

BanditSpec: Adaptive Speculative Decoding via Bandit Algorithms (ICML 2025)

Yunlong Hou*, Fengzhuo Zhang*, Cunxiao Du*, Xuan Zhang*,
Jiachun Pan, Tianyu Pang, Chao Du, **Vincent Y. F. Tan**,
Zhuoran Yang



Yale

August 25, 2025
CNI Seminar, IISc

Background: Canonical Decoding

- **Canonical Decoding:** tokens are generated in an **autoregressive** manner: each token requires a full forward pass through the massive target model $P : \mathcal{X}^* \rightarrow \Delta_{\mathcal{X}}$.

Background: Canonical Decoding

- **Canonical Decoding:** tokens are generated in an **autoregressive** manner: each token requires a full forward pass through the massive target model $P : \mathcal{X}^* \rightarrow \Delta_{\mathcal{X}}$.

Algorithm 2 CANONICAL DECODING

Inputs: initial prompt $\text{pt}_0 = \text{pt} \in \mathcal{X}^*$, target model P .

Procedures:

- 1: Set $t = 0$.
 - 2: **while** $t \neq 0$ and $x_t \neq \text{EOS}$ **do**
 - 3: $t = t + 1$.
 - 4: $x_t \sim P(\cdot \mid \text{pt}_{t-1})$.
 - 5: $\text{pt}_t = \text{concat}(\text{pt}_{t-1}, x_t)$.
 - 6: **end while**
 - 7: **return** $\tau_c = t$, $\text{pt}_{\tau_c} = \text{pt}_t$
-

Background: Canonical Decoding

0 : [BOS]
 1 : [BOS] Adaptive
 2 : [BOS] Adaptive speculative
 3 : [BOS] Adaptive speculative decoding
 4 : [BOS] Adaptive speculative decoding via
 5 : [BOS] Adaptive speculative decoding via bandit
 6 : [BOS] Adaptive speculative decoding via bandit algorithm
 7 : [BOS] Adaptive speculative decoding via bandit algorithm is
 8 : [BOS] Adaptive speculative decoding via bandit algorithm is more
 9 : [BOS] Adaptive speculative decoding via bandit algorithm is more efficient
 10 : [BOS] Adaptive speculative decoding via bandit algorithm is more efficient .
 11 : [BOS] Adaptive speculative decoding via bandit algorithm is more efficient . [EOS]

• Total Time:

$$T_{\text{total}} = T_{\text{target}} \times \tau_c = 30 \text{ ms} \times 11 = 330 \text{ ms}.$$

How can we accelerate this one-by-one token generation process?

$$^0T_{\text{target}} \approx 30 \text{ ms for LLaMA3-8B.}$$

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft → **Verify** → **Accept**

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft \rightarrow **Verify** \rightarrow **Accept**

- **Step 1:** **Draft** the next L tokens via draft model Q *autoregressively*

$$Q(\cdot \mid \text{pt}_{t-1}, \tilde{x}_{t:t+j}), \quad j = 0, \dots, L-1.$$

- Draft models are small and fast (e.g., T5-small (77M)).
- Maybe inaccurate.

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft \rightarrow **Verify** \rightarrow **Accept**

- **Step 1:** **Draft** the next L tokens via draft model Q *autoregressively*

$$Q(\cdot \mid \text{pt}_{t-1}, \tilde{x}_{t:t+j}), \quad j = 0, \dots, L-1.$$

- Draft models are small and fast (e.g., T5-small (77M)).
- Maybe inaccurate.

[BOS] Adaptive speculative decoding **via**

[BOS] Adaptive speculative decoding **via suffix**

[BOS] Adaptive speculative decoding **via suffix tree**

[BOS] Adaptive speculative decoding **via suffix tree way**

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft \rightarrow **Verify** \rightarrow **Accept**

- **Step 2: Verify** the drafted tokens: compute the probabilities of the outputs via the target LLM *in parallel*

$$P(\tilde{x}_{t+j+1} \mid \text{pt}_{t-1}, \tilde{x}_{t:t+j}), \quad j = 0, \dots, L-1.$$

[BOS] Adaptive speculative decoding **via**

[BOS] Adaptive speculative decoding **via** **suffix**

[BOS] Adaptive speculative decoding **via** **suffix** **tree**

[BOS] Adaptive speculative decoding **via** **suffix** **tree** **way**

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft \rightarrow **Verify** \rightarrow **Accept**

- **Step 3: Accept** the drafted tokens **sequentially**.

- Accept \tilde{x}_{t+j+1} with probability

$$\min \left\{ 1, \frac{P(\tilde{x}_{t+j+1} \mid \text{pt}_{t-1}, \tilde{x}_{t:t+j})}{Q(\tilde{x}_{t+j+1} \mid \text{pt}_{t-1}, \tilde{x}_{t:t+j})} \right\}, j = 0, \dots, L-1.$$

[BOS] Adaptive speculative decoding **via**

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft \rightarrow **Verify** \rightarrow **Accept**

- **Step 3: Accept** the drafted tokens **sequentially**.

- Accept \tilde{x}_{t+j+1} with probability

$$\min \left\{ 1, \frac{P(\tilde{x}_{t+j+1} \mid \text{pt}_{t-1}, \tilde{x}_{t:t+j})}{Q(\tilde{x}_{t+j+1} \mid \text{pt}_{t-1}, \tilde{x}_{t:t+j})} \right\}, j = 0, \dots, L-1.$$

[BOS] Adaptive speculative decoding **via** ✓

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft \rightarrow **Verify** \rightarrow **Accept**

- **Step 3: Accept** the drafted tokens **sequentially**.
 - Accept the draft tokens until the first rejected one (e.g. \tilde{x}_i).

[BOS] Adaptive speculative decoding via ✓

[BOS] Adaptive speculative decoding via suffix

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft \rightarrow **Verify** \rightarrow **Accept**

- **Step 3: Accept** the drafted tokens **sequentially**.
 - Accept the draft tokens until the first rejected one (e.g. \tilde{x}_i).

[BOS] Adaptive speculative decoding **via** ✓

[BOS] Adaptive speculative decoding **via** **suffix** ✗

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft \rightarrow **Verify** \rightarrow **Accept**

- **Step 3: Accept** the drafted tokens **sequentially**.
 - Accept the draft tokens until the first rejected one (e.g. \tilde{x}_i).
 - **Correct** \tilde{x}_i by $x_i \sim \text{Norm}\left([P(\cdot \mid \text{pt}, \tilde{x}_{1:i-1}) - Q(\cdot \mid \text{pt}, \tilde{x}_{1:i-1})]_+\right)$.

[BOS] Adaptive speculative decoding **via** ✓

[BOS] Adaptive speculative decoding **via** **suffix** ✗

[BOS] Adaptive speculative decoding **via** **suffix** **bandit**

~~[BOS] Adaptive speculative decoding **via** **suffix** **tree**~~

~~[BOS] Adaptive speculative decoding **via** **suffix** **tree** **way**~~

Preliminary: Speculative Decoding

Speculative Decoding: Accelerate inference of LLMs while maintaining high generation quality (Chen et al., 2023; Leviathan et al., 2023).

Draft → **Verify** → **Accept**

Input: [BOS] Adaptive speculative decoding

[BOS] Adaptive speculative decoding **via suffix tree way**

[BOS] Adaptive speculative decoding **via**

[BOS] Adaptive speculative decoding **via suffix**

[BOS] Adaptive speculative decoding **via suffix tree**

[BOS] Adaptive speculative decoding **via suffix tree way**

[BOS] Adaptive speculative decoding **via ✓**

[BOS] Adaptive speculative decoding **via suffix ✗**

[BOS] Adaptive speculative decoding **via ~~suffix~~ bandit**

⇒ Output: [BOS] Adaptive speculative decoding via bandit

Preliminary: Speculative Decoding

Canonical Decoding:

0 : [BOS]
 1 : [BOS] Adaptive
 2 : [BOS] Adaptive speculative
 3 : [BOS] Adaptive speculative decoding
 4 : [BOS] Adaptive speculative decoding via
 5 : [BOS] Adaptive speculative decoding via bandit
 6 : [BOS] Adaptive speculative decoding via bandit algorithm
 7 : [BOS] Adaptive speculative decoding via bandit algorithm is
 8 : [BOS] Adaptive speculative decoding via bandit algorithm is more
 9 : [BOS] Adaptive speculative decoding via bandit algorithm is more efficient
 10 : [BOS] Adaptive speculative decoding via bandit algorithm is more efficient .
 11 : [BOS] Adaptive speculative decoding via bandit algorithm is more efficient . [EOS]

Speculative Decoding:

0 : [BOS]
 1 : [BOS] Adaptive speculative sampling is decoding
 2 : [BOS] Adaptive speculative decoding via suffix tree way bandit
 3 : [BOS] Adaptive speculative decoding via bandit algorithm is definitely more
 4 : [BOS] Adaptive speculative decoding via bandit algorithm is more efficient . [EOS]

Preliminary: Speculative Decoding

- For one round of speculative decoding:

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

Preliminary: Speculative Decoding

- For one round of speculative decoding:

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

- Total time saved:

$$T_{\text{target}} \times \mathbb{E} \left[\tau_{\text{c}} - \tau_{\text{spec}} \right]$$

τ_{c} : the number of canonical decoding rounds.

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

Preliminary: Speculative Decoding

- For one round of speculative decoding:

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

- Total time saved:

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

τ_c : the number of canonical decoding rounds.

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

- The distribution of the generated token sequence is unbiased (Leviathan et al., 2023; Chen et al., 2023; Yin et al., 2024):

$$\text{pt}_{\tau_{\text{spec}}} \stackrel{d}{=} \text{pt}_{\tau_c}.$$

Preliminary: Speculative Decoding

- Most existing speculative decoding methods use a **fixed** configuration.
 - A **single draft model** or a **fixed set of hyperparameters** is used for all tasks.
 - This is **suboptimal** because the ideal configuration depends heavily on the specific context.

Preliminary: Speculative Decoding

- Most existing speculative decoding methods use a **fixed** configuration.
 - A **single draft model** or a **fixed set of hyperparameters** is used for all tasks.
 - This is **suboptimal** because the ideal configuration depends heavily on the specific context.

Task	Best approach
Code Debugging	Requires precision; a retrieval-based method like Suffix Tree (Oliaro et al., 2024) might be best.
Story Generation	Requires creativity; a smaller draft LLM like Eagle (Li et al., 2024b) might be better.

Preliminary: Speculative Decoding

- Most existing speculative decoding methods use a **fixed** configuration.
 - A **single draft model** or a **fixed set of hyperparameters** is used for all tasks.
 - This is **suboptimal** because the ideal configuration depends heavily on the specific context.

Task	Best approach
Code Debugging	Requires precision; a retrieval-based method like Suffix Tree (Oliaro et al., 2024) might be best.
Story Generation	Requires creativity; a smaller draft LLM like Eagle (Li et al., 2024b) might be better.

Is there any *training-free* method that can *adaptively* choose the hyperparameters such that the latency of speculative decoding can be minimized?

Preliminary: Bandit Model

Algorithm 3 Dynamics of Multi-Armed Bandits

- 1: **Inputs:** K arms, time horizon T .
 - 2: $\mathcal{H}_0 = \emptyset$.
 - 3: **for** $t = 1, 2, \dots, T$ **do**
 - 4: Agent adopts an algorithm ALG to **select** arm I_t based on \mathcal{H}_{t-1} .
 - 5: Environment reveals the reward $X_{I_t,t}$ to the agent.
 - 6: Agent updates the history $\mathcal{H}_t = \text{concat}(\mathcal{H}_{t-1}, (I_t, X_{I_t,t}))$.
 - 7: **end for**
-

Goal: Devise an algorithm ALG to minimize the cumulative regret

$$\max_{i \in [K]} \mathbb{E} \left[\sum_{t=1}^T X_{i,t} \right] - \mathbb{E}_{\text{ALG}} \left[\sum_{t=1}^T X_{I_t,t} \right].$$

Preliminary: Bandit Model

Algorithm 4 Dynamics of Multi-Armed Bandits

- 1: **Inputs:** K arms, time horizon T .
 - 2: $\mathcal{H}_0 = \emptyset$.
 - 3: **for** $t = 1, 2, \dots, T$ **do**
 - 4: Agent adopts an algorithm ALG to select arm I_t based on \mathcal{H}_{t-1} .
 - 5: Environment **reveals** the reward $X_{I_t,t}$ to the agent.
 - 6: Agent updates the history $\mathcal{H}_t = \text{concat}(\mathcal{H}_{t-1}, (I_t, X_{I_t,t}))$.
 - 7: **end for**
-

Goal: Devise an algorithm ALG to minimize the cumulative regret

$$\max_{i \in [K]} \mathbb{E} \left[\sum_{t=1}^T X_{i,t} \right] - \mathbb{E}_{\text{ALG}} \left[\sum_{t=1}^T X_{I_t,t} \right].$$

Preliminary: Bandit Model

Algorithm 5 Dynamics of Multi-Armed Bandits

- 1: **Inputs:** K arms, time horizon T .
 - 2: $\mathcal{H}_0 = \emptyset$.
 - 3: **for** $t = 1, 2, \dots, T$ **do**
 - 4: Agent adopts an algorithm ALG to select arm I_t based on \mathcal{H}_{t-1} .
 - 5: Environment reveals the reward $X_{I_t,t}$ to the agent.
 - 6: Agent **updates** the history $\mathcal{H}_t = \text{concat}(\mathcal{H}_{t-1}, (I_t, X_{I_t,t}))$.
 - 7: **end for**
-

Goal: Devise an algorithm ALG to minimize the cumulative regret

$$\max_{i \in [K]} \mathbb{E} \left[\sum_{t=1}^T X_{i,t} \right] - \mathbb{E}_{\text{ALG}} \left[\sum_{t=1}^T X_{I_t,t} \right].$$

Preliminary: Bandit Model

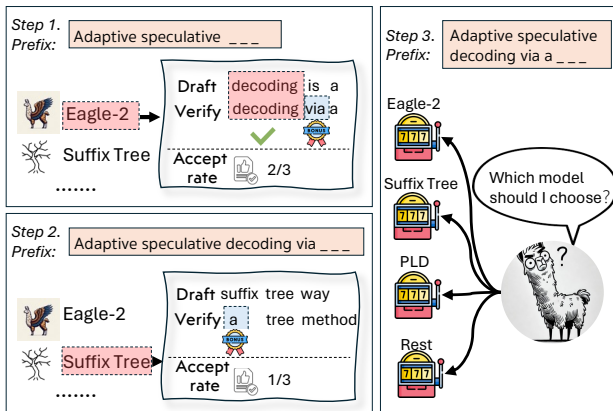
Algorithm 6 Dynamics of Multi-Armed Bandits

- 1: **Inputs:** K arms, time horizon T .
 - 2: $\mathcal{H}_0 = \emptyset$.
 - 3: **for** $t = 1, 2, \dots, T$ **do**
 - 4: Agent adopts an algorithm **ALG** to select arm I_t based on \mathcal{H}_{t-1} .
 - 5: Environment reveals the reward $X_{I_t,t}$ to the agent.
 - 6: Agent updates the history $\mathcal{H}_t = \text{concat}(\mathcal{H}_{t-1}, (I_t, X_{I_t,t}))$.
 - 7: **end for**
-

Goal: Devise an algorithm **ALG** to minimize the cumulative regret

$$\max_{i \in [K]} \mathbb{E} \left[\sum_{t=1}^T X_{i,t} \right] - \mathbb{E}_{\text{ALG}} \left[\sum_{t=1}^T X_{I_t,t} \right].$$

Problem Formulation



Hyperparameter Selection Problem: Each candidate hyperparameter configuration (e.g., draft model) is treated as an “arm”.

Bandit Framework for Speculative Decoding: BANDITSPEC

BANDITSPEC: A *training-free and adaptive* online learning framework for hyperparameter selection in speculative decoding.

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

$\text{ST}(\text{ALG}, \text{pt}, \nu)$: Number of calls to SPECDECSUB .

Bandit Framework for Speculative Decoding: BANDITSPEC

BANDITSPEC: A *training-free and adaptive* online learning framework for hyperparameter selection in speculative decoding.

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

$\text{ST}(\text{ALG}, \text{pt}, \nu)$: Number of calls to SPECDECSUB .

Bandit Framework for Speculative Decoding: BANDITSPEC

BANDITSPEC: A *training-free and adaptive* online learning framework for hyperparameter selection in speculative decoding.

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

$\text{ST}(\text{ALG}, \text{pt}, \nu)$: Number of calls to SPECDECSUB .

Bandit Framework for Speculative Decoding: BANDITSPEC

BANDITSPEC: A *training-free and adaptive* online learning framework for hyperparameter selection in speculative decoding.

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

$\text{ST}(\text{ALG}, \text{pt}, \nu)$: Number of calls to SPECDECSUB .

Bandit Framework for Speculative Decoding: BANDITSPEC

BANDITSPEC: A *training-free and adaptive* online learning framework for hyperparameter selection in speculative decoding.

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

$\text{ST}(\text{ALG}, \text{pt}, \nu)$: Number of calls to SPECDECSUB .

Bandit Framework for Speculative Decoding: BANDITSPEC

Proposition (Property of BANDITSPEC)

For any arm selection algorithm ALG that selects an arm according to the history,

$$\text{pt}_{\text{ST}(\text{ALG})} \stackrel{d}{=} \text{pt}_{\tau_c}.$$

Furthermore, the stopping time $\text{ST}(\text{ALG})$ can be bounded as

$$\frac{\text{len}(\text{pt}_{\text{ST}(\text{ALG})})}{L+1} \leq \text{ST}(\text{ALG}) \leq \text{len}(\text{pt}_{\text{ST}(\text{ALG})}), \text{ a.s.}$$

Bandit Framework for Speculative Decoding: BANDITSPEC

Proposition (Property of BANDITSPEC)

For any arm selection algorithm ALG that selects an arm according to the history,

$$\text{pt}_{\text{ST}(\text{ALG})} \stackrel{d}{=} \text{pt}_{\tau_c}.$$

Furthermore, the stopping time $\text{ST}(\text{ALG})$ can be bounded as

$$\frac{\text{len}(\text{pt}_{\text{ST}(\text{ALG})})}{L+1} \leq \text{ST}(\text{ALG}) \leq \text{len}(\text{pt}_{\text{ST}(\text{ALG})}), \text{ a.s.}$$

- Quality of the token sequence is not compromised.

Bandit Framework for Speculative Decoding: BANDITSPEC

Proposition (Property of BANDITSPEC)

For any arm selection algorithm ALG that selects an arm according to the history,

$$\text{pt}_{\text{ST}(\text{ALG})} \stackrel{d}{=} \text{pt}_{\tau_c}.$$

Furthermore, the stopping time $\text{ST}(\text{ALG})$ can be bounded as

$$\frac{\text{len}(\text{pt}_{\text{ST}(\text{ALG})})}{L+1} \leq \text{ST}(\text{ALG}) \leq \text{len}(\text{pt}_{\text{ST}(\text{ALG})}), \text{ a.s.}$$

- Quality of the token sequence is not compromised.
- The stopping time and the length of the generated tokens sequence are equivalent up to a constant factor.

Bandit Framework for Speculative Decoding: BANDITSPEC

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

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

Bandit Framework for Speculative Decoding: BANDITSPEC

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

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

- $\mathbb{E}[\text{ST}(\text{ALG}, \text{pt}, \nu) \mid \text{pt}, \nu]$: the expected number of speculative decoding rounds for **ALG** to stop.

Bandit Framework for Speculative Decoding: BANDITSPEC

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

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

- $\mathbb{E}[\text{ST}(\text{ALG}, \text{pt}, \nu) \mid \text{pt}, \nu]$: the expected number of speculative decoding rounds for **ALG** to stop.
- $\mathbb{E}[\text{ST}(\text{ALG}_{i^*(\text{pt}, \nu)}, \text{pt}, \nu) \mid \text{pt}, \nu]$: the expected number of rounds to stop if the oracle best hyperparameter $S_{i^*(\text{pt}, \nu)}$ is adopted.

Bandit Framework for Speculative Decoding: BANDITSPEC

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

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

- $\mathbb{E}[\text{ST}(\text{ALG}, \text{pt}, \nu) \mid \text{pt}, \nu]$: the expected number of speculative decoding rounds for **ALG** to stop.
- $\mathbb{E}[\text{ST}(\text{ALG}_{i^*(\text{pt}, \nu)}, \text{pt}, \nu) \mid \text{pt}, \nu]$: the expected number of rounds to stop if the oracle best hyperparameter $S_{i^*(\text{pt}, \nu)}$ is adopted.
- $\text{Reg}(\text{ALG}, \text{pt}, \nu)$ measures the wasted/additional rounds.

Bandit Framework for Speculative Decoding: BANDITSPEC

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

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

- $\mathbb{E}[\text{ST}(\text{ALG}, \text{pt}, \nu) \mid \text{pt}, \nu]$: the expected number of speculative decoding rounds for **ALG** to stop.
- $\mathbb{E}[\text{ST}(\text{ALG}_{i^*(\text{pt}, \nu)}, \text{pt}, \nu) \mid \text{pt}, \nu]$: the expected number of rounds to stop if the oracle best hyperparameter $S_{i^*(\text{pt}, \nu)}$ is adopted.
- $\text{Reg}(\text{ALG}, \text{pt}, \nu)$ measures the wasted/additional rounds.

- **Desired result:**

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

- The additional rounds are negligible for **ALG**, compared to the oracle best hyperparameter $S_{i^*(\text{pt}, \nu)}$.

Challenges

- There is another source of uncertainty, the **stopping time**, in addition to the arm rewards in standard multi-armed bandits

Challenges

- There is another source of uncertainty, the **stopping time**, in addition to the arm rewards in standard multi-armed bandits
 - For instance, even if we know there will be 1 accepted token at each round, the stopping time is

$$\text{ST}(\text{ALG}, \text{pt}, \nu) = \text{len}(\text{pt}_{\tau_c})$$

which depends on the target model P .

Challenges

- There is another source of uncertainty, the **stopping time**, in addition to the arm rewards in standard multi-armed bandits
 - For instance, even if we know there will be 1 accepted token at each round, the stopping time is

$$\text{ST}(\text{ALG}, \text{pt}, \nu) = \text{len}(\text{pt}_{\tau_c})$$

which depends on the target model P .

- The stopping time **depends on the generated tokens**. Mathematically,

$$\psi_t : \mathcal{X}^* \rightarrow \{\text{stop}, \text{continue}\}, \quad \psi_t(\text{pt}_t) = \text{stop or continue}$$

is a function that maps the generated tokens to the stopping decision, which the agent does not know in advance.

Challenges

- There is another source of uncertainty, the **stopping time**, in addition to the arm rewards in standard multi-armed bandits
 - For instance, even if we know there will be 1 accepted token at each round, the stopping time is

$$ST(\text{ALG}, \text{pt}, \nu) = \text{len}(\text{pt}_{\tau_c})$$

which depends on the target model P .

- The stopping time **depends on the generated tokens**. Mathematically,

$$\psi_t : \mathcal{X}^* \rightarrow \{\text{stop}, \text{continue}\}, \quad \psi_t(\text{pt}_t) = \text{stop or continue}$$

is a function that maps the generated tokens to the stopping decision, which the agent does not know in advance.

- The tokens are **generated autoregressively**, thus, the accepted tokens are **not i.i.d.**

Challenges

- There is another source of uncertainty, the **stopping time**, in addition to the arm rewards in standard multi-armed bandits
 - For instance, even if we know there will be 1 accepted token at each round, the stopping time is

$$ST(\text{ALG}, \text{pt}, \nu) = \text{len}(\text{pt}_{\tau_c})$$

which depends on the target model P .

- The stopping time **depends on the generated tokens**. Mathematically,

$$\psi_t : \mathcal{X}^* \rightarrow \{\text{stop}, \text{continue}\}, \quad \psi_t(\text{pt}_t) = \text{stop or continue}$$

is a function that maps the generated tokens to the stopping decision, which the agent does not know in advance.

- The tokens are **generated autoregressively**, thus, the accepted tokens are **not i.i.d.**

Some simplifications are required to make the problem tractable.

Challenges

- There is another source of uncertainty, the **stopping time**, in addition to the arm rewards in standard multi-armed bandits
 - For instance, even if we know there will be 1 accepted token at each round, the stopping time is

$$ST(\text{ALG}, \text{pt}, \nu) = \text{len}(\text{pt}_{\tau_c})$$

which depends on the target model P . \Rightarrow **Anytime-version algorithms**

- The stopping time **depends on the generated tokens**. Mathematically,

$$\psi_t : \mathcal{X}^* \rightarrow \{\text{stop}, \text{continue}\}, \quad \psi_t(\text{pt}_t) = \text{stop or continue}$$

is a function that maps the generated tokens to the stopping decision, which the agent does not know in advance. \Rightarrow **Martingale arguments**

- The tokens are **generated autoregressively**, thus, the accepted tokens are **not i.i.d.** \Rightarrow **Stochastic/Adversarial models**

Some simplifications are required to make the problem tractable.

Arm Selection: Modeling Tokens Stochastically

Assumption (Stationary Mean Values)

There exist K values $\{\mu_i\}_{i \in [K]} \subset [1, L+1]$, such that the expected number of the accepted tokens satisfies

$$\mathbb{E}[Y_{I_t, t} \mid \mathcal{H}_{t-1}, I_t] = \mu_{I_t}.$$

Arm Selection: Modeling Tokens Stochastically

Assumption (Stationary Mean Values)

There exist K values $\{\mu_i\}_{i \in [K]} \subset [1, L+1]$, such that the expected number of the accepted tokens satisfies

$$\mathbb{E}[Y_{I_t, t} | \mathcal{H}_{t-1}, I_t] = \mu_{I_t}.$$

- UCBSPEC: A UCB-type algorithm, whose confidence radius is derived from the self-normalized bound, which holds for all $t \in \mathbb{N}$ (Abbasi-yadkori et al., 2011).
- We design the algorithm as an anytime version, meaning it does not require the time horizon as an input.

Arm Selection: Modeling Tokens Stochastically

Algorithm 7 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

$$n_{i,t} = \sum_{s=1}^t \mathbb{1}\{I_s = i\}, \quad \hat{\mu}_{i,t} = \frac{\sum_{s=1}^t Y_{i,s} \mathbb{1}\{I_s = i\}}{n_{i,t}},$$

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

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

Arm Selection: Modeling Tokens Stochastically

Theorem (Upper Bound for Stationary Case)

Under Assumptions 2 and finite length assumption, given any prompt $\text{pt} \in \mathcal{X}^$ and bandit configuration $\nu = (P, \mathcal{S} = \{S_i\}_{i \in [K]}, L)$, the expected stopping time regret of $\text{ALG} = \text{UCBSPEC}$ is upper bounded as*

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

Here, $\Delta_i := \mu_{i^*} - \mu_i$ and the *hardness parameter*

$$H(\text{pt}, \nu) := \sum_{i \neq i^*} \frac{1}{\mu_{i^*} \Delta_i}.$$

Arm Selection: Modeling Tokens Stochastically

Theorem (Upper Bound for Stationary Case)

Under Assumptions 2 and finite length assumption, given any prompt $\text{pt} \in \mathcal{X}^$ and bandit configuration $\nu = (P, \mathcal{S} = \{S_i\}_{i \in [K]}, L)$, the expected stopping time regret of $\text{ALG} = \text{UCBSPEC}$ is upper bounded as*

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

Here, $\Delta_i := \mu_{i^*} - \mu_i$ and the *hardness parameter*

$$\text{H}(\text{pt}, \nu) := \sum_{i \neq i^*} \frac{1}{\mu_{i^*} \Delta_i}.$$

Theoretical Results: Lower Bound

Truncated geometric distribution (TGD) on $[1, L + 1]$:

$$P_S(x) = \begin{cases} p^{x-1}(1-p), & x = 1, 2, \dots, L, \\ p^L, & x = L + 1. \end{cases} \quad (1)$$

which was considered in the seminal work on speculative decoding (Leviathan et al., 2023).

Proposition (Tightness Result (Informal))

Let $\mathcal{S}_{\text{TGD}} = \{S : P_S \text{ satisfies (1)}\}$. Let $\{S_i\}_{i=1}^K \subset \mathcal{S}_{\text{TGD}}$ and S_i satisfies (1) with p_i , then under the *greedy decoding strategy*,

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

If $p_{i^*} \in (2^{-1/L}, 1)$, the bounds match up absolute constants and L .

Arm Selection: Modeling Tokens Adversarially

Assumption (Adversarial Mean Values)

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.

- This is analogous to the oblivious adversarial case in bandits.
- EXP3SPEC: An EXP3-type algorithm (anytime version).

Arm Selection: Modeling Tokens Adversarially

Assumption (Adversarial Mean Values)

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.

- This is analogous to the oblivious adversarial case in bandits.
- EXP3SPEC: An EXP3-type algorithm (anytime version).

Theorem (Upper Bound for Adversarial Case, Informal)

Under some assumptions and using the **greedy decoding strategy**, given any prompt $\text{pt} \in \mathcal{X}^*$ and bandit configuration $\nu = (P, \mathcal{S} = \{S_i\}_{i \in [K]}, L)$, the expected stopping time regret of $\text{ALG} = \text{EXP3SPEC}$ is

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

Arm Selection: Modeling Tokens Adversarially

Algorithm 8 EXP3SPEC

Inputs: Num. of hyperparams. K , history $\mathcal{H}_t = ((I_s, X_{I_s,s}))_{s=1}^t$.

Procedures:

- 1: Compute the lengths $Y_{I_s,s} = \text{len}(X_{I_s,s})$ for all $s \in [t]$.
- 2: Set the statistics $\hat{Z}_{i,t} = \mathbb{1}\{i = I_t\} \cdot \frac{L+1-Y_{i,t}}{L \cdot p_{t,i}}$ for all $i \in [K]$.
- 3: Set learning rate $\eta_t = \sqrt{\log K / (t \cdot K)}$.
- 4: Set probability vector $p_t \in \Delta_{[K]}$ with for all $i \in [K]$

$$p_{t,i} = \frac{\exp\left(-\eta_t \sum_{s=1}^{t-1} \hat{Z}_{i,s}\right)}{\sum_{j=1}^K \exp\left(-\eta_t \sum_{s=1}^{t-1} \hat{Z}_{j,s}\right)}.$$

- 5: **return** hyperparameter index $I_{t+1} \sim p_t$.
-

Experimental Setup

- **Two Target Models:** LLaMA3-8B-Instruct and Qwen2-7B-Instruct.
- **Four Benchmarks:** Spec Bench (Xia et al., 2024), Alpaca (Taori et al., 2023), Code Editor (Guo et al., 2024) and Debug Bench (Tian et al., 2024).
- **Two Metrics:**
 - *Mean Accepted Tokens (MAT)*: important measurement, mainly focus on the design of the speculative decoding algorithm.
 - *Throughput (Tokens/s)*: closely related to user experience, and it takes all relevant factors into account, e.g., the hardware, the coding, the algorithm, etc.

Experimental Setup

- **Two Target Models:** LLaMA3-8B-Instruct and Qwen2-7B-Instruct.
- **Four Benchmarks:** Spec Bench (Xia et al., 2024), Alpaca (Taori et al., 2023), Code Editor (Guo et al., 2024) and Debug Bench (Tian et al., 2024).
- **Two Metrics:**
 - *Mean Accepted Tokens (MAT)*: important measurement, mainly focus on the design of the speculative decoding algorithm.
 - *Throughput (Tokens/s)*: closely related to user experience, and it takes all relevant factors into account, e.g., the hardware, the coding, the algorithm, etc.
- **Two Settings:**
 - *Batchsize = 1 and four drafting methods:* PLD (Saxena, 2023), Rest (He et al., 2024), Suffix Tree (Oliaro et al., 2024; Hu et al., 2024) and Eagle-2 (Li et al., 2024a).
 - *Simulated real-world scenario:* heterogeneous requests at one time, with different hyperparameters (i.e., speculation length $\gamma \in [4]$) as arms.

Experimental Results: Setting One

Methods	Spec Bench		Alpaca		Code Editor		Debug Bench	
	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)
<i>LLaMA3-8B-Instruct</i>								
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>	<u>4.23</u>	<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
<i>Qwen2-7B-Instruct</i>								
Vanilla	1.00	38.71	1.00	39.32	1.00	39.30	1.00	39.57
PLD	1.55	52.44	1.42	58.41	1.89	64.56	2.15	70.49
Rest	1.31	46.42	1.47	59.01	1.31	53.79	1.22	50.51
Suffix Tree	1.96	68.42	1.46	62.60	2.18	85.75	2.49	101.47
Eagle-2	3.64	97.82	3.61	104.43	4.88	138.58	4.79	126.01
EXP3SPEC	<u>3.76</u>	<u>107.36</u>	<u>3.83</u>	<u>113.90</u>	<u>4.90</u>	<u>160.41</u>	<u>4.86</u>	151.73
UCBSPEC	4.13	112.33	3.93	114.20	4.92	161.35	5.10	<u>151.37</u>

Experimental Results: Setting One

Methods	Spec Bench		Alpaca		Code Editor		Debug Bench	
	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)
<i>LLaMA3-8B-Instruct</i>								
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	<u>3.94</u>	98.15	4.04	110.00	<u>4.79</u>	128.76	4.78	119.12
EXP3SPEC	3.65	<u>102.10</u>	<u>4.23</u>	<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

- The best throuput performance is always achieved by the proposed algorithms.

Experimental Results: Setting One

Methods	Spec Bench		Alpaca		Code Editor		Debug Bench	
	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)
<i>LLaMA3-8B-Instruct</i>								
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	<u>3.94</u>	98.15	4.04	110.00	<u>4.79</u>	128.76	4.78	119.12
EXP3SPEC	3.65	<u>102.10</u>	<u>4.23</u>	<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

- The best throuput performance is always achieved by the proposed algorithms.
- The proposed methods automatically adapt to different prompts, as they are better than the best fixed method.

Experimental Results: Setting One

Methods	Spec Bench		Alpaca		Code Editor		Debug Bench	
	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)	MAT(↑)	Tokens/s(↑)
<i>LLaMA3-8B-Instruct</i>								
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	<u>3.94</u>	98.15	4.04	110.00	<u>4.79</u>	128.76	4.78	119.12
EXP3SPEC	3.65	<u>102.10</u>	<u>4.23</u>	<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

- The best throuput performance is always achieved by the proposed algorithms.
- The proposed methods automatically adapt to different prompts, as they are better than the best fixed method.
- As the empirical performance of UCBSPEC is better than EXP3SPEC, it implies that real-life scenario tends to be benign and may be more aligned with the stationary mean assumption.

Experimental Results: Setting Two

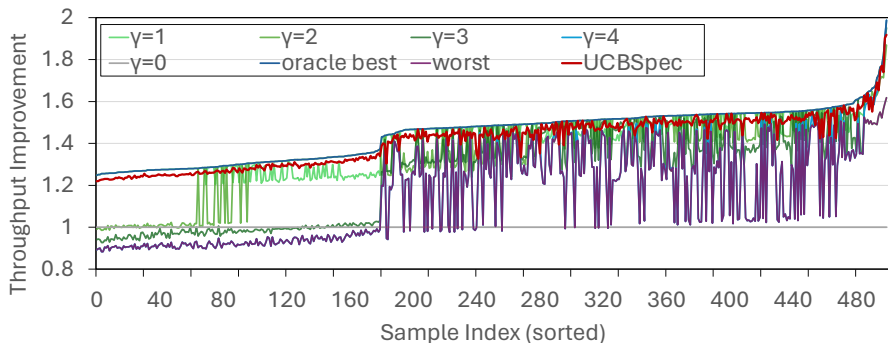


Figure: Target model: LLaMA3-8B-Instruct, Draft Model: Eagle-1

Experimental Results: Setting Two

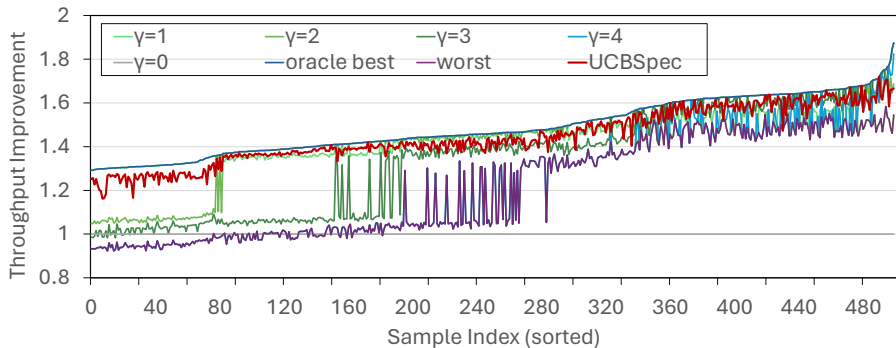


Figure: Target model: Qwen2-7B-Instruct, Draft Model: Eagle-1

Discussions

- The stationary mean assumption (Assumption 2)

$$\mathbb{E}[Y_{I_t,t} | \mathcal{H}_{t-1}, I_t] = \mu_{I_t}.$$

is strictly weaker than the usual i.i.d. assumption in bandits.

- Many variants of UCB (e.g., KL-UCB) algorithms that one can try.
- However, the implementation time of these arm selection algorithm should also be considered, as we aim at reduce the wall-time of the decoding process.

Conclusions and Future Directions

Summary:

- Introduced BANDITSPEC, a bandit-based *training-free* framework for *adaptive* speculative decoding.
- Developed two arm algorithms (UCBSPEC and EXP3SPEC) with theoretical regret guarantees.
- Demonstrated significant empirical improvements in decoding throughput.

Future work:

- Robust bandits and Non-stationary bandits.
- Contextual Bandits.

Reference I

- Abbasi-yadkori, Y., Pál, D., and Szepesvári, C. (2011). Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, volume 24, pages 2312–2320.
- Chen, C., Borgeaud, S., Irving, G., Lespiau, J.-B., Sifre, L., and Jumper, J. (2023). Accelerating large language model decoding with speculative sampling. *arXiv preprint arXiv:2302.01318*.
- Guo, J., Li, Z., Liu, X., Ma, K., Zheng, T., Yu, Z., Pan, D., Li, Y., Liu, R., Wang, Y., Guo, S., Qu, X., Yue, X., Zhang, G., Chen, W., and Fu, J. (2024). Codeeditorbench: Evaluating code editing capability of large language models.
- He, Z., Zhong, Z., Cai, T., Lee, J., and He, D. (2024). REST: Retrieval-based speculative decoding. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1582–1595. Association for Computational Linguistics.
- Hu, Y., Wang, K., Zhang, J., Li, C., and Chen, H. (2024). Sam decoding: Speculative decoding via suffix automaton. *arXiv preprint arXiv:2411.10666*.
- Leviathan, Y., Kalman, M., and Matias, Y. (2023). Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR.
- Li, Y., Wei, F., Zhang, C., and Zhang, H. (2024a). EAGLE-2: Faster inference of language models with dynamic draft trees. In *Empirical Methods in Natural Language Processing*.
- Li, Y., Wei, F., Zhang, C., and Zhang, H. (2024b). Eagle: Speculative sampling requires rethinking feature uncertainty. *arXiv preprint arXiv:2401.15077*.
- Oliaro, G., Jia, Z., Campos, D., and Qiao, A. (2024). Suffixdecoding: A model-free approach to speeding up large language model inference. *arXiv preprint arXiv:2411.04975*.
- Saxena, A. (2023). Prompt lookup decoding.
- Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto, T. B. (2023). Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- Tian, R., Ye, Y., Qin, Y., Cong, X., Lin, Y., Liu, Z., and Sun, M. (2024). Debugbench: Evaluating debugging capability of large language models.

Reference II

- Xia, H., Yang, Z., Dong, Q., Wang, P., Li, Y., Ge, T., Liu, T., Li, W., and Sui, Z. (2024). Unlocking efficiency in large language model inference: A comprehensive survey of speculative decoding. In Ku, L.-W., Martins, A., and Srikumar, V., editors, *Findings of the Association for Computational Linguistics ACL 2024*, pages 7655–7671, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Yin, M., Chen, M., Huang, K., and Wang, M. (2024). A theoretical perspective for speculative decoding algorithm. *arXiv preprint arXiv:2411.00841*.

Thanks for Listening!

