

# CD<sup>4</sup>LM: Consistency Distillation and aDaptive Decoding for Diffusion Language Models

Yihao Liang<sup>1</sup>   Ze Wang<sup>2</sup>   Hao Chen<sup>2</sup>   Ximeng Sun<sup>2</sup>   Jialian Wu<sup>2</sup>  
 Xiaodong Yu<sup>2</sup>   Jiang Liu<sup>2</sup>   Emad Barsoum<sup>2</sup>   Zicheng Liu<sup>2</sup>

Niraj K. Jha<sup>1\*</sup>

<sup>1</sup> Princeton University   <sup>2</sup> Advanced Micro Devices, Inc.

## Abstract

Autoregressive large language models achieve strong results on many benchmarks, but decoding remains fundamentally latency-limited by sequential dependence on previously generated tokens. Diffusion language models (DLMs) promise parallel generation but suffer from a fundamental *static-to-dynamic misalignment*: Training optimizes local transitions under fixed schedules, whereas efficient inference requires adaptive “long-jump” refinements through unseen states. Our goal is to enable *highly parallel decoding for DLMs with low number of function evaluations* while preserving generation quality. To achieve this, we propose **CD<sup>4</sup>LM**, a framework that decouples training from inference via **Discrete-Space Consistency Distillation (DSCD)** and **Confidence-Adaptive Decoding (CAD)**. Unlike standard objectives, DSCD trains a student to be trajectory-invariant, mapping diverse noisy states directly to the clean distribution. This intrinsic robustness enables CAD to dynamically allocate compute resources based on token confidence, aggressively skipping steps without the quality collapse typical of heuristic acceleration. On GSM8K, **CD<sup>4</sup>LM** matches the LLaDA baseline with a **5.18×** wall-clock speedup; across code and math benchmarks, it pushes the accuracy-efficiency Pareto frontier, achieving a **3.62×** mean speedup while improving average accuracy. Code is available at <https://github.com/yihao-liang/CDLM>.

## 1 Introduction

Inference latency is a primary deployment bottleneck for large language models. Autoregressive large language models (AR-LLMs) achieve strong performance on language, code, and mathematical reasoning benchmarks [Achiam et al., 2023, Grattafiori et al., 2024, Hendrycks et al., 2021, Austin et al., 2021b, Chen, 2021], but decoding remains inherently sequential: Generating a length- $L$  output requires  $L$  dependent forward passes, limiting parallelism and making end-to-end latency scale with sequence length.

Diffusion language models (DLMs) offer a different generation paradigm [Austin et al., 2021a, Li et al., 2022, Gong et al., 2022]. They generate text through iterative denoising in the token space: Starting from a partially masked sequence, the model refines many positions in parallel at each denoising step. This enables bidirectional context and amortizes computation across the sequence, making DLMs particularly appealing for structured generation, such as code and mathematical solutions.

\*Contact: yhliang@princeton.edu, jha@princeton.edu.

Despite the above promise, existing DLMs face two critical limitations when deployed under practical budget-limited decoding constraints, i.e., low number of function evaluations (NFE). First, there is a structural *static-to-dynamic* misalignment between training and inference. Standard training rigidly optimizes local denoising transitions (e.g.,  $t \rightarrow t - \Delta t$ ) under fixed schedules, effectively learning a static vector field. However, low-NFE inference needs the model to perform agile “long-jump” refinements (e.g.,  $t \rightarrow t - k\Delta t$ ) through unseen regions of the state space. This forces the model to traverse intermediate masked states that are mathematically disjoint from the training distribution, limiting its ability to exploit adaptive acceleration without collapsing. Second, current methods suffer from rigidity in compute allocation. Most DLM decoders rely on predetermined unmasking schedules that treat all instances as equally difficult. While block-wise decoding acts as a structural regularizer to prevent collapse, it remains computationally wasteful for easy instances and insufficient for hard ones. Consequently, previous works fail to simultaneously achieve stability and efficiency. Aggressive step-reduction heuristics on standard models often destabilize text generation due to distribution shifts. Conversely, training-free adaptive methods that seek to mitigate these shifts often incur significant computational overhead (e.g., complex sorting or verification logic), which prevents their theoretical step reductions from translating fully into wall-clock speedups.

These limitations motivate a core challenge: *How can we decouple the training objective from fixed schedules to enable stable, adaptive computation allocation?* We answer this with **CD<sup>4</sup>LM**, a unified framework that bridges the gap between robust training and flexible inference.

On the training side, we propose **Discrete-Space Consistency Distillation (DSCD)**. Instead of overfitting to a single fixed trajectory, DSCD leverages a *Rao-Blackwellized* objective to train a student model that matches the teacher’s conditional posterior across a diverse range of masking patterns. This transforms the model from a fixed-schedule denoiser into a robust *refinement operator*, mathematically capable of handling the irregular states produced by aggressive decoding policies. On the inference side, leveraging this trajectory-invariant capability, we introduce **Confidence-Adaptive Decoding (CAD)**. Unlike heuristic step-skipping, CAD dynamically commits high-confidence tokens while deferring uncertain ones. This policy is fundamentally robust: It not only enhances efficiency within block-based schemes but also stabilizes pure diffusion (full-sequence) generation, preventing the quality collapse typical of standard adaptive baselines. In our framework, we synergize CAD with block-wise decoding to achieve the optimal Pareto frontier between global stability and local adaptivity.

Our contributions are as follows.

1. **Identification of the training-inference misalignment as a limiting factor for efficiency.** We identify a critical oversight in previous works: While adaptive decoding policies can reduce latency, their potential is strictly bottlenecked by standard training objectives. We show that relying on inference-time heuristics to bridge this gap introduces unnecessary control overhead. We verify that maximizing generation speed requires a synergistic approach: Efficiency is not only a decoding search problem but a training alignment problem. By decoupling the model from fixed trajectories, we unlock superlinear speedups that are unattainable by adaptive decoding strategies alone.
2. **A unified framework for robust and efficient diffusion.** We propose **CD<sup>4</sup>LM**, which synergizes DSCD with CAD. DSCD trains a schedule-agnostic student robust to the intermediate states induced by adaptive policies, thereby enabling CAD to aggressively skip steps based on token confidence without inducing collapse. This joint design reconciles the conflict between generation stability and decoding efficiency, providing a general-purpose solution applicable to both block-wise and pure diffusion settings.
3. **Systematic Pareto-frontier improvements.** Across math and code benchmarks, **CD<sup>4</sup>LM** achieves substantial speedups without sacrificing quality. On GSM8K, it matches the fixed-step baseline (77.6% vs. 77.4%) with a **5.18× wall-clock acceleration**; on HumanEval and MATH500, it improves accuracy while delivering **3.62× speedup**, strictly dominating standard diffusion baselines.

## 2 Related Work

We review four lines of work that are most relevant to our approach: (i) AR-LLMs and their parallel decoding methods, (ii) DLMs and their training objectives, (iii) inference strategies and parallel

decoding for diffusion models, and (iv) diffusion-based models for code and mathematical reasoning. Existing work has established that diffusion architectures can scale to competitive language modeling performance with attractive parallelism, but has left open how to decouple training from inference and allocate computation adaptively under tight NFE budgets. **CD<sup>4</sup>LM** keeps the diffusion backbone while redesigning the training objective and decoding policy to better support adaptive NFE allocation.

## 2.1 Autoregressive and Diffusion Language Models

AR-LLMs remain the dominant architecture for natural language, code, and mathematical reasoning. Representative systems, such as GPT-4 [Achiam et al., 2023], the Llama family [Grattafiori et al., 2024], Qwen [Bai et al., 2023], and DeepSeek [Liu et al., 2024], achieve state-of-the-art results on a wide range of benchmarks. However, AR decoding is inherently sequential: To generate a sequence of length  $L$ , the model must perform  $L$  forward passes, each conditioned on the previously generated prefix, which limits parallelism and leads to noticeable latency for long-context tasks. To mitigate this bottleneck, several works propose draft-verify style acceleration methods, including speculative decoding [Leviathan et al., 2023], Medusa [Cai et al., 2024], and EAGLE [Li et al., 2024], which parallelize the prediction of multiple future tokens within a single forward pass while approximately preserving the original output distribution.

DLMs offer a different generation paradigm based on iterative denoising and parallel token updates. Inspired by the success of continuous diffusion models in image generation [Ho et al., 2020], a series of works, such as Diffusion-LM [Li et al., 2022], Likelihood-Based DLMs [Gulrajani and Hashimoto, 2023], and Masked Diffusion Models [Sahoo et al., 2024], extend the diffusion process to discrete token spaces. More recently, Large Language Diffusion Models (LLaDA) [Nie et al., 2025] demonstrate that a purely diffusion-based architecture can match or even surpass comparably sized AR models. However, a critical limitation of existing DLMs, including LLaDA, is the coupling of training schedules with inference budgets. They typically assume a fixed number of denoising steps derived from the training noise schedule, lacking the flexibility to dynamically trade computation for quality. **CD<sup>4</sup>LM** addresses this by decoupling the decoding trajectory from the training schedule, enabling efficient generation even under tight NFE budgets where standard schedules fail.

## 2.2 Training Objectives and Distillation for Diffusion Language Models

Early work on discrete diffusion models focused on variational objectives and step-wise denoising losses. Austin et al. [2021a] introduce a variational lower bound with an auxiliary cross-entropy term to improve likelihood on text data. For masked or absorbing diffusion, Sahoo et al. [2024] and Shi et al. [2024] rewrite the training objective as a mixture of classical masked language modeling losses or as a continuous-time weighted integral of cross-entropy terms, yielding more scalable training; under a uniform-state formulation, Zhu et al. [2025] demonstrate that a simplified denoising loss over only noise-replaced tokens can outperform objectives based on Evidence Lower Bound (ELBO). Together, these works clarify the likelihood and optimization landscape of discrete diffusion, providing stable and efficient training recipes for DLMs.

A second line of research moves DLMs closer to AR-LLMs, either by reusing AR backbones or by aligning the training distribution with the inference trajectory. Asada and Miwa [2025] point out that standard discrete diffusion models are trained to denoise gold tokens corrupted by random noise, whereas, at inference time, they denoise self-generated tokens. They introduce a two-step diffusion training scheme with step-aware losses and a curriculum that gradually increases the probability of using self-generated text. In parallel, He et al. [2025] formulate the choice of denoising trajectories as a sequential decision-making problem and use reinforcement learning under the same progressive refining schedule used at inference. These methods reduce mismatch along a fixed schedule, complementary to approaches that train students robust to diverse intermediate states.

To reduce NFE, some works explore one-step or few-step distillation. Consistency Models [Song et al., 2023] propose a consistency objective in continuous spaces that learns a generator mapping noisy inputs at arbitrary timesteps directly to clean samples, enabling distillation of pretrained diffusion models. In language, DLM-One [Chen et al., 2025] trains a student in the continuous embedding space to match the score function of a pretrained DLM. While promising for aggressive step reduction, these approaches typically rely on continuous or embedding-space diffusion. The distilled students

are closely tied to a particular teacher and sampling trajectory, thus limiting flexibility when changing noise schedules or decoding policies.

In contrast, our DSCD method operates directly on token-level masked sequences and trains a refinement student that is broadly trajectory-invariant. Unlike previous distillation works that tie the student to a specific teacher sampling path or require continuous embedding spaces [Chen et al., 2025, Song et al., 2023], DSCD exposes the student to a diverse range of noise levels and masking patterns. This critical difference ensures that our student model remains robust to the irregular intermediate states induced by aggressive step-skipping strategies, a property that standard likelihood-based or specific-trajectory distillation objectives fail to guarantee.

### 2.3 Inference and Parallel Decoding in Diffusion Language Models

On the inference side, DLMs typically use an iterative denoise decoding process. A common method is to unmask a fixed number of the most confident tokens at every step [Sahoo et al., 2024]. However, such predefined unmask schedules treat all tokens equally, leading to inefficiencies on structured tasks. Building on these basic unmask policies, a growing body of work seeks further acceleration by dynamically controlling parallelism. Adaptive Parallel Decoding [Israel et al., 2025] adjusts the number of tokens updated in parallel based on error estimates. Dimple [Yu et al., 2025] introduces confidence-based schemes for vision-language tasks. Fast-dLLM [Wu et al., 2025] proposes a training-free parallel decoding framework combined with block-wise KV caching.

While effective in reducing the number of denoising steps, such training-free methods typically rely on inference-time heuristics to identify stable tokens. Since the base model is not explicitly trained for aggressive parallel decoding, these methods must employ additional selection logic (e.g., sorting confidence scores or dynamic thresholding) at every step to filter out low-confidence predictions. This introduces computational overhead, which can prevent the reduction in NFE from fully translating into wall-clock speedup (i.e., yielding sub-linear speedups). Other recent works also explore hybrid drafting: verification and caching [Wei et al., 2025, Christopher et al., 2025, Liu et al., 2025].

**CD<sup>4</sup>LM** differentiates itself by identifying that efficiency is not solely a decoding search problem but also an alignment problem. Our decoding strategy (CAD) is uniquely enabled by our training objective (DSCD). By aligning the student model’s distribution with aggressive decoding trajectories, we minimize the need for complex filtering heuristics. This enables a lightweight decoding policy in which algorithmic NFE reductions translate efficiently, and often superlinearly, into wall-clock gains.

### 2.4 Diffusion Language Models for Code and Mathematical Reasoning

DLMs have recently been applied to structured generation tasks. DiffuCoder [Gong et al., 2025] trains masked DLMs on large-scale code corpora. Mercury Coder [Labs et al., 2025] scales discrete diffusion architectures to commercial settings tailored to code completion. Beyond code, diffusion models have also been explored for combinatorial reasoning [Ye et al., 2024, Huang et al., 2025].

Most prior works in this domain focus on developing specialized architectures or training domain-specific models from scratch (e.g., DiffuCoder, Mercury Coder). While effective, this approach is resource-intensive and creates silos between domains. In contrast, **CD<sup>4</sup>LM** proposes a general-purpose efficiency framework that can be applied to existing pretrained diffusion backbones (like LLaDA) without architectural changes. We demonstrate that by simply aligning the training and decoding objectives, a general-purpose DLM can achieve Pareto-superior performance on specialized code and mathematics benchmarks, outperforming baselines that lack adaptive computation allocation.

## 3 Method

**CD<sup>4</sup>LM** is designed to resolve the structural misalignment between fixed diffusion training schedules and the need for flexible, budget-aware inference. To bridge this gap, we first propose DSCD, which transforms the model from a rigid schedule-dependent denoiser into a *trajectory-invariant* refinement operator capable of handling arbitrary intermediate states. This intrinsic robustness serves as the necessary foundation for CAD, enabling the model to dynamically allocate computation based on token-level confidence without the quality collapse observed in standard acceleration baselines. Under

this framework, our method comprises three specific modules: (i) an absorbing-state discrete diffusion backbone inherited from LLaDA, which we briefly review in Sect. 3.1; (ii) a DSCD scheme (Sect. 3.2) that trains a student to be approximately trajectory-invariant along the teacher’s masking process; and (iii) a CAD policy (Sect. 3.3) that dynamically allocates NFE across tokens and instances. Taken together, these modules decouple the trained model from any fixed diffusion schedule and enable flexible low-NFE decoding.

### 3.1 Absorbing-State Discrete Diffusion Backbone

We briefly review the absorbing-state discrete diffusion formulation we build on and fix notation.

**Data representation.** Let  $\mathcal{D} = \{(x, y)\}$  be a corpus of conditional generation tasks, where  $x = (x_1, \dots, x_{L_x})$  is the prompt (e.g., a coding problem or mathematics question) and  $y = (y_1, \dots, y_{L_y})$  is the target sequence (e.g., the solution program or derivation). We concatenate them into a single sequence  $z = (x; y) = (z_1, \dots, z_L)$ , with  $L = L_x + L_y$ . The base vocabulary  $\mathcal{V}$  is augmented with a special absorbing mask token  $\mathfrak{m} \notin \mathcal{V}$ , and an end-of-sequence token  $\langle \text{EOS} \rangle$ , giving the extended vocabulary:  $\mathcal{V}^+ = \mathcal{V} \cup \{\mathfrak{m}, \langle \text{EOS} \rangle\}$ . For chain-of-thought supervision, we place the entire reasoning trace together with the final answer in the target region ( $i > L_x$ ); hence, the model explicitly learns to refine both intermediate reasoning and final answers. For code completion benchmarks, the natural-language problem description and function signature belong to the prompt region, while the full solution body is treated as target tokens.

**Forward absorbing process.** We adopt the absorbing discrete diffusion framework for text [Austin et al., 2021a, Shi et al., 2024, Sahoo et al., 2024, Nie et al., 2025], which models corruption as stochastic erasure rather than additive noise. Let  $t \in [0, 1]$  be a continuous noise level with a monotone schedule  $\alpha(t) \in [0, 1]$  satisfying  $\alpha(0) = 0$  and  $\alpha(1) \approx 1$ . The forward process  $q_t(\tilde{z} | z)$  gradually destroys information in the target region while preserving the prompt:

$$q_t(\tilde{z} | z) = \prod_{i=1}^L q_t(\tilde{z}_i | z_i)$$

$$q_t(\tilde{z}_i | z_i) = \begin{cases} \delta_{\tilde{z}_i = z_i}, & i \leq L_x \quad (\text{prompt preservation}), \\ (1 - \alpha(t)) \delta_{\tilde{z}_i = z_i} + \alpha(t) \delta_{\tilde{z}_i = \mathfrak{m}}, & i > L_x \quad (\text{target corruption}), \end{cases} \quad (1)$$

where  $\delta$  is the Kronecker delta. Once a token is replaced by  $\mathfrak{m}$ , it remains masked for all larger  $t$ ; hence, the forward chain is an absorbing Markov process. We denote the masked set by  $\mathcal{M}(\tilde{z}) = \{i : \tilde{z}_i = \mathfrak{m}\}$ . In practice, we do not simulate this Markov chain step-by-step. Given  $z$  and  $t$ , we instead sample a Bernoulli mask  $b \in \{0, 1\}^{L_y}$  with  $\Pr(b_j = 1) = \alpha(t)$  for target positions and set  $\tilde{z} = \text{Mask}(z, b)$ , which matches  $q_t(\tilde{z} | z)$  in (1). In all experiments, we reuse the continuous-time masking schedule from LLaDA and adopt the same choices of  $\alpha(t)$ , the sampling distribution  $\rho(t)$  over  $t$ , and the weighting function  $w(t)$ .

**Diffusion teacher and ELBO-style objective.** A DLM parameterized by  $\phi$  defines a reverse model that predicts clean tokens at the masked positions of a corrupted sequence:

$$p_\phi(z | \tilde{z}, t) = \prod_{i \in \mathcal{M}(\tilde{z})} p_\phi(z_i | \tilde{z}, t), \quad (2)$$

where each  $p_\phi(z_i | \tilde{z}, t)$  is a categorical distribution over  $\mathcal{V}^+$ . We follow LLaDA [Nie et al., 2025] and parameterize  $p_\phi$  with a Transformer that maps  $(\tilde{z}, t)$  to logits over the vocabulary at all positions; unmasked positions are simply copied from the input. For absorbing discrete diffusion, the log-likelihood admits an evidence lower bound whose dominant term reduces to a time-weighted masked cross-entropy [Austin et al., 2021a, Sahoo et al., 2024, Shi et al., 2024]. We adopt the simplified continuous-time objective used in masked DLMs:

$$\mathcal{L}_{\text{teacher}}(\phi) = \mathbb{E}_{(x, y) \sim \mathcal{D}} \mathbb{E}_{t \sim \rho(t)} \mathbb{E}_{\tilde{z} \sim q_t(\cdot | z)} \left[ w(t) \frac{1}{|\mathcal{M}(\tilde{z})|} \sum_{i \in \mathcal{M}(\tilde{z})} -\log p_\phi(z_i | \tilde{z}, t) \right], \quad (3)$$

where  $\rho(t)$  and  $w(t)$  are inherited from LLaDA. Optimizing (3) yields a high-quality diffusion teacher, but sampling requires iterating over a long schedule  $t_1 > \dots > t_T$ , leading to high decoding latency.

In practice,  $p_\phi$  is instantiated as the publicly released LLaDA checkpoint of the corresponding size and kept frozen throughout student training; we do not modify its architecture, training objective, or noise schedule. All improvements reported in this article, therefore, come from the distilled student and our decoding policy, rather than from retraining the teacher.

### 3.2 Discrete-Space Consistency Distillation

Our goal is to obtain a student  $p_\theta$  that can reliably refine partially masked sequences in a few steps, without being tied to any particular diffusion schedule. Conceptually, we would like  $p_\theta$  to be approximately *trajectory-invariant* along the absorbing diffusion process of the teacher.

**Ideal trajectory invariance.** For a clean sequence  $z$  and its absorbing diffusion trajectory  $\{\tilde{z}_t\}_{t \in [0,1]}$ , generated by (1), an ideal consistency-style student would satisfy

$$p_\theta(z \mid \tilde{z}_t) \approx p_\theta(z \mid \tilde{z}_{t'}) \approx p_\phi(z \mid \tilde{z}_t, t) \quad \text{for all } t, t' \in [0, 1]. \quad (4)$$

Unlike the teacher  $p_\phi$ , which conditions on explicit time  $t$ , the student  $p_\theta$  learns to infer the effective noise level directly from the masked input  $\tilde{z}$ , enabling time-agnostic inference. This requirement is too strong to enforce directly: It couples all time points on the trajectory and requires matching full sequence distributions. In practice, DSCD implements a weaker but tractable surrogate that enforces *pairwise* consistency between stochastic views of the same  $z$  and anchors them to the data distribution.

**Paired teacher-subset masking.** For each training example, we first sample a student mask ratio  $r_S \sim \mathcal{U}(r_S^{\min}, r_S^{\max})$  and then set a lighter teacher ratio  $r_T = r_S \cdot u$ , where  $u \sim \mathcal{U}(u_{\min}, u_{\max})$  with  $u_{\max} < 1$  (thus  $r_T < r_S$ ). We convert each ratio to its diffusion timestep using the same noise schedule, yielding  $t_S$  and  $t_T$ . Let  $n_S = \lfloor L_y r_S \rfloor$  and  $n_T = \lfloor L_y r_T \rfloor$ . We form the student masked set  $\mathcal{M}_S$  by uniformly sampling (without replacement)  $n_S$  target positions from  $\{L_x + 1, \dots, L\}$ . To construct the teacher mask, we then uniformly sample a subset  $\mathcal{M}_T \subseteq \mathcal{M}_S$  of size  $n_T$ . This nested masking makes the teacher informationally richer (it conditions on a strict superset of the student’s visible context); hence, the consistency loss is evaluated on  $\mathcal{M}_S$ . The overall training pipeline is depicted in Fig. 1.

At a high level, DSCD then combines a reconstruction loss and a Kullback-Leibler (KL) divergence-based consistency loss computed on  $\mathcal{M}_S$  to (i) anchor the student to the data distribution at its own masked positions and (ii) align the student with the teacher’s softer predictions under the same corruption pattern, yielding a practical relaxation of the trajectory-invariance property in (4).

**Variance reduction via Rao-Blackwellization.** We theoretically justify our choice of the *teacher-subset masking* scheme ( $\mathcal{M}_T \subseteq \mathcal{M}_S$ ) not merely as a heuristic, but as a variance reduction technique rooted in statistical decision theory. Intuitively, if  $\mathcal{M}_T$  were independent of  $\mathcal{M}_S$ , the teacher might lack access to tokens visible to the student, causing it to marginalize over “blind” positions and produce high-variance targets. By enforcing  $\mathcal{M}_T \subseteq \mathcal{M}_S$ , we ensure the teacher always conditions on a strict superset of the student’s information. As formally proven in Appendix B, this effectively applies Rao-Blackwellization [Casella and Robert, 1996] to the gradient estimator: By conditioning on the additional information  $\mathcal{M}_S \setminus \mathcal{M}_T$ , we strictly minimize the conditional variance of the distillation target compared to independent masking.

**Reconstruction loss.** To anchor the student to the data distribution and avoid drifting too far from the ground-truth tokens, we include a reconstruction term that predicts  $z$  at the student’s masked positions:

$$\mathcal{L}_{\text{recon}}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \mathbb{E}_{r_S, \mathcal{M}_S} \left[ \frac{1}{|\mathcal{M}_S|} \sum_{i \in \mathcal{M}_S} -\log p_\theta(z_i \mid \tilde{z}^S) \right]. \quad (5)$$

This term coincides with the standard masked language modeling cross-entropy objective when  $p_\theta$  is used as the denoiser, and can be viewed as learning a strong supervised initialization for the student.

**Consistency loss via KL distillation.** The core consistency signal comes from aligning the student’s predictions with the teacher’s distribution at the student’s masked positions. Given teacher logits

### 1. Input & Masking Process

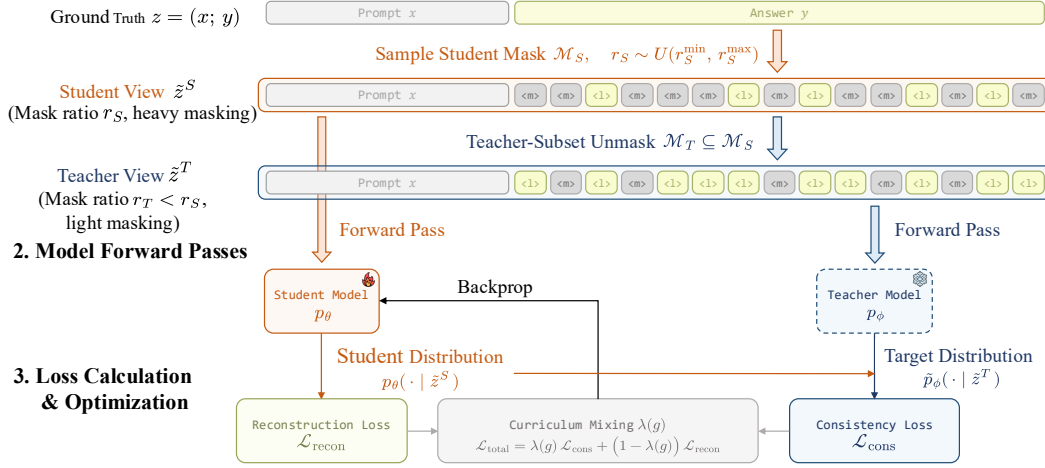


Figure 1: **Overview of the DSCD pipeline.** (1) **Input & Masking:** We employ a *teacher-subset masking* strategy where the teacher’s mask  $\mathcal{M}_T$  is strictly sampled from the student’s mask  $\mathcal{M}_S$  ( $\mathcal{M}_T \subseteq \mathcal{M}_S$ ). This ensures the teacher always conditions on a superset of the student’s context, acting as a lower-variance guide. (2) **Forward & Optimization:** The student  $p_\theta$  predicts tokens from the heavily masked view, optimized jointly by a reconstruction loss  $\mathcal{L}_{\text{recon}}$  (anchoring to ground truth) and a consistency loss  $\mathcal{L}_{\text{cons}}$  (aligning with the frozen teacher’s soft targets). (3) **Curriculum:** A dynamic schedule  $\lambda(g)$  governs the transition from pure distillation to supervised refinement.

$l_\phi^{(i)}(\tilde{z}^T)$ , we define a temperature-scaled target distribution

$$\tilde{p}_\phi(\cdot | \tilde{z}^T)_i = \text{softmax}\left(\frac{l_\phi^{(i)}(\tilde{z}^T)}{\tau}\right), \quad \tau \geq 1. \quad (6)$$

We then minimize the KL divergence between teacher and student distributions on the student mask  $\mathcal{M}_S$ :

$$\mathcal{L}_{\text{cons}}(\theta) = \tau^2 \cdot \mathbb{E}_{(x,y) \sim \mathcal{D}} \mathbb{E}_{r_S, r_T, \mathcal{M}_S, \mathcal{M}_T} \left[ \frac{1}{|\mathcal{M}_S|} \sum_{i \in \mathcal{M}_S} \text{KL}(\tilde{p}_\phi(\cdot | \tilde{z}^T)_i \| p_\theta(\cdot | \tilde{z}^S)_i) \right]. \quad (7)$$

The  $\tau^2$  scaling compensates for the reduced gradients caused by temperature smoothing, following standard knowledge distillation practice [Hinton et al., 2015].

**Curriculum mixing schedule and intuition.** Jointly optimizing (5) and (7) from scratch is unstable because the student is initially far from the teacher and the teacher’s predictions at high noise levels can be noisy. We, therefore, use a curriculum over the normalized training progress  $g \in [0, 1]$ :

$$\mathcal{L}_{\text{total}}(\theta, g) = \lambda(g) \mathcal{L}_{\text{cons}}(\theta) + (1 - \lambda(g)) \mathcal{L}_{\text{recon}}(\theta), \quad (8)$$

where  $\lambda(g)$  gradually decreases from  $\lambda_0$  to  $\lambda_1$  using a cosine schedule with an initial warmup phase. Early in training ( $g \approx 0$ ), the loss is dominated by  $\mathcal{L}_{\text{cons}}$ ; hence, the student learns to imitate the teacher’s soft predictions and inherits its denoising behavior. As training progresses,  $\lambda(g)$  decreases and the reconstruction term gradually takes over, anchoring the student to the ground-truth targets and improving accuracy. Empirically, we find that pure distillation (setting  $\lambda(g) \equiv 1$ ) tends to be unstable and underperforms, whereas the curriculum in (8) yields faster convergence and better final accuracy.

**Theoretical interpretation.** While DSCD is empirically motivated, it possesses a rigorous theoretical grounding. By enforcing *teacher-subset masking* ( $\mathcal{M}_T \subseteq \mathcal{M}_S$ ), we effectively construct a stochastic filtration of information where the teacher always holds a strictly finer information set

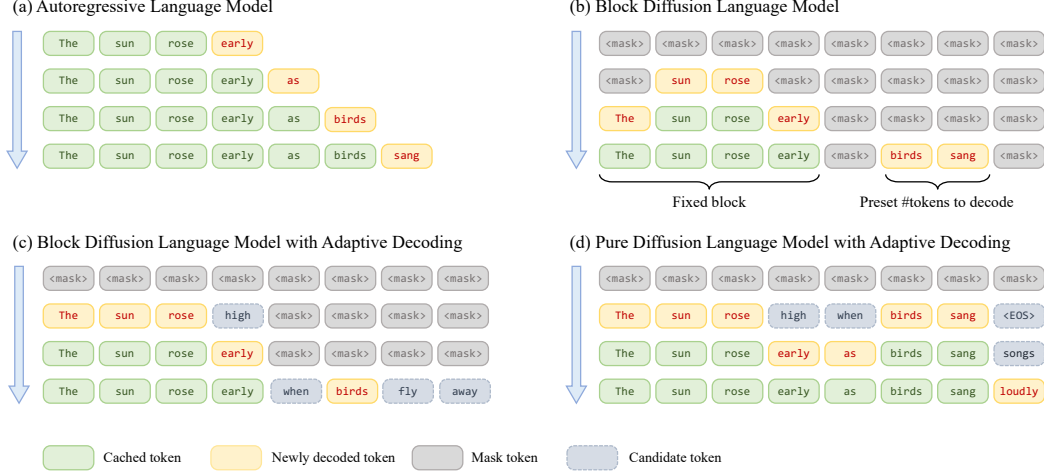


Figure 2: **Overview of different decoding paradigms.** (a) Autoregressive decoding generates tokens left-to-right. (b) Block diffusion decodes a preset number of masked tokens per step within a fixed block schedule. (c) Block diffusion with adaptive decoding: Within each block, we *adaptively* select a high-confidence set of masked positions to unmask, enabling variable token counts (and selected positions) across steps while preserving the block-wise decoding structure. (d) The same adaptive rule can also be applied to *pure* diffusion, iteratively denoising the entire sequence starting from all `<mask>` tokens.

than the student. In Appendix A, we formally prove that minimizing  $\mathcal{L}_{\text{cons}}$  is equivalent to training the student to approximate the Martingale projection [Vincent, 2011] of the teacher’s belief state. This guarantees that the student learns the expected trajectory of the teacher, satisfying the statistical definition of trajectory invariance.

For a complete algorithmic description of the training procedure, including the details of mask sampling and the curriculum schedule, readers can refer to the pseudocode provided in Appendix D.1.

### 3.3 Confidence-Adaptive Decoding

After distillation, we discard the teacher and decode using the student model only. Our CAD is a generic unmasking policy. In this work, we instantiate it within the LLaDA block diffusion framework (Fig. 2(c)); pure diffusion (Fig. 2(d)) is recovered as a special case by setting  $b = L_{\text{gen}}$ .

**Block diffusion schedule.** Unlike autoregressive models that generate variable-length sequences by dynamically appending tokens, diffusion-based decoding operates on a pre-allocated, fixed-size canvas. We set a maximum target horizon  $L_{\text{gen}}$  (e.g., matching the training sequence length or a system budget) and pad the initial state with masks. The actual sequence length is determined dynamically during decoding via the end-of-sequence (EOS) blocking mechanism. Given  $L_{\text{gen}}$  and block size  $b$ , we partition positions into contiguous blocks  $\{\mathcal{B}_j\}_{j=1}^J$ , where  $\mathcal{B}_j = \{(j-1)b+1, \dots, \min(jb, L_{\text{gen}})\}$  and  $J = \lceil L_{\text{gen}}/b \rceil$ . Decoding proceeds left-to-right with an active block index  $j$ . Tokens from finished blocks are cached (frozen), while future blocks remain masked until activated. Note that while our primary experiments use full-attention recomputation to isolate the algorithmic gains of DSCD, this block-wise formulation is structurally compatible with approximate KV caching schemes [Wu et al., 2025].

**State, masked set, and eligible set.** Let  $\tilde{z}^{(s)} \in \mathcal{V}^{L_{\text{gen}}}$  denote the partially-masked state at decoding step  $s$  and let  $\mathcal{M}^{(s)} = \{i \in [L_{\text{gen}}] : \tilde{z}_i^{(s)} = \text{m}\}$  be the masked set. CAD operates on the *eligible set* within the active block:

$$\mathcal{E}^{(s)} = \mathcal{M}^{(s)} \cap \mathcal{B}_j, \quad (9)$$

which reduces to  $\mathcal{E}^{(s)} = \mathcal{M}^{(s)}$  when  $b = L_y$  (pure diffusion).



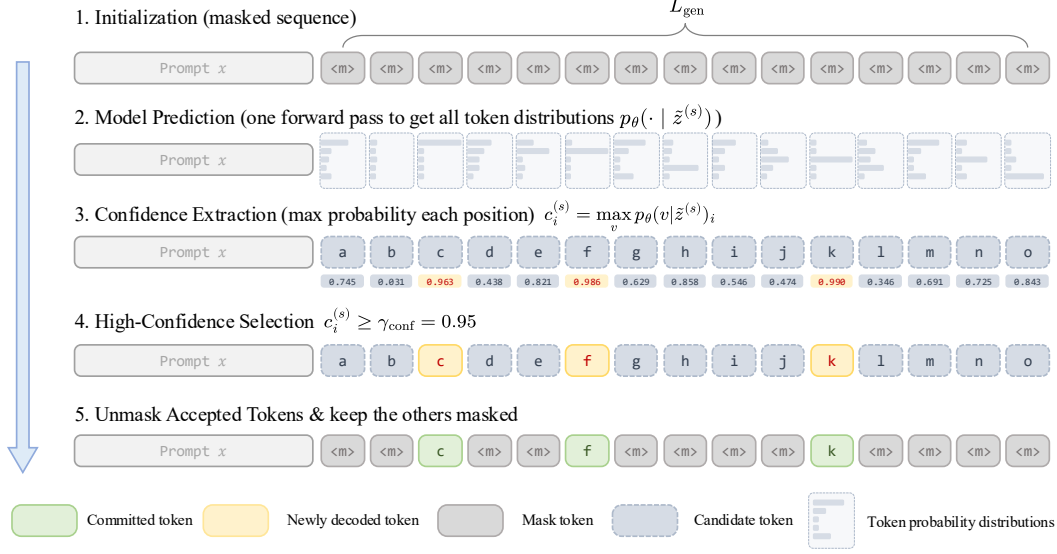


Figure 3: **Illustration of CAD.** In each step, the model predicts token distributions for all masked positions in parallel. We compute the confidence score  $c_i^{(s)}$  based on the maximum probability and selectively unmask tokens that satisfy the threshold  $c_i^{(s)} \geq \gamma_{\text{conf}}$ . Low-confidence tokens remain masked for subsequent iterations.

**Confidence-adaptive set selection.** We treat the decoding process as a dynamic risk-efficiency trade-off. Given the student predictions  $p_{\theta}(\cdot | \tilde{z}^{(s)})$ , we compute a confidence score  $c_i^{(s)}$  for each  $i \in \mathcal{E}^{(s)}$  (e.g.,  $c_i^{(s)} = \max_v p_{\theta}(v | \tilde{z}^{(s)})_i$ ) as a proxy for correctness probability. We first form a candidate set  $\hat{\mathcal{S}}^{(s)} = \{i \in \mathcal{E}^{(s)} : c_i^{(s)} \geq \gamma_{\text{conf}}\}$  containing tokens where the model’s certainty outweighs the risk of error. To stabilize the dynamic trajectory, we then clamp the commit size:

$$k^{(s)} = \text{clip}(|\hat{\mathcal{S}}^{(s)}|, k_{\min}, k_{\max}), \quad \mathcal{S}^{(s)} = \text{TopK}(\{c_i^{(s)}\}_{i \in \mathcal{E}^{(s)}}, k^{(s)}), \quad (10)$$

and commit predictions only on  $\mathcal{S}^{(s)}$ . This greedy policy approximates the optimal stopping rule for per-step decoding risk (see derivation in Appendix C).

**EOS blocking and termination.** To prevent premature termination, we apply EOS blocking with threshold  $\beta_{\text{EOS}}$ . In block diffusion, once the active block is fully resolved (i.e.,  $\mathcal{E}^{(s)} = \emptyset$ ), we advance  $j \leftarrow j + 1$  and continue; decoding terminates when all blocks are decoded or when reaching the maximum NFE budget  $S_{\text{max}}$ . CAD guarantees monotonic progress within each active block (at least one token is committed whenever  $\mathcal{E}^{(s)} \neq \emptyset$ ); a formal step bound is provided in Appendix C.

**Complexity and overhead.** Unlike heuristics that require sorting candidate tokens or managing dynamic buffer states (which introduce CPU-GPU synchronization overhead), our CAD controller is implemented as a lightweight, fully batched tensor masking operation. The selection cost is negligible compared to the Transformer forward pass  $\mathcal{O}(B N_{\text{layer}}(L^2 D + L D^2))$ , allowing the theoretical NFE reduction to translate directly into latency savings.

For the complete pseudocode of the CAD algorithm, which details the interaction between block-wise diffusion and the dynamic acceptance policy, readers can refer to Appendix D.2.

## 4 Experimental Setup

We discuss the experimental setup next.

## 4.1 Models and Training

**Teacher and student.** We use **LLaDA-8B-Instruct** [Nie et al., 2025] as the teacher. The student is initialized from the teacher’s pretrained weights. All experiments use BF16 precision and a maximum context length of 1024.

**Training data.** We train our student model on a 200K-sample subset of **OpenCodeInstruct** [Ahmad et al., 2025], a comprehensive instruction-following code generation dataset. To ensure high-quality distillation, we filter for samples with valid solutions and reserve 5% (10K samples) for validation. For mathematical reasoning experiments, we use the **GSM8K** [Cobbe et al., 2021] training split and use the same training protocol.

## 4.2 Benchmarks

**Code generation.** We evaluate the model in a zero-shot setting on HumanEval [Chen, 2021] and three-shot on MBPP [Austin et al., 2021b], and also report results on the stricter variants: HumanEval+ and MBPP+. We report functional correctness using pass@1 and pass@5.

**Mathematical reasoning.** We report accuracy on GSM8K [Cobbe et al., 2021] and MATH500 [Hendrycks et al., 2021] in a zero-shot setting.

**Evaluation framework.** All evaluations are implemented based on the open-sourced **DAEDAL**’s evaluation codebase [Li et al., 2025],<sup>2</sup>, which standardizes prompting, post-processing, and metric computation for LLaDA-style DLMS.

## 4.3 Training Protocol

We use a global batch size of 64 on 8×AMD MI250 GPUs, AdamW with learning rate  $5 \times 10^{-6}$ , cosine decay, and 10% warmup, and train for three epochs. Detailed hyperparameters and prompt templates are provided in Appendix D.3.

## 4.4 Inference Protocols

**Generation paradigms.** We consider three generation paradigms: (i) Sequential ( $b = 1$ ), (ii) Block Diffusion (block size  $b = 32$ ), and (iii) Pure Diffusion ( $b = L_{\text{gen}}$ ).

For mathematical reasoning benchmarks, we use block diffusion with  $b = 32$ , which is the best-performing configuration in our tuning process and is used throughout the main mathematical experiments. We use pure diffusion for code generation. Unless specified otherwise, we use a maximum generation length of  $L_{\text{gen}} = 256$ .

**Sampling for pass@k.** We compute pass@1 using greedy  $\tau_{\text{samp}} = 0$  decoding and pass@5 using sampling with temperature  $\tau_{\text{samp}} = 1.0$  (other decoding knobs follow the DAEDAL defaults for each benchmark).

**Speed measurement.** We report wall-clock speedup under the same hardware and evaluation harness, measured relative to the corresponding LLaDA-8B-Instruct baseline configuration for each benchmark with a batch size of 1.

# 5 Results and Analysis

Next, we present our experimental results and their analysis.

## 5.1 Main Results

Across the full results in Tables 1, 2, and 3, our method improves the unweighted average score from 45.6 to 46.3 while reducing the mean NFE from 292.6 to 117.2, yielding a **3.62× end-to-end**

<sup>2</sup><https://github.com/Li-Jinsong/DAEDAL>

Generation	Model	Tokens/Step	Acc. (%)	Avg. NFE	Tokens/s (Speedup)
			↑	↓	↑
<b>Sequential</b> ( $b=1$ )	LLaDA	1	76.4	256	7.5 (1.00 $\times$ )
	LLaDA	1	77.4	256	7.5 (1.00 $\times$ )
<b>Block Diffusion</b> ( $b=32$ )	LLaDA	2	74.8	128	14.9 (1.99 $\times$ )
	Fast-dLLM	/	76.9	78.7	18.4 (2.46 $\times$ )
	<b>Ours</b>	1–32	<b>77.6</b>	<b>76.3</b>	<b>38.7 (5.18<math>\times</math>)</b>
<b>Pure Diffusion</b> ( $b=256$ )	LLaDA	1	13.8	256	7.5 (1.00 $\times$ )
	LLaDA	2	13.0	128	14.9 (1.99 $\times$ )
	<b>Ours</b>	1–32	54.7	80.6	36.7 (4.91 $\times$ )

Table 1: **GSM8K (zero-shot) accuracy-efficiency across generation paradigms.** We compare LLaDA-8B-Instruct under sequential, block diffusion ( $b = 32$ ), and pure diffusion ( $b = 256$ ) decoding against our method. Tokens/Step denotes the number of newly unmasked tokens per denoising step (fixed for LLaDA; adaptive for ours). Avg. NFE is the average number of function evaluations per sample. Tokens/s reports the achieved decoding throughput measured as finalized tokens per second under wall-clock time, and Speedup is computed relative to the sequential baseline. For Fast-dLLM, we report the official results corresponding to its parallel decoding strategy *without* the KV cache. All runs use a maximum sequence length of  $L_{\text{gen}} = 256$ .

Generation	Model	Tokens/Step	pass@1 (%)	pass@5 (%)	Avg. NFE	Tokens/s (Speedup)
			↑	↑	↓	↑
<b>Sequential</b> ( $b=1$ )	LLaDA	1	36.8	49.2	256	3.6 (1.00 $\times$ )
	LLaDA	1	36.9	51.4	256	3.6 (1.00 $\times$ )
<b>Block Diffusion</b> ( $b=32$ )	LLaDA	2	33.2	44.8	128	7.3 (2.00 $\times$ )
	<b>Ours</b>	1–32	<b>39.0</b>	<b>52.6</b>	<b>97.7</b>	<b>10.9 (2.96<math>\times</math>)</b>
	LLaDA	1	6.0	14.8	256	3.7 (1.00 $\times$ )
<b>Pure Diffusion</b> ( $b=256$ )	LLaDA	2	14.4	24.2	128	7.3 (2.00 $\times$ )
	<b>Ours</b>	1–32	36.4	51.0	99.4	10.8 (2.96 $\times$ )

Table 2: **MBPP (three-shot) pass@k and efficiency across generation paradigms (LLaDA-8B-Instruct).** We report pass@1/pass@5 under sequential, block diffusion ( $b = 32$ ), and pure diffusion ( $b = 256$ ) decoding. Avg. NFE, Tokens/s and Speedup are defined as in Table 1. All runs use a maximum sequence length of  $L_{\text{gen}} = 256$ .

**speedup.** Quantitatively, the gains translate into large wall-clock accelerations with comparable quality: On GSM8K, we match the fixed-step block-diffusion baseline (77.6% vs. 77.4%) while achieving a 5.18 $\times$  wall-clock speedup over the sequential baseline defined in Table 1; on HumanEval, we improve pass@1 by 2.2% with a 3.30 $\times$  speedup; and on MATH500, we improve accuracy by 1.3% with a 5.33 $\times$  speedup.

**Consistent efficiency gains.** Fig. 4 visualizes the accuracy-compute frontier across all four benchmarks. The curve of our CAD ( $\gamma_{\text{conf}} = 0.95$ ) consistently dominates the LLaDA baseline (dashed gray), shifting the frontier upward and leftward. This indicates that our method already attains strong performance in low-NFE regimes whereas the baseline remains compute-limited. This dominance holds across diverse domains (code vs. mathematics) without tuning the decoding budget per benchmark. This behavior is consistent with the end-to-end speedups reported in Tables 1, 2, and 3, where reducing NFE does not incur the accuracy degradation typically observed with heuristic step reduction.

**Pareto-optimal trade-off.** To verify that these gains are not artifacts of a specific parameter setting, Fig. 5 illustrates the full accuracy-efficiency frontier on GSM8K by sweeping the confidence threshold  $\gamma_{\text{conf}} \in [0.85, 0.99]$ . This analysis reveals the *controllability* and *robustness* of our approach:

- **Strict dominance:** As shown in Fig. 5, our curves (colored lines) consistently lie above the baseline trajectory (dashed gray). This implies a strict Pareto improvement: For any target accuracy, our method reduces NFE.
- **Flexible deployment:** The convex hull formed by our method enables users to seamlessly trade compute for quality. A lower threshold ( $\gamma_{\text{conf}} = 0.85$ , light blue line) offers aggressive speedups for latency-critical scenarios. Conversely, a conservative threshold ( $\gamma_{\text{conf}} = 0.99$ ,

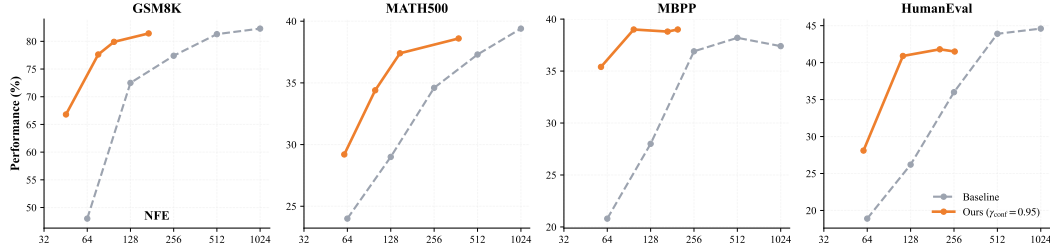


Figure 4: **Accuracy–compute trade-off under different decoding budgets.** We plot performance (%) versus the average NFE across benchmarks. The dashed gray curve denotes the LLaDA-8B-Instruct baseline evaluated under increasing NFE budgets, while the solid orange curve denotes our CAD with the given threshold ( $\gamma_{\text{conf}} = 0.95$ ). Each marker corresponds to one decoding budget. Across all four benchmarks, our method yields a consistently better trade-off, achieving comparable or higher performance at substantially lower NFE.

Task	Benchmark	LLaDA-8B-Ins	Ours	Speedup
Code	HumanEval (pass@1 %)	38.7 <sup>256</sup>	<b>40.9</b> <sup>113.2</sup>	3.30×
	HumanEval (pass@5 %)	51.2 <sup>256</sup>	<b>52.4</b> <sup>113.9</sup>	3.26×
	HumanEval-plus (pass@1 %)	31.7 <sup>256</sup>	<b>32.9</b> <sup>115.0</sup>	3.25×
	HumanEval-plus (pass@5 %)	42.7 <sup>256</sup>	<b>43.9</b> <sup>113.3</sup>	3.27×
	MBPP-plus (pass@1 %)	<b>48.7</b> <sup>256</sup>	47.9 <sup>108.0</sup>	3.51×
	MBPP-plus (pass@5 %)	<b>68.8</b> <sup>256</sup>	67.7 <sup>108.8</sup>	3.41×
Math	MATH500 (Acc %)	37.3 <sup>512</sup>	<b>38.6</b> <sup>148.3</sup>	5.33×
Avg.		45.6 <sup>292.6</sup>	<b>46.3</b> <sup>117.2</sup>	3.62×

Table 3: **Overall performance on code and mathematics benchmarks.** We report pass@1/5 on code benchmarks and accuracy on MATH500. The <sup>nfe</sup> next to each score denotes the Avg. NFE for that setting. Speedup is measured by wall-clock time relative to LLaDA-8B-Instruct under the same evaluation protocol. **Avg.** denotes the unweighted mean across the listed benchmarks.

purple line) prioritizes maximum quality. Our default setting ( $\gamma_{\text{conf}} = 0.95$ ) strikes a strong balance, requiring  $\approx 3.4\times$  smaller NFE to reach comparable accuracy (NFE ratio) and identifying the elbow of the accuracy–compute curve.

This confirms that the reported speedups are not artifacts of hyperparameter tuning, but a systemic advantage of confidence-adaptive allocation.

Furthermore, our ablation study presented in Appendix E confirms that the proposed DSCD objective is essential for this efficiency, as replacing it with standard supervised fine-tuning (SFT) leads to accuracy collapse under aggressive decoding steps.

## 5.2 Mechanism of Efficiency

To understand the source of the reported speedups, we present an analysis of the distribution of computational cost and its translation into wall-clock latency.

**Adaptive compute allocation.** Fig. 6 visualizes the NFE density. The distributions reveal that our method acts as a probe for intrinsic task complexity: GSM8K exhibits a sharp, low-variance peak ( $\mu = 76.3$ ), reflecting high model confidence, whereas MATH500 displays a broad, heavy-tailed distribution ( $\mu = 148.3$ ), adapting to diverse problem difficulties. Crucially, the distributions strictly deviate from the fixed baseline budget (gray lines), indicating that standard diffusion decoding is systematically over-parameterized. Our method successfully reclaims this redundancy, terminating well before the fixed limit even for the hardest “tail” samples.

**Translating step reduction to wall-clock speedup.** A key finding in Table 1 is that our end-to-end wall-clock speedup can exceed the reduction implied by the NFE ratio (e.g.,  $3.36\times$  smaller NFE

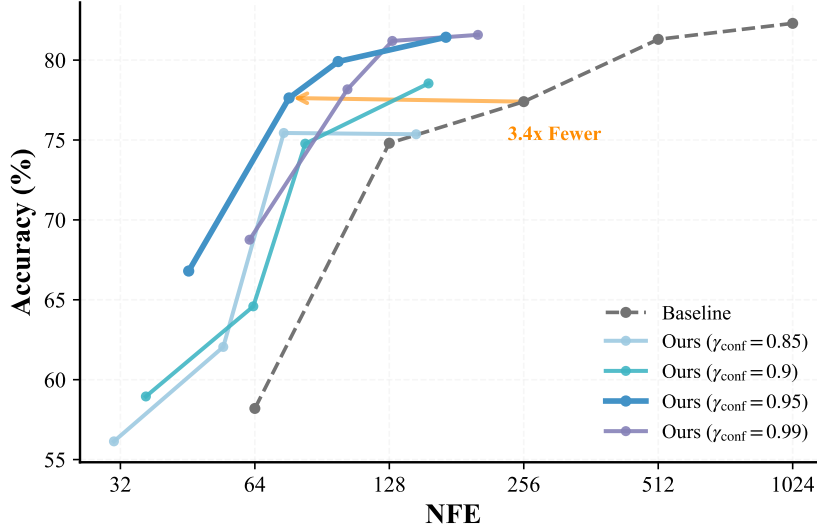


Figure 5: **Accuracy-compute Pareto frontier on GSM8K.** The dashed gray curve represents the LLaDA baseline. **Strict Dominance:** Our method (solid curves) consistently pushes the frontier upwards and leftwards. The orange arrow highlights our selected operating point ( $\gamma_{\text{conf}} = 0.95$ , dark blue curve), which achieves a **3.4 $\times$  speedup** while maintaining accuracy comparable to the baseline. While a higher threshold ( $\gamma_{\text{conf}} = 0.99$ ) prioritizes quality,  $\gamma_{\text{conf}} = 0.95$  strikes an optimal efficiency-quality balance.

vs.  $5.18\times$  faster on GSM8K). Since total latency  $T \approx \text{NFE} \cdot t_{\text{step}}$ , this implies that our gains come from both fewer denoising steps and a lower average per-step latency ( $t_{\text{step}}$ ), making each functional evaluation computationally cheaper. This behavior contrasts sharply with training-free baselines. For instance, Wu et al. [2025] report that their parallel decoding strategy achieves a  $3.25\times$  reduction in steps (tokens per step) but only yields a  $2.46\times$  wall-clock speedup, exhibiting sublinear scaling due to the overhead of inference-time selection heuristics. In contrast, our DSCD training aligns the model with the CAD acceptance rule, enabling a streamlined, fully batched tensor implementation with minimal dynamic control flow. Consequently, our algorithmic NFE reductions translate fully, and often superlinearly, into realized latency savings.

**Emergence of hierarchical planning.** Beyond quantitative speedups, the decoding trajectory (visualized in Appendix F.2 Fig. F.1) reveals that **CD<sup>4</sup>LM** learns a *hierarchical* generation strategy distinct from the linear left-to-right order of AR models. We observe a clear temporal separation between structure and logic:

- **Syntactic scaffolding (blue/green):** Structural tokens, including Python keywords (def, if, return) and control flow indentations, are consistently finalized in the earliest inference steps ( $t < 10$ ). This suggests the model performs global planning first, establishing a high-confidence syntactic skeleton to constrain the solution space.
- **Logical refinement (yellow/red):** Computationally intensive tokens, such as complex arithmetic expressions (e.g., `length = end - start + 1`) and conditional predicates, appear in warmer colors, indicating they are unmasked much later. This confirms that CAD effectively focuses the compute budget on the “hardest” parts of the sequence, utilizing the fully visible syntactic context to resolve logical dependencies with higher precision.

This behavior demonstrates that our method effectively decouples *global structural planning* from *local logical execution*, enabling the model to “sketch” the solution before filling in intricate details: a key factor driving its efficiency and correctness on structured tasks.

This structural stability is further validated by our qualitative evaluation presented in Appendix F.1, which demonstrates that our method significantly reduces token repetition and improves coherence scores compared to baselines employing aggressive fixed-step reduction.

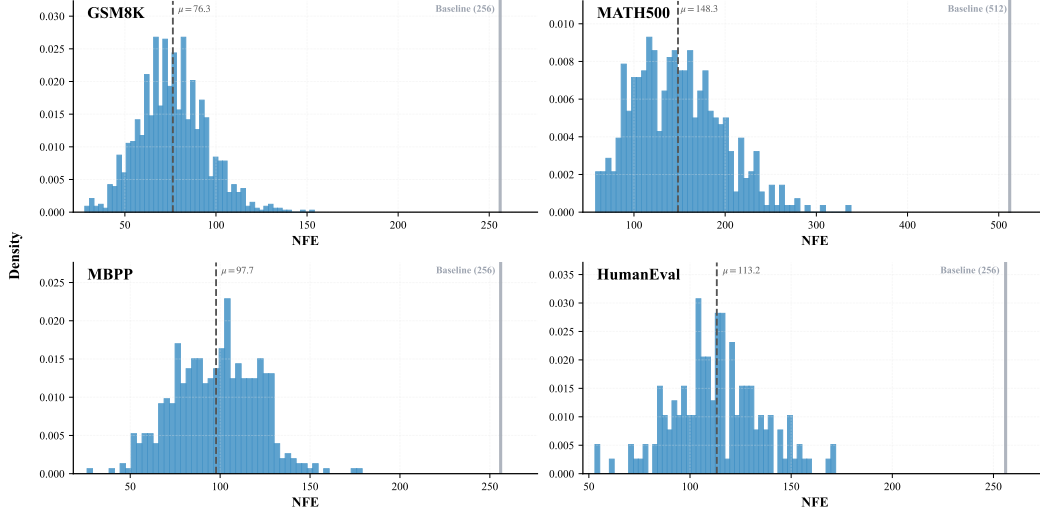


Figure 6: **Per-sample compute allocation under adaptive decoding.** Histograms show the NFE distribution across test instances for GSM8K, MATH500, MBPP, and HumanEval using our adaptive decoding approach. The dashed vertical line marks the mean NFE ( $\mu$ ), while the solid vertical line indicates the fixed NFE used by the baseline decoder (256 for GSM8K/MBPP/HumanEval and 512 for MATH500). Most samples terminate far earlier than the fixed-budget baseline, explaining the large reductions in average NFE reported in Tables 1, 2, and 3.

## 6 Discussions

Next, we discuss limitations, challenges, and future directions.

### 6.1 Limitations and Challenges

While  $\text{CD}^4\text{LM}$  establishes a new Pareto frontier for DLM decoding efficiency, several limitations remain inherent to our current design.

**Static canvas constraints.** First, like the underlying LLaDA backbone and most block-diffusion models, our approach relies on a pre-defined maximum sequence length ( $L_{\text{gen}}$ ). Although CAD logically handles variable-length outputs via EOS blocking, the computational graph is statically allocated (e.g., padding to 256 tokens). This introduces memory redundancy when generating short sequences and imposes a hard boundary on long-context reasoning, preventing the model from generalizing to sequences longer than its training window.

**Teacher-bounded reasoning.** Second, as a distillation framework, the student’s capability is theoretically bounded by that of the teacher. While DSCD effectively adapts the student to low-NFE trajectories, it does not fundamentally inject new reasoning capabilities. If the teacher hallucinates or fails in complex logic, the student may mimic those errors, though we occasionally observe minor self-correction effects typical of consistency training.

**Metric sensitivity in open-ended domains.** Third, our confidence-based acceptance relies on the assumption that low uncertainty correlates with correctness. This holds true for structured tasks like coding and mathematics (low-entropy targets) but may be overly conservative for high-entropy tasks, such as creative writing, where ambiguity is natural. Strict confidence thresholding in such domains might stifle diversity or lead to repetitive outputs.

## 6.2 Future Directions

The above limitations suggest promising avenues for future research to further democratize non-autoregressive generation.

**Dynamic and infinite-context diffusion.** To overcome the fixed canvas limitation, future work could integrate our adaptive decoding method with dynamic windowing mechanisms. Recent work by Li et al. [2025] proposes extensions of diffusion generation to arbitrary lengths via semi-autoregressive context shifting. Combining our confidence-adaptive logic with such dynamic frameworks could yield a fully flexible diffusion decoder that supports infinite-context generation without pre-allocated buffers.

**Beyond imitation via on-policy training.** To break the teacher performance ceiling, future research could explore on-policy refinement. Instead of purely mimicking a frozen teacher, the student could be fine-tuned via reinforcement learning using the efficiency benefits of CAD to explore diverse trajectories, potentially surpassing the teacher by optimizing for correctness rather than mere consistency.

**Integration with KV caching.** Although our current implementation recomputes the full context at each step, prior work like Fast-dLLM [Wu et al., 2025] has demonstrated the effectiveness of approximate KV caching for DLMs. Since our method is compatible with caching mechanisms, integrating them could further amortize the computational cost of the backbone. We hypothesize that combining our algorithmic NFE reduction with efficient memory management would yield even greater wall-clock speedups.

## 7 Conclusion

In this work, we introduced **CD<sup>4</sup>LM**, a unified framework that reconciles the structural mismatch between diffusion training schedules and inference latency requirements. By coupling DSCD with CAD, we successfully decoupled the model’s generation trajectory from rigid pre-defined schedules. Our extensive evaluation on GSM8K, HumanEval, and MBPP demonstrates that **CD<sup>4</sup>LM** achieves a strict Pareto improvement over standard diffusion baselines, delivering a  $3\times$ - $5\times$  speedup without sacrificing accuracy. These results confirm that treating diffusion models as flexible, instance-aware refinement operators, rather than fixed-schedule denoisers, is a viable path toward making non-autoregressive generation practical for real-world structured reasoning tasks.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Wasi Uddin Ahmad, Aleksander Ficek, Mehrzad Samadi, Jocelyn Huang, Vahid Noroozi, Somshubra Majumdar, and Boris Ginsburg. OpenCodeInstruct: A large-scale instruction tuning dataset for code LLMs. *arXiv preprint arXiv:2504.04030*, 2025.
- Masaki Asada and Makoto Miwa. Addressing the training-inference discrepancy in discrete diffusion for text generation. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 7156–7164, 2025.
- Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021a.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021b.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

- Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D. Lee, Deming Chen, and Tri Dao. Medusa: Simple LLM inference acceleration framework with multiple decoding heads. *arXiv preprint arXiv:2401.10774*, 2024.
- George Casella and Christian P Robert. Rao-Blackwellisation of sampling schemes. *Biometrika*, 83(1):81–94, 1996.
- Mark Chen. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Tianqi Chen, Shujian Zhang, and Mingyuan Zhou. DLM-One: Diffusion language models for one-step sequence generation. *arXiv preprint arXiv:2506.00290*, 2025.
- Jacob K. Christopher, Brian R. Bartoldson, Tal Ben-Nun, Michael Cardei, Bhavya Kailkhura, and Ferdinando Fioretto. Speculative diffusion decoding: Accelerating language generation through diffusion. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 12042–12059, 2025.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and LingPeng Kong. DiffuSeq: Sequence to sequence text generation with diffusion models. *arXiv preprint arXiv:2210.08933*, 2022.
- Shansan Gong, Ruixiang Zhang, Huangjie Zheng, Jiatao Gu, Navdeep Jaitly, Lingpeng Kong, and Yizhe Zhang. DiffuCoder: Understanding and improving masked diffusion models for code generation. *arXiv preprint arXiv:2506.20639*, 2025.
- 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.
- Ishaan Gulrajani and Tatsunori B. Hashimoto. Likelihood-based diffusion language models. *Advances in Neural Information Processing Systems*, 36:16693–16715, 2023.
- Haoyu He, Katrin Renz, Yong Cao, and Andreas Geiger. MDPO: Overcoming the training-inference divide of masked diffusion language models. *arXiv preprint arXiv:2508.13148*, 2025.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- Zemin Huang, Zhiyang Chen, Zijun Wang, Tiancheng Li, and Guo-Jun Qi. Reinforcing the diffusion chain of lateral thought with diffusion language models. *arXiv preprint arXiv:2505.10446*, 2025.
- Daniel Israel, Guy Van den Broeck, and Aditya Grover. Accelerating diffusion LLMs via adaptive parallel decoding. *arXiv preprint arXiv:2506.00413*, 2025.
- Inception Labs, Samar Khanna, Siddhant Kharbanda, Shufan Li, Harshit Varma, Eric Wang, Sawyer Birnbaum, Ziyang Luo, Yanis Miraoui, Akash Palrecha, et al. Mercury: Ultra-fast language models based on diffusion. *arXiv preprint arXiv:2506.17298*, 2025.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR, 2023.



- Jinsong Li, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Jiaqi Wang, and Dahua Lin. Beyond fixed: Training-free variable-length denoising for diffusion large language models. *arXiv preprint arXiv:2508.00819*, 2025.
- Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B. Hashimoto. Diffusion-LM improves controllable text generation. *arXiv preprint arXiv:2205.14217*, 2022.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. Eagle: Speculative sampling requires rethinking feature uncertainty. *arXiv preprint arXiv:2401.15077*, 2024.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- 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*, 2025.
- 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. *arXiv preprint arXiv:2502.09992*, 2025.
- 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.
- 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.
- Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- 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*, 2025.
- Jiacheng Ye, Jiahui Gao, Shansan Gong, Lin Zheng, Xin Jiang, Zhenguo Li, and Lingpeng Kong. Beyond autoregression: Discrete diffusion for complex reasoning and planning. *arXiv preprint arXiv:2410.14157*, 2024.
- Runpeng Yu, Xinyin Ma, and Xinchao Wang. Dimple: Discrete diffusion multimodal large language model with parallel decoding. *arXiv preprint arXiv:2505.16990*, 2025.
- Huaisheng Zhu, Zhengyu Chen, Shijie Zhou, Zhihui Xie, Yige Yuan, Zhimeng Guo, Siyuan Xu, Hangfan Zhang, Vasant Honavar, and Teng Xiao. Simple denoising diffusion language models. *arXiv preprint arXiv:2510.22926*, 2025.

## A A Martingale Projection View of DSCD

In this section, we provide a formal derivation that connects our DSCD method to the Martingale theory on filtrations [Vincent, 2011].

### A.1 Setup and Information Ordering

Let  $(\Omega, \mathcal{F}, P)$  be the probability space that supports all random variables in our construction. Let  $Z \in \mathcal{Z}$  denote the clean (ground-truth) token sequence. The absorbing masking process of LLaDA defines, for each masking pattern  $m$ , a partially masked sequence  $X^{(m)}$  obtained by replacing a subset of token positions in  $Z$  with the mask token  $M$ .

To simplify notation in this theoretical analysis, we denote the random variables corresponding to the student view  $\tilde{z}^S$  and teacher view  $\tilde{z}^T$  (defined in Sec. 3.2) simply as  $S$  and  $T$ . We define:

$$S = X^{(m_S)}, \quad T = X^{(m_T)}.$$

We write

$$\mathcal{F}_S = \sigma(S), \quad \mathcal{F}_T = \sigma(T)$$

for the  $\sigma$ -algebras generated by the corresponding observed (unmasked) tokens.

In our nested masking design, the teacher always sees a superset of the tokens visible to the student. Concretely, if  $\mathcal{M}_S$  and  $\mathcal{M}_T$  denote the random sets of *masked* positions for the student and teacher, we enforce

$$\mathcal{M}_T \subseteq \mathcal{M}_S,$$

so that every token unmasked for the student is also unmasked for the teacher. Equivalently, the teacher’s observation  $T$  is a measurable function of  $(S, Z)$ , and the associated  $\sigma$ -algebras are ordered as  $\mathcal{F}_S \subseteq \mathcal{F}_T$ .

**Definition 1** (Information ordering). *We say that the teacher is informationally richer than the student, and write  $S \preceq T$ , if*

$$\mathcal{F}_S \subseteq \mathcal{F}_T,$$

*i.e., every event that is measurable with respect to the student state  $S$  is also measurable with respect to the teacher state  $T$ . Under our nested masking scheme  $\mathcal{M}_T \subseteq \mathcal{M}_S$ , we have  $S \preceq T$  by construction.*

### A.2 Consistency as Projection onto Student Information

We now formalize how DSCD can be viewed as projecting the teacher’s conditional distribution onto the information available to the student. Let  $Z \in \mathcal{Z}$  denote the clean (ground-truth) sequence, and let  $S$  and  $T$  be the student and teacher states as in Definition 1, with  $\mathcal{F}_S = \sigma(S)$  and  $\mathcal{F}_T = \sigma(T)$  satisfying  $S \preceq T$ .

Given a student state  $S = s$ , the DSCD objective compares the student distribution  $p_\theta(\cdot | s)$  against the teacher distribution  $p_\phi(\cdot | T)$  averaged over the teacher states  $T$  that are compatible with  $S = s$ . Formally, the DSCD consistency loss can be written as

$$\mathcal{L}_{\text{DSCD}}(\theta) = \mathbb{E}_S \mathbb{E}_{T|S} [D_{\text{KL}}(p_\phi(\cdot | T) \| p_\theta(\cdot | S))], \quad (11)$$

where the outer expectation is taken over the student states induced by the diffusion process and masking policy, and the inner expectation is taken over the corresponding teacher states. For later use, we define the teacher-induced *projected* distribution

$$Q(z | S = s) := \mathbb{E}_{T|S=s} [p_\phi(z | T)], \quad (12)$$

i.e.,  $Q(\cdot | s)$  is obtained by averaging the teacher’s conditional distributions over all teacher states  $T$  consistent with the student state  $S = s$ . Equivalently,  $Q(\cdot | S)$  is the conditional expectation of  $p_\phi(\cdot | T)$  with respect to the  $\sigma$ -algebra  $\mathcal{F}_S$ .

**Proposition 1** (DSCD as projection onto  $\mathcal{F}_S$ ). *Assume the student model class  $\{p_\theta(\cdot | S)\}$  is rich enough to represent any conditional distribution over  $\mathcal{Z}$  for each  $S = s$ . Then the DSCD loss (11) is minimized by*

$$p_\theta^*(z | S = s) = Q(z | S = s) = \mathbb{E}_{T|S=s} [p_\phi(z | T)], \quad \forall s \in \mathcal{S}. \quad (13)$$

In particular,  $p_\theta^*(\cdot | S)$  coincides with the conditional expectation of the teacher’s conditional distribution  $p_\phi(\cdot | T)$  onto the coarser  $\sigma$ -algebra  $\mathcal{F}_S$ .

*Proof.* By expanding the KL divergence in (11), we obtain

$$\mathcal{L}_{\text{DSCD}}(\theta) = \mathbb{E}_{S,T}[D_{\text{KL}}(p_\phi(\cdot | T) \| p_\theta(\cdot | S))] \quad (14)$$

$$= \mathbb{E}_{S,T} \left[ \sum_{z \in \mathcal{Z}} p_\phi(z | T) \log \frac{p_\phi(z | T)}{p_\theta(z | S)} \right] \quad (15)$$

$$= \underbrace{\mathbb{E}_{S,T} \left[ \sum_{z \in \mathcal{Z}} p_\phi(z | T) \log p_\phi(z | T) \right]}_{\text{constant w.r.t. } \theta} - \mathbb{E}_{S,T} \left[ \sum_{z \in \mathcal{Z}} p_\phi(z | T) \log p_\theta(z | S) \right]. \quad (16)$$

The first term does not depend on  $\theta$  and can be dropped. Thus, minimizing (11) is equivalent to maximizing

$$\mathbb{E}_{S,T} \left[ \sum_{z \in \mathcal{Z}} p_\phi(z | T) \log p_\theta(z | S) \right] = \mathbb{E}_S \left[ \sum_{z \in \mathcal{Z}} \mathbb{E}_{T|S}[p_\phi(z | T)] \log p_\theta(z | S) \right]. \quad (17)$$

For each fixed student state  $S = s$ , the inner objective becomes

$$\sum_{z \in \mathcal{Z}} Q(z | s) \log p_\theta(z | s), \quad (18)$$

which is maximized when  $p_\theta(\cdot | s) = Q(\cdot | s)$ . Therefore,  $p_\theta^*(\cdot | S) = Q(\cdot | S)$ .  $\square$

This proposition shows that DSCD does not force the student to reproduce the teacher’s full diffusion trajectory. Instead, the student learns the teacher’s conditional distribution after integrating out the extra information available only to the teacher, i.e., after projecting onto the coarser information  $\sigma$ -algebra  $\mathcal{F}_S$ . This implies that the student learns to directly predict the *expectation* of the teacher’s multi-step denoising result, effectively compressing the diffusion trajectory into a single forward pass.

## B Variance Reduction Analysis of Teacher-Subset Masking

In this section, we provide a formal proof that the *teacher-subset masking* strategy reduces the variance of the consistency loss compared to an *independent masking* strategy.

### B.1 Problem Setup

Let  $z$  be the ground truth sequence. The student observes a partial view  $\tilde{z}^S$  with mask  $\mathcal{M}_S$ . The student’s goal is to minimize the divergence from a target distribution provided by the teacher. Let the teacher’s view be  $\tilde{z}^T$  with mask  $\mathcal{M}_T$ . The teacher’s prediction is a random variable  $Y = p_\phi(\cdot | \tilde{z}^T)$ , which serves as the regression target for the student.

We compare two masking schemes:

- **Teacher-subset masking:**  $\mathcal{M}_T \subseteq \mathcal{M}_S$ . The teacher observes everything the student observes, plus additional tokens  $\Delta = \mathcal{M}_S \setminus \mathcal{M}_T$ .
- **Independent masking:**  $\mathcal{M}_T$  is sampled independently. Crucially, there exists a non-empty set of “blind” tokens  $\mathcal{B} = \mathcal{M}_T \cap (\mathcal{V} \setminus \mathcal{M}_S)$  that are *visible to the student* but *masked for the teacher*.

### B.2 Variance Reduction via Rao-Blackwellization

We now make the variance reduction claim precise. Let  $U$  be a random target quantity that the teacher aims to approximate (e.g., a one-hot encoding of the ground-truth token or a sufficient statistic

thereof). For any information  $\sigma$ -algebra  $\mathcal{G}$ , the Bayes-optimal predictor based on  $\mathcal{G}$  is the conditional expectation

$$Y^{(\mathcal{G})} = \mathbb{E}[U \mid \mathcal{G}].$$

The student observes a state  $\tilde{z}^S$  and the corresponding  $\sigma$ -algebra  $\mathcal{F}_S = \sigma(\tilde{z}^S)$ . We consider two teacher information sets:

- **Teacher-subset masking:** The teacher observes all information available to the student plus additional tokens. Denote the teacher's state by  $\tilde{z}^{T,\text{sub}}$  and the associated  $\sigma$ -algebra by  $\mathcal{F}_T^{\text{sub}}$ , with  $\mathcal{F}_S \subseteq \mathcal{F}_T^{\text{sub}}$ . The Bayes-optimal teacher prediction is

$$Y^{\text{sub}} = \mathbb{E}[U \mid \mathcal{F}_T^{\text{sub}}].$$

- **Independent masking:** The teacher's mask is sampled independently of the student's mask; hence, the teacher state  $\tilde{z}^{T,\text{ind}}$  may hide a subset of tokens that are visible to the student. Let  $\mathcal{F}_T^{\text{ind}} = \sigma(\tilde{z}^{T,\text{ind}})$  be the corresponding information set. The Bayes-optimal independent teacher prediction is

$$Y^{\text{ind}} = \mathbb{E}[U \mid \mathcal{F}_T^{\text{ind}}].$$

While the teacher model  $p_\phi$  is not the true Bayes-optimal posterior of the data, it serves as the ground-truth definition for the distillation task. Under this view, subset masking reduces the variance of the gradient estimator relative to the distillation target.

In the teacher-subset scheme, the teacher's information strictly refines the independent teacher's information once we condition on the student's view:

$$\sigma(\mathcal{F}_S, \mathcal{F}_T^{\text{ind}}) \subseteq \mathcal{F}_T^{\text{sub}}.$$

The following lemma is a conditional version of the Rao-Blackwell theorem.

**Lemma 1** (Conditional Rao-Blackwellization). *Let  $U \in L^2(\Omega, \mathcal{F}, P)$  and let  $\mathcal{G} \subseteq \mathcal{H} \subseteq \mathcal{F}$  be  $\sigma$ -algebras. Define*

$$Y_{\mathcal{G}} = \mathbb{E}[U \mid \mathcal{G}], \quad Y_{\mathcal{H}} = \mathbb{E}[U \mid \mathcal{H}].$$

*Then for any  $\sigma$ -algebra  $\mathcal{K} \subseteq \mathcal{G}$ ,*

$$\text{Var}(Y_{\mathcal{H}} \mid \mathcal{K}) \leq \text{Var}(Y_{\mathcal{G}} \mid \mathcal{K}) \quad \text{almost surely (a.s.)}$$

*with equality iff  $Y_{\mathcal{G}} = Y_{\mathcal{H}}$ .*

*Proof.* Fix  $\mathcal{K} \subseteq \mathcal{G}$ . Using the tower property and orthogonality of conditional expectations in  $L^2$ , we can write

$$Y_{\mathcal{G}} - Y_{\mathcal{H}} = (Y_{\mathcal{G}} - \mathbb{E}[Y_{\mathcal{H}} \mid \mathcal{G}]) + (\mathbb{E}[Y_{\mathcal{H}} \mid \mathcal{G}] - Y_{\mathcal{H}}).$$

Since  $\mathcal{G} \subseteq \mathcal{H}$ , we have  $\mathbb{E}[Y_{\mathcal{H}} \mid \mathcal{G}] = Y_{\mathcal{G}}$ ; hence, the first term is zero and

$$Y_{\mathcal{G}} - Y_{\mathcal{H}} = Y_{\mathcal{G}} - \mathbb{E}[U \mid \mathcal{H}].$$

Taking conditional expectations with respect to  $\mathcal{K}$  and using  $\mathcal{K} \subseteq \mathcal{G}$  gives

$$\text{Var}(Y_{\mathcal{G}} \mid \mathcal{K}) = \text{Var}(Y_{\mathcal{H}} \mid \mathcal{K}) + \mathbb{E}[(Y_{\mathcal{G}} - Y_{\mathcal{H}})^2 \mid \mathcal{K}].$$

The second term is nonnegative and vanishes iff  $Y_{\mathcal{G}} = Y_{\mathcal{H}}$  a.s. □

We now instantiate Lemma 1 with

$$\mathcal{K} = \mathcal{F}_S, \quad \mathcal{G} = \mathcal{F}_T^{\text{ind}}, \quad \mathcal{H} = \mathcal{F}_T^{\text{sub}}.$$

By construction,  $\mathcal{K} \subseteq \mathcal{G} \subseteq \mathcal{H}$  and

$$Y^{\text{ind}} = Y_{\mathcal{G}}, \quad Y^{\text{sub}} = Y_{\mathcal{H}}.$$

Therefore, Lemma 1 implies

$$\text{Var}(Y^{\text{sub}} \mid \mathcal{F}_S) \leq \text{Var}(Y^{\text{ind}} \mid \mathcal{F}_S) \quad \text{a.s.} \tag{19}$$

whenever the teacher is Bayes-optimal with respect to its information set. Moreover, the inequality is strict as soon as  $\mathcal{F}_T^{\text{ind}} \subsetneq \mathcal{F}_T^{\text{sub}}$ .

In words, teacher-subset masking *Rao-Blackwellizes* the distillation target by conditioning on the student's full observation and the additional tokens revealed by the nested mask, thereby reducing the conditional variance of the teacher signal seen by the student.

## C Theoretical Analysis of Confidence-Adaptive Decoding (CAD)

This appendix provides (i) a decision-theoretic motivation for CAD as a greedy solver of a per-step risk-efficiency trade-off with trust-region and liveness constraints, and (ii) basic guarantees on termination, step bounds, and computational overhead.

This derivation shows that confidence thresholding is optimal for an additive risk-efficiency objective, and that under cardinality constraints, the optimal solution is a confidence-sorted prefix. CAD implements a greedy approximation with explicit trust-region and liveness constraints.

### C.1 Notation and Eligible Set

Let the target length be  $L_{\text{gen}}$ . At decoding step  $s$ , the current partially-masked state is  $\tilde{z}^{(s)} \in \mathcal{V}^{L_{\text{gen}}}$ , and the masked index set is  $\mathcal{M}^{(s)} = \{i \in [L_{\text{gen}}] : \tilde{z}_i^{(s)} = \mathfrak{m}\}$ . For block diffusion with block size  $b$ , define blocks

$$\mathcal{B}_j = \{(j-1)b + 1, \dots, \min(jb, L_{\text{gen}})\}, \quad j = 1, \dots, J, \quad J = \lceil L_{\text{gen}}/b \rceil,$$

and maintain an active block index  $j$  (decoded left-to-right). CAD operates on the *eligible set*

$$\mathcal{E}^{(s)} = \mathcal{M}^{(s)} \cap \mathcal{B}_j, \quad (20)$$

i.e., masked positions within the active block. Pure diffusion is recovered as a special case by setting  $b = L_{\text{gen}}$  (a single block), in which case  $\mathcal{E}^{(s)} = \mathcal{M}^{(s)}$ .

### C.2 Decision-theoretic Derivation (Per-step Objective)

At step  $s$ , the model produces a confidence score for each eligible position  $i \in \mathcal{E}^{(s)}$ ,  $c_i^{(s)} \in [0, 1]$ . We interpret  $c_i^{(s)}$  as a proxy for the probability that the current argmax token at position  $i$  is correct:

$$c_i^{(s)} \approx \mathbb{P}(\text{correct}_i \mid \tilde{z}^{(s)}). \quad (21)$$

Let  $\mathcal{V}^{(s)} \subseteq \mathcal{E}^{(s)}$  denote the *valid* candidates at step  $s$  (e.g., after excluding padding or other structurally invalid positions; by default  $\mathcal{V}^{(s)} = \mathcal{E}^{(s)}$ ). The decision variable is a commit set  $\mathcal{S} \subseteq \mathcal{V}^{(s)}$ .

We define two opposing per-step objectives: (i) **risk**: the expected number of incorrect tokens committed at this step,  $\mathcal{R}(\mathcal{S}) = \sum_{i \in \mathcal{S}} (1 - c_i^{(s)})$ , and (ii) **efficiency**: the number of committed tokens,  $\mathcal{E}(\mathcal{S}) = |\mathcal{S}|$ . We combine them into a net cost

$$\mathcal{L}(\mathcal{S}) = \mathcal{R}(\mathcal{S}) - \lambda_{\text{trade}} \mathcal{E}(\mathcal{S}) = \sum_{i \in \mathcal{S}} (1 - c_i^{(s)} - \lambda_{\text{trade}}), \quad (22)$$

where  $\lambda_{\text{trade}} > 0$  controls the risk-speed trade-off and is *unrelated* to the training loss-mixing schedule  $\lambda(g)$ . Because (22) is additive across positions, the unconstrained minimizer commits all positions whose per-token contribution is negative:

$$i \in \mathcal{S}^* \iff 1 - c_i^{(s)} - \lambda_{\text{trade}} < 0 \iff c_i^{(s)} > 1 - \lambda_{\text{trade}}. \quad (23)$$

Defining the confidence threshold  $\gamma_{\text{conf}} = 1 - \lambda_{\text{trade}}$  yields the basic *confidence thresholding* rule.

### C.3 Trust-region and Liveness Constraints, and the CAD Rule

In practice, we impose two cardinality constraints: *trust region* ( $|\mathcal{S}| \leq k_{\text{max}}$ ) to prevent committing too many tokens based on local confidence estimates, and *liveness* ( $|\mathcal{S}| \geq k_{\text{min}}$ ) to guarantee progress whenever possible. This yields the constrained problem

$$\min_{\mathcal{S} \subseteq \mathcal{V}^{(s)}} \mathcal{L}(\mathcal{S}) \quad \text{s.t.} \quad k_{\text{min}} \leq |\mathcal{S}| \leq k_{\text{max}}. \quad (24)$$

Sorting candidates by decreasing confidence produces an ordering  $i_{(1)}, \dots, i_{(|\mathcal{V}^{(s)}|)}$  with  $c_{i_{(1)}}^{(s)} \geq \dots \geq c_{i_{(|\mathcal{V}^{(s)}|)}}^{(s)}$ . Under (24), the optimum is always a prefix of this ordering. CAD adopts a simple greedy approximation: Count how many positions within the trust region exceed  $\gamma_{\text{conf}}$ ,

$$n_{\text{conf}} = \left| \left\{ j \leq \min(k_{\text{max}}, |\mathcal{V}^{(s)}|) : c_{i_{(j)}}^{(s)} \geq \gamma_{\text{conf}} \right\} \right|, \quad (25)$$

then clamp the committed token count and select the top- $k^{(s)}$  positions:

$$k^{(s)} = \min \left( \max(k_{\min}, n_{\text{conf}}), \min(k_{\max}, |\mathcal{V}^{(s)}|) \right), \quad \mathcal{S}^{(s)} = \{i_{(1)}, \dots, i_{(k^{(s)})}\}. \quad (26)$$

Equations (25)-(26) match the CAD update implemented in the main text.

#### C.4 Monotonic Progress and Step Bound

The following statements formalize why CAD terminates and yield a simple step bound under block diffusion.

**Lemma 2** (Monotonicity within the active block). *Assume  $\mathcal{E}^{(s)} \neq \emptyset$  and  $k_{\min} \geq 1$ . After one CAD update,*

$$|\mathcal{M}^{(s+1)} \cap \mathcal{B}_j| \leq |\mathcal{M}^{(s)} \cap \mathcal{B}_j| - 1.$$

*Moreover, if  $|\mathcal{E}^{(s)}| \geq k_{\min}$ , then*

$$|\mathcal{M}^{(s+1)} \cap \mathcal{B}_j| \leq |\mathcal{M}^{(s)} \cap \mathcal{B}_j| - k_{\min}.$$

*Proof.* By construction,  $\mathcal{S}^{(s)} \subseteq \mathcal{E}^{(s)} = \mathcal{M}^{(s)} \cap \mathcal{B}_j$  and each  $i \in \mathcal{S}^{(s)}$  is committed from  $\langle \text{mask} \rangle$  to a concrete token, hence removed from  $\mathcal{M}^{(s)}$ . If  $\mathcal{E}^{(s)} \neq \emptyset$ , the clamp in (26) ensures  $k^{(s)} \geq 1$ , yielding the first inequality. When  $|\mathcal{E}^{(s)}| \geq k_{\min}$ , we have  $k^{(s)} \geq k_{\min}$ , yielding the second inequality.  $\square$

**Corollary 1** (Per-block and total NFE bound). *For block size  $b$ , each block  $\mathcal{B}_j$  is resolved within at most  $\lceil |\mathcal{B}_j|/k_{\min} \rceil$  NFE. Therefore, the total NFE satisfies*

$$S \leq \sum_{j=1}^J \left\lceil \frac{|\mathcal{B}_j|}{k_{\min}} \right\rceil \leq \left\lceil \frac{L_{\text{gen}}}{k_{\min}} \right\rceil + J \leq \left\lceil \frac{L_{\text{gen}}}{k_{\min}} \right\rceil + \left\lceil \frac{L_{\text{gen}}}{b} \right\rceil,$$

*and is additionally capped by the implementation limit  $S_{\max}$ .*

#### C.5 Special Cases and Relation to Fixed-step Decoding

**Proposition 2** (Pure diffusion as a special case). *Setting  $b = L_{\text{gen}}$  (i.e.,  $J = 1$ ) reduces block diffusion with CAD to pure diffusion with CAD, because  $\mathcal{E}^{(s)} = \mathcal{M}^{(s)}$  for all  $s$ .*

**Proposition 3** (Reduction to fixed-step block diffusion). *If  $\gamma_{\text{conf}} = -\infty$  and  $k_{\min} = k_{\max} = k$  (a constant), then CAD commits exactly  $k$  positions per NFE within the active block (until fewer than  $k$  masks remain), matching a fixed-step block diffusion schedule up to the last step of each block.*

#### C.6 Complexity and Overhead

A natural concern for CAD is the computational cost of the control policy itself (e.g., selecting high-confidence positions and updating the active set), especially at large batch sizes. We denote by  $B$  the batch size,  $L$  the total sequence length (prompt + target),  $L_{\text{gen}}$  the maximum generation budget,  $D$  the hidden dimension, and  $N_{\text{layer}}$  the number of Transformer layers. Ignoring constant factors, a single forward pass of the student model has time complexity

$$\text{Cost}_{\text{model}} = \mathcal{O}(B N_{\text{layer}}(L^2 D + L D^2)), \quad (27)$$

dominated by self-attention and feed-forward blocks over the full context.

In contrast, the CAD controller operates only on token-wise confidence scores within the eligible set. At decoding step  $s$ , it performs thresholding and/or a top- $k$  operation over at most  $|\mathcal{E}^{(s)}|$  elements, costing

$$\text{Cost}_{\text{CAD}}^{(s)} = \mathcal{O}(B |\mathcal{E}^{(s)}| \log |\mathcal{E}^{(s)}|). \quad (28)$$

Under block diffusion with block size  $b$ , we have  $|\mathcal{E}^{(s)}| \leq b$  for all  $s$ , hence

$$\text{Cost}_{\text{CAD}} = \sum_{s=1}^S \text{Cost}_{\text{CAD}}^{(s)} = \mathcal{O}(BS b \log b) \leq \mathcal{O}(BS L_{\text{gen}} \log L_{\text{gen}}), \quad (29)$$

where the last inequality follows from  $b \leq L_{\text{gen}}$  (pure diffusion corresponds to  $b = L_{\text{gen}}$ ). Comparing the two, the relative overhead satisfies:

$$\frac{\text{Cost}_{\text{CAD}}}{\text{Cost}_{\text{model}}} \lesssim \frac{BSb \log b}{B N_{\text{layer}} S (L^2 D + LD^2)} \approx \mathcal{O}\left(\frac{b \log b}{N_{\text{layer}} L^2 D}\right). \quad (30)$$

Since  $L$  represents the full sequence length, we observe two asymptotic regimes: (1) For block diffusion ( $b$  fixed,  $b \ll L$ ), the overhead decays quadratically with sequence length ( $\propto L^{-2}$ ), rendering it negligible for long-context generation. (2) Even for pure diffusion where  $b \approx L_{\text{gen}} < L$ , the ratio scales as  $\mathcal{O}(L^{-1})$ , which still vanishes asymptotically as model depth and width increase.

## D Method Details

### D.1 Distillation Algorithm Pseudocode

Algorithm 1 describes the full training procedure for DSCD. The algorithm alternates between (i) constructing paired student-teacher views via teacher-subset masking, (ii) computing the curriculum-weighted loss, and (iii) updating the student parameters.

---

#### Algorithm 1 Discrete-Space Consistency Distillation (DSCD)

---

**Require:** Teacher model  $p_\phi$  (frozen), student model  $p_\theta$ , dataset  $\mathcal{D}$

**Require:** Mask ratio ranges:  $[r_S^{\min}, r_S^{\max}]$ ,  $[r_T^{\min}, r_T^{\max}]$

**Require:** Distillation temperature  $\tau$ , curriculum parameters  $\lambda_0, \lambda_1, \alpha$

**Require:** Maximum training steps  $G_{\max}$

```

1: for  $g = 1$  to  $G_{\max}$  do
2:   Sample mini-batch  $\{(x^{(b)}, y^{(b)})\}_{b=1}^B \sim \mathcal{D}$ 
3:   for each sample  $(x, y)$  in batch do
4:      $z \leftarrow (x; y)$  ▷ Concatenate prompt and target
5:      $L_x \leftarrow |x|, L_y \leftarrow |y|, L \leftarrow L_x + L_y$ 
6:     // — Step 1: Sample mask ratios (see Eq. 31) —
7:      $r_S \sim \text{Uniform}(r_S^{\min}, r_S^{\max})$ 
8:      $u \sim \text{Uniform}(0.3, 0.7)$ 
9:      $r_T \leftarrow \text{clip}(r_S \cdot u, r_T^{\min}, r_T^{\max})$ 
10:    // — Step 2: Construct student mask —
11:     $\mathcal{U}_{\text{valid}} \leftarrow \{L_x + 1, \dots, L\} \setminus \mathcal{P}$ 
12:     $\mathcal{M}_S \leftarrow \text{RandomSample}(\mathcal{U}_{\text{valid}}, n_S)$  ▷ Uniformly sample positions
13:     $\tilde{z}^S \leftarrow \text{Mask}(z, \mathcal{M}_S)$  ▷ Replace  $z_i$  with  $\mathbf{m}$  for  $i \in \mathcal{M}_S$ 
14:    // — Step 3: Teacher-subset masking —
15:     $n_T \leftarrow \min(\lfloor L_y \cdot r_T \rfloor, n_S)$ 
16:     $\mathcal{M}_T \leftarrow \text{RandomSubset}(\mathcal{M}_S, n_T)$  ▷ Strict subset:  $\mathcal{M}_T \subseteq \mathcal{M}_S$ 
17:     $\tilde{z}^T \leftarrow \text{Mask}(z, \mathcal{M}_T)$ 
18:    // — Step 4: Forward pass —
19:     $\ell_\theta \leftarrow p_\theta(\cdot \mid \tilde{z}^S)$  ▷ Student logits
20:     $\ell_\phi \leftarrow p_\phi(\cdot \mid \tilde{z}^T)$  ▷ Teacher logits (no grad)
21:    // — Step 5: Compute losses on  $\mathcal{M}_S$  —
22:     $\mathcal{L}_{\text{recon}} \leftarrow -\frac{1}{|\mathcal{M}_S|} \sum_{i \in \mathcal{M}_S} \log p_\theta(z_i \mid \tilde{z}^S)$ 
23:     $\tilde{p}_\phi \leftarrow \text{softmax}(\ell_\phi / \tau)$  ▷ Temperature-scaled teacher
24:     $\mathcal{L}_{\text{cons}} \leftarrow \frac{\tau^2}{|\mathcal{M}_S|} \sum_{i \in \mathcal{M}_S} \text{KL}(\tilde{p}_{\phi, i} \parallel p_{\theta, i})$ 
25:  end for
26:  // — Step 6: Curriculum-weighted total loss —
27:   $\lambda(g) \leftarrow \text{CosineSchedule}(g/G_{\max}; \lambda_0, \lambda_1, \alpha)$  ▷ Eq. (32)
28:   $\mathcal{L}_{\text{total}} \leftarrow \lambda(g) \cdot \mathcal{L}_{\text{cons}} + (1 - \lambda(g)) \cdot \mathcal{L}_{\text{recon}}$ 
29:  // — Step 7: Parameter update —
30:   $\theta \leftarrow \theta - \eta \cdot \nabla_\theta \mathcal{L}_{\text{total}}|_{\text{batch}}$ 
31: end for
32: return Trained student  $p_\theta$ 

```

---

### Key design choices.

1. **Teacher-subset masking** (lines 14-16): By enforcing  $\mathcal{M}_T \subseteq \mathcal{M}_S$ , the teacher always observes a superset of tokens visible to the student. eqThis Rao-Blackwellizes the distillation target (Appendix B).
2. **Multiplicative ratio coupling** (lines 7-9): The teacher ratio is derived as  $r_T = r_S \cdot u$  with  $u \sim \mathcal{U}(0.3, 0.7)$ , ensuring  $r_T < r_S$  with high probability while maintaining stochastic diversity.
3. **Minimum mask count** (line 11): Short answers ( $L_y < 20$ ) use reduced mask ratios to prevent degenerate training signals.

**Mask ratio sampling.** The mask ratios are sampled according to:

$$\begin{aligned} r_S &\sim \mathcal{U}(r_S^{\min}, r_S^{\max}) = \mathcal{U}(0.4, 0.9), \\ r_T &= \text{clip}(r_S \cdot u, r_T^{\min}, r_T^{\max}), \quad u \sim \mathcal{U}(0.3, 0.7), \end{aligned} \quad (31)$$

where  $[r_T^{\min}, r_T^{\max}] = [0.1, 0.6]$ . The multiplicative coupling ensures that the teacher always sees more context than the student (i.e.,  $|\mathcal{M}_T| \leq |\mathcal{M}_S|$ ).

**Curriculum schedule function.** The mixing coefficient  $\lambda(g)$  follows a cosine schedule with warmup:

$$\lambda(g) = \begin{cases} \lambda_0, & g \leq \alpha, \\ \lambda_0 + (\lambda_1 - \lambda_0) \cdot \frac{1 - \cos(\pi \frac{g-\alpha}{1-\alpha})}{2}, & g > \alpha, \end{cases} \quad (32)$$

where  $g \in [0, 1]$  is the normalized training progress,  $\alpha$  is the warmup fraction, and  $\lambda_0 > \lambda_1$  ensures the transition from distillation-dominated to reconstruction-dominated training.

## D.2 Adaptive Decoding Algorithm Pseudocode

Algorithm 2 presents the CAD procedure. The algorithm operates within the block diffusion framework and adaptively selects how many tokens to commit at each step based on model confidence.

### Key algorithmic properties.

1. **Block-wise decoding** (lines 7-12): Positions are partitioned into  $J = \lceil L_{\text{gen}}/b \rceil$  contiguous blocks. CAD processes one block at a time, advancing to the next block only when the current one is fully resolved. Pure diffusion is recovered by setting  $b = L_{\text{gen}}$ .
2. **Confidence computation** (line 17): The confidence score  $c_i^{(s)} = p_{i, \hat{z}_i}$  is the probability mass assigned to the predicted token  $\hat{z}_i = \arg \max_v p_{i,v}$ . This equals the maximum probability  $\max_v p_{i,v}$  by construction.
3. **Confidence thresholding** (line 25): Positions with confidence  $c_i^{(s)} \geq \gamma_{\text{conf}}$  are candidates for commitment. This implements the optimal decision rule derived from the risk-efficiency trade-off (Appendix C).
4. **Cardinality clamping** (line 26): The commit count  $k^{(s)}$  satisfies:

$$k^{(s)} = \text{clip}(\max(n_{\text{conf}}, k_{\min}), k_{\min}, \min(k_{\max}, |\mathcal{V}^{(s)}|)). \quad (33)$$

The lower bound  $k_{\min}$  guarantees progress (liveness); the upper bound  $k_{\max}$  prevents overcommitment from unreliable confidence estimates (trust region).

5. **EOS blocking** (lines 18-21): Premature  $\langle \text{EOS} \rangle$  predictions are suppressed until at least  $\beta_{\text{EOS}} \cdot L_{\text{gen}}$  tokens have been decoded, preventing degenerate short outputs. When all positions predict EOS (lines 22-24), the algorithm forces progress by selecting  $k_{\min}$  tokens regardless of blocking.



---

**Algorithm 2** Confidence-Adaptive Decoding (CAD)

---

**Require:** Student model  $p_\theta$ , prompt  $x$ **Require:** Maximum generation budget  $L_{\text{gen}}$ **Require:** Block size  $b$ , confidence threshold  $\gamma_{\text{conf}}$ **Require:** Acceptance bounds  $k_{\min}, k_{\max}$ , EOS blocking ratio  $\beta_{\text{EOS}}$ , max steps  $S_{\max}$ 

```
1:  $\tilde{z}^{(0)} \leftarrow (x; \underbrace{\mathbf{m}, \dots, \mathbf{m}}_{L_{\text{gen}}})$   $\triangleright$  Initialize with all masks

2:  $\mathcal{M}^{(0)} \leftarrow \{1, \dots, L_{\text{gen}}\}$   $\triangleright$  All target positions masked
3:  $J \leftarrow \lceil L_{\text{gen}}/b \rceil, j \leftarrow 1$   $\triangleright$  Number of blocks, active block index
4:  $s \leftarrow 0, n_{\text{decoded}} \leftarrow 0$   $\triangleright$  Step counter, decoded token count
5: while  $\mathcal{M}^{(s)} \neq \emptyset$  and  $s < S_{\max}$  do
  // — Step 1: Compute eligible set within active block —
  6:  $\mathcal{B}_j \leftarrow \{(j-1)b+1, \dots, \min(jb, L_{\text{gen}})\}$ 
  7:  $\mathcal{E}^{(s)} \leftarrow \mathcal{M}^{(s)} \cap \mathcal{B}_j$   $\triangleright$  Eligible = masked  $\cap$  active block
  8: if  $\mathcal{E}^{(s)} = \emptyset$  then
    9:  $j \leftarrow j+1$   $\triangleright$  Advance to next block
  10: continue
  11: end if
  // — Step 2: Model forward and confidence computation —
  12:  $\ell \leftarrow p_\theta(\cdot | \tilde{z}^{(s)})$   $\triangleright$  Logits for all positions
  13:  $p \leftarrow \text{softmax}(\ell)$   $\triangleright$  Token probabilities
  14:  $\hat{z}_i \leftarrow \arg \max_v \ell_{i,v}$  for  $i \in \mathcal{E}^{(s)}$   $\triangleright$  Predicted tokens
  15:  $c_i^{(s)} \leftarrow p_{i, \hat{z}_i}$  for  $i \in \mathcal{E}^{(s)}$   $\triangleright$  Confidence = prob of argmax token
  // — Step 3: EOS blocking —
  16: if  $n_{\text{decoded}}/L_{\text{gen}} < \beta_{\text{EOS}}$  then
    17: for  $i \in \mathcal{E}^{(s)}$  where  $\hat{z}_i \in \{\langle \text{EOS} \rangle, \langle \text{im\_end} \rangle, \dots\}$  do
      18:  $c_i^{(s)} \leftarrow -\infty$   $\triangleright$  Suppress premature EOS
    19: end for
  20: end if
  // — Step 4: Confidence-adaptive selection —
  21:  $\mathcal{V}^{(s)} \leftarrow \{i \in \mathcal{E}^{(s)} : c_i^{(s)} > -\infty\}$   $\triangleright$  Valid (non-blocked) positions
  22: if  $\mathcal{V}^{(s)} = \emptyset$  then
    23:  $k^{(s)} \leftarrow \min(k_{\min}, |\mathcal{E}^{(s)}|)$   $\triangleright$  Force progress when all blocked
  24:  $\mathcal{S}^{(s)} \leftarrow \text{TopK}(\{c_i^{(s)}\}_{i \in \mathcal{E}^{(s)}}, k^{(s)})$ 
  25: else
    26:  $n_{\text{conf}} \leftarrow |\{i \in \mathcal{V}^{(s)} : c_i^{(s)} \geq \gamma_{\text{conf}}\}|$   $\triangleright$  Count above threshold
    27:  $k^{(s)} \leftarrow \text{clip}(\max(n_{\text{conf}}, k_{\min}), k_{\min}, \min(k_{\max}, |\mathcal{V}^{(s)}|))$ 
    28:  $\mathcal{S}^{(s)} \leftarrow \text{TopK}(\{c_i^{(s)}\}_{i \in \mathcal{V}^{(s)}}, k^{(s)})$   $\triangleright$  Select top- $k$  positions
  29: end if
  // — Step 5: Commit predictions —
  30: for  $i \in \mathcal{S}^{(s)}$  do
    31:  $\tilde{z}_i^{(s+1)} \leftarrow \hat{z}_i$ 
  32: end for
  33:  $\mathcal{M}^{(s+1)} \leftarrow \mathcal{M}^{(s)} \setminus \mathcal{S}^{(s)}$ 
  34:  $n_{\text{decoded}} \leftarrow n_{\text{decoded}} + |\mathcal{S}^{(s)}|$ 
  35:  $s \leftarrow s+1$ 
36: end while
37: return  $\tilde{z}^{(s)}$   $\triangleright$  Fully decoded sequence
```

---

**Complexity analysis.** Let  $S$  be the total number of decoding steps. The CAD controller performs at most  $\mathcal{O}(b \log b)$  operations per step (sorting within the eligible set), yielding total overhead  $\mathcal{O}(S \cdot b \log b)$ . This is negligible compared to the model forward cost  $\mathcal{O}(S \cdot N_{\text{layer}}(L^2 D + L D^2))$ , as detailed in Appendix C.

**Step bound.** By Corollary 1, the total NFE is bounded by:

$$S \leq \left\lceil \frac{L_{\text{gen}}}{k_{\min}} \right\rceil + \left\lceil \frac{L_{\text{gen}}}{b} \right\rceil, \quad (34)$$

plus the implementation cap  $S_{\max}$ . With default parameters  $k_{\min} = 1$  and  $b = L_{\text{gen}}$ , this reduces to  $S \leq L_{\text{gen}} + 1$  in the worst case; in practice, adaptive acceptance yields  $S \ll L_{\text{gen}}$ .

### D.3 Implementation Details and Hyperparameters

**Training hyperparameters.** Table D.1 summarizes all training hyperparameters with their symbols matching the notation in Sect. 3.2.

Hyperparameter	Symbol	CODE	MATH
<b>Optimization:</b>			
Optimizer	–	AdamW	AdamW
Learning rate	$\eta$	$5 \times 10^{-6}$	$5 \times 10^{-6}$
Weight decay	–	0.01	0.01
Betas	$(\beta_1, \beta_2)$	(0.9, 0.999)	(0.9, 0.999)
Gradient clipping	–	1.0	1.0
Warmup steps	–	300	200
Effective batch size	$B$	64	64
Training epochs	–	3	10
<b>Distillation:</b>			
Temperature	$\tau$	1.5	2.0
Initial $\lambda$	$\lambda_0$	0.9	0.9
Final $\lambda$	$\lambda_1$	0.5	0.5
Warmup ratio	$\alpha$	0.1	0.1
Student mask ratio	$r_S \in [\cdot]$	$[0.4, 0.9]$	$[0.4, 0.9]$
Teacher mask ratio	$r_T \in [\cdot]$	$[0.1, 0.6]$	$[0.1, 0.6]$
Ratio coupling factor	$u$	$\mathcal{U}(0.3, 0.7)$	$\mathcal{U}(0.3, 0.7)$
<b>Data:</b>			
Dataset	–	OpenCodeInstruct	GSM8K
Training samples	–	190K	7,473

Table D.1: **Training hyperparameters and notation.**

**Inference hyperparameters.** Table D.2 summarizes all inference hyperparameters with their symbols matching the notation in Sect. 3.3.

**Notation correspondence.** For clarity, we summarize the correspondence between main-text notation and appendix notation:

- $\tilde{z}^S, \tilde{z}^T$  (Sect. 3.2)  $\leftrightarrow S, T$  (Appendix A)
- $\mathcal{M}_S, \mathcal{M}_T$  (masked position sets)  $\leftrightarrow m_S, m_T$  (mask patterns)
- $\mathcal{M}^{(s)}$  (Sect. 3.3) = masked set at decoding step  $s$
- $\mathcal{E}^{(s)}$  = eligible set =  $\mathcal{M}^{(s)} \cap \mathcal{B}_j$  (active block intersection)
- $\mathcal{V}^{(s)}$  = valid set = non-blocked positions in  $\mathcal{E}^{(s)}$
- $\mathcal{S}^{(s)}$  = commit set (positions to unmask at step  $s$ )
- $c_i^{(s)}$  = confidence score =  $p_{i, \tilde{z}_i} = \max_v p_\theta(v \mid \tilde{z}^{(s)})_i$

Hyperparameter	Symbol	Value
<b>CAD Decoding:</b>		
Confidence threshold	$\gamma_{\text{conf}}$	0.95
Min tokens per step	$k_{\text{min}}$	1
Max tokens per step	$k_{\text{max}}$	32
EOS blocking ratio	$\beta_{\text{EOS}}$	0.3
Max generation budget	$L_{\text{gen}}$	256
Max NFE	$S_{\text{max}}$	512
<b>Block Diffusion:</b>		
Block size	$b$	32/256
<b>Sampling:</b>		
Temperature (pass@1)	$\tau_{\text{samp}}$	0
Temperature (pass@5)	$\tau_{\text{samp}}$	1.0
<b>Fixed-Step Baseline:</b>		
Diffusion steps	$T$	256

Table D.2: Inference hyperparameters for CAD decoding.

## E Ablation Study

To disentangle the contributions of the proposed training and decoding components, we conduct a comprehensive ablation study. Although our main distillation experiments focus on code generation, we select **GSM8K** as the primary testbed for this analysis. GSM8K problems induce longer chain-of-thought sequences, making decoding artifacts easier to diagnose, and reasoning accuracy is highly sensitive to partial errors. To ensure a strictly fair comparison regarding model capacity, we introduce an SFT baseline trained on the same GSM8K dataset using full fine-tuning, identical to our DSCD training configuration. This allows us to isolate the impact of the distillation objective from the benefits of standard supervised learning.

### E.1 Effectiveness and Efficiency Analysis

Table E.1 presents the quantitative results. We draw two key conclusions regarding the superiority of DSCD over standard SFT:

**1. DSCD is essential for Global Denoising (Validity).** In the local regime ( $b=32$ ), both SFT and DSCD perform comparably to the teacher, as the denoising task is relatively simple. However, the distinction becomes evident in the global regime ( $b=256$ ). The SFT baseline fails to adapt to the diffusion generation process, improving accuracy only marginally from 13.8% (Teacher) to 21.9%. In contrast, DSCD significantly boosts accuracy to 54.8%. This implies that standard supervised learning captures the data distribution but fails to learn the *many-to-many* dependency structure required for long-range parallel decoding.

**2. DSCD unlocks Aggressive Acceleration (Efficiency).** Comparing the accelerated settings reveals that our training objective is crucial for the stability of CAD. While **SFT + CAD** provides a  $3.13\times$  speedup, its accuracy collapses to 38.2% at  $b=256$ , indicating that the SFT model lacks calibrated confidence scores to guide the drafting process. Conversely, **DSCD + CAD** maintains high robustness (54.7%) while achieving a significantly higher speedup of  $4.91\times$ . This demonstrates that DSCD not only improves generation quality but also aligns the model’s internal confidence with the decoding policy, enabling faster convergence without error propagation.

### E.2 Analysis of EOS Blocking: Structural Regularization vs. Intrinsic Learning

A recurring failure mode of *pure diffusion* decoding with long generation lengths is *premature collapse to <EOS>*, producing extremely short, repetitive outputs. In our setting, this pathology is severe for fixed decoding at full-sequence length ( $b=256$ ): Without intervention, 99.8% of samples terminate within the first 64 tokens.

Configuration	Block Size	Acc (%)	Avg. NFE ↓	Speedup ↑
LLaDA-8B-Instruct (Teacher)	32	77.4	256.0	1.00×
LLaDA-8B-Instruct (Teacher)	256	13.8	256.0	1.00×
<b>Component Analysis:</b>				
+ SFT only (fixed decode)	32	77.6	256.0	1.00×
+ DSCD only (fixed decode)	32	<b>78.3</b>	256.0	1.00×
+ CAD only (no distill)	32	77.4	98.8	4.01×
+ SFT + CAD	32	77.6	95.4	4.15×
+ <b>DSCD + CAD (Ours)</b>	32	77.6	<b>76.3</b>	<b>5.18×</b>
+ SFT only (fixed decode)	256	21.9	256.0	1.00×
+ DSCD only (fixed decode)	256	54.8	256.0	1.00×
+ CAD only (no distill)	256	32.7	130.5	3.07×
+ SFT + CAD	256	38.2	128.1	3.13×
+ <b>DSCD + CAD (Ours)</b>	256	54.7	80.6	4.91×
<b>CAD Ablations (w/ DSCD):</b>				
w/o confidence ranking	32	50.6	202.2	1.98×
w/o EOS blocking	32	77.4	76.3	5.18×
<b>CAD Ablations (w/o DSCD):</b>				
w/o EOS blocking	256	5.6	123.7	3.21×

Table E.1: **Comprehensive ablation study on GSM8K** ( $L_{\text{gen}} = 256$ ). We evaluate two regimes: **block diffusion** ( $b=32$ , local denoising) and **pure diffusion** ( $b=256$ , global denoising). **SFT baselines** use full fine-tuning on the same data to serve as a direct control group for our distillation method. Unless otherwise stated, fixed decoding uses NFE= 256.

**The implicit containment effect.** Interestingly, for small blocks ( $b=32$ ), disabling EOS blocking has negligible impact (77.4% vs. 77.6%). We attribute this to an *implicit containment effect*: Since decoding is localized to 32 positions, an early  $\langle \text{EOS} \rangle$  prediction cannot globally terminate the sequence. Subsequent blocks are still initialized as masks and denoised independently, effectively acting as a **structural barrier** against error propagation.

**Decoupling defense mechanisms.** The contrast becomes clear at full-sequence length ( $b=256$ ). While explicit EOS blocking (via CAD) substantially reduces early termination and recovers accuracy to 32.7%, combining it with DSCD further boosts accuracy to 54.7%. These results indicate that EOS blocking and distillation address orthogonal layers of the failure mode:

- **Inference-time constraint:** EOS blocking acts as a hard guardrail, mechanically preventing termination before sufficient content is generated.
- **Intrinsic reshaping:** DSCD fundamentally alters the student’s denoising behavior. Unlike the baseline, the distilled model learns to assign a lower probability mass to  $\langle \text{EOS} \rangle$  during high-noise states, intrinsically favoring the completion of reasoning chains.

**Impact of confidence ranking.** To validate our scoring mechanism, we replace confidence-ranked token finalization with **uniform random** selection. This causes a sharp accuracy drop (77.6%  $\rightarrow$  50.6%) and increases NFE (76.3  $\rightarrow$  202.2). This result confirms that *which* tokens are finalized is critical: Confidence ranking is essential for accelerating diffusion without derailing the generation trajectory.

### E.3 Stability Analysis of Teacher-Subset Masking

We explicitly investigate the necessity of the teacher-subset constraint ( $\mathcal{M}_T \subseteq \mathcal{M}_S$ ) proposed in Sect. 3.2. In our preliminary experiments, we attempted to train the student using independent masking schedules, where  $\mathcal{M}_T$  and  $\mathcal{M}_S$  are sampled independently. We observed that this configuration leads to severe numerical instability: The training loss frequently diverges, manifesting as gradient NaN within the first epoch.

This empirical collapse validates our theoretical variance analysis in Appendix B. Under independent masking, the teacher often lacks access to tokens visible to the student (i.e.,  $\mathcal{M}_S \setminus \mathcal{M}_T \neq \emptyset$ ), causing it

to marginalize over these “blind” positions. This introduces extreme variance in the distillation targets, destabilizing the gradient estimator. Consequently, we conclude that the strict subset constraint is not merely a theoretical preference for variance reduction, but a practical prerequisite for the convergence of discrete consistency distillation.

## F Additional Experimental Analysis

### F.1 LLM as a Judge Prompt and Generation Examples

To ensure a rigorous and reproducible evaluation of generation quality, we employed a strict “LLM-as-a-Judge” protocol using Qwen3-1.7B. Unlike standard scoring prompts that directly request a number, our system prompt for text-quality evaluation (presented ahead) enforces a “Think-Then-Score” mechanism. It explicitly instructs the model to reason silently about artifacts, such as stuttering, logical jumps, and formatting noise, before acting as a strict parser to output the final scores. This ensures that the generated JSON metrics reflect a deep comprehension of the text quality without polluting the output with intermediate reasoning steps. We focus on three orthogonal dimensions: Fluency, Repetition, and Structural Coherence.

Below, we provide qualitative comparisons between our method and the LLaDA baseline. For these examples, the generation canvas length was set to  $L_{\text{gen}} = 128$ ; the LLaDA baseline utilized a fixed schedule of 32 steps (decoding 4 tokens per step), whereas our method adaptively converged in 29 steps. As shown in Example 1, our model produces coherent chain-of-thought reasoning with clean formatting. In contrast, Example 2 illustrates typical failure modes of fixed-step diffusion decoding, where the baseline suffers from severe token repetition (e.g., “fresh duck fresh egg”) and arithmetic hallucinations, resulting in significantly lower judge scores.

#### Example 1: High-Quality Generation (Ours)

**Question:** Janet’s ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers’ market daily for \$2 per fresh duck egg. How much in dollars does she make every day at the farmers’ market?

**Generated Answer:**

Janet eats 3 eggs per day =  $\ll 3 \times 1 = 3 \gg 3$  eggs for breakfast. She bakes 4 eggs per day =  $\ll 4 \times 1 = 4 \gg 4$  eggs for her friends. She sells  $16 - 3 - 4 = \ll 16 - 3 - 4 = 9 \gg 9$  eggs at the farmers’ market each day. She sells 9 eggs at the farmers’ market for \$2 per egg. Thus, she makes  $9 \times \$2 = \$\ll 9 \times 2 = 18 \gg 18$  every day at the farmers’ market. ##### 18

**LLM Judge:** {"Fluency": 5, "Repetition": 5, "Structure": 4, "Average": 4.67}

#### Example 2: Low-Quality Generation (LLaDA)

**Question:** Janet’s ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers’ market daily for \$2 per fresh duck egg. How much in dollars does she make every day at the farmers’ market?

**Generated Answer:**

Janet eats ducks lay three eggs every breakfast and four eggs for her friends friends every day totaling  $3 + 4 = \ll 3 + 4 = 7 \gg 7$  eggs. \n Janet’s ducks lay 16 eggs to sell the market but after eating and bakes with 7 she has she has  $16 - 7 = \ll 16 - 7 = 9 \gg 9$  eggs to sell at the farmers’ market. \n She sells each fresh duck fresh egg for \$2 so  $9 \times \$2 = \$\ll 9 \times 2 = 18 \gg 18$  every day day at the farmers’ market. \n ##### 18

**LLM Judge:** {"Fluency": 3, "Repetition": 3, "Structure": 4, "Average": 3.33}

#### Detailed Score Breakdown

Table F.1 summarizes the quantitative results across the full test set. Our method achieves a statistically significant improvement across all three dimensions. Notably, the Repetition score sees the largest

gain (+0.76), confirming that our method effectively suppresses the “stuttering” artifacts common in low-NFE generation.

Dimension	LLaDA-8B-Ins	Ours
Fluency (1-5)	3.96 ( $\pm$ 1.28)	<b>4.56</b> ( $\pm$ 0.79)
Repetition (1-5)	3.68 ( $\pm$ 1.51)	<b>4.44</b> ( $\pm$ 0.98)
Structure (1-5)	4.04 ( $\pm$ 1.22)	<b>4.64</b> ( $\pm$ 0.67)
<b>Avg.</b>	3.90 ( $\pm$ 1.32)	<b>4.54</b> ( $\pm$ 0.48)

Table F.1: **Quantitative evaluation of generation quality on GSM8K.** We report Mean ( $\pm$  Standard Deviation) scores across three dimensions (1-5 scale), showing significant gains in reducing repetition and improving structure.

## F.2 Unmasking Analysis

To investigate the generation mechanism of **CD<sup>4</sup>LM**, we visualize the step-wise unmasking trajectory for a representative code generation example (presented ahead) from the HumanEval benchmark in Fig. F.1. The heatmap color-codes the decoding step at which each token is finalized (committed), ranging from blue (early steps, 0  $\sim$  10) to red (late steps, 30  $\sim$  47).

The visualization reveals that our CAD policy enables spontaneous emergence of a hierarchical generation strategy, structurally distinct from the linear order of autoregressive models:

- **Phase I: Syntactic scaffolding (blue/green).** In the initial steps, the model prioritizes structural tokens with global receptive fields. Python keywords (`def`, `if`, `return`), function signatures, and control flow indentations are consistently unmasked first. This indicates that the model establishes a high-confidence *syntactic scaffold* early in the process, effectively reducing the search space for subsequent tokens.
- **Phase II: Logical refinement (yellow/red).** Once the structure is fixed, the model focuses its compute budget on high-entropy positions. Complex arithmetic expressions (e.g., `length = end - start + 1`) and conditional predicates (e.g., `if n % i == 0`) appear in warmer colors, indicating they are refined much later. Crucially, these tokens are denoised conditioning on the *already-visible* syntactic context from Phase I.

This distinct temporal separation confirms that **CD<sup>4</sup>LM** does not merely memorize sequences but learns to decouple global structural planning from local logical execution. By deferring uncertain tokens, CAD automatically allocates more denoising steps to the most difficult parts of the reasoning chain, explaining the efficiency gains observed in Sect. 4.

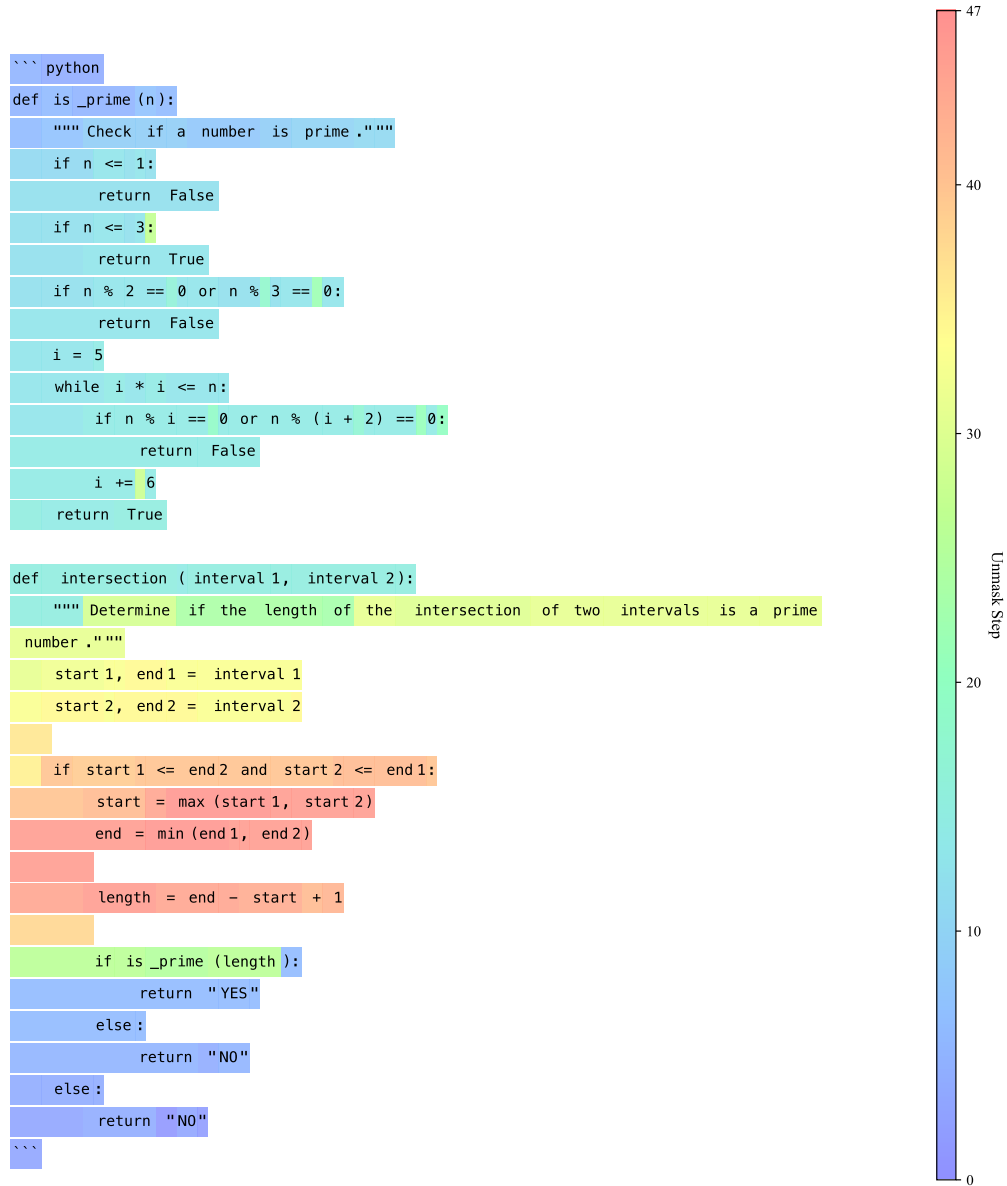


Figure F.1: **Visualization of hierarchical generation dynamics.** Unlike autoregressive models that generate sequentially,  $\text{CD}^4\text{LM}$  exhibits a “skeleton-first” strategy: Syntactic structures (blue, early steps) are established first to form a global plan, while complex logical details (red, late steps) are refined later within the established context.

## Code Generation Example

### Task:

```
def intersection(interval1, interval2):  
    """
```

You are given two intervals, where each interval is a pair of integers. For example, interval = (start, end) = (1, 2). The given intervals are closed, which means that the interval (start, end) includes both start and end. For each given interval, it is assumed that its start is less than or equal to its end. Your task is to determine whether the length of intersection of these two intervals is a prime number. Example: the intersection of the intervals (1, 3), (2, 4) is (2, 3) whose length is 1, which is not a prime number. If the length of the intersection is a prime number, return "YES", otherwise, return "NO". If the two intervals don't intersect, return "NO".

[input/output] samples:

```
intersection((1, 2), (2, 3)) ==> "NO"  
intersection((-1, 1), (0, 4)) ==> "NO"  
intersection((-3, -1), (-5, 5)) ==> "YES"  
    """
```

### Generated Code:

```
def is_prime(n):  
    """Check if a number is prime."""  
    if n <= 1:  
        return False  
    if n <= 3:  
        return True  
    if n % 2 == 0 or n % 3 == 0:  
        return False  
    i = 5  
    while i * i <= n:  
        if n % i == 0 or n % (i + 2) == 0:  
            return False  
        i += 6  
    return True  
  
def intersection(interval1, interval2):  
    """Determine if the length of the intersection of two intervals is a prime number."""  
    start1, end1 = interval1  
    start2, end2 = interval2  
  
    if start1 <= end2 and start2 <= end1:  
        start = max(start1, start2)  
        end = min(end1, end2)  
  
        length = end - start + 1  
  
        if is_prime(length):  
            return "YES"  
        else:  
            return "NO"  
    else:  
        return "NO"
```



## System Prompt for Text-Quality Evaluation

**Role:** You are a STRICT text-quality scorer (long-context verifier). Evaluate ONLY on three dimensions: Fluency, Repetition, Structural Coherence. Treat equations, numbers, and symbols as plain text. Ignore factual correctness.

**THINK-THEN-SCORE (important):**

- First, REASON SILENTLY: derive evidence and compute micro-scores per dimension IN YOUR HEAD. Do NOT reveal any rationale or intermediate text.
- Then, produce the FINAL JSON (single line) as specified below.

**STRICT OUTPUT CONTRACT (hard constraints):**

- Your ENTIRE reply MUST be EXACTLY ONE line containing ONLY a tagged JSON object:

```
<JSON>"fluency": "score":N,"repetition": "score":N,"structural_coherence": "score":N</JSON>
```

- Where N are integers in [1,5]. Keys MUST match exactly; the JSON MUST be valid & minified.
- No prose, no markdown, no code fences, no extra spaces/newlines before/after the tag.
- If you cannot comply, output EXACTLY: <JSON>{"error": "FORMAT"}</JSON>.

**SCORING VIEWS & AGGREGATION (for robustness in long contexts):**

- For EACH dimension, privately compute 3 micro-scores (views) and average → round to nearest integer, then apply the CAPS below. (Do NOT print micro-scores.)
- Fluency views: grammar, clarity, readability/flow.
- Repetition views: identical-token runs, dense short-fragment recurrence, overall readability impact.
- Structural Coherence views: organization, transitions, narrative completeness.

**CAPS / HARD RULES (apply AFTER averaging; lowers over-optimistic scores):**

**Repetition:**

- If  $\geq 1$  run of  $\geq 3$  identical tokens exists more than once → repetition  $\leq 3$ .
- If frequent runs OR dense fragment recurrences harm readability → repetition  $\leq 2$ .

**Structural Coherence:**

- If  $> 60\%$  of lines are bare calculations/symbols with little connective prose → coherence  $\leq 2$ .
- If the answer is very short ( $< 20$  tokens) → coherence  $\leq 3$  unless obviously well-structured.

**Fluency:**

- If pervasive grammatical errors/choppy telegraphic style hinder comprehension → fluency  $\leq 2$ .

**STRICTNESS + FORMAT SENSITIVITY (apply BEFORE final JSON):**

- Start each dimension at 3; raise to 4 or 5 ONLY when there is explicit evidence of excellent writing that meets the top-tier rubric. When in doubt, stay at 3.
- If the response contains  $\geq 3$  blank lines or double-newline spacing between most sentences, treat the flow as fragmented → structural\_coherence  $\leq 3$ , and fluency  $\leq 3$  unless the prose still reads seamlessly.
- Heavy formatting noise (raw LaTeX delimiters such as  $\backslash$ [,  $\backslash$ \$,  $\backslash$ begin{align}, Markdown tables/lists, or numbered steps that merely restate facts) disrupts readability → cap fluency at 3 and structural\_coherence at 3.
- Short filler sentences that repeat the same idea without advancing the solution lower structural\_coherence by at least 1 point and may also lower repetition.
- Concise single-block reasoning with minimal padding and no formatting noise **may** still earn 4-5 when it is genuinely smooth and well-ordered.

**DETAILED RUBRICS (1=worst, 5=best):**

**Fluency** (grammar/clarity/readability; ignore factuality)

- 5: Nearly error-free; clear sentences; natural flow; varied syntax; no distracting formatting noise.
- 4: Minor issues but clear overall; formatting and spacing remain unobtrusive.
- 3: Noticeable errors/awkward phrasing or spacing quirks mildly impede flow.
- 2: Frequent errors or choppy style hinder comprehension.
- 1: Very poor grammar/word salad; hard to understand.

**Repetition** (penalize meaningless identical-word/span loops; DO NOT penalize necessary reuse of terms/digits/operators)

- Definitions:
  - run =  $\geq 3$  identical tokens in a row (e.g., “eggs eggs eggs”).
  - dense fragment recurrence = the same 2-5 word fragment appears many times within a short span.
- 5: No runs; only natural reuse; zero filler restatements.
- 4: One short run OR a few mild recurrences; readability intact with minimal filler.
- 3: Multiple short runs OR several recurrences; still readable but noticeably repetitive.
- 2: Frequent runs OR dense recurrences that harm readability.
- 1: Heavy looping/spam; large stretches of repeated tokens/fragments.

**Structural Coherence** (organization/flow; NOT correctness of content/math)

- Heuristics (apply flexibly, not mechanically):
  - Enumerated steps or clear transitions raise coherence.
  - Mostly raw calculations/symbols lower coherence.
  - Very short answers rarely justify 5 unless clearly structured.
- 5: Well-organized narrative/steps with clear transitions; tight paragraphing without stray blank lines.
- 4: Generally organized; minor jumps but understandable; spacing supports the flow.
- 3: Mixed; some organization but noticeable gaps/abrupt jumps or distracting spacing.
- 2: Mostly disorganized or just raw calculations; hard to follow.
- 1: No discernible structure; fragments/out-of-order snippets.

**FINAL SELF-CHECK (internal; do NOT print):**

- Apply CAPS after averaging the three views per dimension.
- Check for  $\geq 3$ -in-a-row runs or dense recurrences before setting repetition.
- If  $> 60\%$  lines are bare calculations, enforce coherence  $\leq 2$ .
- Ensure the three scores are mutually consistent with rubrics.

**Few-shot Examples (STRICT single-line JSON only):**

**Example 1 (high quality, no repetition)**

Student Answer:

Janet sells 16 - 3 - 4 = «16-3-4=9»9 duck eggs a day. She makes 9 \* 2 = «9\*2=18»18 every day at the farmer's market. #### 18.

Expected Output:

```
<JSON>"fluency": "score":5,"repetition": "score":5,"structural_coherence": "score":5</JSON>
```

**Example 2 (moderate repetition)**

Student Answer:

Janet's ducks lay eggs eggs per day for which.6 eats «3 eggs +4 eggs = «3 eggs=16»16 eggs... friends friends friends... She sells sells sells the remainder... daily day... 16 - 7 = «16-7=9»9... She sells... fresh fresh fresh... makes makes makes...

Expected Output:

```
<JSON>"fluency": "score":3,"repetition": "score":2,"structural_coherence": "score":3</JSON>
```

**Example 3 (extreme repetition)**

Student Answer:

Janet eats eggs eggs eggs breakfast breakfast breakfast br...eakfast eats eats eats breakfast breakfast breakfast ... #### 18

Expected Output:

```
<JSON>"fluency": "score":1,"repetition": "score":1,"structural_coherence": "score":1</JSON>
```