# Retrieval Optimization with Semantic Cache
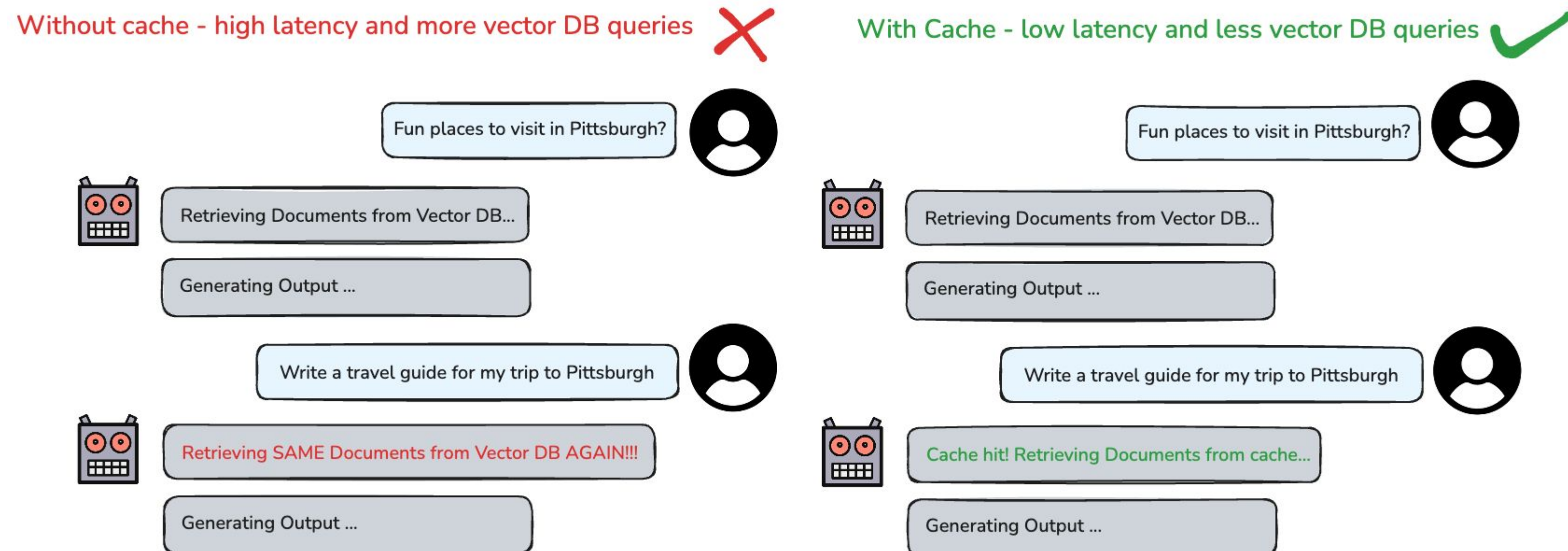
Ziming Wang, Tim Han, Dingfan Zheng, Lei Li
{peterwa2,timhan,dingfanz,leili}@cs.cmu.edu
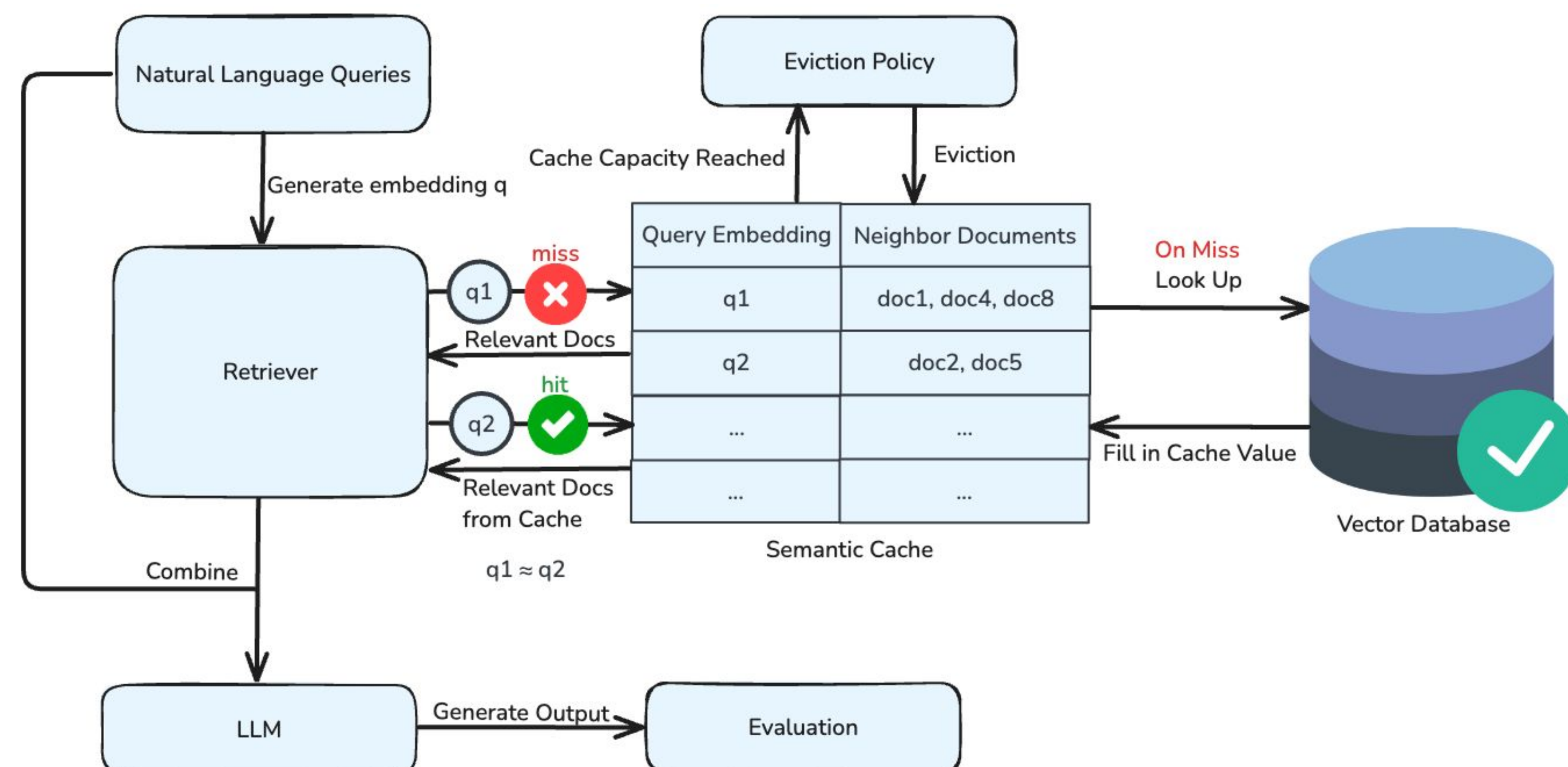
Carnegie Mellon University

## Caching in LLM RAG Applications

- Previous work caches LLM outputs, not fine-grained enough
- Semantically similar question retrieve similar documents, cache the query and <u>retrieved document</u> instead!



Without cache - high latency and more vector DB queries ✗

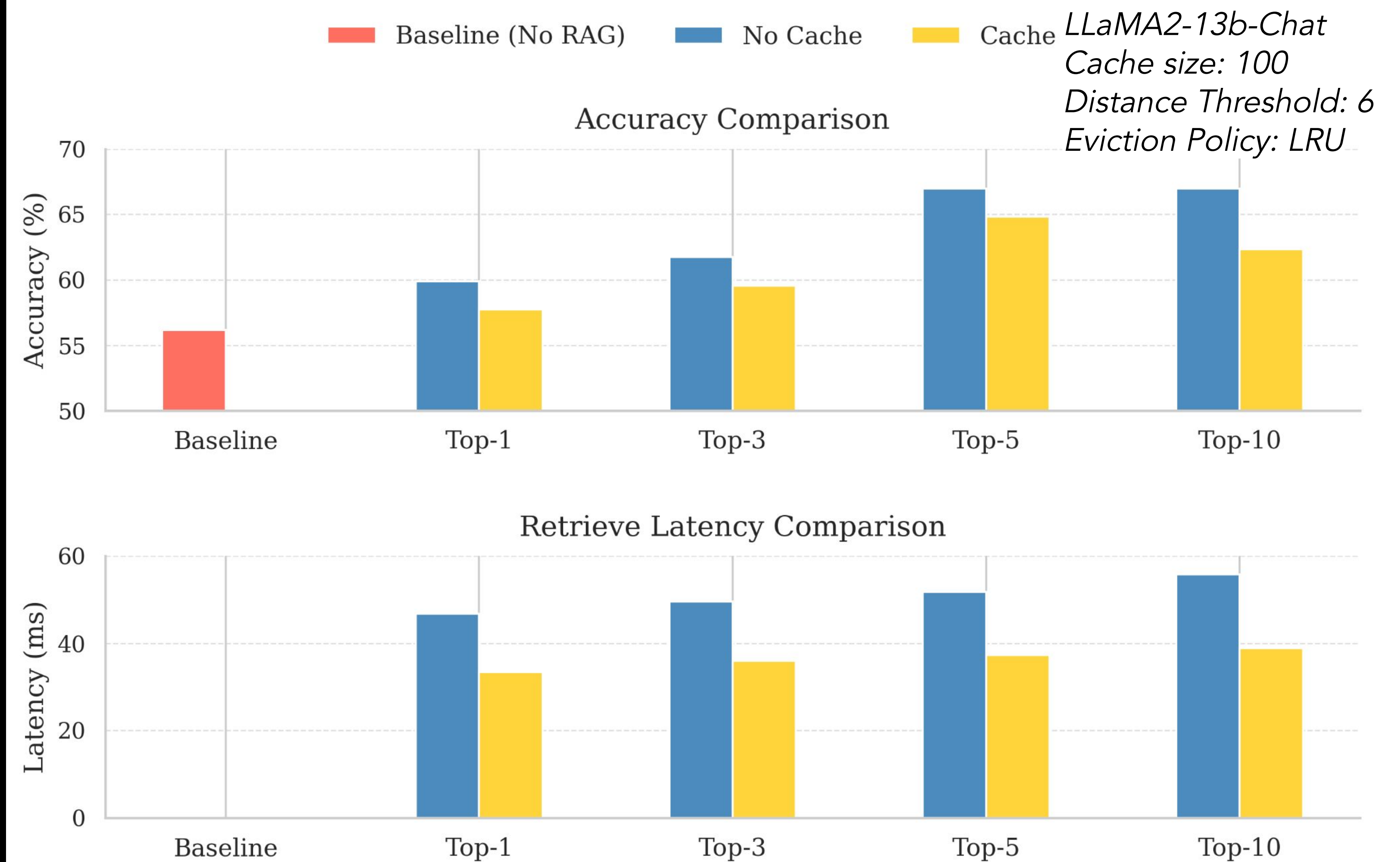With Cache - low latency and less vector DB queries ✓

## Semantic Cache System Design

- Reuse topk results for similar queries based on distance threshold.
- Tunable cache size, threshold, and eviction policies for optimization.
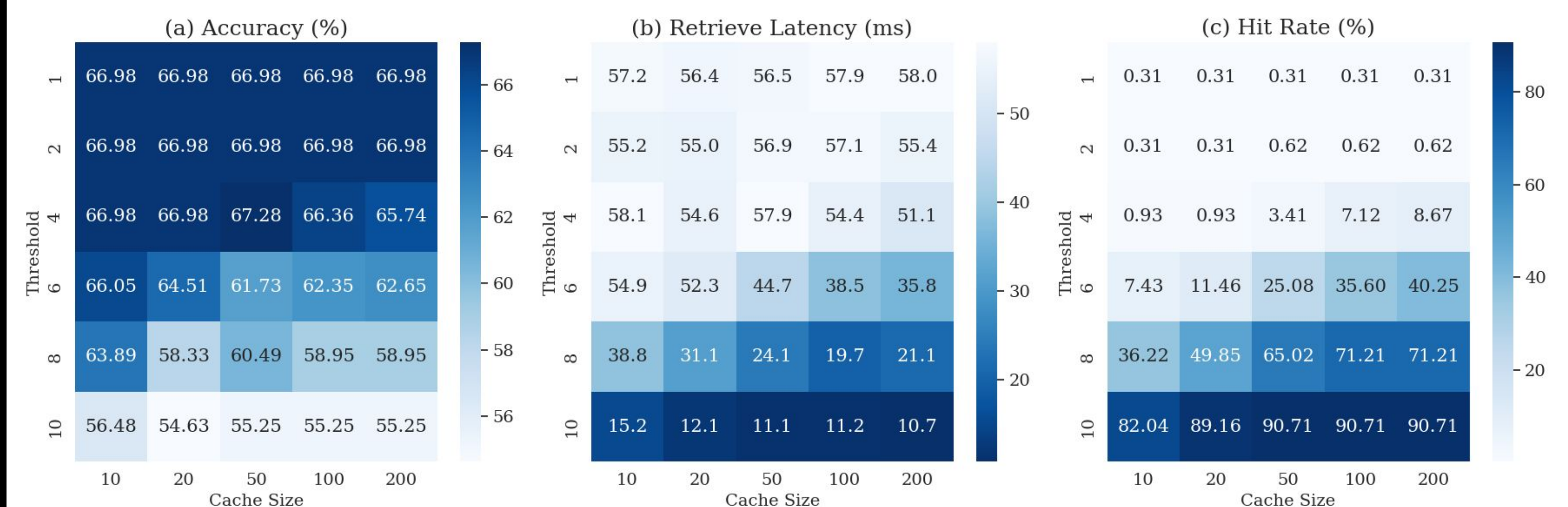


## Results

- Our approach cuts retrieval latency by > 30%, while maintaining high retrieval quality with only a 2–4% drop in accuracy.

Baseline (No RAG) | No Cache | Cache

LLaMA2-13b-Chat
Cache size: 100
Distance Threshold: 6
Eviction Policy: LRU



Accuracy Comparison

Retrieve Latency Comparison

- Impact of Cache Size and Threshold on Performance Metrics



(a) Accuracy (%)

(b) Retrieve Latency (ms)

(c) Hit Rate (%)

## Key Takeaway / Advantages

- Reduces retrieval latency and the number of API calls
- Tuning cache configs based on workloads and resources