



My Contact Details

Reasoning models improve capability but introduce significant memory and efficiency challenges

ThinkKV decomposes CoT into distinct thought types and thought-adaptively performs quantization and eviction

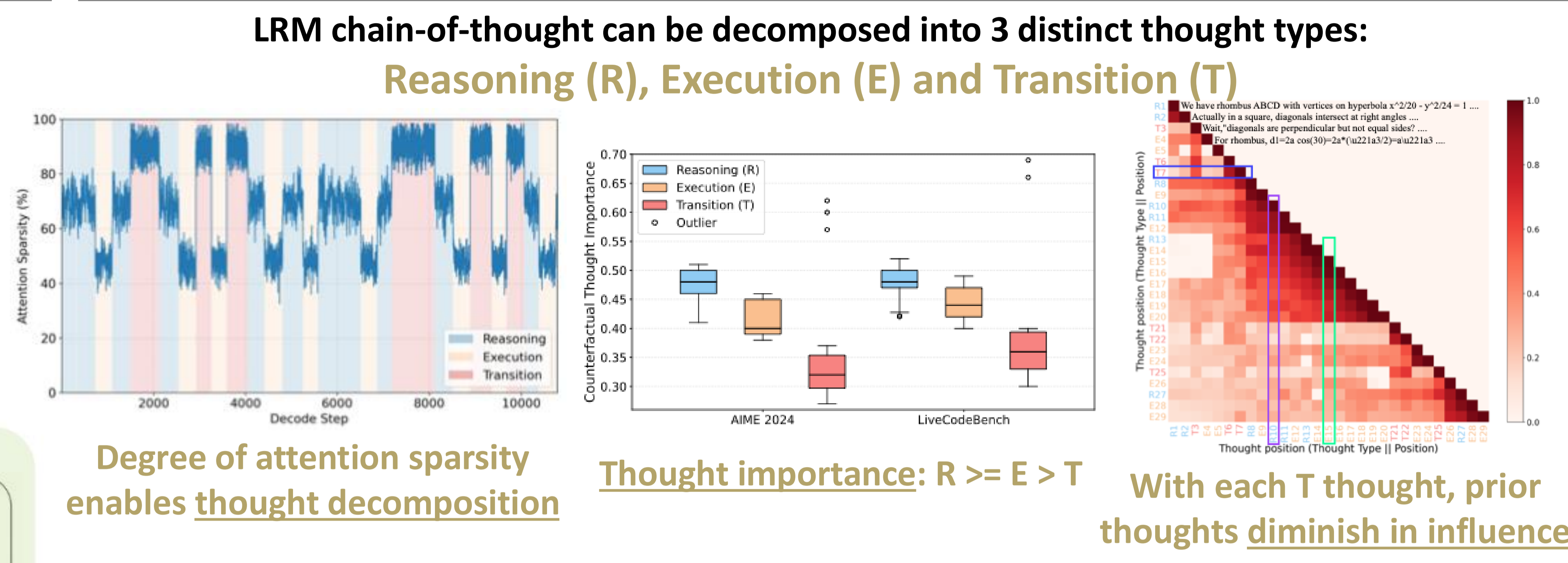
ThinkKV establishes a new pareto frontier

What are Reasoning Models ?

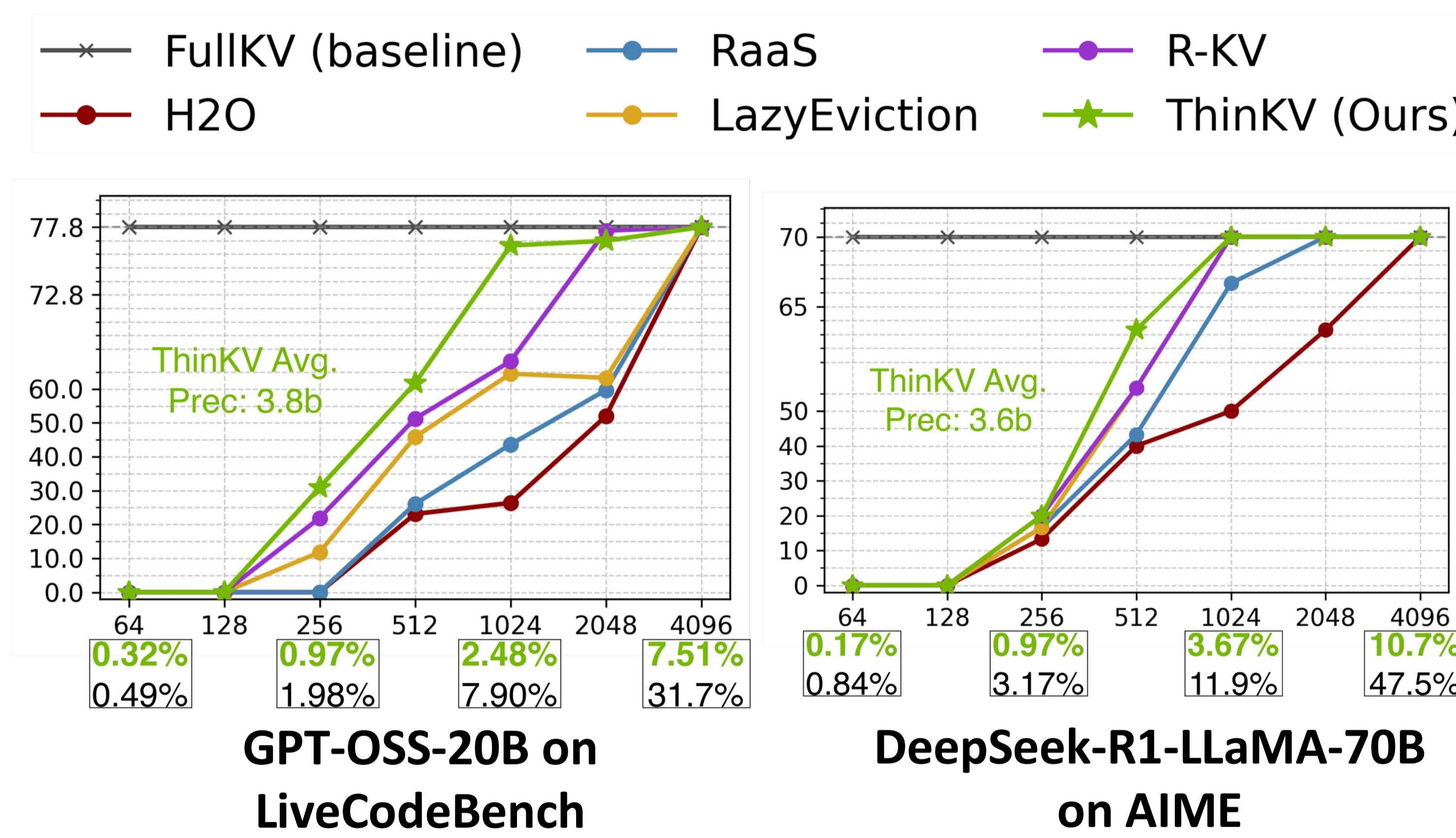
Zero-Shot w/o Reasoning
 Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
 A: The answer (arabic numerals) is
 (Output) 8 ✗

Zero Shot w/ CoT Reasoning (System 2)
 Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
 A: Let's think step by step.
 (Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

ThinkKV's Key Observations



Accuracy and Performance Analysis



Unique Challenges of LRMs

Long output context generation

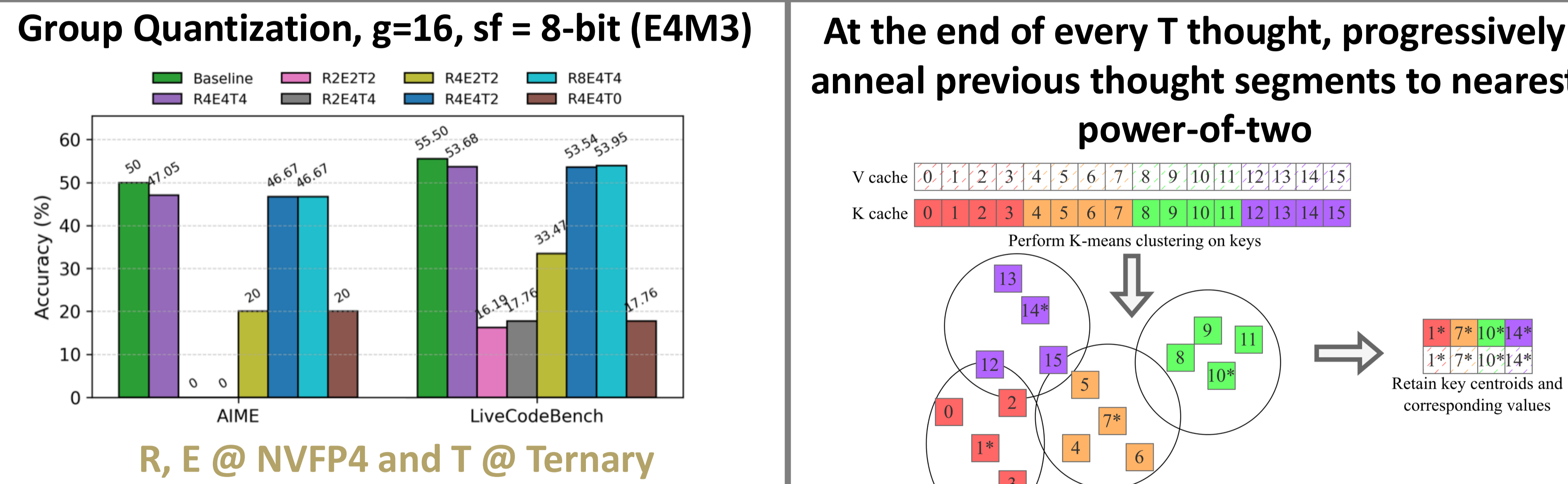
Convolved token importance

Internal memory fragmentation

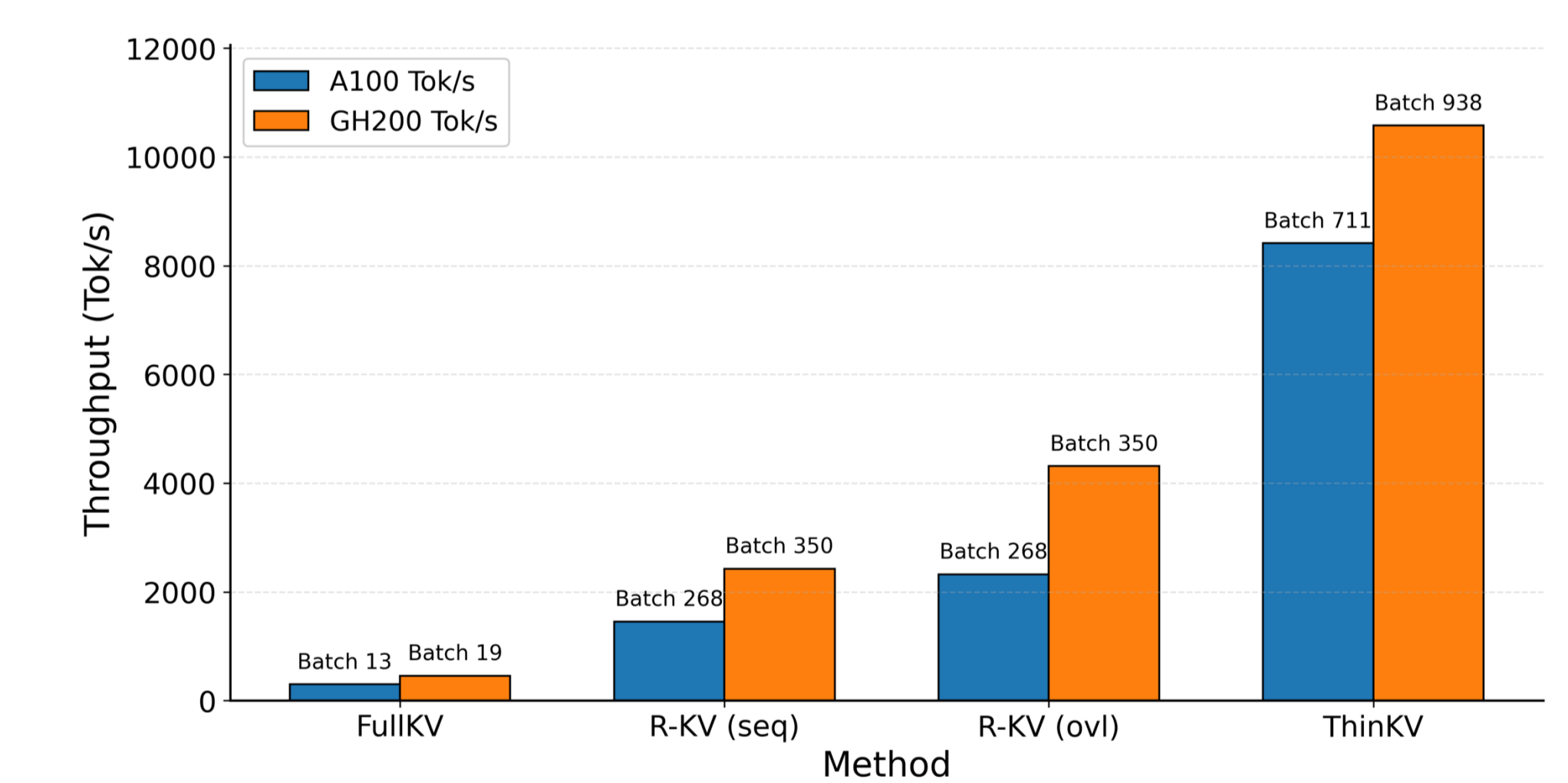
Gather overhead is significant

Model	Weights (GB)	KV Cache (GB)
DeepSeek-LLaMA-8B	16	64
DeepSeek-Qwen-32B	64	128
DeepSeek-LLaMA-70B	140	160

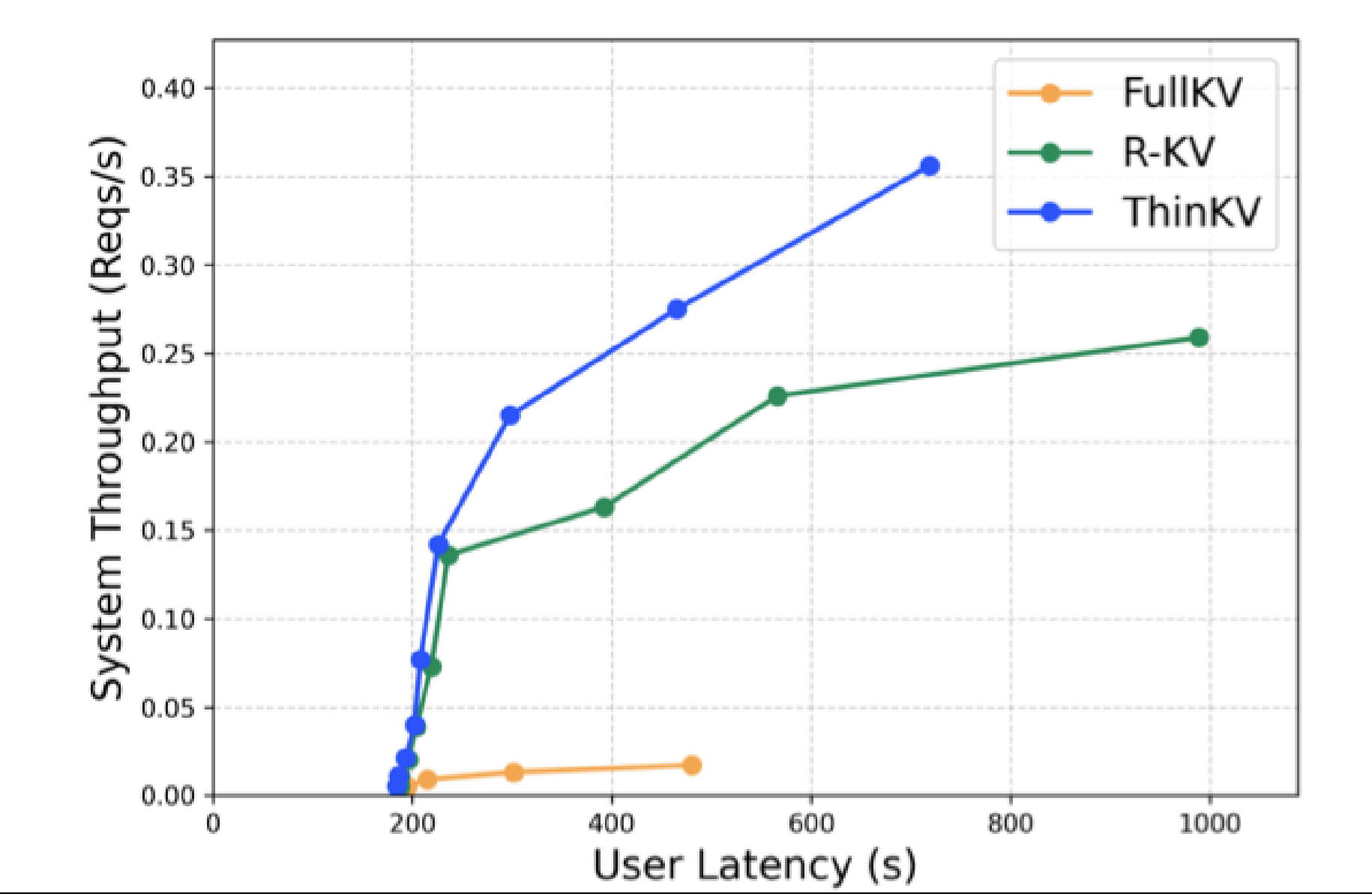
Think Before you Quantize, Think Before you Evict



ThinkKV achieves near lossless accuracy with < 3.67% of FullKV memory



ThinkKV supports 3X higher batch size and up to 5.8X higher TPS



ThinkKV achieves 38% higher reqs/s and 27% lower latency compared to R-KV

Continuous Thinking Kernel: Extending vLLM

New Block Table Fields

- Thought type:** Thought type of tokens in a block.
- Start indices:** Records the start position of the thought segment.
- Segment masks:** If there are multiple start indices, the segment mask is a bit vector.
- Eviction mask:** A bit vector marking positions of tokens evicted.