Published at ICLR 2025 Workshop on Reasoning and Planning for LLMs

# Thinking Slow, Fast: Scaling Inference Compute with Distilled Reasoners

Daniele Paliotta, Junxiong Wang, Metteo Pagliardini, Kevin Y. Li, Aviv Bick, J. Zico Kolter, Albert Gu, Francois Fleuret, Tri Dao

Chanhee Lee1, Yongjun Kim2, Jiwoo Kim1 1 Department of Electrical Engineering 2 Department of Computer Science



# Contents

- I. Introduction
- II. Background
- III. Method
- IV. Experimental Result

## V. Conclusion



### Introduction



# Improvement of recent LLMs inference performance

#### Problem of LLM

- Large Language model (LLM) gains outstanding performance in recent studies, but it has clear limitation on complex inference.
- Simply scaling up the model size is inefficient from latency, memory, computational complexity perspective.

🛑 anthropic 🛛 🛑 chinese 💛 google 🌑 meta 🔵 microsoft 🔵 mistral 🌑 openAl 🎯 other



# Improvement of recent LLM inference performance

#### Overcome the limitation

- Test-time compute is one of the solution of these problems.
- o1 model of OpenAI has proven usefulness of test-time compute
- It gives outstanding performance than original models.





[R1] GPT-o1, OpenAl, 2024



# Optimal architecture for test-time compute

#### Disadvantage of test-time compute

- In transformer architecture, High memory usage and slow inference, Due to their long sequence length.
- This is especially noticeable in transformer structures where a lot of cost is required in the inference step.



Figure 3 | The average response length of DeepSeek-R1-Zero on the training set during the RL process. DeepSeek-R1-Zero naturally learns to solve reasoning tasks with more thinking time.



# Optimal architecture for test-time compute

- Subquadratic models
  - Faster, less memory usage than transformer, which is foundation model of LLMs
  - Inefficient well-trained subquadratic models. → distillation
  - Mamba is well-known subquadratic model.



#### Inference throughput on A100 80GB (prompt length 2048)



## What the author do?

- Problems and solutions summary
  - Test time compute needs long context  $\rightarrow$  Use subquadratic model
- - Lack of well-trained subquadratic model  $\rightarrow$  Distill transformer to subquadratic model



# Background



# Chain of Thought (CoT)

#### Step-by-step inference

- Not compute the result directly, but implement the flow of thought
- Show good performance some specific tasks
  - Mathematical problem
  - Logical determination
  - Analysis of complex sentence





# Chain of Thought (CoT)



Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

[R1] Wei, Jason, et al. "Chain-of-thought prompting elicits reasoning in large language models." Advances in neural information processing systems 35 (2022): 24824-24837.



# Chain of Thought (CoT)

| Math Word Problems (free response)  | Math Word Problems (multiple choice)   | CSQA (commonsense)  |
|---|--|---|
| Q: Roger has 5 tennis balls. He buys<br>2 more cans of tennis balls. Each can<br>has 3 tennis balls. How many tennis<br>balls does he have now? | Q: How many keystrokes are needed<br>to type the numbers from 1 to 500?<br>Answer Choices: (a) 1156 (b) 1392 (c) 1480<br>(d) 1562 (e) 1788   | Q: Sammy wanted to go to where the<br>people were. Where might he go?<br>Options: (a) race track (b) populated areas<br>(c) desert (d) apartment (e) roadblock                |
| A: Roger started with 5 balls. 2 cans<br>of 3 tennis balls each is 6 tennis<br>balls. 5 + 6 = 11. The answer is 11.                             | A: There are 9 one-digit numbers<br>from 1 to 9. There are 90 two-digit<br>numbers from 10 to 99. There are<br>401 three-digit numbers from 100 to<br>500. $9 + 90(2) + 401(3) = 1392$ . The<br>answer is (b). | A: The answer must be a place with a lot of people. Race tracks, desert, apartments, and roadblocks don't have a lot of people, but populated areas do. So the answer is (b). |
| StrategyQA  | Date Understanding   | Sports Understanding  |
| Q: Yes or no: Would a pear sink in water?<br>A: The density of a pear is about 0.6  | Q: The concert was scheduled to be<br>on 06/01/1943, but was delayed by<br>one day to today. What is the date 10<br>days ago in MM/DD/YYYY?  | Q: Is the following sentence<br>plausible? "Joao Moutinho caught the<br>screen pass in the NFC<br>championship."  |
| g/cm^3, which is less than water.<br>Thus, a pear would float. So the<br>answer is no.  | A: One day after 06/01/1943 is<br>06/02/1943, so today is 06/02/1943.<br>10 days before today is 05/23/1943.<br>So the answer is 05/23/1943.   | A: Joao Moutinho is a soccer player.<br>The NFC championship is part of<br>American football, not soccer. So the<br>answer is no.   |
|   |  |   |
| SayCan (Instructing a robot)  | Last Letter Concatenation  | Coin Flip (state tracking)  |
| Human: How would you bring me<br>something that isn't a fruit?  | Q: Take the last letters of the words<br>in "Lady Gaga" and concatenate<br>them.   | Q: A coin is heads up. Maybelle flips<br>the coin. Shalonda does not flip the<br>coin. Is the coin still heads up?  |
| something to eat that isn't a fruit. An<br>energy bar is not a fruit, so I will bring<br>the user an energy bar.                                | A: The last letter of "Lady" is "y". The<br>last letter of "Gaga" is "a".<br>Concatenating them is "ya". So the  | A: The coin was flipped by Maybelle.<br>So the coin was flipped 1 time, which<br>is an odd number. The coin started   |
| Plan: 1. find(energy bar) 2.<br>pick(energy bar) 3. find(user) 4.<br>put(energy bar) 5. done().   | answer is ya.  | heads up, so after an odd number of flips, it will be tails up. So the answer is no.  |

#### Various prompt example

- Arithmetic (green)
- Common sense (orange)
- Symbolic reasoning (blue)

#### Features

- Step-by-steps response
- Generates longer sequence than original cases

Figure 3: Examples of  $\langle input, chain of thought, output \rangle$  triples for arithmetic, commonsense, and symbolic reasoning benchmarks. Chains of thought are highlighted. Full prompts in Appendix G.

[R1] Wei, Jason, et al. "Chain-of-thought prompting elicits reasoning in large language models." Advances in neural information processing systems 35 (2022): 24824-24837.



# Test time compute with CoT

#### Using more resource in inference step

- Generate different result of CoTs
- Select proper response various stratage
  - majority voting
  - Using learnable reward models for each CoTs
- Disadvantage
  - Generate longer CoT sequences than normal cases
    - $\rightarrow$  Require high cost and memory usage



 $\bigcirc$  = Rejected by verifier  $\bigcirc$  = Selected by verifier  $\diamondsuit$  = Intermediate step  $\bigcirc$  = Full solution



## Subquadratic model

### Transformer

- Self-attention mechanism, which is main reason of performance of Transformer architecture, refer all tokens in each steps.
- Its computational complexity is  $O(L^2)$  for input sequence of length L.
- Subquadratic computational complexity
  - Less than  $O(L^2)$  computational complexity for input sequence of length L
  - → Better than transformer to adopt Test-time compute



# Subquadratic model

- Mamba [R2]
  - One of the highlighted subquadratic architecture of LLM
  - State Space Model (SSM) based recurrence model





[R2] Gu, Albert, and Tri Dao. "Mamba: Linear-time sequence modeling with selective state spaces." *arXiv preprint arXiv:2312.00752* (2023). Thinking Slow, Fast: Scaling Inference Compute with Distilled Reasoners

## SSM block of Mamba

Update state vector and discretize





[R2] Gu, Albert, and Tri Dao. "Mamba: Linear-time sequence modeling with selective state spaces." arXiv preprint arXiv:2312.00752 (2023).

## SSM block of Mamba

- Remove Linear Time-Invariant (LTI) constraint to SSM
  - Matrix A contributes major performance of SSM, but it has LTI constraints.
  - Using High-order Polynomial Projection Operators (HiPPO)
  - Construct recent token more precisely than first token.



#### **HiPPO Matrix**



[R2] Gu, Albert, and Tri Dao. "Mamba: Linear-time sequence modeling with selective state spaces." arXiv preprint arXiv:2312.00752 (2023).

## SSM block of Mamba

- Selectiveness
  - Fusion of H3 and Gated MLP





[R2] Gu, Albert, and Tri Dao. "Mamba: Linear-time sequence modeling with selective state spaces." arXiv preprint arXiv:2312.00752 (2023).

# Distillation

### Definition

- Transfer the knowledge large model to small model.
- Efficient compression of LLMs

#### Cross-architecture distillation

- Transformer to RNN[R3], Linear Attention[R4], SSM[R5] etc.
- In this paper, the author distill the knowledge transformer to mamba.



[R3] Kasai, Jungo, et al. "Finetuning pretrained transformers into rnns." arXiv preprint arXiv:2103.13076 (2021).
[R4] Zhang, Michael, et al. "The Hedgehog & the Porcupine: Expressive Linear Attentions with Softmax Mimicry." ICLR, 2024
[R5] Wang, J., Paliotta, D., May, A., Rush, A., and Dao, T. "The mamba in the llama: Distilling and accelerating hybrid models." NeurIPS, 2025

### Method



### MOHAWK: Distill Llama to Mamba

MOHAWK is composed of three stages:





# Sequence transformation & matrix mixer

#### Definition 1: Sequence transformation

- Sequence transformation refers to parameterized map on some sequence
- i.e. sequence transformation combine tokens at various time steps

### Definition 2: Matrix mixer

- Some sequence transformation can be represented by Y = MX
- In this case, *M* is called as matrix mixer
- In a attention mechanism, Softmax(QK<sup>T</sup>) is a matrix mixer





## Causal variants of sequence transformation

#### Definition 3: Causal variants of sequence transformation

 By multiplying a lower triangular matrix filled with 1s (L), we can get casual variants of sequence transformation. That is, if we multiply L to the attention matrix, we can get causal variants of attention matrix.





## Relation between Mamba-2 & causal linear attention

Mamba-2 is a time-varying state-space model, and defined as follow:

$$h_{t+1} = \mathbf{A}_t h_t + \mathbf{B}_t x_t$$
$$y_t = \mathbf{C}_t h_t$$

 Fixing A<sub>t</sub> = I results in the formulation of causal linear attention with the matrices B, C representing the projections of the key and the query, respectively.



# Why Mamba-2 can be the causal linear attention?

• Mamba-2 is a time-varying state-space model, and defined as follow:

$$egin{aligned} h_t &= a_t\,h_{t-1} + B_t\,x_t,\ y_t &= C_t^ op h_t. \end{aligned}$$

• If we unroll the state update with  $h_0 = 0$ , we can get ...

$$h_t = a_t h_{t-1} + B_t x_t = \sum_{i=1}^t \left(\prod_{j=i+1}^t a_j\right) B_i x_i.$$
  
Same function with causal masking



## Why Mamba-2 can be the causal linear attention?

In linear attention, by using kernel-mapping, the output is calculated as follow :

$$y_t \;=\; \phi(q_t)^ op \sum_{i=1}^t \phi(k_i) \, v_i$$

• By computing  $y_t = C_t^{\top} h_t$  with previously unrolled results, we can get

$$y_t = C_t^ op h_t = \sum_{i=1}^t \Big[ C_t^ op (B_i \, x_i) \Big] \prod_{j=i+1}^t a_j$$



# Why Mamba-2 can be the causal linear attention?

 In summarize, attention mechanism's query, key, value matrices can be mapped into Mamba-2 SSM as below:

| Attention mechanism                             | Mamba-2 SSM                        |  |
|---|------------------------------------|--|
| kernelized $key_i = \phi(W_k x_i)$              | $B_i x_i$                          |  |
| $value_i = W_v x_i$                             | (Intrinsically included in $B_i$ ) |  |
| kernelized query <sub>t</sub> = $\phi(W_q x_t)$ | $C_t$                              |  |

#### Please refer to Katharopoulos et al. & Bick et al. for more detailed explanation



# Distillation Llama into Llamba

 Goal: Distillation Llama to Llamba, which means Llama model replacing its attention to Mamba-2 SSMs.





## MOHAWK Stage 1: Matrix Orientation

Goal: Align each student SSM block's SSM to its teacher's attention matrix.

$$egin{aligned} M_\ell^{(T)} &= ext{softmax}(Q_\ell K_\ell^ op) & y_t = C_t^ op \sum_{i=1}^t ig(\prod_{j=i+1}^t A_jig) B_i \, x_i = \sum_{i=1}^t igC_t^ op ig(\prod_{j=i+1}^t A_jig) B_i \, x_i. \ & \underbrace{\mathcal{L}_{ ext{matrix}} = ig\| M_\ell^{(T)}(u_\ell) \, - \, M_\ell^{(S)}(u_\ell) ig\|_F^2} \end{aligned}$$

 Key insight: Matching mixing matrices first, to ensure the student mirrors the teacher's long-range range information flow before any hidden-state alignment



## Distillation Llama into Llamba

Goal: Align each student SSM block's SSM to its teacher's attention matrix.





## MOHAWK Stage 2: Hidden-State Alignment

 Goal: Bring the student's internal representations in each mixer block into close agreement with the teacher's output of attention block.





# MOHAWK Stage 3: Weight transfer & End-to-End KD

- Goal: Finalize the student by
- (a) inheriting compatible teacher weights
- (b) training on final outputs so that its predictions match the teacher's over a small corpus.







## Experimental Results Overview

#### Authors tried to show..

(1) Inference speedup of distilled models is better!

(2) This speedup can result in better scaling for a given inference time budget!



## Inference Time Results

#### Experiment protocols

- Dataset: MATH and GSM8K
  - Realistic setup that matches the prompt and CoTs length
- Varying batch size, tasks of generating 512 tokens from a prompt with 512 tokens
- Prefilling time is not included.
  - It depends on the given prompt
  - Only interested in the time to generate multiple completions given one prompt
- Done on a single NVIDIA H100 GPU

**GSM8K System Prompt:**  $\n\$ *or Given the following problem, reason and give a final answer to the problem.* $\n$ *Your response should end with "The final answer is [answer]" where [answer] is the response to the problem.* $\n$ *Problem:* 



### Inference Time Results

#### Faster generation of distilled models

- Distilled models were faster than Llama baselines.
- MambaInLlama models are slightly faster than Llamba
  - Smaller SSM state size: MambalnLlama(16) < Llamba(64)</p>





#### Experiment Protocol

- Teacher models: Llama-3.2-1B-Instruct & 3B-Instruct
- Distilled Mamba Students: MambaInLlama-1B & 3B, Llamaba-1B & 4B
- 500 sample subset of MATH and GSM8K
- Evaluated coverage and accuracy



#### Experiment Protocol

- Coverage vs Accuracy
  - Coverage: Probability that the generated set contains the correct answer(upper bound)
  - Accuracy: Probability that the selected answer is correct (final output quality)



#### Good Coverage





- **Distilled Models can Cover like Teachers** 
  - "Does distilled model can generate meaningful results?"
    - Previous distilled models: could not achieve similar coverage with teachers...
  - *Time budget: faster models can generate more candidates*Observed scaling of coverage as:
  - - (a) a function of time budget, (b) the number of generation k increases





(a) Scaling with time.

(b) Scaling with number of completions.

- Distilled Models Achieve Competitive Accuracy Under Fixed Time
  - Selection methods
    - Majority voting: Selects the most frequently occurring final answer
    - Weighted Best-of-N: Selects the answer with the highest score based on a reward model or evaluation function. (In this paper, Llama-3.1-8B-based reward model)



## **Result Analysis**

- Larger Students are Better than Smaller Teachers
  - 3B Subquadratic models are faster than 1B Transformer models
  - Proposed MambaInLlama-3B and Llamba-4B model outperforms Llama-1B baseline, in terms of coverage and accuracy while being faster



## **Result Analysis**

- Smaller Models have possibility to achieve better accuracy
  - Smaller models have great coverage, but gap between larger in accuracy test.
  - Smaller models have the ability to generate the correct answer → Developing better reward model can be helpful
  - For the tasks that coverage matters most(e.g. easily verifiable coding, mathematical proofs), smaller models can be preferred



MATH coverage





# Supervised Fine-Tuning

#### SFT improves the models significantly

- While distillation effectively transfers knowledge from the teacher model, SFT further refines and aligns the model's capabilities
- With SFT, subquadratic architectures can surpass teachers!





### Conclusion



V. Conclusion

# Conclusion

- Investigated whether lower-complexity models can leverage their superior generation throughput to outperform similarly sized Transformers
- Focused on reasoning tasks to scale test-time computation
- In fixed memory and computation source, proposed models achieve better coverage and accuracy for most time budgets compared to Transformers
- The paper highlights the potential of Mamba and other attention alternatives



V. Conclusion

# Limitation and Future Research

#### Lack of proposed method

- Propose new distilling method, Develop better reward model, ...
- Need to compare with existing KD models
  - What is the exact limitation of existing models?
  - What is the difference between the existing models and proposed models?

### Need to analyze the reason of improvement

- Why the Mamba model is superior?
- Why was the proposed model able to achieve this performance?



## Thank You!

Chanhee Lee, Yongjun Kim, Jiwoo Kim

