# HW6

Tguo67

**Question(s):**

Read the paper **SpecInfer: Accelerating Large Language Model Serving with Tree-based Speculative Inference and Verification**, summarize the paper, specifically, including the points below:

- What are the motivations/challenges of this work?
- How does the design of this paper address the challenges?
- How does the paper evaluate its design (experiment settings, workloads, metrics)?
- How does the evaluation prove its claims?

Note that it is essential to logically connect the motivation, design, and evaluation, rather than merely listing some points.

Related link:

Paper of SpecInfer: https://dl.acm.org/doi/pdf/10.1145/3620666.3651335

**Answers**:

**Motivation and challenges:**

LLM serving is slow because incremental decoding makes many full passes over model weights to emit a single token. Although speculative decoding helps but it is still challenging to draft many accepted tokens, also acceptance length and verifier efficiency remain the bottlenecks.

**Solutions to the challenges:**

SpecInfer replaces "decode-one-token" with "verify-many at once": (1) it builds a token tree of candidate continuations using multiple small speculative models with diversity/boost-tuning; (2) the target LLM acts as a tree verifier, using tree-based parallel decoding to compute attention for all tree nodes in one pass while reusing KV for shared prefixes.

**Evaluation settings, workloads, metrics:**

Test distributed multi-GPU and offloading regimes on standard prompt workloads, and compare end-to-end latency and number of verifier (LLM) decoding steps against incremental and sequence-speculative baselines.

**Evaluation results:**

SpecInfer cuts end-to-end latency by ~1.5–2.8× (distributed) and ~2.6–3.5× (offloading) while preserving outputs exactly—showing the verifier-as-tree check doesn't trade off quality. Gains coincide with fewer large-model passes and reduced kernel/IO overhead, matching the design's intent (more acceptance via diverse trees; less memory traffic via parallel verification and KV reuse).

**Conclusion:**

LLM serving is bottlenecked by serial, memory/communication heavy decoding, so SpecInfer builds

diverse token trees with small speculative models and then uses the big LLM as a parallel tree verifier (tree-based decoding with KV reuse and a light scheduler). This design achieves about 1.5–2.8× and 2.6–3.5× lower end-to-end latency with identical outputs