

BRiTE: Bootstrapping Reinforced Thinking Process to Enhance Language Model Reasoning



Han Zhong*¹, Yutong Yin*², Shenao Zhang*², Xiaojun Xu*³, Yuanxin Liu*², Yifei Zuo*², Zhihan Liu*²,
Boyi Liu³, Sirui Zheng², Hongyi Guo², Liwei Wang¹, Mingyi Hong⁴, Zhaoran Wang²

¹Peking University, ²Northwestern University, ³ByteDance, ⁴University of Minnesota

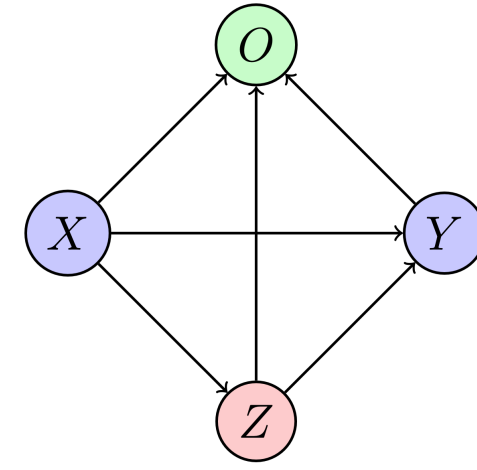


Reasoning as a Graphical Model

Q: What is reasoning in large language models?

A: *Okay, so I need to figure out what reasoning is in large language models (LLMs) is. Let me start by breaking down the question. The user is asking about reasoning ...*

Reasoning in large language models (LLMs) refers to their ability to generate responses that mimic structured logical thought processes to solve problems or answer questions.



$$\mathbb{P}(y, o | x, \theta) = \mathbb{P}(z | x, \theta) \cdot \mathbb{P}(y | x, z, \theta) \cdot \mathbb{P}(o | x, z, \theta)$$

Bootstrap Reinforced Thinking Process

$$\begin{aligned} \mathcal{L}(\theta) &= \log \sum_{(z,y,o) \in \mathcal{Z} \times \mathcal{Y} \times \mathcal{O}} \mathbb{P}(z, y, o | x, \theta) \\ &= \max_{\mathbb{Q} \in \Delta} \left\{ \underbrace{\sum_{(z,y,o)} \log \mathbb{P}(z, y, o | x, \theta) \mathbb{Q}(z, y, o | x, \psi) - \sum_{(z,y,o)} \log \mathbb{Q}(z, y, o | x, \psi) \mathbb{Q}(z, y, o | x, \psi)}_{:= \mathcal{L}_{\psi}(\theta)} \right\} \end{aligned}$$

Maximize $\mathcal{L}(\theta)$ (difficult) \implies Maximize evidence lower bound $\mathcal{L}_{\psi}(\theta)$ (easy)

BRiTE — An EM-type Algorithm

$$\begin{aligned} \mathbb{Q}(z, y, o | x, \psi_{t+1}) &\leftarrow \arg\max_{\mathbb{Q}} \mathcal{L}_{\psi_t}(\theta_t) \\ \text{Thought proposer} &= \frac{\mathbb{P}(z, y, o | x, \theta_t)}{\sum_{(z,y,o)} \mathbb{P}(z, y, o | x, \theta_t)} \\ \text{E} & \\ \theta_{t+1} &= \arg\max_{\theta} \mathcal{L}_{\psi_{t+1}}(\theta) \\ \text{M} &= \arg\max_{\theta} \left\{ \sum_{(z,y,o) \in \mathcal{Z} \times \mathcal{Y} \times \mathcal{O}} \log \mathbb{P}(z, y, o | x, \theta) \cdot \mathbb{Q}(z, y, o | x, \psi_{t+1}) \right\} \end{aligned}$$

Assumptions: $f_{\theta} \in \mathcal{H}$ for a certain RKHS; $\mathbb{P}(z, y | x, \theta) \propto \exp(f_{\theta}(x, z, y))$

Theorem: convergence to optima

$$\min_{1 \leq t \leq T} \left\{ \log \frac{\mathbb{P}(x \in \mathcal{X}, y \in \mathcal{Y}, o \in \mathcal{O} | x, \theta^*)}{\mathbb{P}(x \in \mathcal{X}, y \in \mathcal{Y}, o \in \mathcal{O} | x, \theta_t)} \right\} \leq \frac{\mathbb{D}_{\text{KL}}(\mathbb{P}(\cdot | x, \theta_1) \| \mathbb{P}(\cdot | x, \theta^*))}{T}$$

Concrete Examples of BRiTE

$$\begin{aligned} \text{Scope:} & \quad \begin{aligned} &\bullet \quad o \in \{0,1\}, \mathcal{O} = \{1\} \\ &\bullet \quad \mathcal{Y} \text{ is the response space} \\ &\bullet \quad \mathcal{Z} \text{ is the latent space} \end{aligned} \\ & \quad \begin{aligned} &+ \quad \mathbb{P}(o = 1 | x, z, y) := \exp(R(x, z, y)/\beta) \\ &+ \quad \mathbb{P}(z, y, o = 1 | x, \theta) = \mathbb{P}(z, y | x, \theta) \mathbb{P}(o = 1 | x, z, y) \\ &+ \quad \mathbb{Q}(z, y | x, \psi) := \mathbb{Q}(z, y, o = 1 | x, \psi) \end{aligned} \end{aligned}$$

Example (PPO)

$$\begin{aligned} \mathcal{L}_{\psi}(\theta) &= \sum_{(z,y)} \log \mathbb{P}(z, y, o = 1 | x, \theta) \mathbb{Q}(z, y | x, \psi) \\ &\quad - \sum_{(z,y)} \log \mathbb{Q}(z, y | x, \psi) \mathbb{Q}(z, y | x, \psi) \\ &= \mathbb{E}_{(z,y) \sim \mathbb{Q}} \left[R(x, z, y)/\beta - \log \frac{\mathbb{Q}(z, y | x, \psi)}{\mathbb{P}(z, y | x, \theta)} \right] \end{aligned}$$

$$\begin{aligned} \text{Scope:} & \quad \begin{aligned} &\bullet \quad o \in \{0,1\}, \mathcal{O} = \{1\} \\ &\bullet \quad \mathcal{Y} \text{ is the response space} \\ &\bullet \quad \mathcal{Z} \text{ is the latent space} \end{aligned} \\ & \quad \begin{aligned} &+ \quad \mathbb{P}(o = 1 | x, z, y) := \mathbb{I}(\text{y is correct for } x) \text{ or } \exp(R(x, y)/\beta) \\ &+ \quad \mathbb{P}(z, y, o = 1 | x, \theta) = \mathbb{P}(z, y | x, \theta) \mathbb{P}(o = 1 | x, z, y) \\ &+ \quad \mathbb{Q}(z, y | x, \psi) := \mathbb{Q}(z, y, o = 1 | x, \psi) \end{aligned} \end{aligned}$$

$$\begin{aligned} &\max_{\mathbb{P}} \left\{ \mathbb{E}_{(z,y) \sim \mathbb{P}(\cdot, \cdot | x, \theta_t)} \left[\log \mathbb{P}(z, y | x, \theta) \cdot \mathbb{I}(\text{y is correct for } x) \right] \right\} \\ &\max_{\mathbb{P}} \left\{ \mathbb{E}_{(z,y) \sim \mathbb{P}(\cdot, \cdot | x, \theta_t)} \left[\log \mathbb{P}(z, y | x, \theta) \cdot \exp(R(x, y)/\beta) \right] \right\} \end{aligned}$$

If $\mathcal{Z} = \emptyset$, then it recovers **STaR** and **Reject Sampling Finetuning** or **RestEM**

Learning Intractable Posterior via RL

$$\mathbb{Q}(z, y, o | x, \psi) \leftarrow \arg\max_{\mathbb{Q}} \mathcal{L}_{\psi}(\theta) = \frac{\mathbb{P}(z, y, o | x, \theta)}{\sum_{(z,y,o)} \mathbb{P}(z, y, o | x, \theta)}$$

Intractable

Lemma: the optimal policy for an entropy-regularized token-level MDP

$$\pi^*(a_h \cup \{(s_i, a_i)\}_{i=h+1}^H | s_h) \propto \exp\left(\frac{1}{\beta} \sum_{i=h}^H r(s_i, a_i)\right)$$

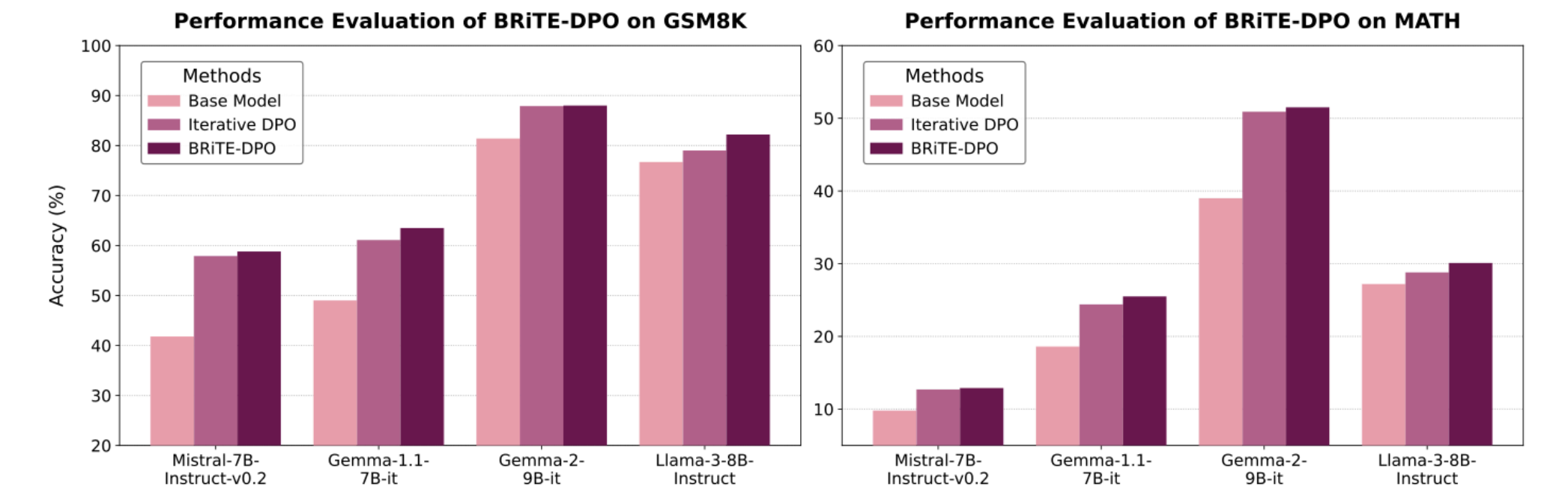
$$\text{Set } \frac{1}{\beta} \sum_{i=0}^H r(s_i, a_i) = \log \mathbb{P}(z, y, o | x, \theta)! \text{ Then } \pi^*(s_H | s_0) = \mathbb{Q}(z, y, o | x, \psi)$$

Experiments

- BRiTE Significantly Improves Existing Rejection Sampling Algorithms.
- BRiTE \geq SFT with Human-Annotated Thinking Process.

Algorithm	Mistral-7B-Instruct-v0.2		Gemma-1.1-7B-it		Gemma-2-9B-it		Llama-3-8B-Instruct	
	GSM8K	MATH	GSM8K	MATH	GSM8K	MATH	GSM8K	MATH
—	41.8	9.8	49.0	18.8	81.3	37.3	79.2	28.3
SFT	52.8	13.6	57.5	19.6	80.1	41.5	72.6	27.1
RS	47.7	10.3	58.4	18.7	87.6	47.5	79.5	28.9
BRiTE	52.2	11.2	59.2	23.7	89.7	50.5	81.0	30.0

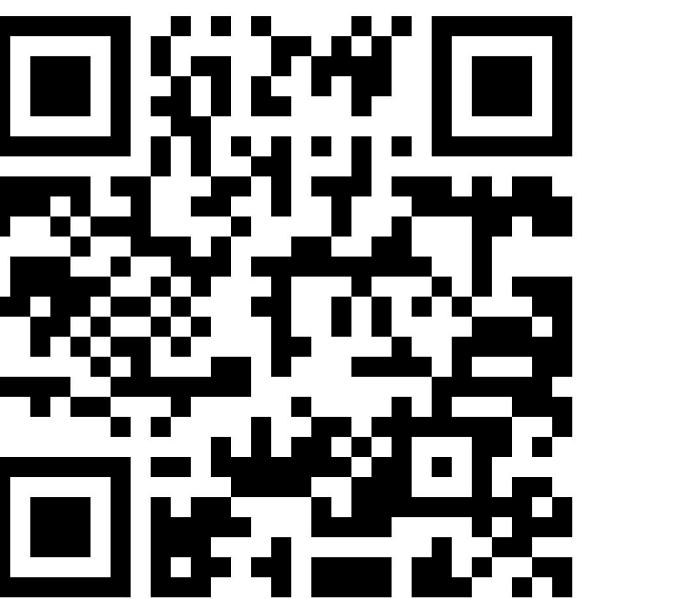
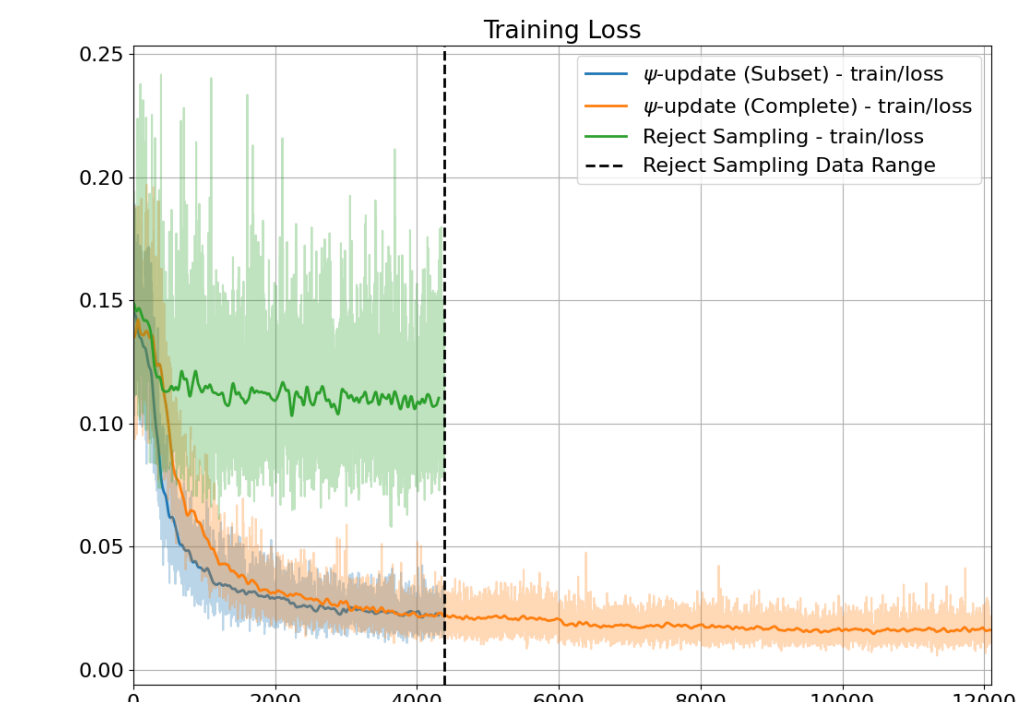
- BRiTE Enhances the Reasoning Capacity in RLHF Stage.



- BRiTE Improves Code Generation Ability.

Algorithm	HumanEval		BCB (Instruct)	
	Basic (%)	Plus (%)	Hard (%)	Full (%)
—	78.0	70.7	10.1	35.5
SFT	78.0	67.7	11.5	37.2
RS	79.3	73.2	11.5	35.6
BRiTE	81.7	72.6	15.5	36.3

- BRiTE Generates High Quality Trajectories for Distillation.



Scan to read our paper!