

Self-Improvement in Language Models: The Sharpening Mechanism

Huang, Block, Foster, Rohatgi, Zhang, Simchowitz, Ash, Krishnamurthy

Some slides & pictures adapted from Akshay Krishnamurthy

LLM Training

Pre-training

The quick brown fox



fox

Post-training

Who proved symmetry implies conservation?



Emmy Noether

Lionel Messi

Inference-time
compute

Think step by step



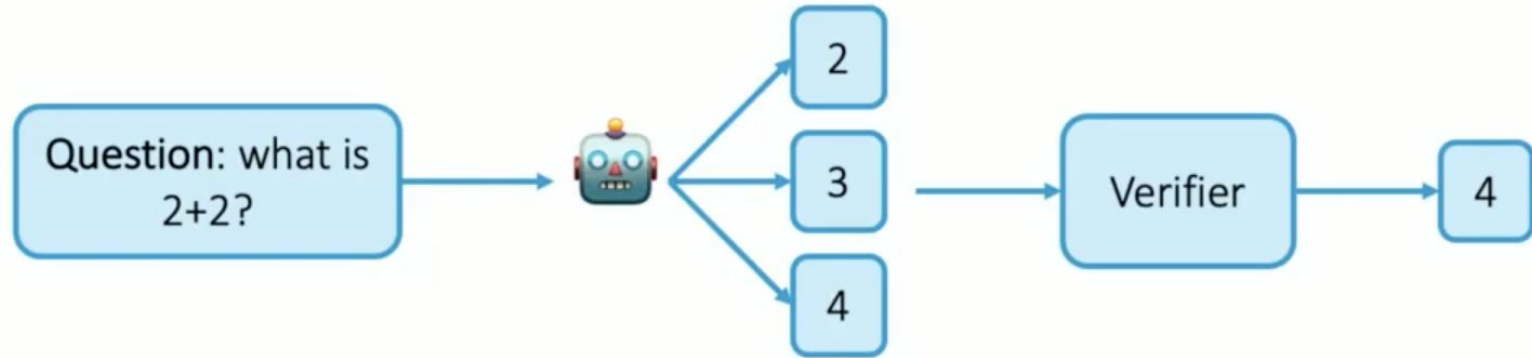
- Chain of thought
- Few-shot prompting
- Model-as-verifier
- Best-of-N

Compute + data heavy

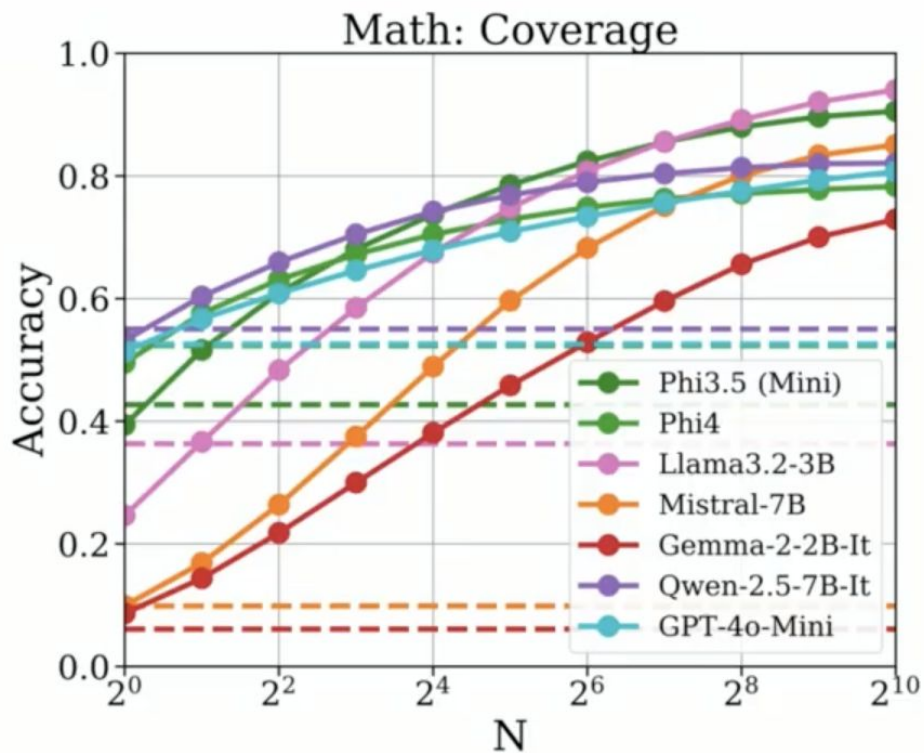
Data scarce

No data, just compute

Inference Time Compute - BoN



Inference Time Compute - BoN



**~32 samples for
good performance**

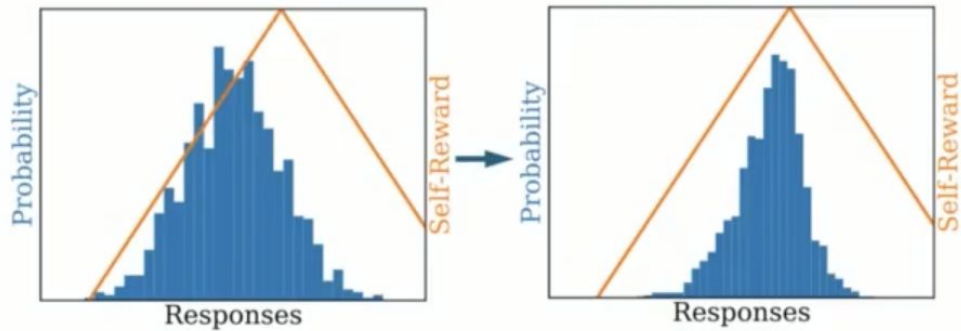
Self Improvement

- Model is improving performance using its own outputs
- Post-processing doesn't provide new information
- Using compute to uncover answer it already knows

Sharpening

Hypothesis: using LLM-as-verifier improves downstream performance

Mechanism: LLM-as-verifier tilts the distribution to output better responses



Self-reward and sharpening

A language model is a distribution $\pi_{\text{base}}(y | x)$ over responses given a prompt

Self reward $r_{\text{self}}(y | x, \pi_{\text{base}})$ is a function of response, prompt, and LM, examples:

- **Sequence level log-probabilities: $\log \pi_{\text{base}}(y | x)$**
- Length-normalized log-probabilities: $\log \pi_{\text{base}}(y | x) / |y|$
- Majority voting
- Prompting LM to score its own response (LM-as-judge)

Self-reward and sharpening

A language model is a distribution $\pi_{\text{base}}(y | x)$ over responses given a prompt

Self reward $r_{\text{self}}(y | x, \pi_{\text{base}})$ is a function of response, prompt, and LM, examples:

- **Sequence level log-probabilities: $\log \pi_{\text{base}}(y | x)$**
- Length-normalized log-probabilities: $\log \pi_{\text{base}}(y | x) / |y|$
- Majority voting
- Prompting LM to score its own response (LM-as-judge)

Sharpening: $\pi^{\dagger}(x) \approx \operatorname{argmax}_y r_{\text{self}}(y; x, \pi_{\text{base}})$

Note: Theory is developed for $\log \pi_{\text{base}}(y | x)$ only as this admits a clean closed-form optimal policy $\pi^*(y|x) \propto \pi_{\text{base}}(y|x)^{(1+1/\beta)}$. Length-normalized and other rewards are tested empirically but lack theoretical guarantees.

(ε, δ) -Sharpened Model

A model $\hat{\pi}$ is (ε, δ) -sharpened if

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(y^*(x) | x) \geq 1 - \delta] \geq 1 - \varepsilon$$

where $y^*(x) = \operatorname{argmax}_y \log \pi_{\text{base}}(y | x)$

Key: if $\delta < \frac{1}{2}$, then $\hat{\pi}(y^(x) | x) > \frac{1}{2} \Rightarrow$ greedy decoding recovers y^**

This is NP hard! Convert computational problem to statistical problem

Inference Time Sharpening with BoN

- Generate N responses
- Return the response with the highest rewards

Training Time Sharpening with BoN

- Generate N responses
- Return the response with the highest rewards
- **Amortize cost by SFT on the best of the N responses**

Sample & Evaluate Framework

- Sample n prompts from the base model
- For each prompt sample N responses
- Query the reward model
- Sample complexity = $n.N$

Coverage

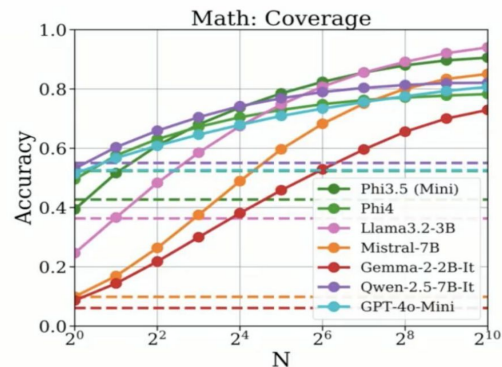
Coverage Coefficient

$$C_{\text{cov}} = \mathbb{E}_{x \sim \mu} [(\pi_{\text{base}}(y^*(x) | x))^{-1}]$$

Small C_{cov} : π_{base} often samples y^ \Rightarrow easy to sharpen*

Large C_{cov} : π_{base} rarely samples y^ \Rightarrow exponentially hard*

Large language monkeys



In practice $\text{CoV} \ll Y = |V|^{\text{length}}$

SFT Sharpening Guarantees

Upper Bound (Theorem 3.2)

SFT-Sharpener achieves (ϵ, δ) -sharpening with sample complexity

$$m = O\left(\frac{C_{\text{cov}} \cdot \log|\Pi| \cdot \log(\delta^{-1})}{\delta \epsilon^2}\right)$$

Lower Bound (Theorem 3.3)

Any algorithm requires

$$m \gtrsim \frac{C_{\text{cov}}}{\epsilon^2} \text{ samples}$$

Greedy Sufficiency (Prop. 3.1)

If $\delta < \frac{1}{2}$, then $\hat{\pi}(y^*(x) | x) > \frac{1}{2}$

\Rightarrow greedy decoding recovers y^* exactly

SFT-Sharpener is minimax optimal

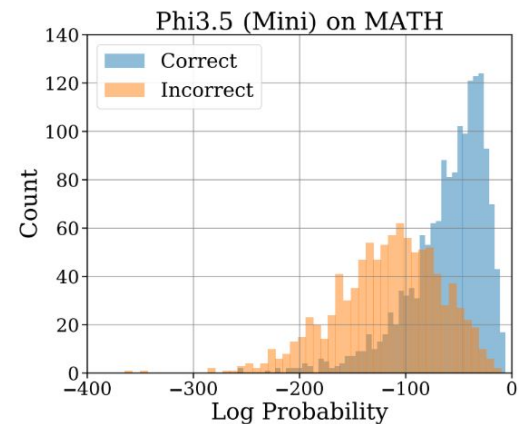
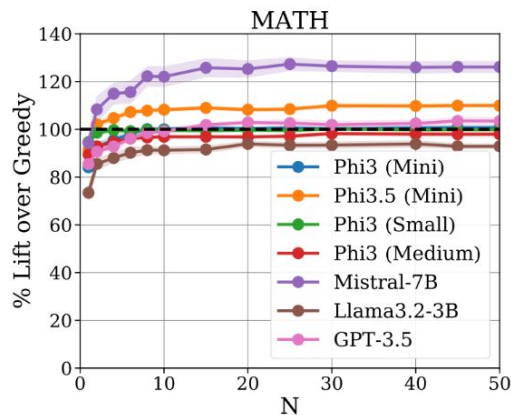
Matches the lower bound up to log factors — but cannot escape dependence on C_{cov}

$$C_{\text{cov}} = \mathbb{E}_{x \sim \mu}[(\pi_{\text{base}}(y^*(x) | x))^{-1}] \quad \text{— small when base model easily samples the best response}$$

Results on Maximum Likelihood Sharpening

BoN: % Lift over Greedy

Phi3 (Mini)	0.7	-2.8	17.0	-4.2	4.5	3.8
Phi3.5 (Mini)	9.9	4.1	3.2	5.9	6.0	-0.9
Phi3 (Small)	-0.2	4.0	-6.8	1.4	0.5	-9.2
Phi3 (Medium)	-2.1	-4.4	1.8	-1.5	3.3	-0.7
Mistral-7B	26.1	10.5	-1.1	11.8	28.6	10.2
Llama3.2-3B	-7.1	-2.6	3.0	-14.9	126.7	6.7
GPT-3.5	3.4	-11.8	-10.4	-0.6	6.7	-7.0
	MATH	GSM8K	ProntoQA	Bio	Phys	Chem



Power of Exploration

Optimization Objective at Iteration t

$$\pi^{(t+1)} \leftarrow \arg \min_{\pi}$$

DPO squared loss

$$-\sum_{D^{(t)}} \left(\beta \log \frac{\pi(y|x)}{\pi_{\text{base}}(y|x)} - \beta \log \frac{\pi(y'|x)}{\pi_{\text{base}}(y'|x)} - (r(x, y) - r(x, y')) \right)^2$$

$$+ \alpha \sum_{D^{(t)}} \log \pi(y'|x)$$

optimism bonus

Dataset Construction

Chosen: $y \sim \pi^{(t)}(\cdot|x)$ (current policy)

Rejected: $y' \sim \pi_{\text{base}}(\cdot|x)$ (base model)

Reward: $r(x, y) = \log \pi_{\text{base}}(y|x)$

What Each Term Does

DPO loss: push π toward high-reward sequences

Optimism: suppress π_{base} samples \Rightarrow explore new regions

Together: find high-reward sequences π_{base} never generates

Power of Exploration

XPO Sample Complexity (Theorem 4.3)

XPO achieves (ϵ, δ) -sharpening with sample complexity

$$m = \tilde{O}\left(\frac{\text{SEC}(\Pi) \cdot \log|\Pi|}{\gamma_{\text{margin}}^2 \delta^2 \epsilon^2}\right) \quad \text{— no dependence on } C_{\text{cov}}$$

SFT-Sharpener vs XPO

SFT-Sharpener

$$\text{Complexity: } \frac{C_{\text{cov}} \cdot \log|\Pi|}{\delta \epsilon^2}$$

Linear softmax: $C_{\text{cov}} = \exp(\Omega(d))$

Offline — stuck with coverage of base model

Cannot escape coverage barrier

XPO

$$\text{Complexity: } \frac{\text{SEC}(\Pi) \cdot \log|\Pi|}{\gamma_{\text{margin}}^2 \delta^2 \epsilon^2}$$

Linear softmax: $\text{SEC}(\Pi) = \tilde{O}(d)$

Online — actively explores beyond base model

Exponential improvement over SFT

For linear softmax models: $\text{SEC}(\Pi) = \tilde{O}(d)$ vs $C_{\text{cov}} = \exp(\Omega(d))$ — **an exponential separation**

SEC measures model class complexity from an exploration perspective, not base model quality