# 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 Liu[3], Sirui Zheng[2], Hongyi Guo[2], Liwei Wang[1], Mingyi Hong[4], Zhaoran Wang[2]

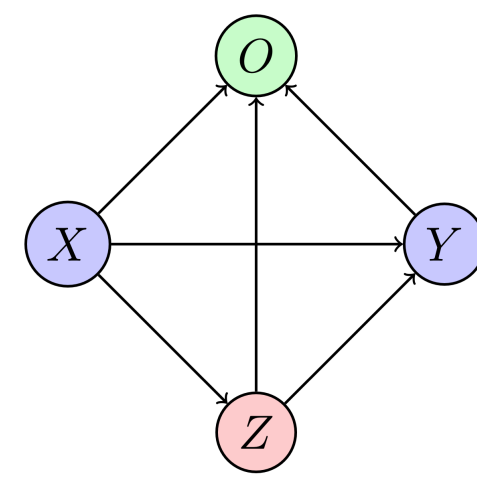[1]Peking University, [2]Northwestern University, [3]ByteDance, [4]University of Minnesota

## 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.



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

## Bootstrap Reinforced Thinking Process

$$\mathcal{L}(\theta) = \log \sum_{(z,y,o) \in \mathcal{Z} \times \mathcal{Y} \times \mathcal{O}} \mathbb{P}(z, y, o \mid x, \theta)$$

$$= \max_{\mathbb{Q} \in \Delta} \left\{ \sum_{(z,y,o)} \log \mathbb{P}(z, y, o \mid x, \theta) \mathbb{Q}(z, y, o \mid x, \psi) - \sum_{(z,y,o)} \log \mathbb{Q}(z, y, o \mid x, \psi) \mathbb{Q}(z, y, o \mid x, \psi) \right\}$$

$$:= \mathcal{L}_\psi(\theta)$$

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

### BRiTE — An EM-type Algorithm

$$\mathbb{Q}(z, y, o \mid x, \psi_{t+1}) \leftarrow \text{argmax}_{\mathbb{Q}} \mathcal{L}_\psi(\theta_t)$$
Thought proposer
$$= \frac{\mathbb{P}(z, y, o \mid x, \theta_t)}{\sum_{(z,y,o)} \mathbb{P}(z, y, o \mid x, \theta_t)}$$
E

$$\theta_{t+1} = \text{argmax}_\theta \mathcal{L}_{\psi_{t+1}}(\theta)$$
$$= \text{argmax}_\theta \left\{ \sum_{z,y,o \in \mathcal{Z} \times \mathcal{Y} \times \mathcal{O}} \log \mathbb{P}(z, y \mid x, \theta) \cdot \mathbb{Q}(z, y, o \mid x, \psi_{t+1}) \right\}$$
M

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

**Theorem**: convergence to optima

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

## Concrete Examples of BRiTE

**Scope**:
- $o \in \{0,1\}$, $\mathcal{O} = \{1\}$
- $\mathcal{Y}$ is the response space
- $\mathcal{Z}$ is the latent space

- $\mathbb{P}(o = 1 \mid x, z, y) := \exp(R(x, z, y)/\beta)$
- $\mathbb{P}(z, y, o = 1 \mid x, \theta) = \mathbb{P}(z, y \mid x, \theta)\mathbb{P}(o = 1 \mid x, z, y)$
- $\mathbb{Q}(z, y \mid x, \psi) := \mathbb{Q}(z, y, o = 1 \mid x, \psi)$

$$\mathcal{L}_\psi(\theta) = \sum_{(z,y)} \log \mathbb{P}(z, y, o = 1 \mid x, \theta)\mathbb{Q}(z, y \mid x, \psi)$$
$$- \sum_{(z,y)} \log \mathbb{Q}(z, y \mid x, \psi)\mathbb{Q}(z, y \mid x, \psi)$$

**Example** (PPO)

$$= \mathbb{E}_{(z,y) \sim \mathbb{Q}} \left[ R(x, z, y)/\beta - \log \frac{\mathbb{Q}(z, y \mid x, \psi)}{\mathbb{P}(z, y \mid x, \theta)} \right]$$

**Scope**:
- $o \in \{0,1\}$, $\mathcal{O} = \{1\}$
- $\mathcal{Y}$ is the response space
- $\mathcal{Z}$ is the latent space

- $\mathbb{P}(o = 1 \mid x, z, y) := \mathbb{I}(y \text{ is correct for } x)$ or $\exp(R(x, y)/\beta)$
- $\mathbb{P}(z, y, o = 1 \mid x, \theta) = \mathbb{P}(z, y \mid x, \theta)\mathbb{P}(o = 1 \mid x, z, y)$
- $\mathbb{Q}(z, y \mid x, \psi) := \mathbb{Q}(z, y, o = 1 \mid x, \psi)$

$$\max_{\mathbb{P}} \left\{ \mathbb{E}_{(z,y) \sim \mathbb{P}(\cdot, \cdot \mid x, \theta_t)} \left[ \log \mathbb{P}(z, y \mid x, \theta) \cdot \mathbb{I}(y \text{ is correct for } x) \right] \right\}$$

$$\max_{\mathbb{P}} \left\{ \mathbb{E}_{(z,y) \sim \mathbb{P}(\cdot, \cdot \mid x, \theta_t)} \left[ \log \mathbb{P}(z, y \mid x, \theta) \cdot \exp(R(x, y)/\beta) \right] \right\}$$

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

## Learning Intractable Posterior via RL

$$\mathbb{Q}(z, y, o \mid x, \psi) \leftarrow \text{argmax}_{\mathbb{Q}} \mathcal{L}_\psi(\theta) = \frac{\mathbb{P}(z, y, o \mid x, \theta)}{\sum_{(z,y,o)} \mathbb{P}(z, y, o \mid x, \theta)}$$
Intractable

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

$$\pi^\star(a_h \cup \{(s_i, a_i)\}_{i=h+1}^H \mid s_h) \propto \exp\left(\frac{1}{\beta} \sum_{i=h}^H r(s_i, a_i)\right)$$
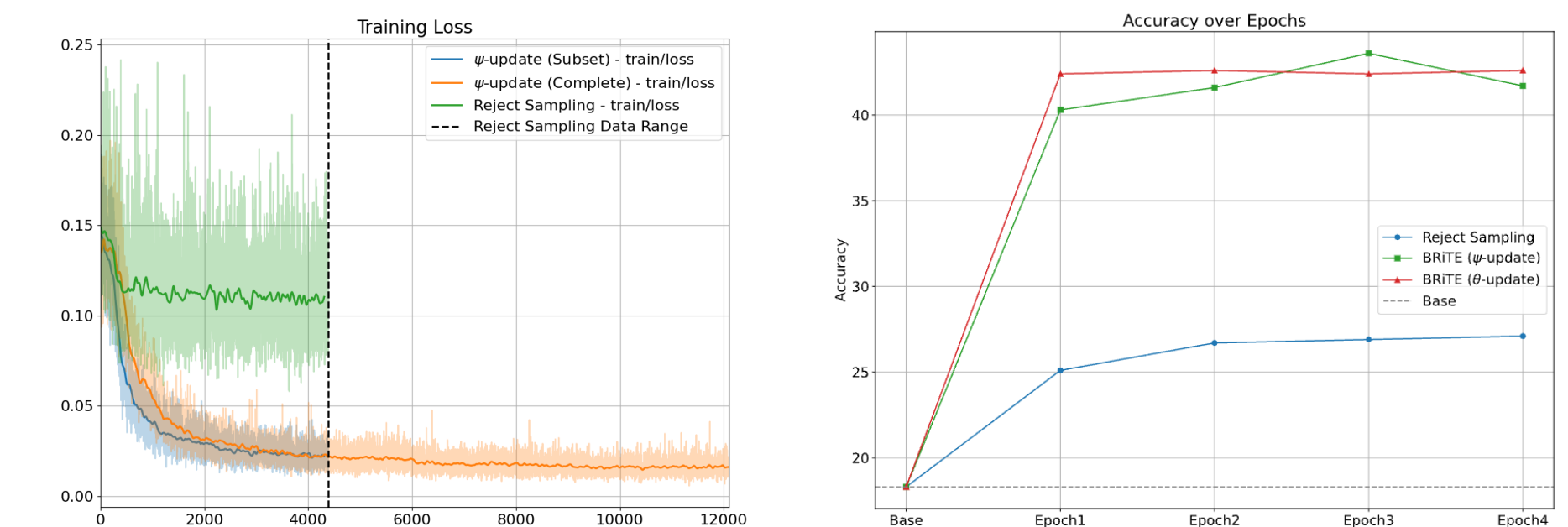
Set $\frac{1}{\beta} \sum_{i=0}^H r(s_i, a_i) = \log \mathbb{P}(z, y, o \mid x, \theta)$! Then $\pi^\star(s_H \mid s_0) = \mathbb{Q}(z, y, o \mid x, \psi)$
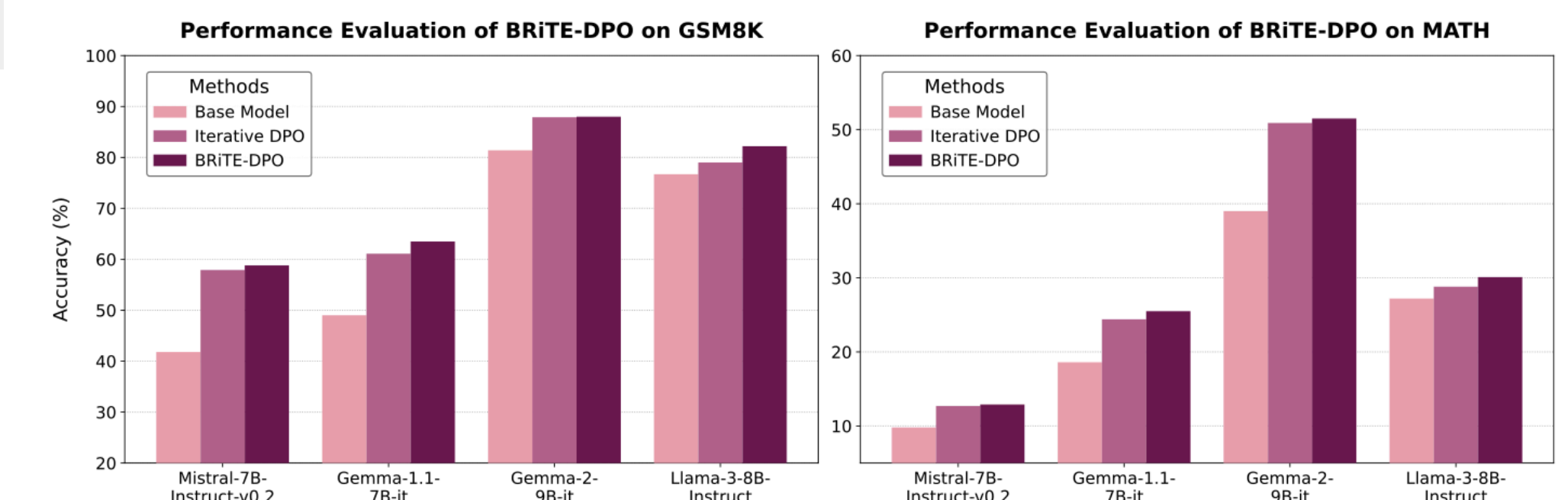
## Experiments

1. BRiTE Significantly Improves Existing Rejection Sampling Algorithms.
2. BRiTE ≥ SFT with Human-Annotated Thinking Process.

| Method | MATH500 | Minerva Math | OlympiadBench | AIME24 | AMC23 | GPQA Diamond |
|---|---|---|---|---|---|---|
| — | 44.1 | 12.9 | 16.1 | 0.9 | 10.1 | 25.9 |
| RS | 54.3 | 21.0 | 23.1 | 5.6 | 31.6 | 26.9 |
| BRiTE ($\psi$-update) | 79.1 | 35.0 | 35.7 | 14.3 | 57.7 | 28.5 |
| BRiTE ($\theta$-update) | 76.9 | 40.6 | 37.0 | 14.4 | 57.1 | 29.8 |
| BRiTE-iter-2 ($\psi$-update) | 80.6 | 41.3 | 37.3 | 14.3 | 57.9 | 29.9 |
| BRiTE-iter-2 ($\theta$-update) | 78.2 | 39.8 | 37.9 | 15.3 | 56.4 | 30.1 |

3. BRiTE Generates High Quality Trajectories for Distillation.



4. BRiTE Enhances the Reasoning and Coding Capacity in RLHF Stage.



| 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 |