# SuffixDecoding: Extreme Speculative Decoding for Emerging AI Applications

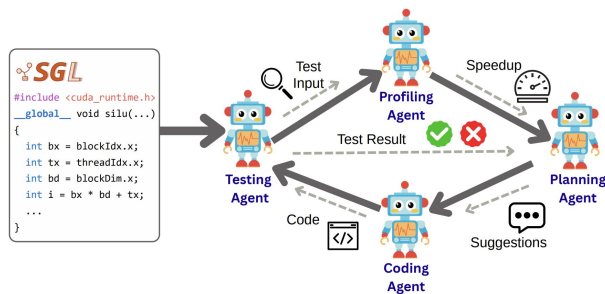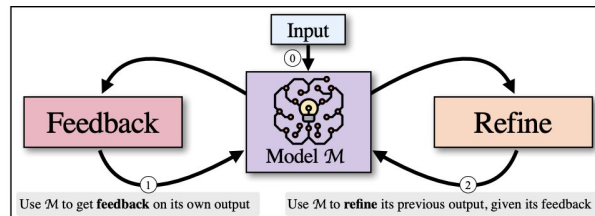Gabriele Oliaro, Zhihao Jia, Daniel Campos, Aurick Qiao

# Agenda

1. **The Problem:** Latency in Modern AI Workloads
2. **Background:** Speculative Decoding 101
3. **Our Solution:** SuffixDecoding
4. **Key Features:** Adaptive, Fast, and Hybrid
5. **Evaluation & Results**
   - End-to-End Speedups
   - Live vLLM Integration
6. **Deeper Dives & Ablations**
7. **Conclusion**

# Agenda

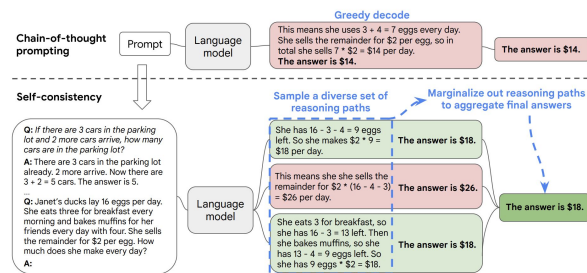# Challenge: Emerging AI Workloads Suffer from High Latency

✅ Inference-Time Scaling techniques can improve the output quality



Multi-Agent Pipelines

Iterative Refinement

Self-Consistency

⏳ But better accuracy comes **at the cost** of more generated tokens

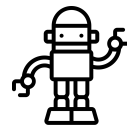# Opportunity: Token Repetitions in Emerging AI Workloads

★ Self-reflection loops

   ✓ Repeated code/text with minor edits

★ Multiple reasoning paths

   ✓ Similar reasoning patterns

★ Multi-agent workflows

   ✓ Shared context across agents

★ Iterative refinement

   ✓ 100+ lines repeated, 2 modified



Change two lines

💡 **To leverage these repetitions, we designed SuffixDecoding!**

# Agenda

# Speculative Decoding 101: Vanilla Decoding

# Speculative Decoding 101: Sequence-Based Speculation

prompt

repeat

Draft model

T1 → T2 → Tx → Ty

LLM

✅ ✅ ❌ 🗑️
T1 T2 ~~Tx~~ T3 Ty'

# Speculative Decoding 101: Tree-Based Speculation

# Speculative Decoding 101: Model-Free Speculation

# Agenda

🚀 SuffixDecoding: Up to **5.3x** end-to-end speedup



Previous Outputs

Index

Suffix Tree Cache

Aggressive Speculation

Fast LLM Agents

# SuffixDecoding: A Tale of Two Trees 🌲 🌳

Request Tree



Ongoing
Inference

| 8 | 5 | 4 | **8** | **5** |

Previous
Outputs

Index

Global Tree

# SuffixDecoding: A Tale of Two Trees 🌲 🌳



Request Tree

Ongoing
Inference

Previous
Outputs

Index

Global Tree

# SuffixDecoding: A Tale of Two Trees 🌲 🌳

# SuffixDecoding: A Tale of Two Trees 🌲 🌳

# SuffixDecoding: A Tale of Two Trees 🌲 🌳

# Agenda

# Efficient Verification via Adaptive Speculation

- Static speculation not viable because <u>verification cost</u> is substantial
- **Adaptive speculation** based on pattern match quality
- Formula: **MAX_SPEC = αp** (where p = pattern match length)

# Blazingly Fast and Memory Efficient Speculation



(a) Memory Scalability

(b) Update and Lookup Performance

# Robustness: SuffixDecoding + EAGLE-3 Hybrid



Mean Accepted Tokens per Step

# Robustness: SuffixDecoding + EAGLE-3 Hybrid

# Robustness: SuffixDecoding + EAGLE-3 Hybrid



Mean Accepted Tokens per Step

# Agenda

# End–To-End Evaluation Benchmarks

AgenticSQL

OpenHands CodeAct + SWE-Bench Verified

SpecBench



Agentic

Open-ended

# Up to 5.3x end-to-end speedup



Speculative Speedups over Vanilla Decoding

Mean Accepted Tokens per Step

# Up to 5.3x end-to-end speedup



Speculative Speedups over Vanilla Decoding

Mean Accepted Tokens per Step

# Up to 5.3x end-to-end speedup



Speculative Speedups over Vanilla Decoding

Mean Accepted Tokens per Step

# Up to 5.3x end-to-end speedup



Speculative Speedups over Vanilla Decoding

Mean Accepted Tokens per Step

# Up to 5.3x end-to-end speedup



Speculative Speedups over Vanilla Decoding

Mean Accepted Tokens per Step

# Up to 4.5x SWE-Bench Verified Task Completion Speedup



Solving SWE-Bench Verified: 1.8x - 4.5x Faster with Suffix Decoding

# Agenda

# Why Two Trees?

We ran an ablation to see what each suffix tree contributes.

- **Per-request Tree Only**: Decent speedup, but low.

- **Global Tree Only**: The long-term history provides most of the speedup.

- **Both Trees**: In almost every case, using both trees is the best.

# How do we know if SuffixDecoding will work?

- Measure the empirical **entropy** (i.e. the "structuredness") of a workload with just 100 example outputs.

- Lower average entropy indicates more predictable outputs and better SuffixDecoding performance.

  - **AgenticSQL Enrich:** 0.171 entropy → 10.41× speedup

  - **SpiderSQL:** 2.50 entropy → 2.19× speedup
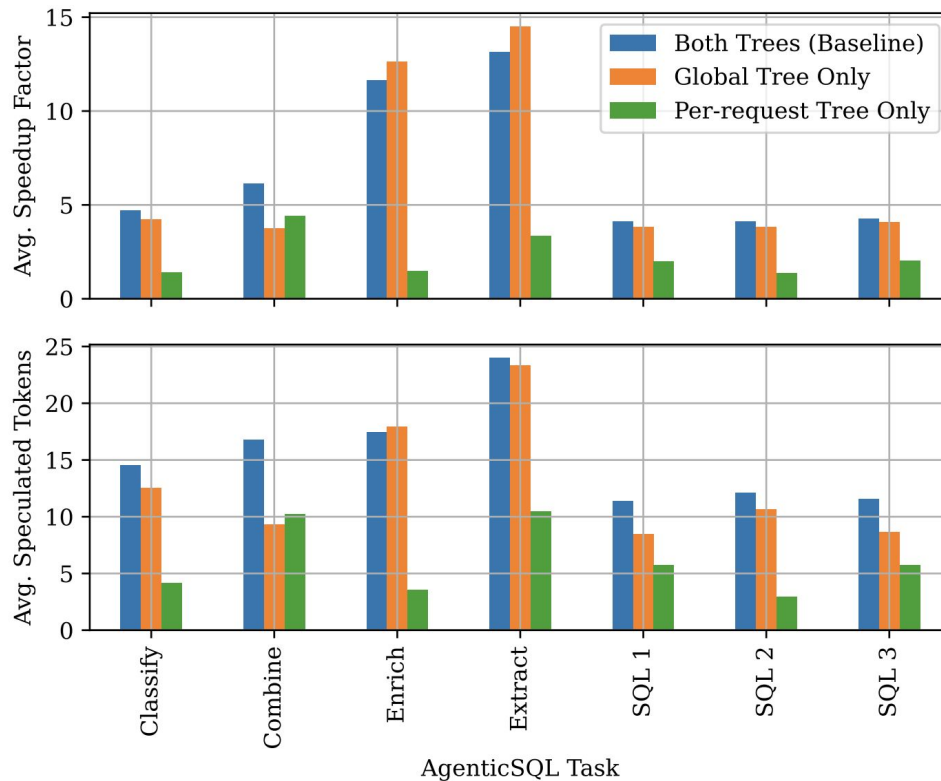
  - **WildChat (open-ended chat):** 3.43 entropy → modest speedup

| Dataset | Average Entropy |
|---|---|
| AgenticSQL (Enrich) | 0.171 |
| AgenticSQL (Classify) | 0.738 |
| AgenticSQL (Extract) | 0.0862 |
| AgenticSQL (SQL1) | 1.52 |
| AgenticSQL (SQL2) | 1.49 |
| AgenticSQL (SQL3) | 1.51 |
| AgenticSQL (Combine) | 1.49 |
| Spider | 2.50 |
| WildChat | 3.43 |
| Magicoder | 2.95 |

# What happens when the input distribution shifts?

🔍 Experiment:

1. **Train Cache** on **WildChat** (open-ended chat).
2. **Shift Workload:** Run inference on **SpiderSQL**
3. **Adapt:** We let SuffixDecoding *add* the new SpiderSQL outputs to the tree



📊 Results:

★ At 0 new examples, speedup is low (~1.5x) but *still better than vanilla*.
★ SuffixDecoding adapts *fast*. After 500 examples, already matching perfect cache

👉 **Takeaway:** The system is robust to distribution shift and adapts online *automatically*.

# Agenda

# Conclusion

✓ SuffixDecoding achieves **5.3x end-to-end speedup** for agentic workloads.

✓ Requires **no model training.**

✓ Maintains **lossless** output quality.

✓ Can be **hybridized** with other speculation methods (e.g. EAGLE-3) for rapid open-ended generation.

✓ Available in vLLM

🏠 [Project Page](Project Page)

Questions? Contact: goliaro@cs.cmu.edu