# *Scalable Billion-point Approximate Nearest Neighbor Search Using SmartSSDs*

**Bing Tian**, Haikun Liu, Zhuohui Duan, Xiaofei Liao, Hai Jin

Huazhong University of Science and Technology

HUAZHONG UNIVERSITY OF SCIENCE AND TECHNOLOGY
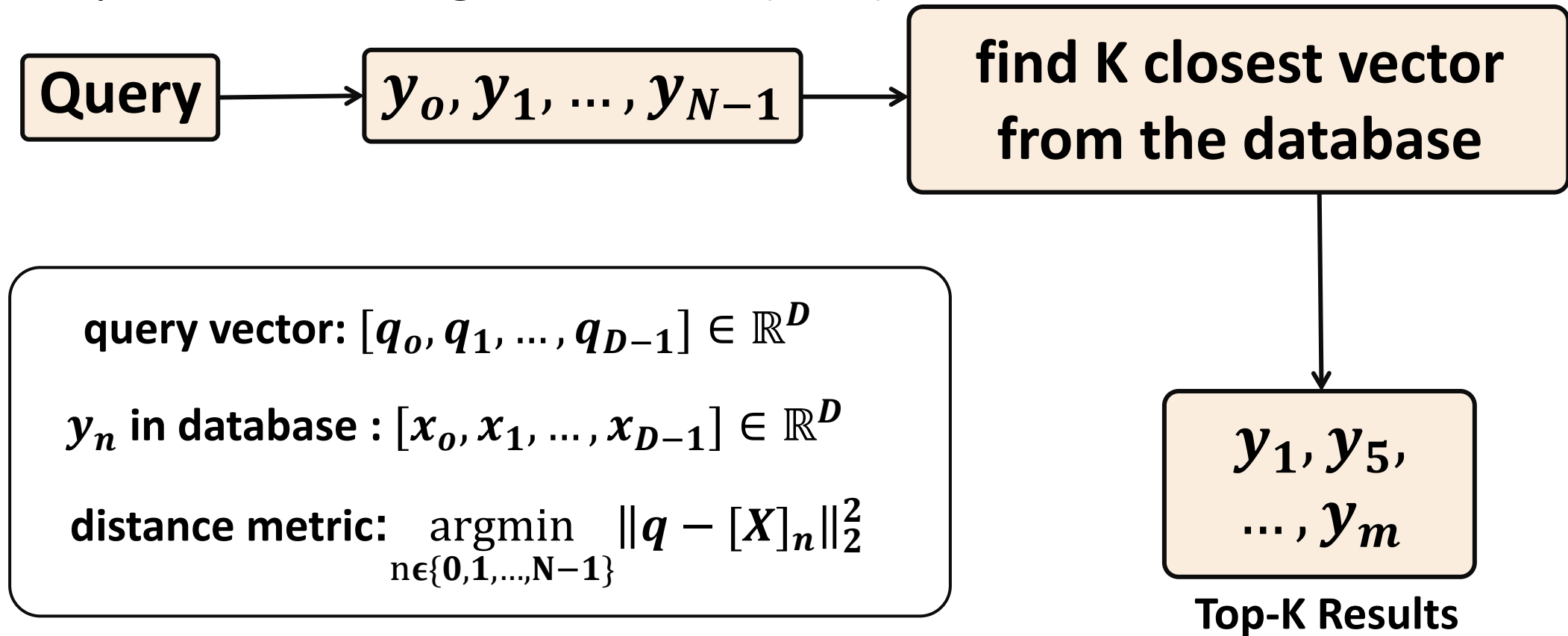
# Outline

❖**Background and Motivation**

❖SmartANNS Design

❖Results

❖Conclusion

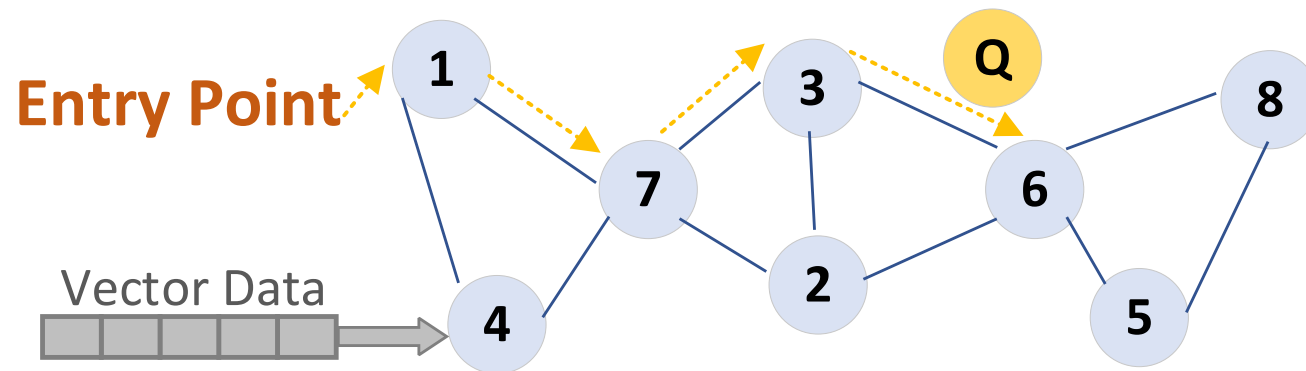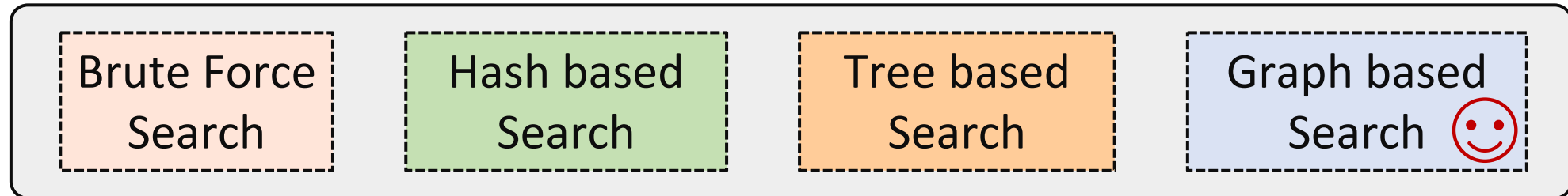# Background of ANNS

❖ Top-K Nearest Neighbor Search (NNS)

$$\boxed{\textbf{Query}} \rightarrow \boxed{y_o, y_1, \dots, y_{N-1}} \rightarrow \boxed{\textbf{find K closest vector from the database}}$$

query vector: $[q_o, q_1, \dots, q_{D-1}] \in \mathbb{R}^D$

$y_n$ in database : $[x_o, x_1, \dots, x_{D-1}] \in \mathbb{R}^D$

distance metric: $\underset{n \in \{0,1,\dots,N-1\}}{\mathrm{argmin}} \|q - [X]_n\|_2^2$

$$y_1, y_5, \dots, y_m$$

**Top-K Results**

# Background of ANNS

❖ Approximate Nearest Neighbor Search (ANNS)

| Brute Force Search | Hash based Search | Tree based Search | Graph based Search 😊 |
|---|---|---|---|



**Entry Point**

Vector Data
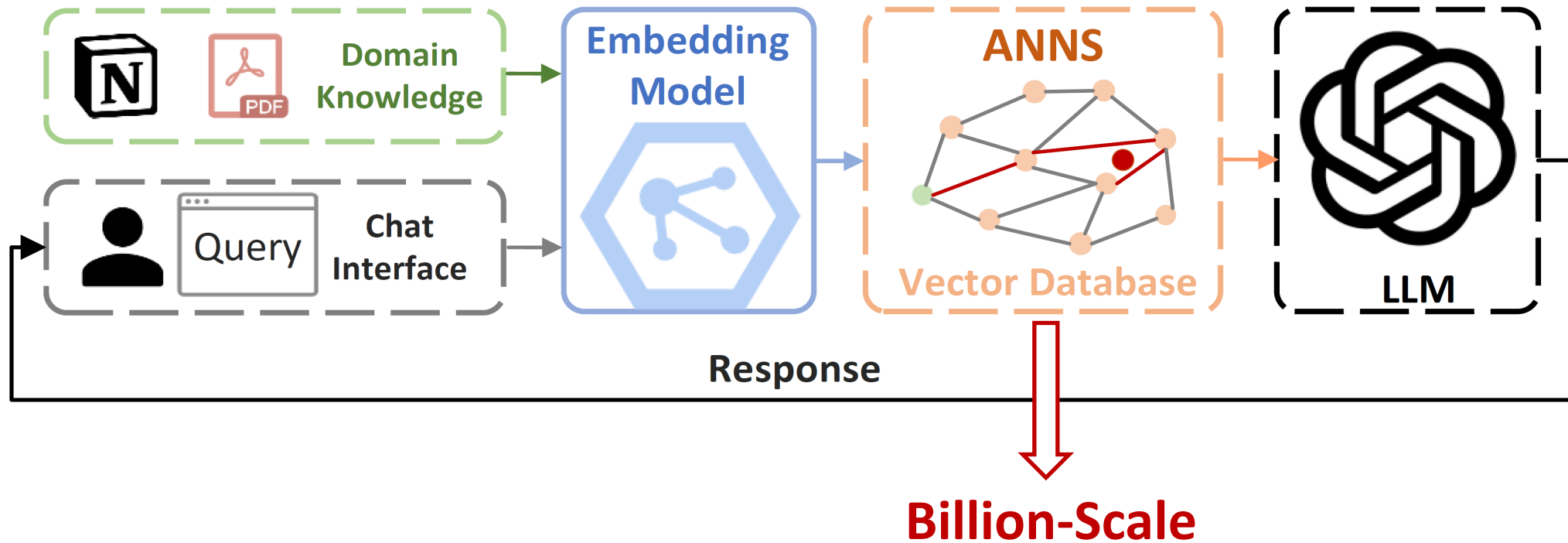
**Traversal Path：**
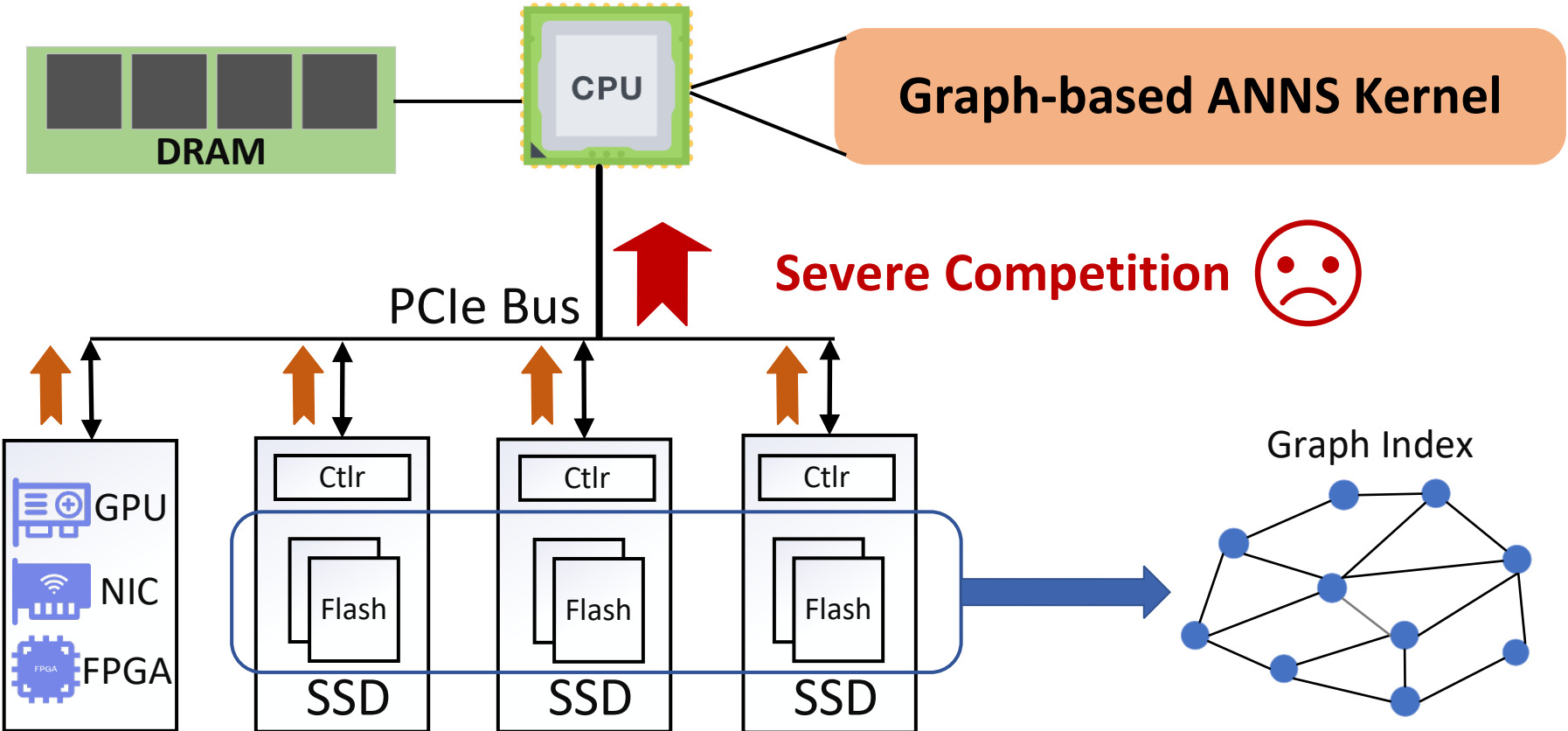
1   7   3   6

**Result:**

3   6

# Background of ANNS

❖ Retrieval-Augmented Generation

# Traditional Computing Architecture for ANNS



**DRAM**

**CPU**

Graph-based ANNS Kernel

PCIe Bus

**Severe Competition**

GPU

NIC

FPGA

Ctlr

Ctlr

Ctlr

Flash

Flash

Flash

SSD

SSD

SSD

Graph Index

# CSD-empowered NDP Architecture



**Large Storage Capacity**

**Less Data Movement**

**Mitigate Host Resource Contention**

**Performance scalability**

# CSD-empowered NDP Architecture

**Large Storage Capacity**

Host

DRAM  CPU

GPU
NIC
FPGA

Ctrl  Ctrl  Ctrl

Flash  Flash  Flash

CSD  CSD  CSD

Controller

NAND

Movement

**Using SmartSSDs to handle large-scale ANNS is promising …**

**Performance scalability**

**Mitigate Host
Resource Contention**

# CSD-based ANNS Solution

❖ **Offline Index Construction**

   ❖ Split dataset

   ❖ Construct graph for each partition

❖ **Online Search**
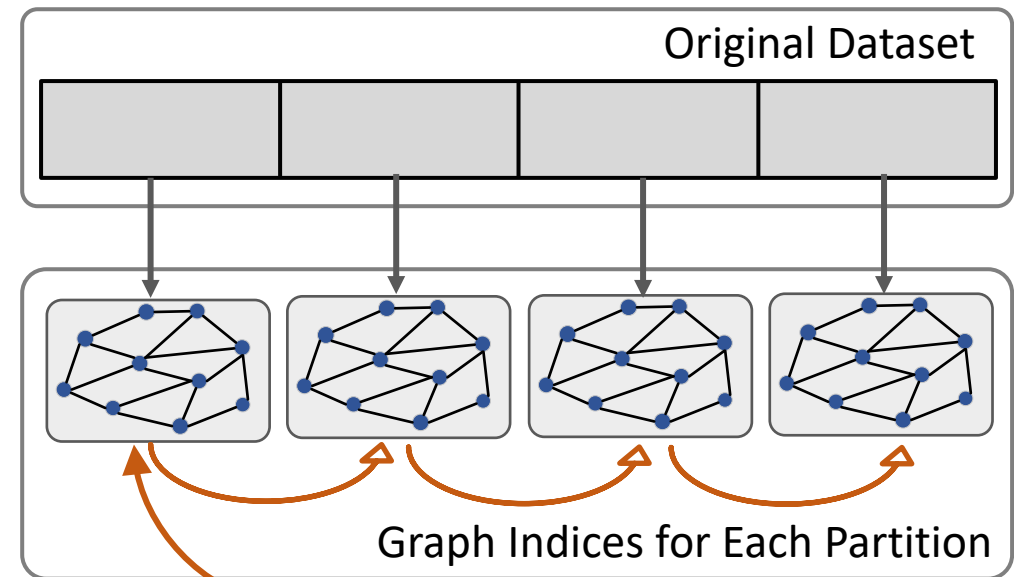
   ❖ <span style="color:red">Traverse all the graph indices</span>

   ❖ Merge all intermediate results

   ❖ Return top-k result

**Significant Computation** **+** **Limited Resource**

**Sub-optimal Performance**

CSDANN [TC'22], SmartSSD-based ANNS



Original Dataset

Graph Indices for Each Partition

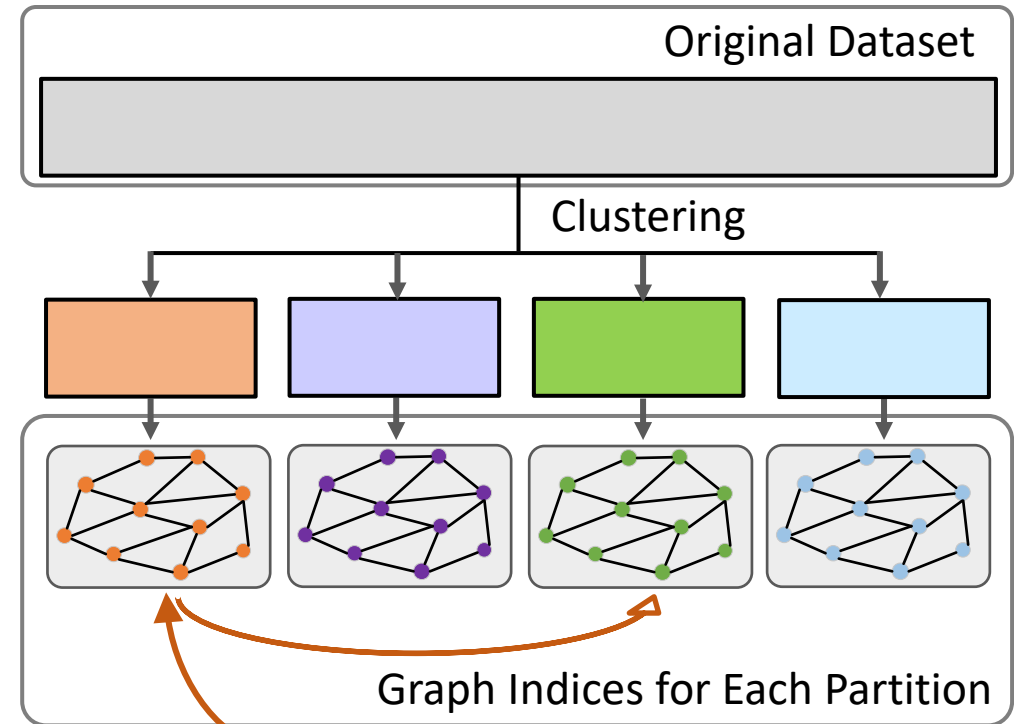# Opportunities of Hierarchical Indexing

❖ **Offline Index Construction**

    ❖ Partition dataset using clustering

    ❖ Construct graph for each shard

❖ **Online Search**

    ❖ Prune irrelevant shards

    ❖ Traverse closest graph of shards

    ❖ Merge and return top-k result

**Less computing overhead**

Original Dataset

Clustering

Graph Indices for Each Partition

# Challenges

## 1. Lack of communication channels



**Global Coordinator**

## 2. Load imbalance across SmartSSDs



**Task Scheduler**

Tasks    Data

Tasks    Data

Tasks    Data

## 3. Differences between queries



Query 1

Query 2

Query 3

**Optimizer**

# Outline

❖ **Background and Motivation**

❖ **SmartANNS Design**

❖ **Results**

❖ **Conclusion**

# Key Designs

Lack of communication channels

**Hierarchical indices in host and SmartSSDs**

Load imbalance across SmartSSDs

**Task scheduling based on the optimized data layout**

Differences between queries

**Learning-based shard pruning algorithm**

# Hierarchical indices

## Construction Step

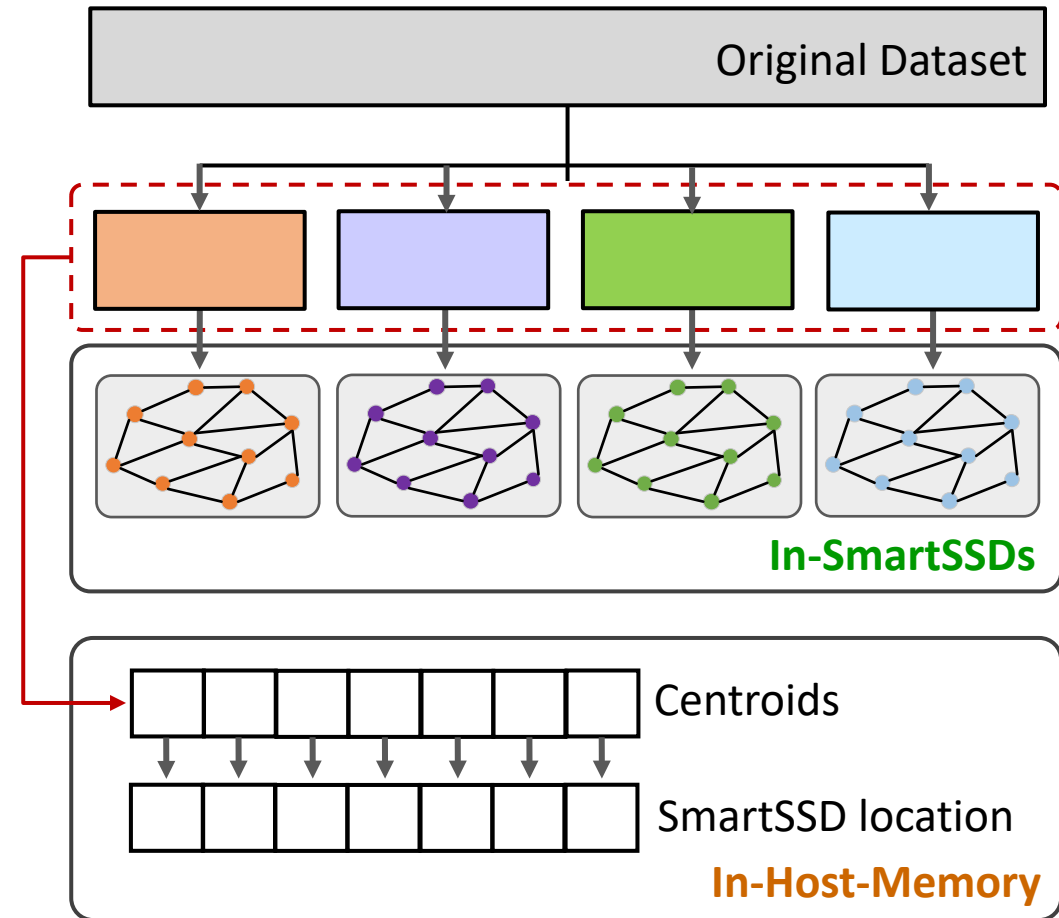❖ Hierarchical Balanced Clustering

❖ Construct HNSW graph indices

For SmartSSDs

❖ Store graph indices in SmartSSDs

For Host CPU

❖ Extract centroids for each shard

❖ Build shard-SmartSSD mapping table

**"Host CPU + SmartSSDs" cooperation**

# Shard Pruning

**Gradient Boosting Decision Trees (GBDT)**

❖ A strong predictive model combing multiple decision trees

Training Stage

| Training size | One million |
|---|---|
| **Learning rate** | 0.05 |
| **Iteration** | 500 |

❖ Iteratively predicting the mean, computing residuals, and fitting weak decision trees to these negative gradients

Inference Stage

❖ Assimilating weighted contributions of all individual weak models

Input Features

**Lightweight and high performance**

❖ The query vector

❖ Distance between the query and the top-k nearest shards ($D_k$) / distance between the query and the top-1 nearest shard ($D_1$)

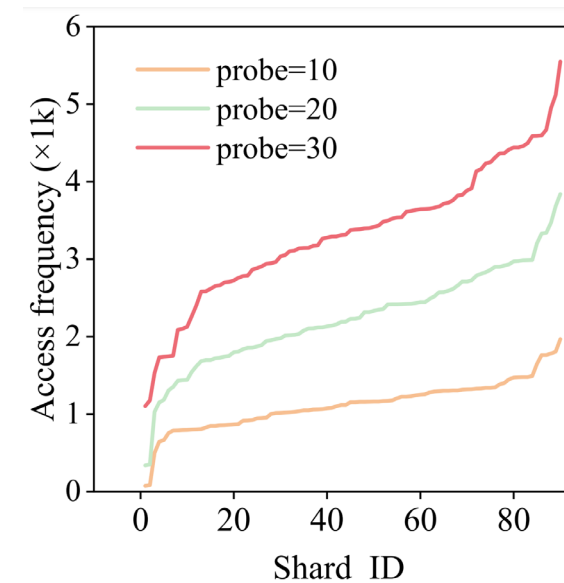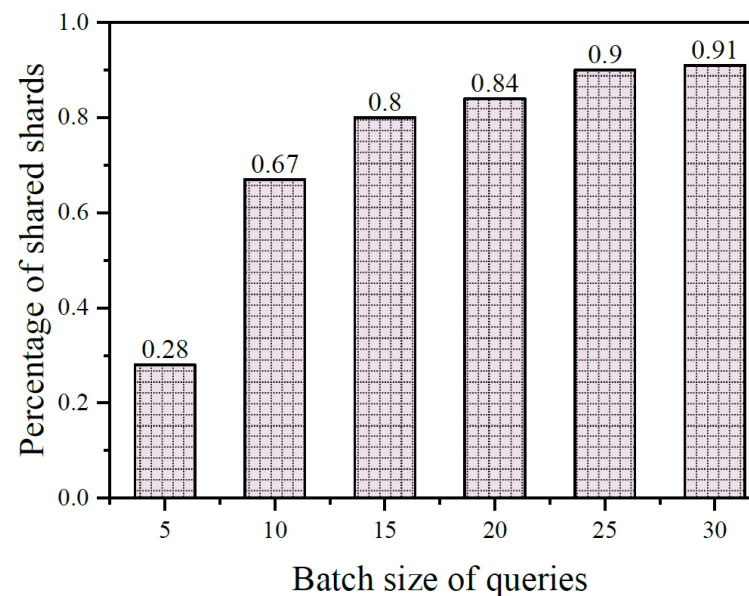❖ The total number of all shards

# Task Scheduling

Data Access Pattern of Hierarchical Indices

**Observation 1**

❖ A large portion of shards are accessed by different queries over a period of time, implying a good data locality

**Observation 2**

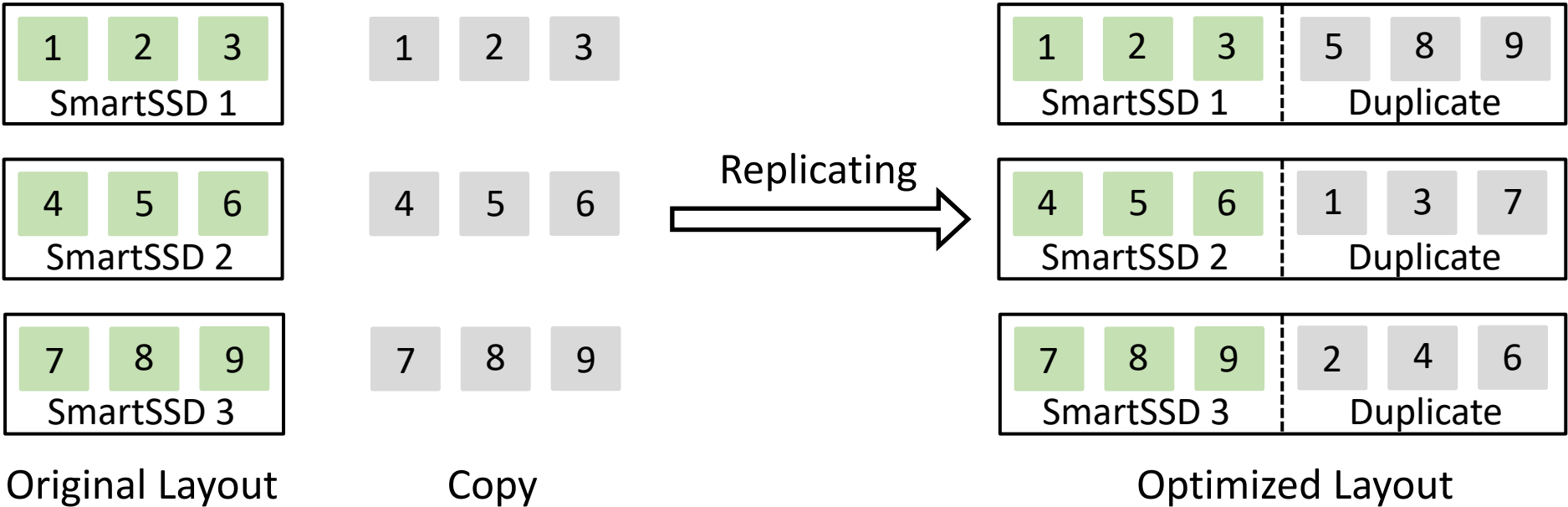❖ The access distribution of different shards are highly skewed



💡 1. Exploiting the data locality among queries

2. Placing hot shards on different SmartSSDs

# Task Scheduling

Optimized Data Layout

**Offering more flexibility for task scheduling**

❖ Iteratively placing shard with highest hotness to the SmartSSD with lowest cumulative hotness

❖ Replicating shards from one SmartSSD to another SmartSSD once



| 1 | 2 | 3 |
|---|---|---|
SmartSSD 1

| 4 | 5 | 6 |
|---|---|---|
SmartSSD 2

| 7 | 8 | 9 |
|---|---|---|
SmartSSD 3

Original Layout

| 1 | 2 | 3 |
|---|---|---|

| 4 | 5 | 6 |
|---|---|---|

| 7 | 8 | 9 |
|---|---|---|

Copy

Replicating

| 1 | 2 | 3 | 5 | 8 | 9 |
|---|---|---|---|---|---|
| SmartSSD 1 | | | Duplicate | | |

| 4 | 5 | 6 | 1 | 3 | 7 |
|---|---|---|---|---|---|
| SmartSSD 2 | | | Duplicate | | |

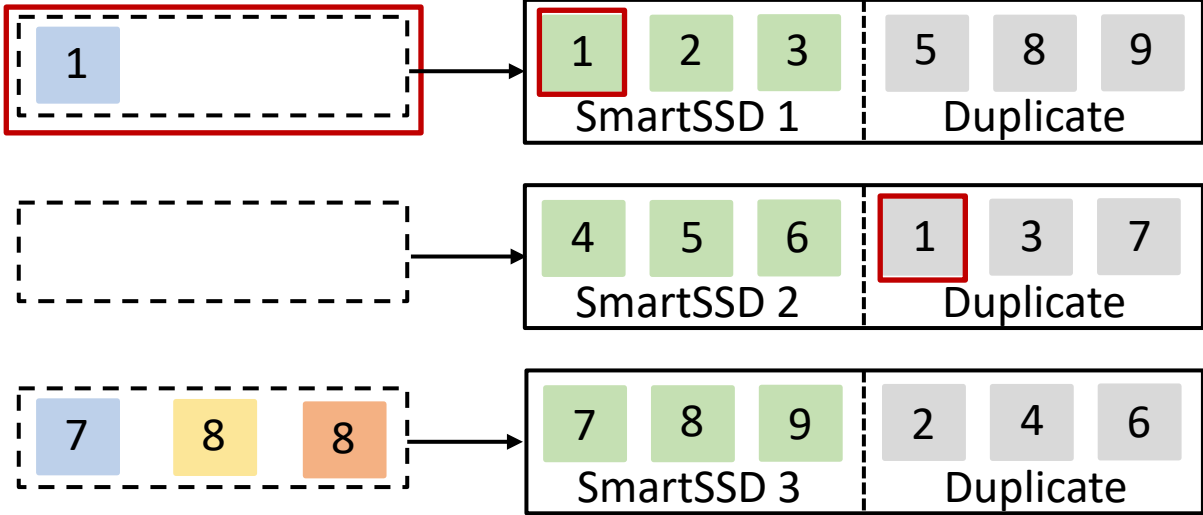| 7 | 8 | 9 | 2 | 4 | 6 |
|---|---|---|---|---|---|
| SmartSSD 3 | | | Duplicate | | |

Optimized Layout

# Task Scheduling

## Scheduling Steps

```
for i = 0; i < Tasks.size(); i++ do
    /* Find devices that store task-relevant shards.    */
    Base[] ← FindBase(Devices, Tasks[i]);
    /* Check whether there is data reuse on device for
       this task.                                       */
    Assign[] ← Check(Device_to_Tasks, Tasks[i]);
    if (Assign.size = 0) ∨ (Assign.size = Base.size)
      then
          Target ← MinWork(DeviceLoad, Base);
          Device_to_Tasks[Target].insert(Tasks[i]);
          DeviceLoad.update(Target);
    else
          TimeInfo[] ← NULL;
          for j = 0; j < Base.size(); j++ do
              Temp[] ← Device_to_Tasks[Base[j]];
              Temp.insert(Tasks[i]);
              TimeInfo[j] ← Estimate(Temp);
          end
          Target ← Min(TimeInfo);
          Device_to_Tasks[Target].insert(tasks[i]);
          DeviceLoad.update(Target);
    end
end
```

Query 1 → 1 3 7    Query 2 → 2 4 8    Query 3 → 1 3 8

1 → Device 1, Device2

| 1 |  →  | 1 | 2 | 3 | : | 5 | 8 | 9 |
SmartSSD 1 : Duplicate

| | → | 4 | 5 | 6 | : | 1 | 3 | 7 |
SmartSSD 2 : Duplicate

| 7 | 8 | 8 | → | 7 | 8 | 9 | : | 2 | 4 | 6 |
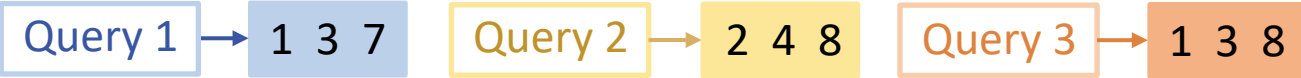SmartSSD 3 : Duplicate

Optimized Layout

# Task Scheduling

## Scheduling Steps

```
for i = 0; i < Tasks.size(); i++ do
    /* Find devices that store task-relevant shards.   */
    Base[ ] ← FindBase(Devices, Tasks[i]);
    /* Check whether there is data reuse on device for
       this task.                                      */
    Assign[ ] ← Check(Device_to_Tasks, Tasks[i]);
    if (Assign.size = 0) ∨ (Assign.size = Base.size)
       then
           Target ← MinWork(DeviceLoad, Base);
           Device_to_Tasks[Target].insert(Tasks[i]);
           DeviceLoad.update(Target);
    else
           TimeInfo[ ] ← NULL;
           for j = 0; j < Base.size(); j++ do
               Temp[ ] ← Device_to_Tasks[Base[j]];
               Temp.insert(Tasks[i]);
               TimeInfo[j] ← Estimate(Temp);
           end
           Target ← Min(TimeInfo);
           Device_to_Tasks[Target].insert(tasks[i]);
           DeviceLoad.update(Target);
    end
end
```
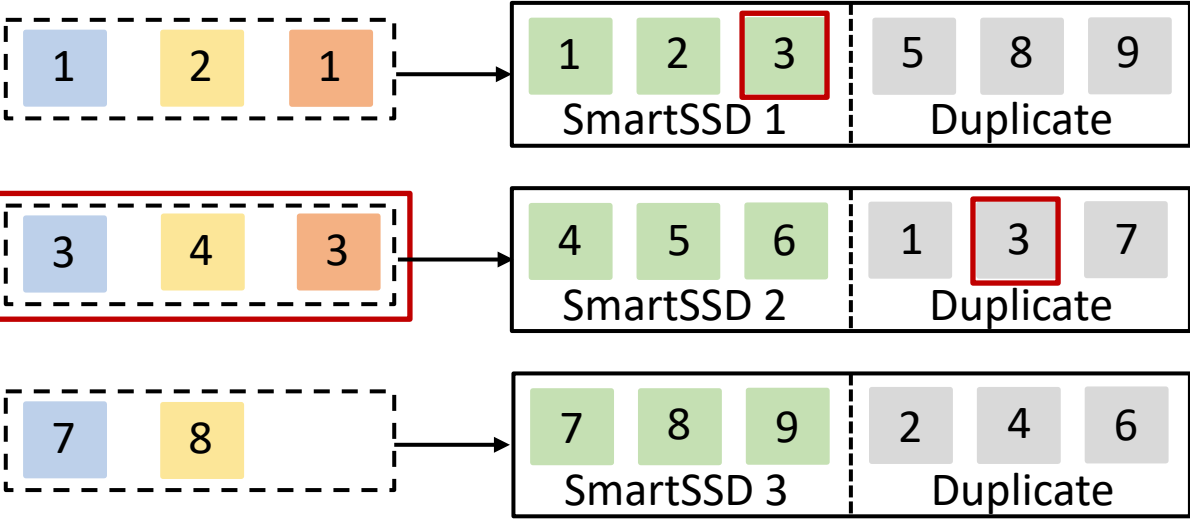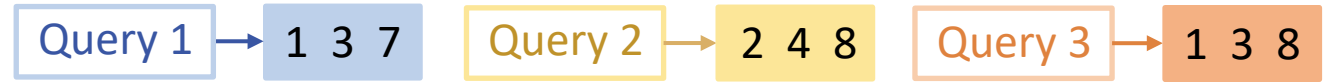


Query 1 → 1 3 7    Query 2 → 2 4 8    Query 3 → 1 3 8

3 → Device 1, Device2

| 1 | 2 | 1 |

| 1 | 2 | 3 | 5 | 8 | 9 |
SmartSSD 1 | Duplicate

| 3 | 4 | 3 |

| 4 | 5 | 6 | 1 | 3 | 7 |
SmartSSD 2 | Duplicate

| 7 | 8 |

| 7 | 8 | 9 | 2 | 4 | 6 |
SmartSSD 3 | Duplicate
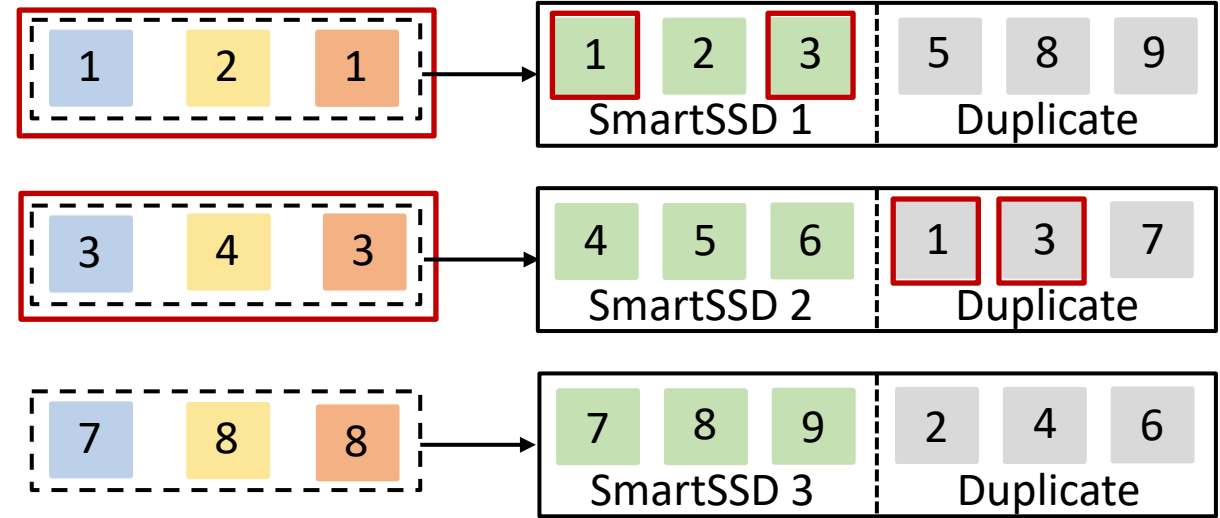
Optimized Layout

# Task Scheduling

## Scheduling Steps

```
for i = 0; i < Tasks.size(); i++ do
    /* Find devices that store task-relevant shards.  */
    Base[ ] ← FindBase(Devices, Tasks[i]);
    /* Check whether there is data reuse on device for
       this task.                                      */
    Assign[ ] ← Check(Device_to_Tasks, Tasks[i]);
    if (Assign.size = 0) ∨ (Assign.size = Base.size)
      then
        Target ← MinWork(DeviceLoad, Base);
        Device_to_Tasks[Target].insert(Tasks[i]);
        DeviceLoad.update(Target);
    else
        TimeInfo[ ] ← NULL;
        for j = 0; j < Base.size(); j++ do
            Temp[ ] ← Device_to_Tasks[Base[j]];
            Temp.insert(Tasks[i]);
            TimeInfo[j] ← Estimate(Temp);
        end
        Target ← Min(TimeInfo);
        Device_to_Tasks[Target].insert(tasks[i]);
        DeviceLoad.update(Target);
end
end
```



Query 1 → 1 3 7     Query 2 → 2 4 8     Query 3 → 1 3 8

**Load balancing and data reusing**

| 1 | 2 | 1 | → | 1 | 2 | 3 | 5 | 8 | 9 |

SmartSSD 1   Duplicate

| 3 | 4 | 3 | → | 4 | 5 | 6 | 1 | 3 | 7 |

SmartSSD 2   Duplicate

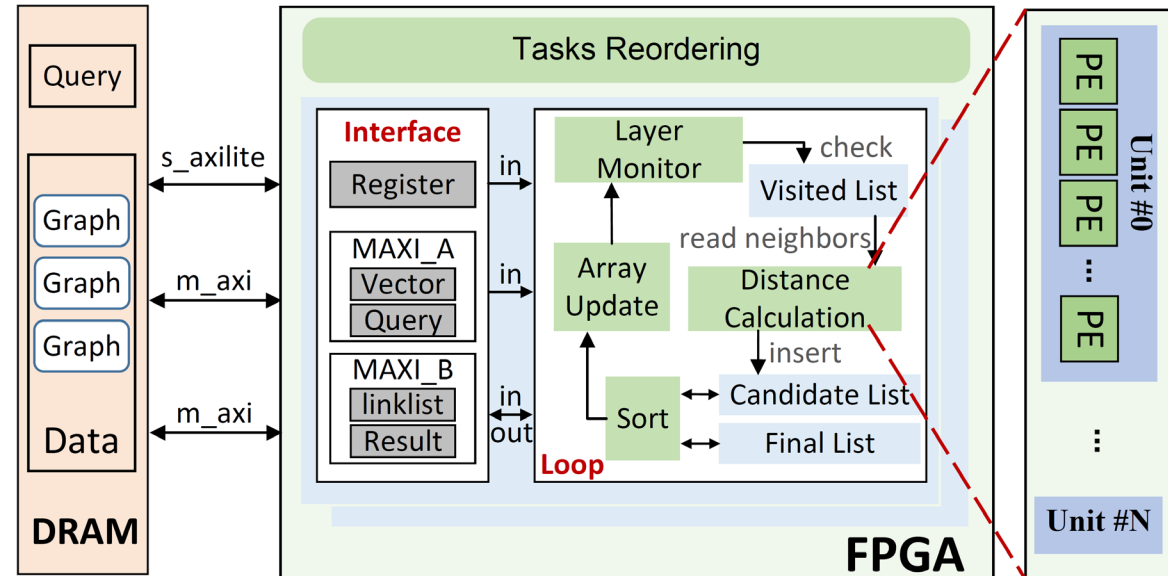| 7 | 8 | 8 | → | 7 | 8 | 9 | 2 | 4 | 6 |

SmartSSD 3   Duplicate

Optimized Layout

# Implementation

## Vector Search Engine

❖ Separate interface for parallel reading

❖ Boolean array as the visited list

❖ Bitonic sort algorithm for lists updating

❖ Loop unrolling and pipelining

❖ Data and kernel pooling



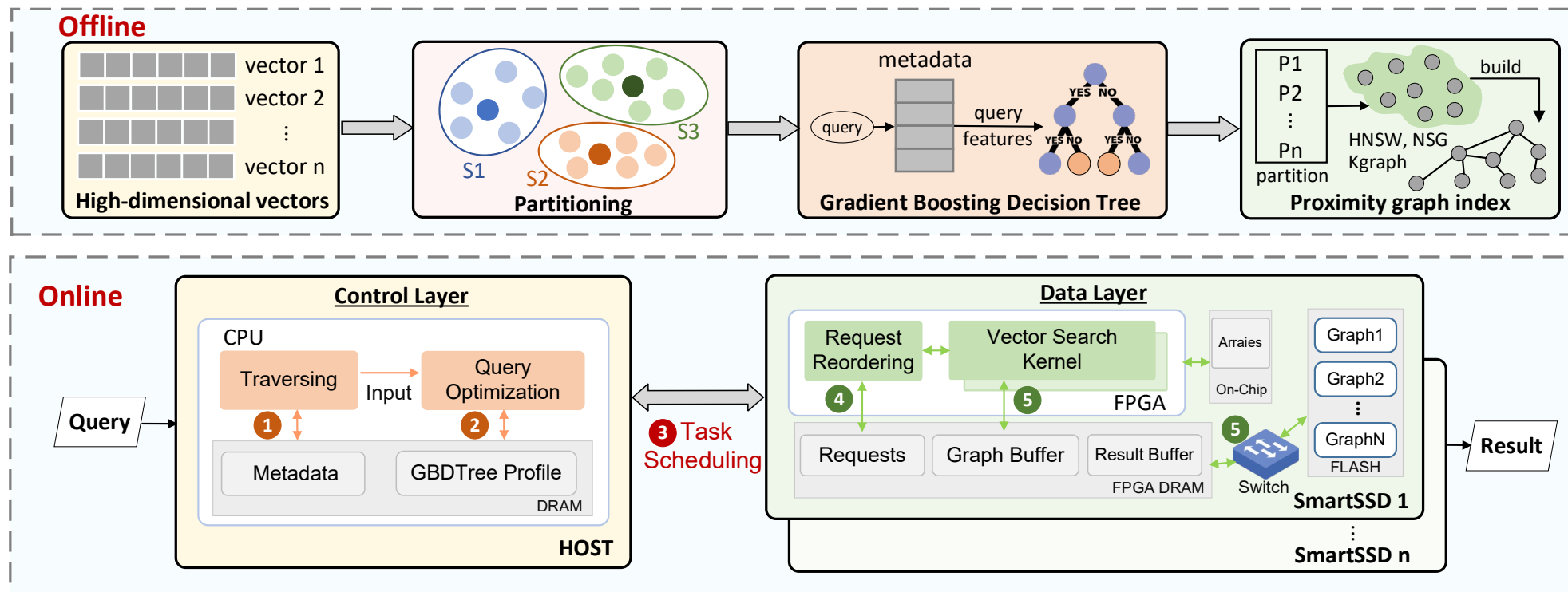## More details: checkout our paper

❖ FPGA kernel details

❖ GBDT implementation (LAET SIGMOD'20)

❖ HNSW search process

# SmartANNS System

❖ "host CPU + SmartSSDs" cooperative processing architecture

# Outline

❖**Background and Motivation**

❖**SmartANNS Design**

❖**Results**

❖**Conclusion**

# Experimental Setup

## Hardware Platform

### Server

| CPU | 2 * Intel Xeon Gold 5220 CPUs |
|-----|-------------------------------|
| GPU | Nvidia Tesla V100 (32GB HBM) |
| DRAM | 128 GB DDR4 |
| OS | Ubuntu 20.04.4 LTS |

### Computational Storage Device

| Type | Samsung SmartSSD |
|------|------------------|
| FPGA | Xilinx Kintex UltraScale+KU15P |
| DRAM | 4 GB DDR4 |
| Flash | 4 TB, 4 GB/s |

## Datasets

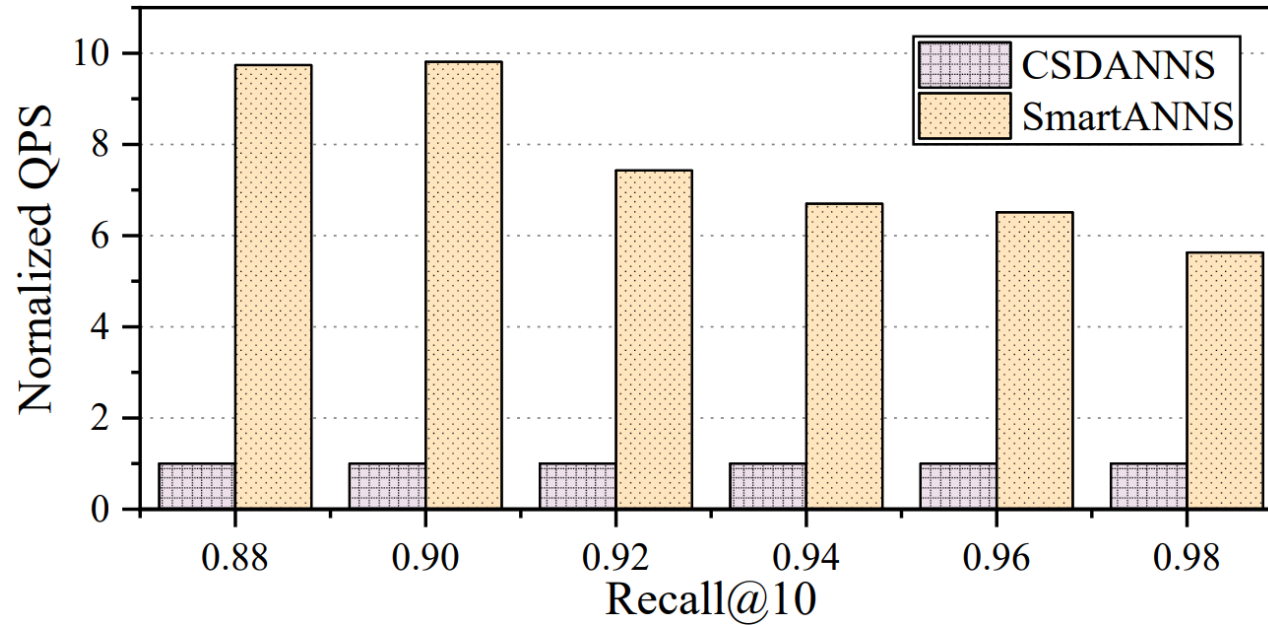| Dataset | Dimension | Data Type | Base Size | Source |
|---------|-----------|-----------|-----------|--------|
| SIFT1B | 128 | Uint8 | 119 GB | Image |
| SPACEV1B | 100 | Int8 | 93 GB | Web Search |
| DEEP1B | 96 | Float32 | 358 GB | Image |
| Turing1B | 100 | Float32 | 373 GB | Web Search |

# Comparison with Baselines

❖ Under different dataset



8.5-10.7X higher QPS compared with the state-of-the-art SmartSSD-based ANNS—CSDANNS

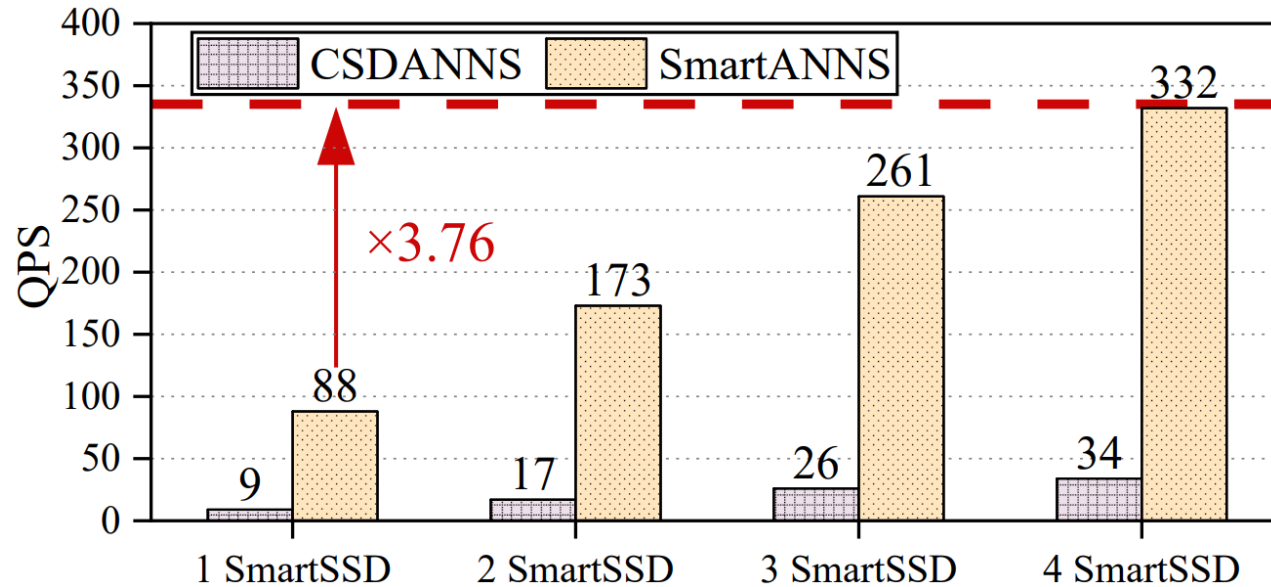# Comparison with Baselines

❖ Under different accuracy



With SIFT1B dataset, SmartANNS achieves 5.6-9.8X higher QPS compared with CSDANNS
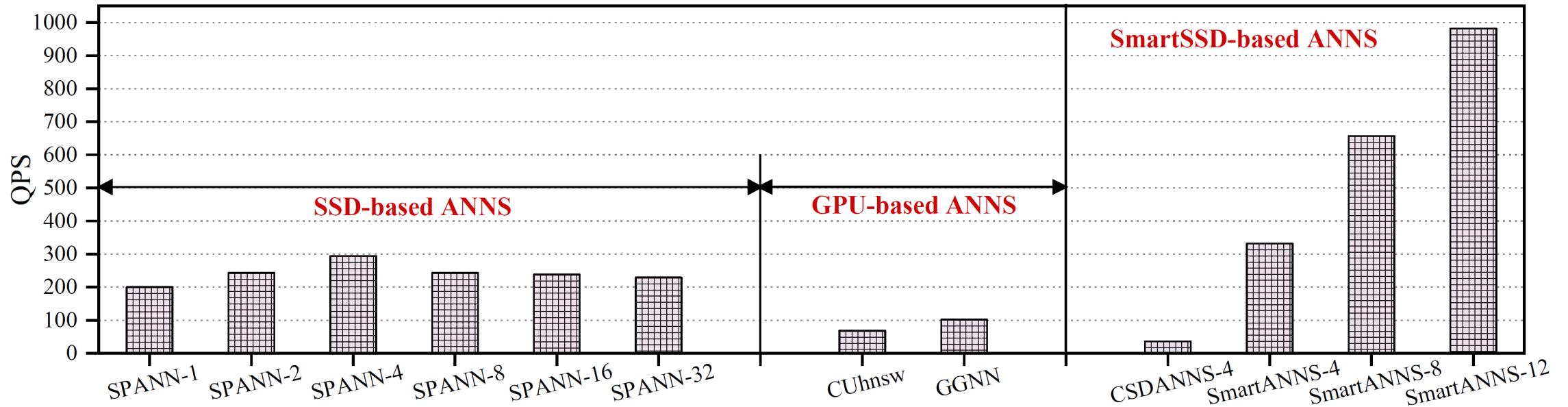
# Comparison with Baselines

❖ Scalability



SmartANNS achieve near-linear performance scalability with the increase of SmartSSDS

# Comparison with SSD/GPU-based ANNS



SmartANNS is more efficient than SSD-based solution and GPU-based solutions

# Outline

❖**Background and Motivation**

❖**SmartANNS Design**

❖**Results**

❖**Conclusion**

❖ SmartANNS:

*A hardware/software co-design architecture using SmartSSDs to support the large-scale and scalable ANNS service*

# *Thanks & QA*