

Main Message

Q: How can we *adaptively* select hyperparameters for speculative decoding?

A: BanditSpec views the text generation process as an *online decision-making problem*. It adaptively selects hyperparameters using *Multi-Armed Bandit* algorithms.

Training-Free with Theoretical Guarantees !



Scan for paper

Problem Formulation

- Accelerate the inference of LLMs while maintaining high generation quality.

Draft \rightarrow Verify \rightarrow Accept

$$T_{\text{spec}} = T_{\text{draft}} \times L + T_{\text{target}} + T_{\text{accept}} \times n_{\text{accepted}} \approx T_{\text{target}}$$

- Total time saved

$$T_{\text{target}} \times \mathbb{E}[\tau_c - \tau_{\text{spec}}]$$

T_{target} : the time of one forward pass of the target model.

τ_c : the number of canonical decoding rounds.

τ_{spec} : the number of speculative decoding rounds.

The BANDITSPEC Framework

Algorithm 1 Speculative Decoding with Bandits (BanditSpec)

Inputs: arm selection algorithm **ALG**, initial prompt $\text{pt}_0 = \text{pt} \in \mathcal{X}^*$, bandit configuration $\nu = (P, \mathcal{S} = \{S_i\}_{i \in [K]}, L)$.

Procedures:

- 1: $t = 0, \mathcal{H}_0 = \emptyset, I_0 = 1, x_{I_0,0} = \emptyset$.
- 2: **while** $\text{EOS} \notin x_{I_t,t}$ **do**
- 3: $t = t + 1$.
- 4: Select a hyperparameter index $I_t = \text{ALG}(\mathcal{H}_{t-1})$.
- 5: $x_{I_t,t} = \text{SpecDecSub}(\text{pt}_{t-1}, P, S_{I_t}, L)$.
- 6: $\text{pt}_t = \text{concat}(\text{pt}_{t-1}, x_{I_t,t})$.
- 7: $\mathcal{H}_t = \text{concat}(\mathcal{H}_{t-1}, (I_t, x_{I_t,t}))$.
- 8: **end while**
- 9: **return** $\text{ST}(\text{ALG}, \text{pt}, \nu) = t, \text{pt}_{\text{ST}(\text{ALG}, \text{pt}, \nu)} = \text{pt}_t$.

- Objective:** Devise an arm selection rule **ALG** to minimize the *stopping time regret*

$$\text{Reg}(\text{ALG}, \text{pt}, \nu) := \mathbb{E}[\text{ST}(\text{ALG}, \text{pt}, \nu) \mid \text{pt}, \nu] - \mathbb{E}[\text{ST}(\text{ALG}_{i^*(\text{pt}, \nu)}, \text{pt}, \nu) \mid \text{pt}, \nu]$$

- Desired result:**

$$\text{Reg}(\text{ALG}, \text{pt}, \nu) = o(\mathbb{E}[\text{len}(\text{pt}_{\tau_c})]) \text{ or } o(\mathbb{E}[\tau_c]).$$

Main Results

Stationary Mean Values assumption: There exist K values $\{\mu_i\}_{i \in [K]} \subset [1, L+1]$, such that conditioned on the history \mathcal{H}_{t-1} and the chosen arm I_t at time t , the expected number of the accepted tokens $\mathbb{E}[Y_{I_t,t} \mid \mathcal{H}_{t-1}, I_t] = \mu_{I_t}$.

Adversarial Mean Values assumption: Let the number of accepted tokens generated by hyperparameter S_i at time step t be $y_{i,t} = \text{len}(X_{i,t})$. We assume $\{y_{i,t}\}_{i \in [K], t \in \mathbb{N}}$ is fixed by the environment before the algorithm starts.

Under the Stationary Mean Values assumption,

ALG = UCBSpec achieves a regret upper bound as

$$\text{Reg}(\text{ALG}, \text{pt}, \nu) = O\left(H(\text{pt}, \nu) \cdot L^2 \cdot \log \mathbb{E}[\text{len}(\text{pt}_{\tau_c})]\right).$$

Additionally, for any algorithm **ALG**, if the acceptance token length follows the truncated geometric distribution, then

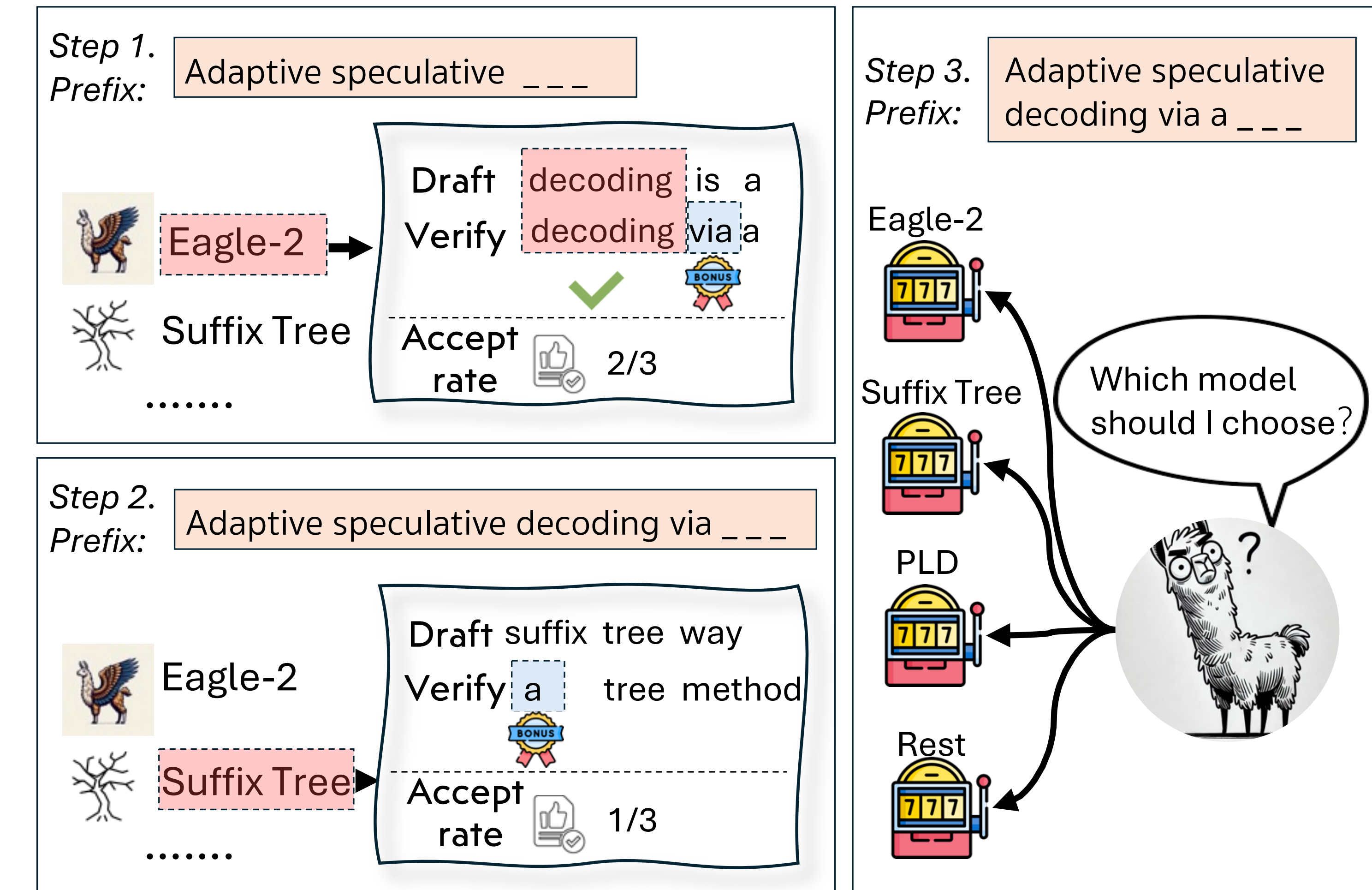
$$\liminf_{m \rightarrow \infty} \frac{\text{Reg}(\text{ALG}, \text{pt}^m, \nu)}{\log(\text{len}(\text{pt}_{\tau_c}^m))} \geq H(\text{pt}^m, \nu) \cdot \frac{p_{i^*}(1 - p_{i^*}^L)}{(1 - p_{i^*})}.$$

Under the Adversarial Mean Values assumption,

ALG = EXP3Spec achieves a regret upper bound as

$$\text{Reg}(\text{ALG}, \text{pt}, \nu) \leq 2L \cdot O\left(\sqrt{\min_{i \in [K]} \text{ST}(\text{ALG}_i) K \log K}\right).$$

Algorithm



Algorithm 2 UCBSpec

Inputs: number of hyperparameter specifications K , history $\mathcal{H}_t = ((I_s, X_{I_s,s}))_{s=1}^t$, confidence parameter δ .

Procedures:

- 1: **if** $t \leq K - 1$ **then return** $I_{t+1} = t + 1$.
- 2: Compute the lengths $Y_{I_s,s} = \text{len}(X_{I_s,s})$ for all $s \in [t]$.
- 3: Set the statistics $\{\hat{\mu}_{i,t}\}_{i \in [K]}, \{\text{UCB}_{i,t} = \hat{\mu}_{i,t} + \text{cr}_{i,t}\}_{i \in [K]}$, where

$$\text{cr}_{i,t} = \frac{L}{2} \sqrt{\frac{1 + n_{i,t}}{n_{i,t}^2} \left(1 + 2 \log \frac{K t^2 (1 + n_{i,t})^{\frac{1}{2}}}{\delta}\right)},$$

- 4: **return** index $I_{t+1} = \arg\max_{i \in [K]} \text{UCB}_{i,t}$.

Experiments

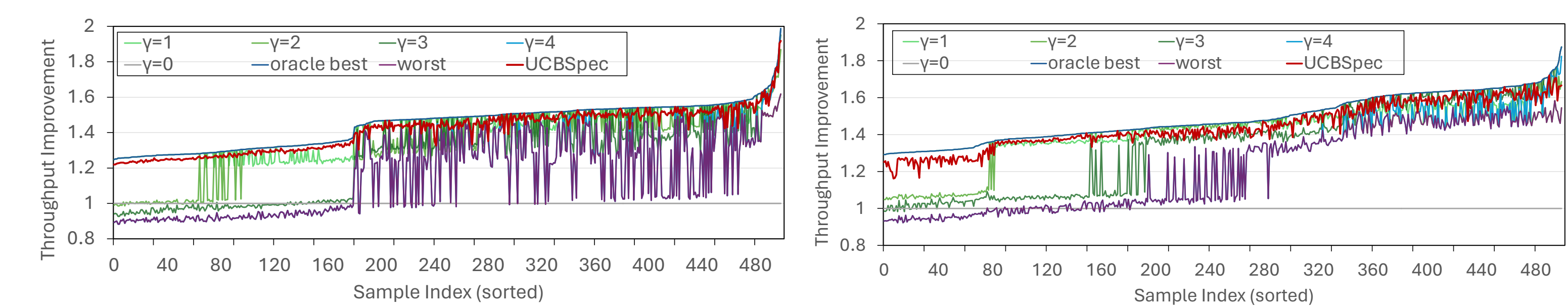


Figure 1. Throughput comparison of speculation lengths $\gamma \in [4]$ and the canonical decoding ($\gamma = 0$).

Table 1. Empirical Comparisons measured by Mean Accepted Tokens (MAT) (†) and Tokens/s (‡).

Methods	Spec Bench		Alpaca		Code Editor		Debug Bench	
	MAT(†)	Tokens/s(‡)	MAT(†)	Tokens/s(‡)	MAT(†)	Tokens/s(‡)	MAT(†)	Tokens/s(‡)
LLaMA3-8B-Instruct								
Vanilla	1.00	35.73	1.00	35.92	1.00	36.32	1.00	36.89
PLD	1.46	43.96	1.53	53.06	2.13	82.61	1.67	82.76
Rest	1.29	40.67	1.48	52.40	1.33	51.32	1.29	48.49
Suffix Tree	1.83	55.10	1.71	64.02	2.30	90.21	2.13	77.56
Eagle-2	3.94	98.15	4.04	110.00	4.79	128.76	4.78	119.12
EXP3Spec	3.65	<u>102.10</u>	4.23	<u>120.38</u>	4.36	<u>137.29</u>	4.50	<u>132.25</u>
UCBSpec	3.98	105.72	4.35	125.78	4.83	138.27	<u>4.60</u>	135.34