STOPPING COMPUTATION FOR CONVERGED TOKENS IN MASKED DIFFUSION-LM DECODING

Anonymous authors

Paper under double-blind review

ABSTRACT

Masked Diffusion Language Models generate sequences via iterative sampling that progressively unmasks tokens. However, they still recompute the attention and feed-forward blocks for every token position at every step—even when many unmasked tokens are essentially fixed, resulting in substantial waste in compute. We propose SURELOCK: when the posterior at an unmasked position has stabilized across steps (our *sure* condition), we *lock* that position—thereafter skipping its query projection and feed-forward sublayers—while caching its attention keys and values so other positions can continue to attend to it. This reduces the dominant per-iteration computational cost from $O(N^2d)$ to O(MNd) where N is the sequence length, M is the number of unlocked token positions, and d is the model dimension. In practice, M decreases as the iteration progresses, yielding substantial savings. On LLaDA-8B, SureLock reduces algorithmic FLOPs by 30–50% relative to the same sampler without locking, while maintaining comparable generation quality. We also provide a theoretical analysis to justify the design rationale of SureLock: monitoring only the local KL at the lock step suffices to bound the deviation in final token probabilities.

1 Introduction

Discrete diffusion language models (DLMs) generate text by iteratively denoising a discrete sequence over T steps (Li et al., 2022; Schiff et al., 2024); While its formulation varies (e.g., token swap, insertion, or masking), a widely used family operates through *masking and unmasking*—Masked Diffusion Language Models (MDLMs) (Sahoo et al., 2024; Shi et al., 2024; Nie et al., 2025; Arriola et al., 2025). Unlike autoregressive (AR) decoding (Vaswani et al., 2017)—whose per-step compute naturally grows by one query token via a KV-cache—standard diffusion-style sampling repeatedly recomputes self-attention and per-token feed-forward sublayers for *all* N token positions at *every* step—even for tokens that have already been unmasked and considered to be stable. Hence, the per-block cost is dominated by computing the attention scores QK^{\top} and applying them to V, i.e., $O(N^2d)$ where d is the model dimension. It leads to substantial waste in compute.

To address the computational challenges of DLM sampling, prior work has mainly progressed along two axes: Temporal approaches shrink the step count—e.g., parallel/dilated unmasking and staged or learned samplers—thereby requiring fewer refinement (Luxembourg et al., 2025; Israel et al., 2025; Wei et al., 2025); Reuse approaches reduce per-step compute by reusing or approximating K/V vectors and partially updating hidden features across steps (Ma et al., 2025; Liu et al., 2025b; Wu et al., 2025a). These choices chiefly either reduce the step count T or amortize work across steps by reusing intermediate states; they do not alter the within-step spatial granularity: each step still issues N query rows, so the attention-dominated cost remains $O(N^2d)$ even in late iterations.

Beyond reducing the step count T or reusing K/V across steps, we pursue a largely orthogonal axis: permanently and monotonically deactivating token positions as the sampling unfolds. We propose SURELOCK: once a token's posterior has stabilized, we lock that position—we cache K/V vectors and therefore skip its Q-projection and per-token feed-forward sublayers. Active token positions can still attend to the locked ones via the cached K/V (Sec. 2). The per-block cost becomes $O(M_tNd)$

 $^{^{1}}$ In this paper, we use *position* to refer to the a fixed token slot i in the length-N sequence; a *position* exists even while its token remains masked.

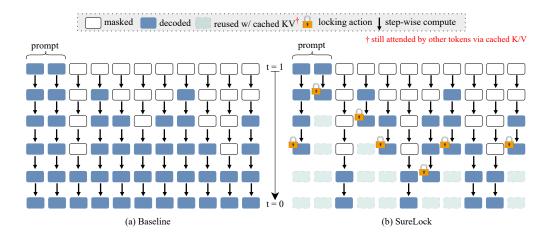


Figure 1: **Iterative sampling by a normal sampler and SureLock.** (a) Baseline consistently recomputes attention scores and FFN sublayers for every token position at every step even after the marginal tokens have become unmasked. (b) SureLock permanently stops recomputing for *locked* positions once these positions are locked. Via cached K/V, other tokens still attend to locked tokens.

for attention and $O(M_t d^2)$ for FFN, instead of $O(N^2 d)$ and $O(N d^2)$, yielding a monotonically decreasing per-step compute profile as M_t shrinks. We lock token positions based on their local stability; a criterion whether the per-step KL divergence of the generative probability distribution falls below a threshold ε . In order to justify using KL as the locking signal, we derive a closed-form bound, which links a per-step KL at locking time to the terminal log-prob deviation². This axis is complementary to the Temporal and Reuse approaches; existing methods continue to operate on the remaining active token positions.

Our work is most closely related to selection-based methods (e.g., DLLM-CACHE (Liu et al., 2025b)) lower per-step compute by updating only a subset of positions. SURELOCK, by contrast, answers a different question—what to remove from compute permanently, instead of what to compute now—so the active position set contracts over time; their selection target of step-wise compute can be tapered down monotonically.

We evaluate SURELOCK on representative MDLMs—LLaDA-8B-Base/-Instruct (Nie et al., 2025)—in Sec. 3. Across diverse decoding settings (e.g., sequence length, locking threshold), perstep algorithmic FLOPs are monotonically decreasing. On continuation generation with WikiText-103 (Merity et al., 2016) and instruction following with MT-Bench (Zheng et al., 2023), SURELOCK reduces algorithmic FLOPs up to 50% in our runs without an expense of generation quality.

2 SureLock

We consider MDLMs that iteratively denoise a length-N sequence, producing intermediate sequences over steps $t=1,\ldots,T$. At step t, an L-block Transformer yields token-wise logits $z_t^{(i)}$ and posteriors $p_t^{(i)}=\operatorname{softmax}(z_t^{(i)})$ at token position i, together with hidden states $h_t^{(i)}\in\mathbb{R}^d$. Let $\mathcal{M}_t\subseteq [N]$ denote the masked positions (prediction targets) and $\bar{\mathcal{M}}_t=[N]\setminus\mathcal{M}_t$ the unmasked ones. In the standard dLLM sampler, each step recomputes self-attention and FFN for all tokens.

In the following, we explain the proposed method **SURELOCK** (Algorithm 1). Once a token's posterior has stabilized at a step, it locks that token position to stop step-wise computations at all subsequent steps by bypassing sublayers and caching K/V values (Sec. 2.1). Our locking criterion is based on the step-wise KL divergence of the posterior (Sec. 2.2). We explain the rationale of using the KL as locking criterion from a theoretical perspective (Sec. 2.3).

²It is deliberately conservative and is intended as a *design rationale*, not as a calibrated predictor for some specific models and datasets

Algorithm 1 SureLock

108

137 138

139 140

141

142

143

144

145

146

147

148 149

150 151

152

153

154 155

156

157 158 159

160

161

```
109
                 Require: sequence length N, total steps T, confidence threshold percentile m, KL threshold \varepsilon
110
                 Require: vocabulary set \mathcal{V}, number of layers L, an unmasking policy UpdateMASK(\cdot)
111
                   1: State: boolean lock \in \{0,1\}^N (init: all 0), caches \mathcal{C} for (k,v) (init: None)
112
                  2: State: previous posteriors p_{t-1} (init: None), masked indices \mathcal{M}_t (init: [N])
                  3: State: frozen block-input \hat{x} \in \mathbb{R}^{N \times d} (init: prompt embeddings)
113
                  4: Notation: X_t \in \mathbb{R}^{N \times d} denotes the model input at t (from embeddings or previous step output)
114
115
                  5: for t = 1 to T do
                                                                                                                                                                          ⊳ one diffusion step
116
                                A_t \leftarrow \{i \mid \text{lock}_i = 0\}, \quad \mathcal{L}_t \leftarrow \{i \mid \text{lock}_i = 1\} \Rightarrow active and locked position at t
                                Initialize block input x^{(1)}[\mathcal{A}_t] \leftarrow X_t[\mathcal{A}_t], x^{(1)}[\mathcal{L}_t] \leftarrow \hat{x}[\mathcal{L}_t]
117
                  7:
                                Compute on active positions i \in A_t:
                  8:
118
                  9:
                                        for \ell = 1 to L
119
                                             Q[\mathcal{A}_t], K[\mathcal{A}_t], V[\mathcal{A}_t] \leftarrow \operatorname{Proj}_{Q,K,V}(\operatorname{LN}(x^{(l)}[\mathcal{A}_t]))
                                                                                                                                                                            10:
120
                                             K^{\mathrm{all}}[\mathcal{A}_t] \leftarrow K[\mathcal{A}_t], \ K^{\mathrm{all}}[\mathcal{L}_t] \leftarrow \mathcal{C}.k[\mathcal{L}_t]; \ \ \text{same for } V^{\mathrm{all}}
                 11:
                                                                                                                                                                                           ▷ assemble
121
                                             h^{(l)}[\mathcal{A}_t] \leftarrow \text{FFN}(\text{Attn}(Q[\mathcal{A}_t], K^{\text{all}}, V^{\text{all}}))
x^{(l+1)}[\mathcal{A}_t] \leftarrow h^{(l)}[\mathcal{A}_t], \quad x^{(l+1)}[\mathcal{L}_t] \leftarrow x^{(l)}[\mathcal{L}_t]
122
                 12:

    b attention with FFN

                 13:
                                                                                                                                                           ⊳ locked rows pass through
123
                 14:
124
                                Obtain posteriors p_t[\mathcal{A}_t] \leftarrow \operatorname{softmax}(\operatorname{Proj}_{out}(x^{(L+1)}[\mathcal{A}_t]))

Unmasking: \mathcal{M}_t \leftarrow \operatorname{UpdateMASK}(p_t, \mathcal{M}_t), \quad \bar{\mathcal{M}}_t \leftarrow [N] \setminus \mathcal{M}_t \quad \triangleright \text{ unmask with posterior}
                 15:
125
                 16:
126
                                Score on active positions i \in A_t:
                 17:
127
                               u_t^{(i)} \leftarrow 1 - \max_{v \in \mathcal{V}} p_t^{(i)}(v) \text{ for } i \in \mathcal{A}_t
D_t^{(i)} \leftarrow \mathrm{KL}(p_t^{(i)} \| p_{t-1}^{(i)}) \text{ if } t > 1 \text{ else } \infty
\mathbf{Locking \ candidates:} \ \mathcal{Z}_t \leftarrow \mathcal{A}_t \ \cap \ \bar{\mathcal{M}}_t
\theta_m \leftarrow \mathrm{Percentile} \big( \{u_t^{(j)}\}_{j \in \mathcal{Z}_t}, m \big)
                 18:
                                                                                                                                                                            ⊳ confidence value
128
                 19:
                                                                                                                                                                                   ⊳ step-wise KL
129
                                                                                                                                                       ⊳ must be active & unmasked
130
                 20:
131
                 21:

    b threshold over candidate positions

132
                 22:
                                      \begin{array}{l} \mathcal{F}_t \leftarrow \{ \ i \in \mathcal{Z}_t \ | \ u_t^{(i)} \leq \theta_m \ \land \ D_t^{(i)} \leq \varepsilon \} \\ & \text{locking evaluation} \\ & \text{lock}[\mathcal{F}_t] \leftarrow 1; \ \ \mathcal{C}_\ell.\{K,V\}[\mathcal{F}_t] \leftarrow \{K_\ell,V_\ell\}[\mathcal{F}_t] \ \ \forall \ell \\ \end{array} \\ & \text{$\lor$ update locked indices and cache} 
                 23:
133
                 24:
134
                                      \hat{x}[\mathcal{F}_t] \leftarrow x^{(1)}[\mathcal{F}_t]
                 25:
                                                                                                                                           135
                 26: end for
136
```

2.1 PERMANENTLY STOPPING STEP-WISE COMPUTE AND CACHING K/V

Once a token i is deemed converged at step t^* , we $permanently\ eliminate$ its position from perstep compute: we **cache** its K/V vectors, and **bypass** its query projection and per-token FFN at all subsequent steps. Let \mathcal{L}_t be the set of indices locked by step t, and define the active set $\mathcal{A}_t := \bar{\mathcal{M}}_t \setminus \mathcal{L}_t$ (unmasked and not yet locked). At step t we assemble the queries $Q[\mathcal{A}_t]$ and run a variable-length attention kernel against the full key/value tables K^{all} and V^{all} where $K^{\mathrm{all}}[\mathcal{L}_t]$, $V^{\mathrm{all}}[\mathcal{L}_t]$ are read from cache. This yields attention outputs only for active positions $\in \mathcal{A}_t$, and we apply the FFN sublayers only to them. Moreover, for locked indices $i \in \mathcal{L}_t$ we keep their predictions fixed, i.e., $p_t^{(i)} \leftarrow p_{t^*}^{(i)}$, and skip their FFN computation. Note that locked positions continue to be attended by other tokens via cached K/V.

2.2 Criterion for Locking: Step-wise KL Divergence

Our locking rule is driven primarily by *local KL*; we will justify the validity of this design in Theorem 1 in Sec. 2.3. We also apply a *confidence* gating as a secondary safeguard to prefer more confident tokens with peaked posteriors. Disabling the confidence gate leaves the theorem unchanged.

Primary: Local KL divergence. For position i at step t, we define the one-step divergence as

$$D_t^{(i)} \triangleq \mathrm{KL}\left(p_t^{(i)} \parallel p_{t-1}^{(i)}\right).$$

Optional: Confidence gate. Let uncertainty $u_t^{(i)} = 1 - \max_{v \in \mathcal{V}} p_t^{(i)}(v)$ and $q_m(u_t)$ the empirical m-th percentile of $\{u_t^{(j)}: j \in \mathcal{A}_t\}$. When enabled, the gate accepts token position i as a locking candidate, if $u_t^{(i)} \leq q_m(u_t)$, i.e., top-m% confidence among active tokens.

Locking rule. We lock token position i at step t^* , if

 $D_{t^{\star}}^{(i)} \leq \varepsilon$ and, if the confidence gate is enabled, $u_{t^{\star}}^{(i)} \leq q_m(u_{t^{\star}})$.

The primary criterion alone suffices Theorem 1. Upon locking we cache $(k_{t^{\star}}^{(i)}, v_{t^{\star}}^{(i)})$ and remove position i from $\mathcal{A}_{>t}$ thereafter.

2.3 DESIGN JUSTIFICATION

The purpose of this section is *design-theoretic*: we justify the use of a local KL threshold as a locking (freezing) criterion by proving that it upper-bounds the terminal error of the log-probability in closed form. The controller based on SureLock locks position i at step t^\star , bypassing Q-projection and FNN sublayers, and reusing its cached K/V thereafter (Alg. 1). We establish that the resulting error of the terminal log-probabilities, relative to an alternative identical sampler without locking, is bounded by $\delta = C_{\text{tail}}\sqrt{\varepsilon}$ under some assumptions, where C_{tail} depends on operator-norm constants of the model and ε is the threshold for determining lock (i.e., $D_{t^\star}^{(i)} \leq \varepsilon$). Therefore, the local KL criterion immediately converts to an explicit design for the terminal error ($\varepsilon(\delta) = \delta^2/C_{\text{tail}}^2$). In the following, $p_t^{(i)}$ and $z_t^{(i)}$ denote the posterior and logits at position i after step t in the no-lock run, while $\hat{p}_t^{(i)}$ denotes the corresponding posterior when i is locked at t^\star .

Note. We emphasize that we here justify the use of a local KL threshold as a locking criterion by proving that it upper-bounds the terminal error in closed form, rather than to predict an empirical error value for a given specific setting.

Standing assumptions for Theorem 1. Let z the logit and $f(z) = \log \operatorname{softmax}(z)$. Fix constants L > 0, $L_{\text{sm}} > 0$, and $\rho \in (0, 1)$. For any position i and step s, the following hold:

- (A1) Locking semantics. Once i is locked at t^* , it is excluded from subsequent re-masking.
- (A2) Geometric tail contraction. $D_s^{(i)} \leq \rho D_{s-1}^{(i)}$ for $s > t^*$.
- (A3) One-step logit smoothness. $||z_s^{(i)} z_{s-1}^{(i)}||_2 \le L\sqrt{D_{s-1}^{(i)}}$ (derivation in Appendix C).
- (A4) Log-softmax Lipschitzness. $||f(z) f(z')||_{\infty} \le L_{\text{sm}}||z z'||_2$.

Theorem 1 (Locking error bound and closed-form threshold). Fix a position i that is unlocked up to t^* and then locked (Alg. 1). Under (A1)–(A4), for any terminal time $T > t^*$,

$$\left\|\log p_T^{(i)} - \log \hat{p}_T^{(i)}\right\|_{\infty} \leq C_{\mathrm{tail}} \sqrt{D_{t^*}^{(i)}}, \qquad C_{\mathrm{tail}} \coloneqq L_{\mathrm{sm}} L/(1 - \sqrt{\rho}).$$

In particular, if the locking test enforces $D_{t^*}^{(i)} \leq \varepsilon$, then the terminal log-probability error is at most $\delta = C_{\text{tail}} \sqrt{\varepsilon}$, so the closed-form threshold is

$$\varepsilon(\delta) = \delta^2 / C_{\text{tail}}^2$$
.

Proof. By (A1), locking token position i at t^* permanently stops the i's step-wise compute, so for all $t \ge t^*$ we have $\hat{p}_t^{(i)} = p_{t^*}^{(i)}$ and $\log \hat{p}_T^{(i)} = \log p_{t^*}^{(i)}$. Therefore

$$\left\|\log p_T^{(i)} - \log \hat{p}_T^{(i)}\right\|_{\infty} = \left\|\log p_T^{(i)} - \log p_{t^{\star}}^{(i)}\right\|_{\infty} = \left\|f(z_T^{(i)}) - f(z_{t^{\star}}^{(i)})\right\|_{\infty},$$

By the triangle inequality and telescoping over $s = t^* + 1, \dots, T$,

$$\|z_T^{(i)} - z_{t^*}^{(i)}\|_2 \le \sum_{s>t^*}^T \|z_s^{(i)} - z_{s-1}^{(i)}\|_2.$$

Applying (A3) term-wise yields

$$||z_T^{(i)} - z_{t^*}^{(i)}||_2 \le L \sum_{s=t^*+1}^T \sqrt{D_{s-1}^{(i)}}.$$

 $\begin{aligned} \text{Under (A2), } D_{s-1}^{(i)} & \leq \rho^{s-1-t^{\star}} D_{t^{\star}}^{(i)} \text{, therefore } \sqrt{D_{s-1}^{(i)}} \leq \rho^{\frac{s-1-t^{\star}}{2}} \sqrt{D_{t^{\star}}^{(i)}} \text{ and } \\ & \sum_{s > t^{\star}} \sqrt{D_{s-1}^{(i)}} \leq \frac{1}{1-\sqrt{\rho}} \sqrt{D_{t^{\star}}^{(i)}}. \end{aligned}$

Finally, by (A4):

$$\left\|\log p_T^{(i)} - \log p_{t^*}^{(i)}\right\|_{\infty} \le L_{\text{sm}} \|z_T^{(i)} - z_{t^*}^{(i)}\|_2 \le C_{\text{tail}} \sqrt{D_{t^*}^{(i)}}$$

which proves the claim and the stated $\varepsilon(\delta)$.

2.4 COMPUTATIONAL COMPLEXITY: ALGORITHMIC FLOPS

This paper reports algorithmic FLOPs, counting only GEMMs (Q/K/V/Out projections, QK^{\top} and AV, and FFN sublayers). Let $N_{\rm gen}$ the number of positions reserved for generation at initial step, $N_{\rm prompt}$ the number of tokens for the prompt, $N_b \coloneqq \max_{b \in B}(N_{\rm prompt}) + N_{\rm gen}$ be the sequence length of batch b; S the number of sampling steps; S the batch size; S the model dimension; S the number of layers; S the number of attention heads; S the number of token positions; S the active index set at S; and S the index of the properties of the sequence length of S the number of attention heads; S the number of token positions; S the active index set at S; and S the properties of the number of token positions; S the active index set at S; and S the properties of the number of token positions; S the active index set at S; and S the properties of the number of token positions; S the active index set at S; and S the properties of the number of token positions; S the active index set at S the properties of the number of token positions; S the active index set at S the properties of the number of token positions; S the active index set at S the properties of the number of token positions; S the active index set at S the number of token positions is S.

Baseline. For a step t, algorithmic FLOPs per batch are constant across t:

$$\mathcal{F}_{\text{base,b}}^t = L \; \underbrace{(4BHN_b{}^2d_h}_{\text{OK}{}^\top + AV} + \underbrace{2BN_bd^2}_{Q} + \underbrace{2BN_bd^2}_{\text{Out}} + \underbrace{4BN_bd \; H_{kv}d_h}_{K,V} + \underbrace{6BN_bd \; d_{\text{ff}}}_{\text{FFN}} \;).$$

SureLock. Algorithmic FLOPs exhibit temporal dynamics since step-wise computation is only for active positions $\in A_{t,b}$, which changes across t:

$$\mathcal{F}_{\text{prop,b}}^{t} = \frac{|\mathcal{A}_{t,b}|}{BN_{b}} \mathcal{F}_{\text{base,b}}^{t} = \frac{M_{t,b}}{T_{b}} \mathcal{F}_{\text{base,b}}^{t}$$

3 EXPERIMENTS

This section reports our empirical evaluation. We consider two kinds of basic tasks—language modeling and instruction following. In each run with different settings, profiling the per-batch sequence length N_b and the number of active positions $M_{t,b} = |\mathcal{A}_{t,b}|$ at each step, we report computational complexity (algorithmic FLOPs). We also assess any quality degradation in these tasks induced by locking behavior of SURELOCK. Unless otherwise noted, we set the number of sampling steps to $S = N_{\rm gen}$, and fix the percentile for the optional confidence gate (Sec. 2.2) to m=20%. The block length for the semi-aggressive generation option is set to $N_g = N_{\rm gen}$, yielding a fully parallel configuration. Temperature and classifier-free guidance scale are set to 0 by default.

3.1 MASKED DIFFUSION LANGUAGE MODELS AND DATASETS.

We evaluate leading open-weight and representative MDLMs with 8 billion parameters: LLaDA-8B-Base and LLaDA-8B-Instruct (Nie et al., 2025).

We use the following datasets: WikiText-103 (Merity et al., 2016) for language modeling (continuation generation) and MT-Bench (Zheng et al., 2023) for instruction following benchmarking. These benchmarks evaluate complementary aspects: WIKITEXT-103 stresses core token-level modeling on broad, category-agnostic text, whereas MT-BENCH evaluates instruction-following ability and response helpfulness across eight categories (e.g., "writing", "coding"); together they cover both core modeling capacity and practical instruction-following behavior. Data usage settings vary by experimental set, so we will explain them in the following as needed. More detailed settings are provided in Appendices D– F.

³Element-wise operations (activations, gating, LN, softmax, RoPE) and cache I/O/packing are lower-order relative to the dominant $O(N_b^2d)$ attention and $O(N_bd^2)$ FFN matmuls and occur similarly across methods, so omitting them does not affect relative comparisons.

| $N_{ m gen}$ | ε | $\downarrow \bar{r}$ | $\downarrow \mathcal{F}_{base}$ | $\downarrow \mathcal{F}_{prop}$ | $\downarrow \mathcal{F}_{prop}/\mathcal{F}_{base}$ |
|--------------|------|----------------------|---------------------------------|---------------------------------|--|
| 64 | 5e-4 | .548 | 9.017e+11 | 4.936e+11 | $0.547 \times$ |
| 64 | 5e-3 | .506 | 9.017e+11 | 4.565e+11 | $0.506 \times$ |
| 64 | 5e-2 | .482 | 9.017e+11 | 4.342e+11 | $0.482 \times$ |
| 256 | 5e-4 | .496 | 3.629e+12 | 1.801e+12 | 0.496× |
| 256 | 5e-3 | .482 | 3.629e+12 | 1.750e+12 | $0.482 \times$ |
| 256 | 5e-2 | .475 | 3.629e+12 | 1.725e+12 | $0.475 \times$ |
| 1024 | 5e-4 | .494 | 1.492e+13 | 7.366e+12 | 0.494× |
| 1024 | 5e-3 | .492 | 1.492e+13 | 7.333e+12 | $0.492 \times$ |
| 1024 | 5e-2 | .490 | 1.492e+13 | 7.315e+12 | $0.490 \times$ |

Table 1: Algorithmic FLOPs per token by setting— \mathcal{F}_{base} and \mathcal{F}_{prop} . \bar{r} is the micro-averaged ratio of the number of active token positions.

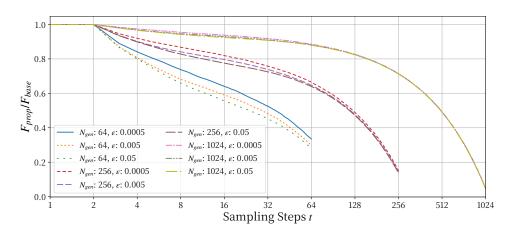


Figure 2: **Step-wise FLOPs ratio.** Ratio of step-wise algorithmic FLOPs, i.e., $\mathcal{F}_{prop}^t/\mathcal{F}_{base}^t$ (and the micro-averaged number of per-step active positions $\bar{r}_t = \sum_b M_{t,b}/\sum_b B N_b$ consistently decreases as steps proceed, explaining later-step savings of computational cost.

3.2 RESULTS

Experiment-1: Computational Complexity. We first evaluate the impact of Surelock on the reduction of algorithmic FLOPs, using LLaDA-8B-Instruct on a fixed set of 32 single-turn prompts from MT-Bench, which we sampled uniformly from each category (Appendix. D). We use a subset of MT-Bench to efficiently sweep diverse inference settings. We set batch size B=4, $N_{\rm gen}=\{64,256,1024\}$, $\varepsilon=\{5e-4,5e-3,5e-2\}$, and run to record the actual number of active token positions $M_{t,b}=|\mathcal{A}_{t,b}|$ sequence length N_b per-batch per-step. Table 1 shows that the computational complexity steadily decreased compared to the same sampler without applying Surelock. The smaller the locking criterion threshold ε , the lower the reduction rate; it aligns with our design intent where ε is set for tightening the locking test (Sec. 2.2). Moreover, on the scale from 5e-2 to 5e-4, ε is less influential to the reduction of computational complexity. We also report the microaveraged ratio of the number of active positions; $\bar{r}=\sum_b\sum_t M_{t,b}/\sum_b\sum_t BN_b$ just as a reference to intuitively grasp how many active rows have been reduced. We can see that \bar{r} is roughly aligned with the reduction ratio of algorithm FLOPs, suggesting that reducing the number of computable token positions leads to reducing computational load.

Experiment-2: Temporal Dynamics of Computational Complexity. Figure 2 shows the dynamics in computational complexity associated with the progression of sampling steps; logging data points was obtained from the same runs as the experiment-1. We can see that, as sampling progresses, the ratio of algorithmic FLOPs decreases at an accelerating rate. This is probably because surrounding most tokens are unmasked in the later steps $(M_t << N)$ leading to the stability of local-KL. We put the results on the runs with different settings for $m = \{10, 40\}\%$ in the Appendix E.

| $N_{ m gen}$ steps $arepsilon$ | $\downarrow GenPPL_{base}$ | $\downarrow \mathcal{F}_{base}$ | \downarrow GenPPL _{prop} | $\downarrow \mathcal{F}_{prop}$ |
|--------------------------------|----------------------------|---------------------------------|-------------------------------------|---------------------------------|
| 64 32 5e-4 | 5.7537 | 4.488e+11 | 6.0782 (1.06×) | 3.138e+11 (0.70×) |
| 64 32 5e-3 | 5.7537 | 4.488e+11 | 6.2177 (1.08×) | $2.788e+11 (0.62\times)$ |
| 64 64 5e-4 | 3.4813 | 8.976e+11 | 4.5722 (1.31×) | $5.203e+11 (0.58\times)$ |
| 64 64 5e-3 | 3.4813 | 8.976e+11 | 4.7596 (1.37×) | 4.664e+11 (0.52×) |
| 128 128 5e-4 | 2.6006 | 1.799e+12 | 2.9202 (1.12×) | 9.585e+11 (0.53×) |
| 128 128 5e-3 | 2.6006 | 1.799e+12 | 3.1042 (1.19×) | $8.964e+11 (0.50\times)$ |
| 128 64 5e-4 | 3.1604 | 8.997e+11 | 3.8371 (1.21×) | $5.620e+11 (0.63\times)$ |
| 128 64 5e-3 | 3.1604 | 8.997e+11 | 3.7353 (1.18×) | $5.083e+11 (0.57\times)$ |
| 256 128 5e-4 | 2.3428 | 1.808e+12 | 2.6092 (1.11×) | $1.012e+12 (0.56\times)$ |
| 256 128 5e-3 | 2.3428 | 1.808e+12 | 2.6687 (1.14×) | $9.497e+11(0.53\times)$ |
| 256 256 5e-4 | 2.0127 | 3.616e+12 | 2.0486 (1.02×) | $1.845e+12 (0.51\times)$ |
| 256 256 5e-3 | 2.0127 | 3.616e+12 | 2.2103 (1.10×) | $1.779e+12 (0.49\times)$ |
| 512 256 5e-4 | 1.6293 | 3.650e+12 | 1.7134 (1.05×) | 1.919e+12 (0.53×) |
| 512 256 5e-3 | 1.6293 | 3.650e+12 | 1.7526 (1.08×) | $1.851e+12 (0.51\times)$ |
| 512 512 5e-4 | 1.4658 | 7.301e+12 | 1.5248 (1.04×) | $3.643e+12 (0.50\times)$ |
| 512 512 5e-3 | 1.4658 | 7.301e+12 | 1.5444 (1.05×) | $3.588e+12 (0.49\times)$ |

Table 2: Generation quality of continuation sequences with and without SURELOCK on LLaDA-8B-Base using WikiText-103. Gen.-PPL is the micro averaged PPL evaluated with LLaMA-3-8B.

| $N_{ m gen}$ | steps | ε | \uparrow Score _{base} | $\downarrow \mathcal{F}_{base}$ | \uparrow Score _{prop} | $\downarrow \mathcal{F}_{prop}$ |
|--------------|-------|------|----------------------------------|---------------------------------|----------------------------------|---------------------------------|
| 64 | 32 | 5e-4 | 3.9 | 4.515e+11 | 3.8 (0.97×) | 3.153e+11 (0.698×) |
| 64 | 32 | 5e-3 | 3.8 | 4.515e+11 | $3.7 (0.97 \times)$ | $2.947e+11 (0.653\times)$ |
| 64 | 64 | 5e-4 | 4.2 | 9.030e+11 | 4.2 (1.00×) | $5.658e+11 (0.627 \times)$ |
| 64 | 64 | 5e-3 | 4.2 | 9.030e+11 | 4.3 (1.02×) | $5.304e+11 (0.587 \times)$ |
| 128 | 64 | 5e-4 | 4.4 | 9.047e+11 | 4.2 (0.96×) | 5.700e+11 (0.630×) |
| 128 | 64 | 5e-3 | 4.3 | 9.047e+11 | $4.2(0.98\times)$ | $5.406e+11 (0.598 \times)$ |
| 128 | 128 | 5e-4 | 4.8 | 1.809e+12 | $4.9(1.02\times)$ | $1.044e+12 (0.577 \times)$ |
| 128 | 128 | 5e-3 | 4.8 | 1.809e+12 | 5.1 (1.06×) | $1.000e+12 (0.553 \times)$ |
| 256 | 128 | 5e-4 | 4.2 | 1.817e+12 | 4.4 (1.05×) | 1.042e+12 (0.573×) |
| 256 | 128 | 5e-3 | 4.3 | 1.817e+12 | $4.2(0.98\times)$ | $1.007e+12 (0.554\times)$ |
| 256 | 256 | 5e-4 | 4.4 | 3.634e+12 | 4.5 (1.02×) | $1.969e+12 (0.542\times)$ |
| 256 | 256 | 5e-3 | 4.4 | 3.634e+12 | 4.7 (1.07×) | $1.920e+12 (0.528 \times)$ |
| 512 | 256 | 5e-4 | 4.0 | 3.667e+12 | 4.3 (1.08×) | 1.971e+12 (0.537×) |
| 512 | 256 | 5e-3 | 4.0 | 3.667e+12 | $4.2(1.05\times)$ | $1.934e+12 (0.528 \times)$ |
| 512 | 512 | 5e-4 | 4.5 | 7.334e+12 | $4.7 (1.04 \times)$ | $3.819e+12 (0.521\times)$ |
| 512 | 512 | 5e-3 | 4.4 | 7.334e+12 | 4.8 (1.09×) | 3.777e+12 (0.515×) |

Table 3: Quality changes in generated responses by LLaDA-Instruct. Score base and Score prop indicate the overall score for the MT-Bench with single-turn settings evaluated with gpt-40.

Experiment-3: Trade-off between Efficiency and Generation Quality. We evaluate LLaDA-8B-Base on Wikitext-103, where we prompt a fixed set of text fragments, requiring the model to generate the continuation, and report Gen.-PPL measured by an external AR model—LLaMA-3-8B (Grattafiori et al., 2024). More details are in Appendix F. We also evaluate LLaDA-8B-Instruct on MT-Bench with single-turn settings; we report the averaged LLM-as-a-judge score from gpt-40 for the generated responses. Other details are described in Appendix F. We set batch size B=4, $N_{\rm gen}=\{64,128,256,512\}$, $\varepsilon=\{5e-4,5e-3\}$ thoughout datasets, and $N_{\rm p}=64$ in WikiText-103. Table 2 and Table 3 show that while computational complexity is indeed decreasing, performance is being maintained at competitive levels. It demonstrates that while computational load is indeed decreasing, performance is maintained at an equivalent level. Upon deeper analysis, the instruct model consistently maintained performance across all settings i.e., -0.1pt at maximum, while reducing computational load about $0.70\times-0.50\times$. The base model also showed the same trend in most settings. However, base Model sometimes revealed relatively large degradation ($\geq 1.21\times$) for short generation lengths ($N_{\rm gen}=\{64,128\}$). It suggests that, in language modeling where prompts are

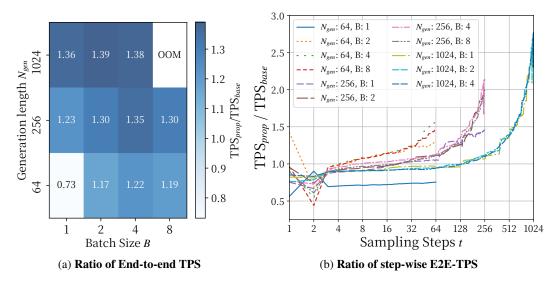


Figure 3: Throughput behavior with SURELOCK: (a) end-to-end TPS ratio across different $N_{\rm gen}$ and batch size B.; (b) per-step TPS ratio increasing as sampling progresses.

Question: Benjamin went to a bookstore and purchased a variety of books. He bought 5 copies of a sci-fi novel, each priced at \$20, 3 copies of a history book priced at \$30 each, and 2 copies of a philosophy book for \$45 each. What was the total cost of his purchases?.

Baseline (score: 10)

Benjamin bought 5 copies of a sci-fi novel, each priced at \$20, so he spent 5 * \$20 = \$100. He bought 3 copies of a history book, each priced at \$30, so he spent 3 * \$30 = \$90. He bought 2 copies of a philosophy book, each priced at \$45, so he spent 2 * \$45 = \$90. The total cost of his purchases is \$100 + \$90 + \$90 = \$280. Summing up,

SureLock (score: 10)

Benjamin bought 5 copies of a sci-fi novel for \$20 each, so the total cost is 5 * \$20 = \$100. He also bought 3 copies of a history book for \$30 each, so the total cost is 3 * \$30 = \$90. Lastly, he bought 2 copies of a philosophy book for \$45 each, so the total cost is 2 * \$45 = \$90. The total cost of his purchases is \$100 + \$90 + \$90 = \$280. Con

Figure 4: Comparison of responses between Baseline vs. SURELOCK on LLaDA-8B-Instrut with $\varepsilon = 5e - 4$, $N_{\rm gen} = 128$, S = 128. The question is sampled from MT-bench with question id= 119.

considered as open questions, setting a short $N_{\rm gen}$ can lead to undesirable final results by SURE-LOCK. While this can affect creative tasks like novel generation, the impact is minimal since such tasks are not typically executed on short sequences e.g., 128 tokens.

Experiment-4: Runtime Performance. We evaluate End-to-end Token/Sec (E2E-TPS), which is the sustained decoding throughput of the model across multiple batches, and temporal dynamics for step-wise TPS. See Appendix G for more details on the metrics. We here follow the settings in Experiment-1 as for the dataset and the model. We set batch size $B = \{1, 2, 4, 8\}$, $N_{\rm gen} = \{64, 256, 1024\}$, $\varepsilon = 5e - 3$. Figure 3a shows that SURELOCK can achieve greater runtime gains in compute-bound areas e.g., $N \ge 256$, B > 1, and no runtime gain was observed under relatively light computational load settings ($N_g = 64$, B = 1); this differs from the trend in computational complexity. SURELOCK reduces the per-step compute with the monotonously decreased number of active positions (M_t) but keeps the costs other than computational complexity. For example, it leaves the KV read intact and accessing weights for specific rows/columns involves discontinuous memory access. Therefore, speedups relatively easily appear in compute-heavy regimes (long $N_{\rm gen}$, larger batch size B) where gains exceeding fixed costs are expected. See Figure 3b for temporal dynamics of step-wise TPS. We can see that the gain in local runtime increases as the step progresses, which is consistent with the computational trends (Figure 2).

Example Analysis. Feeding the same question sampled from MT-Bench to the LLaDA-8B-Instruct, Figure 4 shows the responses from the same sampler with and without SURELOCK. With settings achieving approximately $0.6\times$ of computational complexity on average (See Table 3), we can confirm that not only the quantitative score but also the qualitative quality is comparable to the baseline. While fixing token representations during sampling might seem to introduce disruptions to generated text, we can observe the limited impact. More examples are described in Appendix I.

4 RELATED WORK

 Computational savings for DLMs. Prior efforts to decrease the sampling cost of discrete diffusion LMs (DLMs) largely follow two axes—Temporal and Reuse. Temporal methods shrink the number of sampling steps T via parallel/dilated unmasking and edit-based updates (Ghazvininejad et al., 2019; Gu et al., 2019; Luxembourg et al., 2025; Israel et al., 2025; Huang & Tang, 2025; Wei et al., 2025), or by adopting higher-order / distilled samplers from diffusion literature, e.g., DPM-Solver (Lu et al., 2022). Reuse methods amortize work across steps by reusing or approximating intermediate states—e.g., K/V caching or approximation and selective feature refresh (Ma et al., 2025; Liu et al., 2025b; Wu et al., 2025a), Selection-based partial compute reduces per-step complexity by updating only a fixed number of subset positions, sometimes with periodic full recomputation for stability (Liu et al., 2025b; Wu et al., 2025a). These two axes primarily reduce how many steps we take or how much previous computation we can reuse between steps; they typically leave the within-step active-query count constant up to the end of the sampling step, so per-step compute remains bounded by N even late in sampling. Crucially, these are complementary to Surelock; by focusing only on the gradually contracting active position sets, they could yield even greater savings.

DLMs on longer sequences. Scaling DLMs to longer contexts has been actively explored such as span/block-wise masking (Arriola et al., 2025), adaptive masked token insertion (Arriola et al., 2025), DLMs' suite Rotary Position Embedding (RoPE) Liu et al. (2025a). This line of research directions yields diverse applications such as reasoning on massive code bases. By contrast, decoding longer sequences by DLMs generally requires more sampling steps, making the stationary per-step computational cost more severe. Advancing techniques such as SURELOCK, which monotonically reduces this stationary compute as sampling proceeds, could facilitate DLMs' research on longer contexts, which lags behind compared to AR models (Wu et al., 2025b), i.e., thousands vs. millions.

5 CONCLUSION

We introduced Surelock, a method for iterative decoding in masked-diffusion LMs (MDLMs) that locks converged token positions and thereafter bypasses their query projection and per-token FFN, yielding a monotonic reduction in per-step compute. Locked positions remain fully attendable via cached K/V vectors. The core design is a $local\ KL$ locking test; we justified this approach by deriving a closed-form link between the lock-time local KL and an error bound on the terminal log-probability. On language modeling and instruction-following tasks, Surelock achieves large reductions in algorithmic FLOPs while maintaining task quality. The approach is complementary to step-count reduction and inter-step reuse techniques, as well as position selection methods (Liu et al., 2025b); such methods can operate on the progressively shrinking active token positions induced by Surelock, leaving their integration to future work.

Limitations. We report computational efficiency in terms of algorithmic FLOPs, a kernel-agnostic measure that reflects theoretical savings rather than wall-clock speedups, which isolates the effect of sampling algorithm from hardware-specific factors. Moreover, we use the theorem solely to justify the *algorithm design*, especially the locking criterion based on local KL. While one could estimate L and ρ online to obtain numerical thresholds given a specific model and dataset, we leave it to future work. In our experiment, performance degradation was not particularly noticeable; however, methods for addressing serious errors, e.g., periodic unlocking (Appendix H), are also worth discussing.

ETHICS STATEMENT

This paper aims to reduce computational complexity in inference time for discrete diffusion models, which may have the potential to significantly impact real-world applications. Furthermore, this research does not require diffusion models to generate harmful or sexual contents; we believe our work does not directly contribute to malicious use for such purposes.

7 REPRODUCIBILITY STATEMENT

We will release the code used in our experiments to facilitate reproducibility. Experimental settings are included in Sec. 3. More detailed descriptions are provided in the Appendices.

REFERENCES

- Marianne Arriola, Aaron Gokaslan, Justin T Chiu, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Subham Sekhar Sahoo, and Volodymyr Kuleshov. Block diffusion: Interpolating between autoregressive and diffusion language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. Mask-predict: Parallel decoding of conditional masked language models. *arXiv preprint arXiv:1904.09324*, 2019.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Jiatao Gu, Changhan Wang, and Junbo Zhao. Levenshtein transformer. Advances in neural information processing systems, 32, 2019.
- Chihan Huang and Hao Tang. Ctrldiff: Boosting large diffusion language models with dynamic block prediction and controllable generation. *arXiv* preprint arXiv:2505.14455, 2025.
- Daniel Israel, Guy Van den Broeck, and Aditya Grover. Accelerating diffusion llms via adaptive parallel decoding. *arXiv preprint arXiv:2506.00413*, 2025.
- Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *Advances in neural information processing systems*, 35: 4328–4343, 2022.
- Xiaoran Liu, Zhigeng Liu, Zengfeng Huang, Qipeng Guo, Ziwei He, and Xipeng Qiu. Longllada: Unlocking long context capabilities in diffusion llms. *arXiv preprint arXiv:2506.14429*, 2025a.
- Zhiyuan Liu, Yicun Yang, Yaojie Zhang, Junjie Chen, Chang Zou, Qingyuan Wei, Shaobo Wang, and Linfeng Zhang. dllm-cache: Accelerating diffusion large language models with adaptive caching. arXiv preprint arXiv:2506.06295, 2025b.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in neural information processing systems*, 35:5775–5787, 2022.
- Omer Luxembourg, Haim Permuter, and Eliya Nachmani. Plan for speed–dilated scheduling for masked diffusion language models. *arXiv preprint arXiv:2506.19037*, 2025.
- Xinyin Ma, Runpeng Yu, Gongfan Fang, and Xinchao Wang. dkv-cache: The cache for diffusion language models. *arXiv preprint arXiv:2505.15781*, 2025.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016.
 - 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.

- Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu, Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language models. *Advances in Neural Information Processing Systems*, 37:130136–130184, 2024.
 - Yair Schiff, Subham Sekhar Sahoo, Hao Phung, Guanghan Wang, Sam Boshar, Hugo Dalla-torre, Bernardo P de Almeida, Alexander Rush, Thomas Pierrot, and Volodymyr Kuleshov. Simple guidance mechanisms for discrete diffusion models. *arXiv preprint arXiv:2412.10193*, 2024.
 - Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis Titsias. Simplified and generalized masked diffusion for discrete data. Advances in neural information processing systems, 37: 103131–103167, 2024.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
 - Qingyan Wei, Yaojie Zhang, Zhiyuan Liu, Dongrui Liu, and Linfeng Zhang. Accelerating diffusion large language models with slowfast: The three golden principles. *arXiv preprint arXiv:2506.10848*, 2025.
 - Chengyue Wu, Hao Zhang, Shuchen Xue, Zhijian Liu, Shizhe Diao, Ligeng Zhu, Ping Luo, Song Han, and Enze Xie. Fast-dllm: Training-free acceleration of diffusion llm by enabling kv cache and parallel decoding. *arXiv preprint arXiv:2505.22618*, 2025a.
 - Yuhao Wu, Yushi Bai, Zhiqing Hu, Shangqing Tu, Ming Shan Hee, Juanzi Li, and Roy Ka-Wei Lee. Shifting long-context llms research from input to output. *arXiv preprint arXiv:2503.04723*, 2025b.
 - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023.

A APPENDIX

B USAGE OF LARGE LANGUAGE MODELS

We used LLMs to assist with writing paragraphs and searching literature. However, we retain full responsibility for the content in this paper.

C DETAILED EXPLANATION OF ASSUMPTION A3.

We show that the one-step deviation of the logit vector $z_s^{(i)}$ at position i is controlled by the *previous* step's drift on local posterior measured by KL divergence. Throughout, let $p_t^{(j)} \in \Delta^{V-1}$ be the token posterior at position i and step t, $D_t^{(j)} \coloneqq \mathrm{KL}(p_t^{(j)} \parallel p_{t-1}^{(j)})$, and let $E \in \mathbb{R}^{V \times d_{\mathrm{model}}}$ denote the embedding matrix (rows are token embeddings). We use the block norm $\|\cdot\|_{2,1}$ defined by $\|X\|_{2,1} = \sum_j \|x_j\|_2$ for $X = [x_1, \dots, x_n]^\top$.

Reclaiming Assumption A3. There exists a constant L > 0, depending only on operator norms of the network weights and softmax/LN/activation Lipschitz constants, such that for all steps s and positions i,

$$||z_s^{(i)} - z_{s-1}^{(i)}||_2 \le L\sqrt{D_{s-1}^{(i)}}.$$
 (1)

Below we prove a slightly stronger *global* bound depending on $\sum_j \sqrt{D_{s-1}^{(j)}}$ and then specialize it to the local form equation 1 using a locality factor.

C.1 From posterior drift to expected-embedding drift

For any $q, r \in \Delta^{V-1}$ define the expected input (embedding) $x(q) := E^{\top} q \in \mathbb{R}^{d_{\text{model}}}$. Then

$$||x(q) - x(r)||_2 = ||E^{\top}(q - r)||_2 \le ||E||_2 ||q - r||_2 \le ||E||_2 ||q - r||_1.$$
 (A.3.1)

By Pinsker, $||q-r||_1 \leq \sqrt{2 \operatorname{KL}(q||r)}$, so with $L_{\text{emb}} := ||E||_2 \sqrt{2}$ we have

$$||x(q) - x(r)||_2 \le L_{\text{emb}} \sqrt{\text{KL}(q||r)}.$$
 (A.3.2)

Applying equation A.3.2 to successive steps at position i,

$$\Delta x_s^{(j)} := x(p_{s-1}^{(j)}) - x(p_{s-2}^{(j)}), \qquad \|\Delta x_s^{(j)}\|_2 \le L_{\text{emb}} \sqrt{D_{s-1}^{(j)}}.$$
 (A.3.3)

C.2 SINGLE-HEAD ATTENTION: LIPSCHITZ BOUND

Consider a single-head attention with parameters (W_Q, W_K, W_V) and key/query dimension d_k . For input $X = [x_1, \dots, x_n]^\top$, the head output at position i is

$$o_i(X) = \sum_{j=1}^n \alpha_{ij}(X) W_V x_j, \qquad \alpha_{i\bullet}(X) = \operatorname{softmax}\left(\frac{1}{\sqrt{d_k}}(W_Q x_i)(W_K X)^\top\right).$$
 (A.3.4)

For two inputs X and $X' = X - \Delta X$, a standard three-term decomposition gives

$$\Delta o_i := o_i(X) - o_i(X') = \underbrace{\sum_j \alpha'_{ij} W_V \Delta x_j}_{(I)} + \underbrace{\sum_j (\alpha_{ij} - \alpha'_{ij}) W_V x'_j}_{(II)} + \underbrace{\sum_j (\alpha_{ij} - \alpha'_{ij}) W_V \Delta x_j}_{(III)}.$$

We upper bound (I) and (II) and subsume (III) into them (yielding a conservative bound).

Term (I) By the triangle inequality,

$$||(I)||_2 \le ||W_V||_2 \sum_j \alpha'_{ij} ||\Delta x_j||_2 \le ||W_V||_2 \sum_j ||\Delta x_j||_2 = ||W_V||_2 ||\Delta X||_{2,1}.$$
(A.3.6)

Softmax Jacobian bound. Let $\alpha = \operatorname{softmax}(s)$ and $J_{\operatorname{sm}}(s) = \operatorname{diag}(\alpha) - \alpha \alpha^{\top}$. It is classical that $\sup \|J_{\operatorname{sm}}(s)\|_2 = \frac{1}{2}, \tag{A.3.7}$

achieved at two-point distributions. Hence, with mean value theorem, for score vectors s_i and s_i' ,

$$\|\alpha_{i\bullet} - \alpha'_{i\bullet}\|_2 \le \frac{1}{2} \|s_i - s'_i\|_2, \qquad \|\alpha_{i\bullet} - \alpha'_{i\bullet}\|_1 \le \sqrt{n} \|\alpha_{i\bullet} - \alpha'_{i\bullet}\|_2.$$
 (A.3.8)

Score difference. Write the (i, j) score as

$$s_{ij} = \frac{1}{\sqrt{d_k}} \langle W_Q x_i, W_K x_j \rangle. \tag{A.3.9}$$

Then

$$|\Delta s_{ij}| = \frac{1}{\sqrt{d_k}} |\langle W_Q \Delta x_i, W_K x_j' \rangle + \langle W_Q x_i, W_K \Delta x_j \rangle + \langle W_Q \Delta x_i, W_K \Delta x_j \rangle|$$

$$\leq \frac{\|W_Q\|_2 \|W_K\|_2}{\sqrt{d_k}} \left(\|\Delta x_i\|_2 \|x_j'\|_2 + \|x_i\|_2 \|\Delta x_j\|_2 + \|\Delta x_i\|_2 \|\Delta x_j\|_2 \right). \tag{A.3.10}$$

Assume a uniform radius bound $\|x_j\|_2$, $\|x_j'\|_2 \le R_x$ for all j (true if inputs are bounded and LN is used). Summing equation A.3.10 over j and using Cauchy–Schwarz and $\|\Delta X\|_{2,1} = \sum_j \|\Delta x_j\|_2$ yields

$$\|\Delta s_{ij}\|_{2} = \sqrt{\sum_{j=1}^{n} \Delta(s_{i,j})^{2}} \leq \frac{\|W_{Q}\|_{2} \|W_{K}\|_{2}}{\sqrt{d_{k}}} \left(\sqrt{n} R_{x} \|\Delta x_{i}\|_{2} + R_{x} \|\Delta X\|_{2,1} + \|\Delta x_{i}\|_{2} \|\Delta X\|_{2,1}\right). \tag{A.3.11}$$

Term (II). Using the triangle inequality and operator norm of W_V , we have

$$||(II)||_2 = \left\| \sum_j (\alpha_{ij} - \alpha'_{ij}) W_V x'_j \right\|_2 \le ||W_V||_2 \sum_j |\alpha_{ij} - \alpha'_{ij}| \, ||x'_j||_2. \tag{A.3.12a}$$

Assuming $||x_j'||_2 \le R_x$ for all j, we obtain

$$||(II)||_2 \le ||W_V||_2 R_x ||\alpha_{i\bullet} - \alpha'_{i\bullet}||_1.$$
 (A.3.12b)

Using the norm inequality $\|\cdot\|_1 \leq \sqrt{n} \|\cdot\|_2$,

$$||(II)||_2 \le ||W_V||_2 R_x \sqrt{n} ||\alpha_{i\bullet} - \alpha'_{i\bullet}||_2.$$
 (A.3.12c)

Finally, applying the softmax Jacobian bound $\|\alpha_{i\bullet} - \alpha'_{i\bullet}\|_2 \le \frac{1}{2} \|s_i - s'_i\|_2$ (cf. equation A.3.8), we arrive at

$$||(II)||_2 \le \frac{||W_V||_2 \sqrt{n} R_x}{2} ||s_i - s_i'||_2.$$
 (A.3.12)

Substituting equation A.3.11 into equation A.3.12, we obtain

$$||(II)||_{2} \leq \frac{||W_{V}||_{2}\sqrt{n}R_{x}}{2} \cdot \frac{||W_{Q}||_{2}||W_{K}||_{2}}{\sqrt{d_{k}}} \left(\sqrt{n}R_{x} ||\Delta x_{i}||_{2} + R_{x} ||\Delta X||_{2,1} + ||\Delta x_{i}||_{2} ||\Delta X||_{2,1}\right). \tag{A.3.13a}$$

Using $\|\Delta x_i\|_2 \leq \|\Delta X\|_{2,1}$, this simplifies to

$$||(II)||_2 \le ||W_V||_2 \cdot \frac{||W_Q||_2 ||W_K||_2}{2\sqrt{d_k}} \cdot nR_x^2 ||\Delta X||_{2,1}.$$
 (A.3.13b)

Therefore, combining equation A.3.6 and equation A.3.13b, we obtain the single-head bound

$$\|\Delta o_i\|_2 \le A_{\text{att}} \|\Delta X\|_{2,1}, \qquad A_{\text{att}} := \|W_V\|_2 \left[1 + \frac{\|W_Q\|_2 \|W_K\|_2}{2} \frac{n R_x^2}{\sqrt{d_k}}\right].$$
 (A.3.13)

C.3 BLOCK COMPOSITION: MHA + FFN + RESIDUAL + LN

Let one Transformer block be

$$B(X) = X + FFN(MHA(LN(X))). \tag{A.3.14}$$

702 Assume Lipschitz bounds for each component (w.r.t. $\|\cdot\|_{2,1}$)

$$\|\operatorname{LN}(X) - \operatorname{LN}(X')\|_{2,1} \le L_{\ln} \|X - X'\|_{2,1}, \quad \|\operatorname{MHA}(U) - \operatorname{MHA}(U')\|_{2,1} \le A_{\min} \|U - U'\|_{2,1}, \quad (A.3.15)$$

and

$$\|\text{FFN}(V) - \text{FFN}(V')\|_{2,1} \le A_{\text{ff}} \|V - V'\|_{2,1}$$
 (A.3.16)

For multi-head attention with h heads concatenated and an output projection W_O , one may take

$$A_{\text{mha}} \leq \|W_O\|_2 \left(\max_{\text{head}} A_{\text{att}}^{(\text{head})}\right). \tag{A.3.17}$$

By the residual form equation A.3.14,

$$||B(X) - B(X')||_{2,1} \le ||X - X'||_{2,1} + A_{\text{ff}} A_{\text{mha}} L_{\ln} ||X - X'||_{2,1} = L_{\text{blk}} ||X - X'||_{2,1},$$
 (A.3.18)

with

$$L_{\text{blk}} := 1 + A_{\text{ff}} A_{\text{mha}} L_{\text{ln}}. \tag{A.3.19}$$

For a stack of H blocks,

$$||B^{H}(X) - B^{H}(X')||_{2,1} \le L_{\text{blk}}^{H} ||X - X'||_{2,1}.$$
 (A.3.20)

Finally, with readout $z_i = W_o h_i$, we get

$$||z_i(X) - z_i(X')||_2 \le ||W_o||_2 ||B^H(X) - B^H(X')||_{2,1} \le \underbrace{||W_o||_2 L_{\text{blk}}^H}_{I} ||X - X'||_{2,1}.$$
 (A.3.21)

C.4 PUTTING THINGS TOGETHER

Applying equation A.3.21 to successive sampler states (X_s, X_{s-1}) and using equation A.3.3,

$$||z_s^{(i)} - z_{s-1}^{(i)}||_2 \le L_{\text{net}} ||X_s - X_{s-1}||_{2,1} = L_{\text{net}} \sum_j ||\Delta x_s^{(j)}||_2 \le L_{\text{net}} L_{\text{emb}} \sum_j \sqrt{D_{s-1}^{(j)}}.$$
(A.3.22)

Define $L_{\text{all}} := L_{\text{net}} L_{\text{emb}}$. Equation A.3.22 is the *global* for if equation 1.

$$||z_s^{(i)} - z_{s-1}^{(i)}||_2 \le L_{\text{all}} \sum_j \sqrt{D_{s-1}^{(j)}}.$$
 (A.3.23)

Late in iterative sampling the re-masking rate decreases and many positions become stable. Empirically one can upper bound the tail contribution at step s-1 by

$$\sum_{i \neq i} \sqrt{D_{s-1}^{(j)}} \le \kappa_{s-1}^{(i)} \sqrt{D_{s-1}^{(i)}}, \qquad \kappa_{s-1}^{(i)} \in [0, \infty), \tag{A.3.24}$$

where $\kappa_{s-1}^{(i)}$ is an observable ratio. Substituting equation A.3.24 into equation A.3.23,

$$||z_{s}^{(i)} - z_{s-1}^{(i)}||_{2} \le L_{\text{all}} (1 + \kappa_{s-1}^{(i)}) \sqrt{D_{s-1}^{(i)}} = \underbrace{L_{\text{all}} (1 + \kappa_{s-1}^{(i)})}_{=:L} \sqrt{D_{s-1}^{(i)}}.$$
(A.3.25)

This is exactly Assumption A3 in the main text.

Constants at a glance.

$$L \ = \ \underbrace{\|E\|_2\sqrt{2}}_{L_{\mathrm{emb}}} \times \underbrace{\|W_o\|_2 \left(1 + A_{\mathrm{ff}}\,A_{\mathrm{mha}}\,L_{\ln}\right)^H}_{L_{\mathrm{net}}} \times (1 + \kappa_{s-1}^{(i)}), \quad A_{\mathrm{mha}} \leq \|W_O\|_2 \max_{\mathrm{head}} A_{\mathrm{att}}^{\mathrm{(head)}},$$

with $A_{\mathrm{att}}^{(\mathrm{head})}$ given by equation A.3.13 and $A_{\mathrm{ff}} = \|W_2\|_2 L_{\mathrm{act}} \|W_1\|_2$.

D Dataset Settings for Experiment-1

To efficiently sweep a large configuration space, in experiment-1 to profile algorithm FLOPs, we used a focused subset of MT-Bench rather than the full set. For each of the eight categories e.g., "writing" and "coding," we use the first four single-turn prompts—32 prompts in total—so that diverse phenomena may appear. Prompts are used unedited and are preprocessed with the LLADA-8B-Instruct tokenizer using its default chat template; no system message is injected.

E Additional Results for Experiment-2

We show additional results for different settings of m-percentile used for the optional confidence gate 2.2; we set $m = \{10\%, 40\%\}$, which differ from the one in the main text i.e., m = 20%. The trend in temporal dynamics of FLOPs ratio does not change significantly with m, suggesting m is not a particularly sensitive parameter.

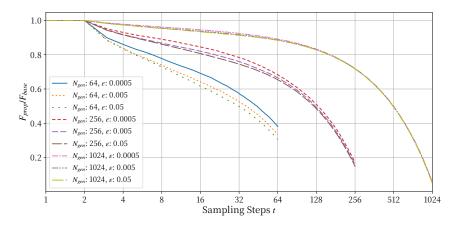


Figure 5: **Step-wise FLOPs ratio.** Algorithmic FLOPs ratio at step t i.e., $\mathcal{F}_{prop}^t/\mathcal{F}_{base}^t$ (and the averaged per-step active row ratio (micro) $\bar{r}_t = \sum_b M_{t,b}/\sum_b B\,N_b$ consistently decreases as steps proceed, explaining later-step savings of computational cost. m-percentile is 10%.

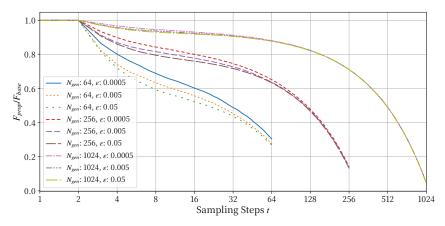


Figure 6: **Step-wise FLOPs ratio.** Algorithmic FLOPs ratio at step t i.e., $\mathcal{F}_{prop}^t/\mathcal{F}_{base}^t$ (and the averaged per-step active row ratio (micro) $\bar{r}_t = \sum_b M_{t,b}/\sum_b B N_b$ consistently decreases as steps proceed, explaining later-step savings of computational cost. m-percentile is 40%.

F DETAILED SETTINGS FOR EXPERIMENT-3

F.1 ON WIKITEXT-103

We extracted a small subset from Wikitext-103 (Merity et al., 2016), which contains diverse text, for the continuation generation task. First, we removed blank lines from the dataset loaded from https://huggingface.co/datasets/Salesforce/wikitext, then extracted the first 120 records at odd indices. Next, we extracted the first 64 tokens from each record to form the prompt x for the continuation. We here used LLaDA-8B-Base tokenizer (Nie et al., 2025).

For each prompt x, we generate a sequence y of fixed length N_g , then compute Gen.-PPL using an external autoregressive language model θ . This work uses Llama-3-8B (Grattafiori et al., 2024), a representative autoregressive language model. We define Gen.-PPL as perplexity microaveraged over the dataset; we first compute the micro-averaged negative log-likelihood $\overline{\text{NLL}}_{\text{micro}} = \sum_i \text{NLL}_i / \sum_i |y_i|$, then $\text{PPL}_{\text{micro}} = \exp(\overline{\text{NLL}}_{\text{micro}})$. In the report, we used natural logarithms.

F.2 ON MT-BENCH

We evaluated LLaDA-8B-Instruct (Nie et al., 2025) using all 80 records in MT-bench Zheng et al. (2023) with a single-turn configuration—we used the first prompt for each record. Essentially, using the LLaDA-8B-Instruct tokenizer, we applied its default chat template to each prompt to form input. We then fed these prompts to the model to generate a responses.

Response evaluation was conducted using the LLM-as-a-judge format. Specifically, we used the official MT-Bench evaluation repository https://github.com/lm-sys/FastChat/tree/main/fastchat/llm_judge as-is. As the judge model, we employed OpenAI's gpt-4o.

G RUNTIME METRICS

E2E-TPS measures the sustained decoding throughput of the model across multiple batches, including all inter-batch gaps and host/device launch overheads incurred during decoding. Concretely, after preparing batches we synchronize the device and start a single wall-clock timer immediately before the first decoding call, and stop it after the last decoding finishes. The numerator is the total number of unmasked tokens produced by the iterative sampling process.

$$E2E\text{-}TPS = \frac{\#\{total\ unmasked\ tokens\}}{wall\text{-}clock\ time}$$

Step TPS at sampling step t is computed as the number of newly unmasked tokens at that step divided by the wall-clock time, aggregating across batches:

$$StepTPS(t) = \frac{\#\{newly unmasked tokens at step t\}}{\Delta t}$$

where Δt is the wall-clock time for that step t. In both baseline and SURELOCK, timing uses CUDA events with synchronization.

H OPTIONAL UNLOCKING

In SURELOCK, once locked, a token position is *excluded from re-masking by construction* (Sec. 2). However, rare context shifts can render a cached posterior stale. Therefore, allowing lightweight *unlocking* to detect context drift within a small finite budget and returning tokens to the active set could also be considered as a useful option. In following, we discuss an implementation for that option.

After every P steps we evaluate a budgeted proxy on locked rows only (e.g., variable-length attention with a reduced subgraph). This yields proxy uncertainty $\tilde{u}_t^{(i)}$ and distributional drift from the locked time $\widetilde{D}_t^{(i)} = \mathrm{KL}(\hat{p}_t^{(i)} \| p_{t_t^{(i)}}^{(i)})$. Let $\theta_t = q_m(u_t)$ be the top-m% uncertainty among *active* tokens at step t. We *unlock* a locked token i if any holds:

$$\tilde{u}_{t}^{(i)} > \theta_{t} \quad \wedge \quad \tilde{D}_{t}^{(i)} > \varepsilon_{\text{unlock}} \quad \wedge \quad t - t_{i}^{\star} > D_{\text{interval}}.$$
 (2)

 D_{interval} is set for suppressing failure to converge due to excessive oscillation in the locking and unlocking behavior. Furthermore, after an unlock, position i cannot be re-locked for D_{cool} steps and must satisfy a stricter KL threshold $\rho\varepsilon$ ($\rho \leq 1$) when re-locking.

On unlock at step t, we invalidate cached K/V for i and reintroduce its Query projection and FFN sublayers from step t+1; the current step is unchanged. Amortized probe overhead on attention-dominant FLOPs is a trade-off with performance. However, the subgraph used for probing can be reduced to a size negligible compared to the dominant cost, depending on its size.

Let $\Delta_t^{(i)} = \|\log p_t^{(i)} - \log p_t^{(i),\mathrm{full}}\|_2$ be deviation from a no-lock decode. Unlocking cannot increase $\Delta_{t'}^{(i)}$ for $t' \geq t$ because the row is computed by the same subgraph as baseline thereafter. Therefore, the error bound (Theorem 1 in Sec. 2.3) holds up to the unlock time; after that, deviation is non-increasing.

I DECODED EXAMPLES

We provide additional response examples for MT-Bench. We randomly selected one question per category.

Question: Compose an engaging travel blog post about a recent trip to Hawaii, highlighting cultural experiences and must-see attractions.

Baseline (score: 1)

Title: "Unveiling the Paradise: A Cultural Journey Through Hawaii"

Introduction:

SURELOCK (score: 3)

Title: "A Journey Through Paradise: Unveiling

the Hidden Gems of Hawaii"

Introduction: As the sun set over the lush land-

scapes of Hawaii, I felt

Figure 7: Quality comparison of responses between Baseline vs. SureLock response. The example question uses the record from MT-bench with question id= 81. Responses are from LLaDA-8B-Instrut with KL threshold $\varepsilon=5e-4$, $N_{\rm gen}=128$, S=128.

Question: Embrace the role of Sheldon from "The Big Bang Theory" as we delve into our conversation. Don't start with phrases like "As Sheldon". Let's kick things off with the following question: "What is your opinion on hand dryers?"

Baseline (score: 2) Hand dryers are useless. **SURELOCK** (score: 1) Hand dryers.

Figure 8: Quality comparison of responses between Baseline vs. SURELOCK response. The example question uses the record from MT-bench with question id= 92. Responses are from LLaDA-8B-Instrut with KL threshold $\varepsilon = 5e - 4$, $N_{\rm gen} = 128$, S = 128.

Question: Thomas is very healthy, but he has to go to the hospital every day. What could be the reasons?

Baseline (score: 3)

here could be several reasons why Thomas has to go to the hospital every day, despite being very healthy. Some possible reasons include:

- 1. Chronic illness: Thomas may have a chronic illness that requires daily medical attention, such as diabetes, cancer, or a heart condition.
- 2. Injury recovery: Thomas may be recovering from a serious injury or surgery that requires daily monitoring and care.
- 3. Medical treatment: Thomas may be undergoing a medical treatment that requires daily visits to the hospital, such as chemotherapy or radiation therapy.
- 4. Medical procedure: Thomas may be undergoing a medical procedure that requires daily visits to

SURELOCK (score: 4)

There could be several reasons why Thomas has to go to the hospital every day despite being healthy. Some possible reasons include:

- 1. Chronic illness: Thomas might have a chronic illness that requires daily medical attention. 2. Medical treatment: He could be undergoing a medical treatment that requires daily visits to the hospital. 3. Recovery period: Thomas might be recovering from an injury or surgery that requires daily monitoring and care.
- 4. Medical condition: He could be managing a medical condition that requires daily monitoring and treatment. 5. Medical procedure: Thomas might be undergoing a medical procedure that requires daily visits to the hospital.

Figure 9: Quality comparison of responses between Baseline vs. SureLock response. The example question uses the record from MT-bench with question id= 103. Responses are from LLaDA-8B-Instrut with KL threshold $\varepsilon=5e-4$, $N_{\rm gen}=128$, S=128.

Question: When rolling two dice, what is the probability that you roll a total number that is at least 3?

Baseline (score: 2)

There are $6 \times 6 = 36$ possible outcomes when rolling two dice. We need to count the number of outcomes where the sum of the dice is at least 3.

The possible sums are 3, 4, 5, and 6.
- Sum of 3: 2 outcomes (1,2), (2,1) - Sum of 4: 3 outcomes (1,3), (2,2), (3,1) - Sum of 5: 4 outcomes (1,4), (2,3), (3,

SURELOCK (score: 3)

There are $6 \times 6 = 36$ possible outcomes when rolling two dice. We will count the number of outcomes where the sum is at least 3.

- Sum = 3: (1,2), (2,1) - Sum = 4: (1,3), (2,2), (3,1) - Sum = 5: (1,4), (2,3), (3,2), (4,1) - Sum = 6: (1,5), (2,4), (3,3), (4,

Figure 10: Quality comparison of responses between Baseline vs. SURELOCK response. The example question uses the record from MT-bench with question id= 114. Responses are from LLaDA-8B-Instrut with KL threshold $\varepsilon=5e-4$, $N_{\rm gen}=128$, S=128.

Question: Write a function to find the highest common ancestor (not LCA) of two nodes in a binary tree.

Baseline (score: 3)

Here is a Python function to find the highest common ancestor (HCA) of two nodes in a binary tree:

"" class TreeNode: def __init__(self, val\overline{\text{0}}, left=None, right=None): self.val = val self.left = left self.right = right

def find_hca(root, p, q): if root is None: return None if root == p or root == q: return root left = find_hca(root.left, p, q) right = find_hca(root.right,

SURELOCK (score: 3)

Here's a Python function to find the highest common ancestor (HCA) of two nodes in a binary tree:

"" class TreeNode: def __init__(self, val=0, left=None, right=None): self.val = val self.left = left self.right = right

def find_hca(root, p, q): if root is None: return None if root == p or root == q: return root left = find_hca(root.left, p, q) right = find_hca(root.right,

Figure 11: Quality comparison of responses between Baseline vs. SURELOCK response. The example question uses the record from MT-bench with question id= 125. Responses are from LLaDA-8B-Instrut with KL threshold $\varepsilon=5e-4$, $N_{\rm gen}=128$, S=128.

Question: Please read the paragraph below and count how many times the words "Amazon", "river", and "you" appear. Please present the results in the format of "word, number of appearances" with each word on a separate line. Sort the lines in order of the number of appearances. The Amazon, a mesmerizing expanse of nature's wonders, is home to the legendary Amazon River. Flowing through awe-inspiring landscapes like the Amazon rainforest, the river weaves its way through Brazil, Colombia, and Peru, giving life to countless creatures. From the mighty jaguars prowling the Amazon jungle to the vibrant macaws soaring above the canopy, this remarkable region teems with biodiversity. Deep within the river's currents, magnificent pink river dolphins gracefully glide alongside piranhas and electric eels. Along the riverbanks, you'll find bustling cities like Manaus, where the urban meets the wild, and Iquitos, a gateway to the heart of the Amazon rainforest. As you venture further, the Amazon River reveals hidden gems like the captivating Anavilhanas Archipelago, a mosaic of islands brimming with rare species. Embark on an adventure, explore the enchanting Amazon River, and immerse yourself in a world teeming with life and untamed beauty.

Baseline (score: 3) Amazon, 2 river, 2 you, 2 **SURELOCK** (score: 3) Amazon, 2 river, 2 you, 2

Figure 12: Quality comparison of responses between Baseline vs. SURELOCK response. The example question uses the record from MT-bench with question id= 136. Responses are from LLaDA-8B-Instrut with KL threshold $\varepsilon=5e-4$, $N_{\rm gen}=128$, S=128.

1029

1030

Question: The city of Vega intends to build a bridge that will span the Vegona River, covering a distance of 1.8 kilometers. The proposed location falls within a seismically active area that has experienced several high-magnitude earthquakes. Given these circumstances, what would be the best approach to constructing the bridge?

1031 1032 1033

1034

1035

1036

1039

1040

1041

1043

1044

1046

1047

1048 1049

1050

1051

1052

Baseline (score: 7)

Given the seismically active nature of the area, the best approach to constructing the bridge would be to use a flexible and resilient design that can withstand the forces of earthquakes. This could involve using a combination of materials such as steel, concrete, and timber, as well as incorporating features such as flexible joints and base isolators to allow the bridge to move in response to seismic activity. Additionally, the bridge should be designed to be lightweight and compact to minimize the risk of damage during an earthquake. Finally, the bridge should be designed and constructed according to seismic codes and regulations to ensure its safety and stability in the event of earthquakes.

SURELOCK (score: 7)

Given the seismically active nature of the area, the best approach to constructing the bridge would be to use a flexible and resilient design that can withstand the effects of earthquakes. This could include using a combination of concrete and steel, as well as incorporating seismic isolation systems such as base isolation or rubber bearings. Additionally, the bridge should be designed to be flexible and able to absorb and dissipate energy from seismic waves, rather than trying to resist them. It is also important to conduct thorough research and testing to evaluate the bridge's performance under various earthquake scenarios to ensure that the bridge remains safe and stable in the face of seismic activity.

Figure 13: Quality comparison of responses between Baseline vs. SureLock response. The example question uses the record from MT-bench with question id= 147. Responses are from LLaDA-8B-Instrut with KL threshold $\varepsilon=5e-4$, $N_{\rm gen}=128$, S=128.

Question: Which methods did Socrates employ to challenge the prevailing thoughts of his time?

1058 1059

1062

1063

1064

1067

1068

1069

1070

1071

1057

Baseline (score: 6)

Socrates employed several methods to challenge the prevailing thoughts of his time. One of the most famous methods was the Socratic method, which involved asking a series of questions to uncover the underlying assumptions and beliefs of his listeners. This method encouraged his listeners to think critically and question their own beliefs.

Another method Socrates used was the concept of the "good life," which he believed was the only way to achieve true happiness. He argued that the good life involved living in accordance with reason and virtue, rather than pursuing material wealth or power.

Finally, Socrates used the concept of the "just city" to challenge the

SURELOCK (score: 6)

Socrates employed several methods to challenge the prevailing thoughts of his time. One of the most famous methods was the Socratic method, which involved asking a series of questions to encourage someone to think critically and examine their own beliefs. This method was used to uncover the underlying assumptions and beliefs of those around him.

Another method Socrates used was the use of irony. He would often use irony to highlight the absurdity of certain beliefs or ideas, and to encourage others to question their own assumptions.

Finally, Socrates also used the power of dialogue to challenge the prevailing beliefs of his time. He would often engage in debates

1074 1075

Figure 14: Quality comparison of responses between Baseline vs. SureLock response. The example question uses the record from MT-bench with question id= 158. Responses are from LLaDA-8B-Instrut with KL threshold $\varepsilon=5e-4$, $N_{\rm gen}=128$, S=128.

1078 1079

1077