

AirIndex: Versatile Index Tuning Through Data and Storage

Supawit Chockchowwat, Wenjie Liu, Yongjoo Park

SIGMOD '24e (Accepted 23 May 2023)

2024. 09. 11

Presentation by ZHU YONGJIE

harasho2015@dankook.ac.kr

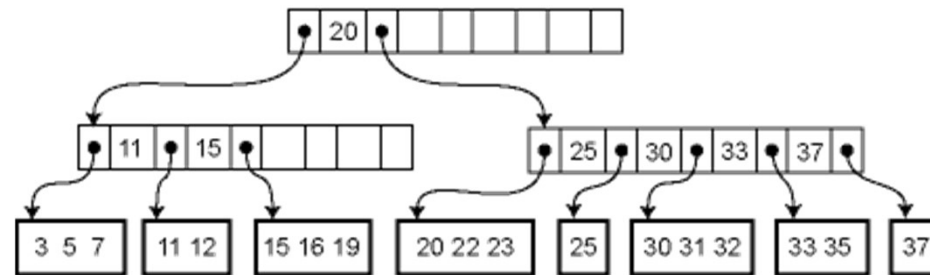
Contents

- 1. Introduction**
2. Motivation
3. AirIndex Overview
4. AirIndex-Model
5. Experiment
6. Conclusion

Introduction

- Introduction to Index Structures

Key \longrightarrow Position

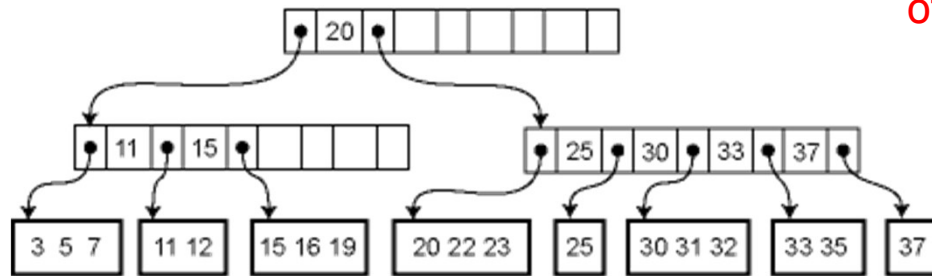


B-Tree

Introduction

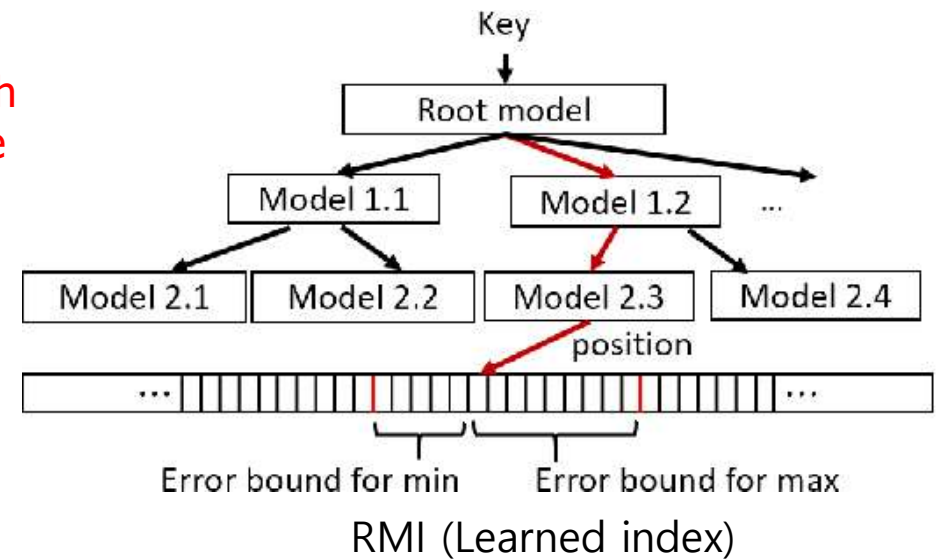
- Introduction to Index Structures

Key \longrightarrow Position



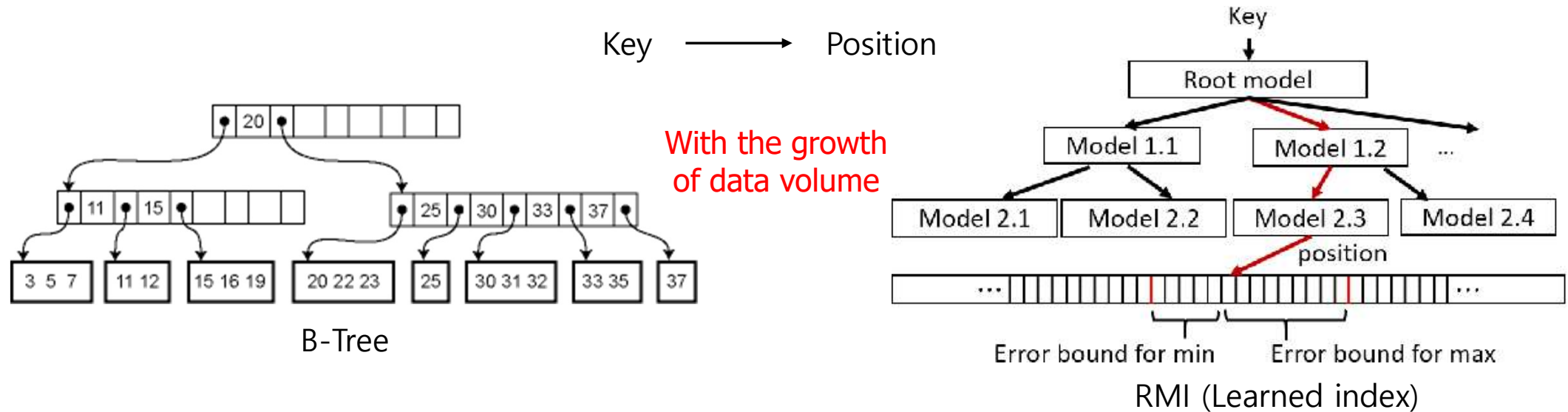
B-Tree

With the growth
of data volume

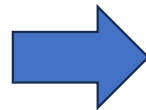


Introduction

- Introduction to Index Structures



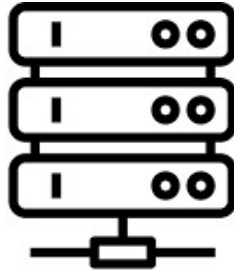
**Different Device environments
&
Different datasets**



Indexes need tuning

Introduction

- Limitations of Existing Indexes



Local SSD

Fast I/O latency
Relatively smaller bandwidth

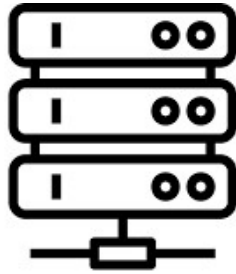


Cloud Server

High I/O latency
Larger bandwidth

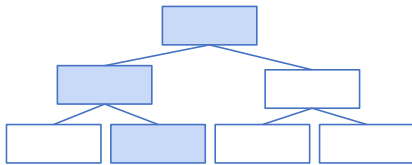
Introduction

- Limitations of Existing Indexes



Local SSD

Fast I/O latency
Relatively smaller bandwidth

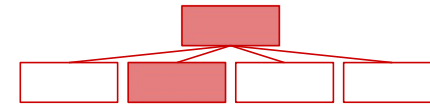


B200: B-tree with 200 child pointers



Cloud Server

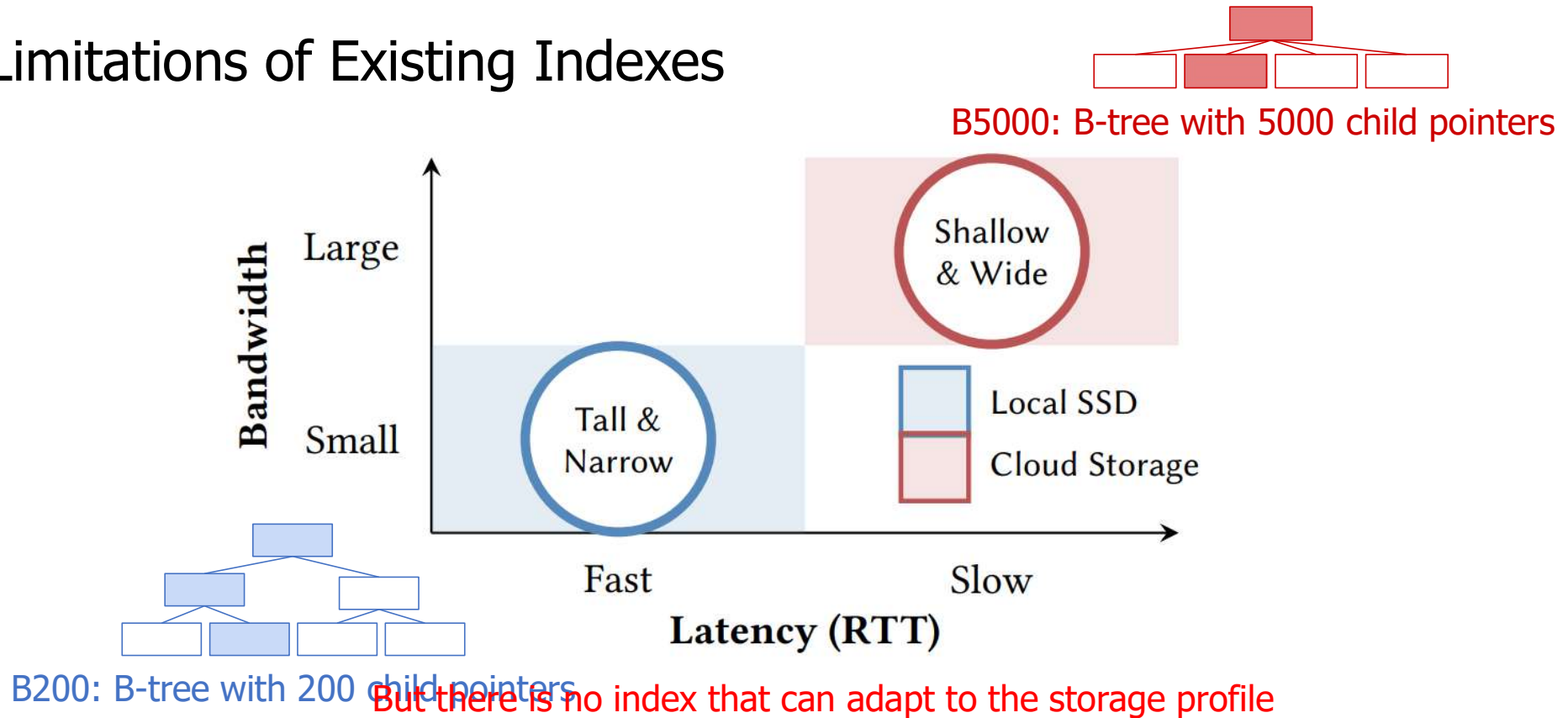
High I/O latency
Larger bandwidth



B5000: B-tree with 5000 child pointers

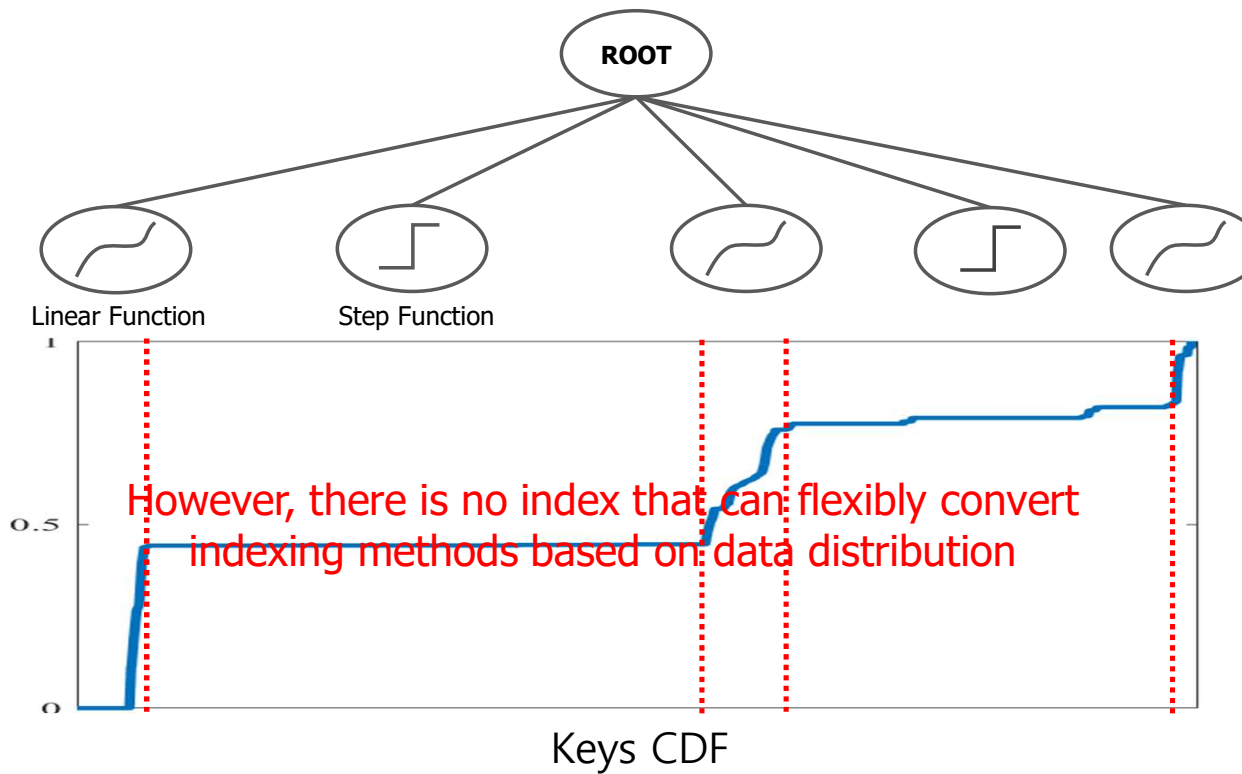
Introduction

- Limitations of Existing Indexes



Introduction

- Limitations of Existing Indexes

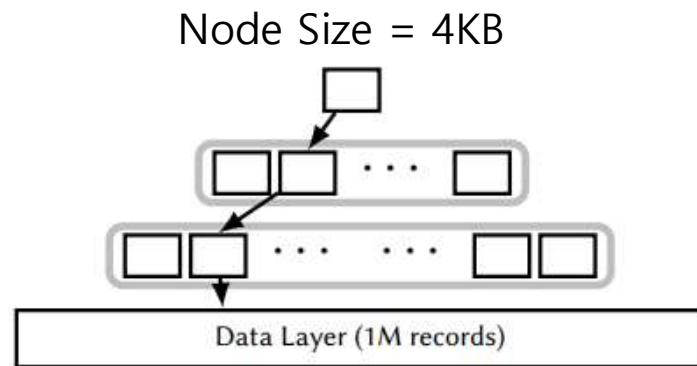


Contents

1. Introduction
- 2. Motivation**
3. AirIndex Overview
4. AirIndex-Model
5. AirTune
6. Experiment
7. Conclusion

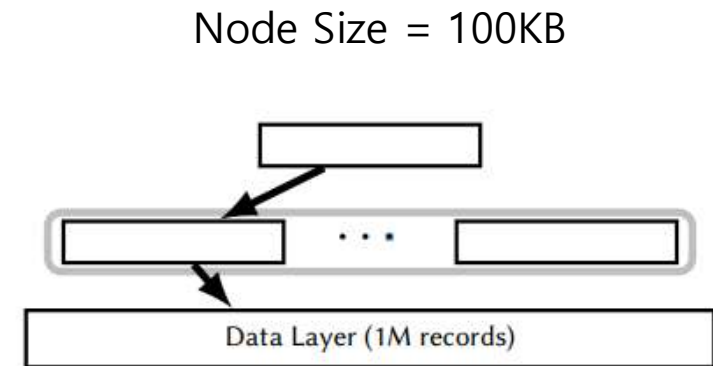
Motivation

- Why does the index need to adapt to different storage profiles



(a) B200: B-tree with 200 child pointers

$$200 * 200 * 200 = 8M > 1M$$



(b) B5000: B-tree with 5,000 child pointers

$$5000 * 5000 = 25M > 1M$$

$$(\text{data transfer time}) = (\text{latency}) + (\text{data size}) / (\text{bandwidth})$$

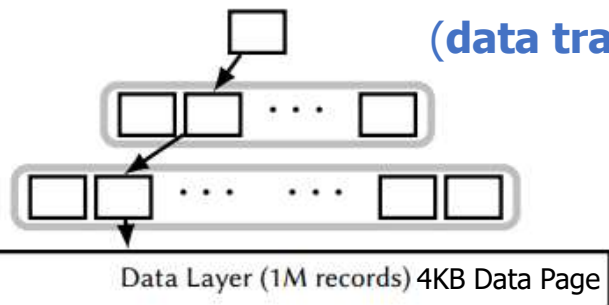
Motivation

- Why does the index need to adapt to different storage profiles

Node Size = 4KB

Node Size = 100KB

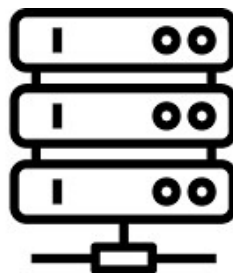
$$(\text{data transfer time}) = (\text{latency}) + (\text{data size}) / (\text{bandwidth})$$



(a) B200: B-tree with 200 child pointers

$$200 \times 200 \times 200 = 8M > 1M$$

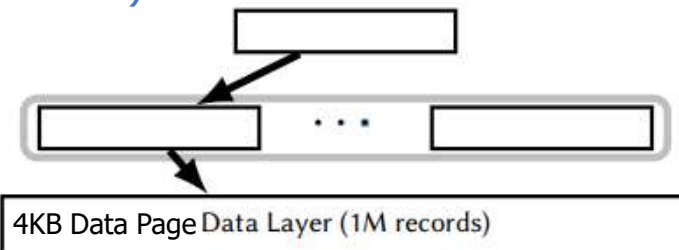
$$3 \times (100\mu s + 4KB / (1GB/s)) + (100\mu s + 4KB / (1GB/s)) = 416\mu s$$



Local SSD

Latency = 100μs

Bandwidth = 1GB/s



(b) B5000: B-tree with 5,000 child pointers

$$5000 \times 5000 = 25M > 1M$$

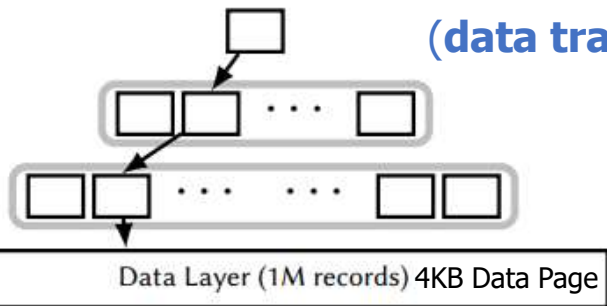
$$2 \times (100\mu s + 100KB / (1GB/s)) + (100\mu s + 4KB / (1GB/s)) = 504\mu s$$

Motivation

- Why does the index need to adapt to different storage profiles

Node Size = 4KB

$$(\text{data transfer time}) = (\text{latency}) + (\text{data size}) / (\text{bandwidth})$$



(a) B200: B-tree with 200 child pointers

$$200 \times 200 \times 200 = 8M > 1M$$

$$3 \times (100\text{ms} + 4\text{KB} / (100\text{MB/s})) + (100\text{ms} + 4\text{KB} / (100\text{MB/s})) = 400.16\text{ms}$$

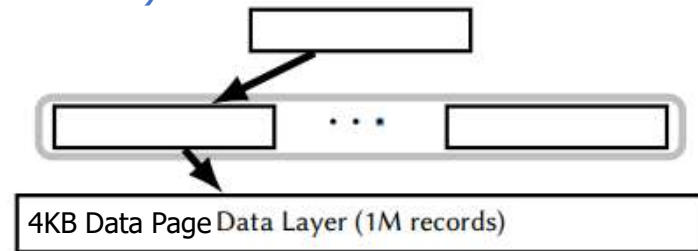


Cloud Server

Latency = 100ms

Bandwidth = 100MB/s

Node Size = 100KB



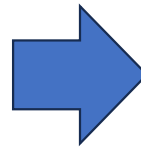
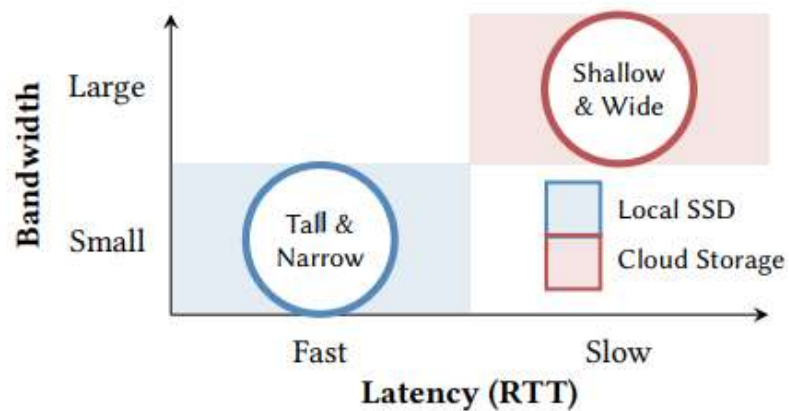
(b) B5000: B-tree with 5,000 child pointers

$$5000 \times 5000 = 25M > 1M$$

$$2 \times (100\text{ms} + 100\text{KB} / (100\text{MB/s})) + (100\text{ms} + 4\text{KB} / (100\text{MB/s})) = 302.04\text{ms}$$

Motivation

- Why does the index need to adapt to different storage profiles



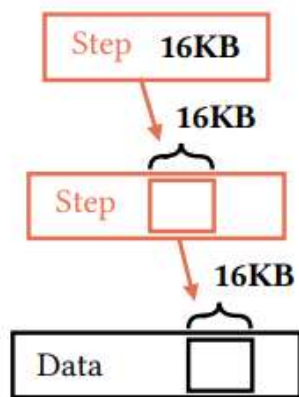
(c) Lookup Speed Comparison

In **local SSDs**, the impact of index depth is not as significant as index width, but the opposite is true in **cloud storage**

Motivation

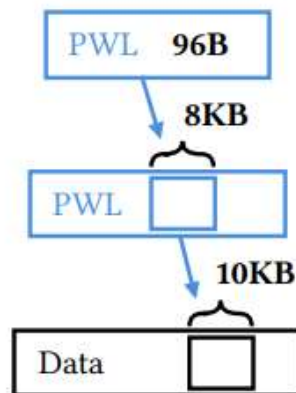
- Why does the index need to adapt to different data distributions

Gmm dataset, SSD(250 μ s latency, 175MB/s bandwidth)



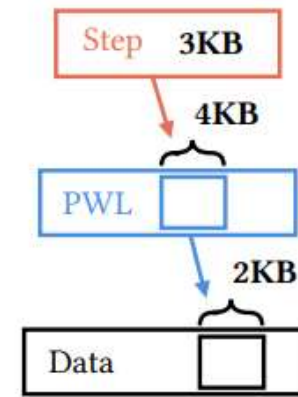
(a) Tuned B-tree

48KB



(b) Tuned PWL Index

18.096KB



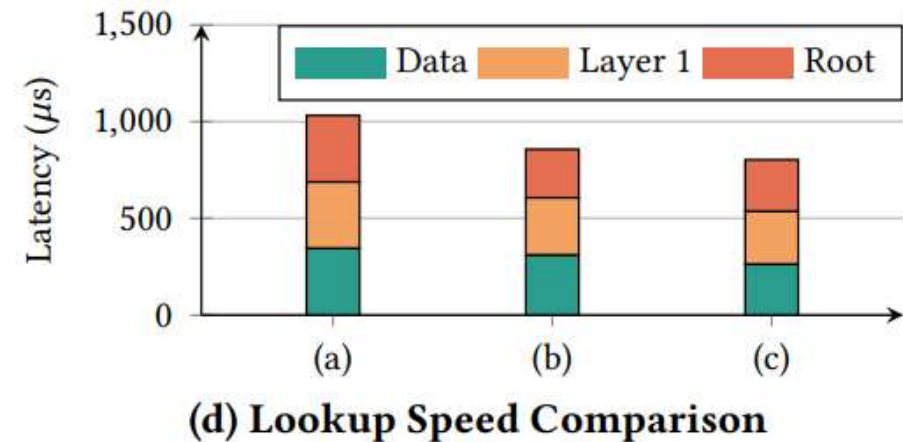
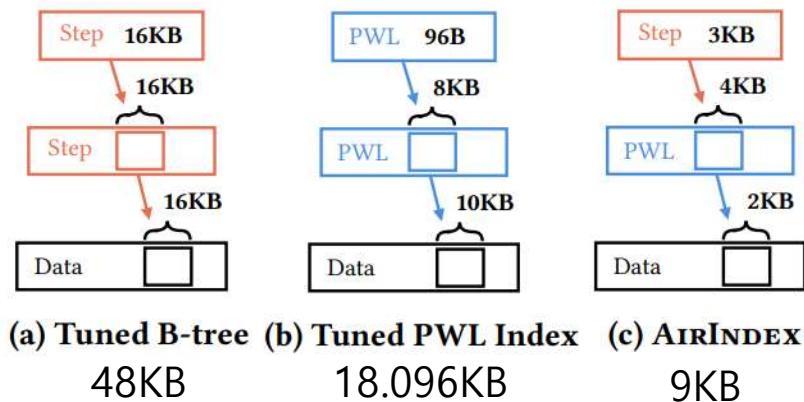
(c) AIRINDEX

9KB

Motivation

- Why does the index need to adapt to different data distributions

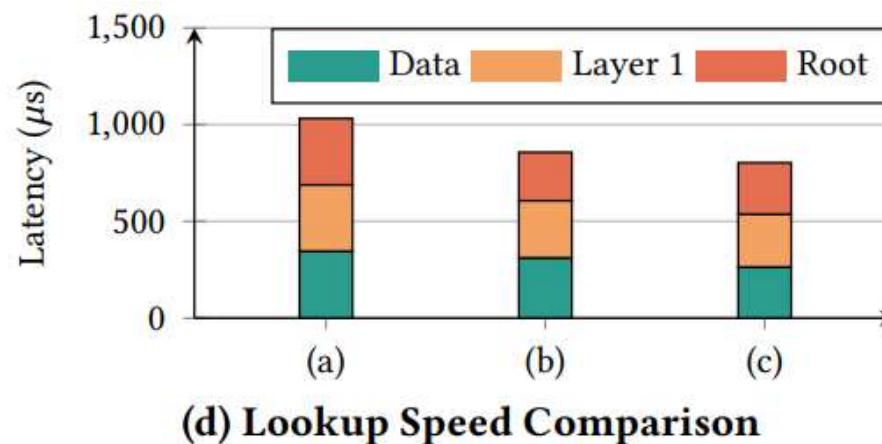
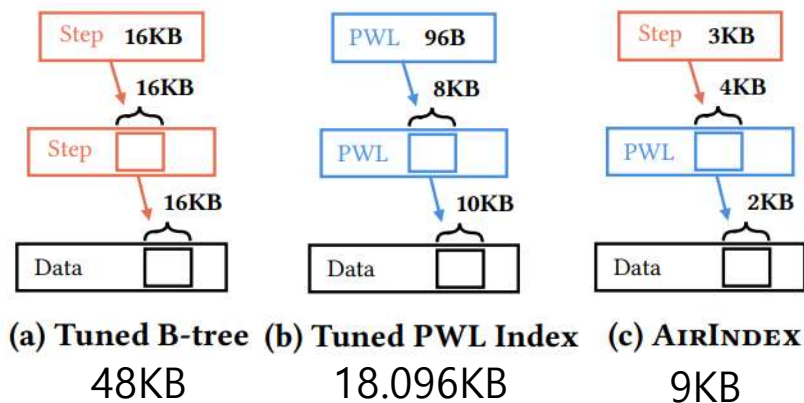
Gmm dataset, SSD(250 μ s latency, 175MB/s bandwidth)



Motivation

- Why does the index need to adapt to different data distributions

Gmm dataset, SSD(250 μ s latency, 175MB/s bandwidth)



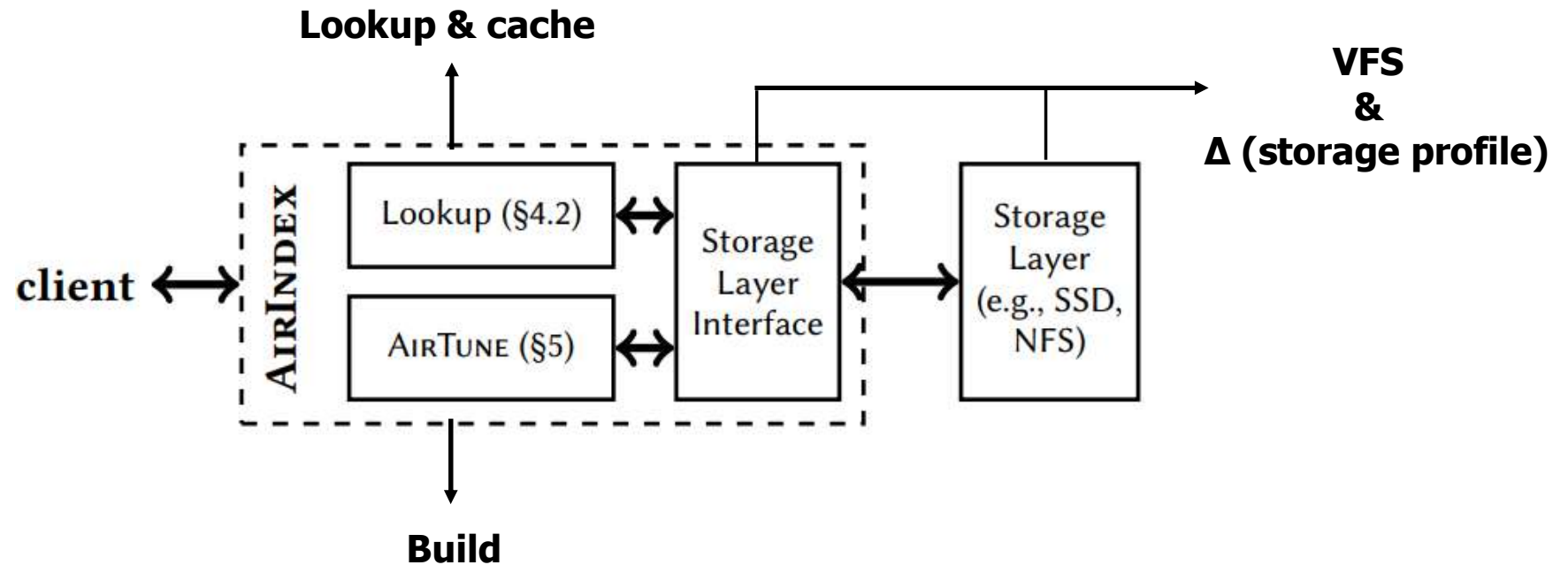
Index needs to be optimized layer by layer based on data distribution

Contents

1. Introduction
2. Motivation
- 3. AirIndex Overview**
4. AirIndex-Model
5. AirTune
6. Experiment
7. Conclusion

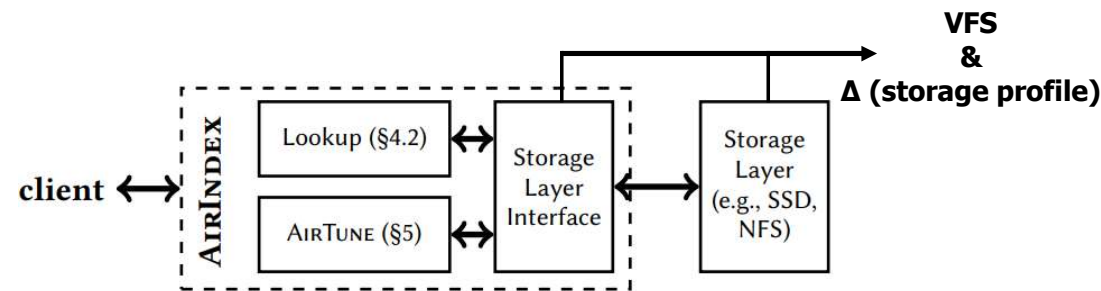
AirIndex Overview

- Architecture



AirIndex Overview

- Storage Model



$$T_{aff}(\Delta) = l + \frac{\Delta}{B}$$

l : latency ; B : bandwidth ;

Δ : Storage layer reads

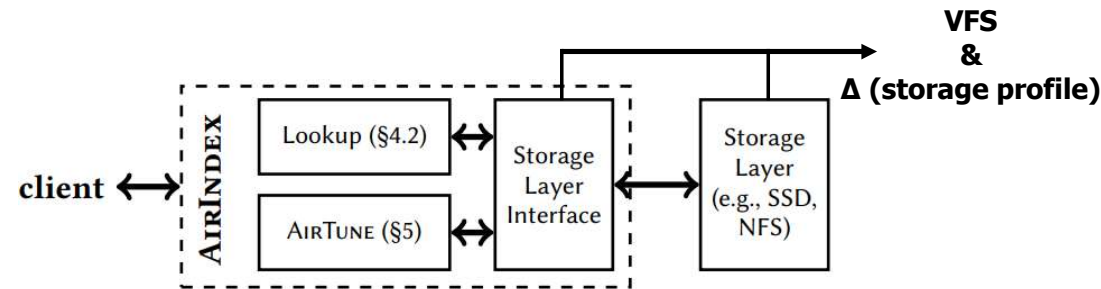
Using l to represent **storage profile**

$$T_{aff-uniform}(\Delta) = \frac{\ell_1 + \ell_0}{2} + \Delta \frac{\ln B_1 - \ln B_0}{B_1 - B_0}$$

The calculation of the expected value of $1/B$ involves logarithms

AirIndex Overview

- Storage Model



$$T_{aff}(\Delta) = l + \frac{\Delta}{B}$$

l : latency ; B : bandwidth ;

Δ : Storage layer reads

Using **latency** to represent **storage profile**

$$T_{\text{aff-uniform}}(\Delta) = \frac{\ell_1 + \ell_0}{2} + \Delta \frac{\ln B_1 - \ln B_0}{B_1 - B_0}$$

The calculation of the expected value of **$1/B$** involves logarithms

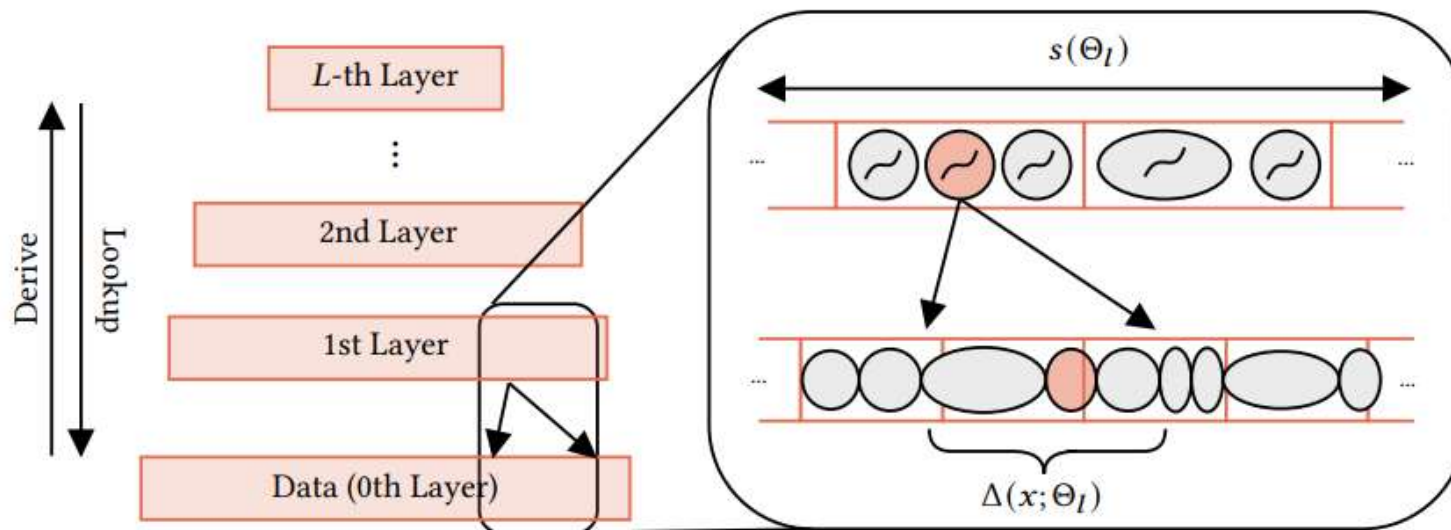
It is more convenient for AirIndex to self tune without considering too many low relational variables

Contents

1. Introduction
2. Motivation
3. AirIndex Overview
- 4. AirIndex-Model**
5. AirTune
6. Experiment
7. Conclusion

AirIndex-Model

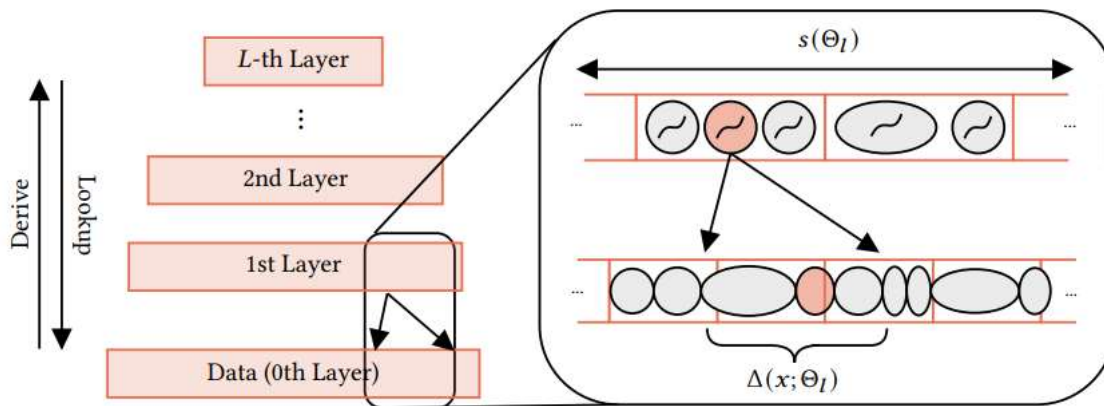
- Hierarchical Indexes



Multi layer & Single layer multiple nodes with multi type nodes

AirIndex-Model

■ Query Process



Algorithm 1: AIRINDEX Query Process, $\text{Lookup}(x; \hat{y}_L, L)$

Input: Query key x , root position \hat{y}_L , number of layers L

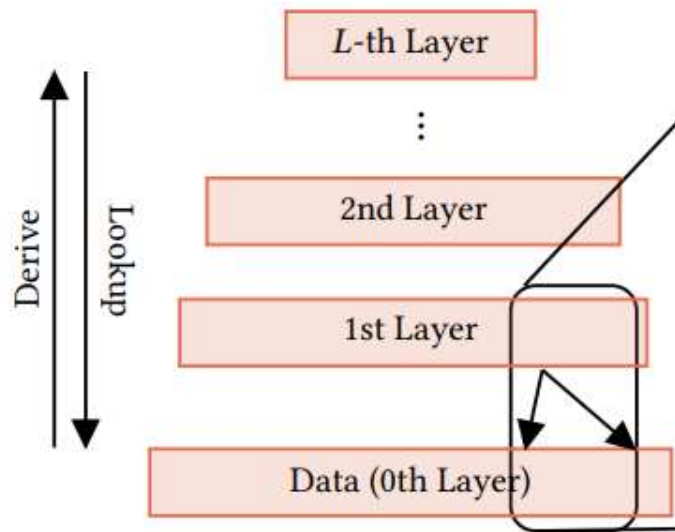
Output: Relevant key-value (x, v)

```

1 for  $l$  from  $L$  to 0 do
2    $\{(x[i], v[i])\}_i \leftarrow \text{Read}(\hat{y}_l(x))$            // Storage access
3    $v_l \leftarrow \text{Search}(x, \{(x[i], v[i])\}_i)$ 
4   if In index layer,  $l \geq 1$  then
5      $\hat{y}_{l-1} \leftarrow \text{ReconstructNode}(v_l)$ 
6 return  $(x, v) = (x, v_0)$ 
    
```

AirIndex-Model

- Latency Under Storage Model



Design variables Θ :

L	Number of layers
NodeType_l	Node type in layer $l \in \{1, \dots, L\}$
n_l	Number of nodes in layer $l \in \{1, \dots, L\}$
$\theta_{l,i}$	Parameters of the i -th node in layer $l, i \in \{1, \dots, n_l\}$

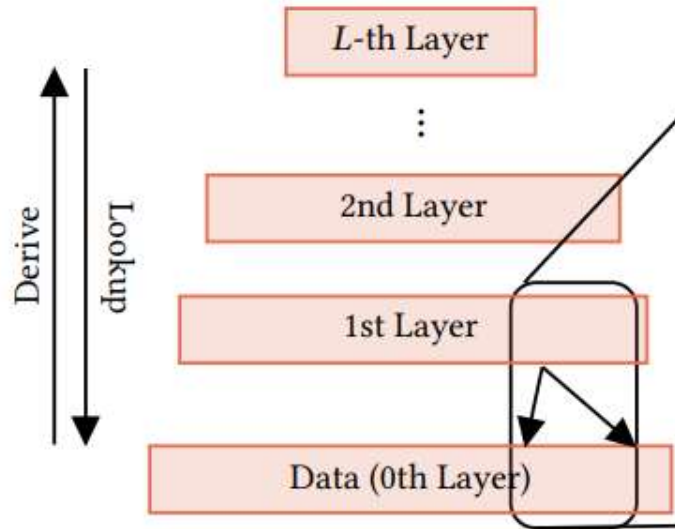
Fixed variables:

T	Storage profile
-----	-----------------

$$\mathcal{L}_{SM}(x; \Theta, T) = T(s(\Theta_L)) + \sum_{l=1}^L T(\Delta(x; \Theta_l))$$

AirIndex-Model

Latency Under Storage Model



Design variables Θ :

L	Number of layers
NodeType_l	Node type in layer $l \in \{1, \dots, L\}$
n_l	Number of nodes in layer $l \in \{1, \dots, L\}$
$\theta_{l,i}$	Parameters of the i -th node in layer $l, i \in \{1, \dots, n_l\}$

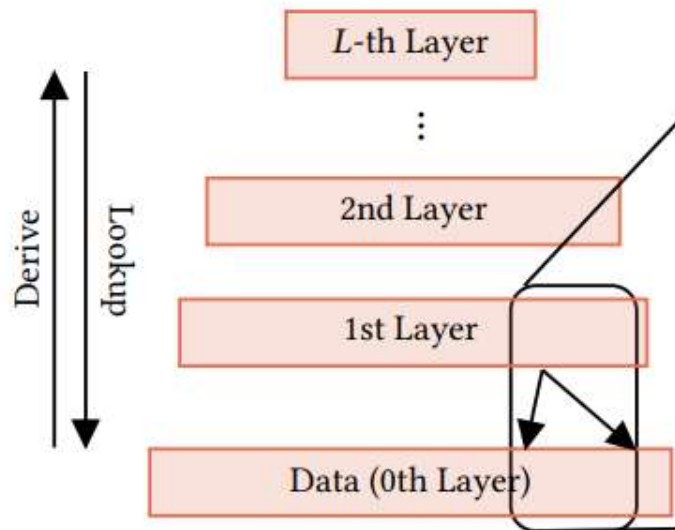
Fixed variables:

T	Storage profile
\mathcal{X}	Query key distribution
D	Key-position collection $D = \{(x_i, y_i)\}_{i=1}^n$

$$\mathcal{L}_{SM}(x; \Theta, T) = \mathbb{E}_{x \sim \mathcal{X}} \left[T(s(\Theta_L)) + \sum_{l=1}^L T(\Delta(x; \Theta_l)) \right]$$

AirIndex-Model

Latency Under Storage Model



Design variables Θ :

L	Number of layers
NodeType_l	Node type in layer $l \in \{1, \dots, L\}$
n_l	Number of nodes in layer $l \in \{1, \dots, L\}$
$\theta_{l,i}$	Parameters of the i -th node in layer $l, i \in \{1, \dots, n_l\}$

Fixed variables:

T	Storage profile
\mathcal{X}	Query key distribution
D	Key-position collection $D = \{(x_i, y_i)\}_{i=1}^n$

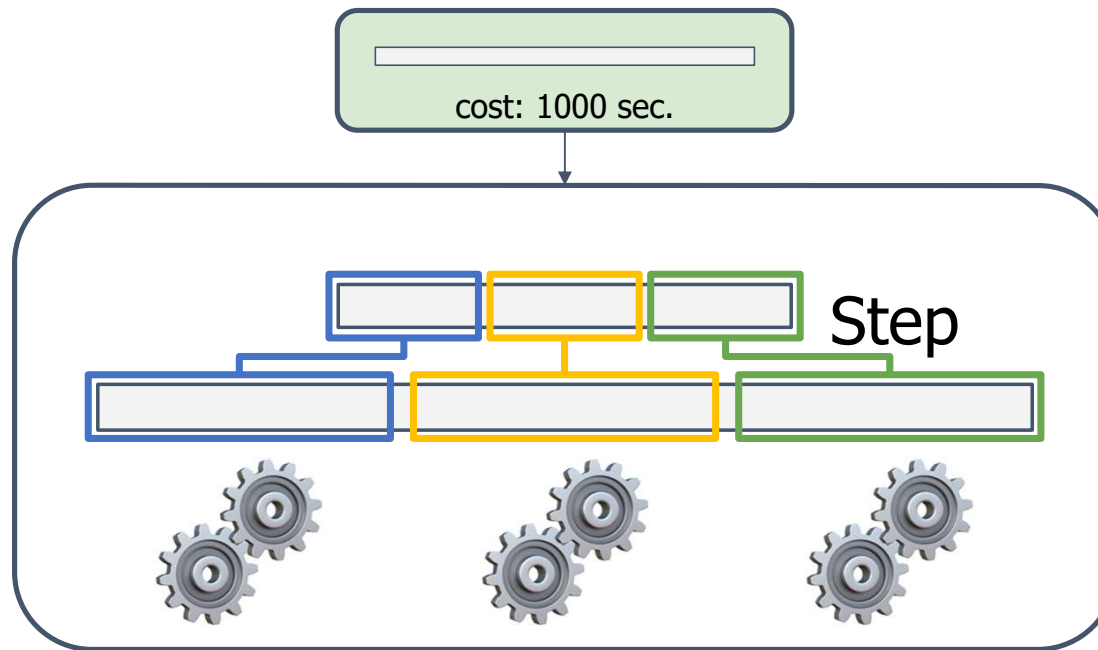
$$\Theta^* = \arg \min_{\Theta} \mathbb{E}_{x \sim \mathcal{X}} \left[T(s(\Theta_L)) + \sum_{l=1}^L T(\Delta(x; \Theta_l)) \right]$$

Contents

1. Introduction
2. Motivation
3. AirIndex Overview
4. AirIndex-Model
- 5. AirTune**
6. Experiment
7. Conclusion

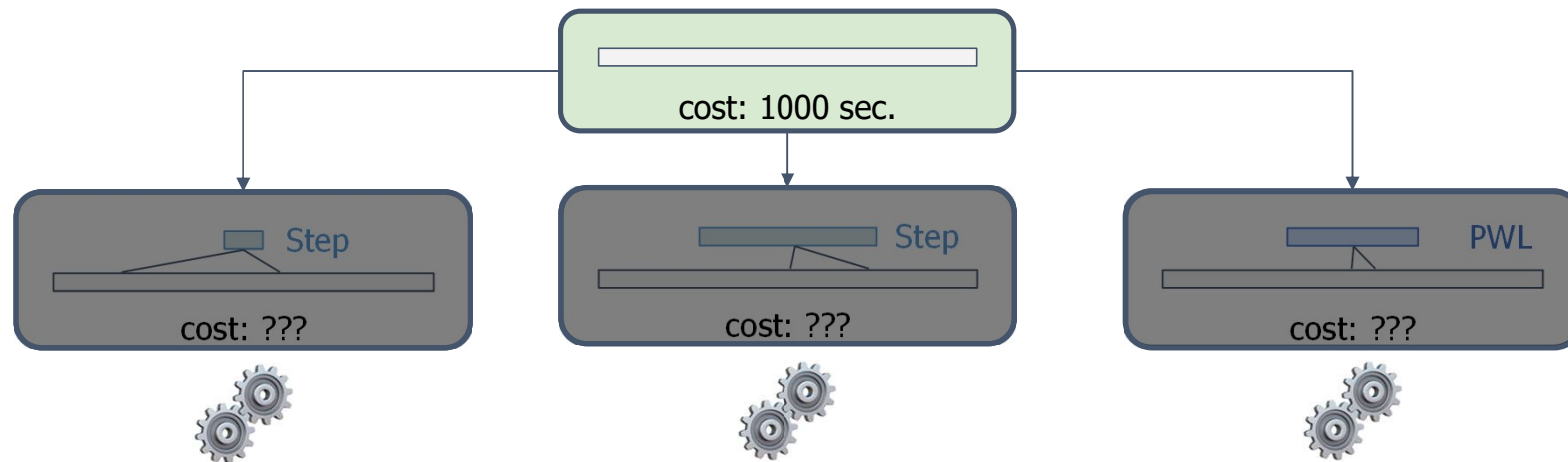
AirTune

- Guided Graph Search



AirTune

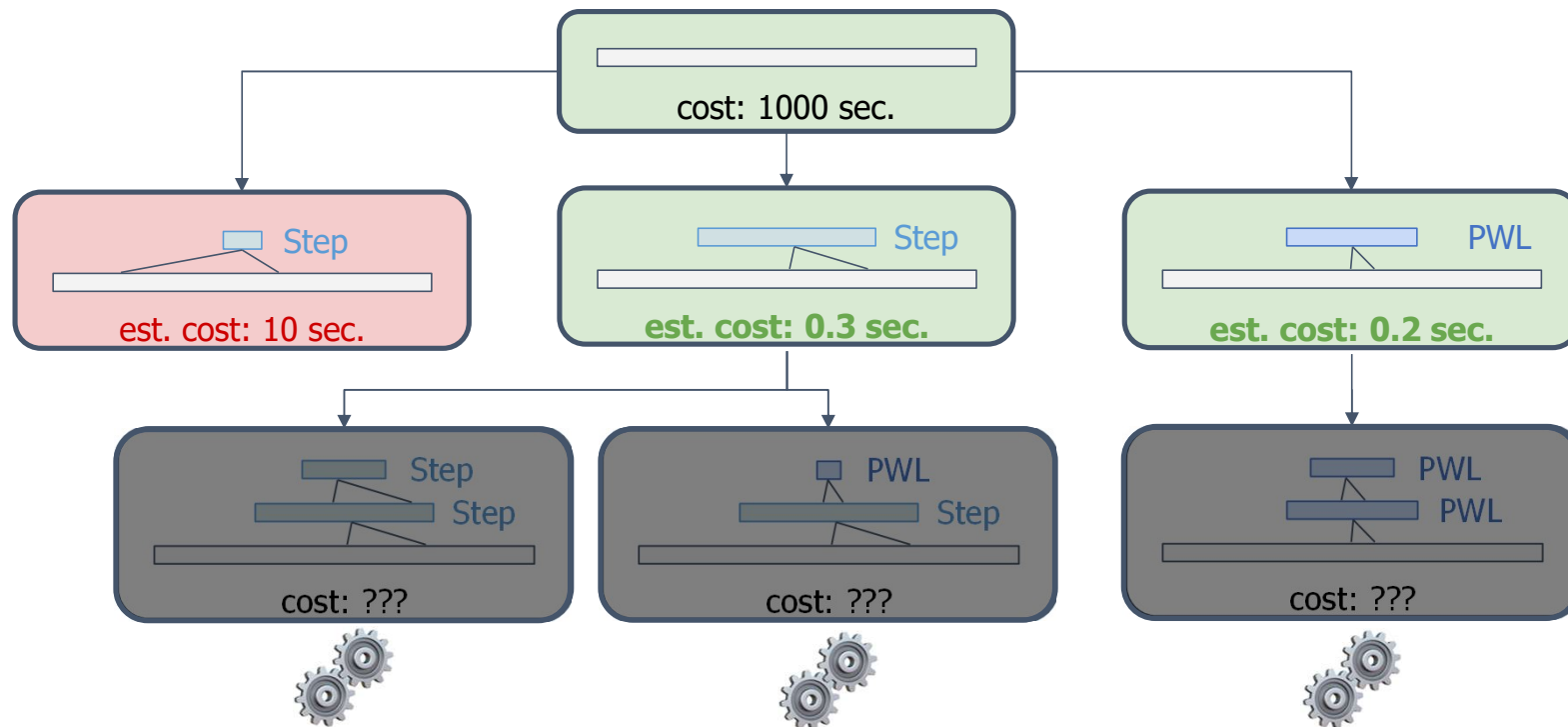
- Guided Graph Search



$$\Theta^* = \arg \min_{\Theta} \mathbb{E}_{x \sim \mathcal{X}} \left[T(s(\Theta_L)) + \sum_{l=1}^L T(\Delta(x; \Theta_l)) \right]$$

AirTune

- Guided Graph Search



Contents

1. Introduction
2. Motivation
3. AirIndex Overview
4. AirIndex-Model
5. AirTune
- 6. Experiment**
7. Conclusion

Experiment

■ Setup

- System Environment :

- Azure cloud platform (8 vCPUs, 32 GiB RAM)
 - NFS: Azure network file system
 - SSD: Azure Premium SSD (256 GiB, 2300 IOPS, 150 MBps, read/write host caching).
 - HDD: Azure Standard HDD (1024 GiB, 500 IOPS, 60 MBps, no host caching)

- Baselines :

- LMDB: B-tree database
- RMI, PGM, ALEX/APEX, PLEX: Learned indexes
- Data Calculator: Index tuner
- B-TREE: AirIndex's tuned B-Tree

Experiment

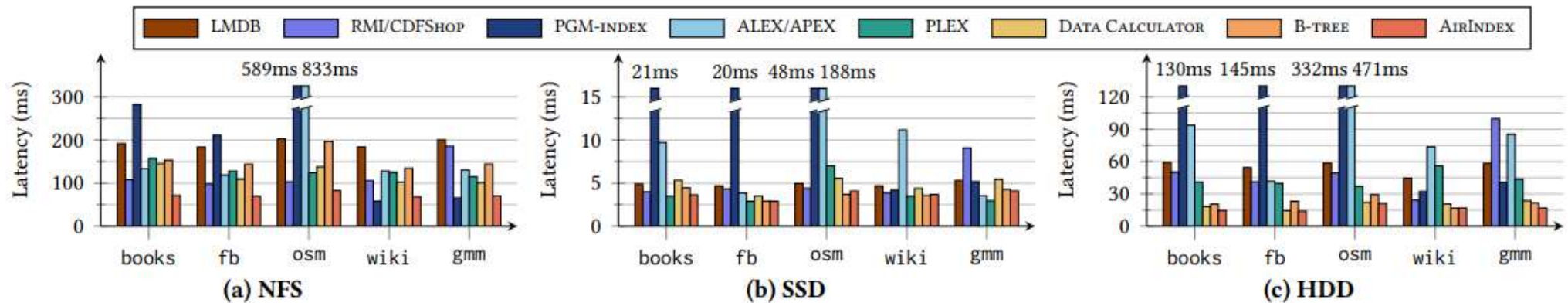
- Setup

- Dataset

- BOOKS(800M)
 - FB(200M)
 - OSM(800M)
 - WIKI(200M)
 - GMM(Gaussian mixture model of 100 normal distribution clusters over 800M keys)

Experiment

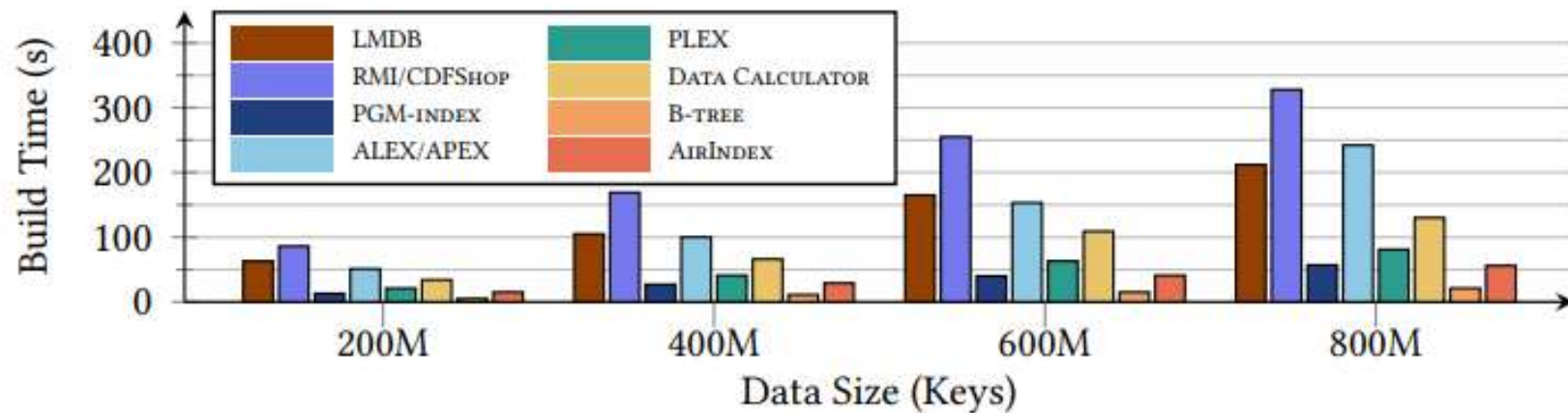
- Faster End-to-end Lookup Speed – Cold-state Latency



AirIndex is consistently one of the fastest methods at searching the first query, across datasets and storage

Experiment

- Competitive Build Time

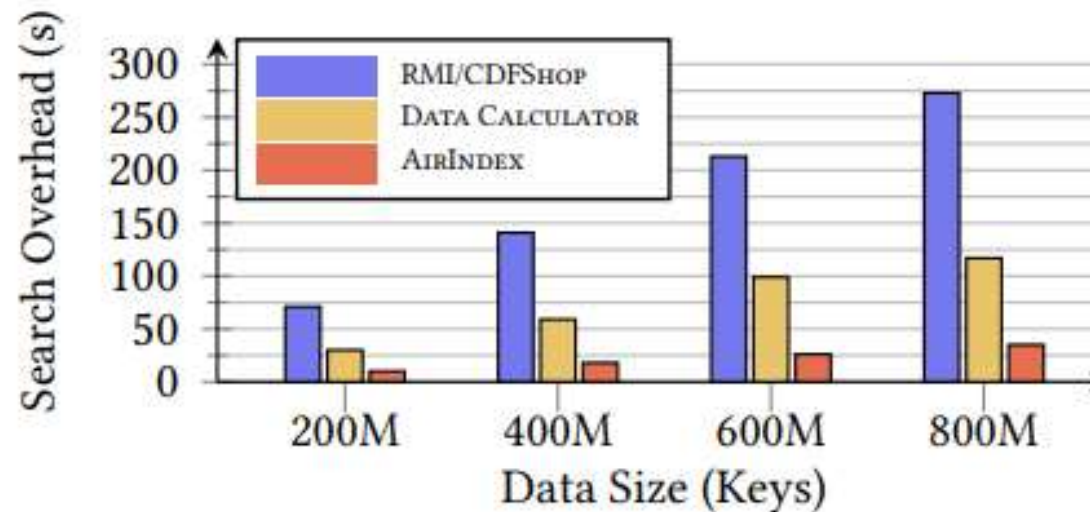


(a) Total Time to Build Index

Despite the default **tuning overhead** of AirIndex's **Airturn**, it still has excellent build time performance

Experiment

- Competitive Build Time



(b) Search Overhead Time

The **tuning cost** of **Airindex** is the **lowest** among existing tuning methods

Contents

1. Introduction
2. Motivation
3. AirIndex Overview
4. AirIndex-Model
5. AirTune
6. Experiment
- 7. Conclusion**

Conclusion

- AirIndex

- AirIndex is the **first** to build high-speed hierarchical indexes by learning data and I/O characteristics
- AirIndex uses a specially constructed **graph search** method (**AirTune**) to explore the search optimization space
- Compared to indexes that have not undergone specific optimization, AirIndex's **data** and **I/O aware optimization** can achieve significantly faster search speeds.

Thank you