The Key Ingredients for Scaling Test-Time Compute (and what is still missing)

Aviral Kumar

Carnegie Mellon University School of Computer Science

Test-Time Scaling for Large Language Models



Test-Time Scaling for Large Language Models



Test-Time Scaling for Large Language Models



Finetuning LLMs is Critical for Test-Time Scaling

Scaling LLM Test-Time Compute Optimally can be More Effective than Scaling Model Parameters

Charlie Snell[•], ¹, **Jaehoon Lee**², **Kelvin Xu**[•], ² and **Aviral Kumar**[•], ² [•]Equal advising, ¹UC Berkeley, ²Google DeepMind, [•]Work done during an internship at Google DeepMind **Key idea:** Finetune LLMs to enable this behavior!



A Family of Test-Time Scaling Algorithms



This training has been done via:

- RL or SFT
- Dense vs sparse supervision
- Different "types" of training data

Training for test-time scaling: Train LLMs to implement these procedures

This Talk: Scaling Test-Time Compute











Amrith Setlur (PhD @ CMU) Matthew Yang (MS @ CMU)



Nived Rajaraman (PhD @ Berkeley)



Formulating a Learning Problem

Main paper and blog post covered:

- Scaling Test-Time Compute without Verification or RL is Suboptimal Setlur, Rajaraman, Levine, Kumar. arXiv 2025
- Optimizing Test-Time Compute requires Solving a Meta-RL Problem Setlur, Qu, Yang, Zhang, Smith, Kumar. CMU MLD Blog, 2025.

Desiderata: What We Want at Test Time

Easy input problems

Don't spend too many tokens in arriving at the solution! (Also referred to as "don't overthink")

Hard input problems

- Make "progress" on hard problems
 - > Explore multiple strategies
 - > Exploit some promising ones, retry in different ways if needed

Do Current Models Enjoy these Desiderata?

Short answer: not really!

Experiment setup

Chop the thinking block in DeepSeek-R1 and ask it to produce best answer

Full trace Partial trace <think> <think> Okay, so I have this Okay, so I have this problem where ... problem where ... Wait, let's parse the Time is up. problem again ... Given the time I've ... </think> Alternatively, perhaps it's better to ... But let me double-check </think> **Solution** **Solution** Step 1: ... Step 1: ... **Final Answer** **Final Answer*

> Hard problems: Make sure to make constant progress

Easy problems: Make sure to be efficient

Do Current Models Enjoy these Desiderata?

> Easy problems: Make sure to be *token*-efficient!

Hard problems: Make sure to make constant progress



Qu*, Yang* et al. Optimizing Test-Time Compute via Meta Reinforcement Finetuning. arXiv 2025.

Short answer: *not really!*

Do Current Models Enjoy these Desiderata?

Short answer: *not really!*

> **Easy problems:** Make sure to be *token*-efficient!

> Hard problems: Make sure to make constant progress



Takeaway: Can make progress by implementing the "algorithm" of running a simple majority vote, *but it does not.*

Formulation: How to Satisfy These Desiderata

Let's start from the final goal

$$\max_{\pi} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{test}}} \left[\mathbb{E}_{\mathbf{z} \sim \pi(\cdot | \mathbf{x})} \left[r(\mathbf{x}, \mathbf{z}) \right] \right]$$

on test problems response sampled from model (longer than typical solution)

Total compute constraint per problem

Qu*, Yang* et al. **Optimizing Test-Time Compute via Meta Reinforcement Finetuning.** arXiv 2025.

Formulation: How to Satisfy These Desiderata

$$\max_{\pi} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{test}}} \left[\mathbb{E}_{\mathbf{z} \sim \pi(\cdot | \mathbf{x})} \left[r(\mathbf{x}, \mathbf{z}) \right] \right] \text{ s.t. } \forall \mathbf{x}, \mathbb{E}_{\pi(\cdot | \mathbf{x})} | \mathbf{z} | \leq C_0$$

But this compute budget is fixed!

$$\max_{\pi} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{train}}} \left[\mathbb{E}_{\mathbf{z} \sim \pi(\cdot | \mathbf{x})} \left[r(\mathbf{x}, \mathbf{z}) \right] \right] \text{ s.t. } \dots$$

Can optimize this via:

- RL (like DeepSeek-R1): outcome-reward RL
- > SFT / STaR: collect data, filter by correctness, maximize likelihood

Why is Outcome Reward + Fixed Budget Bad?



Qu*, Yang* et al. Optimizing Test-Time Compute via Meta Reinforcement Finetuning. arXiv 2025.

Formulation: "Budget-Agnostic" LLMs

Key idea: Incentivize the LLM to make progress regardless of the compute budget





Designing Dense Rewards

Main paper covered:

- Optimizing Test-Time Compute via Meta Reinforcement Finetuning. *Qu*, Yang*, Setlur, Tunstall, Beeching, Salakhutdinov, Kumar. arXiv 2025*
- Also see: Rewarding Progress: Scaling Automated Process Verifiers for LLM Reasoning Setlur, Nagpal, Fisch, Geng, Eisenstein, R. Agarwal, A. Agarwal, Berant, Kumar. ICLR 2025 (Spotlight).

Idea: What is Good Progress on New Problems?



Qu*, Yang* et al. Optimizing Test-Time Compute via Meta Reinforcement Finetuning. arXiv 2025.

Idea: What is Good Progress on New Problems?



We want this area to be as low as possible! (i.e., a notion of cumulative regret)

Inducing a Good Curve That Makes Progress



"Episodes" in the test-time stream

Key idea: Can still push up the performance after every episode!



"Episodes" in the test-time stream

Key idea: Can still push up the performance after every episode!



Key idea: Can still push up the performance after every episode!

Concrete Idea: Progress Reward Design $\max_{\pi} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{train}}} \left[\mathbb{E}_{\mathbf{z} \sim \pi(\cdot | \mathbf{x})} \left[r(\mathbf{x}, \mathbf{z}) \right] \right] \text{ s.t. } \dots$ π over episodes $\mathbb{E}_{\mathbf{x}\sim\mathcal{D}_{\text{train}},\mathbf{z}\sim\pi(\cdot|\mathbf{x})} \left[\sum_{j} \nabla_{\pi} \log \pi(\mathbf{z}_{j}|\mathbf{x},\mathbf{z}_{0:j-1}) \cdot r(\mathbf{x},\mathbf{z}) \right]$ $\mathbb{E}_{\mathbf{x}\sim\mathcal{D}_{\text{train}},\mathbf{z}\sim\pi(\cdot|\mathbf{x})} \left[\sum_{j} \nabla_{\pi} \log \pi(\mathbf{z}_{j}|\mathbf{x},\mathbf{z}_{0:j-1}) \cdot r(\mathbf{x},\mathbf{z}) \right]$ $\mathbb{E}_{\mathbf{x}\sim\mathcal{D}_{\text{train}},\mathbf{z}\sim\pi(\cdot|\mathbf{x})} \left[\sum_{j} \nabla_{\pi} \log \pi(\mathbf{z}_{j}|\mathbf{x},\mathbf{z}_{0:j-1}) \cdot r(\mathbf{x},\mathbf{z}) \right]$ $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{train}}, \mathbf{z} \sim \pi(\cdot | \mathbf{x})} \left[\sum_{j} \nabla_{\pi} \log \pi(\mathbf{z}_{j} | \mathbf{x}, \mathbf{z}_{0:j-1}) \cdot (r(\mathbf{x}, \mathbf{z}) + \alpha \cdot r_{\text{prg}}(\mathbf{x}, \mathbf{z}_{0:j})) \right]$

Progress reward

Qu*, Yang* et al. Optimizing Test-Time Compute via Meta Reinforcement Finetuning. arXiv 2025.

Concrete Idea: Progress Reward Design

$$\mathbb{E}_{\mathbf{x}\sim\mathcal{D}_{\text{train}},\mathbf{z}\sim\pi(\cdot|\mathbf{x})}\left[\sum_{j}\nabla_{\pi}\log\pi(\mathbf{z}_{j}|\mathbf{x},\mathbf{z}_{0:j-1})\cdot(r(\mathbf{x},\mathbf{z})+\alpha\cdot r_{\text{prg}}(\mathbf{x},\mathbf{z}_{0:j}))\right]$$

So far, we were using accuracy of the model that aims to write out the solution given the episodes so far

Progress reward
A "prover" policy that guesses the best answer

$$r_{\text{prg}}(\mathbf{x}, \mathbf{z}_{0:j}) = J_r(\mu(\cdot | \mathbf{x}, \mathbf{z}_{0:j})) - J_r(\mu(\cdot | \mathbf{x}, \mathbf{z}_{0:j-1}))$$



Qu*, Yang* et al. Optimizing Test-Time Compute via Meta Reinforcement Finetuning. arXiv 2025.

(Subset of) Results: Our Approach (MRT)



Open Questions: Dense Rewards

Computational cost

See: RL on Incorrect Synthetic Data Scales the Efficiency of LLM Math Reasoning 8x. NeurIPS 2024.

> Estimating dense rewards requires rollouts, which are costly.

> Can we get more juice out of the same total FLOPs??

The choice of the prover policy

➤ The policy µ determines the progress reward.
 ➤ How should you choose this policy??

See: Rewarding Progress: Scaling Automated Process Verifiers. ICLR 2025.

Other ways of implementing the same principle

Length curriculum and iterative training could be one other way
 Many open-source implementations kinda do this!

Qu*, Yang* et al. Optimizing Test-Time Compute via Meta Reinforcement Finetuning. arXiv 2025.



Comparing RL to SFT for test-time scaling

Main paper covered:

 Scaling Test-Time Compute without Verification or RL is Suboptimal Setlur, Rajaraman, Levine, Kumar. arXiv 2025

Goal: Scaling Test-Time Compute Effectively

Good performance as we allow the model to use more compute

For this analysis, we use a simple notion of binary 0/1 (step) progress

Theory: Using rewards via RL > cloning via SFT

When pre-trained LLM satisfies certain criteria



test-time compute H, #training prompts n is arbitrary $\text{Reward}_{H}(RL) - \text{Reward}_{H}(SFT) > 0$

Result 1: Gap between using rewards vs. not is positive at big H values

Theory: Using rewards via RL > cloning via SFT

When pre-trained LLM satisfies certain criteria



test-time compute *H*, #training prompts $n = \Omega(H)$

40

better than SFT

Criterion 1: Base LLM is Heterogeneous

Intuition: Diversity of traces under dense reward, *per problem*



Criterion 2: Base LLM is "Anticoncentrated"

Intuition: High enough (constant) coverage of better traces

Sample a bunch of traces and plot the *reward histogram*:



Theory: Putting It Together

[Informal] Separation b/w verifier-based & verifier-free

For the notion of step-level progress reward, when training with *compute-budget H* and *data size n, we see:*

1. Lower bound for VF (SFT): SubOpt_H(VF) = $\Omega\left(\frac{\sigma}{\sqrt{n}}\right), \sigma \approx \Omega(H)$

Variance of dense rewards on a given problem

Theory: Putting It Together

[Informal] Separation b/w verifier-based & verifier-free

For the notion of step-level progress reward, when training with *compute-budget H* and *data size n, we see:*

1. Lower bound for VF (SFT): SubOpt_H(VF) = $\Omega\left(\frac{\sigma}{\sqrt{n}}\right)$, $\sigma \approx \Omega(H)$

2. Upper bound for VB (RL): SubOpt_H(VB) = $0\left(\frac{1}{\frac{1}{Chase}},\frac{H}{n}\right)$

Anti-concentration probability mass

Error in learning the reward signal

Theory: Putting It Together

[Informal] Separation b/w verifier-based & verifier-free

For the notion of step-level progress reward, when training with *compute-budget H* and *data size n, we see:*

1. Lower bound for VF (SFT): SubOpt_H(VF) = $\Omega\left(\frac{\sigma}{\sqrt{n}}\right)$, $\sigma \approx \Omega(H)$

2. Upper bound for VB (RL): SubOpt_H(VB) = $0\left(\frac{1}{c_{base}} \cdot \frac{H}{n}\right)$

3. Gap b/w RL and SFT scales as $\Omega\left(\frac{H}{\sqrt{n}}\right)$. Set n = O(H).

Results: Scaling Data and Compute Both



Finding: Gap between verifier-based and verifier-free methods increases

Results: Scaling Compute, Not Data (w/ s1-32B)



Finding: Gap between using reward verification and SFT is positive



Takeaways

Key idea: RL (using reward signal) scales test-time compute better than SFT, with a notion of dense rewards.

>Main findings:

When data is scaled linearly in H, performance gap grows
 When data is kept the same, RL performs better still.
 Base LLMs tend to be heterogeneous & anti-concentrated.

Thank You! Questions?

Scaling Test-Time Compute Without Verification or RL is Suboptimal

Amrith Setlur¹, Nived Rajaraman², Sergey Levine² and Aviral Kumar¹ ¹Carnegie Mellon University, ²UC Berkeley

Optimizing Test-Time Compute via Meta Reinforcement Fine-Tuning

Yuxiao Qu^{*1}, Matthew Y. R. Yang^{*1}, Amrith Setlur¹, Lewis Tunstall², Edward Emanuel Beeching², Ruslan Salakhutdinov¹ and Aviral Kumar¹

¹Carnegie Mellon University, ²Hugging Face, ^{*}Equal contribution

Summary: Scaling Test-Time Compute



Learning "how": Train algorithm $A_{\theta}(x)$ to spend **extra test compute**, search over **responses** & discover the **final answer**.



MRT incorporates dense rewards Dense reward reward outcome 0/1 reward U Test compute budget (H) Summary: Use dense rewards w/ progress

