

# Empirical Analysis of Decoding Biases in Masked Diffusion Models

Pengcheng Huang<sup>1</sup>, Tianming Liu<sup>1</sup>, Zhenghao Liu<sup>1</sup>, Yukun Yan<sup>2</sup>, Shuo Wang<sup>2</sup>, Tong Xiao<sup>1</sup>, Zulong Chen<sup>3</sup>, Maosong Sun<sup>2</sup>

1 Northeastern University, China 2 Tsinghua University-NLP Lab 3 Alibaba Group



> 7%

avg. gain over the best decoding baseline

44.7 ≈ 45.3

LLaDA-1.5 + UNCODE vs. Qwen-2.5-7B (AR)

3 × 7

MDM backbones × reasoning & planning tasks

0

extra training — plug-and-play, decoding-only

## Motivation and Background

- Masked Diffusion Models (MDMs) decode by iterative unmasking — any-order, multi-token, non-autoregressive.
- The unmasking order is decisive: standard MDMs greedily unmask the least-uncertain positions first.

## The Two Decoding Biases

(1) **Rigid Boundary Bias**: boundary tokens are always decoded first, collapsing inward — the answer is fixed before the rationale exists.

(2) **Trivial Token Bias**: high-frequency, low-information tokens (punctuation, fillers) are over-prioritized, wasting decoding steps.

**Question**: Jill earns \$20/hour teaching and \$30/hour coaching. What is her weekly salary if she works 35 hours teaching and 15 coaching?

**Ground Truth**:  $(20 * 35) + (30 * 15) = \$1150$

**A: Rigid Boundary Bias** Decode earlier → later: ■■■■■  
Jill makes  $\$20 * 35 = \$700$  as a teacher. Then, adding her \$30 coaching pay gives a total weekly salary of  $700 + 30 = \$730$ .  
**Confidence**: ★★★★ **Boundary-First Decoding**

**B: Trivial Token Bias** Decode earlier → later: ■■■■■  
First, for the teacher salary, she earns  $\$20 * 35 = \$3500$ . ; Then, adding the \$450 coaching pay, the total is \$3500.  
**Confidence**: ★★★★ **Trivial Token Force 4-digit Format (Target: 700)**

**C: Ideal Decoding** Decode earlier → later: ■■■■■  
First, teaching pays  $\$20 * 35 = \$700$ . Next, coaching pays  $\$30 * 15 = \$450$ . In total, she makes  $\$700 + 450 = \$1150$  per week.  
**Confidence**: ★★★★ **Principled Reasoning Trajectory**

Uncertainty-based decoding prioritizes (A) boundaries and (B) trivial tokens, deviating from the (C) ideal reason-then-answer path.

## UNCODE: Unmasking Calibration

A lightweight, training-free calibration of the unmasking priority  $s_t^i$  via two complementary priors:

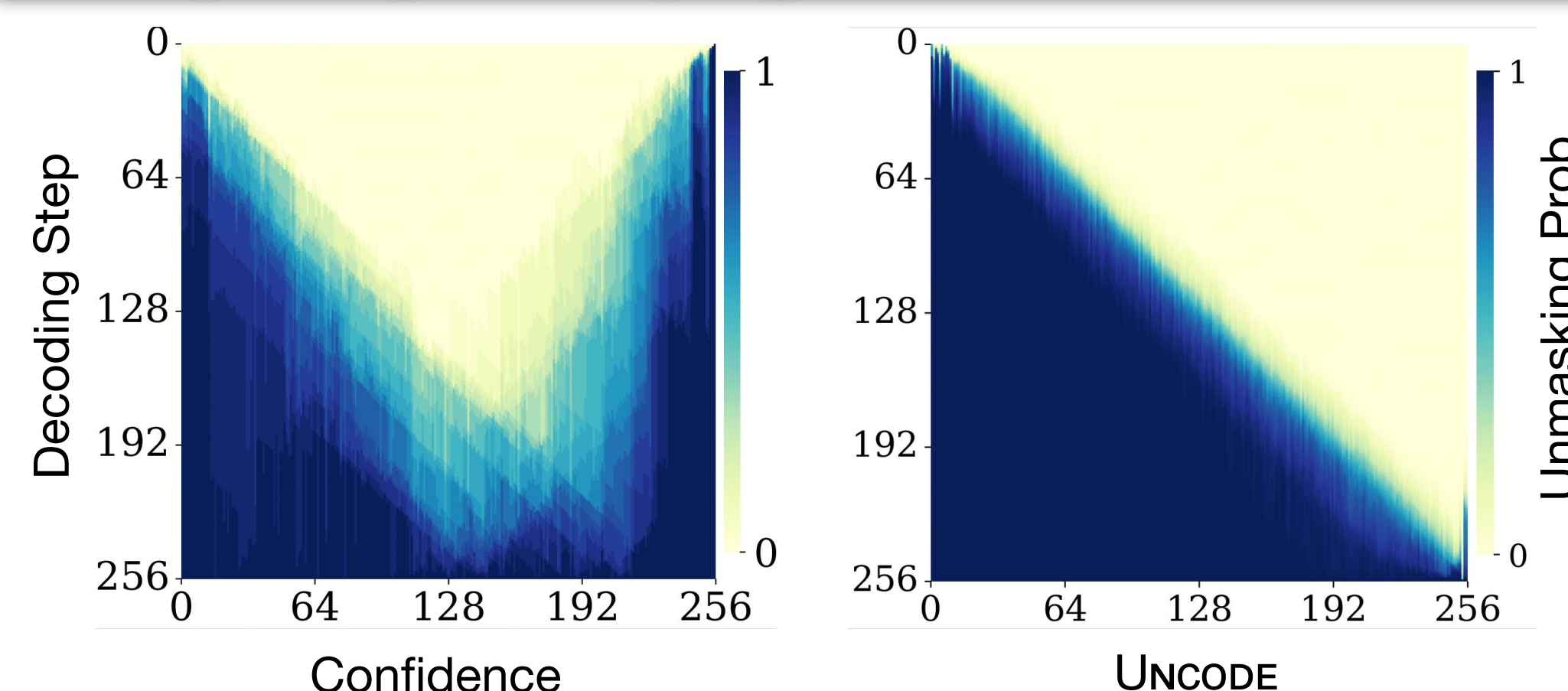
$$\tilde{s}_t^i = \underbrace{P^i}_{\text{position}} \cdot \underbrace{S_t^i}_{\text{semantics}} \cdot \underbrace{\mathcal{F}(p_\theta(\cdot | x_t, i))}_{\text{raw uncertainty}}$$

**Positional Trajectory Prior**:  $P^i = e^{-\lambda i}$  — discourages boundary-first decoding;  $\lambda$  interpolates between flexible any-order and left-to-right generation.

**Semantic Informativeness Prior**:  $S_t^i = \min(-\log_{P_D}(\bar{x}_t^i), \alpha)$ , down-weights frequent trivial tokens via corpus self-information (clipped at  $\alpha$ ).

⇒ Keeps parallel any-order decoding while **globally reshaping the trajectory** and **promoting informative content**.

## Trajectory Reshaping



## Experimental Setup

- 3 MDM backbones: LLaDA-8B-Instruct, LLaDA-1.5-8B, Dream-7B.
- 7 benchmarks: HumanEval, MBPP (code); GSM8K, MATH500 (math); GPQA (science); Countdown, Sudoku (planning).
- 8 decoding baselines: Uniform, Confidence, Entropy, Margin, EB-Sampler, Semi-AR, Fast-dLLM.

## Experimental Results

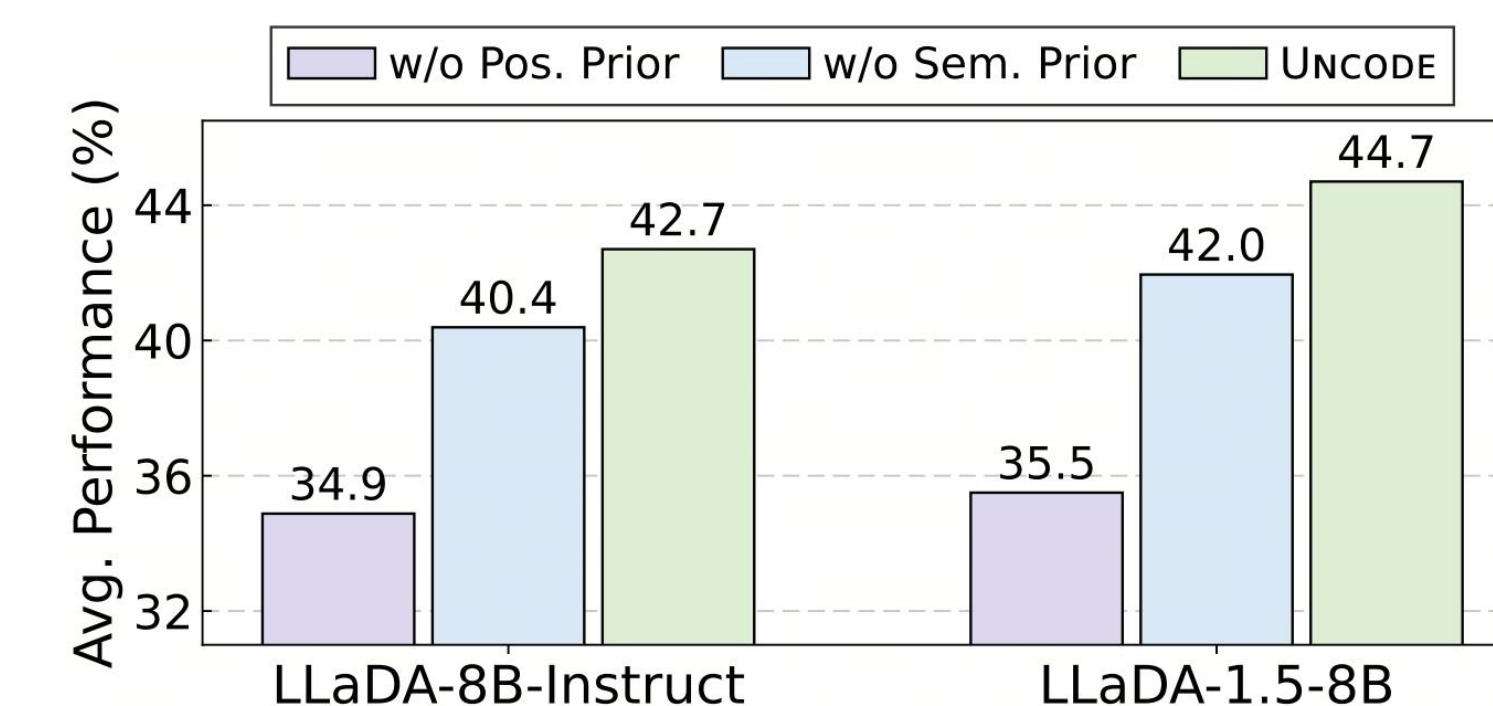
### Overall Performance

- UNCODE outperforms all baseline models across datasets.

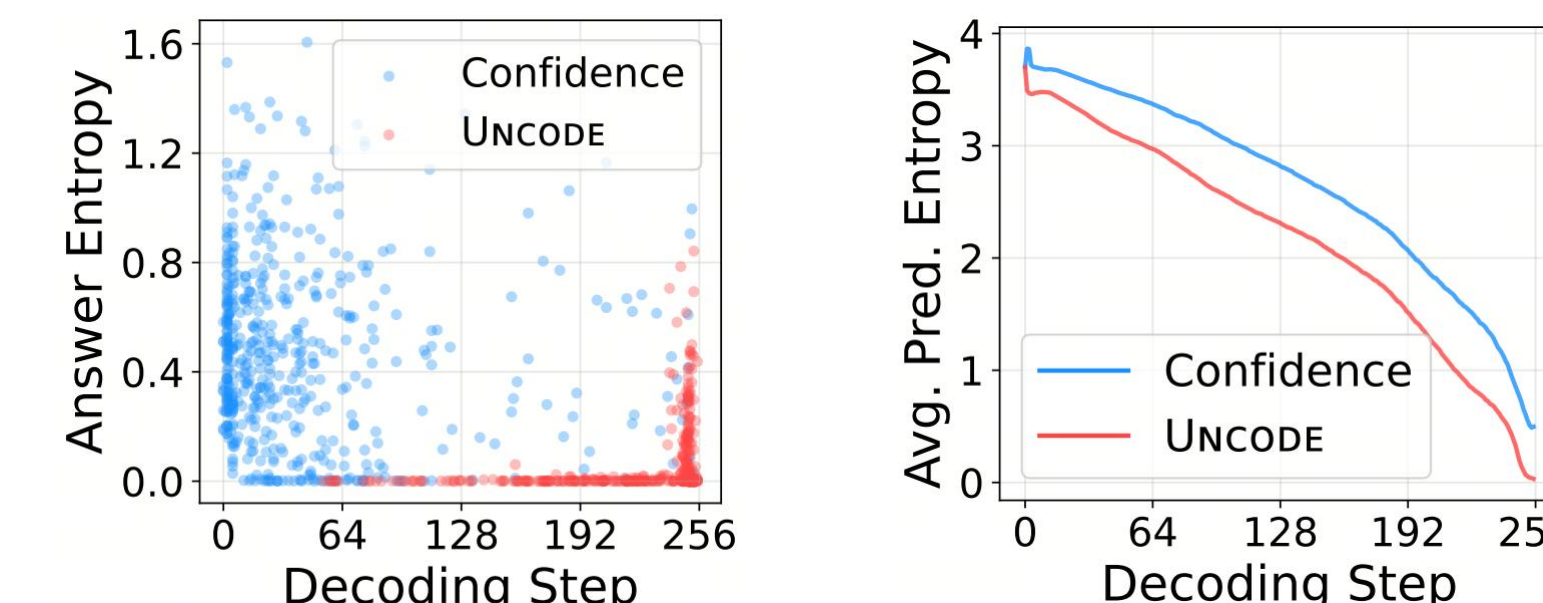
Method	GSM8K	MATH	Sudoku	C.down	Avg.
<i>Autoregressive reference</i>					
Qwen-2.5-7B	71.9	64.2	0.0	7.7	45.3
<i>LLaDA-8B-Instruct</i>					
Semi-AR (best baseline)	77.9	27.6	0.0	32.6	35.7
UNCODE	79.2	34.8	29.8	36.3	42.7
<i>LLaDA-1.5-8B</i>					
Semi-AR (best baseline)	80.7	34.2	0.0	32.4	37.1
UNCODE	82.2	37.4	33.4	35.0	44.7

### Ablation and Analysis

- Both priors matter



- Defers answer tokens to later steps with higher confidence, and reduces global uncertainty faster than the baseline.



## Resources



Explore more interesting experiments in our paper and code, and feel free to contact us at pengcheng.neu@outlook.com or via WeChat