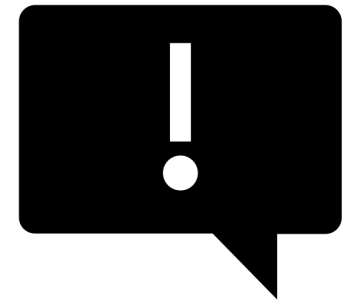# Index Structure Journey 1st Week
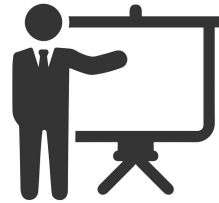# Traditional Index Structure Overview

Presented by Hojin Shin

System Software Lab.

# Contents

- Introduction: Data and In-memory Index

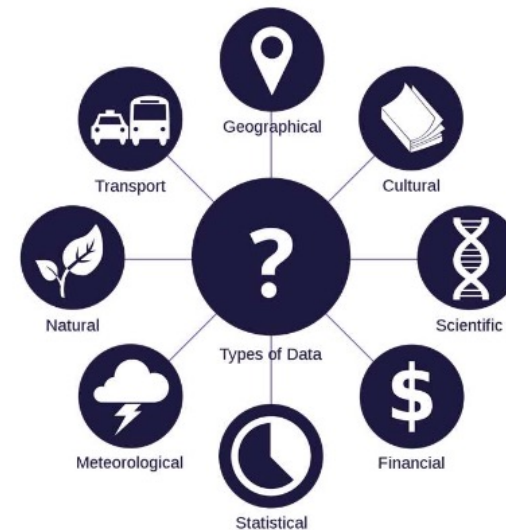- Various In-memory Index

- Study Github: Index Structure Journey

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# Introduction:
# Data and In-memory Index

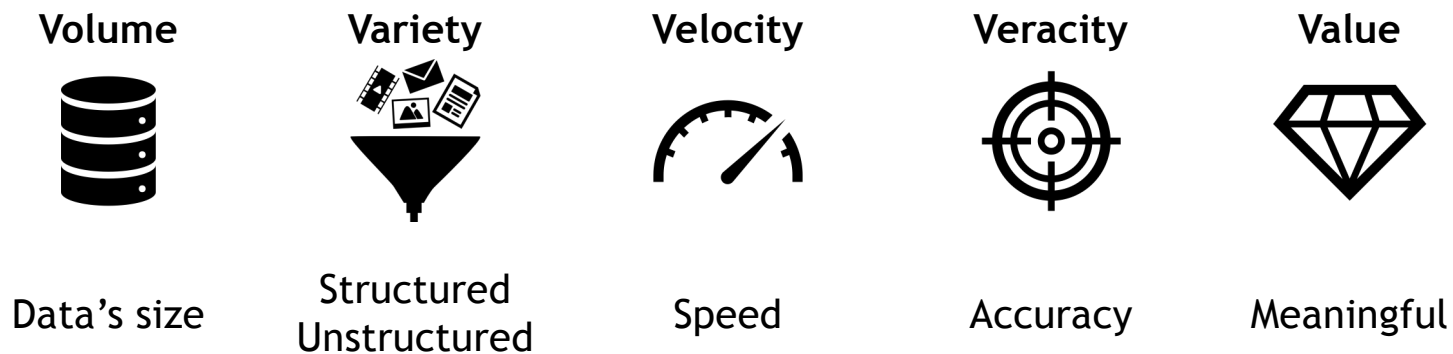# Introduction: Data and In-memory Index

- **What is data?**
  - Units of information, often numeric, that are collected through observation
  - Fact on which a theory is based
  - Data in the form of letters, numbers, sounds, pictures that a computer can process

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# Introduction: Data and In-memory Index

- **What is BigData?**
  - A large amount of structured data that exceed existing DB management tools
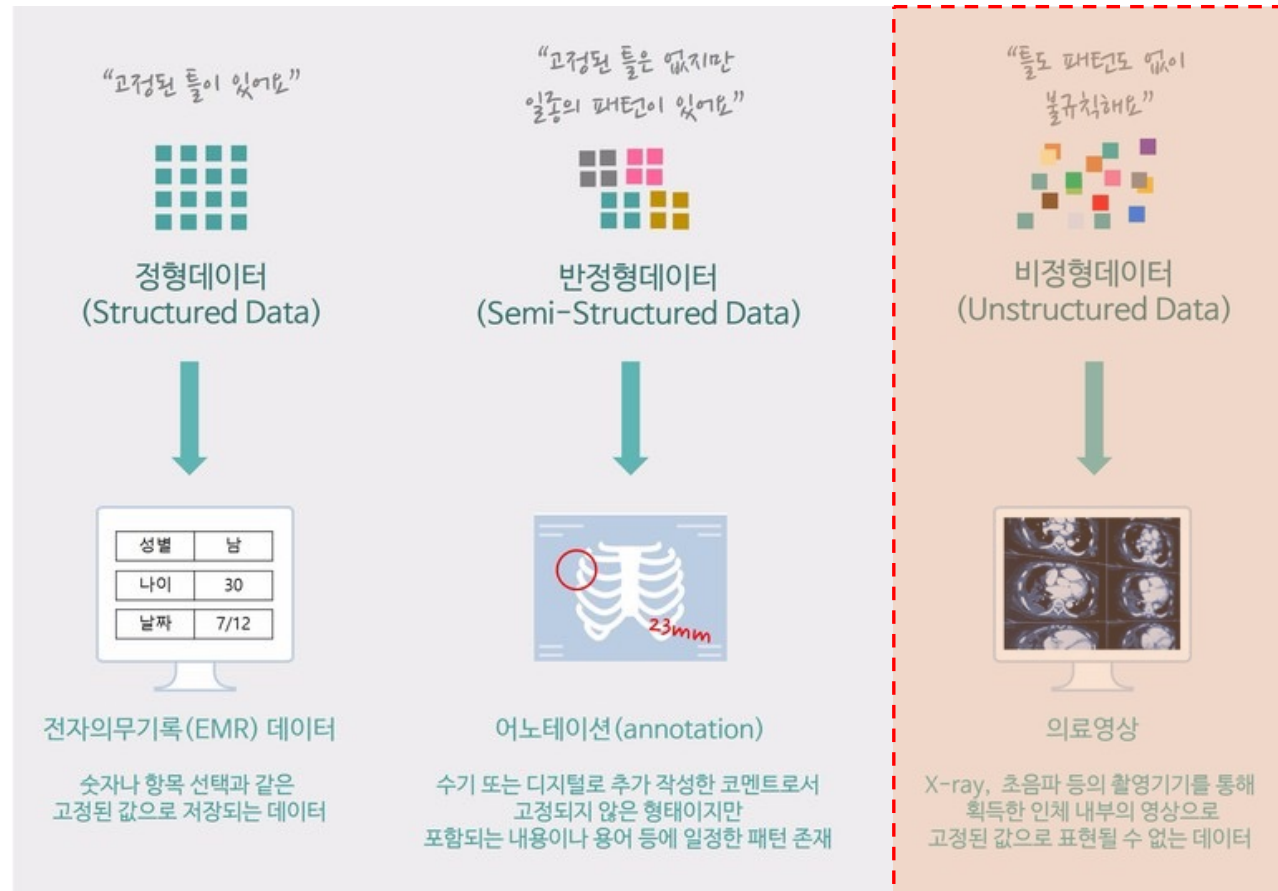  - Set of unstructured data that is not in the form of data
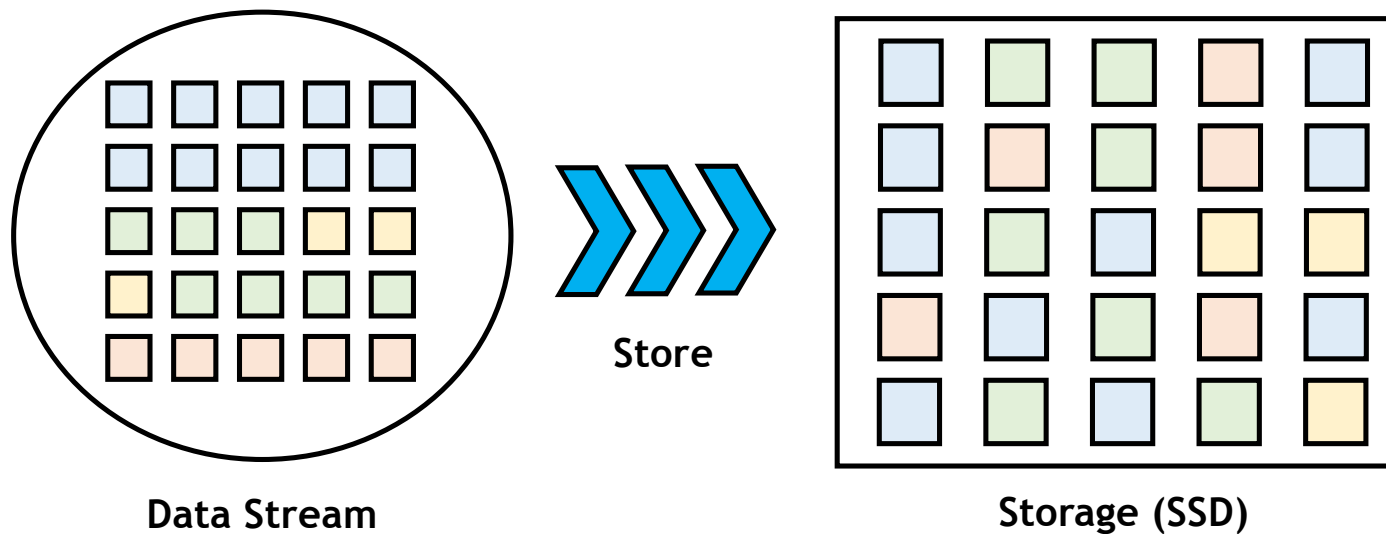
**BigData 5V**

| Volume | Variety | Velocity | Veracity | Value |
|--------|---------|----------|----------|-------|
| Data's size | Structured Unstructured | Speed | Accuracy | Meaningful |

# Introduction: Data and In-memory Index

- **Kind of Data**

  - Structured Data

    - Data organized and processed into a form suitable for immediate statistical analysis

    - Data stored in fixed fields

  - Unstructured Data

    - One piece of data, not a set of data, is objectified as collected data

    - Difficult to understand the meaning of a value because there is no set rule

  - Semi-structured Data

    - File type, metadata (schema of structured data inside data)

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# Introduction: Data and In-memory Index

# Introduction: Data and In-memory Index

- What is In-memory Index?



Data Stream

Store

Storage (SSD)

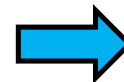DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# Introduction: Data and In-memory Index

- What is In-memory Index?



User  →  **Blue Color**  →  Storage (SSD)

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# Introduction: Data and In-memory Index

- What is In-memory Index?

**User** → **Blue Color** → 

**Index (Memory)**

- location (block#)
- location (block#)
- location (block#)
- location (block#)

→ **Storage (SSD)**

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# Introduction: Data and In-memory Index

- **The reason to use In-memory indexing**

  - Memory has a low latency for processing requests

    - If we process each request to disk, it goes slow down

    - DRAM: 100ns < SATA SSD: ~70us

  - To effectively index the input data

  - To reduce disk-based I/O

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# Introduction: Data and In-memory Index

- **The reason to use In-memory indexing**

  - Memory has a low latency for processing requests

    – If we process each request to disk, it goes slow down

    – DRAM: 100ns < SATA SSD: ~70us

  - To effectively index the input data

  - To reduce disk-based I/O

```
    0.5 ns - CPU L1 dCACHE reference
    1   ns - speed-of-light (a photon) travel a 1 ft (30.5cm) distance
    5   ns - CPU L1 iCACHE Branch mispredict
    7   ns - CPU L2  CACHE reference
   71   ns - CPU cross-QPI/NUMA best  case on XEON E5-46*
  100   ns - MUTEX lock/unlock
  100   ns - own DDR MEMORY reference
  135   ns - CPU cross-QPI/NUMA best  case on XEON E7-*
  202   ns - CPU cross-QPI/NUMA worst case on XEON E7-*
  325   ns - CPU cross-QPI/NUMA worst case on XEON E5-46*
 10,000   ns - Compress 1K bytes with Zippy PROCESS
 20,000   ns - Send 2K bytes over 1 Gbps NETWORK
250,000   ns - Read 1 MB sequentially from MEMORY
500,000   ns - Round trip within a same DataCenter
10,000,000   ns - DISK seek
10,000,000   ns - Read 1 MB sequentially from NETWORK
30,000,000   ns - Read 1 MB sequentially from DISK
150,000,000   ns - Send a NETWORK packet CA -> Netherlands
|  |  |  |
|  |  | ns|
|  | us|
| ms|
```

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# Various In-memory Index

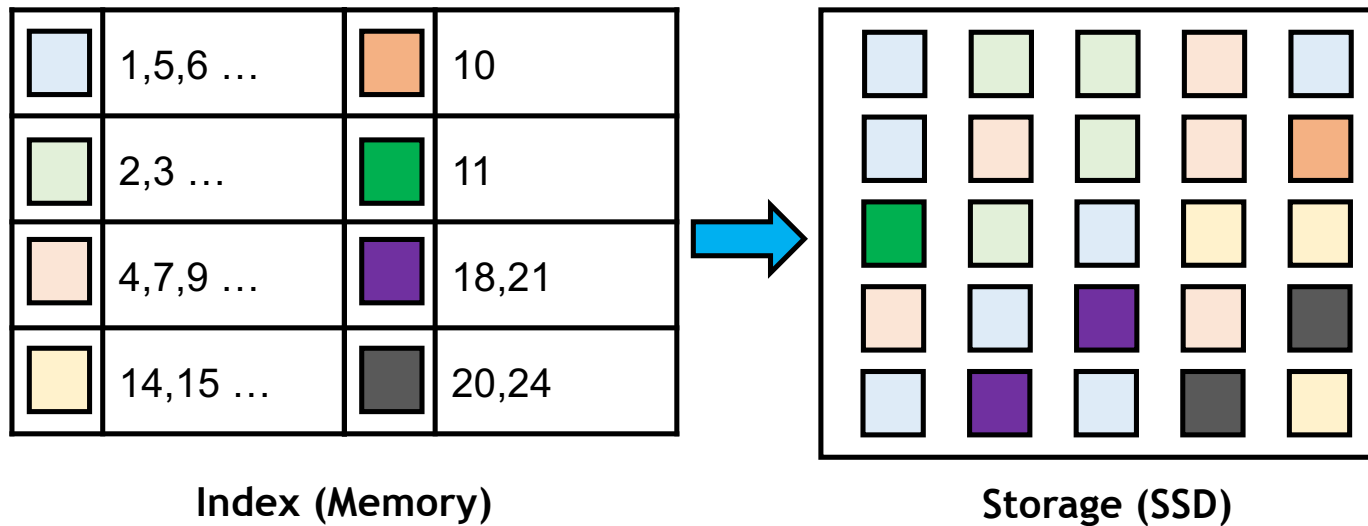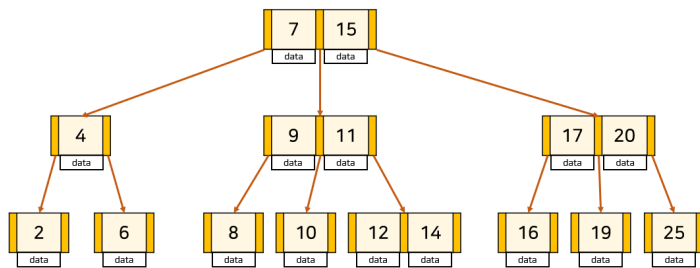# What I've seen: Observation

- Why do various in-memory indexes exist?
  - Aren't arrays enough?
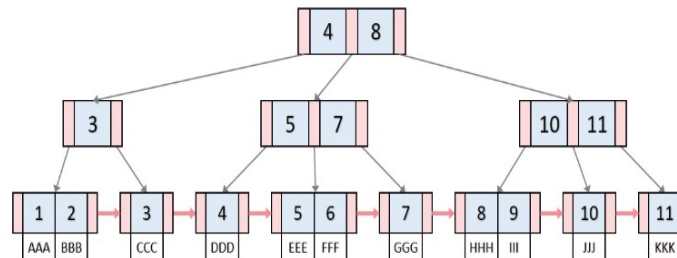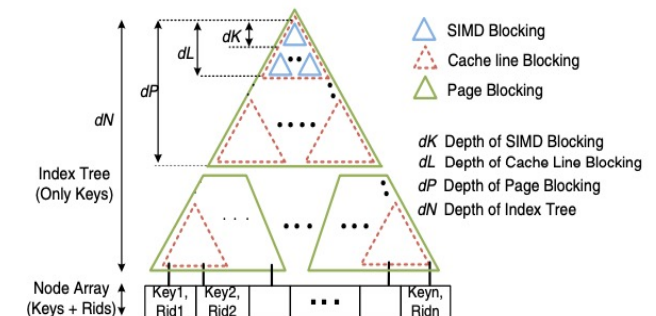    - Problems arise when there is a lot of data!



**Index (Memory)**  →  **Storage (SSD)**

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# What I've seen: Observation

- Why do various in-memory indexes exist?
  - Aren't arrays enough?
    - Problems arise when there is a lot of data!



**Index (Memory)**                    **Storage (SSD)**

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# What I've seen: Observation

- **In-memory structure: Tree-based**
  - B-Tree, B+-Tree, FAST (Fast Architecture Sensitive Tree)
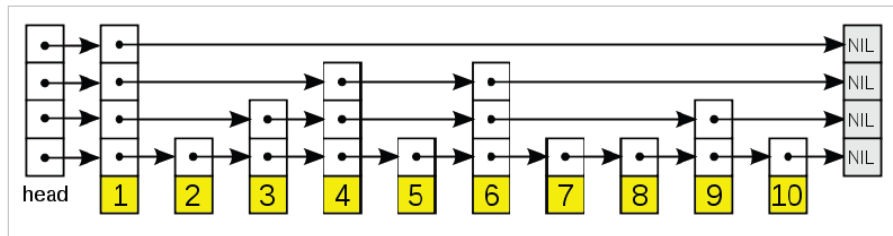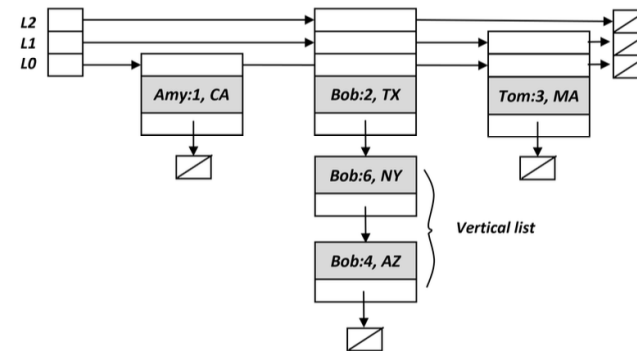  - AVL Tree, RB-Tree



**B-Tree[1]**



**B+-Tree[2]**



**FAST[3]**

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# What I've seen: Observation

- In-memory structure: List-based
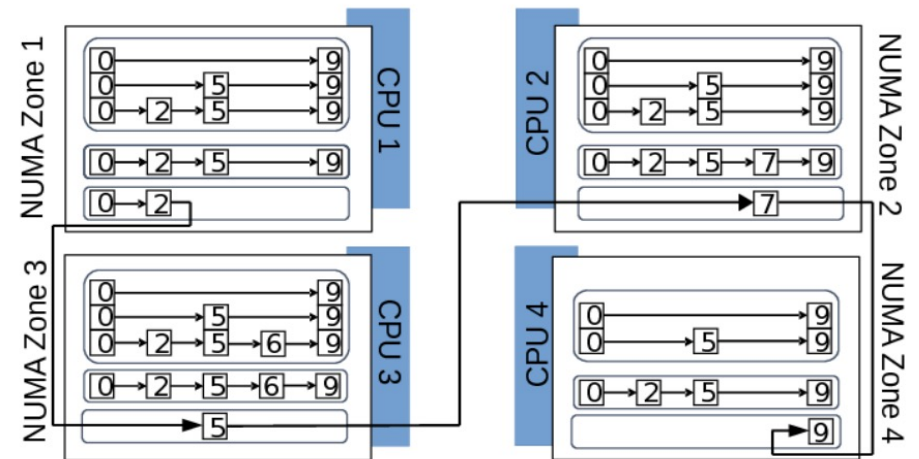  - Skiplist, JellyFish
  - No Hot Spot Skiplist, NUMA Skiplist



**Skiplist[4]**



**JellyFish[5]**

# What I've seen: Observation

- **In-memory structure: List-based**
  - Skiplist, JellyFish
  - No Hot Spot Skiplist, NUMA Skiplist



Fig. 2. The contention-friendly non blocking skip list structure

**NHS Skiplist[6]**



**NUMA Skiplist[7]**

DANKOOK UNIVERSITY

Dankook University
System Software Laboratory

# What I've seen: Observation

- In-memory structure: Trie-based

  - Radix Tree, ART (Adaptive Radix Tree)



**Radix Tree**



Fig. 1. Adaptively sized nodes in our radix tree.

**Adaptive Radix Tree[8]**

# What I've seen: Observation

- **In-memory structure: Hybrid**

  - S3, CSSL (Cache Sensitive Skip List)



S3[9]



CSSL[10]

# Study Github:
# Index Structure Journey

# Index Structure Journey

- Study Github

  - Link: https://github.com/DKU-StarLab/IndexStructureJourney

# Index Structure Journey

- **Open Source**

  - Traditional Index benchmark

    - Index-microbench: https://github.com/wangziqi2016/index-microbench

  - Learned Index benchmark

    - GRE bench: https://github.com/gre4index/GRE

# Q & A

## Thank you!

# Reference

- Reference List
  - [1] Douglas Comer, "Ubiquitous B-Tree", ACM Computing Surveys, 1979
  - [2] R. Bayer, et al. "Organization and maintenance of large ordered indices", SIGFIDET '70
  - [3] Changkyu Kim et al. "FAST: fast architecture sensitive tree search on modern CPUs and GPUs", SIGMOD '10
  - [4] William Pugh, "Skip lists: a probabilistic alternative to balanced trees", Communications of the ACM 1990
  - [5] Jeseong Yeon, et al. "JellyFish: A Fast Skip List with MVCC", Middleware '20
  - [6] Tyler Crain, et al. "No Hot Spot Non-blocking Skip List", ICDCS 2013
  - [7] Henry Daly, et al. "NUMASK: High Performance Scalable Skip List for NUMA", DISC 2018
  - [8] Viktor Leis, et al. "The adaptive radix tree: ARTful indexing for main-memory databases", ICDE 2013
  - [9] Jingtian Zhang, et al. "S3: a scalable in-memory skip-list index for key-value store", VLDB 2019
  - [10] Sprenger, et al. "Cache-Sensitive Skip List: Efficient Range Queries on Modern CPUs", Data Management on New Hardware 2016

**DANKOOK UNIVERSITY**

Dankook University
System Software Laboratory