# Accelerating Retrieval-Augmented Generation

**Derrick Quinn**, Mohammad Nouri, Neel Patel,
Alireza Salimi[1], Sukhan Lee[2], Hamed Zamani[1], and Mohammad Alian

**Cornell, UMass Amherst[1], Samsung Electronics[2]**

**CornellEngineering**
Electrical and Computer Engineering

CSL

# Retrieval Augmented Generation (RAG)

- Majority of LLM-powered applications use RAG

- Access to fresh data is a must have!
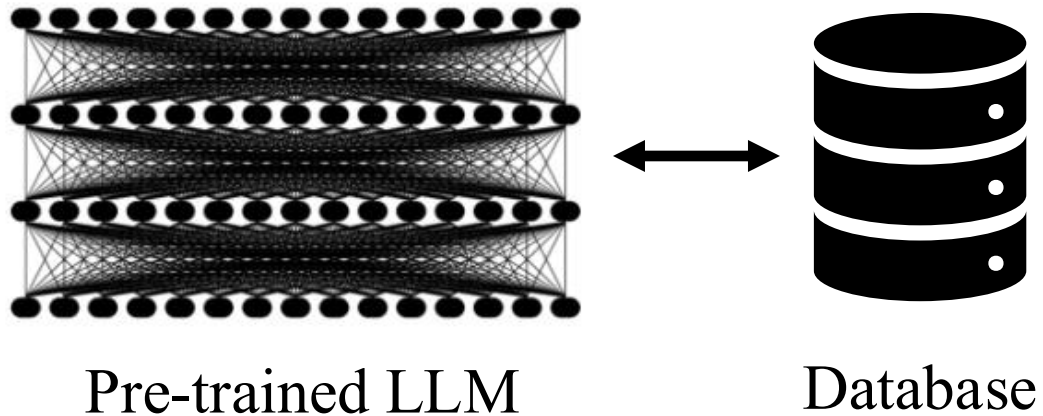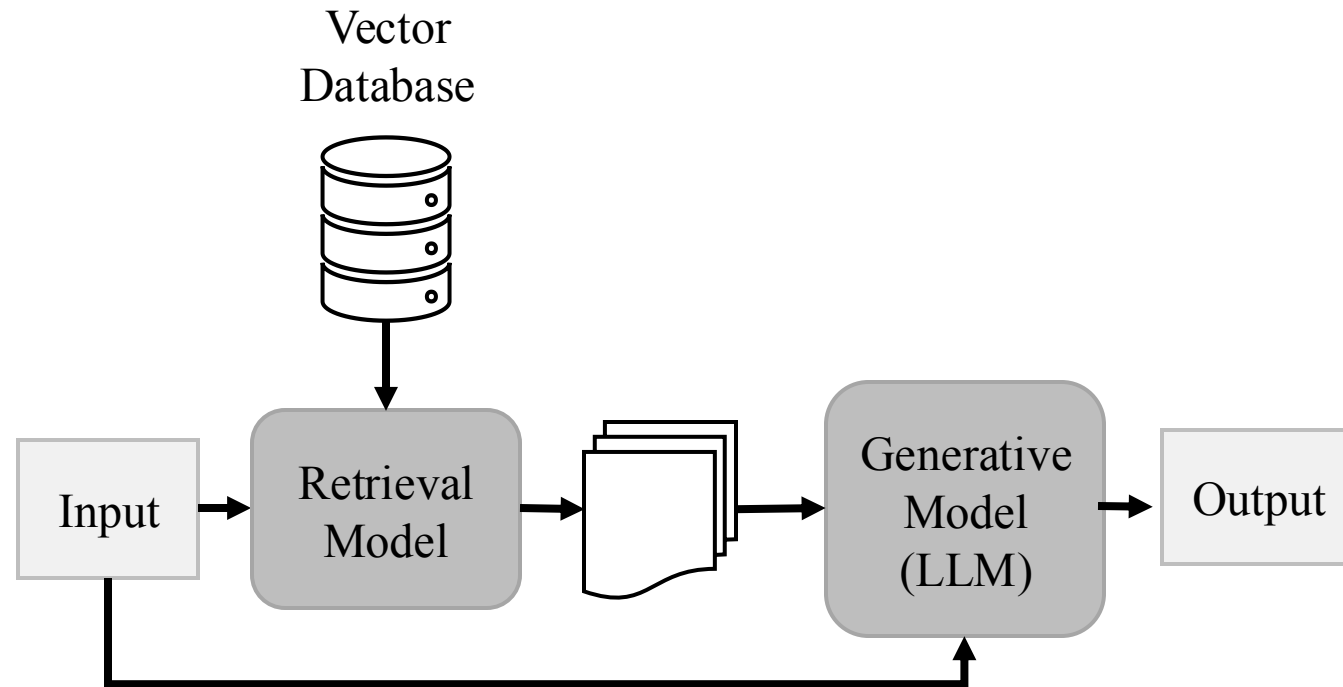  - Updated knowledge, factual grounding

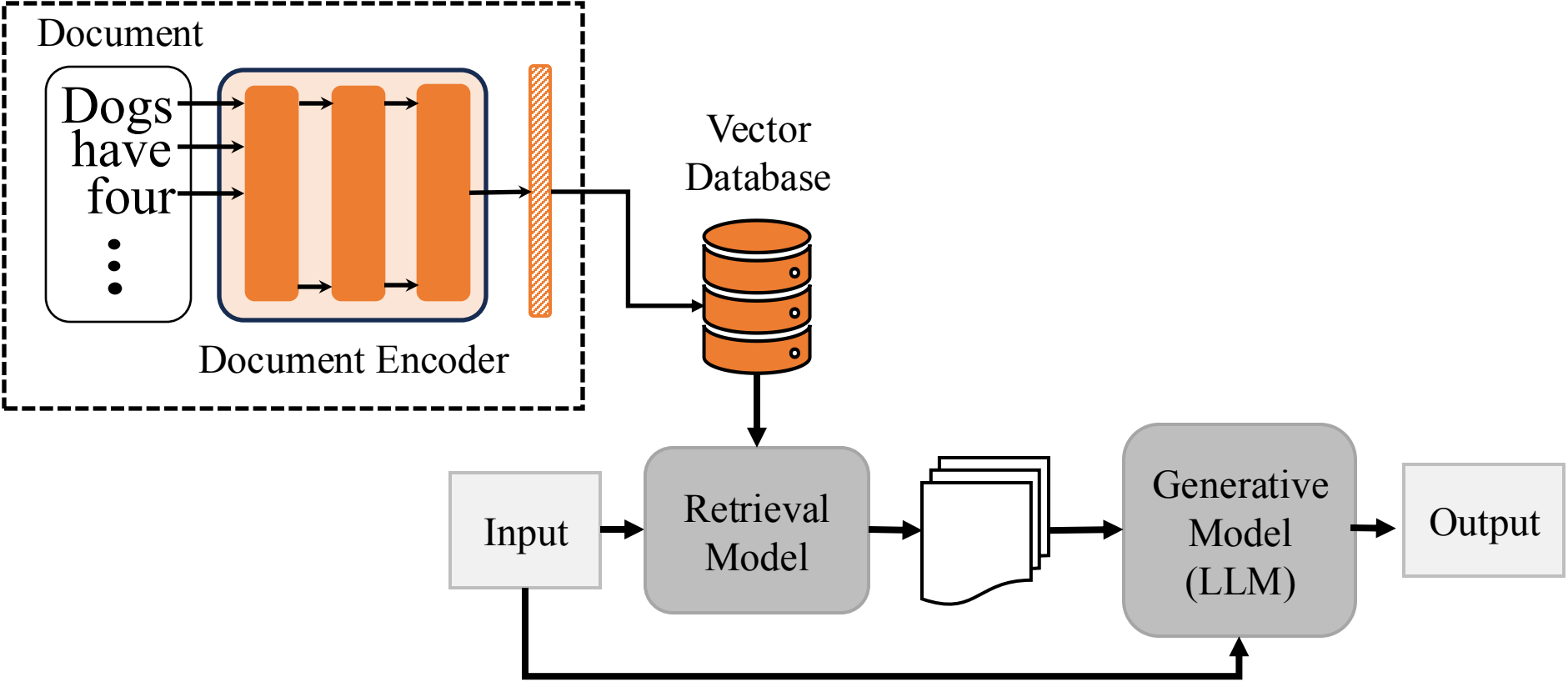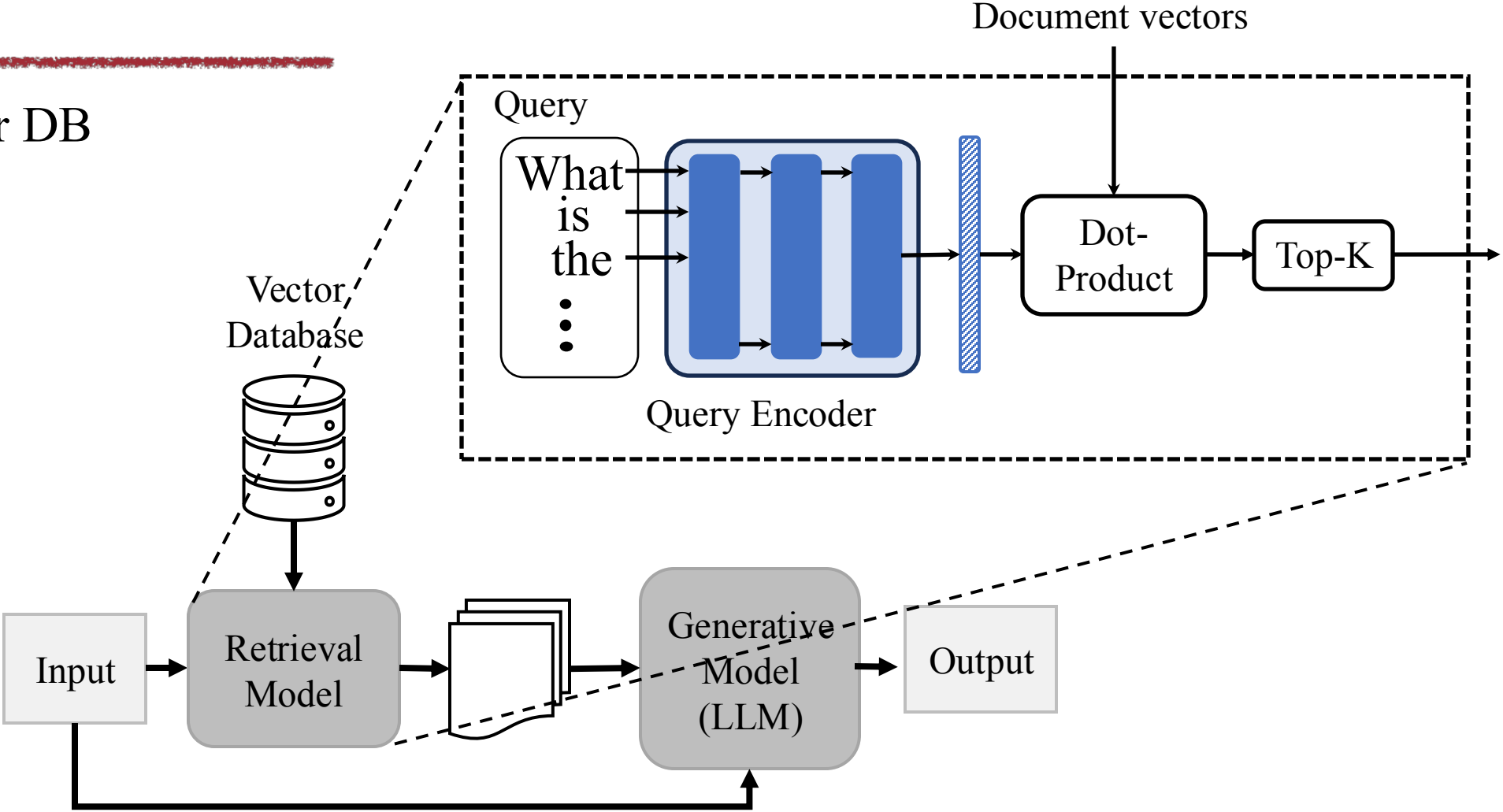Pre-trained LLM            Database

Image generated by ChatGPT!

# RAG Pipeline Overview

# RAG Pipeline Overview

Offline: Populate Vector DB

# RAG Pipeline Overview

Online: Search Vector DB
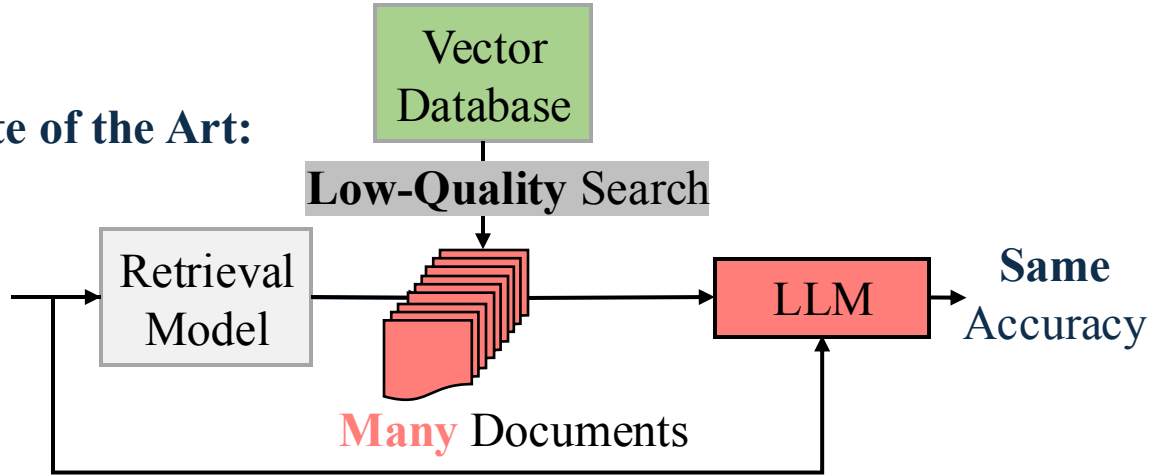
Document vectors
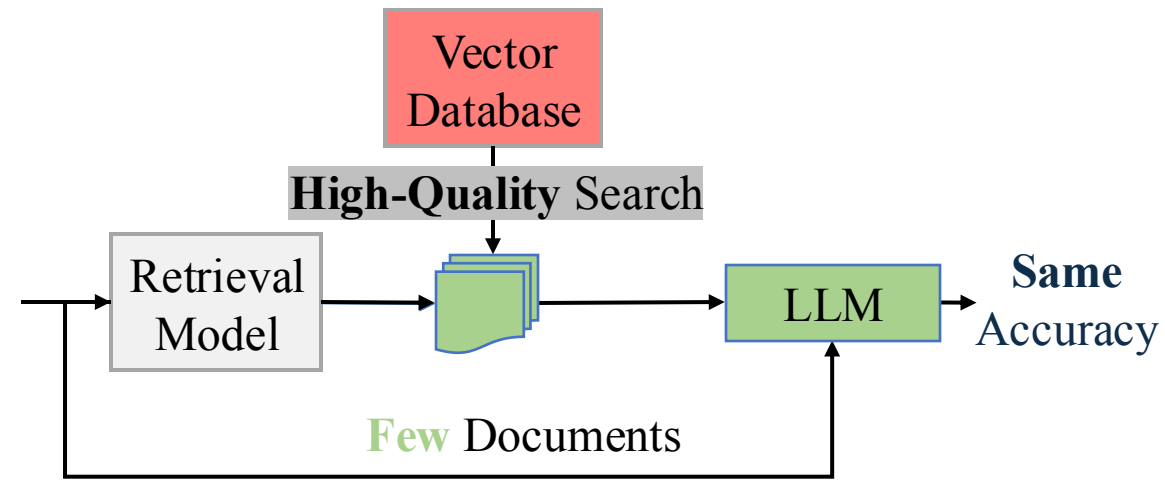
Query

What
is
the
⋮

Query Encoder

Dot-Product

Top-K

Vector Database

Input → Retrieval Model → Generative Model (LLM) → Output

# Insight: Interplay of Retrieval and Generation



|  | State of the Art | This Work |
|---|---|---|
| Accuracy | +++ | +++ |
| Retrieval Speed | +++ | ++ |
| Generation Speed | -- | + |
| Overall Speed | + | +++ |

**State of the Art:**

Vector Database → **Low-Quality** Search → Retrieval Model → **Many** Documents → LLM → **Same** Accuracy

**This Work (Near-Memory Acceleration of Exact Search)**

Vector Database → **High-Quality** Search → Retrieval Model → **Few** Documents → LLM → **Same** Accuracy

# Insight: Interplay of Retrieval and Generation



Vector
Database

Our contributions:

1. We showcase the system-level interactions of retrieval quality

2. We leverage algorithm-hardware co-design to build a high-quality retrieval accelerator, accelerating RAG

3. We showcase the utility of CXL for the interface of high-capacity accelerators

Accuracy

Retrieval
Speed

Generation
Speed

Overall
Speed

**Same**
Accuracy

**Same**
Accuracy

**Few** Documents

**This Work (Near-Memory Acceleration of Exact Search)**

# Content

- Background
- Profiling RAG Applications
  - A Case for Near-Memory Exhaustive Search Acceleration
- Introducing Intelligent Knowledge Store (IKS)
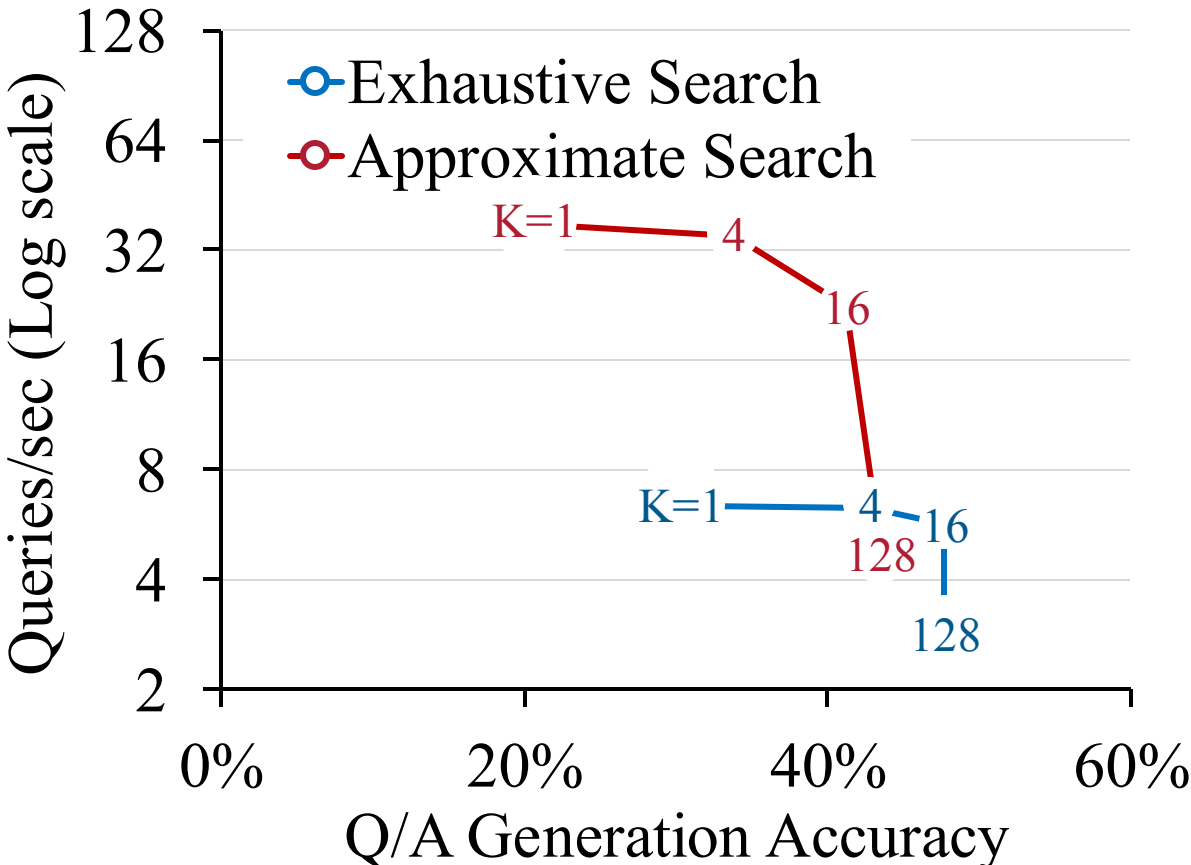- Evaluating IKS in a RAG pipeline

# Profiling RAG Applications

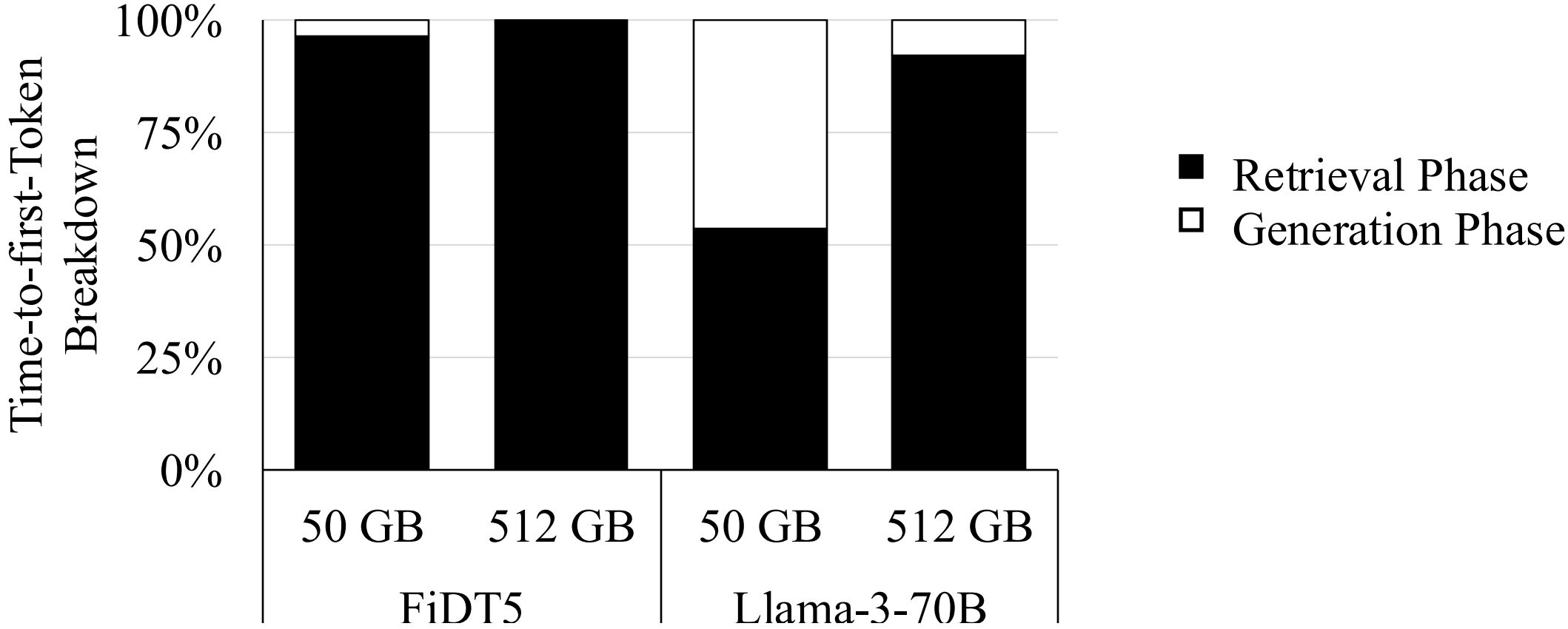# Evaluating RAG Applications: Methodology

- Question answering RAG applications
  - LLM: FiDT5 (T5 Fusion in Decoder), Llama-3-8B, Llama-3-70B
  - Retrieval model: BERT uncased to generate embeddings
  - 50-512GB vector database
- Meta FAISS library for similarity search
  - **Approximate** Nearest Neighbor **Search**: HNSW
  - **Exact/Exhaustive** Nearest Neighbor **Search**: Search the entire vector DB
- Experimental setup
  - Retrieval on CPU (Intel Xeon 4[th] gen 4416+)
  - Generation on GPU (H100 80GB)
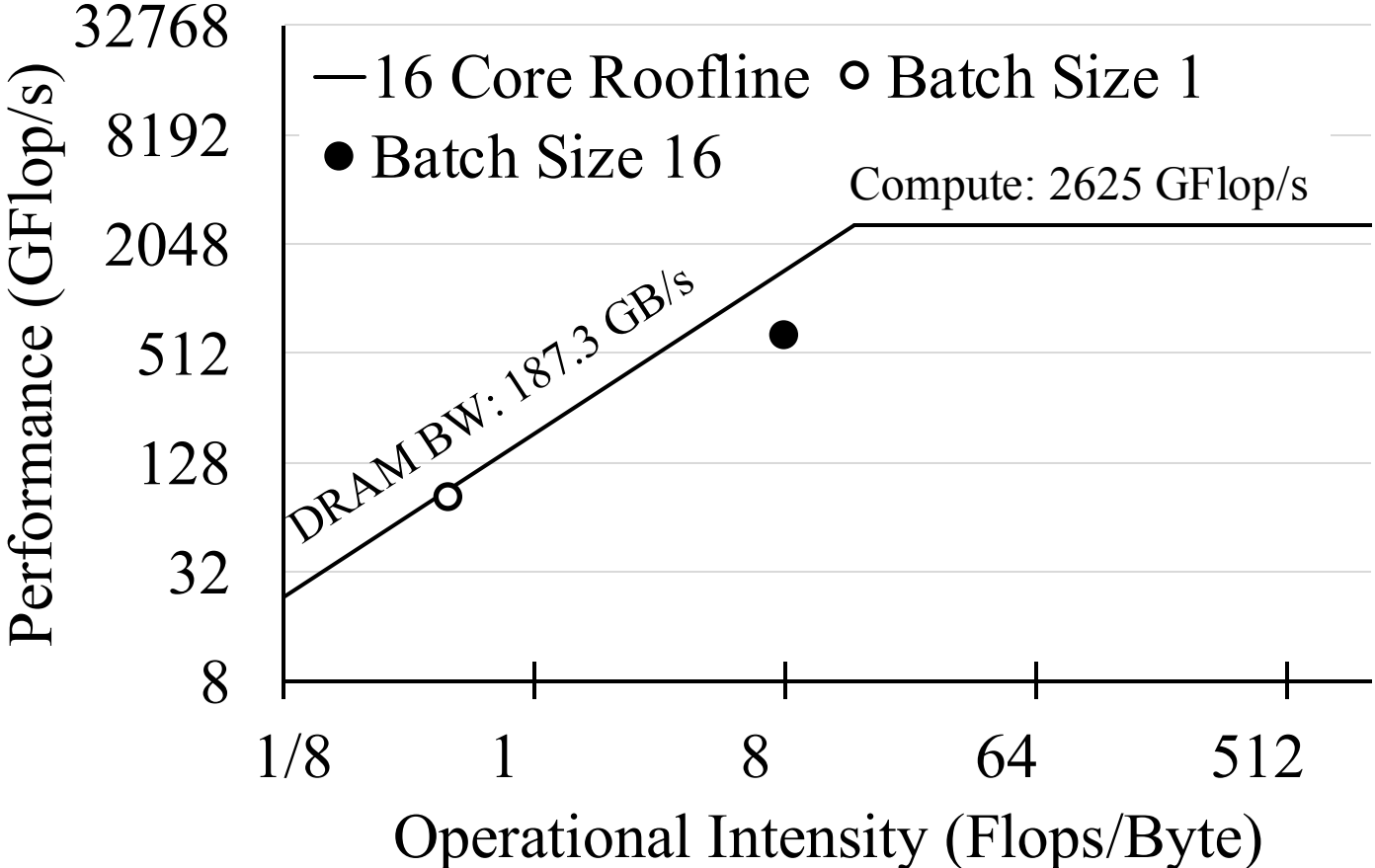  - Cycle-approximate simulator for IKS (Ours)

# Observation 1: Trade-off Space

# Observation 2: Exhaustive Search Bottlenecks RAG

# Observation 3: Exhaustive Search is Memory-Bound

# Case for Near-Memory Exhaustive Search Acceleration

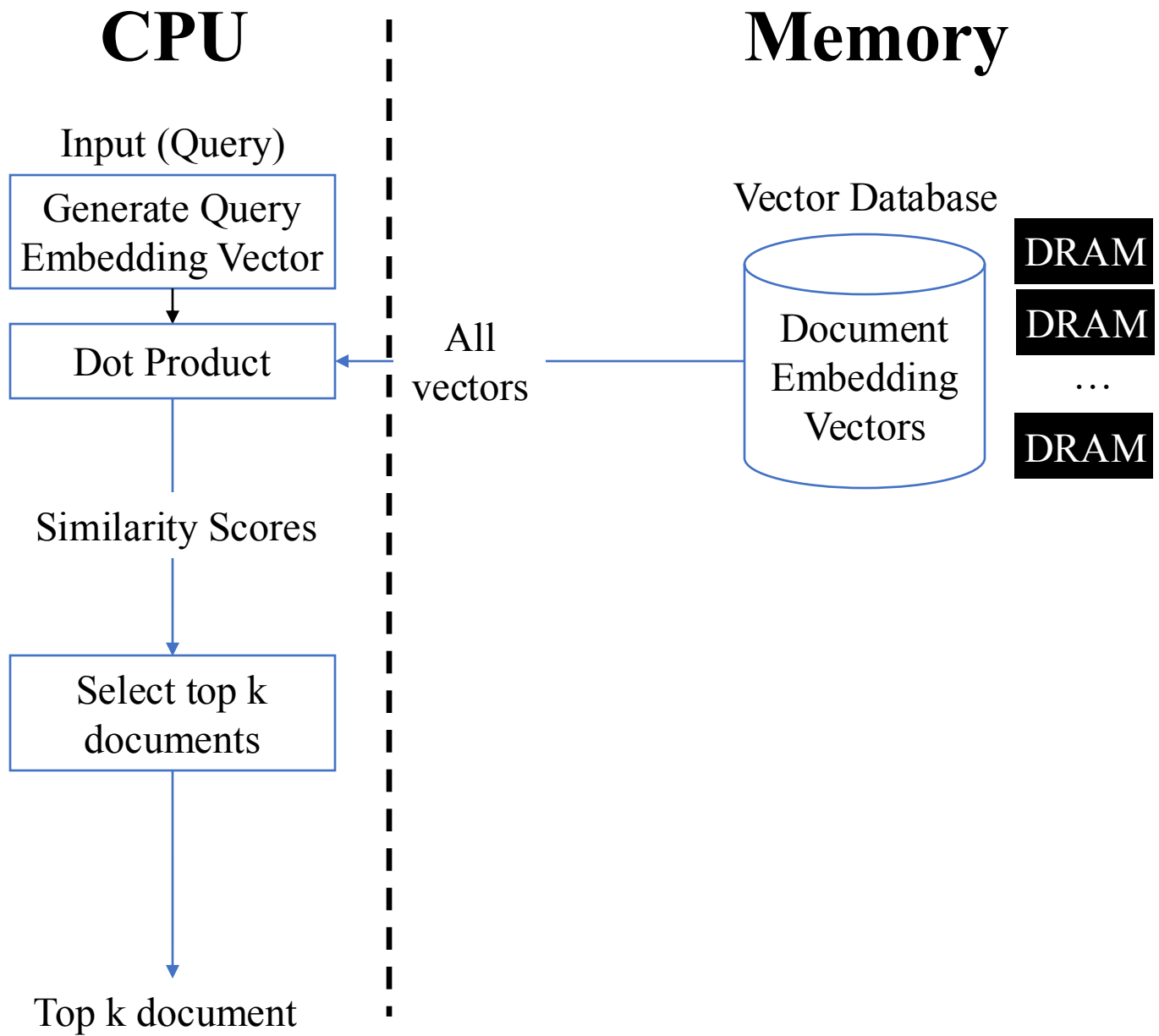# Near-Memory Exhaustive Search Acceleration is Justified

Observations:

1. Exhaustive search provides useful accuracy for RAG

2. Exhaustive search can dominate time-to-first-token

3. Exhaustive search is memory-bound

- Therefore, we apply algorithm-hardware co-design:

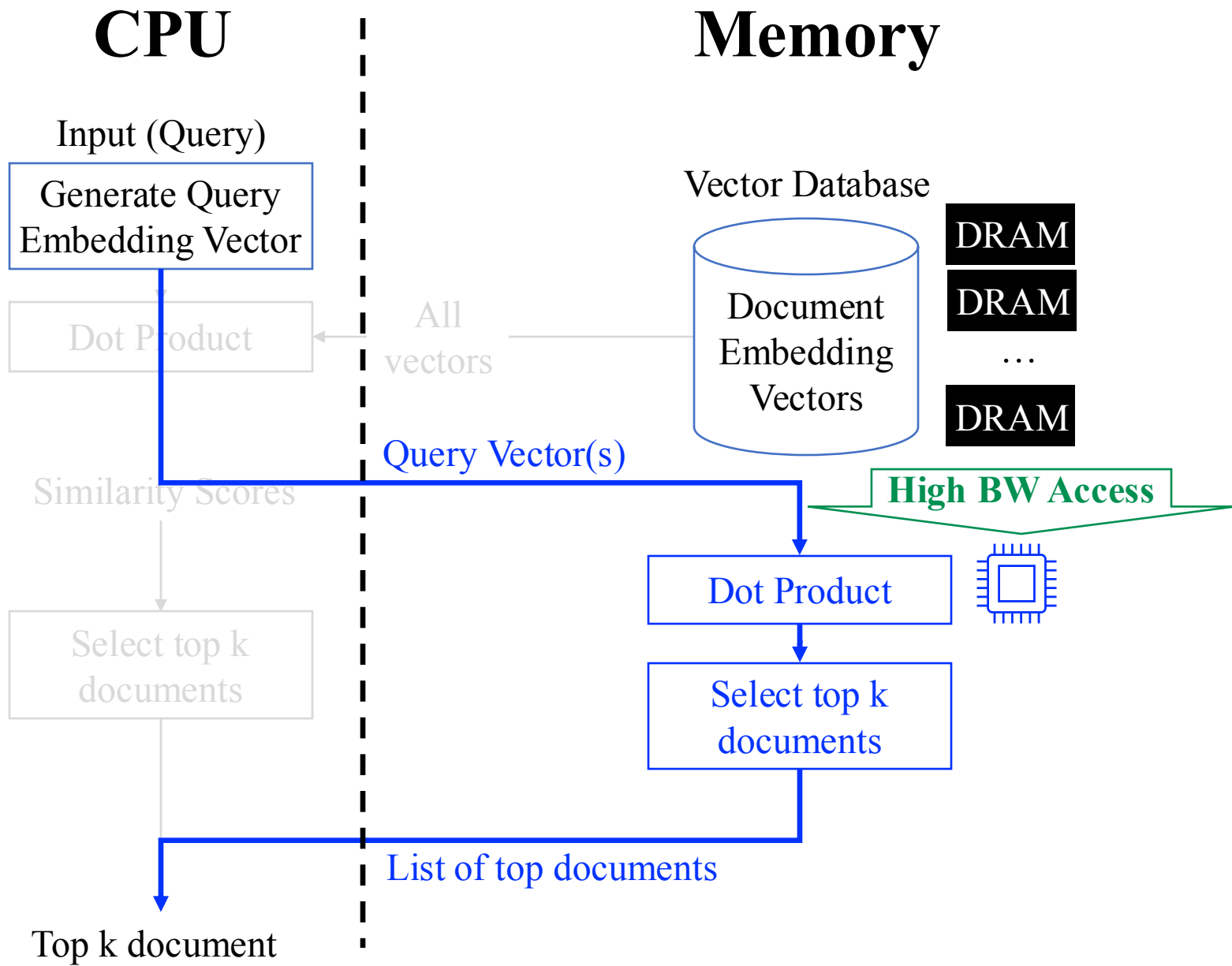**RAG Motivates Near-Memory Acceleration of Exhaustive Search**

# Intelligent Knowledge Store (IKS)

**CPU**

**Memory**

Baseline:
CPU Retrieval

Input (Query)

Generate Query Embedding Vector

Dot Product

Similarity Scores

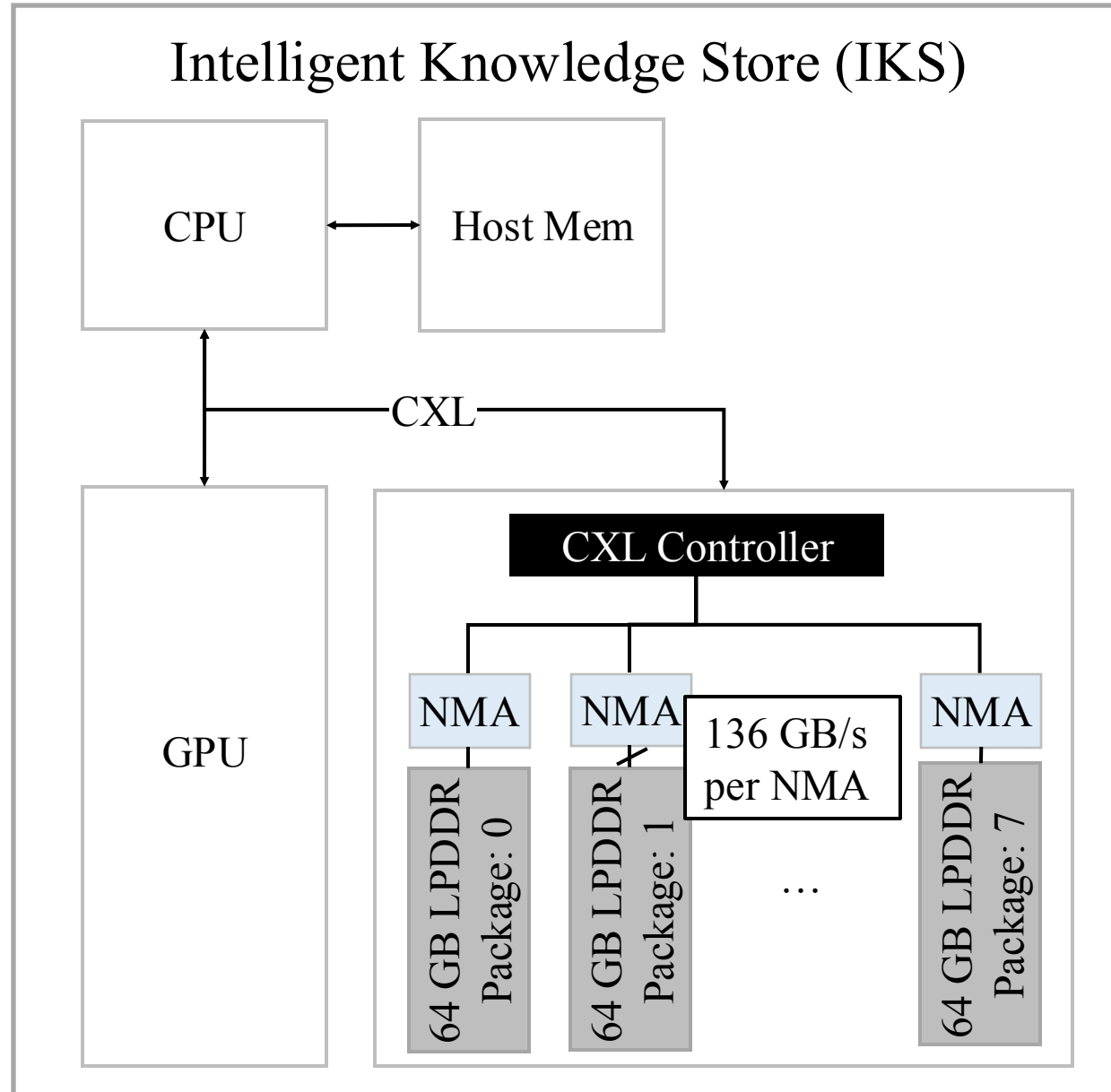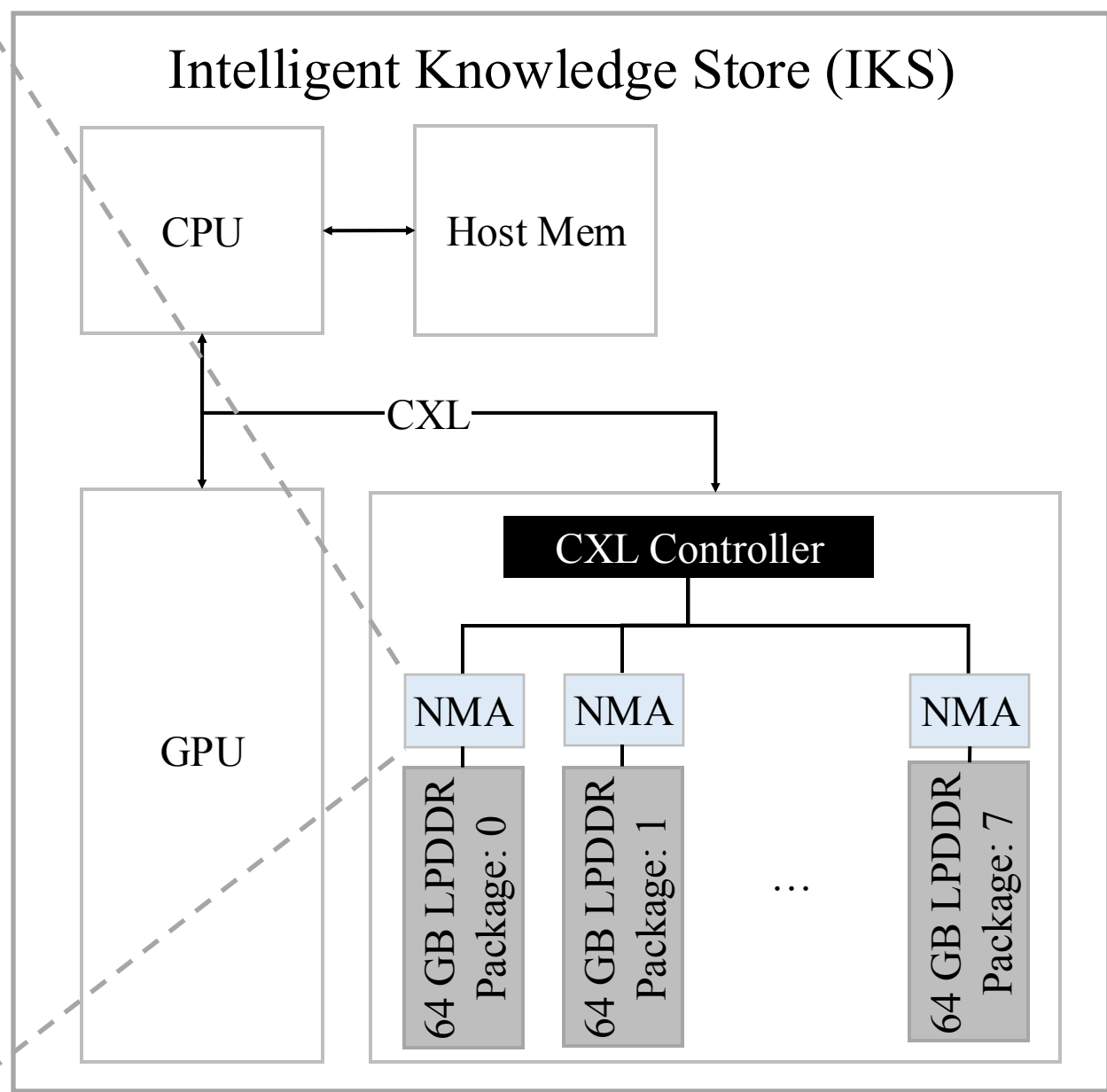Select top k documents

Top k document

All vectors

Vector Database

Document Embedding Vectors

DRAM

DRAM

…

DRAM

This work: CPU + IKS Retrieval

**CPU**

Input (Query)

Generate Query Embedding Vector

Dot Product

Similarity Scores

Select top k documents

Top k document

**Memory**

Vector Database

Document Embedding Vectors

DRAM
DRAM
…
DRAM

High BW Access

Query Vector(s)

Dot Product

Select top k documents

List of top documents

# IKS Overview

- IKS provisioned with:
  - 8x Near-Memory Accelerators (NMA)
  - 512 GB total capacity
  - 1.1 TB/sec internal bandwidth
- IKS supports:
  - Usermode polling
  - Multi-tenancy with host applications
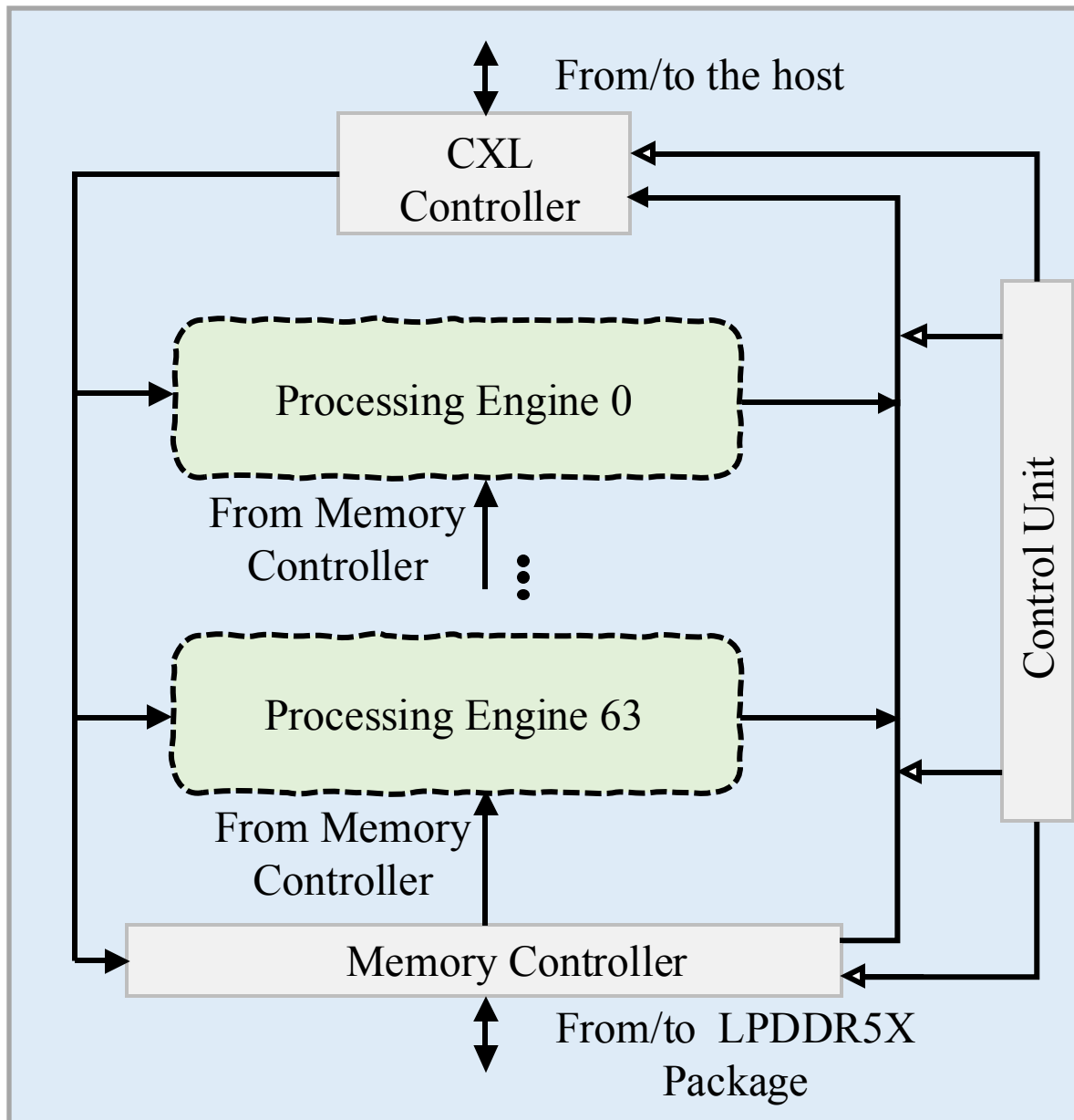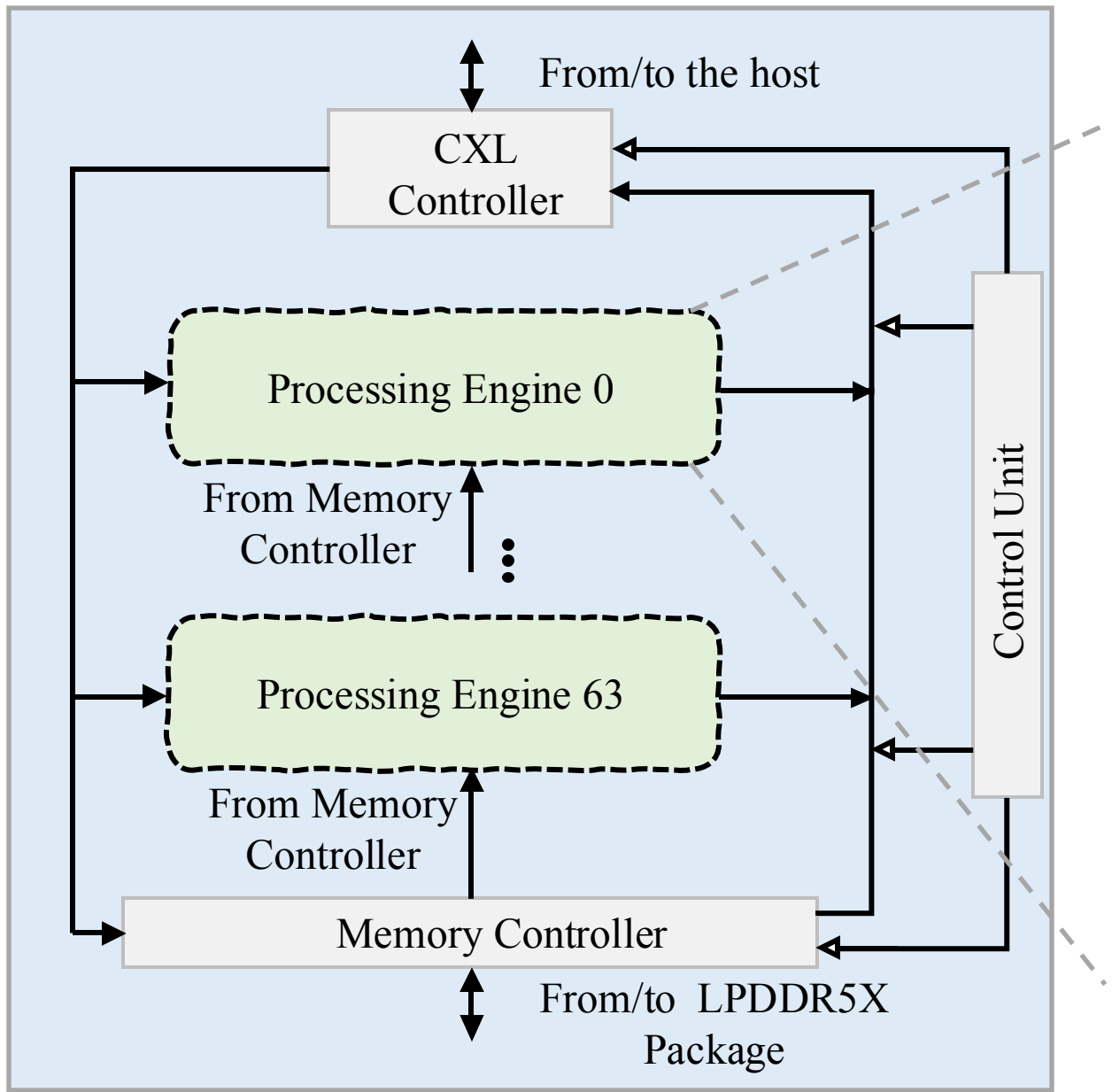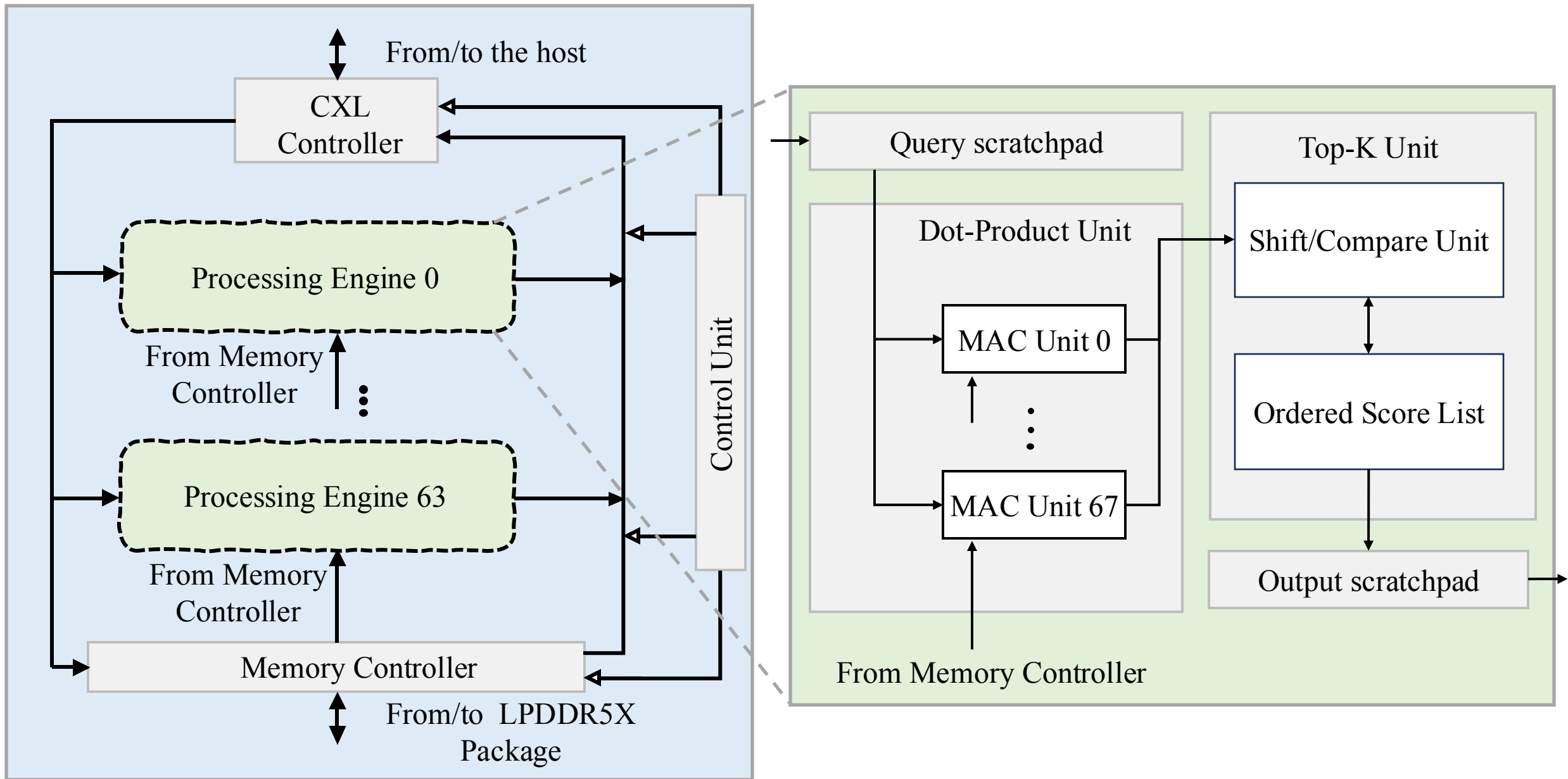  - Batch Size up to 64



Intelligent Knowledge Store (IKS)

CPU ↔ Host Mem

CXL

GPU

CXL Controller

NMA  NMA  136 GB/s per NMA  NMA

64 GB LPDDR Package: 0  64 GB LPDDR Package: 1  …  64 GB LPDDR Package: 7
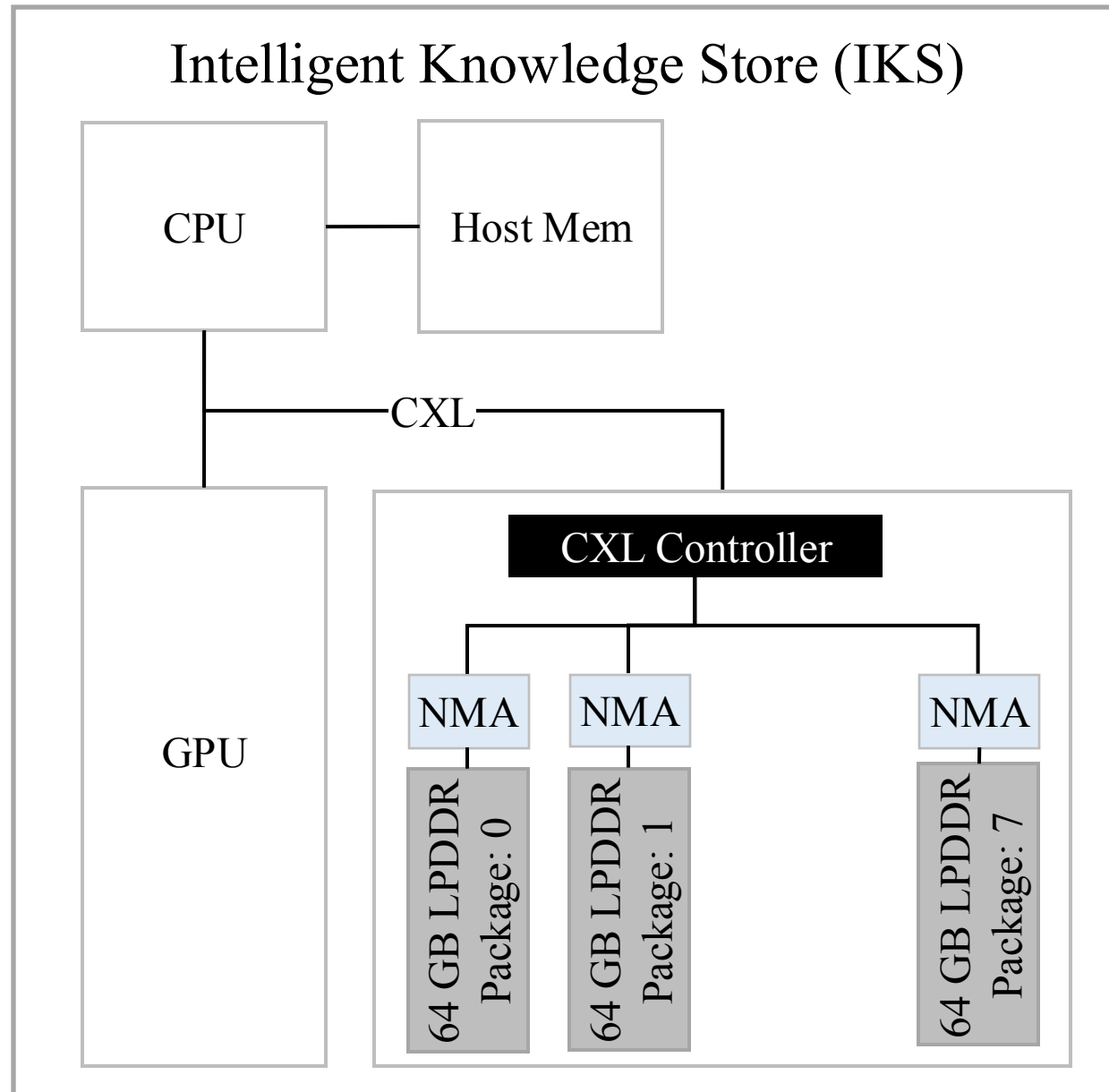
Intelligent Knowledge Store (IKS)

# Leveraging CXL to implement IKS

- IKS *collaborates* with CPU!

- Option 1: DMA
  - High initial overhead

- Option 2: PIO (CXL.io)
  - Low bandwidth

- Our approach: Use CXL.cache
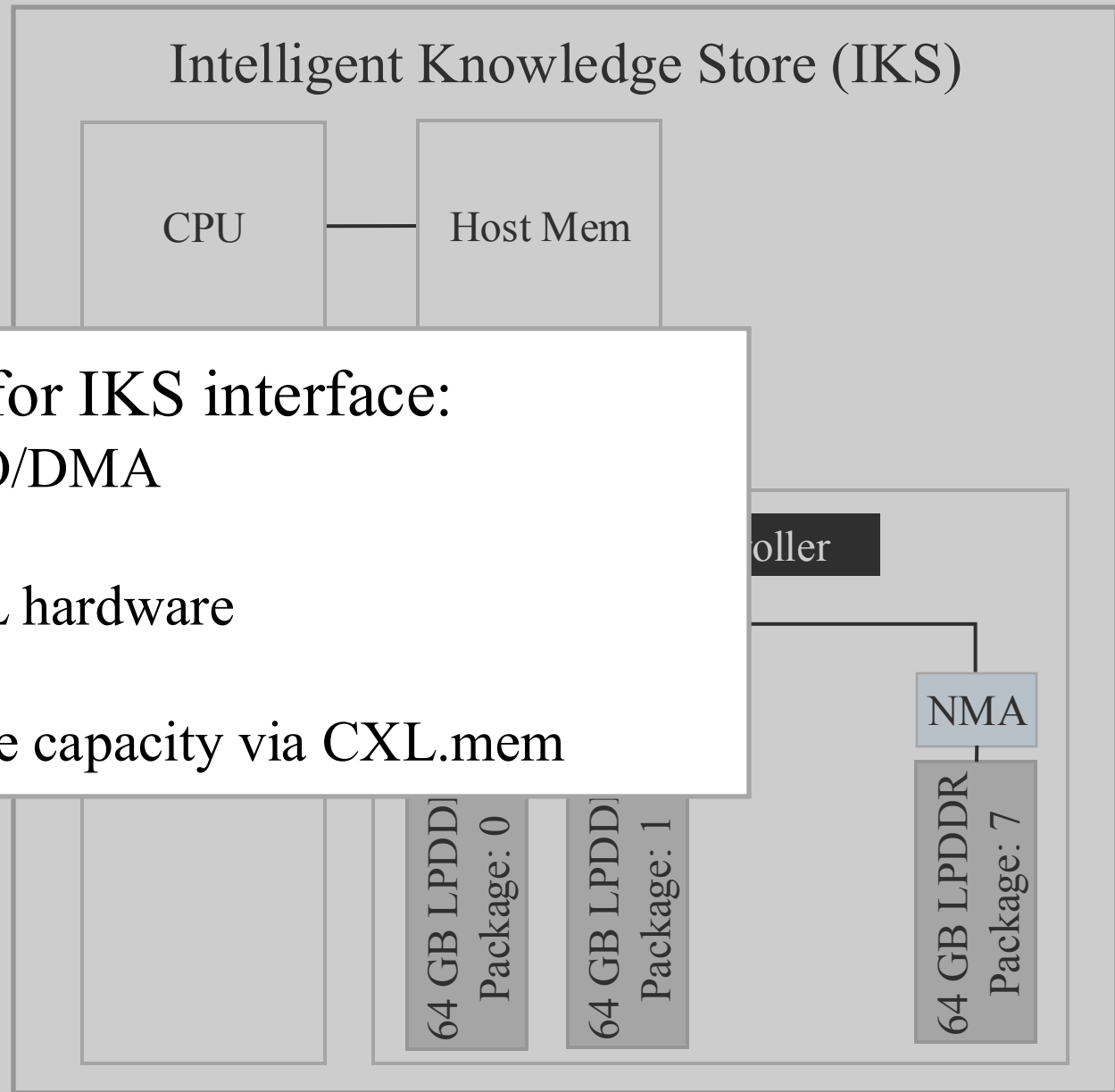  - Build request in cache: high BW!
  - No host mem hop: low latency!

Intelligent Knowledge Store (IKS)

CPU — Host Mem

CXL

GPU

CXL Controller

NMA    NMA    NMA

64 GB LPDDR Package: 0    64 GB LPDDR Package: 1    64 GB LPDDR Package: 7

# Leveraging CXL to implement IKS

- IKS *collaborates* with CPU!

- Option 1
  - High

- Option 2
  - Low

- Our app
  - Build request in cache: high BW!
  - No host mem hop: low latency!

## Intelligent Knowledge Store (IKS)

CPU —— Host Mem

oller

NMA

64 GB LPDDR Package: 0

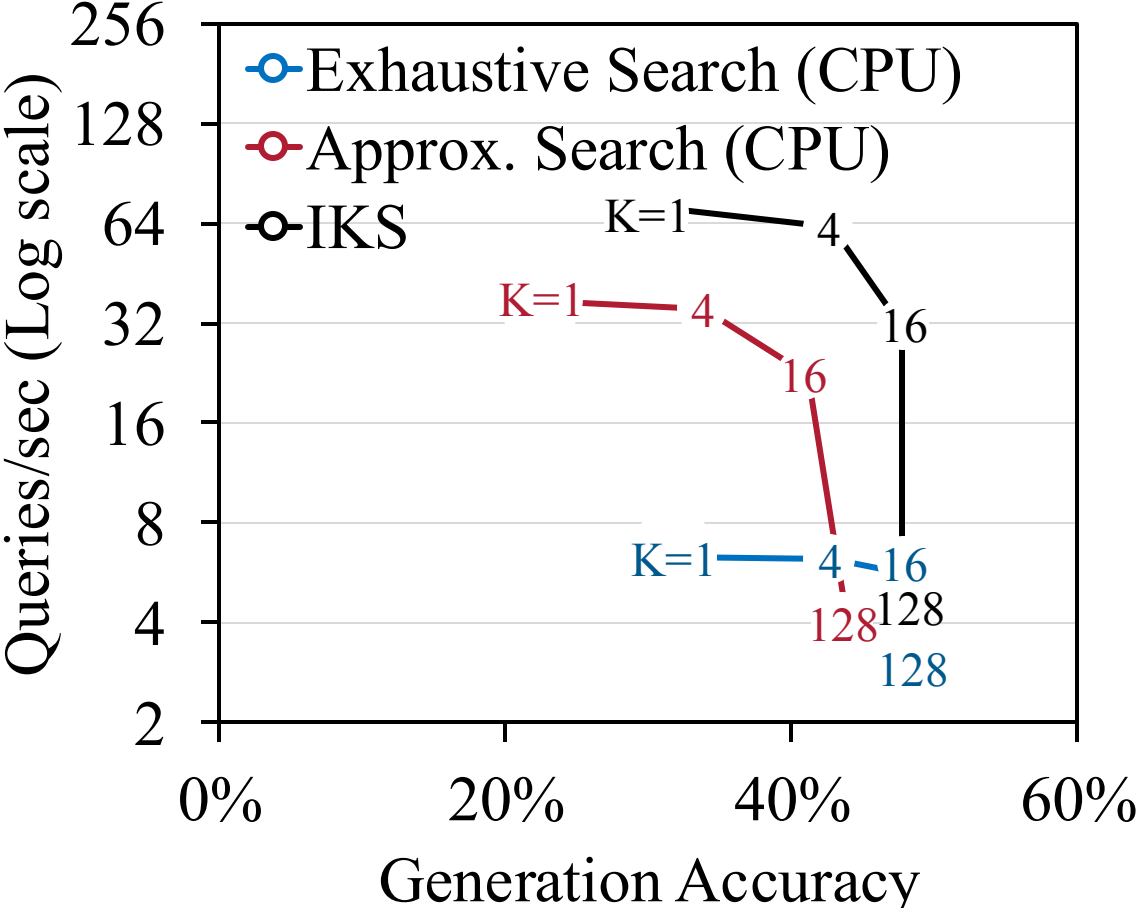64 GB LPDDR Package: 1

64 GB LPDDR Package: 7

---

## CXL.mem/cache for IKS interface:

- Improved performance vs. PIO/DMA

- Uses existing commodity CXL hardware

- Enables disaggregation of huge capacity via CXL.mem
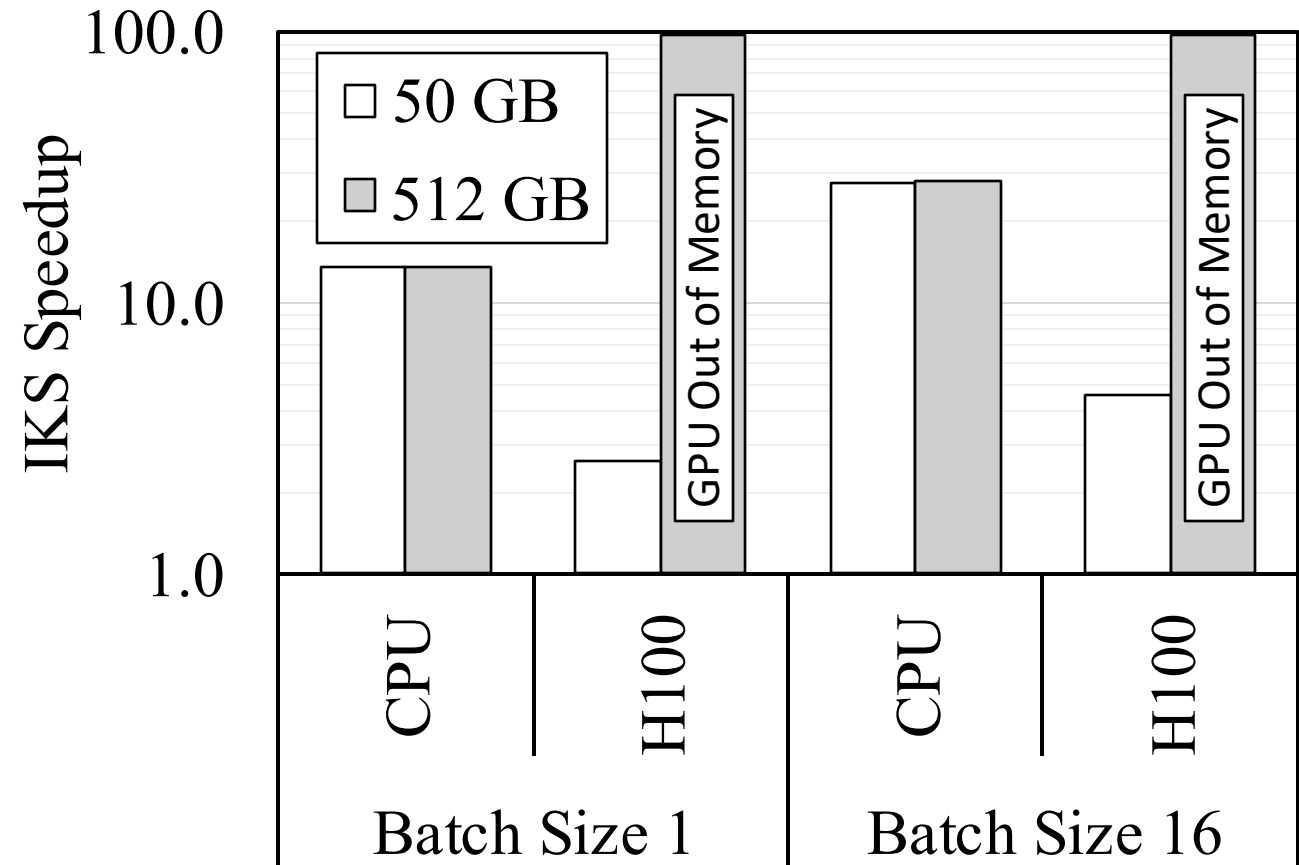
# Evaluating IKS

# IKS Accelerates RAG
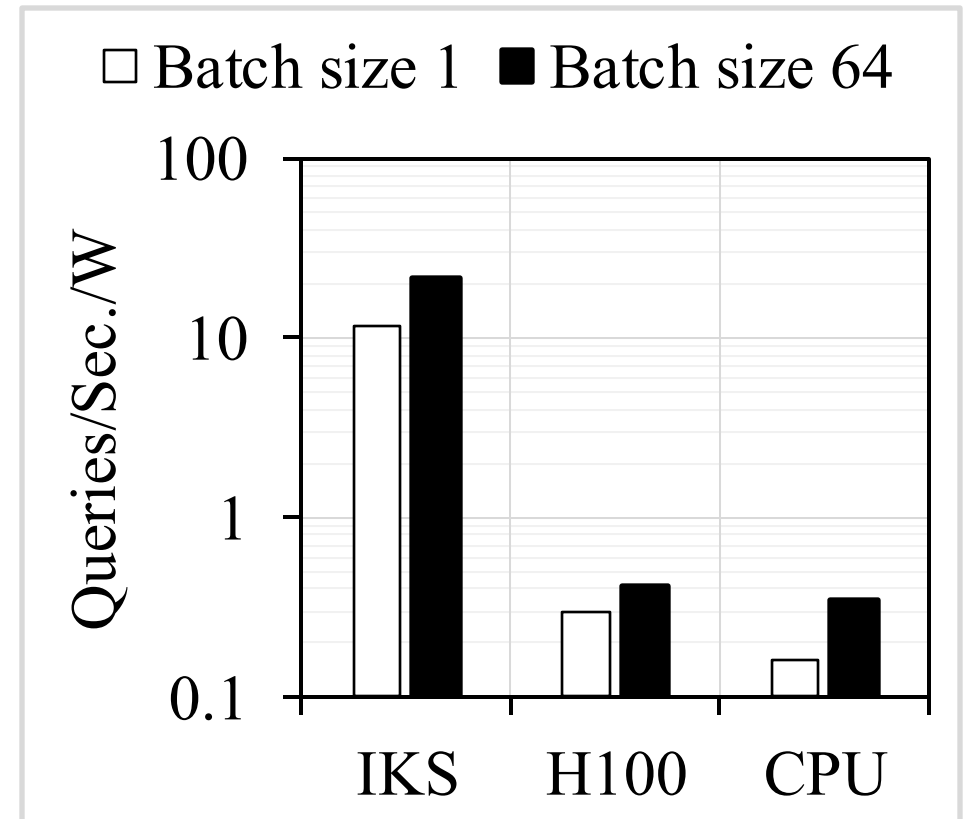
# Effectiveness of IKS Retrieval

- IKS speedup:
  - Up to 13.5x over CPU
  - Up to 4.6x over 1x GPU

- 512 GB capacity enables large-scale retrieval

IKS vs. CPU and H100 GPU (Exhaustive Search)

# Area and Power Efficiency

- Each NMA: 3.4 mm$^2$ (TSMC 16 nm)
- Total area: 220 mm$^2$
- LPDDR: 34.7 *W*
- Entire IKS power consumption
  - Batch size 1: 35.2 *W*
  - Batch size 64: 65.0 *W*



Throughput Retrieval Efficiency Comparison
(Higher is Better)

# Contributions

- We showcased the system-level interactions of retrieval quality
- We leveraged algorithm-hardware co-design to build a high-quality retrieval accelerator, accelerating RAG
- We showcased the utility of CXL for the interface of high-capacity accelerators

For more information visit Alian Research Group at
https://arg.csl.cornell.edu

**Cornell**Engineering
Electrical and Computer Engineering