

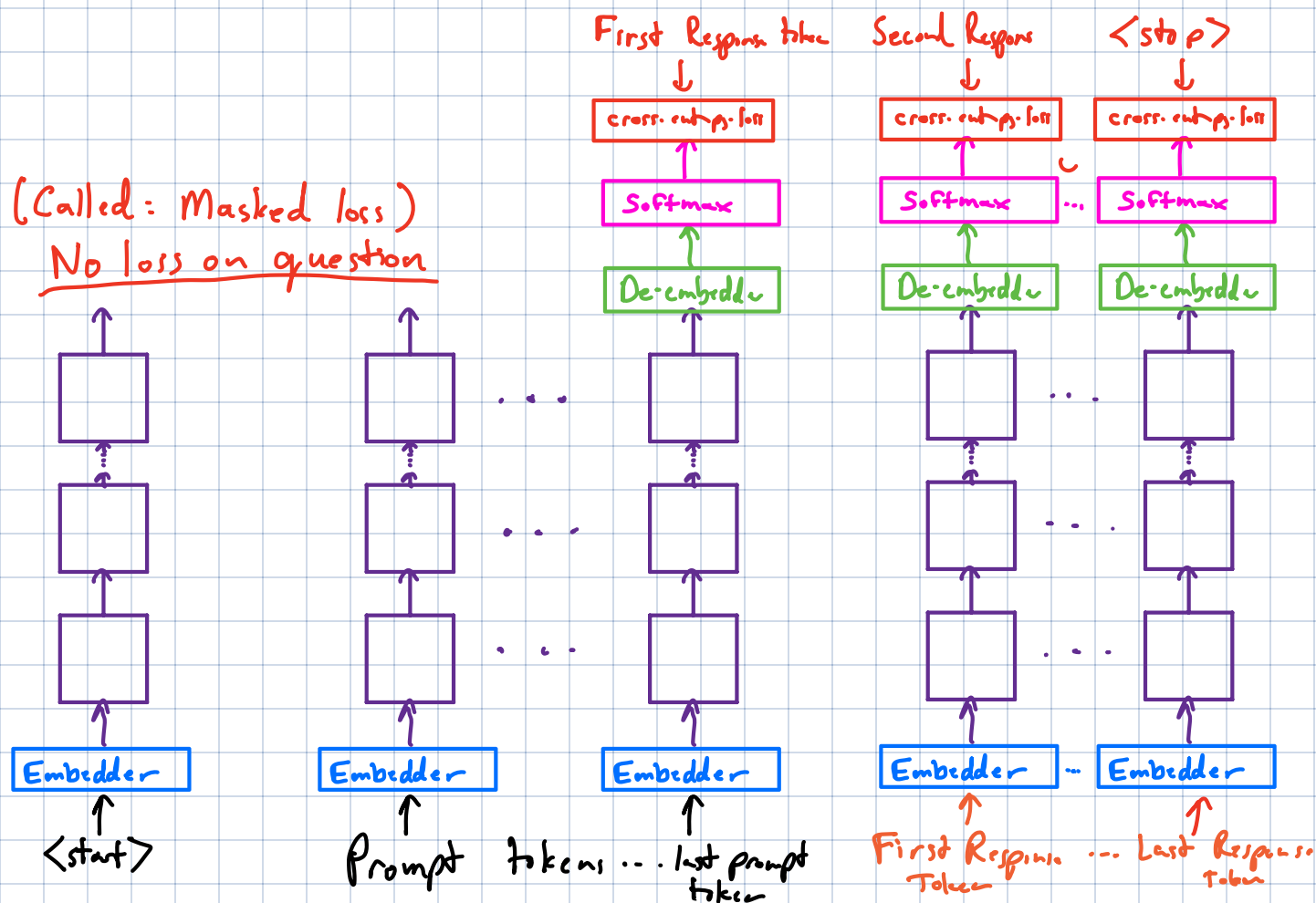
Today: Post-training and test-time compute

(We'll return to generative models)
after post-training interlude

Announce: Fill out survey
Extra Credit for everyone (3%)
IF 75% of class does the survey
Due: Fri of Thanksgiving week.
This Friday!!!
Finish your peer review
Today!!!

Post-training including RLVR.

Recall basic SFT for instruction following.



For LLMs, one key advance was "Chain-of-thought." So it is useful to have examples that show their work while solving problems.

But how can we make a model better at solving problems?

Two parts to the answer:

A) Be willing to spend more compute while answering.
 ↗ test time compute

B) Train it to be better.

Test-time compute: 0) Oldest Approach: pure prompting.
"Think step by step. Be careful."

1) Repeated generation. Instead of one try, generate N responses.

At Inference/Test time: A) Take Majority/Plurality Vote
B) Take highest prob
C) Train a reward model to grade.

2) Sample better...

Recent Paper: Karan & Du, "Reasoning with Sampling: Your base model is smarter than you think" Oct 16, 2025.

We want a good answer. That includes thinking. If we sample from original model, what can go wrong? Moving forward one token at a time, we can make a mistake... and then flounder.

Older Approach: Beam-Search with some k .

To do what? Sample from higher-likelihood sequences.

Alternative: sample not from $p(\vec{x})$ but from $\text{normalized}(p(\vec{x}))^\alpha$ where $\alpha > 1$

Is this just low temperature sampling?

No. Consider position t .

$$P_{\text{desired}}(x_t | x_0, \dots, x_{t-1}) = \frac{\sum_{x_{>t}} (p(x_0, \dots, x_{t-1}, x_t, x_{t+1}, \dots, x_T))^\alpha}{\sum_{\tilde{x}_t} \left[\sum_{x_{>t}} (p(x_0, \dots, x_{t-1}, \tilde{x}_t, x_{t+1}, \dots, x_T))^\alpha \right]}$$

Sequence. (pointing to \vec{x})
"integrate out" (pointing to the summation over $x_{>t}$)

The corresponding low-temp sampling doesn't raise the future completions to the desired power and then sum them.

Instead

$$P_{\text{lowTemp}}(x_T | x_0, \dots, x_{T-1}) = \frac{(p(x_T | x_0, \dots, x_{T-1}))^\alpha}{\sum_{\tilde{x}_T} (p(\tilde{x}_T | x_0, \dots, x_{T-1}))^\alpha}$$

$$= \frac{\left(\sum_{x_T} p(x_0, \dots, x_{T-1}, x_T, x_{T+1}, \dots, x_T) \right)^\alpha}{\sum_{\tilde{x}_T} \left[\sum_{x_T} p(x_0, \dots, x_{T-1}, \tilde{x}_T, x_{T+1}, \dots, x_T) \right]^\alpha}$$

Raising a sum to a power is not the same as the sum of powers

Key Idea in the paper: Use MCMC Techniques to sample

Algorithm 1: Power Sampling for Autoregressive Models

Input : base p ; proposal p_{prop} ; power α ; length T

Hyperparams: block size B ; MCMC steps N_{MCMC}

Output : $(x_0, \dots, x_T) \sim p^\alpha$

1 **Notation**: Define the unnormalized intermediate target

$$\pi_k(x_{0:kB}) \propto p(x_{0:kB})^\alpha$$

Goal

2 **for** $k \leftarrow 0$ **to** $\lceil \frac{T}{B} \rceil - 1$ **do**

3 Given prefix $x_{0:kB}$, we wish to sample from π_{k+1} . Construct initialization \mathbf{x}^0 by extending autoregressively with p_{prop} :

$$x_t^{(0)} \sim p_{\text{prop}}(x_t | x_{<t}), \quad \text{for } kB + 1 \leq t \leq (k+1)B.$$

Set the current state $\mathbf{x} \leftarrow \mathbf{x}^0$.

4 **for** $n \leftarrow 1$ **to** N_{MCMC} **do**

5 Sample an index $m \in \{1, \dots, (k+1)B\}$ uniformly.

6 Construct proposal sequence \mathbf{x}' with prefix $x_{0:m-1}$ and resampled completion:

$$x'_t \sim p_{\text{prop}}(x_t | x_{<t}), \quad \text{for } m \leq t \leq (k+1)B.$$

7 Compute acceptance ratio (9)

$$A(\mathbf{x}', \mathbf{x}) \leftarrow \min \left\{ 1, \frac{\pi_k(\mathbf{x}')}{\pi_k(\mathbf{x})} \frac{p_{\text{prop}}(\mathbf{x} | \mathbf{x}')}{p_{\text{prop}}(\mathbf{x}' | \mathbf{x})} \right\}.$$

Draw $u \sim \text{Uniform}(0, 1)$;

8 **if** $u \leq A(\mathbf{x}', \mathbf{x})$ **then accept** and set $\mathbf{x} \leftarrow \mathbf{x}'$

9 **end**

10 Set $x_{0:(k+1)B} \leftarrow \mathbf{x}$ to fix the new prefix sequence for the next stage.

11 **end**

12 **return** $x_{0:T}$

← Take random steps towards more likely generations

→ Traditional way to write in MCMC. In this case, reduces to $\frac{p_{\text{prop}}(\mathbf{x})}{p_{\text{prop}}(\mathbf{x}')}$ because of AR generation in p_{prop} .

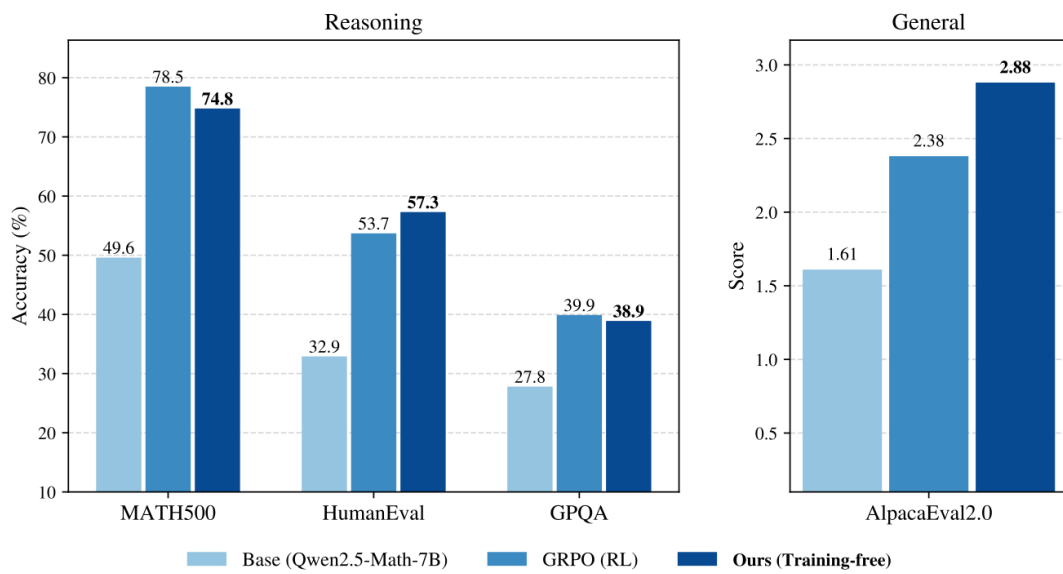


Fig 9 from Paper.

Exercise: What simpler alternative should also have been in this paper?

1) Try "best-of-N" with highest-pr.p seg chosen

1') Try beam search

2) Per prompt optimization

As distinct from RLHF: — with human feedback.

RLVR: Reinforcement Learning with Verifiable Rewards

Key new paper: Khatri, et al "The Art of Scaling Reinforcement Learning Compute for LLMs," Oct 15, 2025.

Treatment here influenced by this, along with DeepSeek paper.

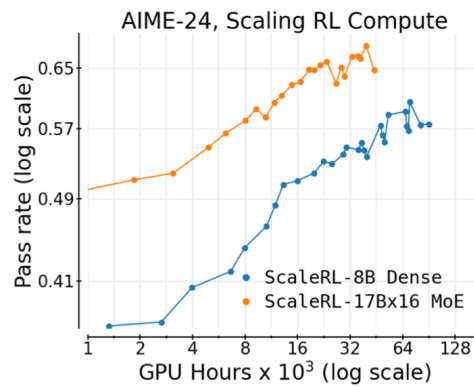
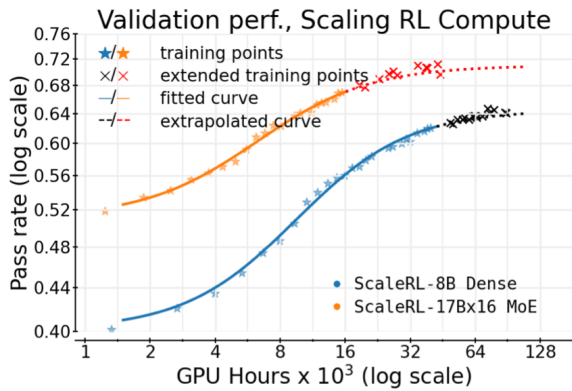


Fig 1 in ScaleRL...

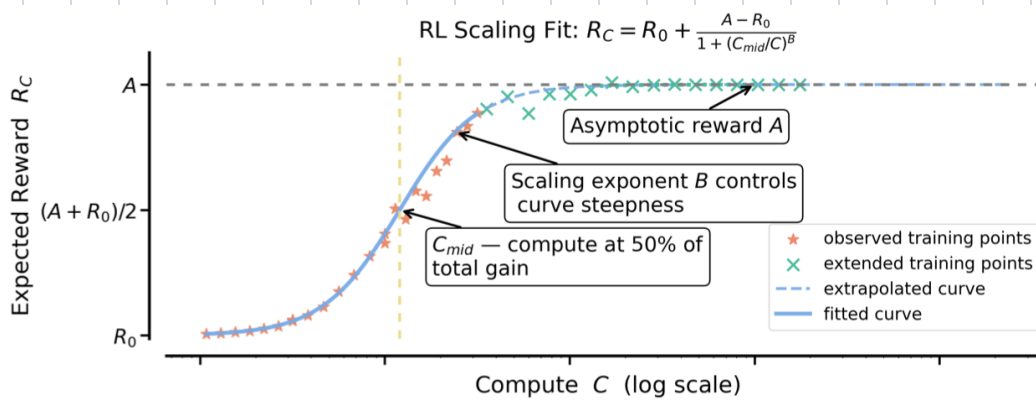


Fig 3 in ScaleRL

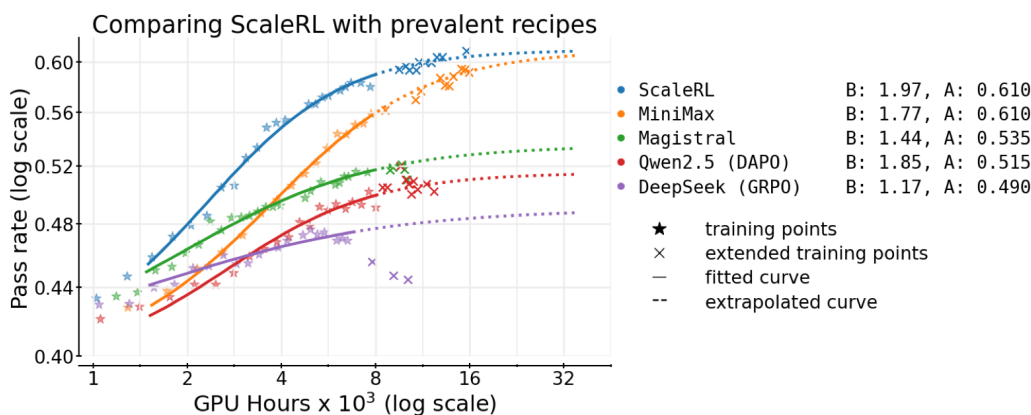


Fig 2 in ScaleRL...

"Bay of tokens perspective"

Unpacking what to do...

$$\mathcal{J}_{\text{ScaleRL}}(\theta) = \mathbb{E}_{x \sim D, \{y_i\}_{i=1}^G \sim \pi_{\text{gen}}^{\theta_{\text{old}}}(\cdot|x)} \left[\frac{1}{\sum_{g=1}^G |y_g|} \sum_{i=1}^G \sum_{t=1}^{|y_i|} \text{sg}(\min(\rho_{i,t}, \epsilon)) \hat{A}_i^{\text{norm}} \log \pi_{\text{train}}^{\theta}(y_{i,t}) \right],$$

$$\rho_{i,t} = \frac{\pi_{\text{train}}^{\theta}(y_{i,t})}{\pi_{\text{gen}}^{\theta_{\text{old}}}(y_{i,t})}, \quad \hat{A}_i^{\text{norm}} = \hat{A}_i / \hat{A}_{\text{std}}, \quad 0 < \text{mean}(\{r_j\}_{j=1}^G) < 1, \quad \text{pass_rate}(x) < 0.9,$$

"Importance Sampling"
Distribution Mismatch.

Demom across G generations recall

VR gives this.

Don't bother if we get 90% mtd.

Recall HW3 Trick (Reparameterization Gradient Estimator)

$$\nabla_{\theta} \mathbb{E}_{y \sim p_{\theta}} [F(y)] = \mathbb{E}_{y \sim p_{\theta}} [F(y) \nabla_{\theta} \log p_{\theta}]$$

$\text{sg}(\dots) \leftarrow$ To stop the gradient from inadvertently going into $\rho_{i,t}$.

$$\min(\rho_{i,t}, \epsilon) = \begin{cases} \epsilon & \text{if } \rho_{i,t} \geq \epsilon \\ \rho_{i,t} & \text{if } \rho_{i,t} < \epsilon \end{cases}$$

→ Caps the overweighting of this sample in case it is strongly preferred by π_{train} as opposed to π_{gen}

In general $\pi_{\text{gen}}^{\theta_{\text{old}}}$ closely tracks $\pi_{\text{train}}^{\theta}$ since it gets updated with a lag.