

# Overcoming Joint Intractability with Lossless Hierarchical Speculative Decoding

Yuxuan Zhou<sup>1\*</sup>, Fei Huang<sup>2</sup>, Heng Li<sup>1</sup>, Fengyi Wu<sup>3</sup>, Tianyu Wang<sup>3</sup>, Jianwei Zhang<sup>2</sup>, Junyang Lin<sup>2†</sup>, Zhi-Qi Cheng<sup>3†</sup>

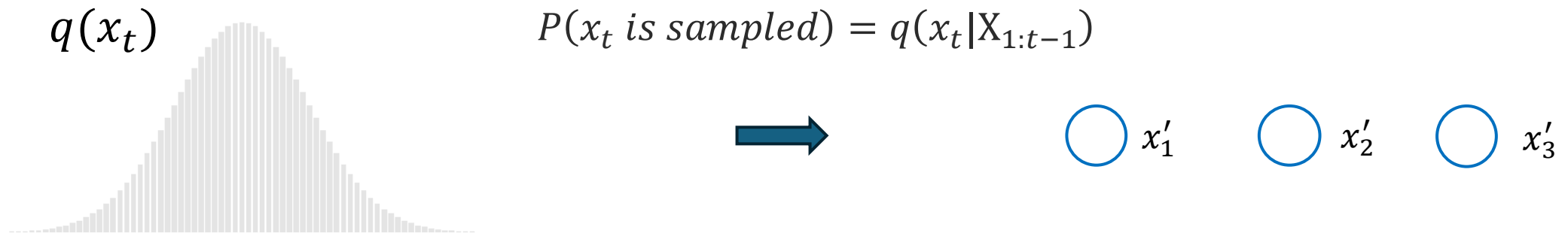
<sup>1</sup>Independent Researcher    <sup>2</sup>Qwen Team, Alibaba Inc.    <sup>3</sup>University of Washington  
zhouyuxuanyx@gmail.com, junyang.ljy@alibaba-inc.com, zhiqics@uw.edu

\*Work done during internship at Qwen Team, Alibaba Inc.

†Corresponding author.

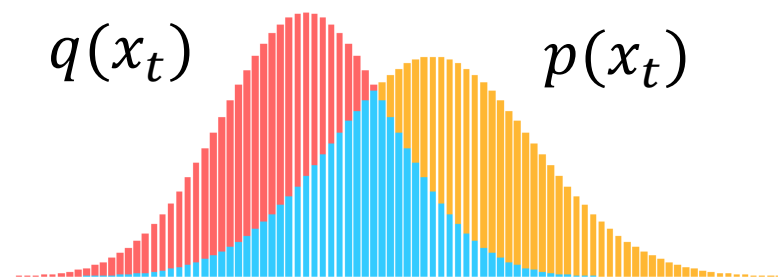
## Tokenwise Verification [1]

- **Step 1:** sample draft tokens from the draft distribution



## Tokenwise Verification

➤ **Step 2:** verify a draft token based on likelihood ratio

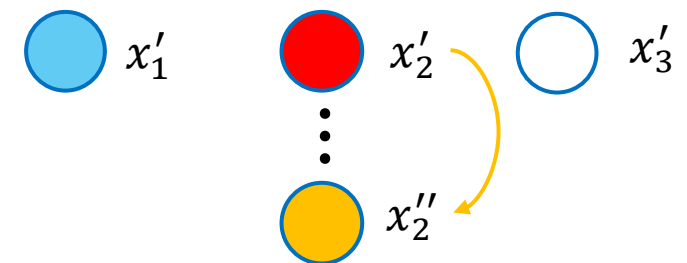


$$\text{Accept Probability: } \min\left\{1, \frac{p(x_t)}{q(x_t)}\right\}$$

**Case 1:** ✓ accept  $\Rightarrow$  check next token

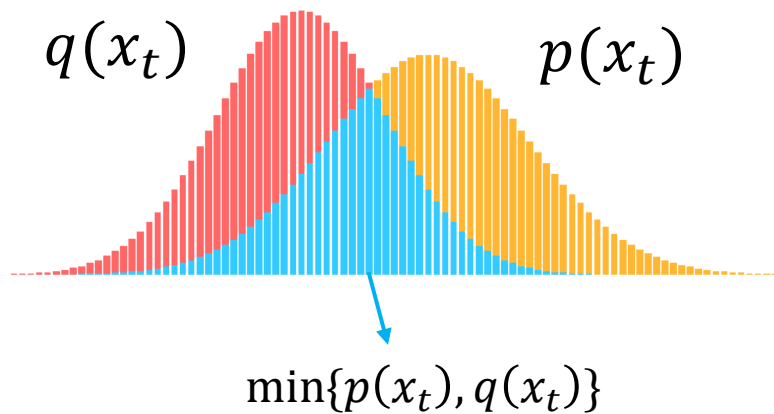


**Case 2:** ✗ reject  $\Rightarrow$  resample



## Tokenwise Verification

- **Step 2:** verify a token based on likelihood ratio
  - **Case 1:** token is accepted



Accept Probability:

$$h = \min\left\{\frac{p(x_t)}{q(x_t)}, 1\right\}$$



- **Output Probability:**

$$P(x_t \text{ is sampled, } x_t \text{ is accepted})$$

$$= q(x_t) \min\left\{\frac{p(x_t)}{q(x_t)}, 1\right\}$$

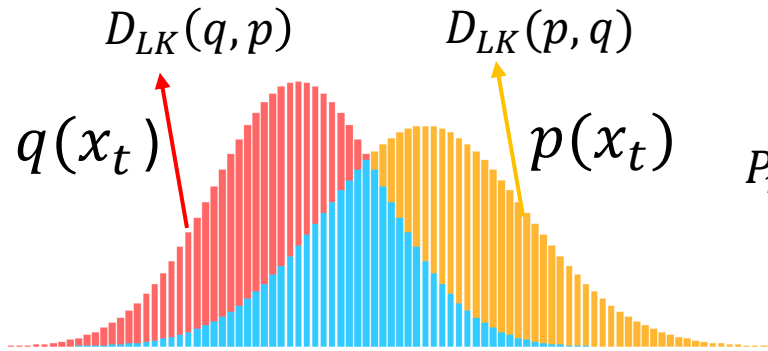
$$= \min\{p(x_t), q(x_t)\}$$

blue area

# Tokenwise Verification

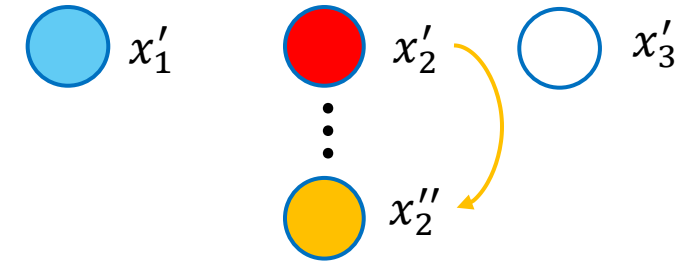
➤ **Step 2:** resample a token using likelihood ratio

- **Case 2:** token is rejected



Resample Probability:

$$P_{res}(x_t) = \frac{p(x_t) - \min\{p(x_t), q(x_t)\}}{\sum_{\tilde{x}_t \in \mathcal{V}} p(\tilde{x}_t) - \min\{p(\tilde{x}_t), q(\tilde{x}_t)\}}$$



➤ **Output Probability:**

$P(\tilde{x}_t \neq x_t \text{ is sampled and rejected, } x_t \text{ is resampled})$

$$= \sum_{\tilde{x}_t \in \mathcal{V}} q(\tilde{x}_t) \left(1 - \min\left\{\frac{p(x_t)}{q(x_t)}, 1\right\}\right) \frac{p(x_t) - \min\{p(x_t), q(x_t)\}}{\sum_{\tilde{x}_t \in \mathcal{V}} p(\tilde{x}_t) - \min\{p(\tilde{x}_t), q(\tilde{x}_t)\}}$$

$$= \underbrace{D_{LK}(q, p)}_{\text{red area}} \frac{p(x_t) - \min\{p(x_t), q(x_t)\}}{\underbrace{D_{LK}(p, q)}_{\text{orange area}}}$$

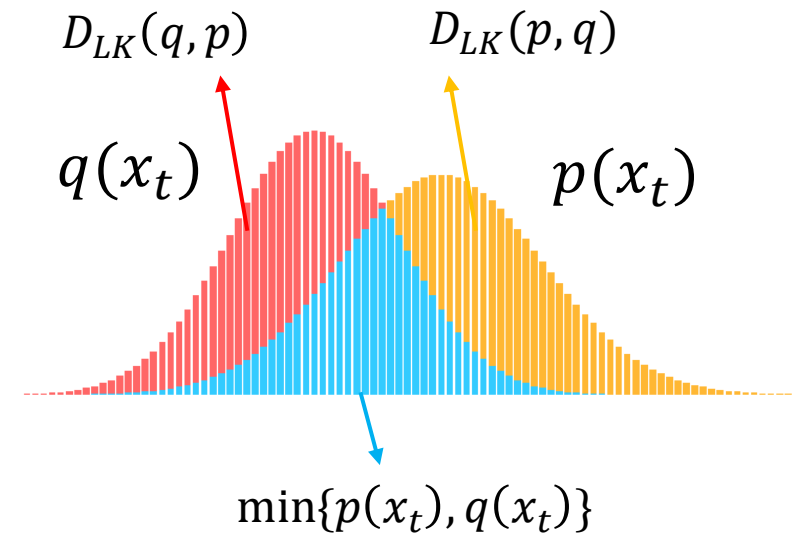
red area = orange area

## Tokenwise Verification

### ➤ Output Distribution:

$$\begin{aligned} P(x_t \text{ is yielded}) &= P(x_t \text{ is sampled and accepted}) + P(\tilde{x}_t \text{ is sampled and rejected, } x_t \text{ is resampled}) \\ &= \min\{p(x_t), q(x_t)\} + p(x_t) - \min\{p(x_t), q(x_t)\} \\ &= p(x_t) \end{aligned}$$

➔ **Target distribution is preserved!**



## Joint Verification

### ➤ Accept Probability [2]:

$$\text{tokenwise: } h(x_t) = \min\left\{1, \frac{p(x_t)}{q(x_t)}\right\} \longrightarrow \text{joint: } h(X_{1:t}) = \min\left\{1, \frac{p(X_{1:t})}{q(X_{1:t})}\right\}$$

### ➤ Illustrative Example:

“Eliza worked a total of 45 hours, so {she, work, ed, 45}”

input prompt                      draft tokens

### Draft and target probs:

$$\{p(x_t)\} = \{0.72, 1.00, 1.00, 1.00\}, \quad \{q(x_t)\} = \{0.88, 0.79, 0.65, 0.26\},$$

### Accept probs

$$\{h(x_t)\} = \{0.82, 1.00, 1.00, 1.00\}, \quad h(X_{1:4}) = 1$$

[2] Zongyue Qin, et, al. Optimized multitoken joint decoding with auxiliary model for llm inference. ICLR 2025.

## Joint Verification

### ➤ Output Distribution:

$$P(X_{1:t} \text{ is yielded}) = P(X_{1:t} \text{ is sampled and accepted}) + P(\tilde{X}_{1:t} \text{ is sampled and rejected, } X_{1:t} \text{ is resampled})$$

- **Case 1:**  $P(X_{1:t} \text{ is sampled and accepted}) = q(X_{1:t}) \min\left\{1, \frac{p(X_{1:t})}{q(X_{1:t})}\right\} = \underbrace{\min\{p(X_{1:t}), q(X_{1:t})\}}_{\text{blue area}}$

- **Case 2:**  $P(\tilde{X}_{1:t} \neq X_{1:t} \text{ is sampled and rejected, } X_{1:t} \text{ is resampled})$

$$= D_{Branch}(q, p | X_{1:t-1}) p_{res}(X_{1:t})$$

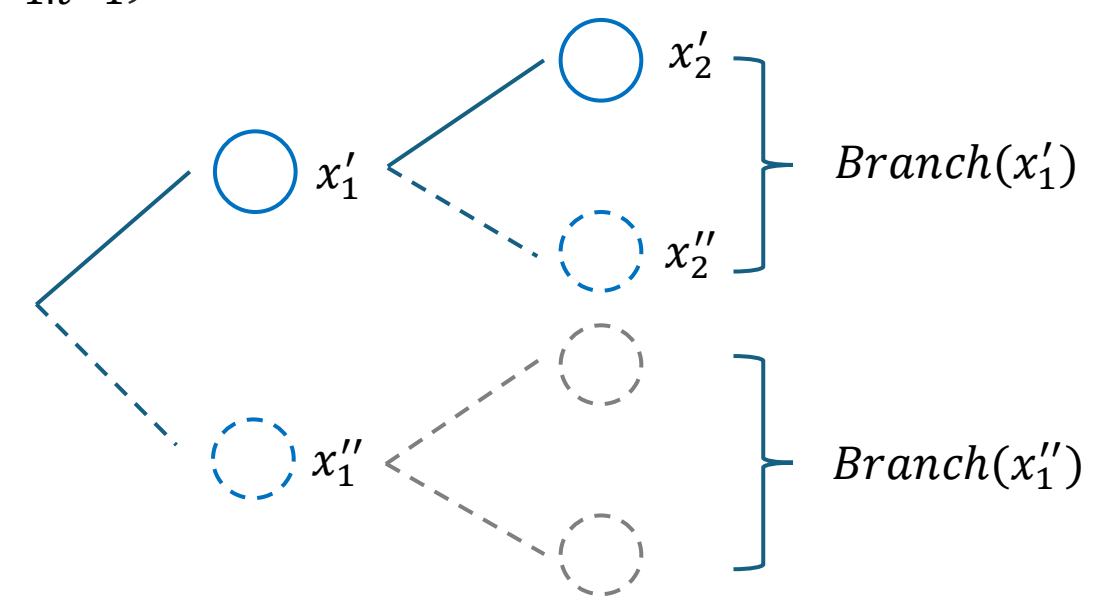
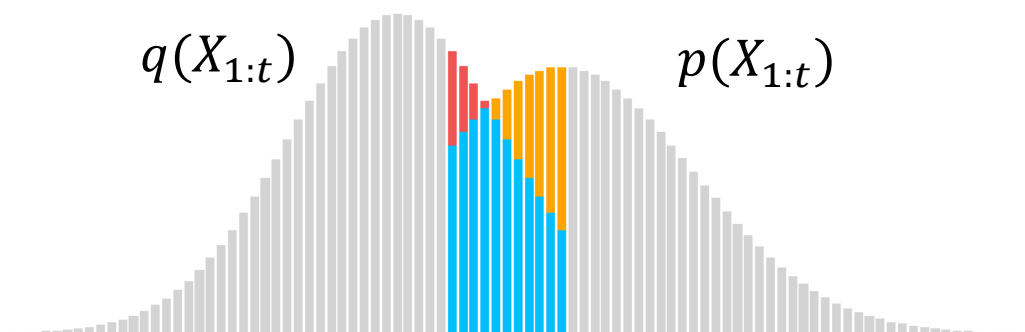
$$= D_{Branch}(q, p | X_{1:t-1}) \frac{p(X_{1:t}) - \min\{p(X_{1:t}), q(X_{1:t})\}}{D_{Branch}(p, q | X_{1:t-1})}$$

# Joint Verification

➤ **Case 2:**

$$\text{Issue 1: } \underbrace{D_{Branch}(q, p|X_{1:t-1})}_{\text{red area}} \neq \underbrace{D_{Branch}(p, q|X_{1:t-1})}_{\text{orange area}}$$

**Issue 2:**  $p(\cdot | \tilde{X}_{1:t-1} \neq X'_{1:t-1})$  and  $q(\cdot | \tilde{X}_{1:t-1} \neq X'_{1:t-1})$  are unknown



➡  $P(X_{1:t} \text{ is yielded}) \neq p(X_{1:t})$

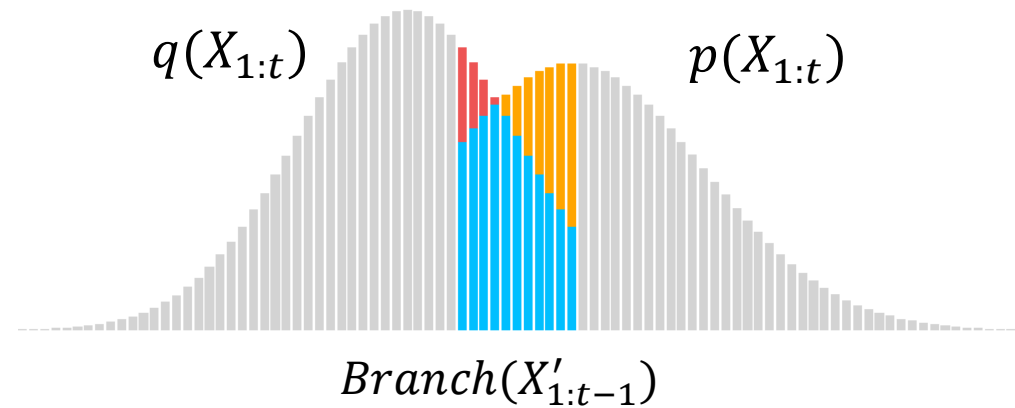
# Method

## Theory

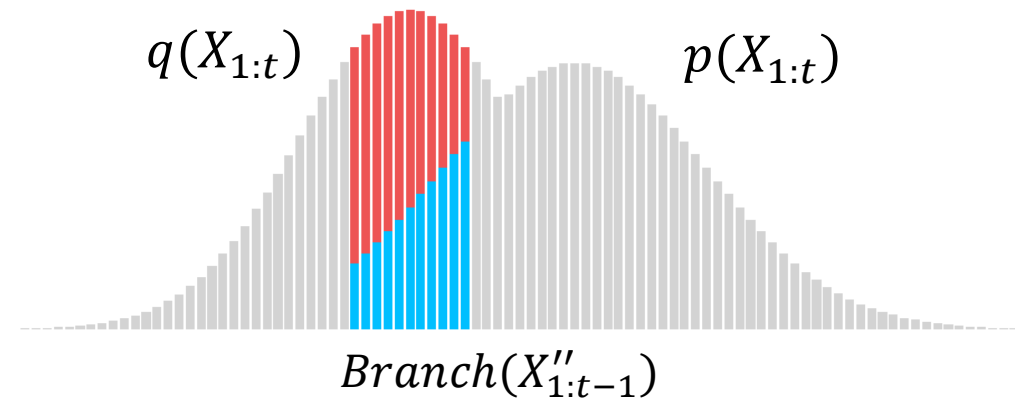
### ➤ Definition: Asymmetry of Branch Divergence

$$\Delta_{Branch}(X_{1:t-1}) = D_{Branch}(p, q|X_{1:t-1}) - D_{Branch}(q, p|X_{1:t-1})$$

- **negative asymmetry:** excess probability mass  $\Rightarrow$  target dist. locally recoverable



- **positive asymmetry:** deficient probability mass  $\Rightarrow$  target dist. locally unrecoverable



# Method

## Theory

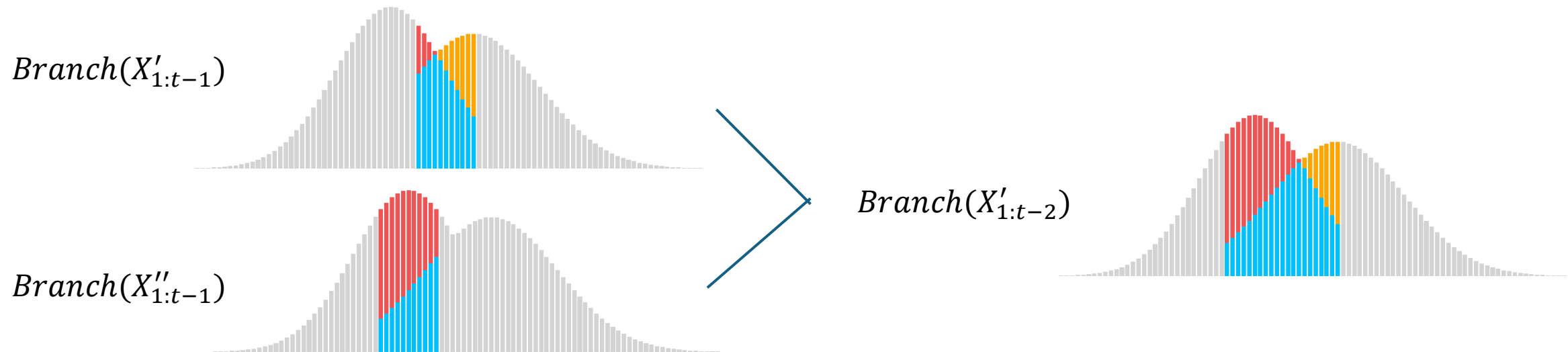
### ➤ Theorem: Hierarchy of Branch Divergence

$$\sum_{\Delta_{Branch}(X_{1:t-2}, \tilde{x}_{t-1}) > 0} \Delta_{Branch}(X_{1:t-2}, \tilde{x}_{t-1}) = D_{Branch}(p, q | X_{1:t-2})$$

### ➤ Key Insight:

Target distribution is recoverable by **recursively** compensating deficient mass with excess mass at each branch, because:

$$D_{Branch}(p, q | X_0) = 0$$



# Naive Hierarchical Speculative Decoding

## ➤ Method:

$$\sum_{\tau=0}^{\gamma} \sum_{\tilde{X}_{\tau+1:\gamma}} P(\tilde{X}_{1:\gamma} \neq X_{1:\gamma} \text{ is sampled and rejected, } X_{1:\gamma} \text{ is resampled}) =$$

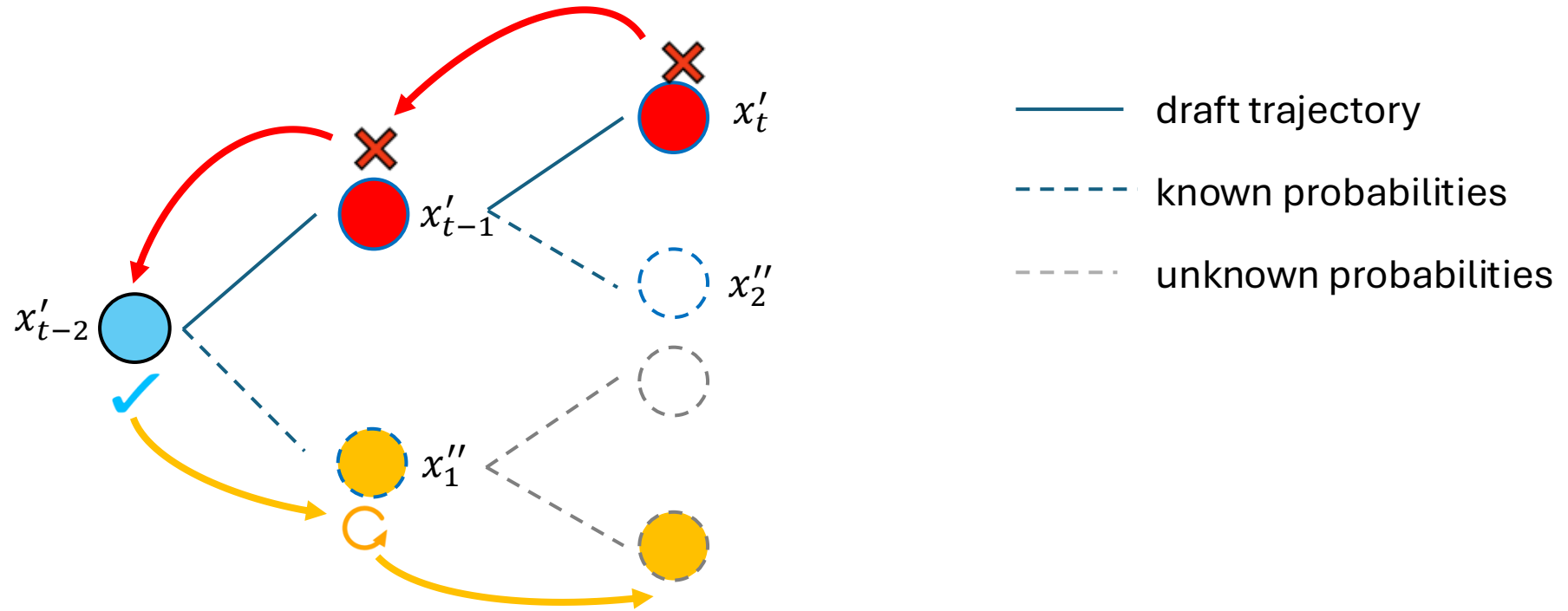
$$\sum_{\tau=0}^{\gamma} \sum_{\tilde{X}_{\tau+1:\gamma}} q(X_{1:\tau} \tilde{X}_{\tau+1:\gamma}) \cdot \prod_{t=\tau+1}^{\gamma} (1 - h_t) \cdot h_{\tau} \cdot \prod_{t=\tau+1}^{\gamma} P_{res}(x_t)$$

Explanation of probability terms:

- **Sampling**  $q(X_{1:\tau} \tilde{X}_{\tau+1:\gamma})$ : generate the initial draft sequence.
- **Backward Scan**  $\prod_{t=\tau+1}^{\gamma} (1 - h_t)$ : reject tokens until the first accepted prefix is found.
- **Acceptance**  $h_{\tau}$ : accept the longest prefix  $X_{1:\tau}$ .
- **Resampling**  $\prod_{t=\tau+1}^{\gamma} P_{res}(x_t)$ : recursively resample to recover the target probability.

# Naive Hierarchical Speculative Decoding

## ➤ Illustration:



## ➤ Issue:

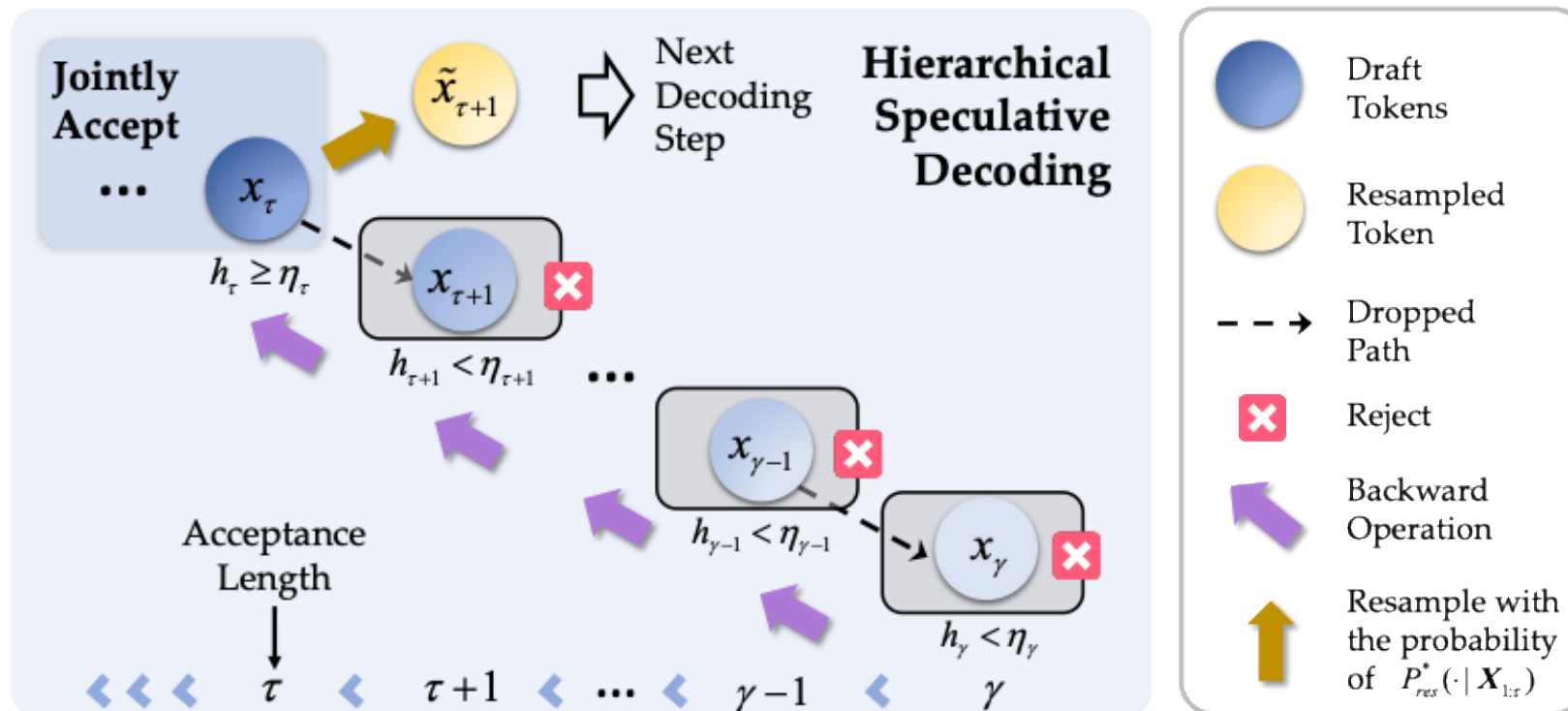
$P_{res}(x_t)$  requires the target and draft probability at the newly resampled trajectory,

➡  $\gamma - \tau + 1$  calls to the target model

# Hierarchical Speculative Decoding

## ➤ Method: Capped Branch Resampling

By clipping the maximum prefix ratio, we derive the corresponding accept and resample rule, which recovers the target distribution with a **single resampling step** within the accessible branches.



Thank you!