

# T\*: Progressive Block Scaling for MDM Through Trajectory Aware RL

Anonymous ACL submission

## Abstract

We present **T\***, a simple TRACERL-based training curriculum for progressive block-size scaling in masked diffusion language models (MDMs). Starting from an AR-initialized small-block MDM, **T\*** transitions smoothly to larger blocks, enabling higher-parallelism decoding with minimal performance degradation on math reasoning benchmarks. Moreover, further analysis suggests that **T\*** can converge to an alternative decoding schedule  $\hat{S}$  that achieves comparable performance.

## 1 Introduction

Before the current wave of large language models (LLMs), bidirectional Transformers trained with masked language modeling were a widely adopted backbone for NLP systems, with BERT and its optimized variants as canonical examples (Devlin et al., 2019; Liu et al., 2019). Today, autoregressive (AR) modeling via next-token prediction dominates both scaling practice and deployed systems (Brown et al., 2020; Touvron et al., 2023).

In parallel, diffusion language models have begun to emerge as viable alternatives or complements to the autoregressive decoding paradigm. Masked diffusion models stochastically mask a subset of tokens under a ratio-parameterized corruption process and optimize cross-entropy on masked positions to recover the original sequence (Sahoo et al., 2024). For scalability, recent work initializes diffusion LMs from pretrained autoregressive LLMs and trains them with random-mask diffusion objectives (Ye et al., 2025; Cheng et al., 2025). At inference time, they often adopt blockwise decoding that denoises tokens within each block while generating blocks autoregressively to preserve global coherence (Arriola et al., 2025). In this setting, the block size is a control parameter that interpolates between stronger AR-like causality and higher-parallel masked updates.

Within each block, the denoising schedule is typically determined by model confidence. Let  $\mathcal{M}^{(s)}$  be the set of masked positions at denoising step  $s$ , and let the model predict a token distribution at each masked position. A common heuristic defines a confidence score

$$c_i^{(s)} = \max_{v \in \mathcal{V}} p_\theta(x_i = v \mid x^{(s)}, Q), \quad i \in \mathcal{M}^{(s)},$$
$$U^{(s)} = \{i \in \mathcal{M}^{(s)} : c_i^{(s)} \geq \eta\}, \quad (1)$$

and then materializes tokens in  $U^{(s)}$  (e.g., via argmax or sampling), while leaving the rest masked for subsequent refinement.

When examining the **SDAR** series models across scales (1.7B–30B) and block sizes (4–64), we find that math-centric reasoning becomes increasingly sensitive to larger blocks: accuracy generally degrades as block size  $B$  grows, with more pronounced drops for smaller models, which is also reported by Cheng et al. (2025). Under an absorbing-state corruption, the negative-ELBO objective reduces to a reweighted cross-entropy on masked positions:

$$\mathcal{L}(\theta) = \mathbb{E}_{x_0 \sim p_{\text{data}}, x_t \sim q(x_t | x_0), t \sim U(0,1)} \left[ -\frac{1}{t} \sum_{\ell=1}^L \mathbf{1}_{\{x_{t,\ell} = [/MASK]\}} \cdot \log p_\theta(x_{0,\ell} \mid x_t) \right], \quad (2)$$

where  $t \sim U(0, 1)$  controls the masking ratio and  $x_t$  is obtained by independently masking tokens. For a block of size  $B$ , the expected number of  $[/MASK]$  positions is  $tB$ , so larger blocks contain more masked tokens to resolve within each denoising stage. Since we scale block sizes in powers of two ( $B = 2^n$ ), the number of tokens involved per step—and hence the degree of within-block reordering—grows exponentially with the stage index  $n$ . Moreover, standard corpora for supervised fine-tuning (SFT) only specify token targets but do not provide supervision for “correct” unmasking

073 schedule.

074 In this work, we propose  $\mathbf{T}^*$ , an easy-to-  
 075 implement yet effective strategy for progressive  
 076 block-size scaling that increases block size with  
 077 minimal performance degradation.  $\mathbf{T}^*$  offers a prac-  
 078 tical route for masked diffusion models (MDMs)  
 079 to preserve the strong reasoning capability inher-  
 080 ited from AR-initialized small-block models while  
 081 moving toward higher-parallel decoding. Further  
 082 analysis suggests that  $\mathbf{T}^*$  can induce an alterna-  
 083 tive decoding schedule, rather than reverting to the  
 084 canonical left-to-right schedule.

## 085 2 Methodology

### 086 2.1 Trajectory-aware RL

087 We adopt TRACERL as our trajectory-aware rein-  
 088 forcement learning backbone, and build our method  
 089 on top of it (Wang et al., 2025). TraceRL views  
 090 diffusion decoding as a multi-step denoising tra-  
 091 jectory and performs policy optimization on the  
 092 same trajectory used at inference. Given a prompt  
 093  $Q$ , a diffusion LM produces a trajectory  $\tau =$   
 094  $\tau(1) \cup \dots \cup \tau(T)$ , where  $T$  is the number of de-  
 095 noising steps and  $\tau(t)$  denotes the set of tokens  
 096 decoded (unmasked) at step  $t$ . For brevity, we de-  
 097 note the trajectory prefix by  $\tau_{<t} := \tau(1:t-1)$  and  
 098 suppress the dependence on  $Q$  when it is clear.

099 TraceRL applies a PPO-style objective over all  
 100 decoded tokens along the trajectory:

$$\begin{aligned}
 J(\theta) = & \mathbb{E}_{\tau \sim \pi_{\theta_{\text{old}}}} \left[ \sum_{t=1}^T \frac{1}{|\tau(t)|} \sum_{o \in \tau(t)} C_{\epsilon}(\rho_t(o), A(o)) \right] \\
 & - \beta \text{KL}(\pi_{\theta} \parallel \pi_{\theta_{\text{old}}}),
 \end{aligned}
 \tag{3}$$

101 where  $C_{\epsilon}(r, A) = \min\{rA, \text{clip}(r, 1-\epsilon, 1+\epsilon)A\}$   
 102 is the clipped surrogate and  
 103

$$\rho_t(o) = \frac{\pi_{\theta}(o \mid \tau_{<t}, Q)}{\pi_{\theta_{\text{old}}}(o \mid \tau_{<t}, Q)}.
 \tag{4}$$

104 In the simplest verifiable-reward setting, a single  
 105 sequence-level reward (e.g., correctness of the final  
 106 answer) is broadcast to the trajectory and used to  
 107 form the advantages in Eq. 3.

108 To enable finer credit assignment over denoising  
 109 steps, TRACERL aggregates token-level rewards  
 110 (and value predictions) into step-level quantities  
 111 by averaging within each denoising step, and com-  
 112 putes step-wise advantages via TD/GAE (Schul-  
 113 man et al., 2015). These step advantages are then  
 114 assigned back to all tokens decoded at the corre-  
 115 sponding step, so that learning signals propagate  
 116

---

### Algorithm 1: $\mathbf{T}^*$ : Progressive Block Scal- ing

---

**Input:** Base model  $\theta_0, \mathcal{D}$ , Initial  $B_0$ , Target  $\hat{B}$   
**Output:** Optimized model  $\theta$   
*// Initialization*  
 1  $\theta \leftarrow \theta_0; B \leftarrow B_0;$   
*// Progressive Scaling Training*  
 2 **while**  $B \leq \hat{B}$  **do**  
 3     Sample a batch  $d$  from  $\mathcal{D}$ ;  
    *// Data Preparation*  
 4      $d_1, d_2 \leftarrow \text{SPLIT}(d);$   
    *// First Update Step*  
 5      $\theta \leftarrow \text{TRACERL}(\theta, d_1, B);$   
    *// Shift Mechanism*  
 6      $\Delta \leftarrow B/2;$   
 7      $d'_2 \leftarrow \text{SHIFT}(d_2, \Delta);$   
    *// Second Update Step*  
 8      $\theta \leftarrow \text{TRACERL}(\theta, d'_2, B);$   
 9      $B \leftarrow 2B;$   
 10 **return**  $\theta$

---

through the entire denoising trajectory rather than  
 only the final output (Lightman et al., 2023b).

### 2.2 Progressive Block Scaling

We propose a progressive block-scaling strategy,  $\mathbf{T}^*$ , which uses trajectory-aware RL as a catalyst to adapt the denoising policy under the current block size and then relaxes the constraint by enlarging blocks. Concretely, for each stage with block size  $B$ , we run one epoch consisting of three steps. (I) RL: on the first 50% of the training batches, we perform TRACERL updates under the current block partition. (II) SHIFTING: on the remaining 50%, we keep the same block size  $B$  but shift block boundaries by an offset  $\Delta = B/2$ , so that each block straddles two neighboring blocks from the original partition; this reduces boundary artifacts and helps positional representations adapt to cross-block interactions. (III) EXPANDING: we merge adjacent blocks and double the block size,  $B \leftarrow 2B$ , and proceed to the next stage. Algorithm 1 summarizes the training loop.

## 3 Experiments

### 3.1 Setup

We conduct experiments with the masked diffusion models **SDAR-1.7B-Chat** and **SDAR-4B-Chat**. These models are trained via block-diffusion strategy using different block sizes  $B \in \{4, 8, 16, 32\}$ . Our training dataset consists of 8K high-quality mathematical problems with difficult levels 3-5 from Open1math. In each step, we randomly sample 128 problems from the

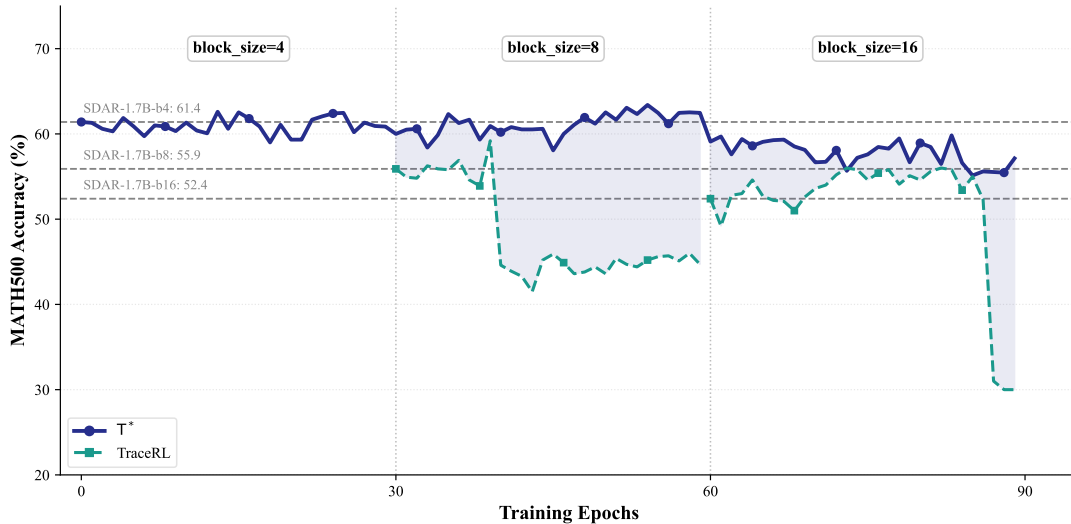


Figure 1: **Validation accuracy during block scaling (1.7B)**. MATH500 validation accuracy over training epochs for  $T^*$  and a direct TRACERL baseline (dashed). Vertical dotted lines indicate stage transitions ( $B=4 \rightarrow 8 \rightarrow 16$ ). Horizontal dashed lines show the accuracies of the original **SDAR** checkpoints trained at each block size.

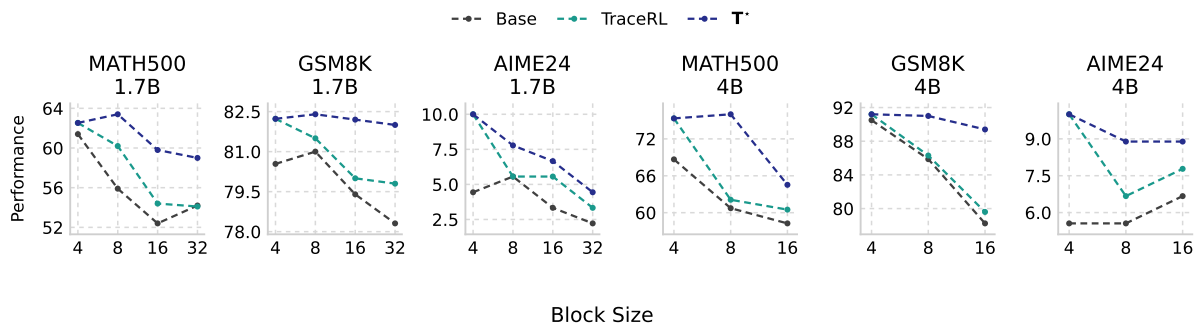


Figure 2: **Performance vs. block size across model scales**. Performance on MATH500, GSM8K, and AIME24 as a function of block size  $B$  for **SDAR** models at 1.7B (left) and 4B (right), comparing the Base model, TRACERL trained at the same  $B$ , and our progressive strategy  $T^*$ .

dataset and generate 16 responses per problem using the static sampling strategy. Training process is performed on an 8-GPU H200 cluster using the AdamW optimizer with a learning rate of  $1 \times 10^{-6}$ . To prevent the policy from collapsing or drifting far from the base model, we apply a KL-divergence penalty with  $\beta = 0.01$ .

**Baselines:** for each block size  $B \in \{8, 16, 32\}$  of **SDAR**, we directly apply 30 epochs of TRACERL training.

**Evaluation** We evaluate the models on MATH500 (Hendrycks et al., 2021; Lightman et al., 2023a), GSM8K (Cobbe et al., 2021) and AIME24 (Art of Problem Solving, 2024a,b). During sampling, we use the same block size for inference as used in training stage and report the Pass@3 accuracy.

### 3.2 General Performance

Figure 1 plots MATH500 validation accuracy throughout training for the 1.7B model. While  $T^*$  remains relatively stable across stages, the direct TRACERL baseline exhibits abrupt collapses: a sharp drop occurs during the  $B=8$  stage (from  $\sim 56\%$  to the low-40% range), and another collapse appears near the end of the  $B=16$  stage (down to  $\sim 30\%$ ). We find this instability is highly sensitive to initialization at the target block size: applying TRACERL directly on the **SDAR-1.7B-Chat-b8** checkpoint collapses, whereas continuing TRACERL at  $B=8$  starting from a TRACERL-trained  $B=4$  diffusion policy (our stage transition) remains stable. A plausible explanation is that larger-block **SDAR** checkpoints operate under weaker conditioning contexts (cf. Eq. 2) and thus start from a lower-

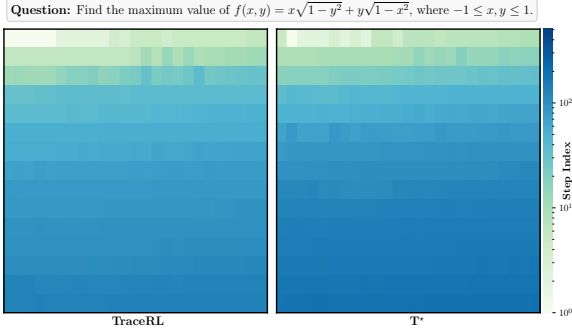


Figure 3: Decoding schedule under TRACERL vs.  $\mathbf{T}^*$ . More results can be found in Appendix A.2

confidence regime, yielding noisier rollouts and higher-variance advantage estimates; when such advantages are broadcast to many tokens per denoising step, ratio-based updates can trigger likelihood drift and collapse, consistent with the Lazy Likelihood-Displacement “death spiral” analysis for GRPO-style training (Deng et al., 2025; Gao et al., 2025).

Figure 2 shows that, across different model sizes,  $\mathbf{T}^*$  consistently matches or exceeds the performance of the base models and TRACERL at the same block size on MATH500, GSM8K, and AIME24. When we expand the block size, the base models generally show a downward trend, while  $\mathbf{T}^*$  remains more stable and achieves the strongest results at most evaluated block sizes; TRACERL often improves over the base model at smaller blocks but is typically below  $\mathbf{T}^*$  at larger blocks. Unless otherwise noted, all scores are reported for the checkpoint that attains the best validation accuracy on MATH500 during training (and are then evaluated on MATH500, GSM8K, and AIME24). The exact scores can be found in Table 2 and Appendix A.1.

### 3.3 Schedule

We compute LOCALSTRICT (Gong et al., 2025). Let  $\pi = (\pi_1, \dots, \pi_n)$  denote the linearized unmasking order obtained by sorting token positions by their first-unmask step (ties broken by smaller positions). LOCALSTRICT is defined as the fraction of events that decode the leftmost remaining masked position:

$$\text{LOCALSTRICT} = \frac{1}{n} \sum_{k=1}^n \mathbb{1} \left[ \pi_k = \min_{j \geq k} \pi_j \right]. \quad (5)$$

Higher values indicate a schedule closer to the canonical left-to-right order  $S_0$ , while lower values

reflect more non-monotone masked updates.

Model	LocalStrict	Accuracy
Qwen3-1.7B	1.000	70.2
Qwen2.5-1.5B	1.000	55.0
SDAR-1.7B-b32	0.743	54.2
+ TRACERL	0.704	54.1
+ $\mathbf{T}^*$	0.730	59.0
SDAR-1.7B-b16	0.766	52.4
+ TRACERL	0.824	54.4
+ $\mathbf{T}^*$	0.804	59.8
SDAR-1.7B-b8	0.915	55.9
+ TRACERL	0.984	60.2
+ $\mathbf{T}^*$	0.854	63.4

Table 1: **LocalStrict vs. accuracy on MATH500.** LocalStrict is computed by Eq. 5; higher values indicate a decoding order closer to the canonical left-to-right schedule.

Figure 3 visualizes token-level first-unmask step indices under TRACERL and  $\mathbf{T}^*$ . Table 1 reports LOCALSTRICT and accuracy under different block sizes. Overall, both methods retain largely monotone unmasking behavior (i.e., LOCALSTRICT remains high), but neither collapses to a strictly deterministic left-to-right schedule; instead, the learned step-wise schedules differ under the target block size (see Appendix A.2 for more examples).

## 4 Conclusion

Experiments show that  $\mathbf{T}^*$  stably scales block size with minimal performance loss, providing a practical recipe for diffusion language models to inherit strong reasoning ability from AR-initialized small-block checkpoints. We analyze the collapse of direct TraceRL at larger block size models and present a potential mitigation perspective related to Lazy Likelihood Displacement. Finally, our schedule analysis suggests that trajectory-aware RL can induce a non-canonical denoising schedule under a fixed block size.

Recent work encourages non-linear reasoning via explicit external scaffolds such as tree/graph-structured search over intermediate thoughts (Yao et al., 2023; Besta et al., 2024; Yao et al., 2024). In contrast, our experiments show that trajectory-aware RL can modify the model’s internal denoising policy (i.e., token-finalization schedule) and improve reasoning performance without introducing an external search procedure, suggesting internal schedule learning as a complementary direction.

## 250 Limitations

251 The limitations of this work can be summarized as:

- 252 •  $T^*$  mitigates but does not fully eliminate  
253 degradation under block expansion; we sus-  
254 pect residual drops are partly due to the lack  
255 of a high-quality “cold-start” stage.
- 256 • We did not scale to very large blocks (e.g.,  
257  $B=64$  or  $128$ ) in our  $T^*$  curriculum, because  
258 the inference engine becomes unstable at large  
259 block sizes.

## 260 Ethical Statements

261 In this paper, we propose strategies to improve the  
262 SQL generation capabilities of LLMs. To the best  
263 of our knowledge, we do not expect our system  
264 would have negative impacts on society.

## 265 References

266 Marianne Arriola, Aaron Gokaslan, Justin T. Chiu, Zhi-  
267 han Yang, Zhixuan Qi, Jiaqi Han, Subham Sekhar  
268 Sahoo, and Volodymyr Kuleshov. 2025. [Block diffusion: Interpolating between autoregressive and diffusion language models](#). In *International Conference on Learning Representations*.

272 Art of Problem Solving. 2024a. [2024 AIME i](#). AoPS  
273 Wiki. Accessed: 2026-01-05.

274 Art of Problem Solving. 2024b. [2024 AIME ii](#). AoPS  
275 Wiki. Accessed: 2026-01-05.

276 Maciej Besta, Nils Blach, Ales Kubicek, Robert Ger-  
277 stenberger, Michal Podstawski, Lukas Gianinazzi,  
278 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadow-  
279 ski, Piotr Nyczyk, and Torsten Hoeffler. 2024. [Graph of thoughts: Solving elaborate problems with large language models](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, Vancouver, Canada*, pages 17682–17690.

284 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
285 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
286 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
287 Askeell, Sandhini Agarwal, Ariel Herbert-Voss,  
288 Gretchen Krueger, Tom Henighan, Rewon Child,  
289 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,  
290 Clemens Winter, Christopher Hesse, Mark Chen, Eric  
291 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,  
292 Jack Clark, Christopher Berner, Sam McCandlish,  
293 Alec Radford, Ilya Sutskever, and Dario Amodei.  
294 2020. [Language models are few-shot learners](#).

295 Shuang Cheng, Yihan Bian, Dawei Liu, Yuhua Jiang, Yi-  
296 hao Liu, Linfeng Zhang, Wenhai Wang, Qipeng Guo,  
297 Kai Chen, Biqing Qi, and Bowen Zhou. 2025. [SDAR: A synergistic diffusion–autoregression paradigm for scalable sequence generation](#).

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 300  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 301  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 302  
Nakano, Christopher Hesse, and John Schulman. 303  
2021. [Training verifiers to solve math word prob- 304  
lems](#). *arXiv preprint arXiv:2110.14168*. 305

Wenlong Deng, Yushu Li, Boying Gong, Yi Ren, Chris- 306  
tos Thrampoulidis, and Xiaoxiao Li. 2025. On 307  
grp collapse in search-r1: The lazy likelihood- 308  
displacement death spiral. *arXiv preprint 309  
arXiv:2512.04220*. 310

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 311  
Kristina Toutanova. 2019. [BERT: Pre-training of 312  
deep bidirectional transformers for language under- 313  
standing](#). In *Proceedings of the 2019 Conference of 314  
the North American Chapter of the Association for 315  
Computational Linguistics: Human Language Tech- 316  
nologies, Volume 1 (Long and Short Papers)*, pages 317  
4171–4186, Minneapolis, Minnesota. Association for 318  
Computational Linguistics. 319

Chang Gao, Chujie Zheng, Xiong-Hui Chen, Kai Dang, 320  
Shixuan Liu, Bowen Yu, An Yang, Shuai Bai, Jingren 321  
Zhou, and Junyang Lin. 2025. [Soft adaptive policy 322  
optimization](#). *arXiv preprint arXiv:2511.20347*. 323

Shansan Gong, Ruixiang Zhang, Huangjie Zheng, Ji- 324  
atao Gu, Navdeep Jaitly, Lingpeng Kong, and Yizhe 325  
Zhang. 2025. [Diffucoder: Understanding and im- 326  
proving masked diffusion models for code genera- 327  
tion](#). 328

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul 329  
Arora, Steven Basart, Eric Tang, Dawn Song, and Ja- 330  
cob Steinhardt. 2021. [Measuring mathematical prob- 331  
lem solving with the MATH dataset](#). In *Advances in 332  
Neural Information Processing Systems*. 333

Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri 334  
Edwards, Bowen Baker, Teddy Lee, Jan Leike, 335  
John Schulman, Ilya Sutskever, and Karl Cobbe. 336  
2023a. [Let’s verify step by step](#). *arXiv preprint 337  
arXiv:2305.20050*. Defines the nonstandard MATH- 338  
500 evaluation split released in the PRM800K repos- 339  
itory. 340

Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri 341  
Edwards, Bowen Baker, Teddy Lee, Jan Leike, John 342  
Schulman, Ilya Sutskever, and Karl Cobbe. 2023b. 343  
[Let’s verify step by step](#). 344

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man- 345  
dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, 346  
Luke Zettlemoyer, and Veselin Stoyanov. 2019. 347  
[Roberta: A robustly optimized BERT pretraining 348  
approach](#). 349

Subham Sekhar Sahoo, Marianne Arriola, Yair Schiff, 350  
Aaron Gokaslan, Edgar Marroquin, Justin T. Chiu, 351  
Alexander Rush, and Volodymyr Kuleshov. 2024. 352  
[Simple and effective masked diffusion language mod- 353  
els](#). In *Advances in Neural Information Processing 354  
Systems*. 355

356 John Schulman, Philipp Moritz, Sergey Levine,  
357 Michael I. Jordan, and Pieter Abbeel. 2015. [High-](#)  
358 [dimensional continuous control using generalized](#)  
359 [advantage estimation](#).

360 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier  
361 Martinet, Marie-Anne Lachaux, Timothée Lacroix,  
362 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal  
363 Azhar, Aurélien Rodriguez, Armand Joulin, Edouard  
364 Grave, and Guillaume Lample. 2023. [LLaMA: Open](#)  
365 [and efficient foundation language models](#).

366 Yinjie Wang, Ling Yang, Bowen Li, Ye Tian, Ke Shen,  
367 and Mengdi Wang. 2025. [Revolutionizing reinforce-](#)  
368 [ment learning framework for diffusion large language](#)  
369 [models](#).

370 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,  
371 Thomas L. Griffiths, Yuan Cao, and Karthik  
372 Narasimhan. 2023. [Tree of thoughts: Deliberate](#)  
373 [problem solving with large language models](#). In *Ad-*  
374 *vances in Neural Information Processing Systems*.

375 Yao Yao, Zuchao Li, and Hai Zhao. 2024. [GoT: Effec-](#)  
376 [tive graph-of-thought reasoning in language models](#).  
377 In *Findings of the Association for Computational Lin-*  
378 *guistics: NAACL 2024*, pages 2901–2921, Mexico  
379 City, Mexico. Association for Computational Lin-  
380 guistics.

381 Jiacheng Ye, Zihui Xie, Lin Zheng, Jiahui Gao, Zirui  
382 Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong.  
383 2025. [Dream 7b: Diffusion large language models](#).

Base model	Method	MATH500 $\uparrow$	GSM8K $\uparrow$	AIME24 $\uparrow$
<b>SDAR-1.7B-Chat-b4</b>	–	61.40	80.54	4.44
	TRACERL	62.53	82.23	10.00
<b>SDAR-1.7B-Chat-b8</b>	–	55.90	81.00	5.56
	TRACERL	60.20	81.50	5.56
	<b>T*</b>	63.40	82.40	7.78
<b>SDAR-1.7B-Chat-b16</b>	–	52.40	79.40	3.33
	TRACERL	54.40	80.00	6.66
	<b>T*</b>	59.80	82.20	6.66
<b>SDAR-1.7B-Chat-b32</b>	–	54.20	78.31	2.22
	TRACERL	54.10	79.80	3.33
	<b>T*</b>	59.00	82.00	4.44
<b>SDAR-4B-Chat-b4</b>	–	68.67	90.50	5.56
	TRACERL	75.33	91.20	10.00
<b>SDAR-4B-Chat-b8</b>	–	60.73	85.87	5.56
	TRACERL	62.10	86.30	6.67
	<b>T*</b>	76.00	91.00	8.89
<b>SDAR-4B-Chat-b16</b>	–	58.26	78.24	6.67
	TRACERL	60.50	79.60	7.78
	<b>T*</b>	64.53	89.40	8.89

Table 2: **Reasoning performance under different block sizes.** “–” denotes the original SDAR--Chat-b $B$  checkpoint, TRACERL applies trajectory-aware RL at the same block size  $B$ , and **T\*** denotes our progressive block-size scaling.

## A.2 Case Study

**Question:** Find the product  $CD$  of the integers  $C$  and  $D$  for which  $\frac{C}{x-3} + \frac{D}{x+8} = \frac{4x-23}{x^2+5x-24}$  for all real values of  $x$  except  $-8$  and  $3$ .

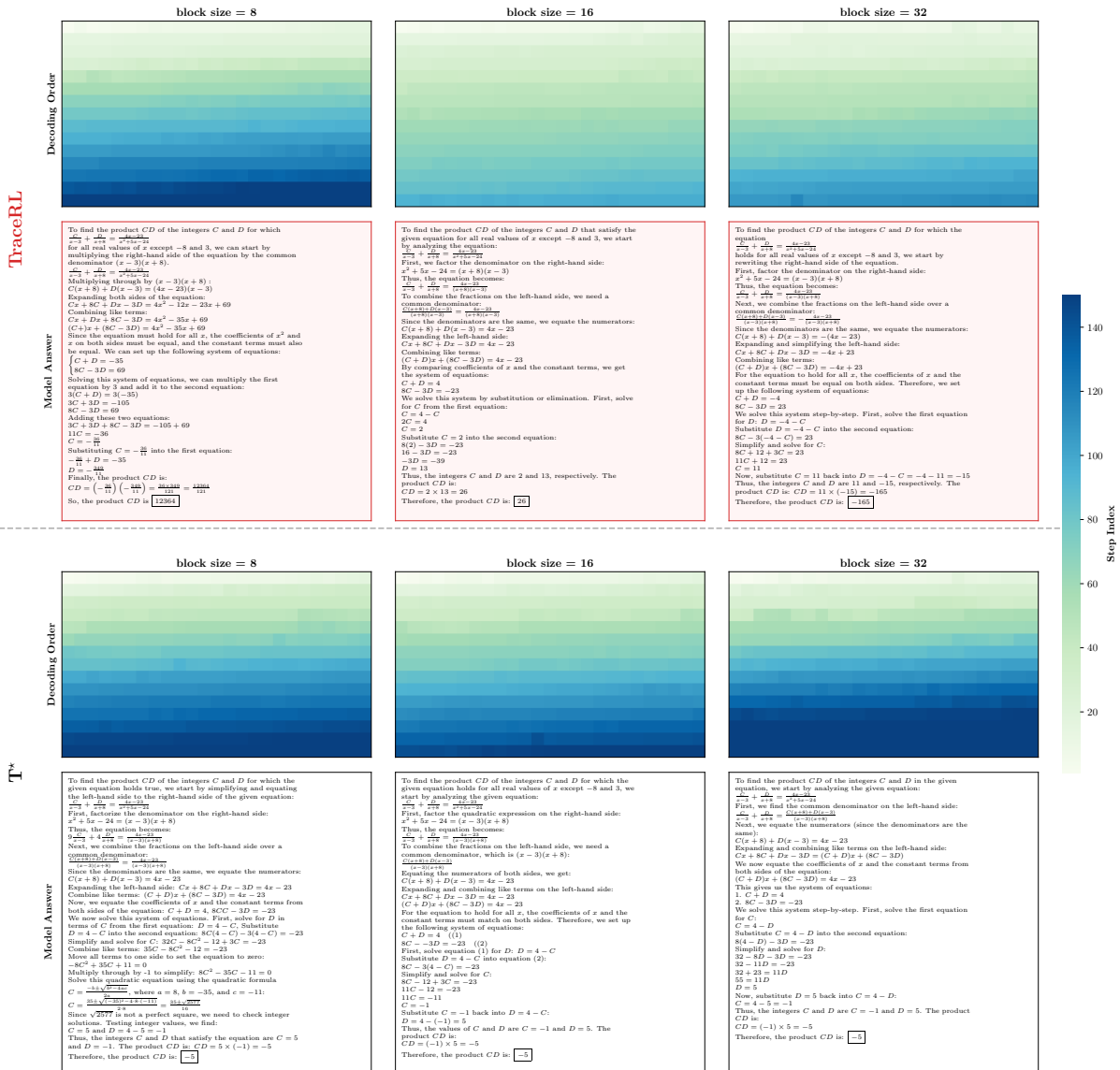


Figure 4: Case study: decoding schedule under TRACERL vs. T\*. We visualize the token-level first-unmask step index (heatmaps; darker means decoded later) and the corresponding model solutions for a representative algebra problem, evaluated with block sizes  $B \in \{8, 16, 32\}$ . The top row shows a model trained with direct TRACERL at the same block size, and the bottom row shows the model obtained by T\*.