# Reasoning Models Can Be Effective Without Thinking - Ma et. al.

**Keno Bürger**
Scientific Seminar Machine Learning
11 July 2025

Associate Professorship of Machine Learning
Technical University of Munich

# Motivation

- Reasoning models have significantly improved the performance of LLM systems
- Include an explicit, lengthy Thinking process at the cost of increased token usage and latency
- Paper questions the necessity of this Thinking process
- Ongoing research to improve efficiency and effectiveness of reasoning models
- Investigates the effectiveness of reasoning models without the Thinking process
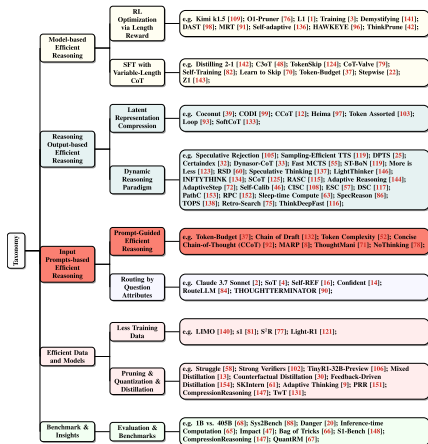
# Background Knowledge: Reasoning Models

- Training Reasoning Models:
    - Chain-of-Thought (CoT) Fine-Tuning: Fine-tune models on datasets with explicit reasoning steps to improve interpretability and accuracy
    - Instruction Tuning: Train models on task-specific instructions to guide reasoning for complex tasks
    - Reinforcement Learning (RL): Use RL with human feedback (RLHF) to reward correct and interpretable reasoning steps
- Key Difference from non-Reasoning Models:
    - Reasoning models are explicitly trained to generate intermediate steps, while non-reasoning models rely on implicit patterns in the data

Introduction
○●○○○

Methodology
○○○○○

Key Results
○○○

Analysis & Discussion
○○○

3

# Background Knowledge: Efficiency in LLMs

- Efficiency Trade-offs:
  - Explicit reasoning improves accuracy but increases latency and token usage
  - Token-constrained settings require balancing reasoning depth and efficiency
- Implicit Reasoning:
  - Models generate answers without explicit reasoning steps
- Scaling Approaches:
  - Sequential Scaling: Processes reasoning steps one at a time, increasing latency but ensuring step-by-step accuracy
  - Parallel Scaling: Executes multiple reasoning paths simultaneously, reducing latency but requiring more computational resources

# Related Work

- Contrasting prior work, this paper demonstrates competetive performance by disabling the Thinking process
- Techniques like CoT and instruction tuning to elicit reasoning abilities in LLMs [Huang+ 2023; Ott+ 2023; Xu+ 2025]
- Ongoing debate about whether LLMs truly reason or simply mimic reasoning patterns from training data [Hochlehnert+ 2025; Xu+ 2025]
  → necessity of explicit reasoning steps?

# Main Contribution

- Identification of NoThinking's **low-budget superiority**, particularly in token-constrained settings
- Development of a **parallel scaling framework** combining NoThinking with best-of-N sampling, achieving **7× lower latency** and **4× less token usage** compared to sequential Thinking approaches
- comparable results can be achieved through **implicit reasoning**

Introduction
○○○○●

Methodology
○○○○○

Key Results
○○○

Analysis & Discussion
○○○

6

# Methods and Experimental Setup

| **Thinking (6099 Completion Tokens)** |
|---|
| **User:** Return your final response within than . There exist real numbers $x$ and $y$, both greater than 1 such that $\log_x(y^x) = \log_y(x^{4y}) = 10$. Find $xy$. You should carefully think about the problem and reason step by step. |
| **Assistant:** |
| \<think\> |
| Alright, so I have this logarithmic problem here, and I need to find the product $xy$ where $x$ and $y$ are real numbers greater than 1 |
| ● ● ● (vertical) |
| \</think\> |
| ... thus the value of $xy$ is 25. |

| **NoThinking (373 Completion Tokens)** |
|---|
| **User:** Return your final response within than . There exist real numbers $x$ and $y$, both greater than 1 such that $\log_x(y^x) = \log_y(x^{4y}) = 10$. Find $xy$. Please write the answer for this math problem directly without any thinking process. |
| **Assistant:** |
| \<think\> |
| Okay I have finished thinking. |
| \</think\> |
| ... thus the value of $xy$ is 25. |

- DeepSeek-R1-Distill-Qwen-32B (additional baseline: Qwen-32B-Instruct)
- Budget forcing technique [Muennighoff+ 2025]: average token usage from NoThinking applied to Thinking → at this limit model is prompted to produce a final answer

# DeepSeek-R1-Distill-Qwen-32B

- First-generation reasoning model from DeepSeek
- Distilled from DeepSeek-R1 based on Qwen2.5-32B architecture
- Created through sophisticated knowledge distillation process
- Two RL stages in teacher model to discover improved reasoning patterns
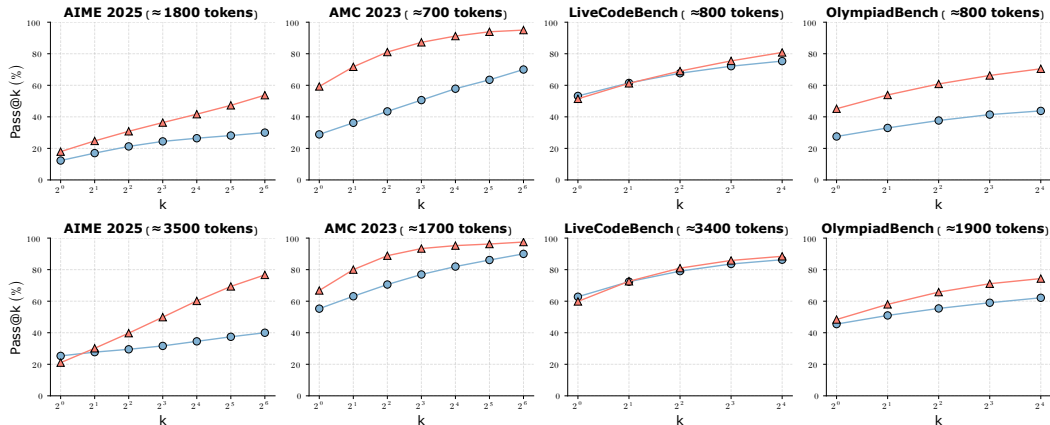- Two SFT stages to improve reasoning and non-reasoning capabilities

Introduction
○○○○○

**Methodology**
○●○○○

Key Results
○○○

Analysis & Discussion
○○○

8

# Benchmarks

| Use Case | Benchmark | Details |
|---|---|---|
| Mathematical Problem Solving | AIME 2024/AIME 2025 | ■ 15 questions, 3 hours<br>■ Answers are integers from 0 to 999 |
| | AMC 2023 | ■ 25 multiple choice questions, 40 minutes<br>■ One correct answer out of 5 choices |
| | Olympiad Bench | ■ Advanced reasoning benchmark containing math and physics problems<br>■ 8476 problems with an sophisticated automated scoring pipeline to verify solutions |
| Coding | LiveCodeBench | ■ Contamination-free evaluation of coding abilities<br>■ Code execution framework evaluates generated programs against test case |
| Formal Theorem Proving | MiniF2F | ■ Statements from olympiads (AMC, AIME, IMO) and high-school/undergraduate math classes<br>■ Formal system proof checkers |
| | ProofNet | ■ Logic and theorem proving benchmark<br>■ 371 examples with formal theorem statements, natural language theorem statements, and proofs. |

Introduction
○○○○○

Methodology
○○●○○

Key Results
○○○

Analysis & Discussion
○○○

9

# Metrics

- **pass@k** measuring the probability of obtaining at least one correct output
    - $k$ randomly selected samples, n generated completions, c correct outputs
    - $\text{pass@}k = \mathbb{E}_{\text{problems}}\left[1 - \frac{\binom{n-c}{k}}{\binom{n}{k}}\right]$
- **Mean and standard deviation** of entropy
    - entropy: measure of uncertainty in a probability distribution (model prediction)
    - higher mean entropy (close to the theoretical maximum): high uncertainty, greater diversity across questions
    - low variance (small compared to the the mean): more consistent diversity

Introduction
○○○○○

**Methodology**
○○○●○

Key Results
○○○

Analysis & Discussion
○○○

10

# Scaling

- **Parallel Scaling**: Distributes computations across multiple devices to process tasks simultaneously
- **Sequential Scaling**: Executes computations one after another
- **Metrics**: Latency, maximum number of tokens.
- **Perfect verifiers not available**:
    - *Confidence-based*: Model's self-certainty (confidence score) in predictions. Methods include highest confidence selection and weighted voting across answers.
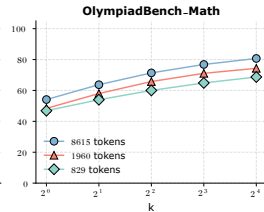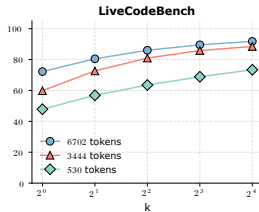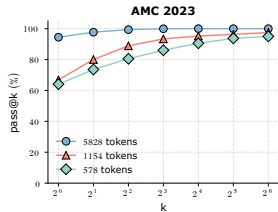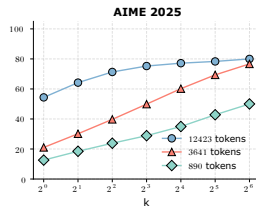    - *Majority voting*: Selection of the most common response (consensus)

# Token Efficiency

Introduction
○○○○○

Methodology
○○○○○

Key Results
●○○

Analysis & Discussion
○○○

12

# Accuracy

Introduction
○○○○○

Methodology
○○○○○

Key Results
○●○

Analysis & Discussion
○○○

13

# Parallel Scaling

| Task | Thinking | BF (tokens) | Pass@K | Selection Methods (Pass@1) | | |
|------|----------|-------------|--------|--------------------|----------------------|---------------------|
| | | | | **Majority Voting** | **Confidence + Highest** | **Confidence + Voting** |
| AIME 2024 | *Thinking* | 3500 | 73.33 | 43.33 | 40.00 | **46.67** |
| | NoThinking | 3500 | 77.30 | 46.67 | 20.00 | **50.00** |
| AIME 2025 | *Thinking* | 3500 | 40.00 | **30.00** | **30.00** | **30.00** |
| | NoThinking | 3500 | 53.73 | **33.33** | 20.00 | **33.33** |
| AMC 2023 | *Thinking* | 2400 | 92.50 | **77.50** | 65.00 | **77.50** |
| | NoThinking | 2400 | 95.00 | 77.50 | 57.50 | **85.00** |

- Improved pass@1 results at similar or lower latency
- Reduced token usage for tasks with perfect verifiers

Introduction
○○○○○

Methodology
○○○○○

Key Results
○○●

Analysis & Discussion
○○○

14

# Analysis



- In low-budgt settings, NoThinking outperforms Thinking
- Parallel scaling $\rightarrow$ improved accuracy-latency tradeoff

## Discussion

| Strengths | Limitations |
| --- | --- |
| <ul><li>Broad evaluation</li><li>Analysis of different scaled models</li><li>Efficiency gains</li><li>Practical Applications with parallel scaling</li></ul> | <ul><li>Poor performance at pass@1 and Coding</li><li>Results limited to one model</li><li>No interpretive/subjective/humaties-style tasks</li><li>No ablation studies on different components</li><li>Raw performance still favors large models</li></ul> |

$\Rightarrow$ Peer review and test on general conversation/common sense reasoning tasks

$\Rightarrow$ Applications with strict latency or cost constraints

# References I

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu and Maosong Sun.
**OlympiadBench: A Challenging Benchmark for Promoting AGI with Olympiad-Level Bilingual Multimodal Scientific Problems**. 2024. eprint: 2402.14008.

Andreas Hochlehnert, Hardik Bhatnagar, Vishaal Udandarao, Samuel Albanie, Ameya Prabhu and Matthias Bethge.
**A Sober Look at Progress in Language Model Reasoning: Pitfalls and Paths to Reproducibility**. Apr. 2025. DOI: 10.48550/arXiv.2504.07086.

Jie Huang and Kevin Chen-Chuan Chang. **Towards Reasoning in Large Language Models: A Survey**. May 2023. DOI: 10.48550/arXiv.2212.10403.

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen and Ion Stoica.
**LiveCodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code**. June 2024. DOI: 10.48550/arXiv.2403.07974.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès and Tatsunori Hashimoto. **s1: Simple test-time scaling**. Mar. 2025. DOI: 10.48550/arXiv.2501.19393.

Simon Ott, Konstantin Hebenstreit, Valentin Liévin, Christoffer Egeberg Hother, Milad Moradi, Maximilian Mayrhauser, Robert Praas, Ole Winther and Matthias Samwald. **ThoughtSource: A central hub for large language model reasoning data**.
In: *Scientific Data* 10.1 (Aug. 2023). ISSN: 2052-4463. DOI: 10.1038/s41597-023-02433-3.

# References II

Fengli Xu, Qianyue Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, Chenyang Shao, Yuwei Yan, Qinglong Yang, Yiwen Song, Sijian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao and Yong Li. **Towards Large Reasoning Models: A Survey of Reinforced Reasoning with Large Language Models**. Jan. 2025. DOI: 10.48550/arXiv.2501.09686.

Kunhao Zheng, Jesse Michael Han and Stanislas Polu. **MiniF2F: a cross-system benchmark for formal Olympiad-level mathematics**. Feb. 2022. DOI: 10.48550/arXiv.2109.00110.