
One-Pass to Reason: Token Duplication and Block-Sparse Mask for Efficient Fine-Tuning on Multi-Turn Reasoning

| | | |
|--|---|--|
| Ritesh Goru DevRev Austin, USA ritesh.goru@devrev.ai | Shanay Mehta DevRev Bengaluru, India shanay.mehta@devrev.ai | Prateek Jain DevRev Austin, USA prateek.jain@devrev.ai |
|--|---|--|

Abstract

Fine-tuning Large Language Models (LLMs) on multi-turn reasoning datasets requires N (number of turns) separate forward passes per conversation due to reasoning token visibility constraints, as reasoning tokens for a turn are discarded in subsequent turns. We propose duplicating response tokens along with a custom attention mask to enable single-pass processing of entire conversations. We prove our method produces identical losses to the N -pass approach while reducing time complexity from $O(N^3)$ to $O(N^2)$ and maintaining the same memory complexity for a transformer based model. Our implementation is available online¹.

1 Introduction

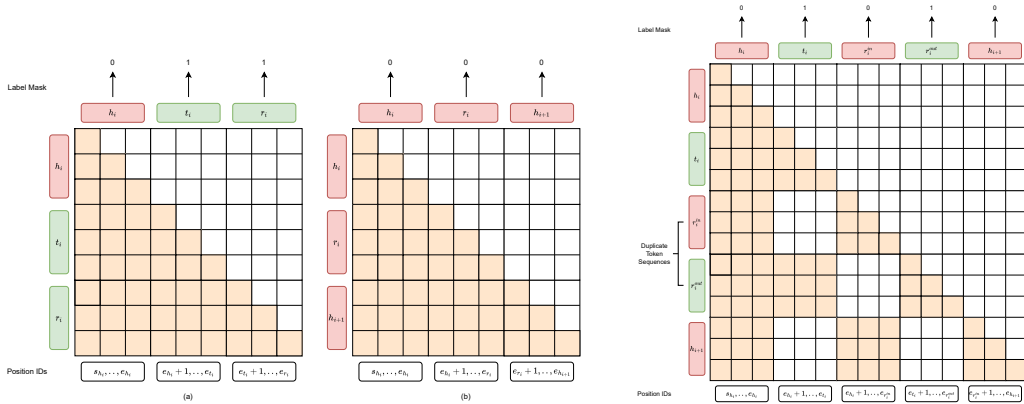
Recent advances in large language models (LLMs) have catalyzed a shift toward reasoning models that generate intermediate steps before producing final responses. While numerous studies fine-tune LLMs for multi-turn dialogues [15, 11], these approaches assume non-reasoning dialogues.

Training LLMs for multi-turn reasoning conversations presents two fundamental constraints: **Visibility Constraints**: Following industry-standard practices for multi-turn conversations [8, 1], reasoning tokens must be visible during generation but hidden from subsequent conversation turns, requiring conditional visibility that static attention masks cannot satisfy. (2) **Position ID Discrepancy**: Response tokens follow reasoning tokens during generation but directly follow human messages in a later context, creating positional misalignment. While prior works have explored masking techniques and position ID assignment to control information flow [14, 5, 10], none address these constraints.

Our contributions. (1) We present a theoretical framework featuring a block-sparse visibility mask and strategic position ID assignment scheme that enables processing an entire multi-turn reasoning conversation in a single forward pass while maintaining training correctness (Theorem 2.1). (2) Due to the absence of a publicly available multi-turn reasoning dataset (to the best of our knowledge), we create and release a novel dataset, MathChat_{sync}Reasoning, in which each assistant message is augmented with synthetically generated reasoning. (3) We provide comprehensive empirical validation for the proposed framework on Qwen3 models.

Notation. We use \mathcal{D} to denote a multi-turn reasoning dataset where each conversation $c \in \mathcal{D}$ consists of alternating human messages h_i and assistant messages a_i such that $c = (h_i, a_i)_{i=1}^N$ for N turns. Each assistant message a_i comprises thinking tokens t_i and response tokens r_i . We denote $\mathcal{H}_{<i} = (h_j, r_j)_{j=1}^{i-1}$ as conversation history before turn i . For token sequence x , s_x , and e_x represent starting and ending position IDs. The notation $x \rightarrow \mathcal{A}(\cdot)$ indicates sequences that x attends to, and $\mathcal{L}(\cdot)$ denotes language modeling loss (detailed in Appendix A.1).

¹<https://github.com/devrev/One-Pass-to-Reason>



(i) N-Pass causal masks; \square denotes non-zero attention. Panels: (a) generating response, (b) response in context. (ii) Custom attention mask for 1-Pass approach; \square represents non-zero attention.

2 Single Pass Fine-tuning on Multi-Turn Reasoning

In this section, we highlight the challenges associated with fine-tuning on multi-turn reasoning data and present a single-pass training approach. In multi-turn reasoning data, response tokens r_i must attend to reasoning tokens t_i when generating a_i , yet reasoning tokens must be hidden in later turns $a_{j>i}$. As a result, a single static mask cannot support both conditions in a conversation within a single forward pass—a capability that is often feasible with non-reasoning datasets.

2.1 N-Pass Approach

A straightforward baseline performs a separate forward pass for each turn $(\mathcal{H}_{<i}, h_i, a_i)$. This is correct but inefficient: N turns $\Rightarrow N$ training examples per conversation. Fig. 1i illustrates the masks for generation and for later context.

2.2 1-Pass Approach

The primary challenge in applying a single forward pass during training due to discrepancy in the attention behavior of r_i can be illustrated as follows²:

$$r_i \rightarrow \begin{cases} \mathcal{A}(\mathcal{H}_{<i}, h_i, t_i) & \text{generation} \\ \mathcal{A}(\mathcal{H}_{<i}, h_i) & \text{context} \end{cases}$$

We can resolve this issue through the following steps:

Duplicate response tokens. Split r_i into r_i^{out} (used in generation; attends to t_i) and r_i^{in} (context only; does not attend to t_i).

Custom Attention Mask. Duplication of response tokens makes it possible to have a single attention mask that satisfies visibility constraints. We define a custom masking strategy for each type of token sequence $(h_i, t_i, r_i^{in}, r_i^{out})$, ensuring that each token only attends to the appropriate subsequence:

$$\begin{aligned} h_i &\rightarrow \mathcal{A}(\mathcal{H}_{<i}^{in}) & r_i^{in} &\rightarrow \mathcal{A}(\mathcal{H}_{<i}^{in}, h_i) \\ t_i &\rightarrow \mathcal{A}(\mathcal{H}_{<i}^{in}, h_i) & r_i^{out} &\rightarrow \mathcal{A}(\mathcal{H}_{<i}^{in}, h_i, t_i) \end{aligned}$$

Assigning Consistent Position IDs. After duplication of response tokens, we need to assign consistent position IDs to tokens to maintain the correct relative positions—as if multiple forward passes were performed for each turn in the conversation. The following assignment of the first

²For ease of understanding, we omit the detail that each token within a token sequence also attends to all its preceding tokens, which must be encoded in the attention mask.

position ID for each token sequence ensures the relative positions are correct and equivalent to N-Pass approach³:

$$s_{t_i} = s_{r_i^{in}} = e_{h_i} + 1 \quad s_{r_i^{out}} = e_{t_i} + 1 \quad s_{h_{i+1}} = e_{r_i^{in}} + 1$$

Label Mask. Include loss only on t_i and r_i^{out} :

$$h_i \leftarrow 0 \quad t_i \leftarrow 1 \quad r_i^{in} \leftarrow 0 \quad r_i^{out} \leftarrow 1$$

Fig. 1ii shows the combined mask with consistent positions and labels.

Theorem 2.1. *Consider a language model with output distributions determined solely by attention patterns, positional encodings, and input representation. For any conversation c as input to the model, the sum of the N-Pass language modeling losses is equivalent to the 1-Pass loss:*

$$\mathcal{L}^{1-Pass}(c) = \sum_{i=1}^N \mathcal{L}_i^{N-Pass}(\mathcal{H}_{<i}, h_i, a_i)$$

Proof is in the Appendix B.1

2.3 Complexity Analysis

We compare the computational complexity of our 1-Pass method against N-Pass approach for transformer-based models with hidden dimension d [13]. Table 1 summarizes the time and memory complexities for a conversation c , where ℓ denotes its characteristic turn length.

Table 1: Time and Memory Complexity for N-Pass and 1-Pass Approach

| | N-PASS | 1-PASS |
|--------|-------------------|-------------------|
| $T(c)$ | $O(N^3 \ell^2 d)$ | $O(N^2 \ell^2 d)$ |
| $M(c)$ | $O(N^2 \ell^2)$ | $O(N^2 \ell^2)$ |

The 1-Pass approach reduces asymptotic time complexity by one order in N while retaining the same asymptotic memory as N-pass (with a larger constant). Full derivations are in Appendix C.

3 Experiments

We evaluate our single-pass fine-tuning on Qwen-3 models (4B, 8B, 32B) with QLoRA [3] using a $8 \times H100$ instance (CUDA 12.8, PyTorch 2.7.0). Our method implemented in LLaMA-Factory [16] and benchmarked against multi-pass baselines. We use FlashAttention2 (FA2) [2] and FlexAttention [4] backends. Our 1-Pass method uses FlexAttention only, as FA2 doesn’t support custom attention mask. We compare our **1-Pass method** against the **N-Pass baseline**. Both are evaluated with and without sequence packing⁴ [6]. See Appendix D.3 for more details.

Dataset Creation. Addressing lack of a public multi-turn dataset with explicit per-turn reasoning, we introduce MathChat_{sync}Reasoning⁵, derived from MathChat_{sync} [7]. Assistant turns are augmented with explicit reasoning generated using gpt-4.1-mini, conditioned on dialogue history and current assistant response. Refer to Appendix D.1 for more details. All our experiments are conducted on this dataset.

3.1 Results:

Training Speedup. Figure 2a shows training speedups. Our 1-Pass method with packing(Flex-Pack-1-Pass) is $1.05 \times$, $1.21 \times$, and $1.22 \times$ faster than FA2-N-Pass baseline with packing(FA2-Pack-N-Pass)

³Position IDs are assigned sequentially based on the order of tokens within each sequence.

⁴We set the cutoff length to the maximum number of tokens in any datapoint in the dataset for all our experiments.

⁵<https://huggingface.co/datasets/devrev-research/MathChatSync-reasoning>

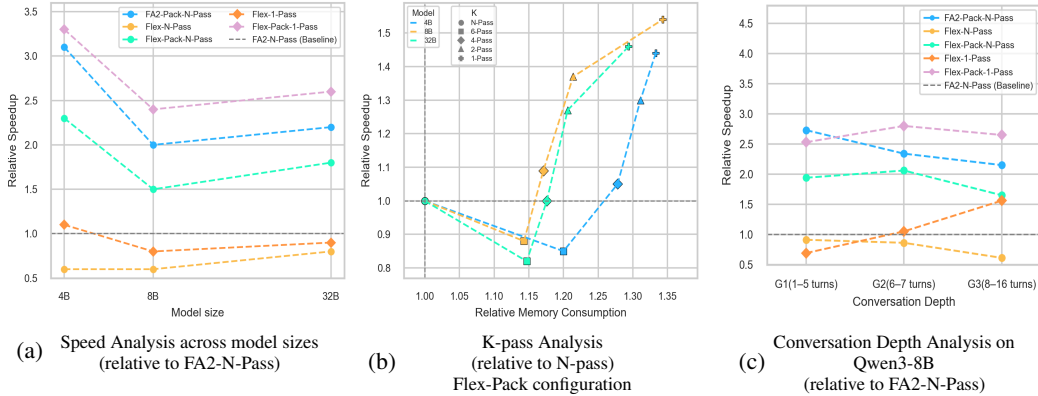


Figure 2: Training-time experiments

on 4B, 8B, and 32B models, respectively. Despite FlexAttention’s inherent slowness versus FA2, our method’s single-pass efficiency compensates. Compared to N-Pass FlexAttention with packing (Flex-Pack-N-Pass), our Flex-Pack-1-Pass yields $1.44\times$, $1.54\times$, and $1.46\times$ speedups for 4B, 8B, and 32B models, respectively. Without packing, our 1-pass method (Flex-1-Pass) lags FA2-N-Pass baseline for 8B and 32B models. We hypothesize that this is because response-token duplication widens the length disparity between conversations, making the method more sensitive to the absence of packing than the N-Pass baseline. Across all experiments, the 1-Pass variants consume roughly 33% more GPU memory than their N-Pass counterparts.

K-Pass Trade-offs. The 1-Pass and N-Pass approaches represent two extremes: processing the entire conversation in a single pass or in as many passes as there are turns. We also investigate intermediate settings, processing each conversation in K passes. We split every dialogue into K contiguous chunks and apply our single-pass mask only to the current chunk, duplicating response tokens and computing loss exclusively for that portion (see Appendix D.4.1 for full details). Figure 2b reveals a speed-memory trade-off for $K \in \{1, 2, 4, 6, N\}$. Our 1-Pass method maximizes speed with $\sim 33\%$ more memory (vs. N-Pass). $K=2$ offers a balance ($1.30\times$ – $1.37\times$ speedups, $\sim 20\%$ extra memory). Gains diminish for $K > 4$ because, beyond $K = 4$, the extra time incurred by the longer sequences created through token duplication outweighs the savings from processing a few turns together.

Conversation Scalability. The dataset contains conversations with depths from 1 to 16 turns. To analyse the effect of depth, we partition it into three groups: G1 (1–5 turns), G2 (6–7 turns), and G3 (8–16 turns)⁶. Figure 2c shows our Flex-Pack-1-Pass speedups (vs. FA2-Pack-N-Pass) grow with conversation depth ($0.93\times$, $1.19\times$, $1.23\times$ for G1, G2, G3 respectively). A similar trend appears when comparing our method without packing (Flex-1-Pass) to the FA2-N-Pass baseline: speedups of $0.69\times$, $1.05\times$, and $1.56\times$ for G1, G2, and G3, respectively. This supports the theoretical complexity reduction from $O(N^3)$ to $O(N^2)$, as efficiency gains become more pronounced with depth.

These results confirm the computational savings of single-pass training, making multi-turn reasoning fine-tuning practical at scale. See Appendix D.4 for full results.

4 Conclusion

We presented an optimized 1-Pass training method for multi-turn reasoning that reduces time complexity from $O(N^3)$ to $O(N^2)$ via strategic token duplication and custom attention mask. Our theoretical analysis confirms loss equivalence with the N-Pass method, enabling efficient training for longer conversations. As multi-turn reasoning becomes central to complex AI tasks, our method offers a scalable and broadly applicable solution. Future work includes exploring adaptive strategies to balance memory-efficiency trade-offs. Additionally, we aim to benchmark performance on latest back-ends such as FlashAttention3 [12] and port our masking logic to these faster implementations.

⁶This uneven distribution originates from the underlying MathChat_{sync} dataset, which is heavily skewed toward 5–7 turn conversations, a bias that propagates to our reasoning corpus.

References

- [1] Anthropic. Anthropic extended thinking. <https://docs.anthropic.com/en/docs/build-with-claude/extended-thinking>, 2025. Accessed: 2025-04-17.
- [2] Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In *Proceedings of the 12th International Conference on Learning Representations*, 2024.
- [3] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: efficient finetuning of quantized llms. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023.
- [4] Juechu Dong, Boyuan Feng, Driss Guessous, Yanbo Liang, and Horace He. Flex attention: A programming model for generating optimized attention kernels. *arXiv preprint arXiv:2412.05496*, 2024.
- [5] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. GLM: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335, Dublin, Ireland, May 2022.
- [6] Mario Michael Krell, Matej Kosec, Sergio P. Perez, and Andrew Fitzgibbon. Efficient sequence packing without cross-contamination: Accelerating large language models without impacting performance. *arXiv preprint arXiv:2107.02027*, 2022.
- [7] Zhenwen Liang, Dian Yu, Wenhao Yu, Wenlin Yao, Zhihan Zhang, Xiangliang Zhang, and Dong Yu. Mathchat: Benchmarking mathematical reasoning and instruction following in multi-turn interactions. *arXiv preprint arXiv:2405.19444*, 2024.
- [8] OpenAI. Openai reasoning. <https://platform.openai.com/docs/guides/reasoning>, 2024. Accessed: 2025-04-17.
- [9] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. Technical report, OpenAI, June 2018. URL https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf. Accessed: 10-11-2023.
- [10] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [11] Traian Rebedea, Makesh Sreedhar, Shaona Ghosh, Jiaqi Zeng, and Christopher Parisien. CantTalkAboutThis: Aligning language models to stay on topic in dialogues. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 12232–12252, Miami, Florida, USA, November 2024.
- [12] Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. Flashattention-3: Fast and accurate attention with asynchrony and low-precision. In *Advances in Neural Information Processing Systems*, 2024.
- [13] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [14] Franklin Wang and Sumanth Hegde. Accelerating direct preference optimization with prefix sharing. *arXiv preprint arXiv:2410.20305*, 2024.
- [15] Zezhong Wang, Kingshan Zeng, Weiwen Liu, Liangyou Li, Yasheng Wang, Lifeng Shang, Xin Jiang, Qun Liu, and Kam-Fai Wong. ToolFlow: Boosting LLM tool-calling through natural and coherent dialogue synthesis. 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 4246–4263, Albuquerque, New Mexico, April 2025.

- [16] Y. Zheng et al. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024.

A Background

A.1 Language Modeling Loss

For a token sequence $(\mathcal{H}_{<i}, h_i, a_i)$, the language modeling loss [9] for assistant message a_i can be expressed as:

$$\mathcal{L}(\mathcal{H}_{<i}, h_i, a_i) = -\log(P_{\Theta}(a_i | (\mathcal{H}_{<i}, h_i))) \quad (1)$$

where language model is parameterized by Θ .

B Proofs

B.1 Proof for Theorem 2.1

We establish the equivalence by demonstrating that both approaches yield identical probability distributions over sequences, which directly implies equal language modeling losses.

The proof proceeds in three parts: we show that (1) position encodings are equivalent, (2) attention patterns are identical, and (3) the resulting loss functions are mathematically equivalent.

Part I: Position Encoding Equivalence. Consider the position ID assignments for turn i as defined in Section 2.2. In the 1-Pass approach, output tokens receive positions:

$$\begin{aligned} s_{t_i} &= e_{h_i} + 1 \\ s_{r_i^{out}} &= e_{t_i} + 1 \end{aligned}$$

while input tokens from previous turns $j < i$ receive:

$$\begin{aligned} s_{r_j^{in}} &= e_{h_j} + 1 \\ s_{h_{j+1}} &= e_{r_j^{in}} + 1 \end{aligned}$$

This assignment ensures that tokens maintain the same relative positional relationships as in the N-Pass approach, where each turn processes tokens sequentially within separate forward passes.

Part II: Attention Pattern Preservation. The custom attention mask defined in Section 2.2 ensures causal dependencies are preserved. For turn i , the attention patterns are:

Output tokens:

$$\begin{aligned} t_i &\rightarrow \mathcal{A}(\mathcal{H}_{<i}^{in}, h_i) \\ r_i^{out} &\rightarrow \mathcal{A}(\mathcal{H}_{<i}^{in}, h_i, t_i) \end{aligned}$$

Input tokens from previous turns $j < i$:

$$\begin{aligned} h_j &\rightarrow \mathcal{A}(\mathcal{H}_{<j}^{in}) \\ r_j^{in} &\rightarrow \mathcal{A}(\mathcal{H}_{<j}^{in}, h_j) \end{aligned}$$

These patterns exactly replicate the causal attention available in the N-Pass approach.

Part III: Loss Function Equivalence. The language modeling loss for turn i in the N-Pass approach is:

$$\mathcal{L}_i^{N-Pass}(\mathcal{H}_{<i}, h_i, a_i) = -\log P_{\theta}(t_i, r_i | \mathcal{H}_{<i}, h_i) \quad (2)$$

By the autoregressive factorization:

$$\mathcal{L}_i^{N\text{-Pass}}(\mathcal{H}_{<i}, h_i, a_i) = -\log P_\theta(t_i | \mathcal{H}_{<i}, h_i) - \log P_\theta(r_i | \mathcal{H}_{<i}, h_i, t_i) \quad (3)$$

The total loss across all turns is:

$$\mathcal{L}^{N\text{-Pass}}(c) = \sum_{i=1}^N \mathcal{L}_i^{N\text{-Pass}}(\mathcal{H}_{<i}, h_i, a_i) \quad (4)$$

For the 1-Pass approach, the loss is computed as:

$$\mathcal{L}^{1\text{-Pass}}(c) = -\sum_{i=1}^N [\log P_\theta(t_i | \mathcal{H}_{<i}^{\text{in}}, h_i) + \log P_\theta(r_i^{\text{out}} | \mathcal{H}_{<i}^{\text{in}}, h_i, t_i)] \quad (5)$$

Key insight: Since $r_j = r_j^{\text{in}} = r_j^{\text{out}}$ (identical content in different positions) and the position encodings and attention patterns are equivalent as established in Parts I and II, the internal representations are identical. Therefore:

$$P_\theta(t_i | \mathcal{H}_{<i}, h_i) = P_\theta(t_i | \mathcal{H}_{<i}^{\text{in}}, h_i) \quad (6)$$

$$P_\theta(r_i | \mathcal{H}_{<i}, h_i, t_i) = P_\theta(r_i^{\text{out}} | \mathcal{H}_{<i}^{\text{in}}, h_i, t_i) \quad (7)$$

Combining equations (4), (5), (6), and (7):

$$\mathcal{L}^{N\text{-Pass}}(c) = \mathcal{L}^{1\text{-Pass}}(c) \quad (8)$$

C Complexity Analysis

C.1 Input Length

C.1.1 N-Pass Approach

In the N-Pass approach, each turn i is processed in a separate forward pass. The input to the model at turn i is:

$$\mathcal{H}_{<i}, h_i, t_i, r_i$$

because human and assistant response tokens from previous turns remain in the conversation history, while earlier reasoning tokens are discarded.

Let $L_{N\text{-Pass}}$ denote the maximum input length possible for the N-Pass approach for a conversation c . It can be defined by:

$$L_{N\text{-Pass}} = \sum_{i=1}^N (|h_i| + |r_i|) + \max_{i=1}^N |t_i|, \quad (9)$$

which is sum of all the human messages and response tokens for entire conversation and maximum length of thinking tokens across turns. To simplify further, assume:

$$|h_i|, |t_i|, |r_i| \in O(\ell).$$

where ℓ denote the characteristic turn component length, defined as $\ell = P_{95}(|h_i|, |t_i|, |r_i| : i \in [1, N], c \in \mathcal{D})$, where P_{95} is the 95th percentile operator. Then:

$$L_{N\text{-Pass}} \in O((2N + 1)\ell) = O(N\ell). \quad (10)$$

C.1.2 1-Pass Approach

Our 1-Pass approach processes the entire conversation c in a single forward pass. The input length $L_{1\text{-Pass}}$ can be calculated as:

$$L_{1\text{-Pass}} = \sum_{i=1}^N (|h_i| + |t_i| + 2|r_i|) \in O(4N\ell) = O(N\ell). \quad (11)$$

C.2 Time Complexity Analysis

For a transformer with hidden dimension d and context length n , each layer requires $O(n^2d)$ operations when $n \gg d$ [13].

N-Pass Approach: Under the N-Pass approach, each of the N turns requires a forward pass, each operating on $O(L_{N-Pass}) = O(N\ell)$ tokens. Thus, for conversation c :

$$T_{N-Pass}(c) \in O(N \times (N\ell)^2d) = O(N^3\ell^2d). \quad (12)$$

1-Pass Approach: In the 1-Pass approach, all the conversation tokens are given as input at once, thus operating on L_{1-Pass} tokens yielding a cost of:

$$T_{1-Pass}(c) \in O((4N\ell)^2d) = O(N^2\ell^2d). \quad (13)$$

This represents a factor of N improvement in asymptotic complexity, with substantial gains for large N .

C.3 Memory Complexity Analysis

A transformer layer with input context length n has memory complexity $O(n^2)$ assuming $n \gg d$.

N-Pass Approach: Peak Memory requirement for N-Pass approach is at L_{N-Pass} input. Thus for conversation c :

$$M_{N-Pass}(c) \in O((2N + 1)^2\ell^2) = O(N^2\ell^2). \quad (14)$$

1-Pass Approach: Memory requirement for 1-Pass approach can be given by:

$$M_{1-Pass}(c) \in O((4N)^2\ell^2) = O(N^2\ell^2). \quad (15)$$

Though 1-Pass incurs a higher constant factor due to response token replication, both approaches exhibit identical asymptotic memory complexity.

D Experiments

D.1 Dataset Creation

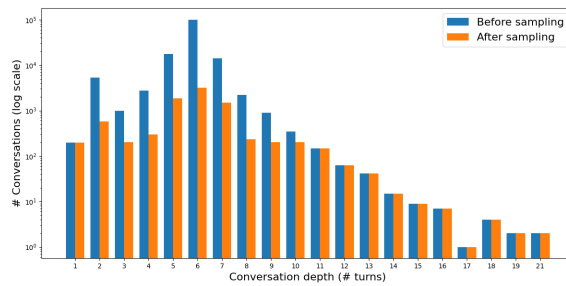


Figure 3: Dataset depth distribution: before vs. after sampling

To enable supervised training with explicit step-by-step reasoning, we construct and release MathChat_{sync} Reasoning along with its generation script. The dataset is obtained by augmenting the original MathChat_{sync} corpus [7] with a synthetically-generated rationale for every assistant turn. The procedure comprises three stages.

1. Source corpus. MathChat_{sync} is a synthetic, dialogue-based mathematics tutoring dataset containing 144,978 conversations with alternating human and assistant messages but no reasoning traces.

2. Depth-balanced sampling. Conversation depth in MathChat_{sync} is highly skewed toward six-turn dialogues (69 % of all conversations; see Figure 3). To mitigate this bias, we first down-sample depth-6 dialogues from 100,443 to 30,000 instances. From the resulting pool we draw a stratified sample of 8,000 conversations.

- For each depth d , we calculate the proportion of the pool that depth represents.
- We allocate to that depth the corresponding proportion of the 8,000-conversation budget, rounding up to the nearest whole conversation.
- If the resulting number is below 200, we raise it to (i) 200 or (ii) the total number of conversations available at that depth, whichever is smaller. This guarantees broad coverage across conversation depths.

The final split contains 8,797 assistant turns. Figure 3 compares the depth distribution before and after sampling.

3. Reasoning augmentation. For every assistant turn we generate an intermediate reasoning string using gpt-4.1-mini. The model is provided with (i) the dialogue history up to the current human utterance and (ii) the assistant’s reply, and is instructed to output only the hidden rationale that could have produced that reply. These rationales are concatenated to the original conversations to form MathChat_{sync}Reasoning.

D.2 Efficient mask generation

We present an efficient algorithm for generating the custom attention mask required by our 1-Pass training method. The algorithm leverages vectorized GPU operations to compute visibility patterns without explicit loops.

Algorithm 1 Efficient Custom Attention Mask Generation

Require: Role IDs tensor $\mathbf{R} \in \{0, 1, 2, 3, 4\}^{B \times L}$ where B is batch size, L is sequence length

Ensure: 4D attention mask $\mathbf{M} \in \mathbb{R}^{B \times 1 \times L \times L}$

```

1: // Step 1: Compute turn IDs via cumulative sum
2:  $\mathbf{R}_{\text{shift}} \leftarrow \text{roll}(\mathbf{R}, \text{shift} = 1, \text{dim} = 1)$ 
3:  $\mathbf{R}_{\text{shift}}[:, 0] \leftarrow 0$ 
4:  $\text{turn\_increment} \leftarrow (\mathbf{R} \neq 0) \wedge (\mathbf{R} = 1) \wedge (\mathbf{R}_{\text{shift}} \neq 1)$ 
5:  $\mathbf{T} \leftarrow \text{cumsum}(\text{turn\_increment}, \text{dim} = 1)$ 
6:  $\mathbf{T}[\mathbf{R} = 0] \leftarrow 0$  {Zero out padding positions}
7:
8: // Step 2: Create base causal non-padding mask
9:  $\mathbf{i} \leftarrow [0, 1, \dots, L - 1]$ 
10:  $\text{non\_pad} \leftarrow (\mathbf{R} \neq 0)$ 
11:  $\mathbf{M}_{\text{base}} \leftarrow (\mathbf{i}[:, \text{None}] \geq \mathbf{i}[\text{None}, :]) \wedge \text{non\_pad}[:, :, \text{None}] \wedge \text{non\_pad}[:, \text{None}, :]$ 
12:
13: // Step 3: Apply role-specific visibility constraints (K-map optimized)
14:  $\text{turn\_equal} \leftarrow (\mathbf{T}[:, :, \text{None}] = \mathbf{T}[:, \text{None}, :])$ 
15:  $\mathbf{R}_i \leftarrow \mathbf{R}[:, :, \text{None}]; \mathbf{R}_j \leftarrow \mathbf{R}[:, \text{None}, :]$ 
16:  $\mathbf{M}_{\text{final}} \leftarrow \mathbf{M}_{\text{base}} \wedge [(\mathbf{R}_j = 1) \vee (\mathbf{R}_j = 4 \wedge \text{turn\_equal})$ 
17:  $\vee (\mathbf{R}_j = 3 \wedge \mathbf{R}_i \neq 4) \vee (\mathbf{R}_j = 3 \wedge \neg \text{turn\_equal})$ 
18:  $\vee (\mathbf{R}_j = 2 \wedge \text{turn\_equal} \wedge \mathbf{R}_i \neq 3)]$ 
19:
20: // Step 4: Convert to 4D attention weights
21:  $\mathbf{M} \leftarrow \text{where}(\mathbf{M}_{\text{final}}.\text{unsqueeze}(1), 0, -\infty)$ 
22: return  $\mathbf{M}$ 

```

Implementation Notes:

- All operations are performed on GPU using PyTorch’s vectorized tensor operations
- Role IDs: 0 = padding, 1 = human, 2 = thinking, 3 = response (first copy), 4 = response (second copy)

- The boolean expression in Step 3 is optimized using Karnaugh map reduction to minimize logical operations
- The algorithm avoids explicit loops by leveraging broadcasting and logical operations
- For CPU tensors, we temporarily move computation to GPU before returning results to the original device

D.3 Experimental Setup

All training runs are initiated using llamafactory-cli in SFT mode. We apply QLoRA with 4-bit NF4 quantization, using a LoRA rank of 32 and a scaling factor of $\alpha = 64$. Training is performed for three epochs with bfloat16 (bf16) precision.

We enable the Liger kernel for improved efficiency. Each GPU processes a batch size of 2, with gradient accumulation over 4 steps. This setup yields an effective batch size of 64 across the 8-GPU node.

D.4 Comprehensive Results

We report the complete numerical results that support the figures in Section 3 in Tables 2, 3 and 4. We report two metrics for every configuration:

- **Throughput** (“samples per sec.”) — the average number of *full conversations* processed per second.
- **Peak GPU memory** — the peak memory recorded during training.

| Model Size | Run Setting | Samples per sec. | Peak Memory(GB) | Relative Speedup | Relative Peak Memory |
|------------|---------------------------|------------------|-----------------|------------------|----------------------|
| 4B | FA2-N-Pass(Baseline) | 1.985 | 9 | 1.0 | 1.00 |
| | FA2-Pack-N-Pass | 6.241 | 9 | 3.1 | 1.00 |
| | Flex Atten-N-Pass | 1.286 | 9 | 0.6 | 1.00 |
| | Flex Atten+Packing-N-Pass | 4.550 | 9 | 2.3 | 1.00 |
| | Flex-1-Pass | 2.107 | 12 | 1.1 | 1.33 |
| | Flex-Pack-1-Pass | 6.552 | 12 | 3.3 | 1.33 |
| 8B | FA2-N-Pass(Baseline) | 2.307 | 14 | 1.0 | 1.00 |
| | FA2-Pack-N-Pass | 4.522 | 14 | 2.0 | 1.00 |
| | Flex-N-Pass | 1.365 | 14 | 0.6 | 1.00 |
| | Flex-Packing-N-Pass | 3.561 | 14 | 1.5 | 1.00 |
| | Flex-1-Pass | 1.736 | 18.8 | 0.8 | 1.34 |
| | Flex-Pack-1-Pass | 5.484 | 18.8 | 2.4 | 1.34 |
| 32B | FA2-N-Pass(Baseline) | 0.601 | 34 | 1.0 | 1.00 |
| | FA2-Pack-N-Pass | 1.299 | 34 | 2.2 | 1.00 |
| | Flex-N-Pass | 0.465 | 34 | 0.8 | 1.00 |
| | Flex-Packing-N-Pass | 1.078 | 34 | 1.8 | 1.00 |
| | Flex-1-Pass | 0.521 | 44 | 0.9 | 1.29 |
| | Flex-Pack-1-Pass | 1.578 | 44 | 2.6 | 1.29 |

Table 2: **Throughput and peak memory across execution strategies.** FA2 = FlashAttention 2; Flex = FlexAttention. Pack denotes dynamic sequence-packing; “1-Pass” is our proposed approach. Relative columns are computed with respect to the corresponding FA2–N-Pass baseline.

D.4.1 Implementing K-Pass Processing

To obtain the results in Table 3 we extend our Optimised 1-Pass scheme to an intermediate *K-Pass* schedule. Assume a conversation contains N assistant turns $(h_1, t_1, r_1), \dots, (h_N, t_N, r_N)$.

- Chunking the dialog.** We partition the conversation into K contiguous chunks, each containing $\lceil N/K \rceil$ turns (the last chunk may be shorter).

| Model Size | K | Samples per sec. | Peak Memory(GB) | Relative Speedup | Relative Peak Memory |
|------------|------------------|------------------|-----------------|------------------|----------------------|
| 4B | N-Pass(baseline) | 4.55 | 9 | 1.00 | 1.00 |
| | 6-Pass | 3.89 | 10.8 | 0.85 | 1.20 |
| | 4-Pass | 4.76 | 11.5 | 1.05 | 1.28 |
| | 2-Pass | 5.91 | 11.8 | 1.30 | 1.31 |
| | 1-Pass | 6.55 | 12 | 1.44 | 1.33 |
| 8B | N-Pass(baseline) | 3.56 | 14 | 1.00 | 1.00 |
| | 6-Pass | 3.13 | 16 | 0.88 | 1.14 |
| | 4-Pass | 3.87 | 16.4 | 1.09 | 1.17 |
| | 2-Pass | 4.87 | 17 | 1.37 | 1.21 |
| | 1-Pass | 5.48 | 18.8 | 1.54 | 1.34 |
| 32B | N-Pass(baseline) | 1.08 | 34 | 1.00 | 1.00 |
| | 6-Pass | 0.88 | 39 | 0.82 | 1.15 |
| | 4-Pass | 1.08 | 40 | 1.00 | 1.18 |
| | 2-Pass | 1.37 | 41 | 1.27 | 1.21 |
| | 1-Pass | 1.58 | 44 | 1.46 | 1.29 |

Table 3: **Speed-memory trade-off as a function of K** . Each dialogue is split into K equal-length chunks that are processed sequentially in a *single* forward/backward pass. $K = N$ corresponds to the per-turn baseline, while $K = 1$ is our single-pass method. All experiments use the FlexAttention backend with sequence packing (Flex-Pack), the configuration that achieved the best overall speed in our primary evaluation.

| | Run Setting | Samples per sec. | Peak Memory(GB) | Relative Speedup | Relative Peak Memory |
|---------|----------------------|------------------|-----------------|------------------|----------------------|
| Group 1 | FA2-N-Pass(Baseline) | 2.54 | 14 | 1.00 | 1 |
| | FA2-Pack-N-Pass | 6.93 | 14 | 2.73 | 1 |
| | Flex-N-Pass | 2.32 | 14 | 0.91 | 1 |
| | Flex-Packing-N-Pass | 4.94 | 14 | 1.94 | 1 |
| | Flex-1-Pass | 1.74 | 18.8 | 0.69 | 1.34 |
| | Flex-Pack-1-Pass | 6.43 | 18.8 | 2.53 | 1.34 |
| Group 2 | FA2-N-Pass(Baseline) | 1.02 | 14 | 1 | 1 |
| | FA2-Pack-N-Pass | 2.39 | 14 | 2.34 | 1 |
| | Flex-N-Pass | 0.87 | 14 | 0.86 | 1 |
| | Flex-Packing-N-Pass | 2.10 | 14 | 2.06 | 1 |
| | Flex-1-Pass | 1.07 | 18.8 | 1.05 | 1.34 |
| | Flex-Pack-1-Pass | 2.86 | 18.8 | 2.80 | 1.34 |
| Group 3 | FA2-N-Pass(Baseline) | 1.06 | 14 | 1 | 1 |
| | FA2-Pack-N-Pass | 2.28 | 14 | 2.15 | 1 |
| | Flex-N-Pass | 0.65 | 14 | 0.61 | 1 |
| | Flex-Packing-N-Pass | 1.75 | 14 | 1.65 | 1 |
| | Flex-1-Pass | 1.66 | 18.8 | 1.56 | 1.34 |
| | Flex-Pack-1-Pass | 2.81 | 18.8 | 2.65 | 1.34 |

Table 4: **Impact of conversation depth (Qwen-3 8B)**. Group 1 (1–5 turns), Group 2 (6–7 turns), and Group 3 (8–16 turns). Our 1-Pass approach gains more speed as depth increases, in line with the theoretical $O(N^2)$ vs. $O(N^3)$ complexity gap.

- (b) **Selective token duplication.** Within the *current* chunk we apply the same response-token duplication as in Section 2.2: $r_i^{\text{in}}, r_i^{\text{out}}$. All earlier chunks act purely as context and therefore retain their original, non-duplicated responses. This progressively lowers the number of duplicated tokens as K increases, which is the main source of the memory savings reported in Table 3.
- (c) **Attention and position IDs.** The custom attention mask and position-ID assignment described in Section 2.2 are applied *only* to the duplicated tokens of the active chunk. Context tokens keep the standard causal mask.
- (d) **Loss computation.** The label mask is set to 1 for t_i and r_i^{out} *inside* the active chunk and 0 elsewhere, so each pass trains only on the new turns while reusing earlier content as fixed context.

Conceptually, the K -Pass schedule interpolates between the extremes:

- $K = N$ reproduces the per-turn baseline (no response duplication, minimal memory, maximal passes);
- $K = 1$ is our 1-Pass method (maximum duplication, single pass, fastest).