

# RetroInfer: A Vector-Storage Approach for Scalable Long-Context LLM Inference

*Yaoqi Chen, Jinkai Zhang, Baotong Lu , Qianxi Zhang, Chengruidong Zhang, Jingjia Luo, Di Liu, Huiqiang Jiang, Qi Chen, Jing Liu, Bailu Ding, Xiao Yan, Jiawei Jiang, Chen Chen, Mingxing Zhang, Yuqing Yang, Fan Yang, Mao Yang*

**Presented by Xiaoqi Li  
2025-06-03**



# Contexts

**Background**

**Motivation**

**RetroInfer**

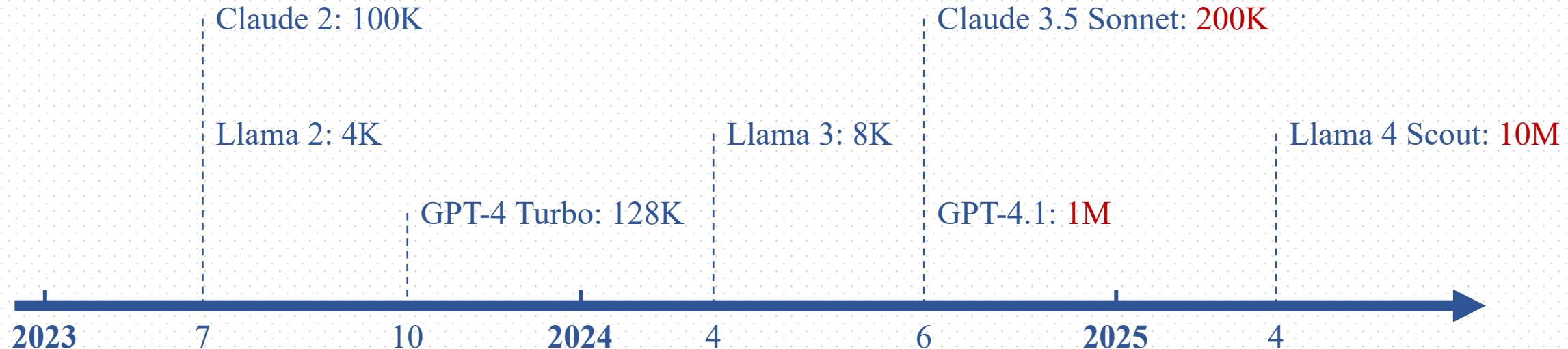
**Evaluations**



# Longer and Longer Context Window in LLMs

## ◆ Longer context window

- ◆ Improved understanding and factual accuracy
- ◆ Enhanced summarization and information retrieval
- ◆ Advanced in-context learning





# Expensive Long Context Serving

- ❑ **O( $n^2$ )** Inference latency to the context size n
- ❑ Most of the time is spent on **attention** operations

Prompt Length	128K	256K	512K	1M
Total Latency (s)	32.8	111	465	1,765
FFN (s)	7.6	15	31	70
Attention (s)	25.2	96	434	1,695

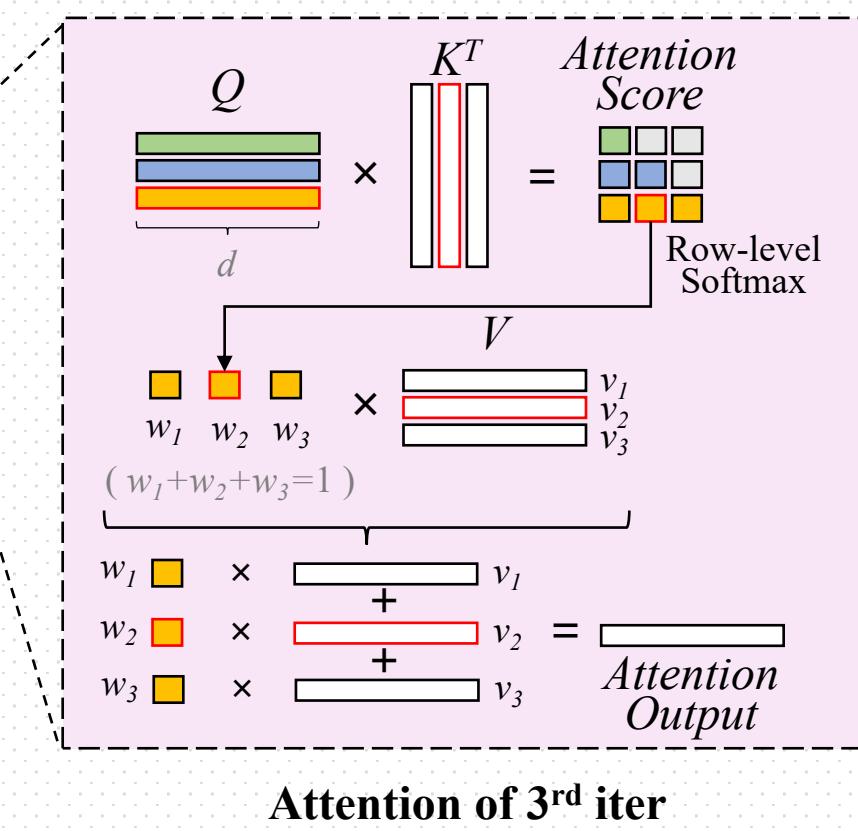
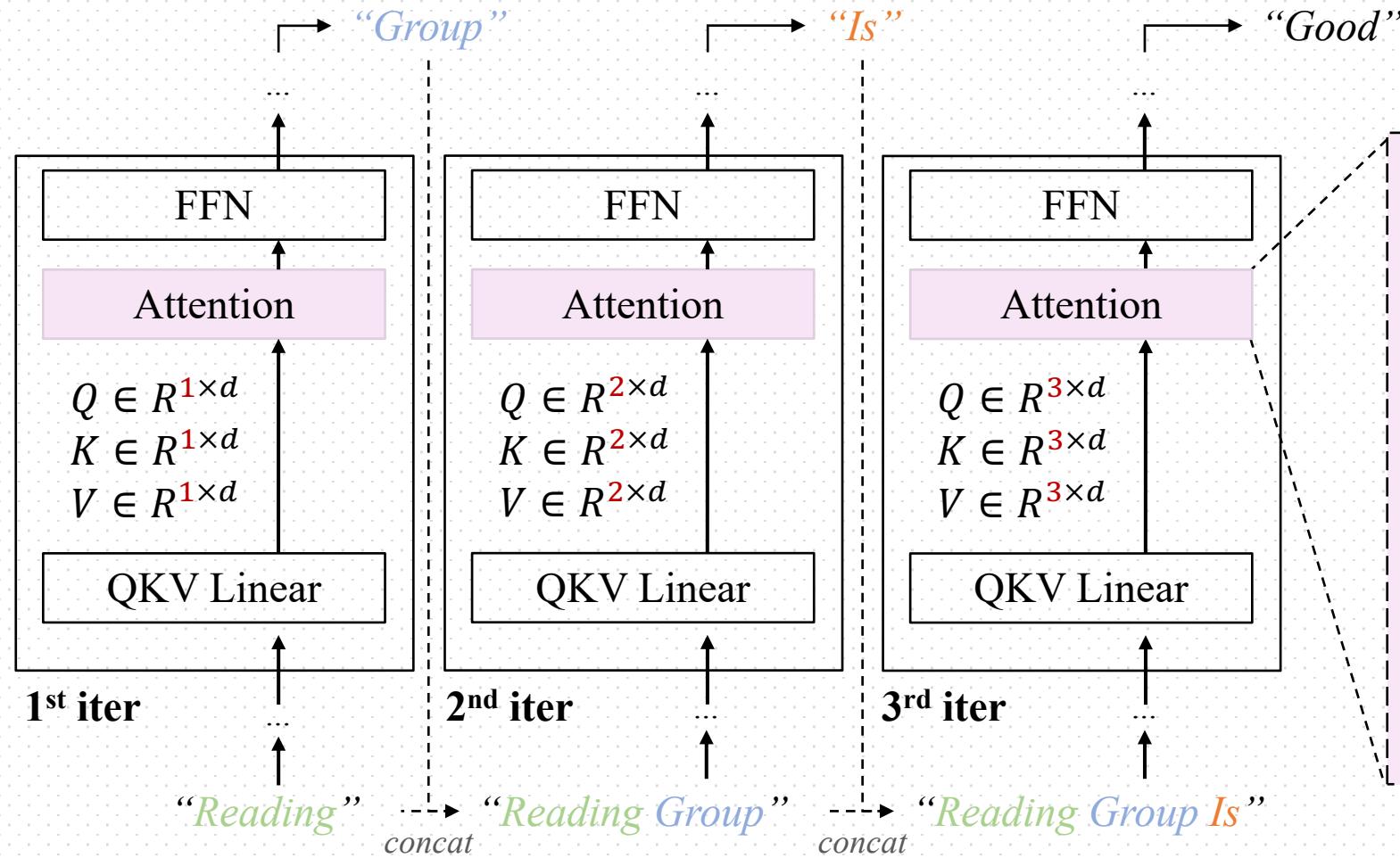
**1695 / 1765 ≈ 96%**

Decoding latency of Llama-3-8B  
across different context lengths on one A100 GPU<sup>[1]</sup>

[1] RetrievalAttention: ACCELERATING LONG-CONTEXT LLM INFERENCE VIA VECTOR RETRIEVAL

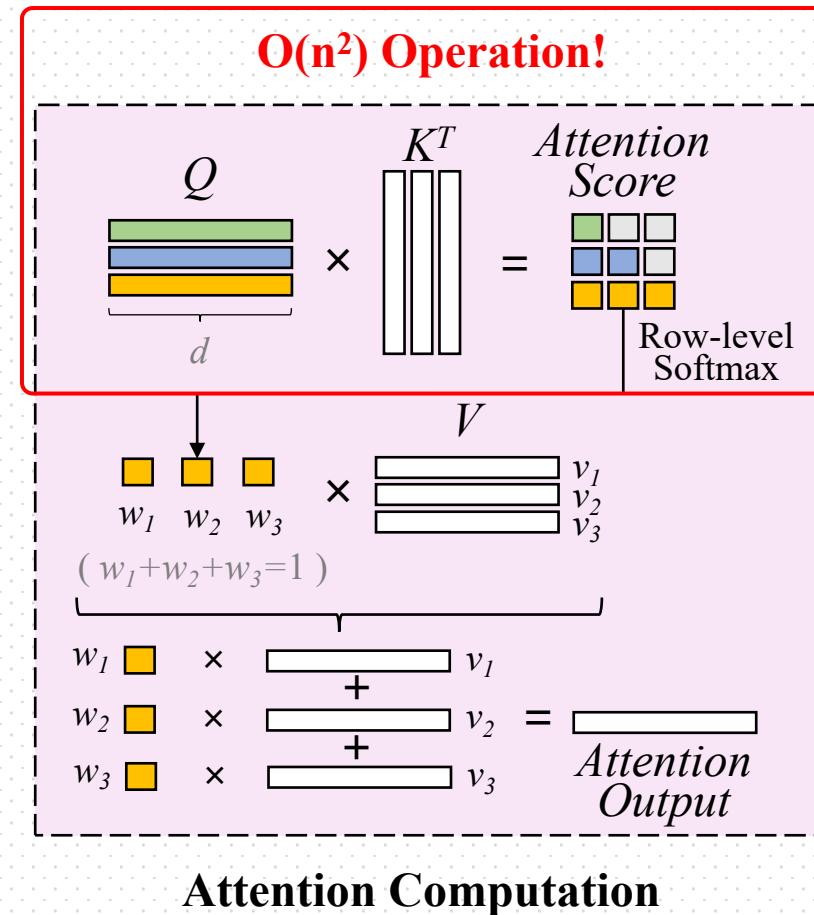


# Vanilla Attention: $O(n^2)$ Complexity





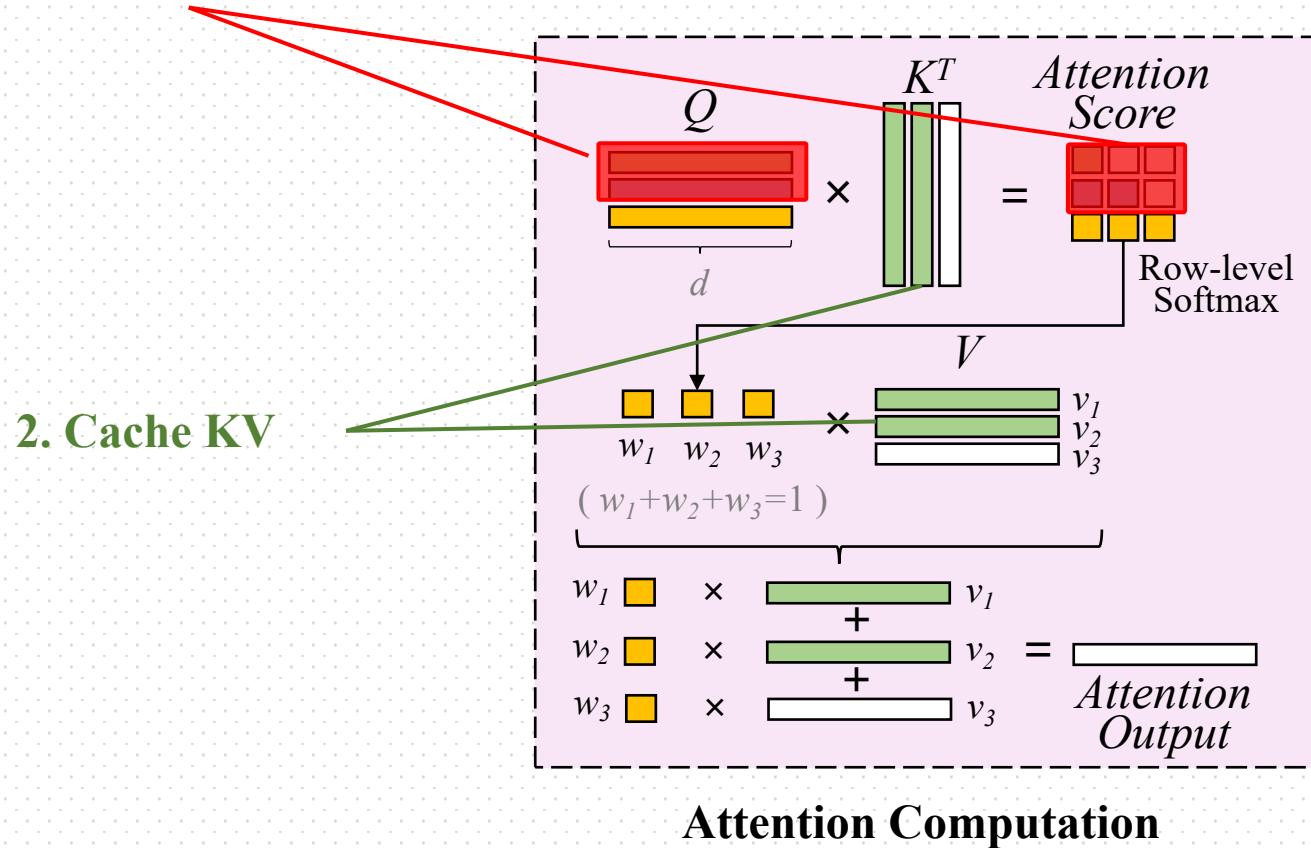
# Vanilla Attention: $O(n^2)$ Complexity





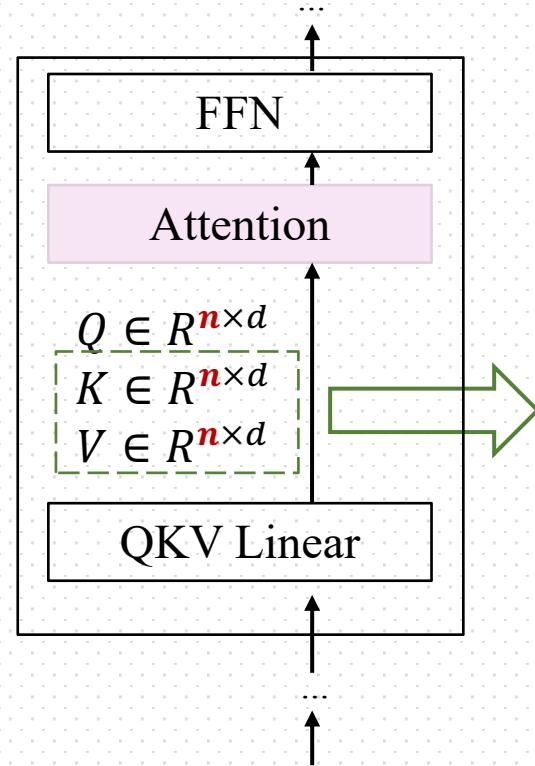
# Optimize Vanilla Attention

## 1. Remove Unneeded Operation



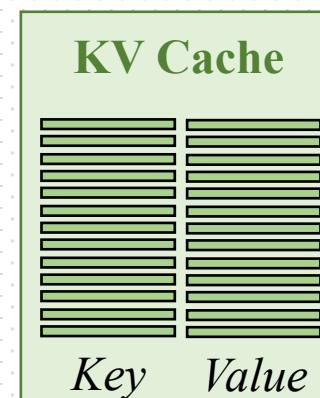


# Using KV Cache to Accelerate Decoding

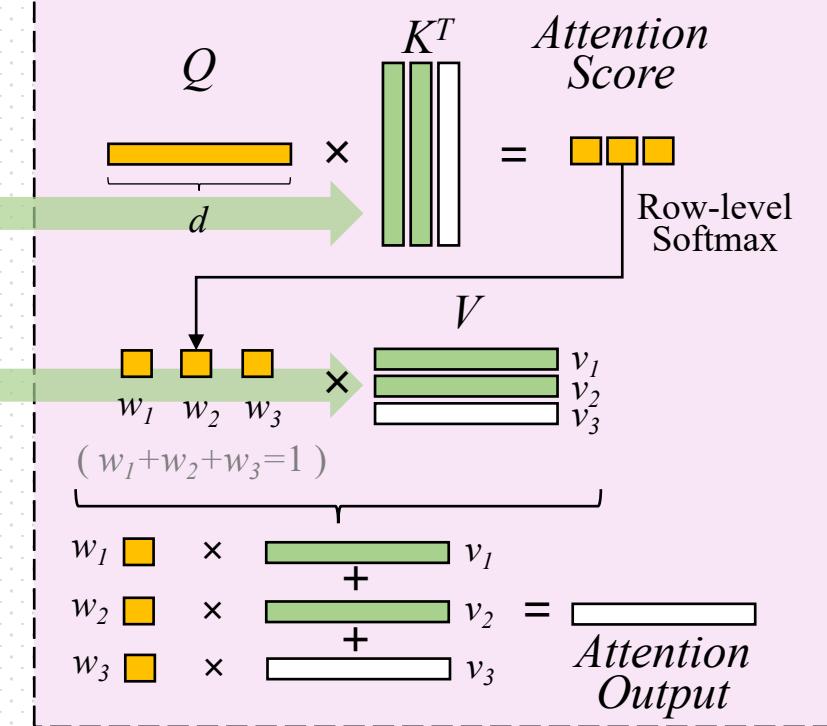


*"This is a very long prompt..."*

## Prefilling Stage



## Decoding Stage Attention Computation





# But KV Cache Is Not Enough!

❑ KV cache grows linearly with context window length!

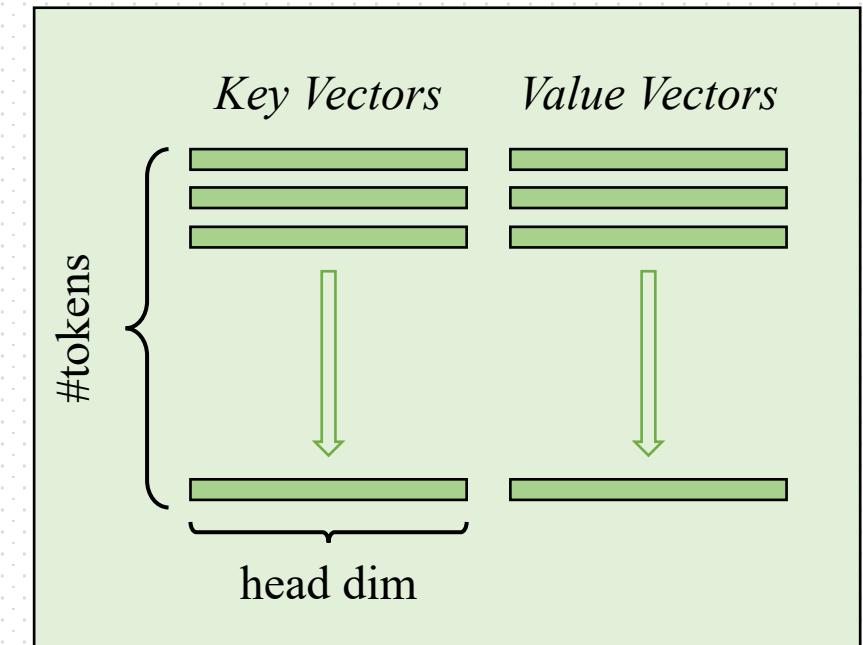
LLM Model: Llama3-8B			
#Layers	Hidden Dim	GQA Compress Ratio	Data Type
32	4096	4	BF16

KV Cache size:

- A single token:  $2 * 32 * 4096 / 4 * 2B = 128KB$
- 1M tokens:  $128KB * 1M = 128GB$

GPU memory capacity:

- RTX5090: 32GB
- A100: 40/80 GB



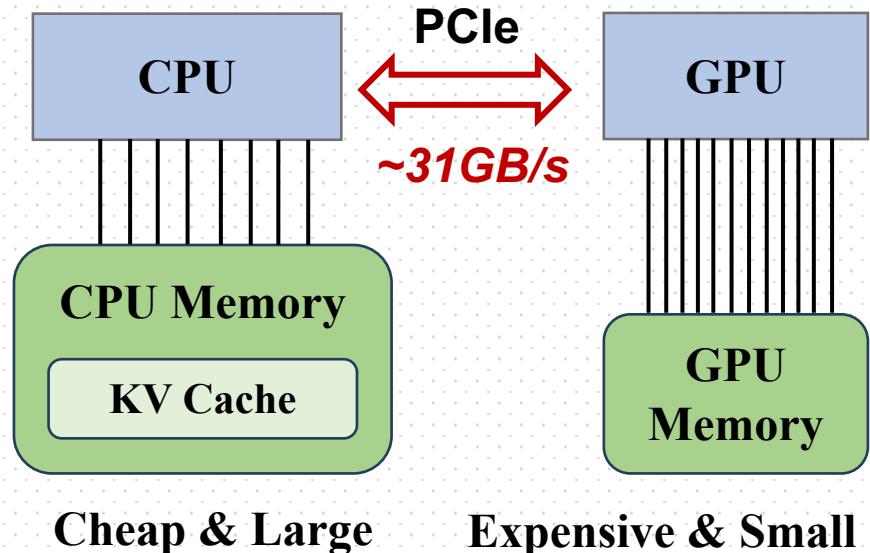
KV Cache



# KV Cache Optimizations

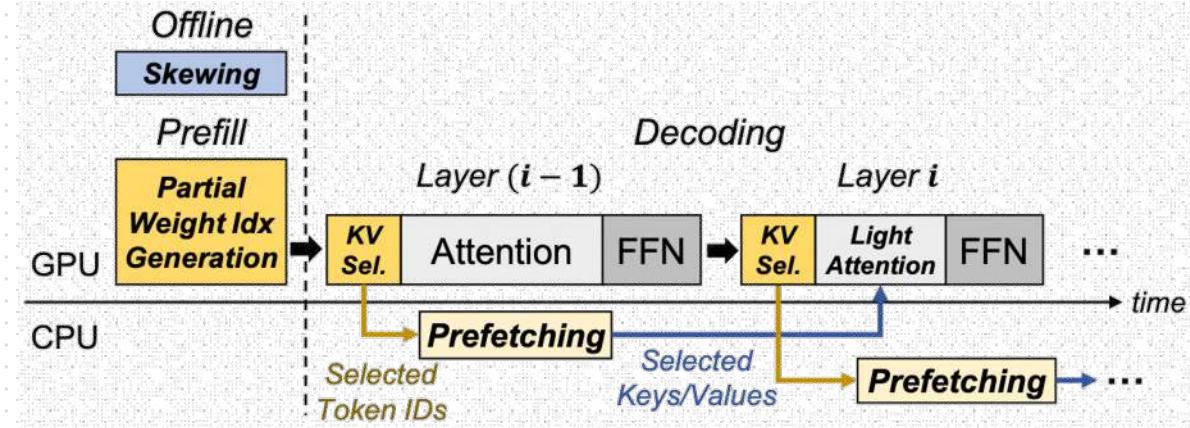
## □ Opt. 1: Offload

- ❖ Offload KV cache to CPU DRAM
- ❖ PCIe bandwidth is limited



## □ InfiniGen (OSDI 24)

- ❖ Prefetching next layer's KV vectors



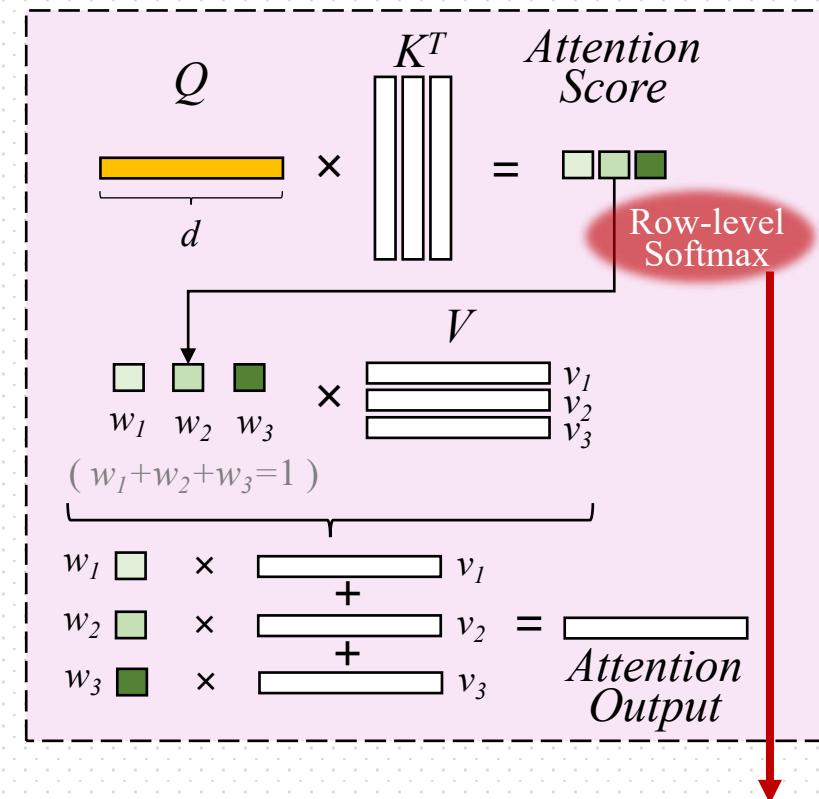
InfiniGen leverages prefetching to overlap data transfer



# KV Cache Optimizations

## Opt. 2: TopK Attention

- ❖ Attention is sparse<sup>[1]</sup>.
- ❖ In most transformer layers, n% top KV can carry enough information to maintain model accuracy.
- ❖ Only use top n% token with largest attention weights



Before softmax:	-0.61	-0.02	0.31
After softmax:	0.19	0.34	0.47

[1] Memory-efficient Transformers via Top-k Attention

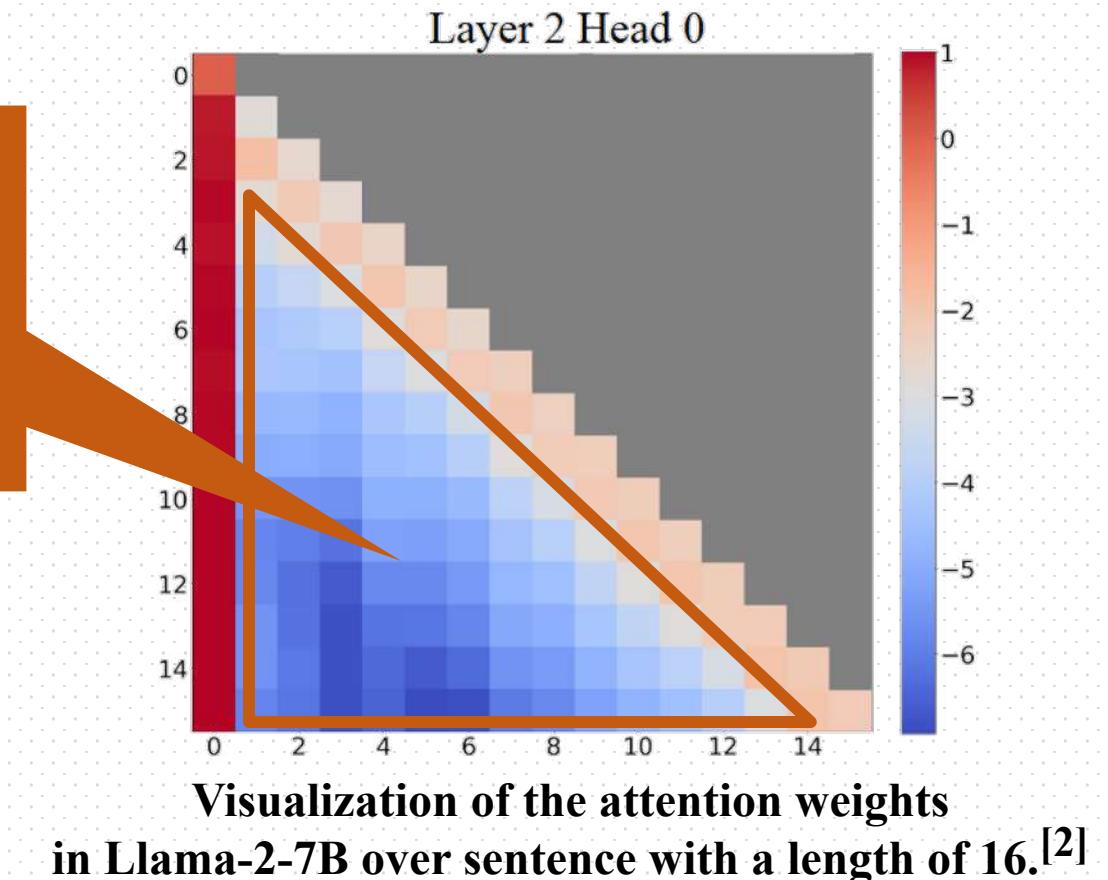


# KV Cache Optimizations

## □ Opt. 2: TopK Attention

- ❖ Attention is sparse
  - ❖ In most cases, we only need top KV pairs to get good information and maintain accuracy.
  - ❖ Only use top n% token with largest attention weights
- Drop unimportant keys  
(StreamingLLM, H2O, ...)
- Easy to implement
  - Memory saving

Sounds good, but...



[1] Memory-efficient Transformers via Top-k Attention

[2] StreamingLLM: Efficient Streaming Language Models with Attention Sinks

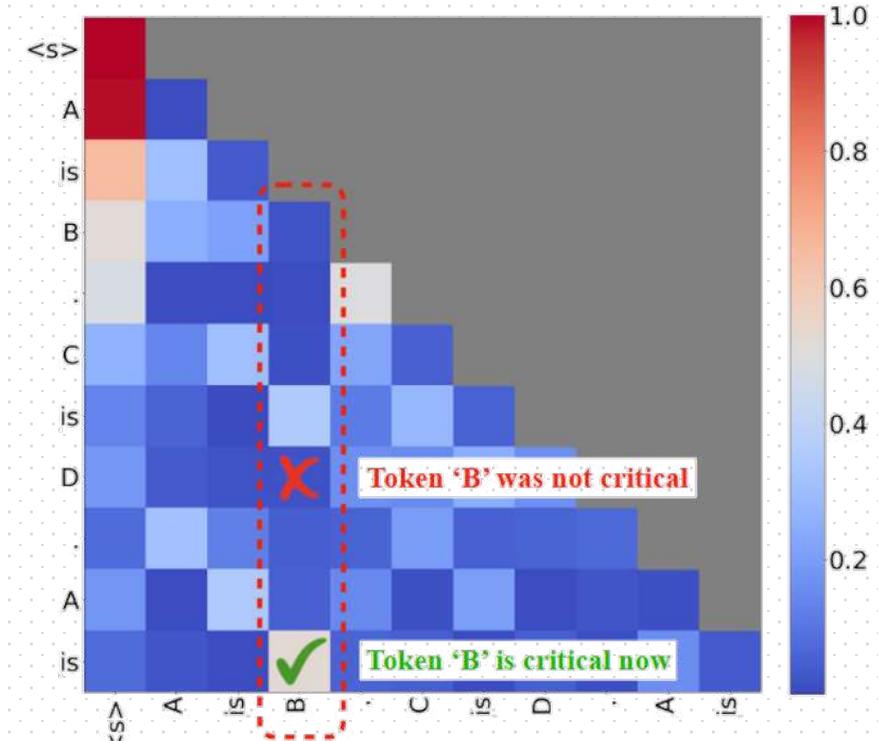


# KV Cache Optimizations

## □ Opt. 2: TopK Attention

- ❖ Attention is sparse<sup>[1]</sup>.
- ❖ In most transformer layers, n% top KV can carry enough information to maintain model accuracy.
- ❖ Only use top n% token with largest attention weights

Critical tokens depend on the query



Visualization of the attention weights of prompt “A is B. C is D. A is”.<sup>[2]</sup>

[1] Memory-efficient Transformers via Top-k Attention

[2] QUEST: query-aware sparsity for efficient long-context LLM inference

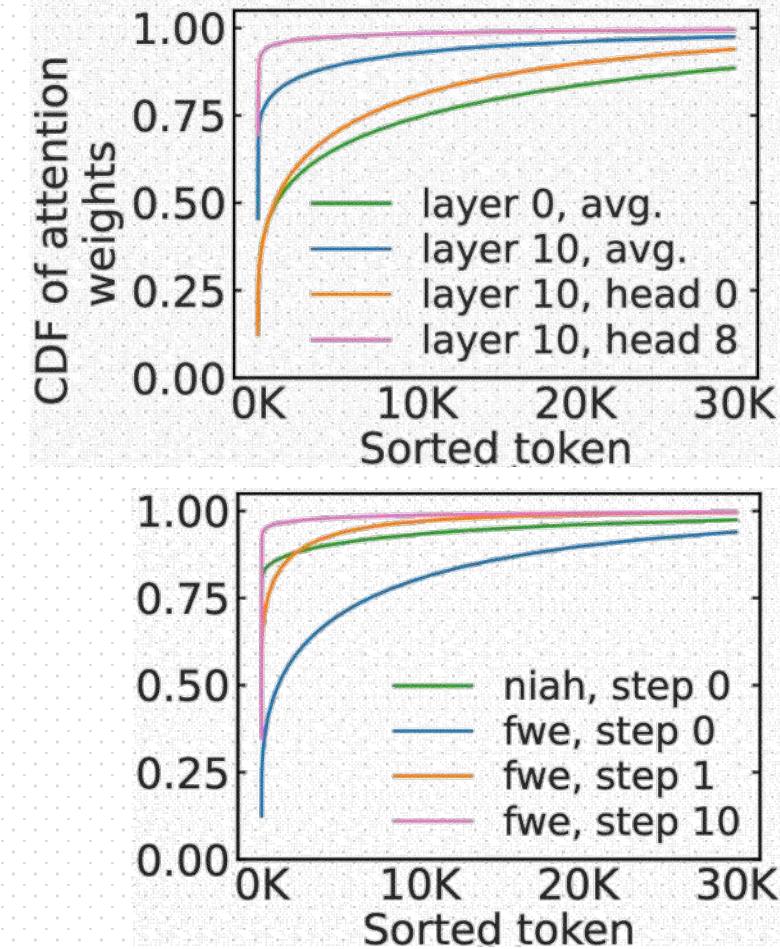


# KV Cache Optimizations

## □ Opt. 2: TopK Attention

- ❖ Attention is sparse<sup>[1]</sup>.
- ❖ In most transformer layers, n% top KV can carry enough information to maintain model accuracy.
- ❖ Only use top n% token with largest attention weights

The sparsity exhibits **high variability**



Variety of Attention Sparsity

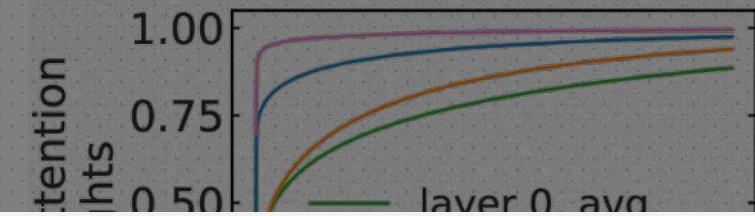
[1] Memory-efficient Transformers via Top-k Attention



# KV Cache Optimizations

## □ Opt. 2: TopK Attention

- ❖ Attention is sparse<sup>[1]</sup>.

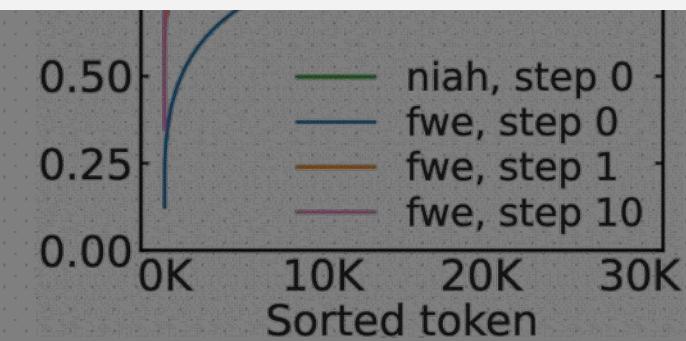


Static methods is inadequate

We need some strategy to find topk token from **full KV** efficiently

Only use top K/N token with  
largest attention weights

The sparsity exhibits **high variability**



Variety of Attention Sparsity

[1] Memory-efficient Transformers via Top-k Attention

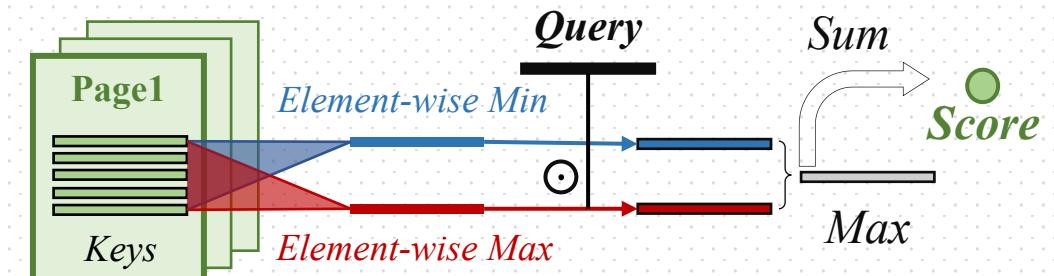


# KV Cache Optimizations

## □ Opt. 2: TopK Attention

## □ Quest (ICLR 24)

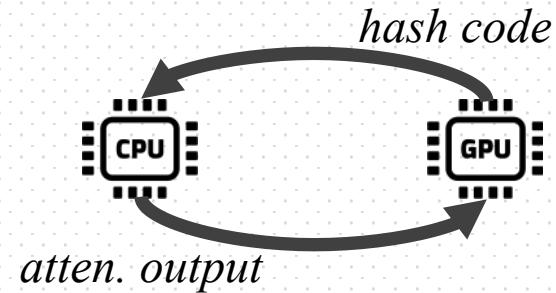
- ❖ KV vectors stored in pages, each with a pair of min/max key vector



Quest use min/max key to find topk pages

## □ MagicPIG (ICLR 25)

- ❖ Use Locality Sensitive Hashing (LSH) to sample important Key/Value



MagicPIG leverages LSH to do sample



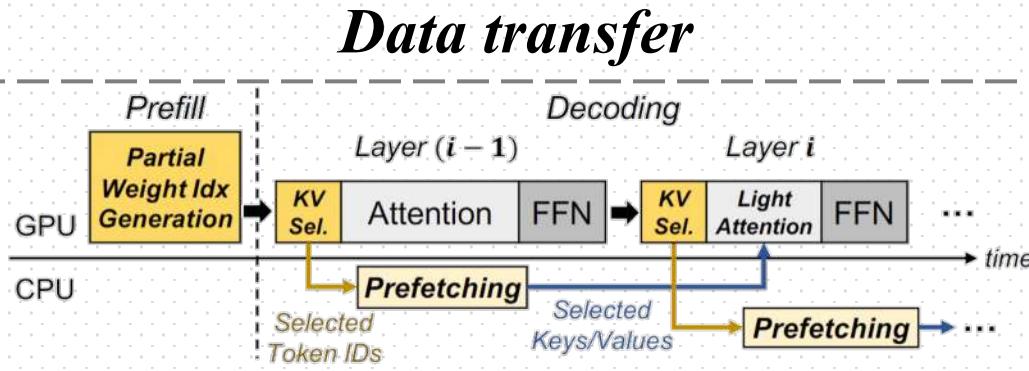
# Existing Long-Context Inference

- ❑ Existing approaches face challenges in accuracy and efficiency
- ❑ Drop-out based (*e.g. StreamingLLM*)
  - ❖ Assuming static token importance → Accuracy loss
- ❑ Other Leading Methods
  - ❖ See next page...



# Existing Long-Context Inference: SOTAs

InfiniGen  
OSDI 24

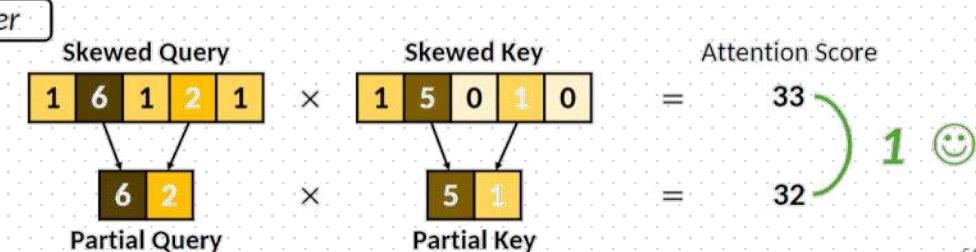


Prefetching, but frequent data transfer

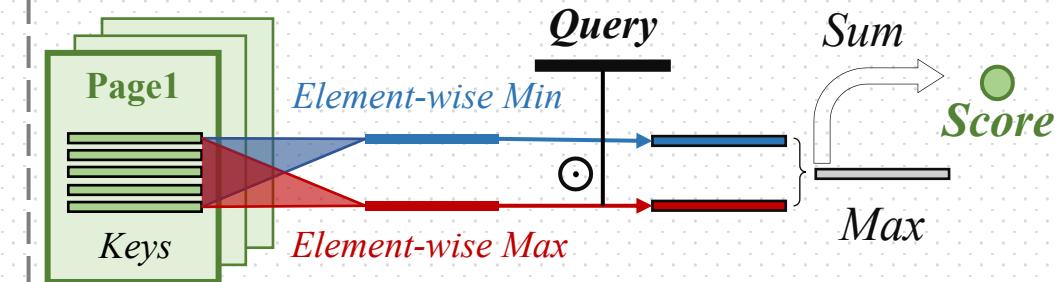
Quest  
ICLR 24



## TopK Selection



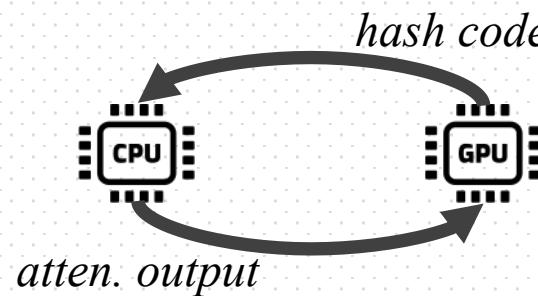
Use reduced dimension to find topk, still costly



Use page level reduced keys, accuracy loss

MagicPIG  
ICLR 25

Attention is  
done on CPU



Use Locality Sensitive Hashing (LSH)  
to sample important Key/Value,  
accuracy loss



# Contexts

Background

Motivation

RetroInfer

Evaluations

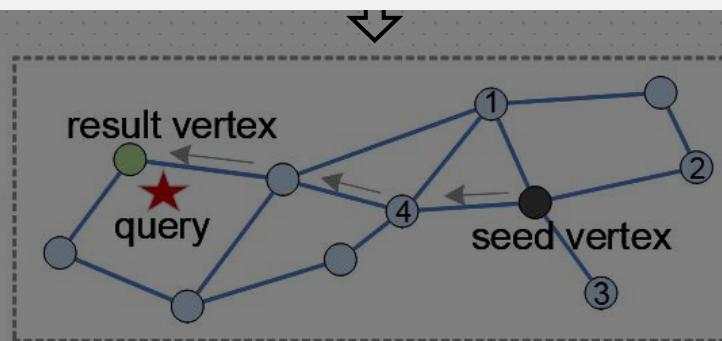


# Top-K Selection Is Natively Supported in ANNS

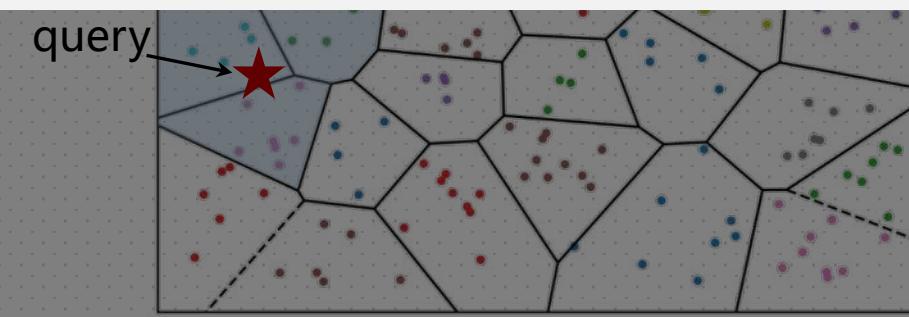
## □ ANNS (Approximate Nearest Neighbor Search)

- ❖ Say, <5% search yields >90% recall

## Can We Leverage ANNS?



Graph-based



Cluster-based



# Challenges: When ANNS Meets KV Cache

## □ Index traversal and construction costs

❖ Quest latency<sup>[1]</sup>: **3.61 s** (prefilling) and **14.86 ms** (decoding) on Llama-2-7B using A100 GPU, 32k context

Algo.	ANNS Test Results		KV Cache Estimations	
	Index construction	Avg. latency	Index construction (in prefilling stage)	Latency per token (in decoding stage)
IVF FLAT	562 ms	2.58 us	575.5 s	2.64 ms
HNSW	300 ms	1.02 us	307.2 s	1.04 ms

ANNS test achieving 0.9 recall  
on SIFT subset (32k samples, 128-dim, float16)  
CPU: Intel(R) Xeon(R) 6982P-C

[1] LServe: Efficient Long-sequence LLM Serving with Unified Sparse Attention



# Challenges: When ANNS Meets KV Cache

## □ Collaboration between GPU and CPU

### ❖ GPU

- Highspeed computation, fast memory access
- **Limited memory capacity**

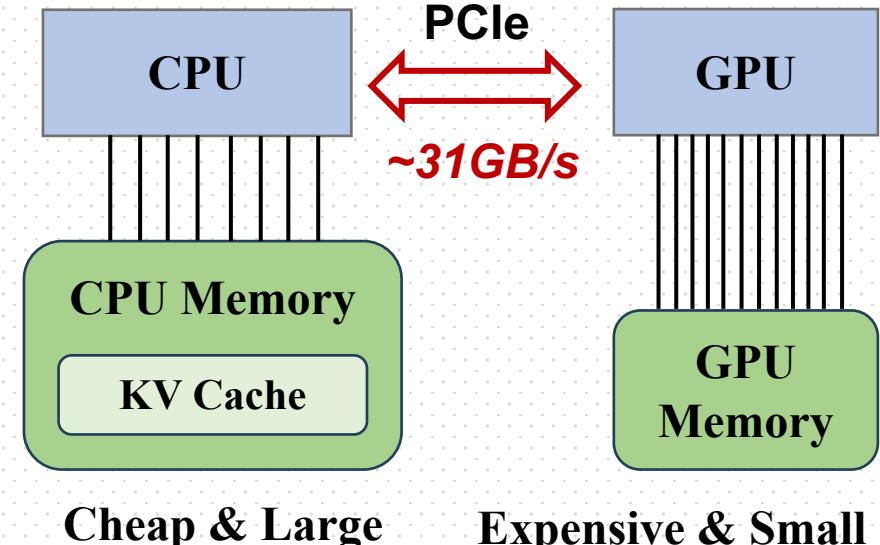
### ❖ CPU

- Abundant memory
- **Lower computational throughput**
- The control logic for managing index structures is often better suited for CPU execution

## □ Questions

### ❖ Where to place data and computation?

### ❖ How to leverage parallelism and overlap to fully utilize both processors?





# Contexts

❑ Background

❑ Motivation

❑ RetroInfer

❑ Evaluations



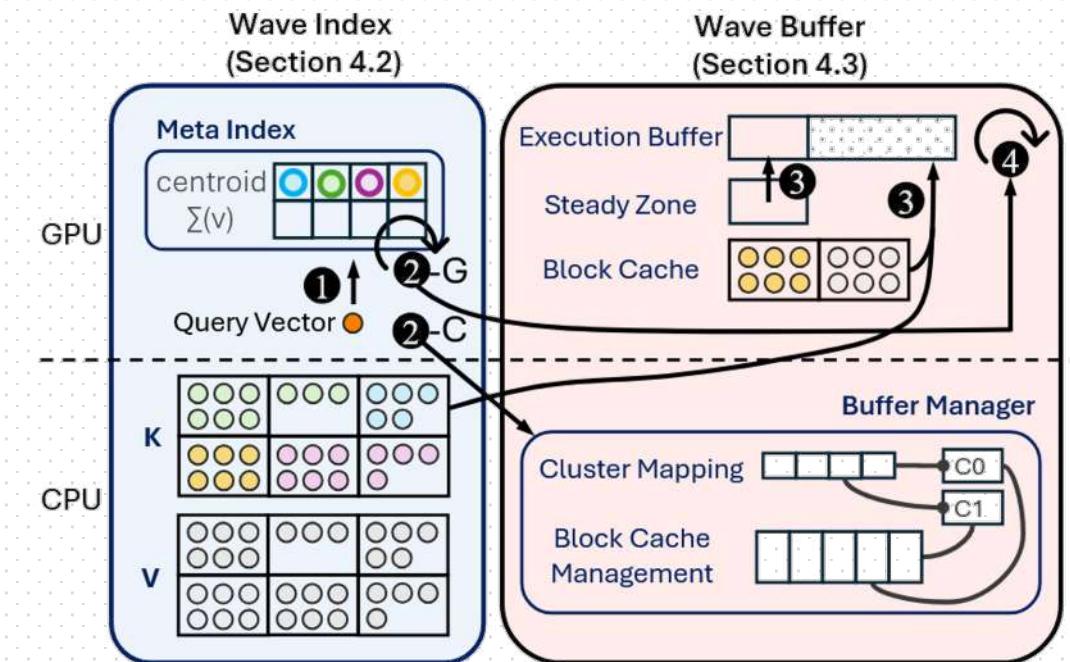
# RetroInfer

## ❑ RetroInfer

- ❖ A Vector-Storage System for KV Cache

## ❑ Key designs

- ❖ Wave index
  - Wave: Attention-aWare VEctor index
- ❖ Wave buffer



RetroInfer Overview



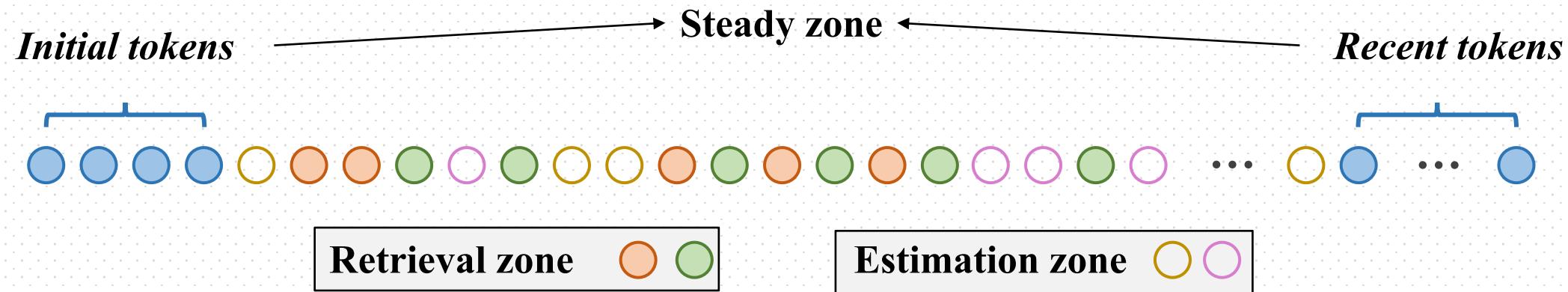
# Wave Index: 3 Zones

## □ Distinguish 3 zones within the token sequence

- ❖ Steady zone: initial + recent tokens<sup>[1]</sup>
- ❖ Retrieval zone: retrieve using cluster-based index
- ❖ Estimation zone: unretrieved clusters

Physical

} Logical



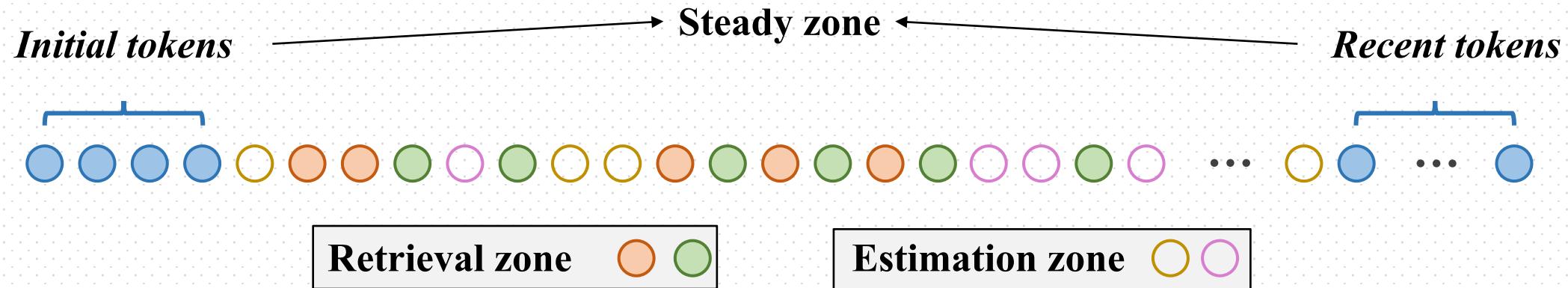
[1] StreamingLLM: Efficient Streaming Language Models with Attention Sinks



# Wave Index: 3 Zones

## □ Tripartite Attention Approximation

- ❖ Steady zone  $o^0$ : precise attention
- ❖ Retrieval zone  $o^1$ : precise attention
- ❖ Estimation zone  $o^2$ : estimated attention

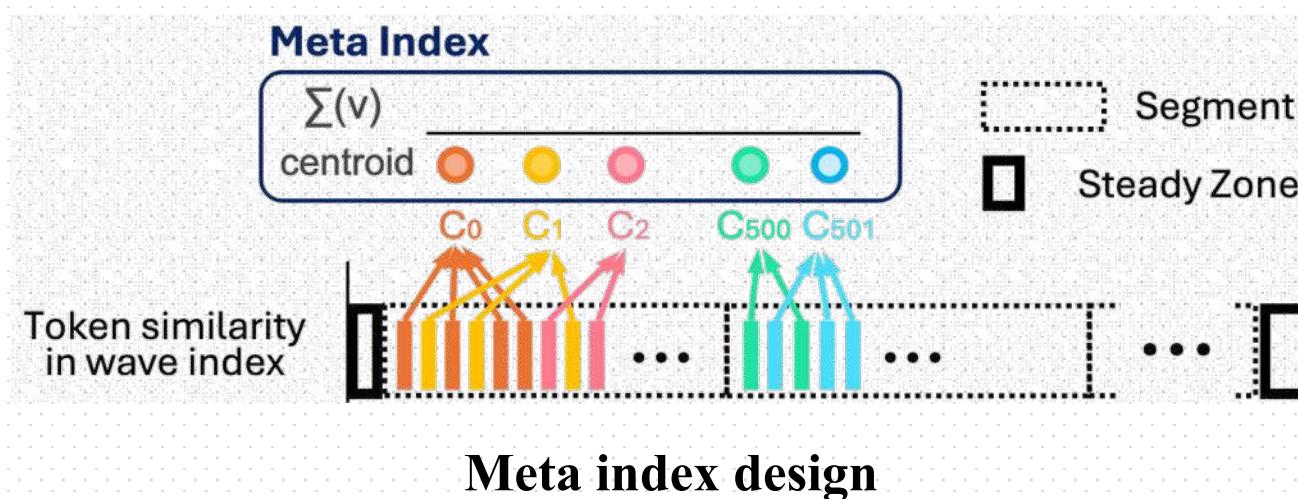




# Wave Index: Meta Index

## □ Meta index

- ❖ Spherical k-means clustering
- ❖ Segment-level clustering
  - To save index construction overhead
  - Based on the coarse-grained spatial locality ?





# Wave Index: Estimation Zone

## □ Accuracy-Bounded Attention Estimation

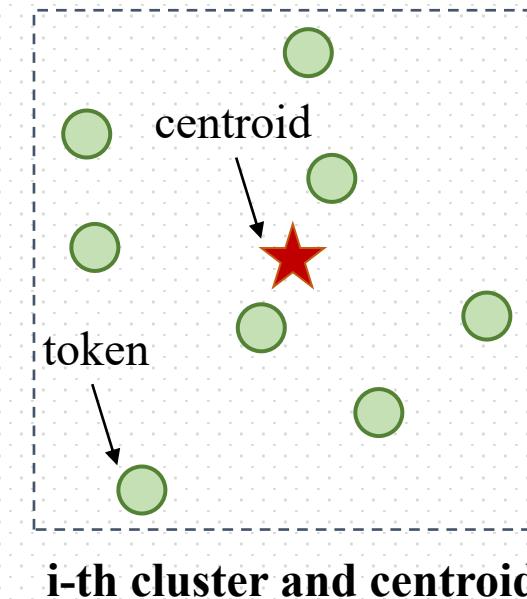
- ❖ For un-retrieved clusters, use centroids to generate estimate attention

use centroid's weight

$$\tilde{a}_i = \frac{e^{\frac{q \cdot c_i^T}{\sqrt{d}}}}{\sum_{j=1..p} e^{\frac{q \cdot K_j^T}{\sqrt{d}}} + \sum_{j=1..m} (s_j \cdot e^{\frac{q \cdot C_j^T}{\sqrt{d}}})}, i = 1..m$$

i-th centroid weight      p tokens in steady zone      all m clusters with  $s_j$  tokens each

$$token \text{ atten.} = \text{softmax} \left( \frac{qK^T}{\sqrt{d}} \right)_{token} v_{token}$$





# Wave Index: Estimation Zone

## □ Accuracy-Bounded Attention Estimation

- ❖ For un-retrieved clusters, use centroids to generate estimate attention

use centroid's weight

$$\tilde{a}_i = \frac{e^{\frac{q \cdot C_i^T}{\sqrt{d}}}}{\sum_{j=1..p} e^{\frac{q \cdot K_j^T}{\sqrt{d}}} + \sum_{j=1..m} (s_j \cdot e^{\frac{q \cdot C_j^T}{\sqrt{d}}})}, i = 1..m$$

p tokens  
in steady zone      all m clusters  
with  $s_j$  tokens each

$$token \text{ atten.} = \text{softmax} \left( \frac{qK^T}{\sqrt{d}} \right)_{token} v_{token}$$

*i*th cluster atten.

$$= \sum_{token} \text{softmax} \left( \frac{qK^T}{\sqrt{d}} \right)_{token} v_{token}$$

$$\approx \sum_{token} \tilde{a}_i v_{token} = \tilde{a}_i \sum_{token} v_{token}$$

$$= \tilde{a}_i \sum v$$



# Wave Index: Estimation Zone

## □ Accuracy-Bounded Attention Estimation

- ❖ For un-retrieved clusters, use centroids to generate estimate attention

i-th centroid weight  $\tilde{a}_i = \frac{e^{\frac{q \cdot C_i^T}{\sqrt{d}}}}{\sum_{j=1..p} e^{\frac{q \cdot K_j^T}{\sqrt{d}}} + \sum_{j=1..m} (s_j \cdot e^{\frac{q \cdot C_j^T}{\sqrt{d}}})}$ ,  $i = 1..m$

$p$  tokens in steady zone      all  $m$  clusters with  $s_j$  tokens each

$$\left[ \begin{array}{l} C_i = \frac{\sum_{j=1..s_i} K_j}{s_i} \\ e^{\frac{q \cdot C_i^T}{\sqrt{d}}} \leq \frac{\sum_{j=1..s_i} e^{\frac{q \cdot K_j^T}{\sqrt{d}}}}{s_i} \end{array} \right]$$

Jensen's inequality

The estimation is a lower bound



# Wave Index: Data Placement and Workflow

## □ Physical placement

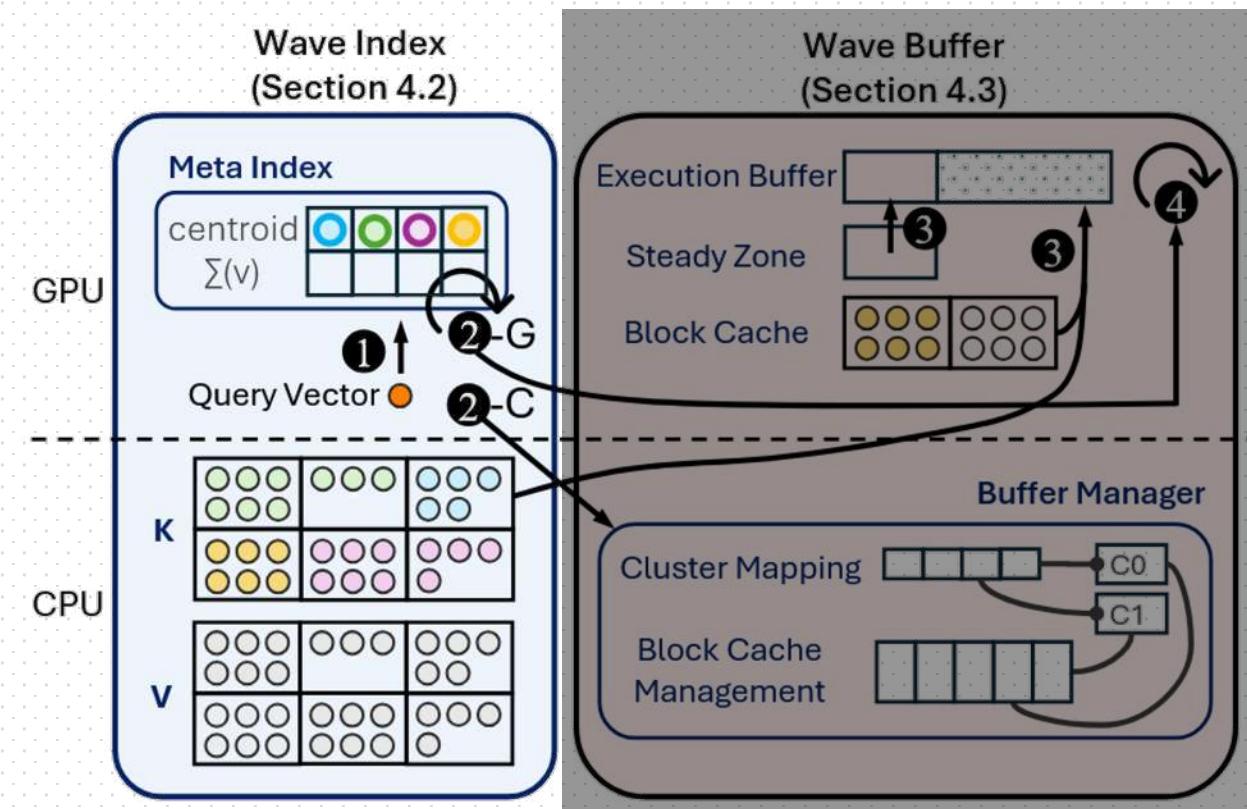
- ❖ CPU: KV vectors (stored in blocks)
- ❖ GPU: centroids +  $\sum v$

## □ After prefilling

- ❖ GPU performs segmented clustering
- ❖ KV vectors offload **asynchronously**
- ❖ CPU construct related data structure **asynchronously**

## □ Decoding

- ❖ GPU search centroids
- ❖ CPU send retrieved KV blocks



RetroInfer Overview

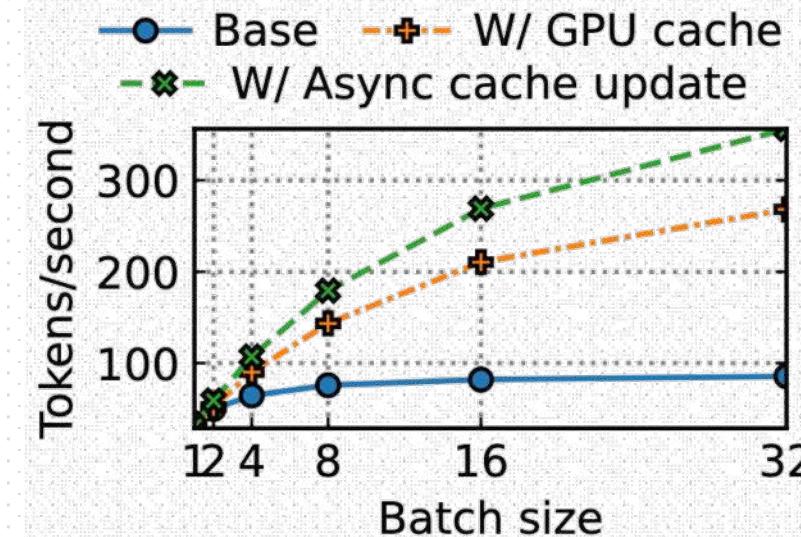


# Wave Buffer: Why and How Is it Possible

- ❑ Data traffic is still high
  - ❖ Though wave index reduce it to <2% of full attention

## ❑ A GPU buffer is designed

- ❖ Based on a phenomenon “*adjacent decoding steps for neighboring query vectors show strong temporal locality*”



Throughput with batch sizes  
Base is constrained by PCIe bandwidth



# Wave Buffer: Overview

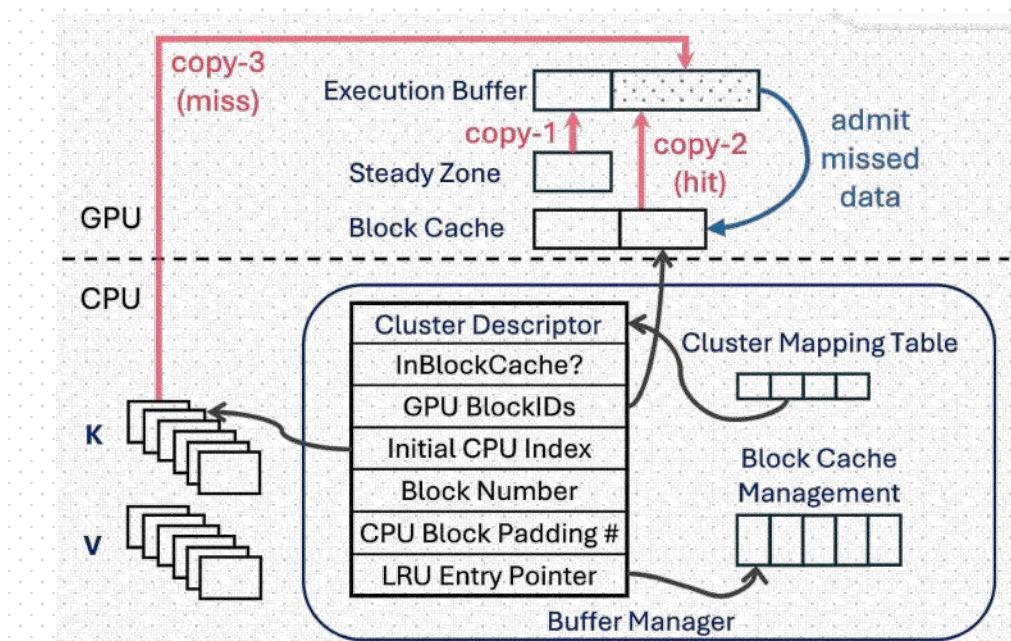
## □ Buffer contents on GPU

- ❖ Execution buffer
- ❖ Steady zone
- ❖ Block cache

## □ During decoding: 3 Data copy

## □ Control plane on CPU

- ❖ A table containing cluster-id to infos
- ❖ Cache policy: just LRU





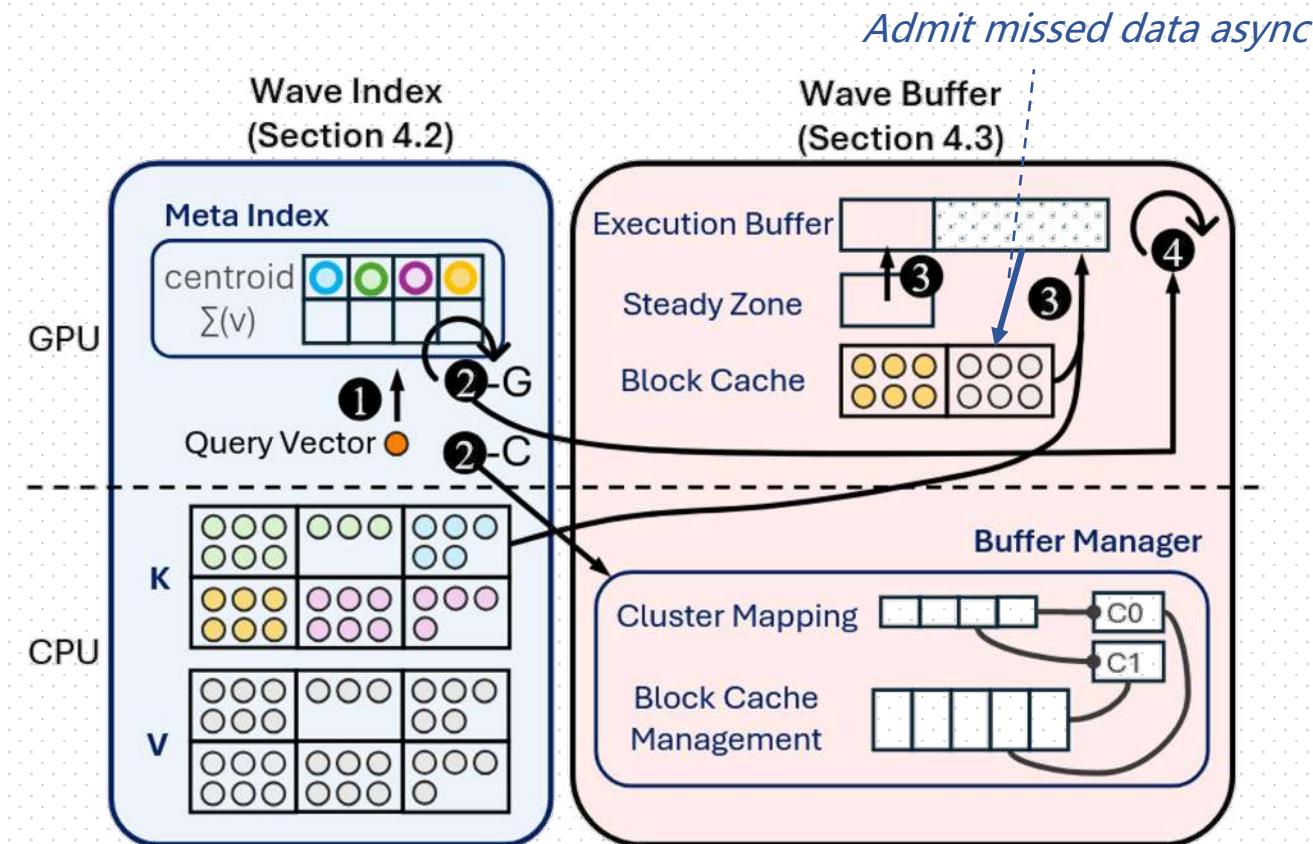
# Coordination between CPU/GPU

## □ Synchronous cache access

- ❖ The KV vectors are ready as soon as possible

## □ Asynchronous cache update

- ❖ Cache update
- ❖ Admit missed data
- ❖ Removed from critical path



RetroInfer Overview



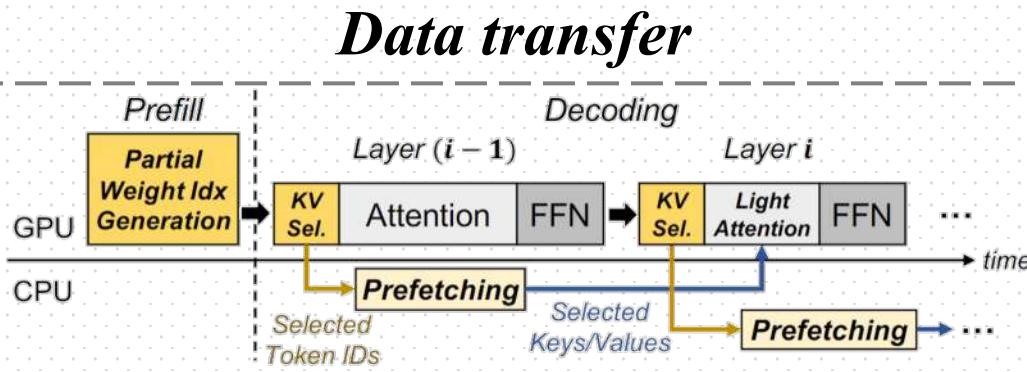
# Contexts

- Background
- Motivation
- RetroInfer
- Evaluations



# Baselines

InfiniGen  
OSDI 24



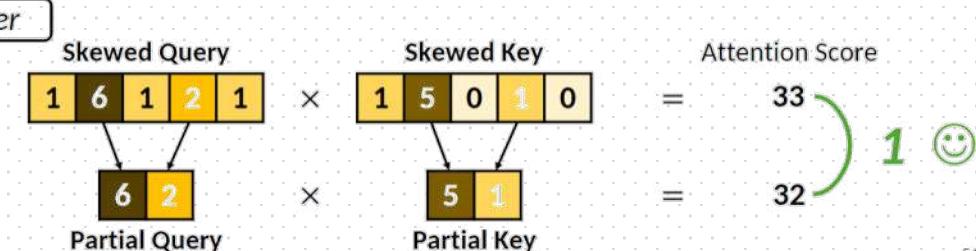
Prefetching, but frequent data transfer

Quest  
ICLR 24

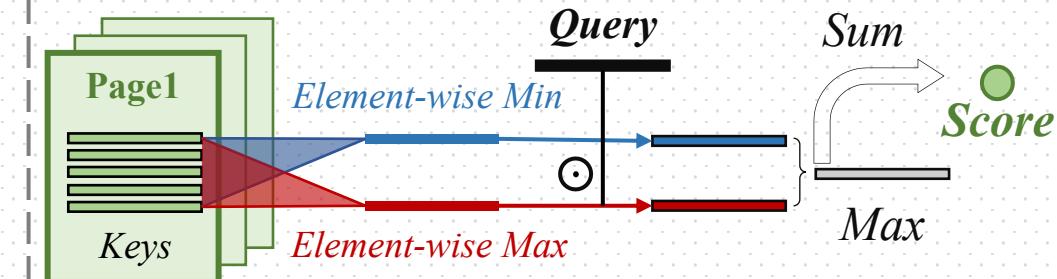


Quest-Offload is implemented  
in this paper

## TopK Selection



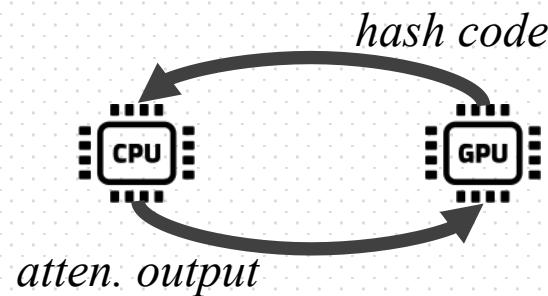
Use reduced dimension to find topk, still costly



Use page level reduced keys, accuracy loss

MagicPIG  
ICLR 25

Attention is  
done on CPU



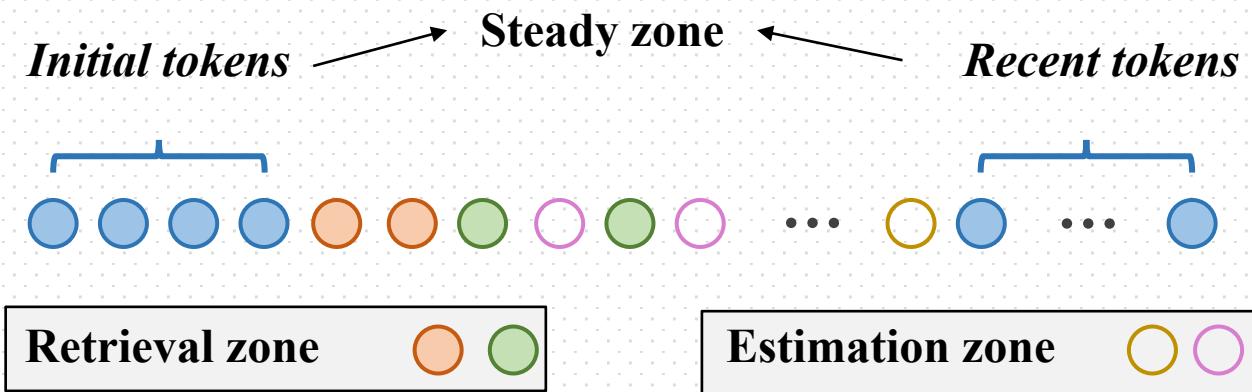
Use Locality Sensitive Hashing (LSH)  
to sample important Key/Value,  
accuracy loss



# Parameters

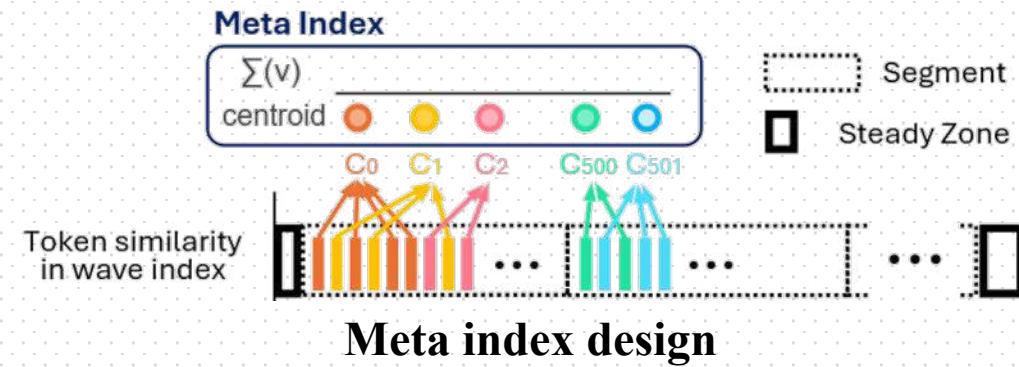
## ❑ Zones setup

- ❖ Steady zone: 4 + 64
- ❖ Retrieval zone: 1.8%
- ❖ Estimation zone: 23%



## ❑ Segment clustering

- ❖ 8K tokens in a segment
- ❖ #clusters = 1/16 #tokens





# Parameters

## ❑ Cache Setup

- ❖ 5% of all KV vectors
- ❖ Page size: 2KB

## ❑ Baseline setup

- ❖ Follow the paper statements



# Other Configuration

## ❑ Hardware

- ❖ 1 × NVIDIA A100 GPU (80GB memory)
- ❖ AMD EPYC 7V12 64-core CPU (1900GB memory)
- ❖ PCIe 4.0 × 16 (~31GB/s)

## ❑ Models

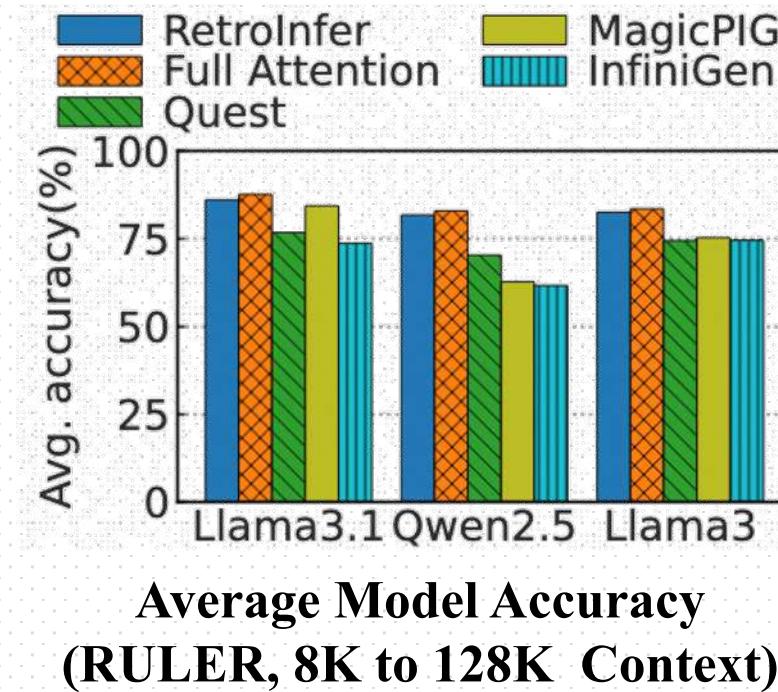
- ❖ Llama3.1-8B      ]    128K context
- ❖ Qwen2.5-7B      ]
- ❖ Llama38B-1048K



# Accuracy

❑ RetroInfer always approaches full attention accuracy

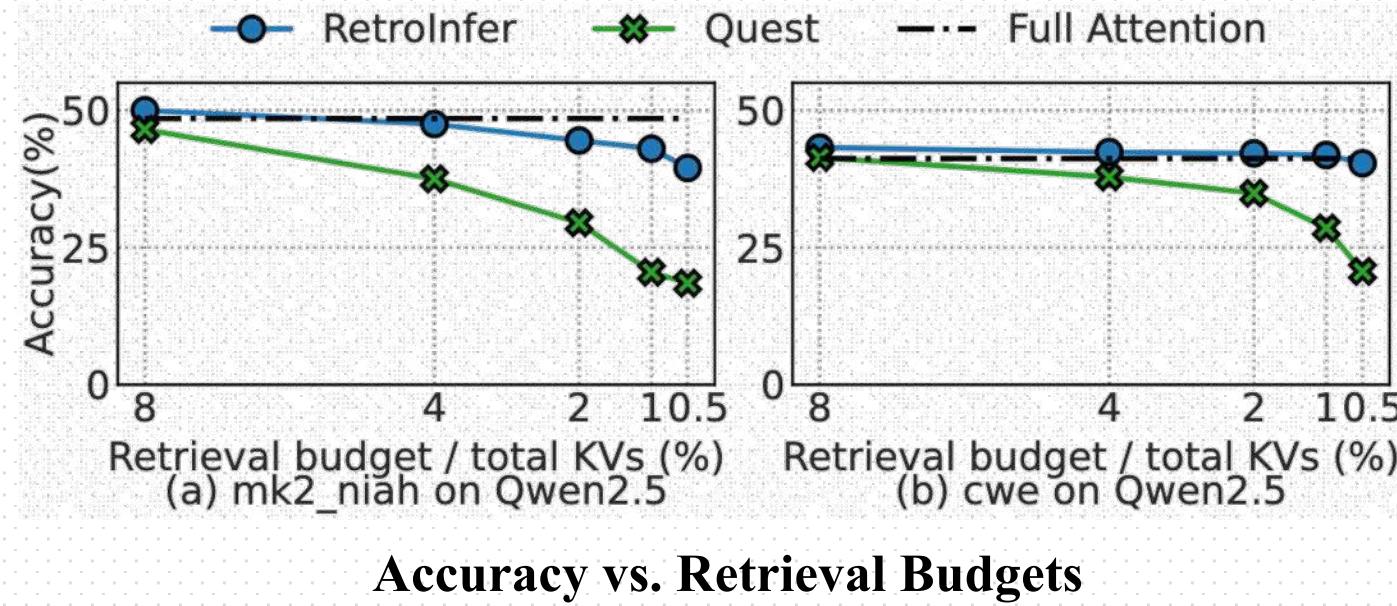
- ❖ Drop **0.73%/0.78%/1.46%** compared to full attention
- ❖ Achieves **5.48%/15.51%/3.33%** accuracy improvements compared to the best-performing baselines on each model





# Accuracy

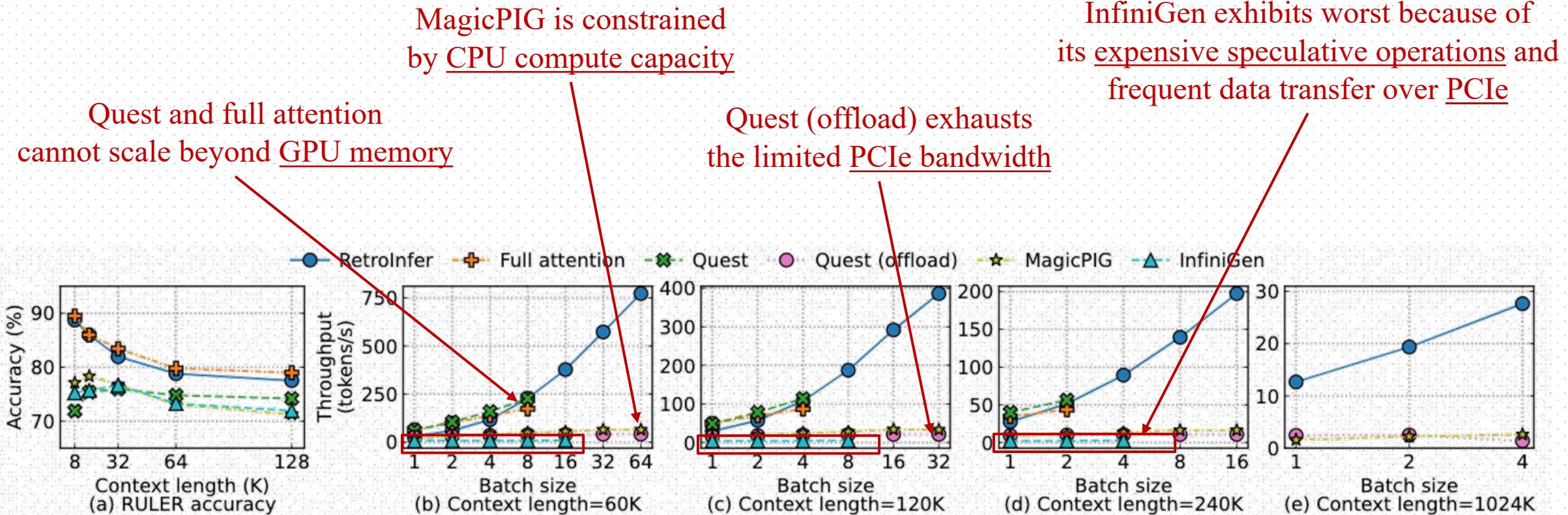
- ❑ RetroInfer maintains accuracy close to full attention, even with a low retrieval budget
- ❑ RetroInfer sometimes outperforms full attention





# Efficiency

## ❑ RetroInfer scales well and achieves the best throughput

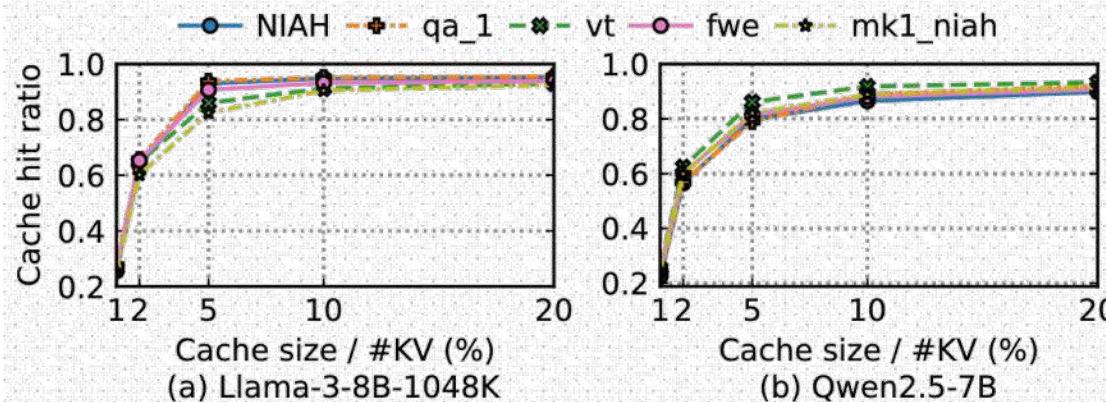


**Accuracy and Decoding throughput vs. Context Length and Batch Size**  
**Benchmark: RULER, model: Llama3-8B-1048K**

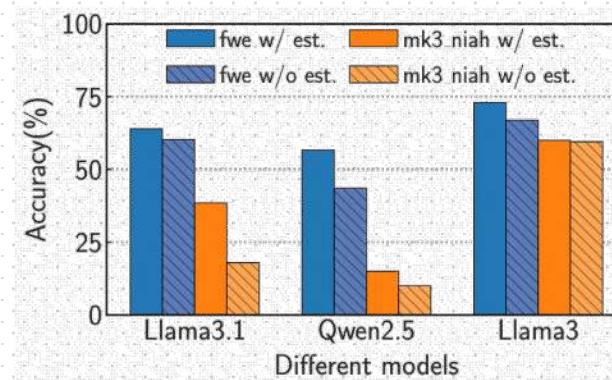


# Micro Analysis

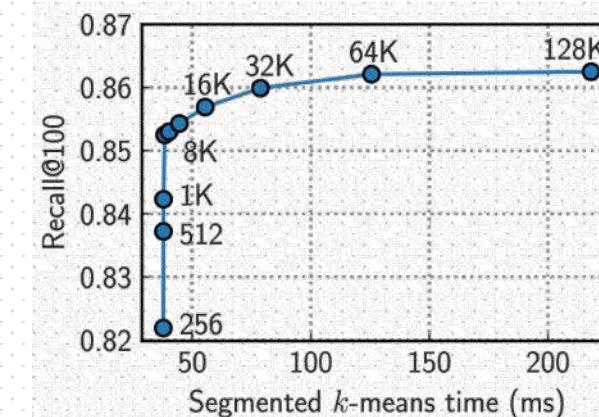
## □ Cache ratio: 5% is good



## □ Estimation helps a lot



## □ Segment size: 8K is good



## □ Decoding Latency Breakdown (ms)

Batch size	1	32
Access meta index	0.1447	0.5132
Estimation zone attention	0.1010	0.5398
Access cluster mapping table	0.2405	0.5902
Update cluster mapping table	0.0142	0.5176
Data copy to execution buffer	0.1281	0.9008
Steady and retrieval zone attention	0.0777	0.2950
Admit missed data	0.0572	0.1038



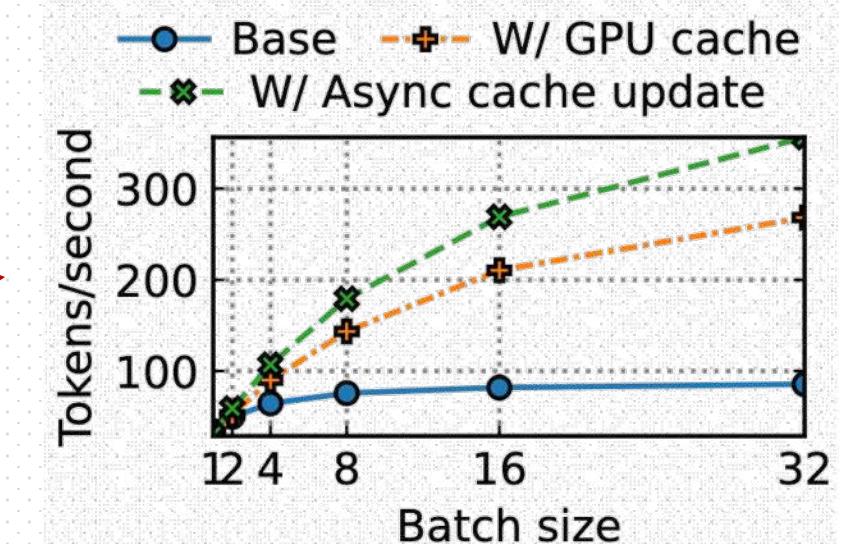
# Conclusions

## □ Pros

- ❖ Innovative technique: estimation via centroids and  $\sum v$
- ❖ High completeness of system design
- ❖ Excellent experimental results
- ❖ Good story: vector storage system

## □ Cons

- ❖ Buffer contributes a lot... 
- ❖ Disparity between motivation and solution
- ❖ Unclear explanation or lack of deep analysis
  - Impact of Zone ratio (retrieval / estimation)
  - ANNS OOD issue (see next page)



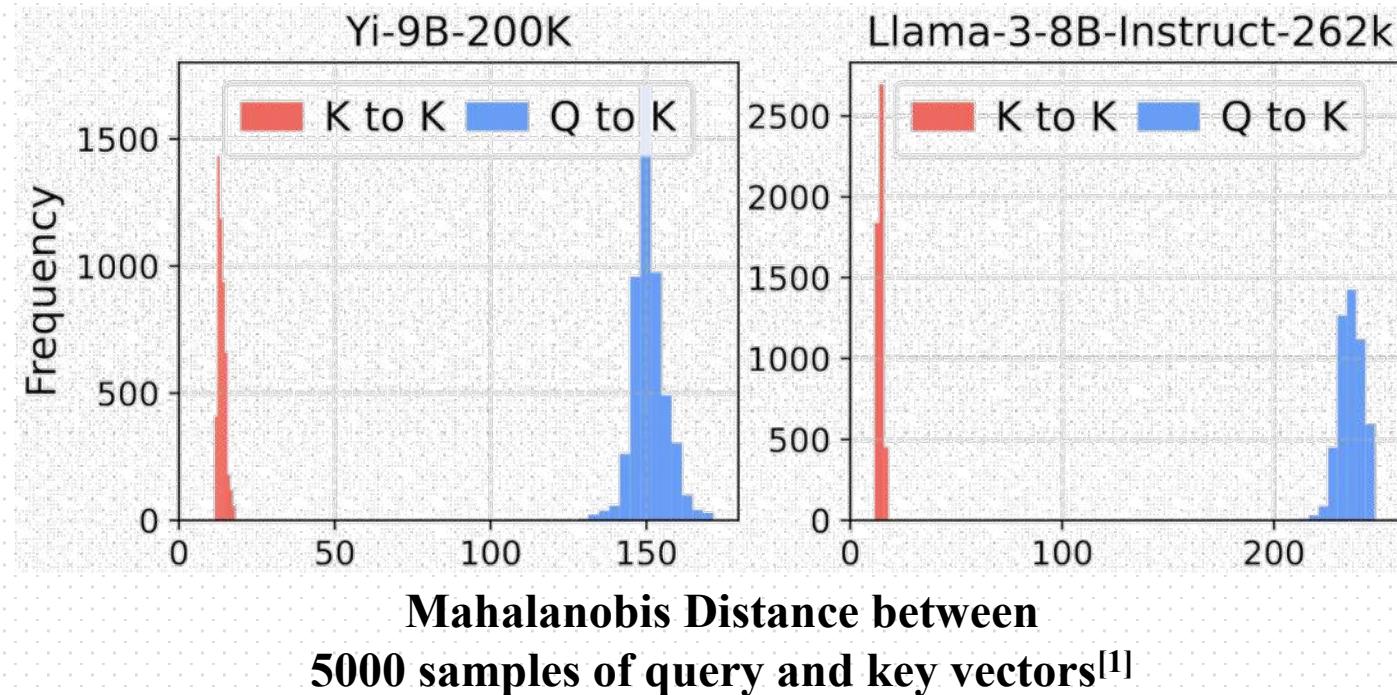


# OOD Issue: When ANNS Meets KV Cache

## ❑ OOD (out-of-distribution)

- ❖ The distribution of queries differs from that of the base data

## ❑ Search in the KV cache exhibits strong OOD characteristics

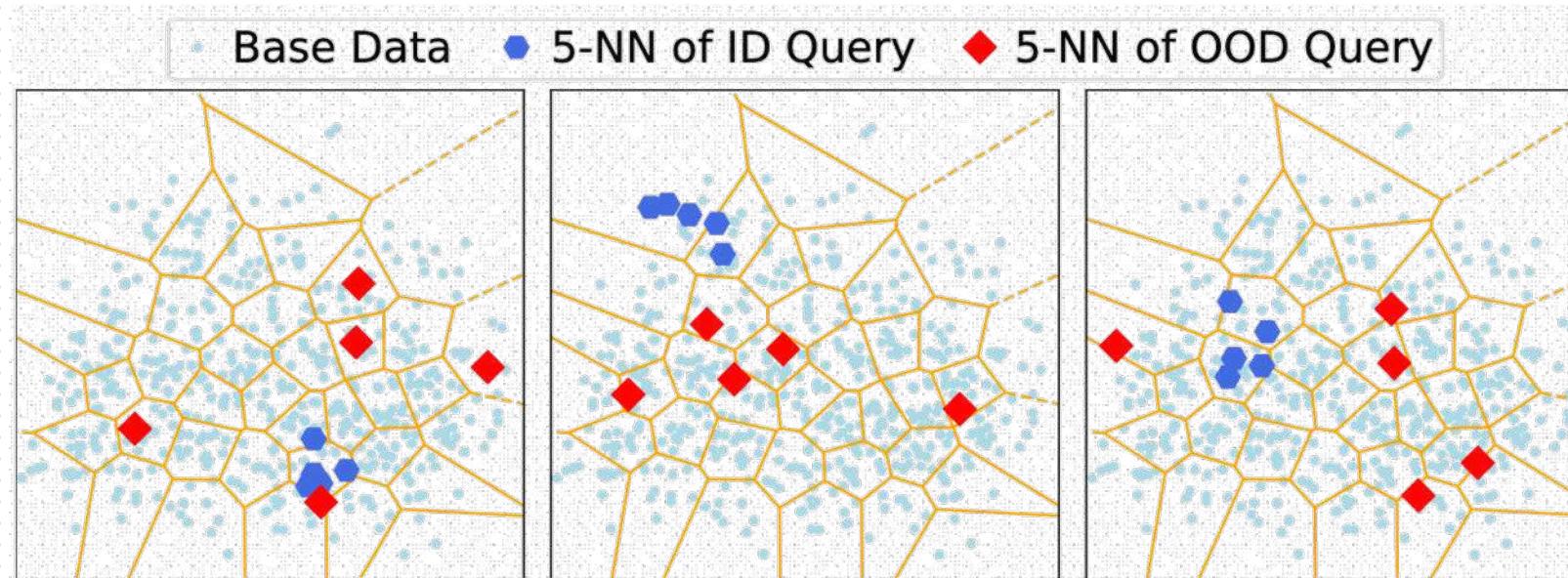


[1] RetrievalAttention: Accelerating Long-Context LLM Inference via Vector Retrieval



# OOD Increases the Searching Difficulty

- ❖ Most ANNS indexes are built for in-distribution search
- ❖ OOD queries challenge the assumptions of standard ANNS methods.



**5 Nearest neighbors' distribution in cluster-based index  
for ID query and OOD query<sup>[1]</sup>**

[1] RoarGraph: A Projected Bipartite Graph for Efficient Cross-Modal Approximate Nearest Neighbor Search 47



**End**

Thank you!

Q&A