

# DIP: Dynamic In-Context Planner For Diffusion Language Models

Anonymous ACL submission

## Abstract

Diffusion language models (DLMs) have shown strong potential for general natural language tasks with in-context examples. However, due to the bidirectional attention mechanism, DLMs incur substantial computational cost as context length increases. This work addresses this issue with a key discovery: unlike the sequential generation in autoregressive language models (ARLMs), the diffusion generation paradigm in DLMs allows *efficient dynamic adjustment of the context* during generation. Building on this insight, we propose **Dynamic In-Context Planner (DIP)**, a context-optimization method that dynamically selects and inserts in-context examples during generation, rather than providing all examples in the prompt upfront. Results show DIP maintains generation quality while achieving up to  $12.9\times$  inference speedup over standard inference and  $1.17\times$  over KV cache-enhanced inference.

## 1 Introduction & Related Work

Diffusion language models have emerged as a new paradigm for language modeling, demonstrating strong potential across general natural language tasks (Nie et al., 2025; Ye et al., 2025; Austin et al., 2021; DeepMind, 2025; Arriola et al., 2025; Lou et al., 2024; Sahoo et al., 2024; Lou et al., 2024). However, the DLM’s decoding mechanism is fundamentally different. By utilizing bidirectional attention and iterative refinement, DLMs remove the strict left-to-right constraint inherent to AR generation. This paradigm shift offers fine-grained control over the generation process, enabling flexible editing and global planning (Zhang et al., 2025; Jin et al., 2025; Li et al., 2025; Chen et al., 2025).

Similar to ARLMs, DLMs also benefit from In-Context Learning (ICL), in which providing in-context examples improves performance on many retrieval-augmented generation tasks (Gao et al., 2023). Existing ICL approaches for DLMs largely

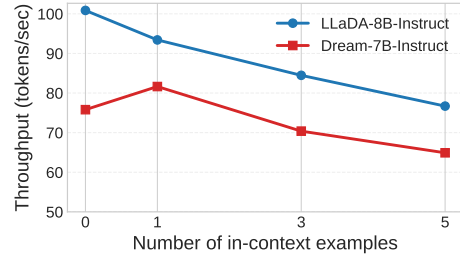


Figure 1: DLMs using Fast-dLLM (Wu et al., 2025) shows a negative correlation between throughput and the number of examples in the prompt.

inherit the AR convention (Radford et al., 2018; Bi et al., 2025), which keeps the prompt fixed and uses all in-context examples throughout generation (Radford et al., 2018; Dong et al., 2024). As shown in Figure 1, the static all-in-context example approach could lead to slower inference speed. Popular DLMs (Ye et al., 2025; Nie et al., 2025) utilize bidirectional attention, which scales poorly with context length. While recent optimizations, such as KV cache (Liu et al., 2025; Song et al., 2025), parallel decoding (Jiang et al., 2025; Wu et al., 2025), and block-wise decoding (Wu et al., 2025), can reduce the computation of static prompts, they do not resolve the fundamental inefficiency: the generative process still requires computing full attention over context during the generation process, which limits inference throughput (Wu et al., 2025).

This work addresses this issue with a **key discovery**: unlike the sequential generation, which requires a fixed context in ARLMs, the diffusion generation paradigm in DLMs is free from these constraints and allows *efficient dynamic adjustment of the context* during generation. Take the 3-shot GSM8k task as an example. In ARLMs, all examples must be provided at once in the prompt before generation begins. On the contrary, as illustrated in the middle column of Figure 2, the diffusion generation paradigm allows generation with 1 to 3

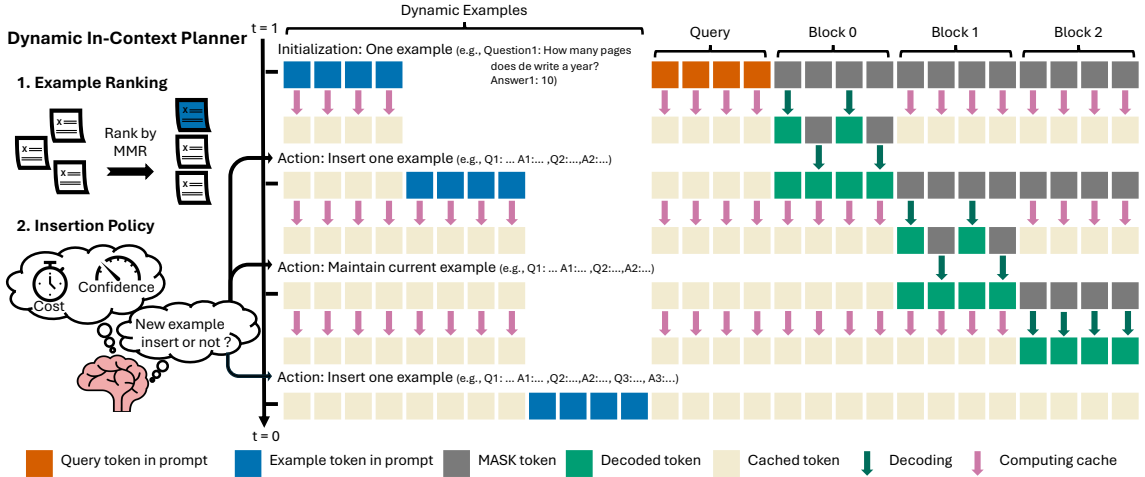


Figure 2: Dynamic In-Context Planner (DIP): (1) Example ranking stage uses MMR to rank the candidate examples by their marginal utility. (2) Insertion policy progressively inserts new examples into the context between blocks.

examples across different stages of the generation process (initialized with 1-shot at block 1, insert 1 example at block 2, ...). By dynamically selecting in-context examples based on the evolving generation state, we could increase inference throughput, ideally without compromising the generation quality.

Building on these insights, we propose **Dynamic In-Context Planner (DIP)**, a context-optimization method that dynamically selects in-context examples during generation. DIP consists of a Maximal Marginal Relevance (MMR) based example ranking stage and a confidence- and generation-process-aware example insertion policy that performs insertion decisions during generation. DIP is training-free and can serve as a plug-in method that efficiently selects examples and updates the prompt on top of prominent DLM frameworks such as Fast-dLLM (Wu et al., 2025).

Our main contributions include: (1) We make a key discovery that the diffusion generation paradigm of DLMs allows efficient dynamic context adjustment, opening a new avenue for efficient DLMs inference based on dynamic context optimization. (2) We propose DIP, a practical context-optimization method that dynamically updates in-context examples during generation, which consistently yields inference speedup. (3) Empirical results on the 5-shot GSM8k task show that DIP can maintain generation quality while achieving up to  $12.9\times$  inference speedup over LLaDA’s standard inference (Nie et al., 2025) and  $1.17\times$  over KV cache-enhanced inference like Fast-dLLM (Wu et al., 2025).

## 2 Masked DLMs with Dynamic Context

In this section, we briefly discuss the modeling details of DLMs that enable dynamic ICL with minimal overhead. Masked DLMs (the absorbing state discrete diffusion models) (Sahoo et al., 2024; Zhu et al., 2025; Ye et al., 2025) consist of a forward and reverse process, where an absorbing state is a mask token, denoted as [MASK], which represents missing or corrupted information (Lou et al., 2024; Nie et al., 2025).

In the *forward* process,  $q$  with an initial sequence  $y_0$  of length  $N$  is organized into discrete *time steps*, where each time step corresponds to a single masking stage. At each time step, independent tokens are gradually replaced by [MASK] tokens with probability  $t \in [0, 1]$ , which is defined as:

$$q(y_t|y_0) = \prod_{i=0}^{N-1} q(y_t^i|y_0^i), \quad (1)$$

where  $i \in [0, N]$  and  $y^i$  stands for the  $i$ -th token. Once a token transitions to [MASK], it remains in this state throughout the forward process. As  $t$  increases to 1, the original sequence will eventually become a fully masked sequence.

Conversely, in the *reverse* process  $p_\theta$ , at each time step, masked DLMs aim to predict the probability distribution over all masked tokens simultaneously given  $y_t$ , which can be formulated as:

$$p_\theta(y_s|y_t) = \prod_{i=0}^{N-1} p_\theta(y_s^i|y_t^i). \quad (2)$$

Based on these predictions, the additional sampling method  $S$  decides which tokens to "unmask" at each step. As  $t$  decreases to 0, the masked sequence will eventually recover the original sequence.

Masked DLMs usually employ a bidirectional attention mechanism (e.g., BERT (Devlin et al., 2019)) and scale poorly with context length, as shown in Figure 1. To address this, recent KV cache inference frameworks often use block diffusion, dividing the overall generation length into  $b$  equal-length discrete blocks and performing semi-AR decoding on the block level. For instance, Fast-dLLM (Wu et al., 2025) employs a KV cache for prefix and masked-suffix tokens in non-decoding-blocks, achieving significant speedups.

### 3 Problem Statement

We formally define the new dynamic example optimization task: instead of using all examples as context throughout DLM’s generation process, we aim to decide when and which examples to insert during generation to accelerate inference while maintaining inference quality.

**Definition 1 (ICL Setting)** Given a prompt  $p$  with a pool of in-context examples  $E_{pool} = \{(q_i, a_i)\}_{i=1}^K$ , where  $q_i$  and  $a_i$  denote the query and answer of the  $i$ -th example, respectively.

**Definition 2 (Efficient Inference Setting)** Given a KV cache-based inference process (e.g., Fast-dLLM (Wu et al., 2025)) divided into  $N$  non-KV cache steps (or blocks), indexed by  $b \in \{0, \dots, N - 1\}$ .

**Goal:** The goal is to find an optimal policy  $\pi$  that selects a subset of examples  $E_b \subseteq \{E_{pool}, \emptyset\}$  for each block  $b$ , which maintains the generation quality while improving the overall throughput.

### 4 Dynamic In-Context Planner

As shown in Figure 2, the Dynamic In-Context Planner (DIP) is a two-stage framework that decouples selection complexity into: (1) example ranking stage and (2) example insertion policy. The *ranking stage* orders the candidate examples in  $E_{pool}$  by their marginal utility. Then, an *example insertion policy* is a monotonic-increasing policy that progressively inserts these examples into the context at different generation blocks  $b$  based on generation confidence and generation stage. A formal algorithm is outlined in Appendix A.

#### 4.1 Block-wise Decoding with KV Cache

As DIP aims to optimize inference throughput, it is essential to ensure compatibility with the KV cache-enhanced inference framework. We build DIP on top of Fast-dLLM (Wu et al., 2025), a KV cache

enhanced block-wise inference framework, which requires a necessary non-KV cache call at the start of discrete block decoding to update the KV cache, providing an opportunity to efficiently adjust the prompt during generation (details in Algorithm 1 of Appendix A). This update usually introduces minimal overhead to the generation process.

In principle, DIP is a plug-in method that works with any masked DLM KV cache enhanced inference framework, whereas a non-KV-cache call is required to update the KV cache.

#### 4.2 Example Ranking

In the ranking stage, DIP aims to assign ranks to examples in  $E_{pool}$  based on their importance (Carbonell and Goldstein, 1998; Pickett et al., 2024; Nafee et al., 2025). Then we can use the top-related examples to guide our generation. To achieve this and balance the examples’ relevance and diversity with the test input, DIP deploys Maximal Marginal Relevance (MMR) on the query sentence embeddings to derive a ranked list.

$$\text{MMR}(q_i, \lambda) = \lambda \cdot \text{Sim}(q_p, q_i) - (1 - \lambda) \max_{j \in S} \text{Sim}(q_i, q_j), \quad (3)$$

where  $q_p$  denotes the representation of the actual query LLM aims to solve,  $S$  is the index set of already selected examples, and  $\lambda \in [0, 1]$  controls the trade-off between relevance and redundancy.

#### 4.3 Example Insertion Policy

Given ranked examples, DIP needs to decide when to insert the next-best example into the prompt and perform a prompt update. Intuitively, a good example-insertion policy should balance inference efficiency alone with generation confidence.

**Inference efficiency.** As illustrated in Figure 1, there is an inverse correlation between model inference throughput and the number of in-context examples. Since the cost of bidirectional attention scales poorly with sequence length, the policy should aim to minimize the context length to achieve greater inference gain. However, this is not always beneficial, as fewer examples can degrade performance in some cases. As the number of examples of our policy is non-decreasing, we will also gain greater inference speedup if we introduce new examples at later stages. Guided by those principles, we introduce a generation process penalty term  $G$  to scale the probability of inserting

an example as follows:

$$G(n, N, \epsilon) = (1 - \epsilon) + \epsilon \cdot (n/N), \quad (4)$$

where  $n$  stands for the current number of non-KV cache call,  $N$  stands for total number of non-KV cache call, and  $\epsilon \in [0, 1]$  is a hyperparameter controlling the ratio of penalty.

**Generation confidence guided.** While our goal is to improve inference speed, we should also maintain the same generation quality. It is essential to detect changes in generation quality during the decoding process. How can we effectively measure the shift in generation quality? The answer is the confidence of the generated tokens. Confidence is a critical signal for sampling methods to decide which tokens to unmask at each step, thus a suitable measure without additional cost.

DIP’s example insertion policy uses the average confidence across all generated tokens as a reference and compares it with the tokens generated in the current block or step. Then we define the probability for inserting a new example as the scaled ratio between two averages as follows:

$$P(\text{insert}) = \frac{1 - \mu}{2(1 - \bar{\mu})}, \quad (5)$$

where  $\mu$  stands for the average confidence of tokens from the current block and  $\bar{\mu}$  stands for the running average of all generated tokens.

To incorporate the time penalty, DIP’s final example insertion policy can be defined as a Bernoulli distribution as follows:

$$\begin{aligned} P(a = \text{insert}) &= P(\text{insert}) \cdot G(n, N, \epsilon) \\ P(a = \text{keep}) &= 1 - P(\text{insert}) \cdot G(n, N, \epsilon) \end{aligned} \quad (6)$$

## 5 Experiment

### 5.1 Experiment Setup

We evaluate DIP on the representative LLaDA-Instruct model, with a generation length of 256 and block size of 32, on the 5-shot GSM8k task (Cobbe et al., 2021). For embedding used in the ranking stage, we use Qwen3-Embedding-0.6B throughout the experiment. To ensure reproducibility, our experiment uses the *lm-eval* library (Gao et al., 2024). For metrics, we compare the flexible match accuracy provided by the *lm-eval* library along with throughput, measured as the average number of output tokens generated per second. For baseline, we selected Fast-dLLM (Wu et al., 2025) and decoding from LLaDA (Zhu et al., 2025). We run all experiments on a single NVIDIA 5090 GPU to ensure a fair, consistent comparison.

Method	Throughput	Acc
LLaDA	6.94 (1.0×)	78.2
Fast-dLLM	75.87 (10.9×)	77.7
DIP( $\lambda=0.1, \epsilon=0.2$ )	86.54 (12.5×)	77.8
DIP( $\lambda=0.1, \epsilon=1$ )	89.47 (12.9×)	76.7

Table 1: Results on 5-shot GSM8k with LLaDA-8B-Instruct (Zhu et al., 2025). Blue: relative speedup.

Parameters	Throughput	Accuracy
$\epsilon = 0.2, \lambda = 0.0$ ,	87.45	76.6
$\epsilon = 0.2, \lambda = 0.25$	86.49	77.6
$\epsilon = 0.2, \lambda = 0.5$	86.44	76.3
$\epsilon = 0.2, \lambda = 0.75$	86.68	76.5
$\epsilon = 0.2, \lambda = 1.0$	86.46	77.3
$\lambda = 0.1, \epsilon = 0.0$	85.54	76.9
$\lambda = 0.1, \epsilon = 0.25$	86.27	77.8
$\lambda = 0.1, \epsilon = 0.5$	87.51	77.0
$\lambda = 0.1, \epsilon = 0.75$	88.81	76.0
$\lambda = 0.1, \epsilon = 1.0$	89.47	76.7

Table 2: Ablation study on DIP’s insertion policy over choice of MMR parameter  $\lambda$  and generation process penalty parameter  $\epsilon$ .

## 5.2 Results

In Table 1, we present the accuracy and throughput using DIP. Results show that DIP outperforms LLaDA’s decoding and previous KV cache-only methods in terms of throughput. DIP can achieve up to 12.9× inference speedup compared to LLaDA’s decoding and 1.17× over Fast-dLLM, while maintaining on-par accuracy. In Table 2, we present the results over the choice of hyperparameter. Results show that: (1) increasing the MMR parameter  $\lambda$ , which emphasizes embedding similarity over redundancy, has little effect on throughput but slightly affects accuracy; (2) increasing the generation process penalty parameter  $\epsilon$ , which puts a larger penalty on early insertion, will increase the throughput but shows performance degradation.

## 6 Conclusion

In this paper, we made a key discovery: the diffusion generation paradigm of DLMs enables efficient context adjustment. Inspired by this, we propose DIP, a context optimization method that dynamically updates context examples during generation. Results show that DIP can maintain generation quality while achieving up to 12.9× inference speedup compared to standard inference and 1.17× compared to KV-cache enhanced inference.

## 301 Limitations

302 There are a few limitations to our work. First,  
303 due to computational constraints and time limita-  
304 tions, our experiment was conducted on the GSM8k  
305 benchmarks using LLaDA-8B-Instruct models. We  
306 encourage future work to explore DIP’s perfor-  
307 mance across more tasks and models. Second,  
308 DIP is built and tested on the Fast-dLLM inference  
309 framework. While theoretically DIP is compatible  
310 with standard KV cache methods, we encourage  
311 future work to explore different inference frame-  
312 works with DIP. Lastly, we adopted the inference  
313 setup from Fast-dLLMs. We encourage future work  
314 to explore different inference settings with dynamic  
315 context.

## 316 References

317 Marianne Arriola, Subham Sekhar Sahoo, Aaron  
318 Gokaslan, Zhihan Yang, Zhixuan Qi, Jiaqi Han,  
319 Justin T Chiu, and Volodymyr Kuleshov. 2025. Block  
320 diffusion: Interpolating between autoregressive and  
321 diffusion language models. In *The Thirteenth Inter-  
322 national Conference on Learning Representations*.

323 Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel  
324 Tarlow, and Rianne Van Den Berg. 2021. Structured  
325 denoising diffusion models in discrete state-spaces.  
326 *Advances in neural information processing systems*,  
327 34:17981–17993.

328 Zhenyu Bi, Meng Lu, Yang Li, Swastik Roy, Wei-  
329 jie Guan, Morteza Ziyadi, and Xuan Wang. 2025.  
330 Optagent: Optimizing multi-agent llm interactions  
331 through verbal reinforcement learning for enhanced  
332 reasoning. *arXiv preprint arXiv:2510.18032*.

333 Jaime Carbonell and Jade Goldstein. 1998. The use of  
334 mmr, diversity-based reranking for reordering doc-  
335 uments and producing summaries. In *Proceedings  
336 of the 21st annual international ACM SIGIR confer-  
337 ence on Research and development in information  
338 retrieval*, pages 335–336.

339 Zigeng Chen, Gongfan Fang, Xinyin Ma, Ruonan  
340 Yu, and Xinchao Wang. 2025. dparallel: Learn-  
341 able parallel decoding for dllms. *arXiv preprint  
342 arXiv:2509.26488*.

343 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
344 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
345 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
346 Nakano, and 1 others. 2021. Training verifiers  
347 to solve math word problems. *arXiv preprint  
348 arXiv:2110.14168*.

349 DeepMind. 2025. Gemini Diffusion — deep-  
350 mind.google. [Accessed 01-01-2026].

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
351 Kristina Toutanova. 2019. Bert: Pre-training of deep  
352 bidirectional transformers for language understand-  
353 ing. In *Proceedings of the 2019 conference of the  
354 North American chapter of the association for com-  
355 putational linguistics: human language technologies,  
356 volume 1 (long and short papers)*, pages 4171–4186.  
357

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan  
358 Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu,  
359 Baobao Chang, and 1 others. 2024. A survey on  
360 in-context learning. In *Proceedings of the 2024 con-  
361 ference on empirical methods in natural language  
362 processing*, pages 1107–1128.  
363

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Bider-  
364 man, Sid Black, Anthony DiPofi, Charles Foster,  
365 Laurence Golding, Jeffrey Hsu, Alain Le Noac’h,  
366 Haonan Li, Kyle McDonell, Niklas Muennighoff,  
367 Chris Ociepa, Jason Phang, Laria Reynolds, Hailey  
368 Schoelkopf, Aviya Skowron, Lintang Sutawika, and  
369 5 others. 2024. The language model evaluation har-  
370 ness.  
371

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,  
372 Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen  
373 Wang, and Haofen Wang. 2023. Retrieval-augmented  
374 generation for large language models: A survey.  
375 *arXiv preprint arXiv:2312.10997*, 2(1).  
376

Yuchu Jiang, Yue Cai, Xiangzhong Luo, Jiale Fu, Jiarui  
377 Wang, Chonghan Liu, and Xu Yang. 2025. d<sup>2</sup> cache:  
378 Accelerating diffusion-based llms via dual adaptive  
379 caching. *arXiv preprint arXiv:2509.23094*.  
380

Xiangqi Jin, Yuxuan Wang, Yifeng Gao, Zichen Wen,  
381 Biqing Qi, Dongrui Liu, and Linfeng Zhang. 2025.  
382 Thinking inside the mask: In-place prompting in  
383 diffusion llms. *arXiv preprint arXiv:2508.10736*.  
384

Pengxiang Li, Yefan Zhou, Dilxat Muhtar, Lu Yin,  
385 Shilin Yan, Li Shen, Yi Liang, Soroush Vosoughi,  
386 and Shiwei Liu. 2025. Diffusion language models  
387 know the answer before decoding. In *NeurIPS 2025  
388 Workshop on Efficient Reasoning*.  
389

Zhiyuan Liu, Yicun Yang, Yaojie Zhang, Junjie Chen,  
390 Chang Zou, Qingyuan Wei, Shaobo Wang, and Lin-  
391 feng Zhang. 2025. dllm-cache: Accelerating diffu-  
392 sion large language models with adaptive caching.  
393 *arXiv preprint arXiv:2506.06295*.  
394

Aaron Lou, Chenlin Meng, and Stefano Ermon. 2024.  
395 Discrete diffusion modeling by estimating the ratios  
396 of the data distribution. In *Proceedings of the 41st In-  
397 ternational Conference on Machine Learning*, pages  
398 32819–32848.  
399

Mahmud Wasif Nafee, Maiqi Jiang, Haipeng Chen, and  
400 Yanfu Zhang. 2025. Dynamic retriever for in-context  
401 knowledge editing via policy optimization. In *Pro-  
402 ceedings of the 2025 Conference on Empirical Meth-  
403 ods in Natural Language Processing*, pages 16755–  
404 16768.  
405

406 Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang,  
407 Jingyang Ou, Jun Hu, JUN ZHOU, Yankai Lin, Ji-  
408 Rong Wen, and Chongxuan Li. 2025. Large language  
409 diffusion models. In *The Thirty-ninth Annual Con-  
410 ference on Neural Information Processing Systems*.

411 Marc Pickett, Jeremy Hartman, Ayan Kumar Bhowmick,  
412 Raquib-ul Alam, and Aditya Vempaty. 2024. Better  
413 rag using relevant information gain. *arXiv preprint  
414 arXiv:2407.12101*.

415 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya  
416 Sutskever, and 1 others. 2018. Improving language  
417 understanding by generative pre-training.

418 Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron  
419 Gokaslan, Edgar Marroquin, Justin Chiu, Alexan-  
420 der Rush, and Volodymyr Kuleshov. 2024. Simple  
421 and effective masked diffusion language models. *Ad-  
422 vances in Neural Information Processing Systems*,  
423 37:130136–130184.

424 Yuerong Song, Xiaoran Liu, Ruixiao Li, Zhigeng Liu,  
425 Zengfeng Huang, Qipeng Guo, Ziwei He, and Xipeng  
426 Qiu. 2025. Sparse-dllm: Accelerating diffusion  
427 llms with dynamic cache eviction. *arXiv preprint  
428 arXiv:2508.02558*.

429 Chengyue Wu, Hao Zhang, Shuchen Xue, Zhijian Liu,  
430 Shizhe Diao, Ligeng Zhu, Ping Luo, Song Han, and  
431 Enze Xie. 2025. Fast-dllm: Training-free accelera-  
432 tion of diffusion llm by enabling kv cache and parallel  
433 decoding. *arXiv preprint arXiv:2505.22618*.

434 Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui  
435 Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong.  
436 2025. Dream 7b: Diffusion large language models.  
437 *arXiv preprint arXiv:2508.15487*.

438 Lingzhe Zhang, Liancheng Fang, Chiming Duan,  
439 Minghua He, Leyi Pan, Pei Xiao, Shiyu Huang, Yun-  
440 peng Zhai, Xuming Hu, Philip S Yu, and 1 others.  
441 2025. A survey on parallel text generation: From par-  
442 allel decoding to diffusion language models. *arXiv  
443 preprint arXiv:2508.08712*.

444 Fengqi Zhu, Rongzhen Wang, Shen Nie, Xiaolu Zhang,  
445 Chunwei Wu, Jun Hu, Jun Zhou, Jianfei Chen,  
446 Yankai Lin, Ji-Rong Wen, and 1 others. 2025. Llada  
447 1.5: Variance-reduced preference optimization for  
448 large language diffusion models. *arXiv preprint  
449 arXiv:2505.19223*.

450 **A Appendix: Algorithm for Fast-dLLM**  
451 **and DIP**

452 Algorithm 1 presents the decoding process of Fast-  
453 dLLM (Wu et al., 2025). Building on top of Fast-  
454 dLLM, we present the algorithm of DIP in Algo-  
455 rithm 2.

456 **B Appendix: Use Of Ai Assistants**

457 In this paper, we used Gemini solely to help with  
458 grammar and fluency. We also use Copilot for  
459 syntax assistance.

460 **C Appendix: Code Implementation**

461 Please refer to the supplementary materials for the  
462 code implementation.

---

**Algorithm 1: Fast-dLLM Decoding with Threshold**

---

**Input:** model  $p_\theta$ , prompt  $p$ , answer length  $L$ , blocks  $N$ , block size  $B$ , step per block  $T$ , threshold  $\tau$

- 1  $x \leftarrow [p; [\text{MASK}], \dots, [\text{MASK}]]$
- 2 Initialize KV Cache for  $x$ . ▷▷▷ KV Cache Init
- 3 **for**  $n = 1$  to  $N$  **do**
- 4      $s \leftarrow [ |p| + (n - 1)B ], e = |p| + nB$
- 5     **Update KV Cache with non-KV cache call** ▷▷▷ KV Cache Update
- 6     **for**  $t = 1$  to  $T$  **do**
- 7         Use cache, call  $p_\theta$  on  $x^{[s,e]}$  ▷▷▷ KV Cache Reuse
- 8         For each masked  $x^i$ , confidence  $c^i = \max_x p_\theta(x^i)$
- 9         Unmask all  $i$  in  $[s, e)$  with  $c^i \geq \tau$ , with index  $\max c^i$
- 10         **if all**  $x^{[s,e]} \neq [\text{MASK}]$  **then**
- 11             **break, move to next block**
- 12 **return**  $x$ ;

---

---

**Algorithm 2: DIP with Fast-dLLM**

---

**Input:** model  $p_\theta$ , answer length  $L$ , blocks  $N$ , block size  $B$ , step per block  $T$ , threshold  $\tau$ , example insertion policy  $\pi$ , examples set  $E$ , number of examples  $K$ , query  $q$

- 1 Rank  $E$  by MMR,  $E_{rank} = [E_1, E_2, \dots, E_K]$  ▷▷▷ **Example Ranking Stage**
- 2  $p \leftarrow [E_1; q], k = 1$  ▷▷▷ **Initialize Prompt With  $\pi$**
- 3  $x \leftarrow [p; [\text{MASK}], \dots, [\text{MASK}]]$
- 4 Initialize KV Cache for  $x$ . ▷▷▷ KV Cache Init
- 5  $\bar{\mu} = 0$  ▷▷▷ Running Confidence Mean
- 6  $\mu = 0$  ▷▷▷ Confidence Mean Of Current Block
- 7 **for**  $n = 1$  to  $N$  **do**
- 8      $s \leftarrow [ |p| + (n - 1)B ], e = |p| + nB$
- 9     **if**  $n! = 1$  **then**
- 10         Insertion =  $\pi(\mu, \bar{\mu})$  ▷▷▷ **Get Action Following  $\pi$**
- 11         **if** Insertion == *True* **then**
- 12              $p \leftarrow [E_1, E_2, \dots, E_k; q], k += 1$  ▷▷▷ **Update Prompt**
- 13              $lp = \text{length of } p$
- 14              $x \leftarrow [p; [\text{MASK}], \dots, [\text{MASK}]]$
- 15         **Update KV Cache with non-KV cache call** ▷▷▷ KV Cache Update
- 16         **for**  $t = 1$  to  $T$  **do**
- 17             Use cache, call  $p_\theta$  on  $x^{[s,e]}$  ▷▷▷ KV Cache Reuse
- 18             For each masked  $x^i$ , confidence  $c^i = \max_x p_\theta(x^i)$
- 19             Unmask all  $i$  in  $[s, e)$  with  $c^i \geq \tau$ , with index  $\max c^i$
- 20             **if all**  $x^{[s,e]} \neq [\text{MASK}]$  **then**
- 21                  $\mu = \text{mean}(c^i \text{ of } x^{[s,e]})$  ▷▷▷ Save Confidence Mean Of Current Block
- 22                  $\bar{\mu} = \text{mean}(c^i \text{ of } x^{[lp,e]})$  ▷▷▷ Save Confidence Mean Of Generated Tokens
- 23                 **break, move to next block**
- 24              $\mu = \text{mean}(c^i \text{ of } x^{[s,e]})$  ▷▷▷ Save Confidence Mean Of Current Block
- 25              $\bar{\mu} = \text{mean}(c^i \text{ of } x^{[lp,e]})$  ▷▷▷ Save Confidence Mean Of Generated Tokens
- 26 **return**  $x$ ;

---