
DINGO: Constrained Inference for Diffusion LLMs

Tarun Suresh*, Debangshu Banerjee*, Shubham Ugare, Sasa Misailovic, Gagandeep Singh

Department of Computer Science,
University of Illinois Urbana-Champaign

Abstract

Diffusion LLMs have emerged as a promising alternative to conventional autoregressive LLMs, offering substantial potential for improving runtime efficiency. However, existing diffusion models fail to provably enforce user-specified formal constraints, such as regular expressions, which makes them unreliable for tasks that require structured outputs, such as fixed-schema JSON generation. Unlike autoregressive models, which generate tokens sequentially, diffusion LLMs predict a block of tokens in parallel. This parallelism makes traditional constrained decoding algorithms, designed to enforce constraints with sequential token prediction, ineffective at preserving the true output distribution. To address this limitation, we propose DINGO, a dynamic programming-based constrained decoding strategy that is both efficient and provably distribution-preserving. DINGO enables sampling of output strings with the highest probability under the model’s predicted distribution while strictly adhering to any user-specified regular expression. On standard symbolic math and JSON generation benchmarks, DINGO achieves up to a 68% points of improvement over unconstrained inference. The code is available at [DINGO](#).

1 Introduction

Autoregressive LLMs demonstrate impressive performance across a wide range of tasks, including logical reasoning [Pan et al., 2023], theorem proving [Yang et al., 2023], and code generation [et. al., 2021]. However, because they generate one token at a time, they can be slow when producing long responses. Recent work has explored using diffusion models to accelerate token generation by predicting blocks of tokens in parallel. For tasks such as logical reasoning, where the LLM output is fed into symbolic solvers like Z3 [Fedoseev et al., 2024], syntactic correctness of the output is essential. Prior works [Poesia et al., 2022, Ugare et al., 2024a, Loula et al., 2025] have shown that LLMs frequently make syntactic and semantic errors, often generating structurally invalid outputs that cause downstream tasks to fail due to unparsable input. To mitigate this issue, constrained decoding has emerged as a promising approach that provably ensures structural correctness by projecting the LLM output onto a set of valid strings, typically defined by a regular grammar or, more generally, a context-free grammar (CFG). However, existing constrained decoding techniques are designed specifically for autoregressive LLMs and rely on their step-by-step generation process to prune invalid tokens that cannot lead to structurally valid outputs. At each generation step, the decoder selects the highest-probability token from the set of valid options, based on the LLM’s output distribution.

In contrast, diffusion LLMs predict blocks of tokens in parallel without sequential dependencies, making existing constrained decoding algorithms incompatible. Furthermore, greedy token selection in autoregressive models maximizes the probability locally at each step but can be suboptimal over an entire sequence, potentially leading to structurally valid yet lower-quality outputs that fail to maximize the overall probability of valid strings. Lew et al. [2023], Park et al. [2024b] have reported this distortion in output distribution for autoregressive LLMs under constrained decoding. Therefore,

*Equal contributing authors ordered randomly

any constrained decoding algorithm for diffusion LLMs should also ensure that enforcing formal constraints does not come at the cost of distorting the true output distribution.

Key Challenges: Diffusion LLMs generate a block of tokens starting from a fully masked string composed of special mask tokens \perp , and iteratively unmask one or more tokens at each step until producing a fully unmasked output. Each unmasking step (referred to as a diffusion step) can unmask tokens at arbitrary positions in the block, with no left-to-right sequential dependency across steps. As a result, designing constrained decoding for diffusion LLMs requires addressing the following:

- **RQ1:** Efficiently detecting invalid tokens and restricting token choices at each diffusion step to ensure the final unmasked string is always structurally correct.
- **RQ2:** Ensuring the generated token block maximizes the probability under the output distribution.

Contributions: We present the first constrained decoding algorithm for diffusion LLMs, making the following contributions:

- We introduce DINGO, the first constrained decoding algorithm for diffusion LLMs that supports any user-specified regular expression. DINGO provably ensures that the output string is always a valid prefix of some string in the target regular language.
- DINGO uses dynamic programming to ensure that the output string achieves the maximum probability among all valid strings over the output block with respect to the true output distribution. This approach guarantees scalability while maintaining optimality (e.g., maximizing the probability), in contrast to existing methods such as [Park et al., 2024b], which rely on repeated resampling. Resampling-based methods are computationally expensive and unsuitable for practical deployment.
- Extensive experiments on multiple open-source diffusion LLMs and benchmarks show that DINGO significantly outperforms standard unconstrained decoding, achieving up to a 68% improvement on challenging tasks such as the GSM-symbolic benchmark for symbolic reasoning [Mirzadeh et al., 2024] and a JSON generation benchmark [NousResearch, 2024].

Roadmap: We provide the necessary background in Section 2, formalize constrained decoding for diffusion LLMs in Section 3, describe the DINGO algorithm along with its correctness and optimality proofs in Section 4, and present experimental results in Section 5.

2 Background

Notation: In the rest of the paper, we use small case letters x for constants, bold small case letters (\mathbf{x}) for strings, capital letters X for functions, \cdot for string concatenation, $|\mathbf{x}|$ to denote the length of \mathbf{x} .

Diffusion LLM: The diffusion LLM $\mathcal{L}_{m,n} : V^m \rightarrow V^n$ processes finite strings $\mathbf{p} \in V^m$ over a finite alphabet V including the special mask symbol \perp and produces the output string $\mathbf{o} \in V^n$. Typically $\mathbf{o} = \mathbf{p} \cdot \mathbf{r}$ with length n represents the entire output string of \mathcal{L} where \mathbf{p} is the input prompt, \mathbf{r} is the response, and $m + |\mathbf{r}| = n$. \mathcal{L} can compute the response \mathbf{r} over a single block Austin et al. [2021], Ye et al. [2025], Nie et al. [2025] in pure diffusion setup or over multiple blocks i.e. $\mathbf{r}_1 \cdot \mathbf{r}_2 \cdots \mathbf{r}_k$ in a semi-autoregressive setup where different blocks are computed sequentially from left to right [Han et al., 2023, Arriola et al., 2025].

At a high level, to compute a block of tokens of size d , \mathcal{L} pads the prompt \mathbf{p} with a fully masked suffix, resulting in $\mathbf{p} \cdot \perp^d$, where \perp^d denotes a sequence of d special mask tokens \perp . The model then iteratively unmasks a subset of these tokens at each step, ultimately producing a fully unmasked output string \mathbf{o} . Each such step is referred to as a diffusion step, and \mathcal{L} typically applies T diffusion steps to compute \mathbf{o} . The number of steps T is usually a fixed, predetermined constant satisfying $T < d$, which enables greater scalability compared to their autoregressive counterparts.

Definition 2.1 (Diffusion step). A diffusion step $f_n : V^n \times \mathbb{N} \rightarrow V^n$ applies a single unmasking step to a masked (or, a partially masked) string of length n to compute a new masked (or, possibly unmasked) string of the same length. The first argument represents the input string appended with the output block while the second argument dictates the number of masked tokens in the output string.

Each diffusion step f_n consists of two components: a transformer step $\mathcal{N}_n : V^n \rightarrow \mathbb{R}_+^{|V| \times n}$, which predicts the token probability distribution at each output position, and a mask prediction step

$\mathcal{M}_n : \mathbb{R}_+^{|V| \times n} \times \mathbb{N} \rightarrow \mathbb{R}_+^{|V| \times n}$, which determines which token positions to remark. Typically, for each position, the mask prediction step identifies the token with the highest probability and compares these maximum probabilities across positions. \mathcal{M}_n then greedily remarks positions with relatively lower max-probability scores [Nie et al., 2025] and produces the modified token distribution. Further details about \mathcal{N}_n and \mathcal{M}_n are in Appendix A.

Formally, the diffusion step is defined as $f_n(\mathbf{x}_{i-1}, i) = D_{m,n}(\mathcal{M}_n(\mathcal{N}_n(\mathbf{x}_{i-1}), i))$ where $D_{m,n} : \mathbb{R}_+^{|V| \times n} \rightarrow V^n$ is the decoder. We now use the diffusion step to formally define the diffusion LLM for generating strings of length n in either a single-block or multi-block setting.

Definition 2.2 (Single block diffusion LLM). A diffusion LLM that outputs a block of d tokens given an input $\mathbf{p} \in V^m$ using T diffusion steps is a function $\mathcal{L}_{m,n} : V^m \rightarrow V^n$, where $n = m + d$, and the output is $\mathbf{o} = \mathbf{p} \cdot \mathbf{r} = \mathcal{L}_{m,n}(\mathbf{p})$. Let $f_n : V^n \times \mathbb{N} \rightarrow V^n$ denote a single diffusion step, and let $P_{m,n} : V^m \rightarrow V^n$ be the padding function. Then the output is computed as $\mathbf{o} = \mathcal{L}_{m,n}(\mathbf{p}) = \mathbf{x}_T$, where: $\mathbf{x}_0 = P_{m,n}(\mathbf{p}) = \mathbf{p} \cdot \perp^d$ and $\mathbf{x}_i = f_n(\mathbf{x}_{i-1}, i)$ for $1 \leq i \leq T$.

Definition 2.3 (Semi Autoregressive diffusion LLM). 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 k blocks, where each block is computed via a call to the single block diffusion model. The output of 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 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 all $1 \leq i \leq k$.

DFA and regular expression: We provide necessary definitions regarding regular expression.

Definition 2.4. (DFA) A DFA $D_{\mathcal{R}} = (Q, \Sigma, \delta, q_0, F)$ for a regular expression \mathcal{R} is a finite-state machine that deterministically processes input strings to decide membership in the language $L(\mathcal{R}) \subseteq \Sigma^*$ defined by \mathcal{R} . It consists of states Q , a start state q_0 , a set of accepting states F , and transition rules $\delta : Q \times \Sigma \rightarrow Q$ and the input alphabet Σ .

Definition 2.5 (extended transition function). The extended transition function $\delta^* : \Sigma^* \times Q \rightarrow Q$ maps an input (\mathbf{w}, q) to the resulting state q_r , obtained by sequentially applying δ to each character c_i in $\mathbf{w} = c_1 \cdots c_{|\mathbf{w}|}$, starting from state q .

Definition 2.6 (Live DFA states). Given a DFA $(Q, \Sigma, \delta, q_0, F)$, let Q_l represent the set of live states such that $q \in Q_l$ iff $\exists \mathbf{w} \in \Sigma^*$ s.t. $\delta^*(\mathbf{w}, q) \in F$.

3 Optimal Constrained Decoding

We formalize the correctness and optimality of constrained decoding for any diffusion LLM with respect to a user-defined regular expression \mathcal{R} . Given \mathcal{R} , let $L(\mathcal{R}) \subseteq \Sigma^* \subseteq (V \setminus \perp)^*$ denote the set of all finite strings that satisfy the expression \mathcal{R} .

Correctness: A valid constrained decoding algorithm must ensure that the output string always remains a valid *prefix* of some string in $L(\mathcal{R})$, effectively eliminating any output that cannot be extended into valid completions. By treating the output string as a prefix rather than a fully completed string, we can accommodate the semi-autoregressive setup, where blocks of tokens are appended to the right of the current output. This approach avoids prematurely rejecting strings that may lead to valid completions in subsequent blocks and also aligns with the notion of correctness adopted in existing constrained decoding algorithms for the autoregressive LLM [Ugare et al., 2024b, Banerjee et al., 2025a]. We denote the set of all valid prefixes of $L(\mathcal{R})$ as $L_P(\mathcal{R})$.

Each diffusion step f_n produces a string over the vocabulary V , which may include one or more special mask tokens \perp . These tokens act as placeholders for actual (non-mask) tokens that will be filled in during future diffusion steps. To account for these future substitutions, we define a masked (or partially masked) string as valid if there exists a replacement for all mask tokens such that the resulting fully unmasked string is a valid prefix of some string in $L(\mathcal{R})$. To formalize this notion, we first define the *substitution set*, which represents the set of fully unmasked strings obtained by replacing all mask tokens in a masked or partially masked string. We then use substitution sets to define the correctness of the constrained decoder.

Definition 3.1 (Substitution Set). Given a masked (or, partially masked) string $\mathbf{x} \in V^n$, the *substitution set* $\mathcal{S}(\mathbf{x}) \subseteq (V \setminus \{\perp\})^n$ is the set of all fully unmasked strings obtained by replacing each occurrence of \perp in \mathbf{x} with a token from $V \setminus \{\perp\}$. For unmasked strings with no \perp , $\mathcal{S}(\mathbf{x}) = \{\mathbf{x}\}$

Definition 3.2 (Correctness of Constrained decoder). Any deterministic decoder $D_{m,n,\mathcal{R}} : \mathbb{R}_+^{|V| \times n} \rightarrow V^n$ is a valid constrained decoder if, for all $n \in \mathbb{N}$, input prompt \mathbf{p} and for any output distribution \mathcal{D}_n provided as n probability vectors each of size $|V|$, there exists an unmasked string \mathbf{x} in the substitution set $\mathcal{S}(D_{m,n,\mathcal{R}}(\mathcal{D}_n))$ of the decoded output such that actual response $\mathbf{p} \cdot \mathbf{r} = \mathbf{x}$ is a valid prefix i.e., $\mathbf{r} \in L_P(\mathcal{R})$.²

Optimality: Given a distribution \mathcal{D}_n and a regular expression \mathcal{R} , the set of decodings that are valid prefixes for \mathcal{R} (as defined in Definition 3.2) may not be unique. An optimal constrained decoder selects, among all valid strings, the string that maximizes the probability under \mathcal{D}_n . The output distribution \mathcal{D}_n is represented as n vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$, each of size $|V|$, where the i -th vector \mathbf{v}_i captures the token distribution at position i . For any masked position j , \mathbf{v}_j assigns probability 1 to the mask token \perp and 0 to all other tokens. Assuming the input prompt has length m , the token distribution of the actual response is given by $\mathbf{v}_{m+1}, \dots, \mathbf{v}_n$. For any output string $\mathbf{r} = t_{m+1} \dots t_n$, let $P(\mathbf{r} \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n)$ denote the probability of the string \mathbf{r} under the output distribution. Then, the optimal constrained decoding can be formalized as follows:

$$\mathbf{r}^* = \arg \max_{\mathbf{r}} P(\mathbf{r} \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n) \text{ s.t. } \exists \mathbf{x} \in V^*. (\mathbf{x} \in \mathcal{S}(\mathbf{r})) \wedge (\mathbf{x} \in L_P(\mathcal{R})) \quad (1)$$

Since the token distributions $\mathbf{v}_{m+1}, \dots, \mathbf{v}_n$ are independent across positions, the probability of the string \mathbf{r} can be written as $P(\mathbf{r} \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n) = \prod_{i=m+1}^n \mathbf{v}_i[t_i]$ where $\mathbf{v}_i[t_i]$ denotes the probability assigned to token t_i by the vector \mathbf{v}_i . Using this, we can rewrite the optimization problem from Eq. 1 as follows:

$$\mathbf{r}^* = \arg \max_{\mathbf{r}=t_{m+1} \dots t_n} \prod_{i=m+1}^n \mathbf{v}_i[t_i] \text{ s.t. } \exists \mathbf{x} \in V^*. (\mathbf{x} \in \mathcal{S}(\mathbf{r})) \wedge (\mathbf{x} \in L_P(\mathcal{R})) \quad (2)$$

4 DINGO Algorithm

The search space for Eq. 2 is exponential— $|V|^d$, where $d = n - m$ denotes the block length, making naive enumeration-based methods impractical. To efficiently retrieve the optimal output string \mathbf{r}^* from Eq. 2, DINGO leverages dynamic programming. Given a regular expression \mathcal{R} , it first modifies the transition function to handle the mask symbol \perp , which is then utilized during inference.

4.1 Precomputation

For a user-provided \mathcal{R} and the corresponding DFA $D_{\mathcal{R}} = (Q, \Sigma, \delta_{\mathcal{R}}, q_0, F)$ (referred to as character-level DFA) with $\Sigma \subseteq (V \setminus \perp)$, we first construct the token-level DFA $D_t = (Q, (V \setminus \perp), \delta_t, q_0, F)$ recognizing $L(\mathcal{R})$ over strings generated by \mathcal{L} . A single token $\mathbf{t} \in (V \setminus \perp)$ can span across multiple characters in Σ i.e. $\mathbf{t} = c_1 \dots c_l$ where $c_i \in \Sigma$. To construct the token-level transition function $\delta_t : Q \times (V \setminus \perp) \rightarrow Q$, we process each token $\mathbf{t} \in (V \setminus \perp)$ and state $q \in Q$ by executing the character-level DFA $D_{\mathcal{R}}$ on the sequence of constituent characters $c_1 \dots c_l$, starting from state q , and record the resulting state q_r . We then define the token-level transition as $\delta_t(q, \mathbf{t}) = q_r$.

To handle the special mask token $\perp \in V$, we define the transition function $\delta_{\perp} : Q \rightarrow 2^Q$. For each state $q \in Q$, $\delta_{\perp}(q)$ returns the set of states $Q_r \subseteq Q$ that are reachable via a single token transition using δ_t . Formally, $\delta_{\perp}(q) = \{q' \mid q' = \delta_t(q, \mathbf{t}); \mathbf{t} \in (V \setminus \perp)\}$. Since δ_{\perp} may return multiple states, it resembles the transition function of a non-deterministic finite automaton (NFA). The precomputation step combines δ_t and δ_{\perp} to define $\delta : Q \times V \rightarrow 2^Q$, which is used in the dynamic programming step. Using the token-level DFA D_t , we also construct the set of live states $Q_l \subseteq Q$ (Definition 2.6).

$$\delta(q, \mathbf{t}) = \begin{cases} \{\delta_t(q, \mathbf{t})\} & \text{if } \mathbf{t} \in (V \setminus \perp), \\ \delta_{\perp}(q) & \text{if } \mathbf{t} = \perp. \end{cases}$$

4.2 DINGO Dynamic Programming

Before going into details, we present two key observations that lead to the decoding algorithm.

²More precisely, if there exists at least one \mathbf{r} that is a valid prefix (i.e., $\mathbf{r} \in L_P(\mathcal{R})$), then constrained decoding is always capable of retrieving one of them.

Observation 1: Determining whether a fully unmasked string $\mathbf{r} = t_1 \cdots t_d \in (V \setminus \perp)^*$ is a valid prefix is equivalent to checking whether the resulting state q_r , obtained by applying δ to the sequence $t_1 \cdots t_d$ starting from q_0 , is live. Similarly, for a partially (or fully) masked string \mathbf{r}_\perp , applying δ to $t_1 \cdots t_d$ yields a set of resulting states Q_r . In this case, \mathbf{r}_\perp is a valid prefix if and only if any state $q \in Q_r$ is live (Definition 3.2).

Observation 2: For optimality, it is sufficient to track the maximum probability path from the start state q_0 to each resulting state in Q_r . Once these paths are computed, we select the one with the highest probability that leads to a live state. The corresponding string is the optimal string \mathbf{r}^* (or one of the optimal strings in case of multiple optimal solutions) for the optimization problem in Eq. 2.

Based on these observations, the main challenge is to efficiently maintain the maximum probability path to each reachable state in Q_r . We address this using a dynamic programming (DP) approach, similar to traditional graph-based DP algorithms such as [Forney, 1973].

DP states: For each token position $1 \leq i \leq d$ in the block, the DP maintains: a) $W[i, q]$, which records the maximum probability with which a state $q \in Q$ can be reached from the start state q_0 via transitions on some token sequence with length i ; and b) $Pr[q, i]$, which stores the last transition, i.e., the previous state and the corresponding token, that led to the maximum probability stored in $W[i, q]$. If a state q is unreachable, then $W[i, q] = 0$. Formally, given the probability vectors $\mathbf{v}_1, \dots, \mathbf{v}_i$, $W[i, q]$ is defined as follows where δ_t^* is extended transition function (Definition 2.5).

$$W[i, q] = \max_{t_1 \dots t_i} \prod_{j=1}^i \mathbf{v}_j[t_j] \quad \text{s.t. } q = \delta_t^*(t_{m+1} \cdots t_n, q_0)$$

DP state update: Given the states at token position i , we describe the computation for position $i + 1$. Initially, $W[i, q] = 0$ for all $q \neq q_0$, and $W[i, q_0] = 1$ (lines 1 – 3 in Algo. 1). To compute $W[i + 1, q]$ for each $q \in Q$, we consider all tokens $t \in V$ (including the mask token \perp) that can transition to q from some previous state q' at step i . Among all such transitions, we select the one with the highest probability and add it to the maximum probability path reaching q' at step i . The value $Pr[i + 1, q]$ stores the previous state and token that lead to the maximum probability path to q at step $i + 1$ (lines 12 – 15 in Algo. 1). Formally,

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

Path construction: We consider all reachable states q at the end of the block with $W[d, q] > 0$. Among the live states $q_l \in Q_l$ satisfying this condition, we select the state q_{\max} with the highest value of $W[d, q_l]$. We then use Pr to iteratively reconstruct the token sequence backward that forms the maximum probability path starting from q_{\max} and ending at q_0 (lines 20 – 22 in Algo. 1).

Semi-autoregressive setup: In semi-autoregressive setup, we may not start from DFA start state q_0 since one or more blocks of tokens $\mathbf{r}_1 \cdots \mathbf{r}_l$ may have been generated in left the current block. Provide the string $\mathbf{r}_1 \cdots \mathbf{r}_l$ ends at a live state q_l , we can apply dynamic programming approach with the initialization $W[0, q_l] = 1$ and $W[0, q] = 0$ for all state $q \neq q_l$. Details are in Appendix D.

4.3 Correctness of DINGO

Proposition 4.1. [Correctness] 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 $\exists \mathbf{x} \in V^*. (\mathbf{x} \in \mathcal{S}(\mathbf{r})) \wedge (\mathbf{x} \in L_P(\mathcal{R}))$ holds.

Proof sketch: DINGO ensures that if a state $q \in Q$ is reachable in i tokens, then $W[i, q] > 0$ for all $1 \leq i \leq d$. Since $L_P(\mathcal{R}) \cap (V \setminus \perp)^d \neq \{\}$, there exists a state $q_l \in Q_l$ that is reachable in d steps. Therefore, $W[d, q_{\max}] > 0$ (see line 16 in Alg.1). Consequently, there exists a sequence $\mathbf{x} \in \mathcal{S}(\mathbf{r})$ such that $\delta^*(\mathbf{x}, q_0) = q_{\max} \in Q_l$, implying that $\mathbf{x} \in L_P(\mathcal{R})$. Formal proof is in AppendixB.

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 $\exists \mathbf{x} \in V^*. (\mathbf{x} \in \mathcal{S}(\mathbf{r}')) \wedge (\mathbf{x} \in L_P(\mathcal{R}))$, $P(\mathbf{r}' \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n) \leq P(\mathbf{r}^* \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n)$.

Proof Sketch: Formal proof is in Appendix B.

Algorithm 1 DINGO DP

Require: q_0 , block length d , probability vectors $\mathbf{v}_1, \dots, \mathbf{v}_d$ for the current block, Q_t, 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**
- 7: **for** $(q \in Q)$ **do**
- 8: **for** $t \in V$ **do**
- 9: $q' \leftarrow \delta(q, t)$
- 10: $V_i(q, q'), T_i(q, q') \leftarrow \text{MaxTransition}(\mathbf{v}_i, t, q, q')$
- 11: **for** $i \in \{1, \dots, d\}$ **do** ▷ DP computation loop
- 12: **for** $(q \in Q) \wedge (q' \in Q)$ **do**
- 13: **if** $W[i, q] < W[i-1, q'] \times V_i(q, q')$ **then**
- 14: $W[i, q] \leftarrow W[i-1, q'] \times V_i(q, q')$ ▷ Update maximum probability path to q
- 15: $Pr[i, q] \leftarrow (q', T_i(q, q'))$ ▷ Update the parents accordingly
- 16: $q_{max} \leftarrow \arg \max_{q \in Q_t} W[d, q]$
- 17: **if** $W[d, q_{max}] = 0$ **then** ▷ No valid prefixes
- 18: **return** None, q_{max}
- 19: $\mathbf{r}^* \leftarrow \{\}, q_{curr} \leftarrow q_{max}$
- 20: **for** $i \in \{d, \dots, 1\}$ **do** ▷ Decoding the optimal string \mathbf{r}^*
- 21: $q_{curr}, t \leftarrow Pr[i, q_{curr}]$
- 22: $\mathbf{r}^* \leftarrow \mathbf{r}^* \cdot t$
- 23: **return** $\text{reverse}(\mathbf{r}^*), q_{max}$

4.4 DINGO algorithm

Algorithm 1 presents DINGO steps. The two main loops dominating its computational complexity involve calculating transition costs and performing the DP updates respectively.

First, for each of the d time steps, the algorithm computes the optimal single-token transition costs $V_i(q_s, q_t)$ between all source states $q_s \in Q$ and target states $q_t \in Q$. This is achieved by iterating through each source state q_s , each token $t \in V$, and then for each state q_t reached from q_s via t (i.e., $q_t \in \delta(q_s, t)$), updating the cost $V_i(q_s, q_t)$ with $\mathbf{v}_i[t]$ if it is better. The complexity for this part is $O(d \cdot (|Q|^2 + \sum_{q_s \in Q} \sum_{t \in V} |\delta(q_s, t)|))$. The sum $\sum_{q_s} \sum_t |\delta(q_s, t)|$ represents the total number of transitions, $N_{\text{trans}} = O(|Q| \cdot |V| + |Q| \cdot N_{\perp})$, where N_{\perp} is the maximum number of states reachable via the \perp token. Thus, this part takes $O(d \cdot (|Q|^2 + |Q| \cdot |V|))$.

Second, the core dynamic programming update calculates $W[i, q]$ for each diffusion step i and state q . This involves iterating over d diffusion steps, $|Q|$ current states q , and for each q , considering all $|Q|$ possible previous states q' . This leads to a complexity of $O(d \cdot |Q|^2)$.

Combining these dominant parts, the total complexity is $O(d \cdot (|Q|^2 + |Q| \cdot |V|) + d \cdot |Q|^2)$, which simplifies to $O(d \cdot (|Q|^2 + |Q| \cdot |V|))$. This can be expressed as $O(d \cdot |Q| \cdot (|Q| + |V|))$.

4.5 Context free grammars and reasoning

Currently, DINGO operates with regular grammars. For context-free grammars (CFGs), existing works Ugare et al. [2024b] perform greedy pruning of invalid tokens in a left-to-right manner. The same greedy generation strategy can also be applied to diffusion models, albeit at the cost of optimality. We already have an implementation of greedy left-to-right generation (referred to as Greedy Constrained in Section 5), which naturally extends to CFGs. However, as demonstrated by our experiments, greedy decoding (also employed in more recent works Mündler et al. [2025]) yields significantly worse performance compared to the optimal constrained decoder DINGO for regular expressions. Moreover, DINGO already supports reasoning-augmented grammars Banerjee et al. [2025b], allowing the inclusion of reasoning tokens that do not conform to the output grammar.

5 Experiments

In this section, we evaluate DINGO on a math reasoning task (GSM-Symbolic Mirzadeh et al. [2024]) and a schema-based text-to-JSON task (JSONModeEval [NousResearch, 2024]) and demonstrate significant improvement over baselines. In both tasks, we use the LLaDA-8B-Base (LLaDA-8B-B) Nie et al. [2025], LLaDA-8B-Instruct (LLaDA-8B-I) Nie et al. [2025], Dream-v0-Base-7B (Dream-B-7B) Ye et al. [2025], and Dream-v0-Instruct-7B (Dream-I-7B) Ye et al. [2025] models.

Experimental Setup. We run experiments on a 48-core Intel Xeon Silver 4214R CPU with 2 Nvidia RTX A5000 GPUs. DINGO is implemented using PyTorch Paszke et al. [2019] and the HuggingFace transformers library Wolf et al. [2020]. The token-level DFA is implemented in Rust using a highly efficient regex-DFA library to minimize overhead during DFA construction and LLM inference. We report the mean number of DFA states and transitions as well as the offline pre-computation time in Appendix E.

Baselines. We compare DINGO against unconstrained diffusion LLM generation. Furthermore, to highlight the benefit of optimal constrained decoding with DINGO, we implement a constrained decoding strategy Greedy Constrained that mirrors existing autoregressive constrained generation methods Willard and Louf [2023], Ugare et al. [2024b]. Greedy Constrained iterates over the diffusion block and at each position i computes a binary mask $m \in \{0, 1\}^{|V|}$ based on the DFA, specifying valid tokens ($m = 1$) and excluded tokens ($m = 0$). Decoding is then performed on the masked probability distribution $m \odot v_i$, where \odot denotes element-wise multiplication. Since in some cases, Unconstrained outperforms Greedy Constrained, we also report Best of Greedy + Unconstrained, which takes the better result of the two approaches for each problem in the dataset.

Math Reasoning: We evaluate DINGO on GSM-Symbolic Mirzadeh et al. [2024] dataset, which consists of reasoning-based math world problems where numerical values and names are replaced by symbolic variables. Diffusion LLMs are tasked with generating correct symbolic expression solutions to those word problems. We evaluate correctness by using the Z3 solver [De Moura and Bjørner, 2008] to check if the final expressions from the LLM generations are functionally equivalent to the ground truth expressions. We set the generation length to 128, number of blocks to 8, and total diffusion steps to 64 and prompt the LLMs with 4-shot examples from GSM-Symbolic [Mirzadeh et al., 2024] (the prompts can be found in Appendix F.1). We initialize DINGO and Greedy Constrained with a regex (shown in Appendix F.2) that permits math expressions wrapped in \langle and \rangle and natural language text outside these expressions for reasoning as done in CRANE Banerjee et al. [2025a].

Table 1 compares the performance of DINGO with the baseline methods. The Accuracy (%) column reports the percentage of functionally correct LLM-generated expressions, Parse (%) indicates the percentage of syntactically valid responses (i.e., expressions without invalid operations), and Time provides the average time in seconds taken to generate a completion.

As displayed in the table, DINGO significantly improves functional correctness over the baselines. For instance, for LLaDA-8B-I, DINGO outperforms unconstrained generation by 13 percentage points and Greedy Constrained generation by 5 percentage points. Furthermore, DINGO achieves 100% syntactic accuracy across all models evaluated. On the other hand, unconstrained and Greedy Constrained generation make many syntactic errors, especially for non-instruct tuned models. For these cases, generation with Greedy Constrained results in responses that are syntactically valid prefixes but not syntactically valid by themselves. We present case studies in Appendix F.3. Importantly, DINGO is extremely efficient, introducing marginal overhead compared to unconstrained generation.

JSON Generation: We further evaluate DINGO on a text-to-JSON generation task JSON-ModeEval, which consists of zero-shot problems specifying a JSON schema and a request to generate a JSON object that contains specified contents. Generating JSON that adheres to a specified schema is extremely important for applications like tool use and function calling Ugare et al. [2024b], Willard and Louf [2023]. We evaluate the correctness of JSON generated by an LLM by first evaluating whether the JSON string can be parsed and converted to a valid JSON object. We further evaluate whether the generated JSON is valid against the schema specified in the prompt. We set the generation length to 128, number of blocks to 1, and the total diffusion steps to 64. For the constrained generation methods, we convert each problem’s JSON schema into its corresponding regular expression and guide the diffusion LLM to generate output conforming to that regex.

Table 1: Comparison of constrained and unconstrained generation methods on GSM-Symbolic.

Model	Method	Acc. (%)	Parse (%)	Time (s)
LLaDA-8B-B	Unconstrained	25	54	9.06
	Greedy Constrained	30	75	9.31
	Best of Greedy + Unconstrained	30	75	9.08
	DINGO	31	100	9.22
LLaDA-8B-I	Unconstrained	19	35	23.78
	Greedy Constrained	27	98	23.97
	Best of Greedy + Unconstrained	27	98	23.8
	DINGO	32	100	23.92
Dream-B-7B	Unconstrained	17	33	16.02
	Greedy Constrained	21	41	16.13
	Best of Greedy + Unconstrained	21	41	16.04
	DINGO	23	100	16.19
Dream-I-7B	Unconstrained	32	61	23.89
	Greedy Constrained	34	93	24.01
	Best of Greedy + Unconstrained	34	93	23.9
	DINGO	36	100	23.91

Table 2 presents the results of our experiment. The Parse (%) column reports the percentage of syntactically valid LLM generations while the Accuracy (%) column reports the percentage of generations that are both syntactically valid and valid against their respective schemas. Notably, DINGO achieves 100% schema validation and syntactic accuracy, while baseline methods struggle in many cases to generate valid JSON. We attribute this to the fact that Greedy Constrained may distort the distribution through its greedy approximation and can only generate a valid prefix, not a full parsable generation Park et al. [2024a].

Table 2: Comparison of constrained and unconstrained generation methods for JSON Schema.

Model	Method	Acc. (%)	Parse (%)	Time (s)
LLaDA-8B-B	Unconstrained	57	59	6.37
	Greedy Constrained	80	80	6.47
	Best of Greedy + Unconstrained	88	90	6.41
	DINGO	100	100	6.43
LLaDA-8B-I	Unconstrained	87	91	6.7
	Greedy Constrained	78	79	6.81
	Best of Greedy + Unconstrained	99	99	6.73
	DINGO	100	100	6.78
Dream-B-7B	Unconstrained	15	18	5.31
	Greedy Constrained	23	23	5.41
	Best of Greedy + Unconstrained	32	35	5.34
	DINGO	100	100	5.45
Dream-I-7B	Unconstrained	85	87	6.4
	Greedy Constrained	30	30	6.51
	Best of Greedy + Unconstrained	91	93	6.43
	DINGO	100	100	6.55

Comparison with finetuning and rejection sampling: We compare DINGO with two other possible baselines: (a) fine-tuning on a specific dataset’s JSON schema, and (b) rejection sampling, where we iteratively sample responses and check them against a regular expression (referred to as the Sample-Filter approach). For the fine-tuning approach, we perform supervised fine-tuning (SFT) on a dataset of 20k natural language + JSON schema to JSON instance pairs Azinn [2024]. We used the d1 codebase Zhao et al. [2025] for the SFT implementation. While for the rejection sampling approach,

Table 3: Comparison of DINGO with rejection sampling and finetuning.

Model	Method	Acc. (%)	Parse (%)	Time (s)
LLaDA-8B-B	Unconstrained	57	59	6.37
	Greedy Constrained	80	80	6.47
	Sample-Filter	68	69	13.00
	Format Fine-Tuning	60	60	6.41
	DINGO	100	100	6.43
LLaDA-8B-I	Unconstrained	87	91	6.70
	Greedy Constrained	78	79	6.81
	Sample-Filter	93	97	9.80
	Format Fine-Tuning	89	92	6.72
	DINGO	100	100	6.78

we sampled upto 1000 responses per question. Table 3 shows DINGO significantly outperforms both the baselines.

Ablation Study on The Number of Diffusion Blocks: We analyze the performance of DINGO on GSM-Symbolic using different numbers of diffusion blocks. We run generation with a response length of 128, using 64 total diffusion steps, and each of 1, 2, and 8 blocks. As shown in Figure 1, DINGO performs well across all block settings, outperforming baselines in both functional and syntactic correctness. The ablations on the number of diffusion blocks are presented in Appendix I.

Impact of varying output lengths and different unmasking schedules: We present an ablation study evaluating different output lengths and three different unmasking strategies: a) random, b) top2 margin, and c) entropy-based unmasking strategies Nie et al. [2025] in Appendix L.

6 Related Works

To the best of our knowledge, our work is the first to provide provable guarantees on constrained adherence for inference in diffusion language models. We next discuss the broader set of related works on diffusion language models and constrained language model decoding.

Diffusion Language Models: Diffusion Language Models Austin et al. [2021] have emerged as a promising alternative to traditional autoregressive architectures Radford et al. [2019], offering advantages in parallel processing and controllability while addressing limitations in sequential generation. Recent advances in semi-autoregressive diffusion models Han et al. [2023], Nie et al. [2025], Ye et al. [2025], Arriola et al. [2025] have significantly narrowed the performance gap with autoregressive counterparts. SSD-LM [Han et al., 2023] introduced a semi-autoregressive approach that performs diffusion over the natural vocabulary space, enabling flexible output length and improved controllability by iteratively generating blocks of text while facilitating local bidirectional context updates. More recently, several breakthrough models have advanced the field: LLaDA (Large Language Diffusion with mAsking) achieved competitive performance with SOTA open-source autoregressive models of a similar size like LLaMA3-8B through a forward data masking process and

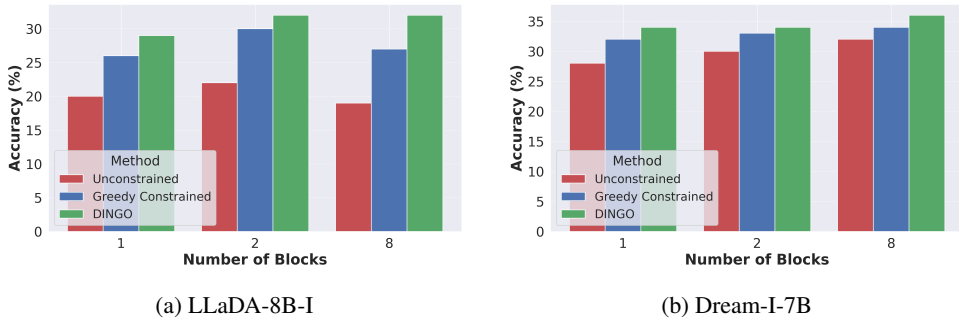


Figure 1: Ablation Study on The Number of Diffusion Blocks For GSM-Symbolic

a reverse process, parameterized by a vanilla Transformer to predict masked tokens [Nie et al., 2025]. BD3-LMs (Block Discrete Denoising Diffusion Language Models) Arriola et al. [2025] introduced a novel approach that interpolates between discrete denoising diffusion and autoregressive models while supporting flexible-length generation and improving inference efficiency with KV caching. Most recently, Dream-7B Ye et al. [2025] emerged as a strong open diffusion large language model that matches state-of-the-art autoregressive (AR) language models of similar size.

Constrained Decoding with Autoregressive LLMs: Constrained decoding has shown promising results in augmenting autoregressive language models. Researchers have developed efficient techniques for ensuring syntactic correctness in regular [Deutsch et al., 2019, Willard and Louf, 2023, Kuchnik et al., 2023] or context-free [Koo et al., 2024, Ugare et al., 2024a, Dong et al., 2024, Banerjee et al., 2025a] languages. Other works have focused on semantically constrained decoding through Monte Carlo sampling [Lew et al., 2023, Loula et al., 2025] or backtracking [Poesia et al., 2022, Ugare et al., 2025]. Lew et al. [2023], Park et al. [2024a] demonstrated that all these approaches that perform greedy constrained approximation for inference can distort the sampling distribution. DINGO addresses this challenge by performing optimal constrained sampling on blocks of tokens in a diffusion language model, which partially mitigates distribution distortion issues.

Concurrent to our work, Cardei et al. [2025] performs constrained sampling from diffusion language models by minimizing a loss function defined using a surrogate model used for scoring constraints. However, their proposed method does not guarantee convergence to the constraint and necessitates a differentiable surrogate model. In contrast, our work focuses on providing provable guarantees for constraint satisfaction during inference without the need of an additional surrogate model.

Limitations DINGO is optimal for per-block generation, making it ideal for pure diffusion settings. However, this optimality may not hold in semi-autoregressive setups involving multiple blocks. Currently, our approach is limited to regular language constraints, while programming languages often belong to context-free or context-sensitive classes. As a result, our method cannot directly enforce these more expressive constraints, which have been addressed in prior work on autoregressive constrained generation. Nonetheless, we believe the core dynamic programming framework behind DINGO can be extended to support richer language classes in future work. Moreover, important constraints like toxicity mitigation fall outside formal language classes, highlighting directions for further research.

7 Conclusion

We presented DINGO, a novel dynamic programming approach that enables diffusion LLMs to generate outputs that strictly adhere to regular language constraints while preserving the model’s underlying distribution. Our method overcomes the limitations of traditional constrained decoding algorithms that fail with parallel token prediction. Our experimental results on symbolic math and JSON generation tasks demonstrate significant improvements over unconstrained inference, demonstrates that DINGO is an effective solution for structured output generation with diffusion models. Our work bridges an important gap in making diffusion LLMs reliable for applications requiring formal guarantees.

Acknowledgement

We thank the anonymous reviewers for their insightful comments. This work was supported by funding through NSF Grants No. CCF-2238079, CCF-2316233, CNS-2148583, CCF-2313028, CCF-2217144 and a Research Gift from Amazon AGI Labs.

References

- Marianne Arriola, Subham Sekhar Sahoo, Aaron Gokaslan, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Justin T Chiu, and Volodymyr Kuleshov. Block diffusion: Interpolating between autoregressive and diffusion language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=tyEyYT267x>.
- Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=h7-XixPCAL>.
- Christian Azinn. Json training dataset. <https://huggingface.co/datasets/ChristianAzinn/json-training>, 2024. Accessed: 2025-10-23.
- Debangshu Banerjee, Tarun Suresh, Shubham Ugare, Sasa Misailovic, and Gagandeep Singh. CRANE: Reasoning with constrained LLM generation. *arXiv preprint arXiv:2502.09061*, 2025a. URL <https://arxiv.org/pdf/2502.09061>.
- Debangshu Banerjee, Tarun Suresh, Shubham Ugare, Sasa Misailovic, and Gagandeep Singh. Crane: Reasoning with constrained llm generation, 2025b. URL <https://arxiv.org/abs/2502.09061>.
- Michael Cardei, Jacob K Christopher, Thomas Hartvigsen, Brian R. Bartoldson, Bhavya Kailkhura, and Ferdinando Fioretto. Constrained language generation with discrete diffusion models, 2025. URL <https://arxiv.org/abs/2503.09790>.
- Leonardo De Moura and Nikolaj Bjørner. Z3: an efficient smt solver. In *Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, TACAS’08/ETAPS’08*, page 337–340, Berlin, Heidelberg, 2008. Springer-Verlag. ISBN 3540787992.
- Daniel Deutsch, Shyam Upadhyay, and Dan Roth. A general-purpose algorithm for constrained sequential inference. In *Proceedings of the Conference on Computational Natural Language Learning*, 2019. URL <https://aclanthology.org/K19-1045/>.
- Yixin Dong, Charlie F Ruan, Yaxing Cai, Ruihang Lai, Ziyi Xu, Yilong Zhao, and Tianqi Chen. XGrammar: Flexible and efficient structured generation engine for large language models. *arXiv preprint arXiv:2411.15100*, 2024. URL <https://arxiv.org/pdf/2411.15100>.
- Chen et. al. Evaluating large language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- Timofey Fedoseev, Dimitar Iliev Dimitrov, Timon Gehr, and Martin Vechev. LLM training data synthesis for more effective problem solving using satisfiability modulo theories. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS’24*, 2024. URL <https://openreview.net/forum?id=hR4Hskr4GX>.
- G. D. Forney. The viterbi algorithm. *Proc. of the IEEE*, 61:268 – 278, March 1973.
- Xiaochuang Han, Sachin Kumar, and Yulia Tsvetkov. Ssd-lm: Semi-autoregressive simplex-based diffusion language model for text generation and modular control, 2023. URL <https://arxiv.org/abs/2210.17432>.
- Terry Koo, Frederick Liu, and Luheng He. Automata-based constraints for language model decoding. In *Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=BDBdblmyzY>.
- Michael Kuchnik, Virginia Smith, and George Amvrosiadis. Validating large language models with RELM. *Proceedings of Machine Learning and Systems*, 5, 2023. URL https://proceedings.mlsys.org/paper_files/paper/2023/file/93c7d9da61ccb2a60ac047e92787c3ef-Paper-mlsys2023.pdf.

- Alexander K Lew, Tan Zhi-Xuan, Gabriel Grand, and Vikash Mansinghka. Sequential Monte Carlo steering of large language models using probabilistic programs. In *ICML 2023 Workshop: Sampling and Optimization in Discrete Space*, 2023. URL <https://openreview.net/pdf?id=U12K0qXxXy>.
- João Loula, Benjamin LeBrun, Li Du, Ben Lipkin, Clemente Pasti, Gabriel Grand, Tianyu Liu, Yahya Emara, Marjorie Freedman, Jason Eisner, Ryan Cotterell, Vikash Mansinghka, Alex Lew, Tim Vieira, and Tim O’Donnell. Syntactic and semantic control of large language models via sequential Monte Carlo. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/pdf?id=xoXn62FzD0>.
- Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models, 2024. URL <https://arxiv.org/abs/2410.05229>.
- Niels Mündler, Jasper Dekoninck, and Martin Vechev. Constrained decoding of diffusion llms with context-free grammars, 2025. URL <https://arxiv.org/abs/2508.10111>.
- Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models, 2025. URL <https://arxiv.org/abs/2502.09992>.
- NousResearch. json-mode-eval, 2024. URL <https://huggingface.co/datasets/NousResearch/json-mode-eval>.
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning, 2023. URL <https://arxiv.org/abs/2305.12295>.
- Kanghee Park, Jiayu Wang, Taylor Berg-Kirkpatrick, Nadia Polikarpova, and Loris D’Antoni. Grammar-aligned decoding. *Advances in Neural Information Processing Systems*, 37:24547–24568, 2024a. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/2bdc2267c3d7d01523e2e17ac0a754f3-Paper-Conference.pdf.
- Kanghee Park, Jiayu Wang, Taylor Berg-Kirkpatrick, Nadia Polikarpova, and Loris D’Antoni. Grammar-aligned decoding. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b. URL <https://openreview.net/forum?id=5G7ve8E1Lu>.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*. 2019.
- Gabriel Poesia, Alex Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. Synchromesh: Reliable code generation from pre-trained language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=KmtVD97J43e>.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI*, 2019. URL https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf. Accessed: 2024-11-15.
- Shubham Ugare, Tarun Suresh, Hangoo Kang, Sasa Misailovic, and Gagandeep Singh. SynCode: Improving LLM code generation with grammar augmentation. *arXiv preprint arXiv:2403.01632*, 2024a. URL <https://arxiv.org/pdf/2403.01632>.
- Shubham Ugare, Tarun Suresh, Hangoo Kang, Sasa Misailovic, and Gagandeep Singh. Syncode: Llm generation with grammar augmentation, 2024b. URL <https://arxiv.org/abs/2403.01632>.
- Shubham Ugare, Rohan Gumaste, Tarun Suresh, Gagandeep Singh, and Sasa Misailovic. IterGen: Iterative structured LLM generation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/pdf?id=ac93gRzxxV>.

- Brandon T Willard and Rémi Louf. Efficient guided generation for large language models. *arXiv preprint arXiv:2307.09702*, 2023. URL <https://arxiv.org/pdf/2307.09702>.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In *Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2020. URL <https://aclanthology.org/2020.emnlp-demos.6>.
- Kaiyu Yang, Aidan M. Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan Prenger, and Anima Anandkumar. Leandojo: Theorem proving with retrieval-augmented language models, 2023. URL <https://arxiv.org/abs/2306.15626>.
- Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong. Dream 7b, 2025. URL <https://hkunlp.github.io/blog/2025/dream>.
- Siyan Zhao, Devaansh Gupta, Qinqing Zheng, and Aditya Grover. d1: Scaling reasoning in diffusion large language models via reinforcement learning, 2025. URL <https://arxiv.org/abs/2504.12216>.

NeurIPS Paper Checklist

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

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

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

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

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

IMPORTANT, please:

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

1. Claims

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

Answer: [Yes]

Justification: Yes, the abstract and introductions claims are validated through the proposed approach and experimental results.

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: Limitations are discussed in the Related Works section.

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: [\[Yes\]](#)

Justification: Yes, refer to Section 4.

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: All experiment details are in Section 5.

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: Code is provided.

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, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

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

Justification: All experiment details are mentioned in Section 5.

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: [\[No\]](#)

Justification: All experiments used greedy decoding, which is deterministic. The results will not change between runs.

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: Yes, the details are provided in Section 5.

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: The paper conforms to 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: The paper discusses the impact of the work in the Introduction and validate through the experiments.

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: 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: All citations for assets used are given.

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 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 details about the algorithm and limitations are described in the paper

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: No 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: No human subject involved.

Guidelines:

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

16. Declaration of LLM usage

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

Answer: [No]

Justification: LLM is not used for important, original, or non-standard component of the research.

Guidelines:

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

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

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

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

B Proofs

Proposition 4.1. [Correctness] 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 $\exists \mathbf{x} \in V^* . (\mathbf{x} \in \mathcal{S}(\mathbf{r})) \wedge (\mathbf{x} \in L_P(\mathcal{R}))$ holds.

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 string \mathbf{r} satisfy the soundness property (see Definition 3.2). In other words, if there is at least one fully unmasked valid prefix with non-zero probability then DINGO retrieves a valid string.

We show this by induction on the position of tokens. Before moving to the proof, we first define extended transition function δ^* when $\delta : Q \times V \rightarrow 2^Q$ outputs a set of states instead of single 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 reachable states starting from q_0 . This can be defined as $\delta^*(\{\}, q_0) = \{q_0\}$ and $\delta^*(t_1 \dots t_{m+1}, q_0) = \cup_{q \in \delta^*(t_1 \dots t_m, q_0)} \delta(q, t_{m+1})$.

1. Let $0 \leq i \leq d$, and let $t_1 \dots t_i \in V^i$ denote any token sequence with positive probability 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 induction on i .
 - (a) Base case $i = 0$: For empty strings only start state q_0 is reachable. DINGO initializes $W[0, q_0] = 1 > 0$ and for all $q \neq q_0$, $W[0, q] = 0$. (lines 1 – 3 in Algo. 1).
 - (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 $q' \in \delta^*(t_1 \dots t_i, q_0)$ and $q \in \delta(q', t_{i+1})$. By the inductive hypothesis, for all such q'

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

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

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

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 $\arg \max$, $W[d, q_{max}] \geq W[d, q_l] > 0$.
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, $W[d, q_{max}] > 0$.
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 $\mathbf{x} \in L_P(\mathcal{R})$.

□

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 $\exists \mathbf{x} \in V^*. (\mathbf{x} \in \mathcal{S}(\mathbf{r}')) \wedge (\mathbf{x} \in L_P(\mathcal{R}))$, $P(\mathbf{r}' \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n) \leq P(\mathbf{r}^* \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n)$.

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 $W[i, q] = \prod_{j=1}^i \mathbf{v}_{m+j}[\mathbf{r}_j^*]$.

- (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, $W[i, q] = \prod_{j=1}^{i-1} \mathbf{v}_{m+j}[\mathbf{r}_j^*] \times \mathbf{v}_{m+i}(\mathbf{r}_i^*) = \prod_{j=1}^i \mathbf{v}_{m+j}[\mathbf{r}_j^*]$.

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

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 $\prod_{j=1}^i \mathbf{v}_{m+j}[\mathbf{r}_j'] \leq W[d, q]$ using induction on i .

- (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,

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

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

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 line 16 in Algo. 1, $q_{max} = \arg \max_{q \in Q_l} W[d, q]$. Thus, by the definition of $\arg \max$, $W[d, q_d] \leq W[d, q_{max}]$. From the inductive hypothesis above, $\prod_{j=1}^d \mathbf{v}_{m+j}[\mathbf{r}_j'] \leq W[d, q_d] \leq W[d, q_{max}]$.

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

□

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  $\text{None}, q_{max}$ 
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  $\text{reverse}(\mathbf{r}^*), q_{max}$ 

```

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

D Semi-Autoregressive

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 k blocks, where each block is computed via a call to the single block diffusion model. The output of 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 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 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  $\mathbf{r}$ 
```

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

E Token Transitions Statistics

Table 4: 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

In Table 4, we report the precomputation time and the number of states in the DFA for both tasks. For JSON generation, different regular expressions are used for different schemas; therefore, we report the mean precomputation time and mean number of states. The maximum number of states and precomputation times across all questions are 455 and 17.7 (Dream) 21.3 (LLaDA) seconds, respectively.

F GSM-Symbolic

F.1 GSM-Symbolic Prompt

```

You are an expert in solving grade school math tasks. You will be presented
with a grade-school math word problem with symbolic variables and be
asked to solve it.

Before answering you should reason about the problem (using the <reasoning>
field in the response described below). Intermediate symbolic expressions
generated during reasoning should be wrapped in << >>.

Only output the symbolic expression wrapped in << >> that answers the
question. The expression must use numbers as well as the variables
defined in the question. You are only allowed to use the following
operations: +, -, /, //, %, *, and **.

You will always respond in the format described below:
Let's think step by step. <reasoning> The final answer is <<symbolic
expression>>

There are {t} trees in the {g}. {g} workers will plant trees in the {g} today
. After they are done, there will be {tf} trees. How many trees did the {
g} workers plant today?

Let's think step by step. Initially, there are {t} trees. After planting,
there are {tf} trees. The number of trees planted is <<tf - t>>. The
final answer is <<tf - t>>.

If there are {c} cars in the parking lot and {nc} more cars arrive, how many
cars are in the parking lot?

Let's think step by step. Initially, there are {c} cars. {nc} more cars
arrive, so the total becomes <<c + nc>>. The final answer is <<c + nc>>.

{p1} had {ch1} {o1} and {p2} had {ch2} {o1}. If they ate {a} {o1}, how many
pieces do they have left in total?

Let's think step by step. Initially, {p1} had {ch1} {o1}, and {p2} had {ch2}
{o1}, making a total of <<ch1 + ch2>>. After eating {a} {o1}, the
remaining total is <<ch1 + ch2 - a>>. The final answer is <<ch1 + ch2 - a
>>.

{p1} had {l1} {o1}. {p1} gave {g} {o1} to {p2}. How many {o1} does {p1} have
left?

Let's think step by step. {p1} started with {l1} {o1}. After giving {g} {o1}
to {p2}, {p1} has <<l1 - g>> {o1} left. The final answer is <<l1 - g>>.

{question}

```

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


[illegible]

F.3 GSM-Symbolic Case Studies


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

Greedy Constrained: Let's think step by step.
 ``Reasoning Text ...`` of emails requiring require a
 response is <<b - (c/100* b ** b * a) ** a>>. The
 final answer is <<b - (c/100* b ** b * a) ** a>>.

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 

Case Study 2:

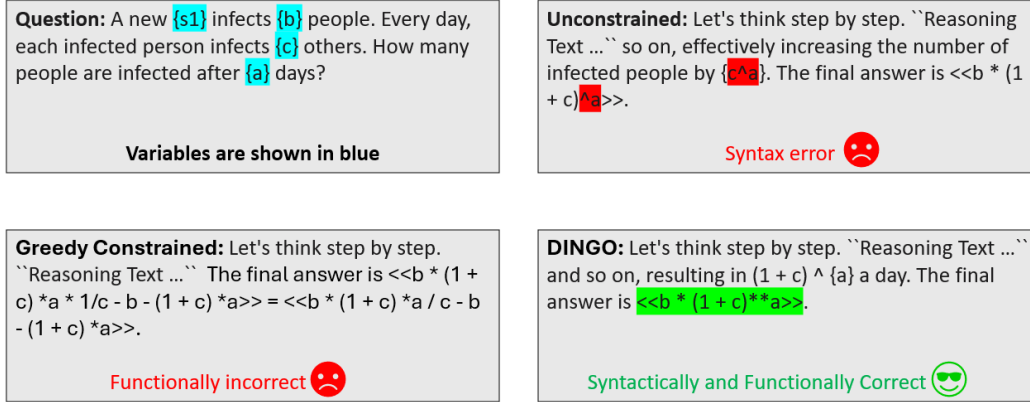


Figure 3: An example from the GSM-symbolic dataset (variables in blue), where unconstrained generation produces syntactically incorrect output, and greedy constrained generation yields a syntactically valid but incorrect answer. In contrast, DINGO generates the correct answer.

G JSON-Mode

G.1 JSON-Mode Example Prompt

```
You are a helpful assistant that answers in JSON. Here's the json schema you must adhere to:
<schema>
{'title': 'PromotionalCampaign', 'type': 'object', 'properties': {'campaignID': {'title': 'Campaign ID', 'type': 'string'}, 'productID': {'title': 'Product ID', 'type': 'string'}, 'startDate': {'title': 'Start Date', 'type': 'string', 'format': 'date'}, 'endDate': {'title': 'End Date', 'type': 'string', 'format': 'date'}, 'discountDetails': {'title': 'Discount Details', 'type': 'string'}}, 'required': ['campaignID', 'productID', 'startDate', 'endDate']}
</schema>

I'm organizing a promotional campaign for our new eco-friendly laundry detergent, which is part of our household products line. The campaign will start on June 1, 2023, and end on June 30, 2023. We're planning to offer a 15% discount on all purchases during this period. The campaign ID is CAMP123456, the product ID is PROD7891011, and the discount details are 15% off on all purchases.

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

```
\\{[ ]?"campaignID"[ ]?:[ ]?"([~\\\\\\\\\\\\\\\\x00-\\\\x1F\\\\x7F-\\\\x9F]|\\\\\\\\[~\\\\\\\\\\\\\\\\])*" [ ]?,[ ]?"productID"[ ]?:[ ]?"([~\\\\\\\\\\\\\\\\\\\\x00-\\\\x1F\\\\x7F-\\\\x9F]|\\\\\\\\[~\\\\\\\\\\\\\\\\])*" [ ]?,[ ]?"startDate"[ ]?:[ ]?"(?:\\\\d{4})-(?:0[1-9]|1[0-2])-(?:0[1-9]|1[1-2][0-9]|3[0-1])"[ ]?,[ ]?"endDate"[ ]?:[ ]?"(?:\\\\d{4})-(?:0[1-9]|1[0-2])-(?:0[1-9]|1[1-2][0-9]|3[0-1])"([ ]?,[ ]?"discountDetails"[ ]?:[ ]?"([~\\\\\\\\\\\\\\\\\\\\x00-\\\\x1F\\\\x7F-\\\\x9F]|\\\\\\\\[~\\\\\\\\\\\\\\\\])*" )?[ ]?\\}
```

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: <pre><schema> ... </schema></pre> <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}]}
Greedy Constrained: { "orderId" : "4", 'items': [{ 'name': 'Wireless Mouse', 'price': 15.99 }, { 'name': 'Bluetooth Keyboard', 'price': 29.99 }, { 'name': 'Webcam HD', 'price': 45.50 }] }	DINGO: {'orderId': 'ORD10234', 'items': [{ 'name': 'Wireless Mouse', 'price': 15.99 }, { 'name': 'Bluetooth Keyboard', 'price': 29.99 }, { 'name': 'Webcam HD', 'price': 45.50 }] }}

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.

<p>Question: Here's the json schema you must adhere to: <code><schema> ... </schema></code> 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.</p>	<p>Unconstrained: <code>\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."}</code></p> <p>Syntax error (missing '{') 😞</p>
<p>Greedy Constrained: <code>{ "customerID": "1UST29382", "feedbackScore": 4, \n\n\n\n"comments"\n\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</code></p> <p>Syntax error (incomplete output) 😞</p>	<p>DINGO: <code>{'customerID': 'CUST29382', 'feedbackScore': 4, 'comments': 'The service was very prompt and professional, but I wish there was more variety in the product options.'}</code></p> <p>Syntactically Correct 😊</p>

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.

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

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

Table 5: 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
Dream-I-7B	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
	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

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

We run generation with a response length of 128, using 64 total diffusion steps, and each of 1, 2, and 8 blocks. Table 6 presents the result.

Table 6: 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

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

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

Table 7: 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

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

We run generation with a response length of 128, 1 block, and each of 16, 32, 64, and 128 total diffusion steps. Table 8 presents the result.

Table 8: 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

L Impact of varying output lengths and different unmasking schedules

Table 9: Ablation Study for Different Output Lengths for JSON-Mode. All experiments use a single block, 64 diffusion steps, and generation length of 128.

Model	Output Length	Method	Acc. (%)	Parse (%)	Time (s)
LLaDA-8B-I	128	Unconstrained	87	91	6.70
	128	Greedy Constrained	78	79	6.81
	128	Best of Greedy + Unconstrained	99	99	6.73
	128	DINGO	100	100	6.78
	256	Unconstrained	84	87	20.25
	256	Greedy Constrained	74	76	20.62
	256	Best of Greedy + Unconstrained	95	96	20.29
	256	DINGO	100	100	20.39
	512	Unconstrained	84	84	45.77
	512	Greedy Constrained	79	80	46.24
	512	Best of Greedy + Unconstrained	96	96	45.83
	512	DINGO	100	100	46.05
Dream-I-7B	128	Unconstrained	85	87	6.40
	128	Greedy Constrained	30	30	6.51
	128	Best of Greedy + Unconstrained	91	93	6.43
	128	DINGO	100	100	6.55
	256	Unconstrained	88	89	17.73
	256	Greedy Constrained	26	26	18.02
	256	Best of Greedy + Unconstrained	93	94	17.75
	256	DINGO	100	100	17.86
	512	Unconstrained	90	90	36.88
	512	Greedy Constrained	25	25	37.28
	512	Best of Greedy + Unconstrained	95	95	36.90
	512	DINGO	100	100	36.09

Table 10: Ablation Study for Different Unmasking Schedules for JSON-Mode.

Model	Unmasking Type	Method	Acc. (%)	Parse (%)	Time (s)
LLaDA-8B-I	Low-confidence	Unconstrained	87	91	6.70
	Low-confidence	Greedy Constrained	78	79	6.81
	Low-confidence	Best of Greedy + Unconstrained	99	99	6.73
	Low-confidence	DINGO	100	100	6.78
	Random	Unconstrained	80	89	6.52
	Random	Greedy Constrained	93	95	6.74
	Random	Best of Greedy + Unconstrained	98	98	6.53
	Random	DINGO	100	100	6.76
	Top2 Margin	Unconstrained	80	81	7.03
	Top2 Margin	Greedy Constrained	98	98	7.17
	Top2 Margin	Best of Greedy + Unconstrained	100	100	7.06
	Top2 Margin	DINGO	100	100	7.11
	Entropy	Unconstrained	86	89	6.89
	Entropy	Greedy Constrained	71	72	7.03
	Entropy	Best of Greedy + Unconstrained	97	98	6.91
	Entropy	DINGO	100	100	7.01
Dream-I-7B	Low-confidence	Unconstrained	78	82	6.32
	Low-confidence	Greedy Constrained	41	41	6.46
	Low-confidence	Best of Greedy + Unconstrained	90	92	6.33
	Low-confidence	DINGO	100	100	6.42
	Random	Unconstrained	55	68	6.17
	Random	Greedy Constrained	30	30	6.28
	Random	Best of Greedy + Unconstrained	68	77	6.18
	Random	DINGO	100	100	6.30
	Top2 Margin	Unconstrained	76	80	6.71
	Top2 Margin	Greedy Constrained	40	40	6.85
	Top2 Margin	Best of Greedy + Unconstrained	89	91	6.73
	Top2 Margin	DINGO	100	100	6.83
	Entropy	Unconstrained	85	87	6.40
	Entropy	Greedy Constrained	30	30	6.51
	Entropy	Best of Greedy + Unconstrained	91	93	6.43
	Entropy	DINGO	100	100	6.55