

The Art of Scaling Reinforcement Learning Compute for LLMs

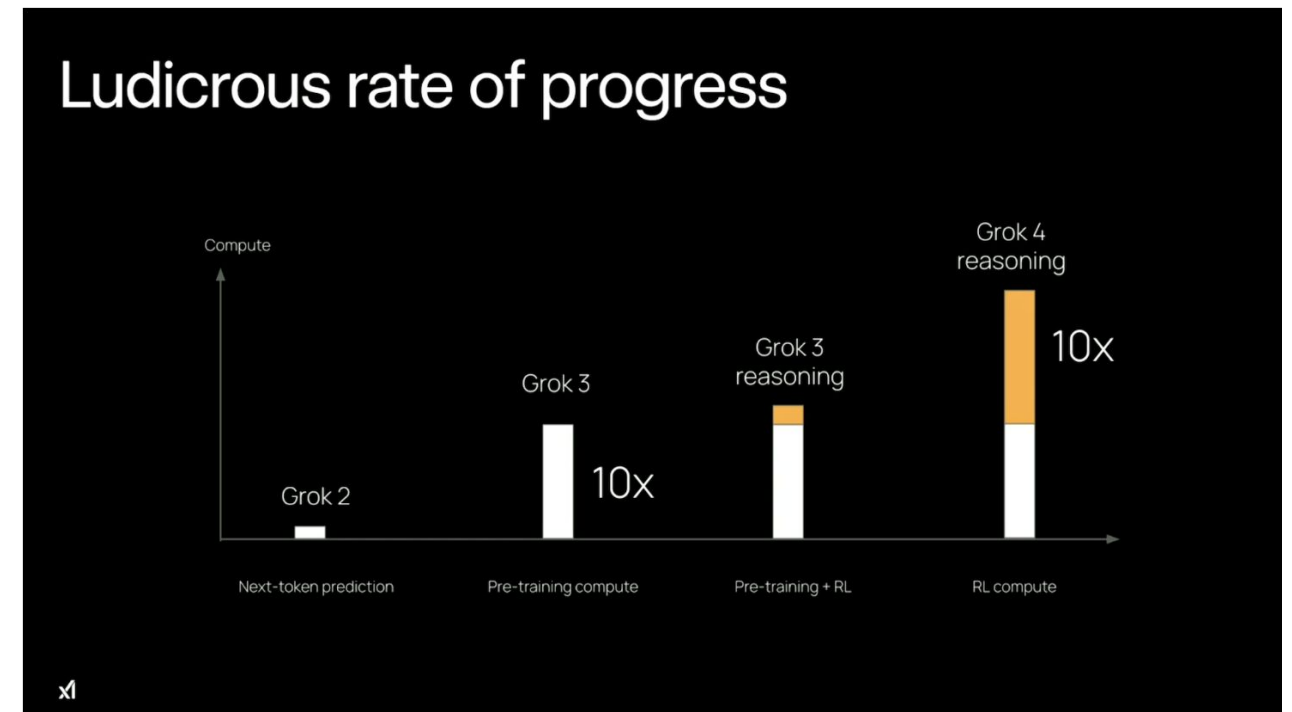
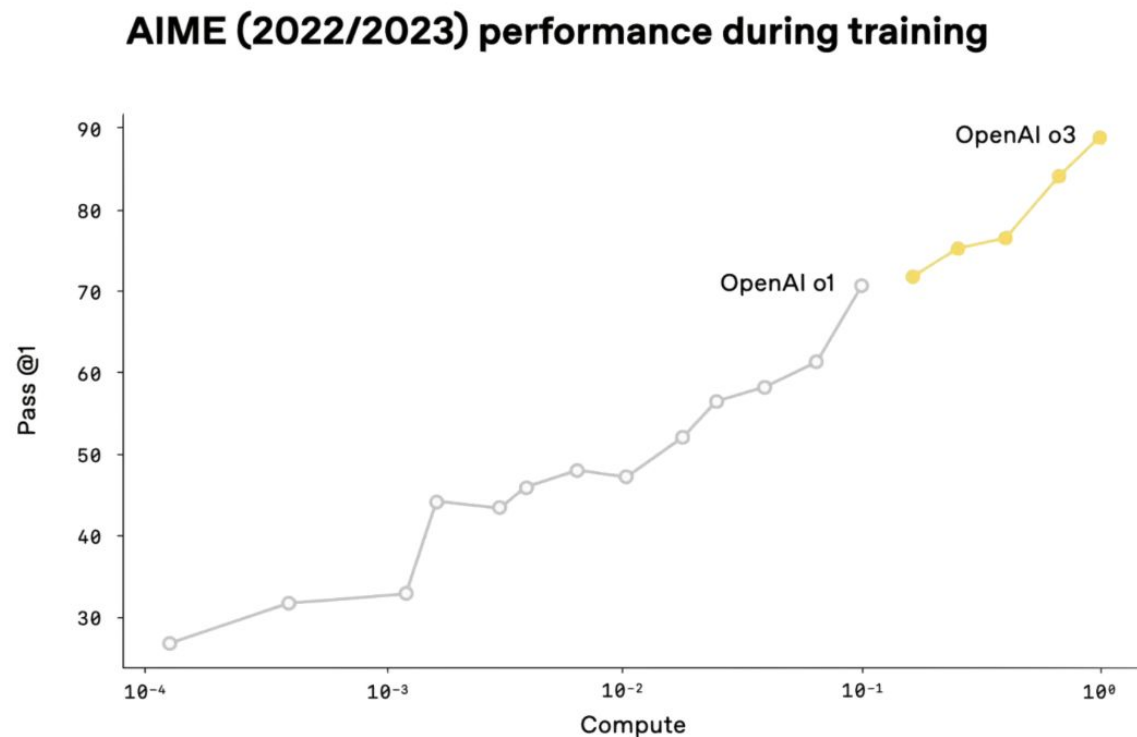
Devvrit Khatri*, Lovish Madaan*, Rishabh Tiwari, Rachit Bansal, Sai Surya Duvvuri, Manzil Zaheer, Inderjit Dhillon, David Brandfonbrener, Rishabh Agarwal

* - equal contribution



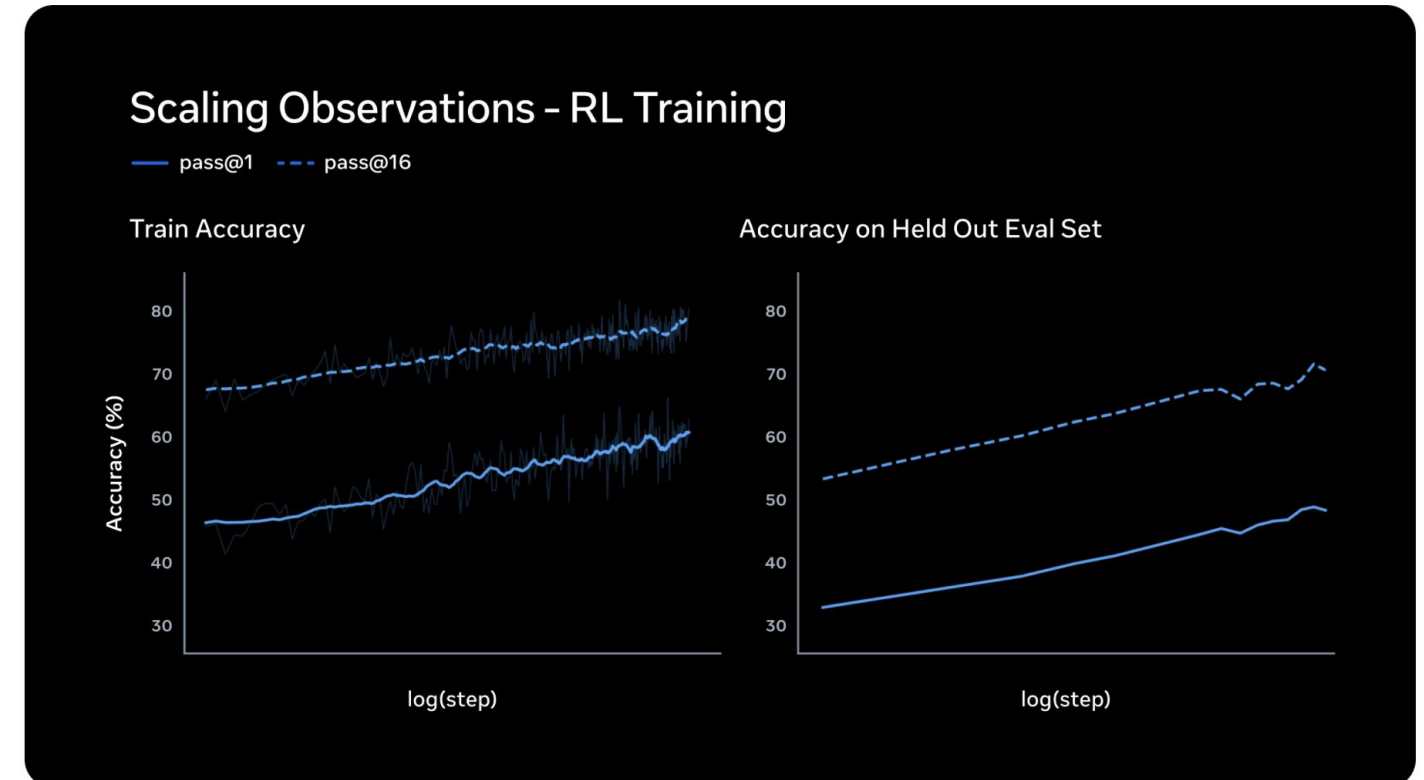
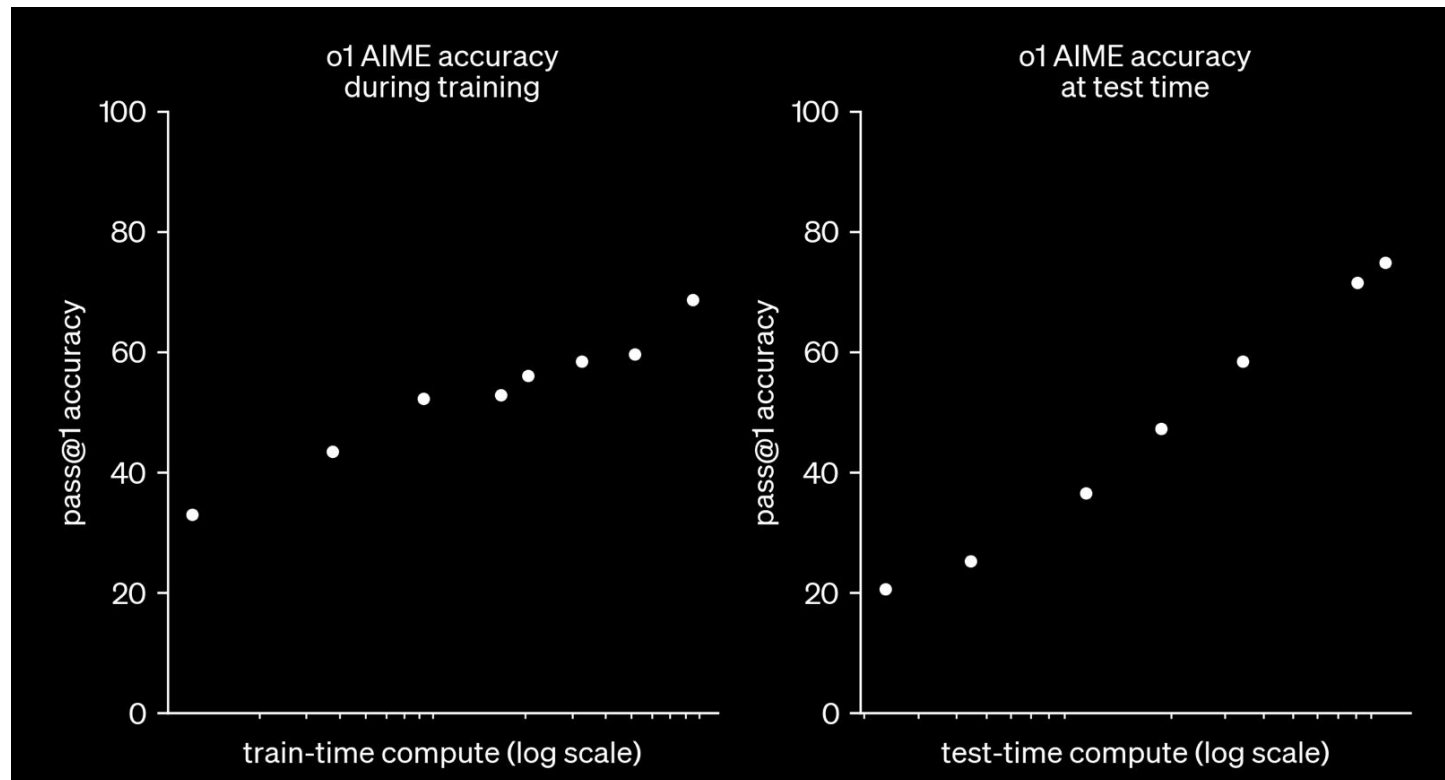
Motivation

- Reinforcement Learning (RL) is becoming the new scaling frontier for large language models (LLMs)



Motivation

- RL unlocks new capabilities in a base model like test-time thinking and agentic capabilities



Motivation

- Many design choices while building the RL infra stack
 - Off-policyness, loss, clipping, numerics, adaptive prompt filtering, generation length, batch size, etc.

Many algorithms

GRPO, DAPO, PPO,
CISPO, GSPO, IcePop, ...

Many design choices

Normalization, curriculum...

No scaling framework

Can't predict from
small runs

Misleading early results

Best at 5K \neq best
at 100K GPU-hrs

- Question: What design choices actually work best at scale? And can you achieve smooth scaling like pre-training?

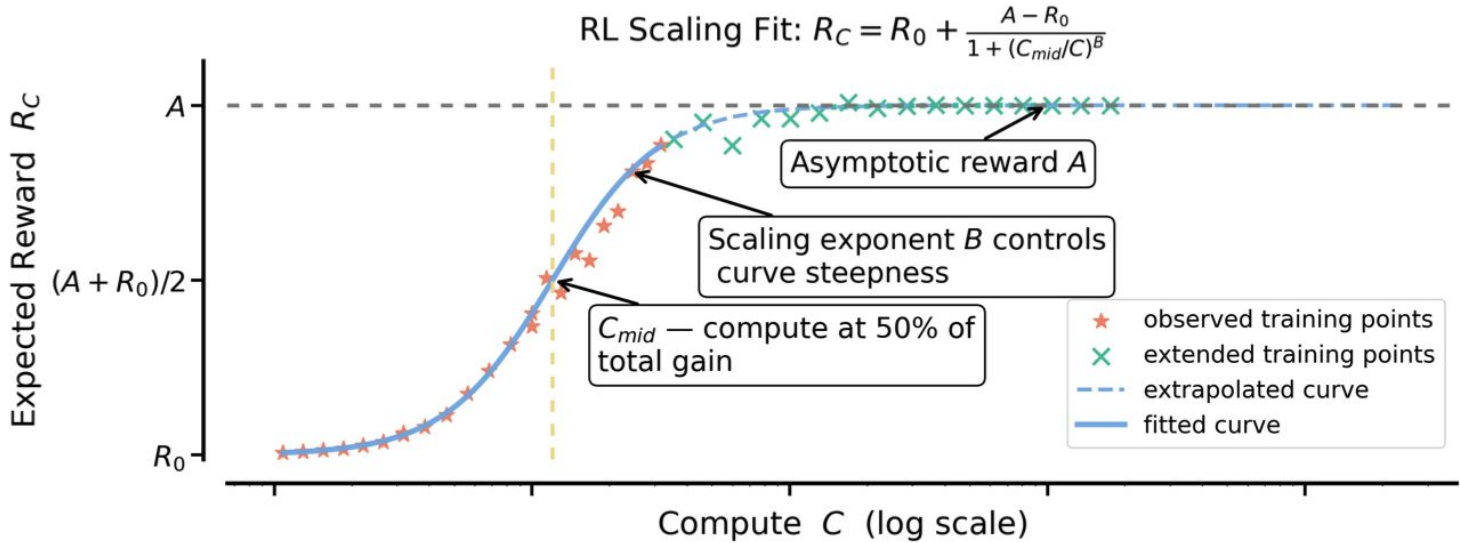
Predictable Scaling in RL

$$\underbrace{R_C - R_0}_{\text{Reward Gain}} = \underbrace{(A - R_0)}_{\text{Asymptotic Reward Gain}} \times \underbrace{\frac{1}{1 + (C_{\text{mid}}/C)^B}}_{\text{Compute Efficiency}}$$

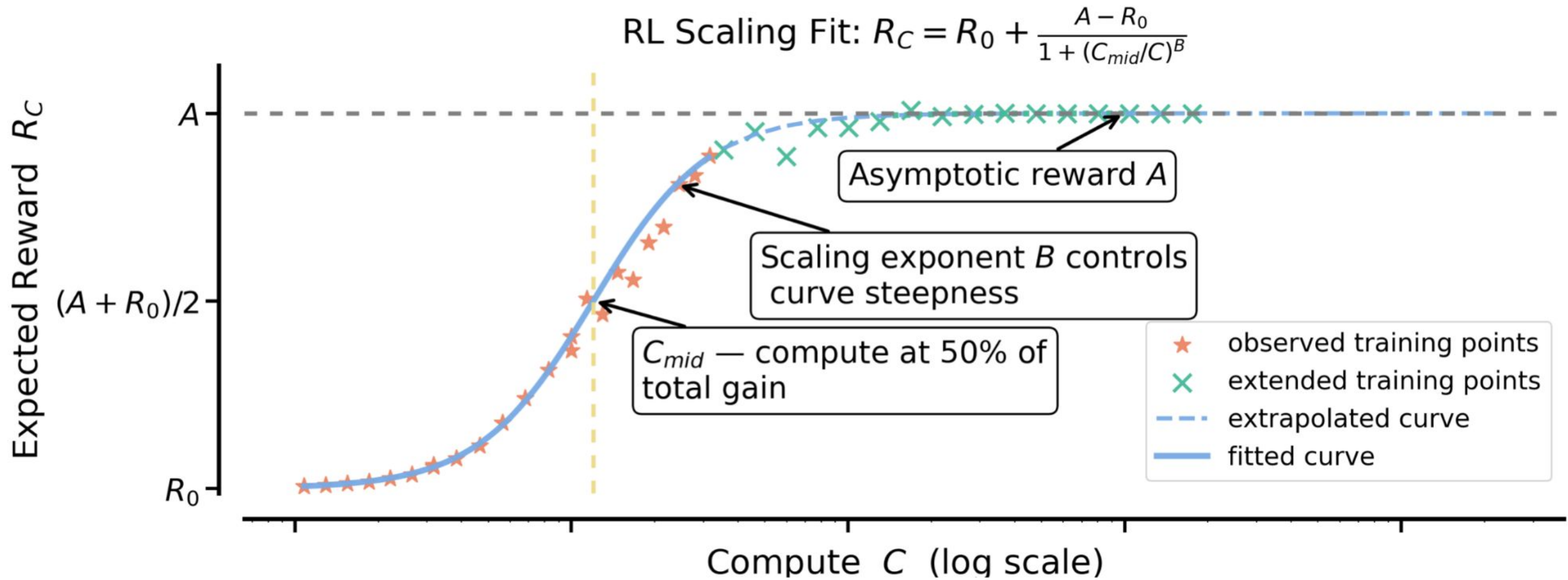
(fixed model and training data)

- Function of compute
- Bounded (because of accuracy metric)

Predictable Scaling



Predictable Scaling in RL



Asymptotic Ceiling
How high?

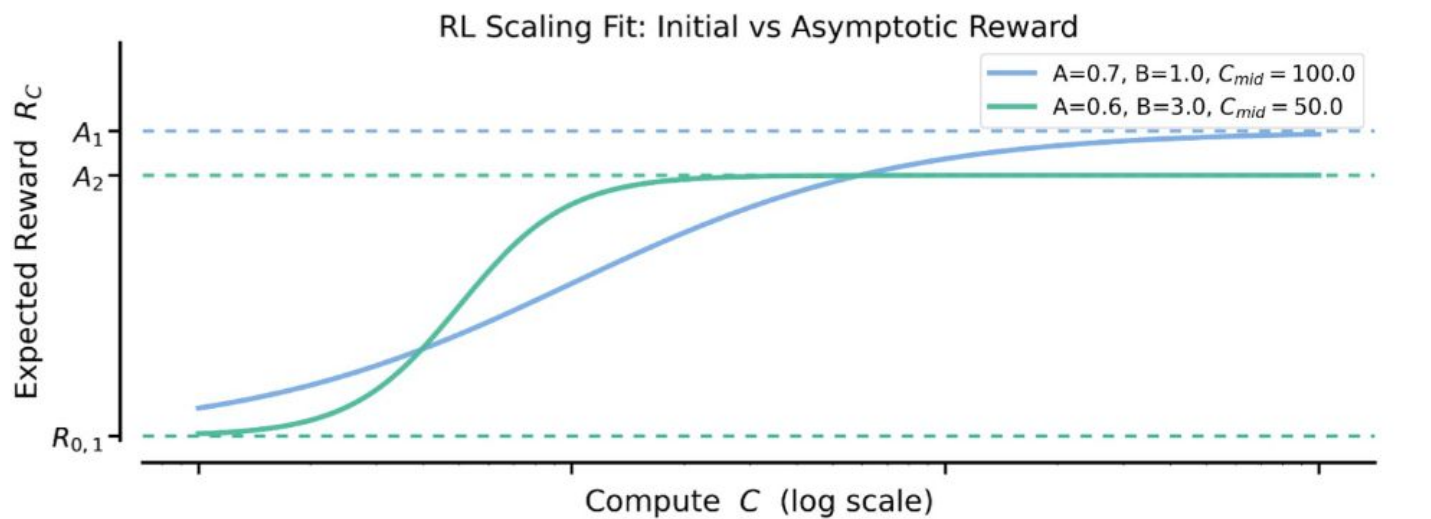
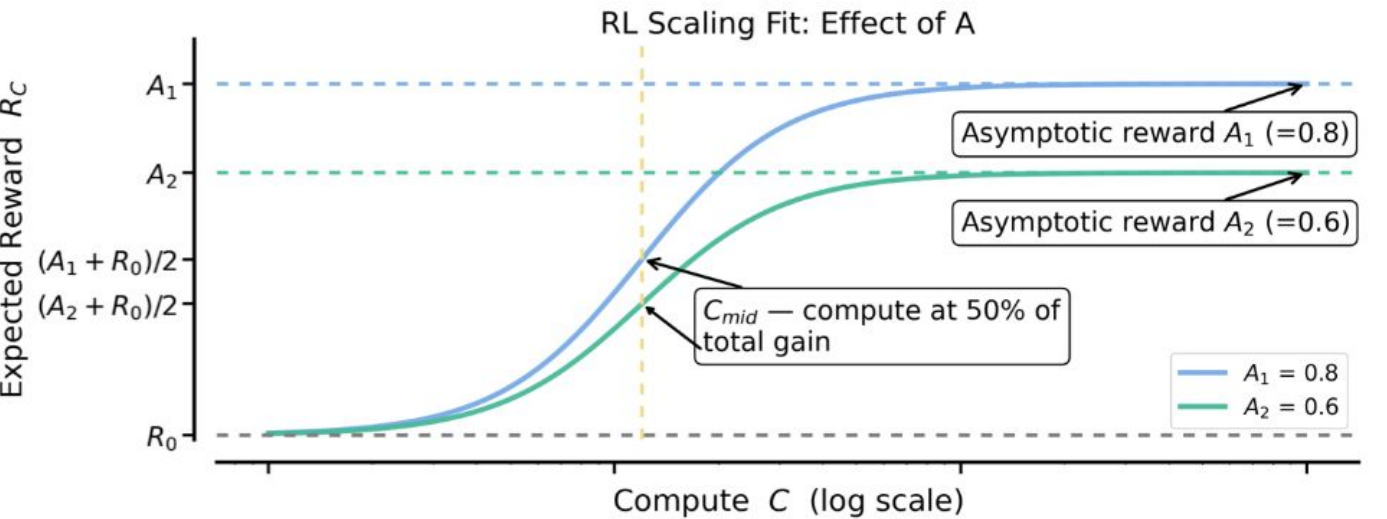
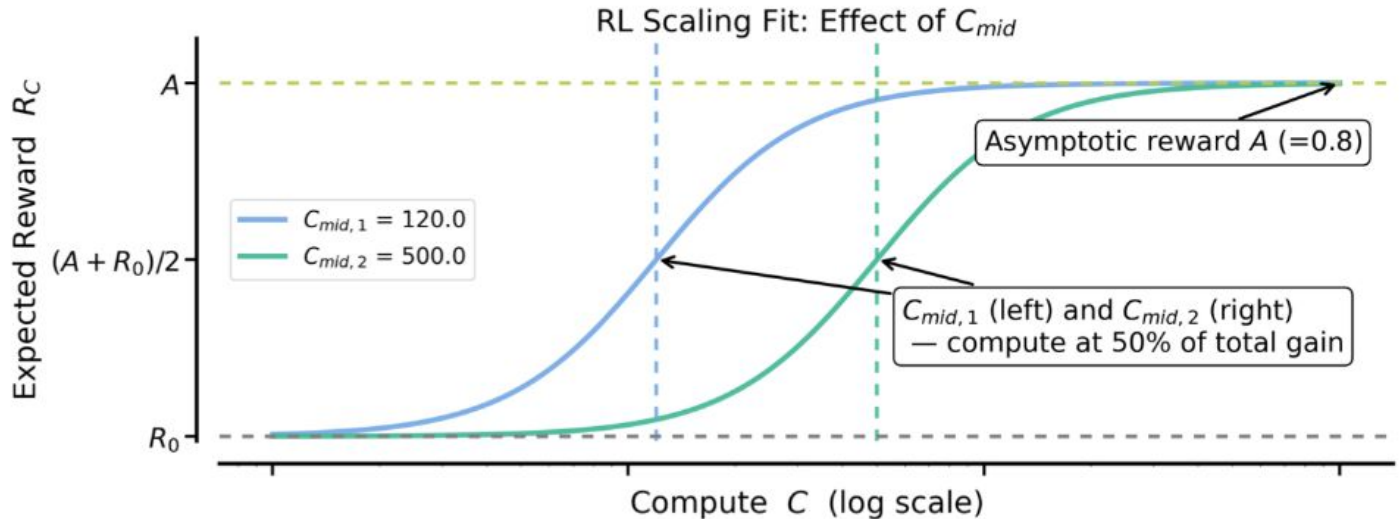
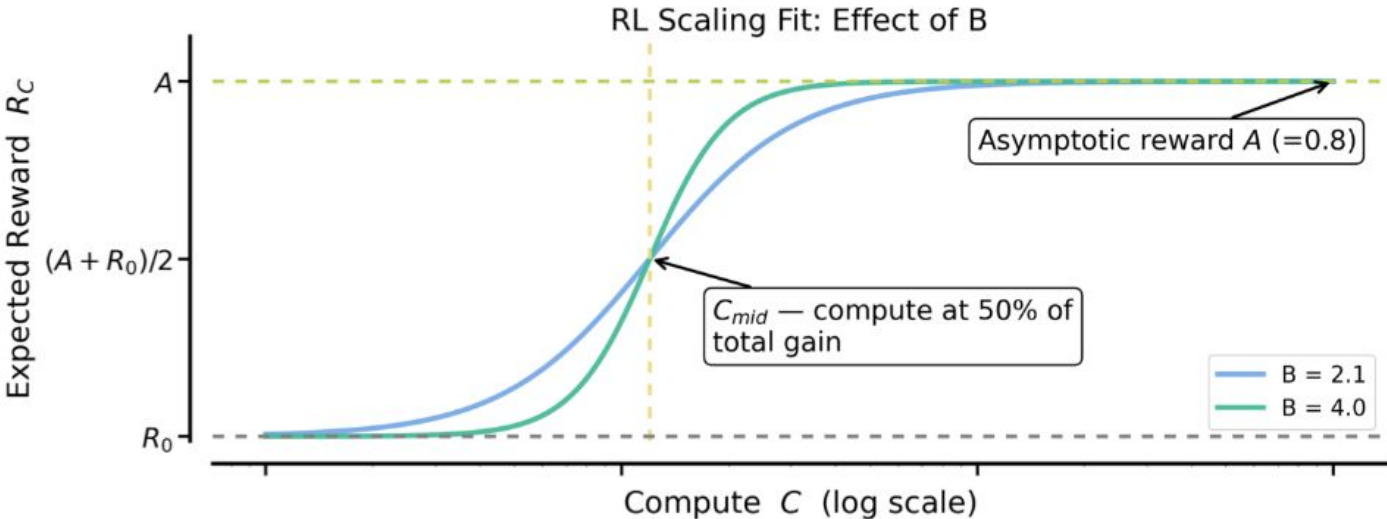


Compute Efficiency
How fast?

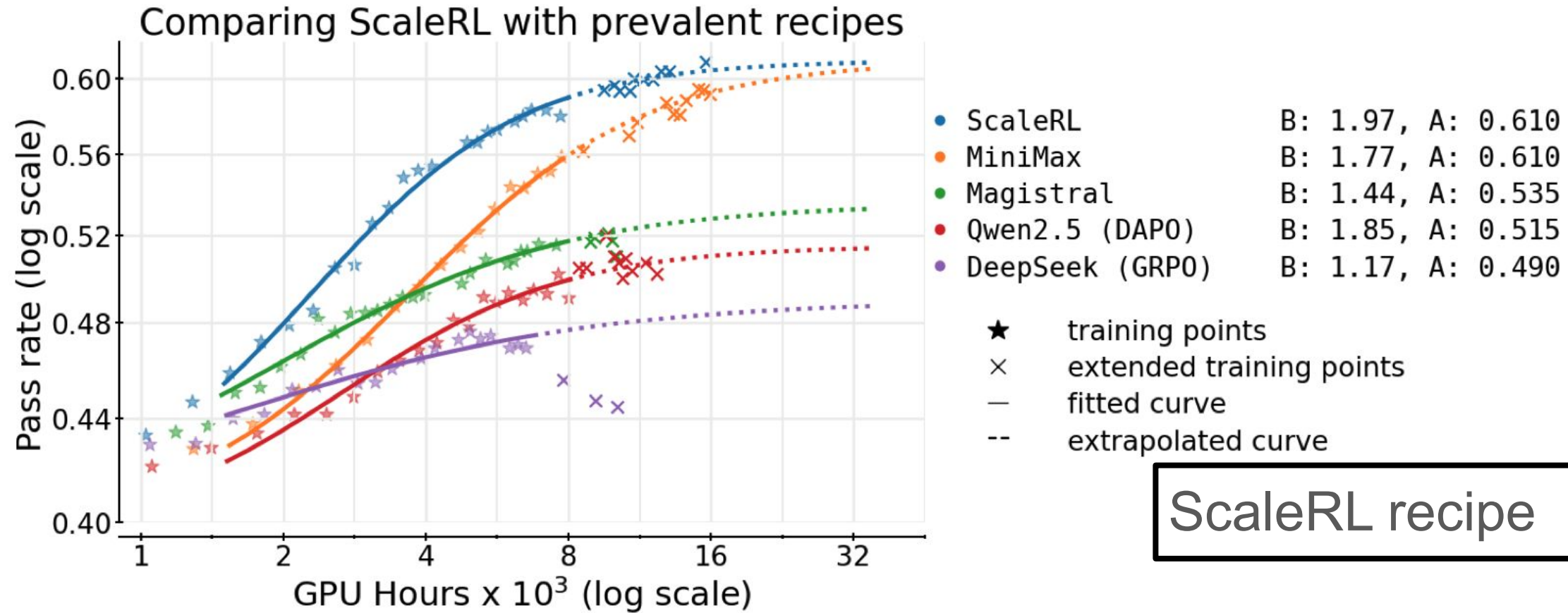


Inflection
When saturate?

Predictable Scaling in RL



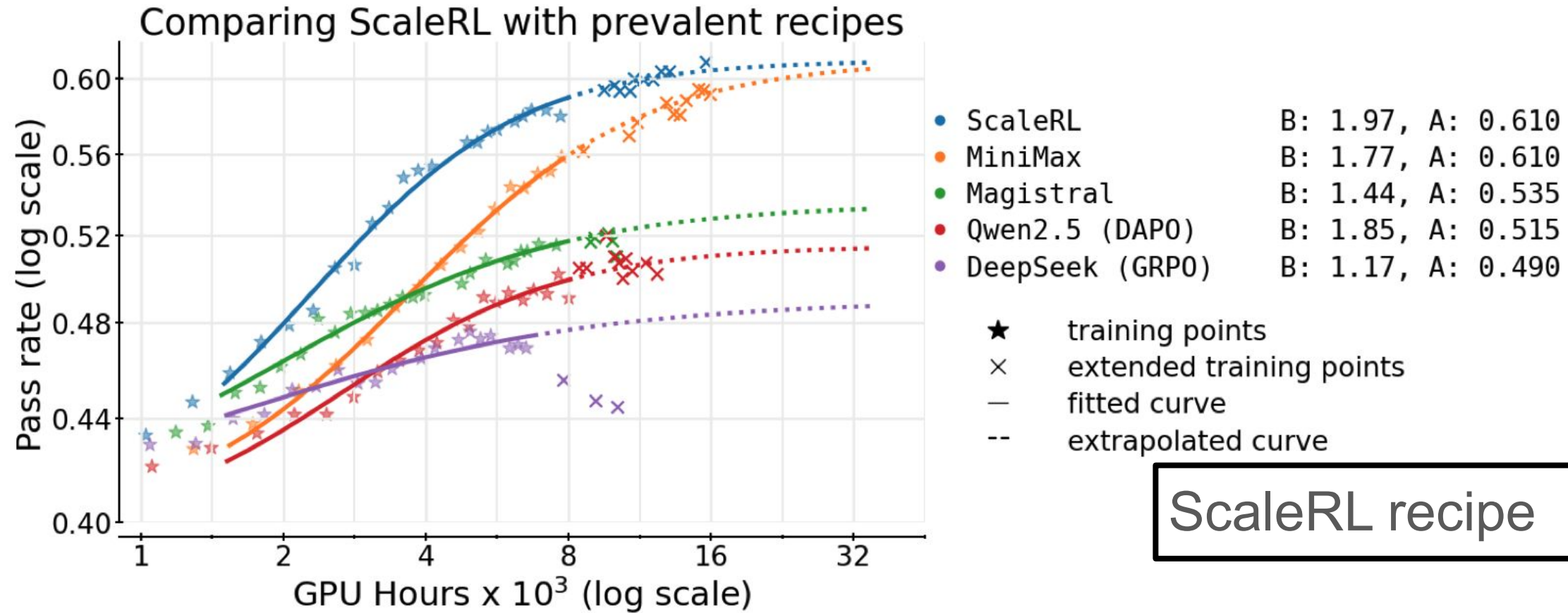
Different recipes have different ceilings



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

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

Different recipes have different ceilings



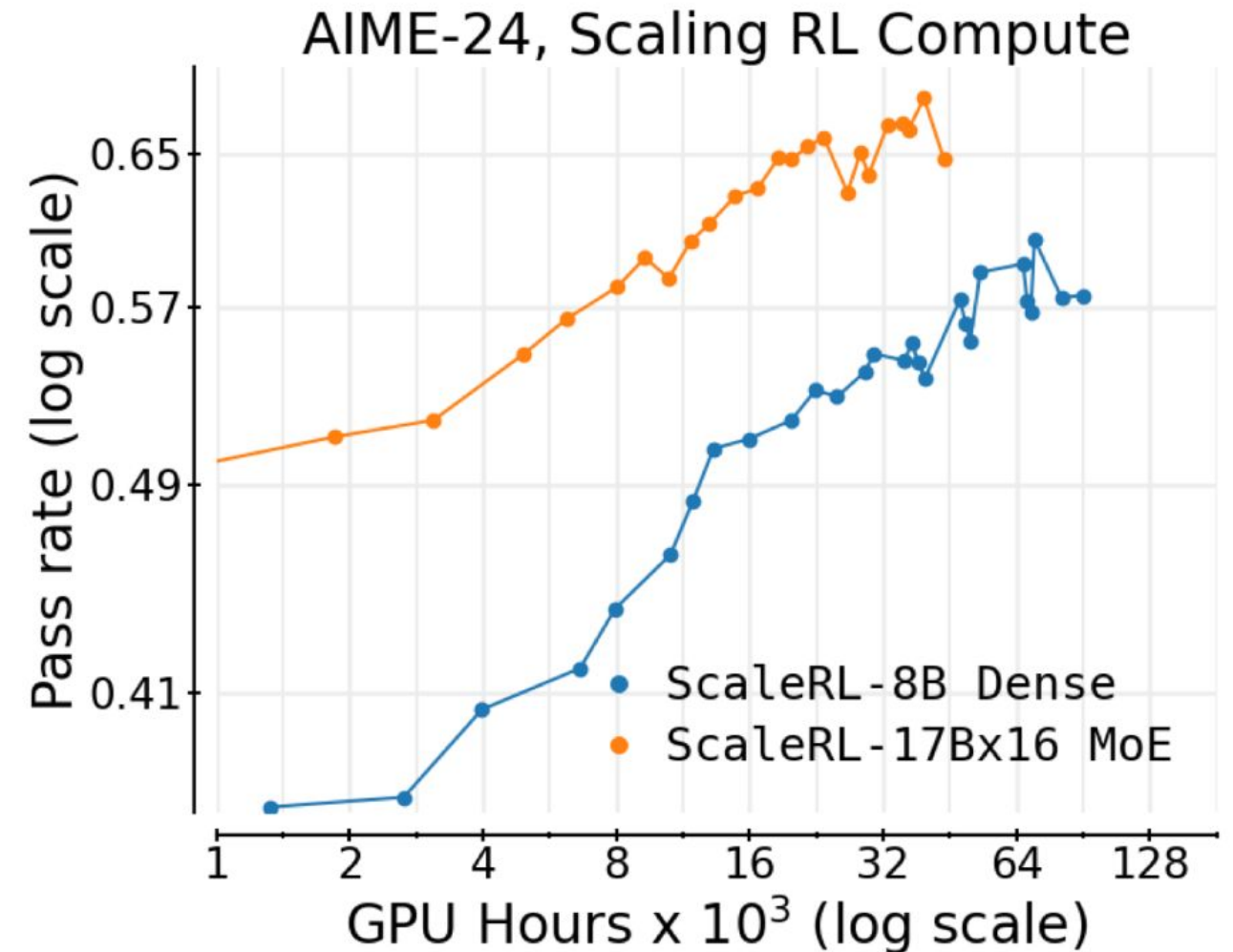
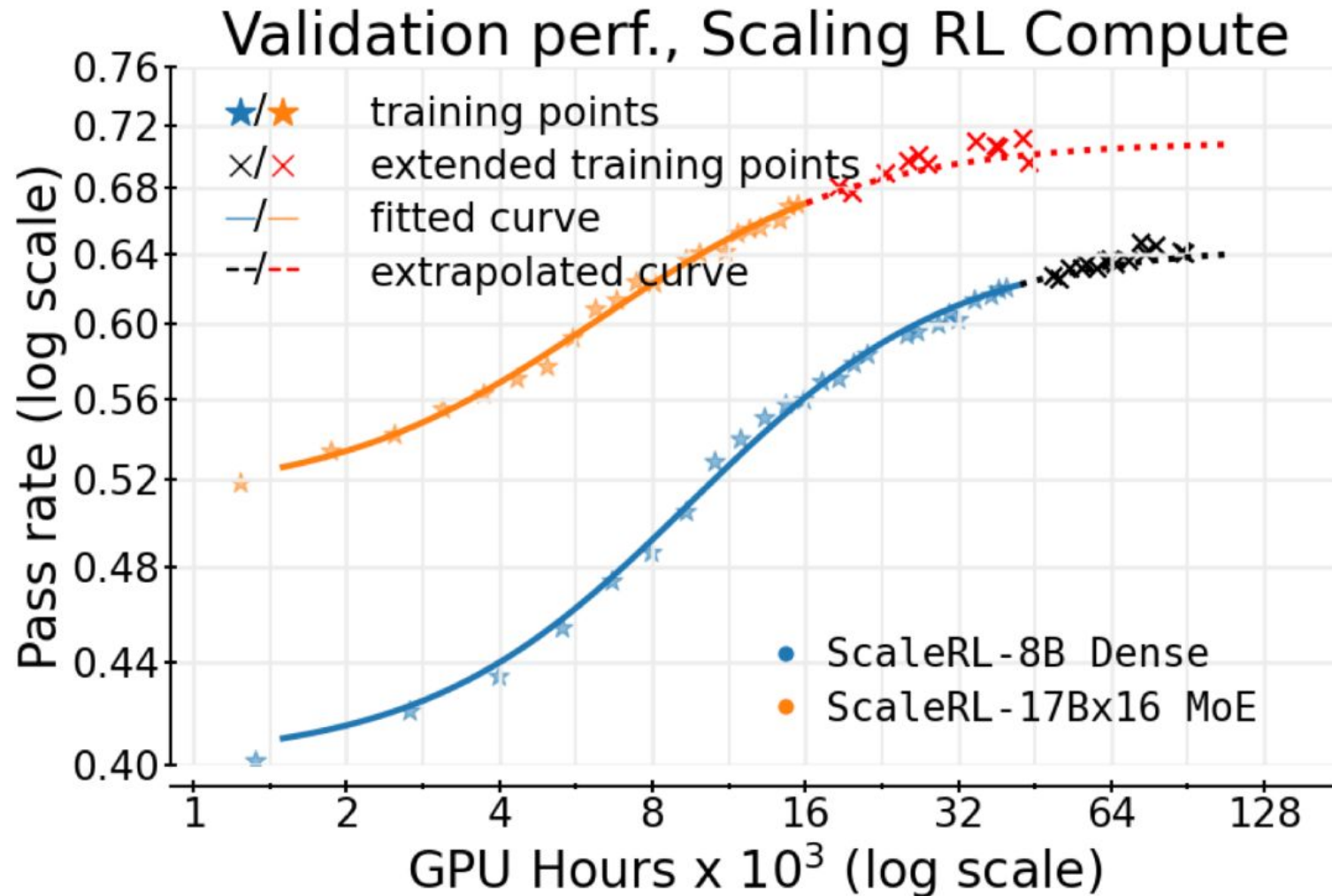
ScaleRL recipe

Cannot fix a low-ceiling recipe by training for longer!

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

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

Predictions at O(100k) GPU-hours

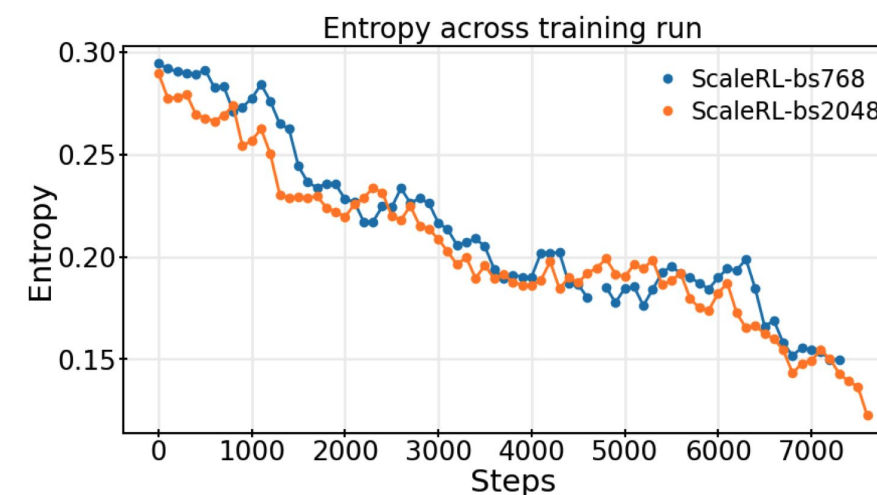
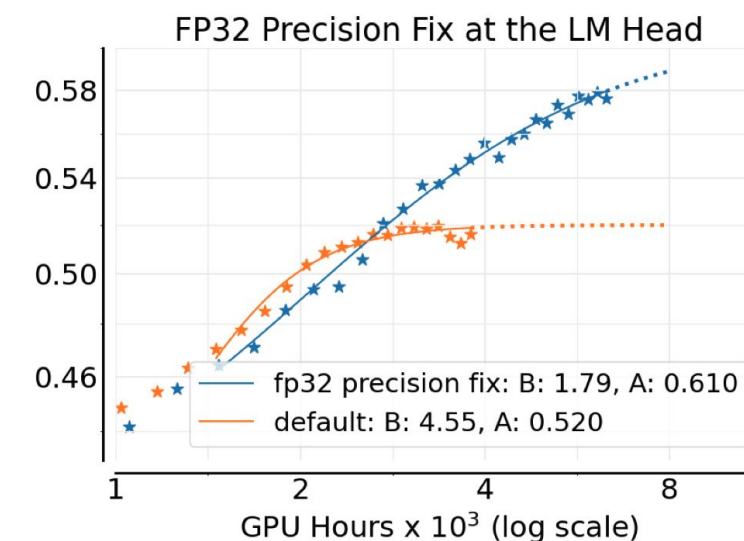
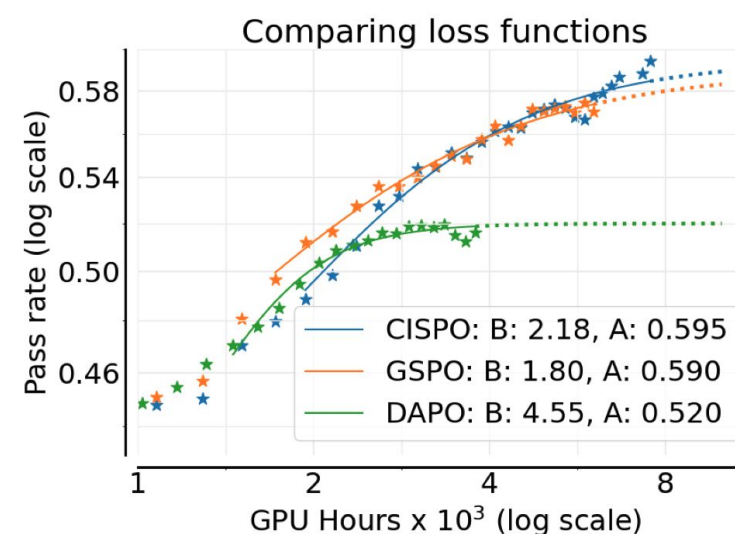


Granular Findings

- Higher asymptote
 - Loss function (CISPO over DAPO/GRPO), numerics (FP32 at the output head)

- Training efficiency
 - Off-policy (PipelineRL), data curriculum

- Scaling obvious axes
 - Generation length, batch size, model size



Thank You!

<https://arxiv.org/abs/2510.13786>

