Supawit Chockchowwat, Wenjie Liu, Yongjoo Park {supawit2,wenjie3,yongjoo}@illinois.edu CreateLab, University of Illinois at Urbana-Champaign

ABSTRACT

The end-to-end lookup latency of a hierarchical index—such as a B-tree or a learned index—is determined by its structure such as the number of layers, the kinds of branching functions appearing in each layer, the amount of data we must fetch from layers, etc. Our primary observation is that by optimizing those structural parameters (or *designs*) specifically to a target system's I/O characteristics (e.g., latency, bandwidth), we can offer a faster lookup compared to the ones that are not optimized. Can we develop a systematic method for finding those optimal design parameters? Ideally, the method must have the potential to generate almost any existing index or a novel combination of them for the fastest possible lookup.

In this work, we present a new data-and-I/O-aware index builder (called AIRINDEX) that can find high-speed hierarchical index designs in a principled way. Specifically, AIRINDEX minimizes an objective expressing the end-to-end latency in terms of various *designs*—the number of layers, types of layers, and more—for given data and a *storage profile*, using a graph-based optimization method purpose-built to address the computational challenges rising from the inter-dependencies among index layers and the exponentially many candidate parameters in a large search space. Our evaluations confirm that AIRINDEX can find optimal index designs, build them within the times comparable to existing methods, and deliver up to 4.1× faster lookup than a lightweight B-tree library (LMDB), 3.3×–46.3× faster than state-of-the-art learned indexes (RMI/CDFSHOP, PGM-INDEX, ALEX/APEX, PLEX), and 2.0× faster than DATA CAL-CULATOR's suggestion on various dataset and storage settings.

ACM Reference Format:

Supawit Chockchowwat, Wenjie Liu, Yongjoo Park. 2024. AIRINDEX: Versatile Index Tuning Through Data and Storage. In *Proceedings of ACM SIG-MOD/PODS International Conference on Management of Data (SIGMOD '24)*. ACM, New York, NY, USA, 18 pages. https://doi.org/XXXXXXXXXXXXXXXXX

1 INTRODUCTION

Hierarchical indexes (e.g., binary search tree, B-tree) allow us to quickly locate a relevant item by fetching (only) a small fraction of data inside each index layer. The B-tree [17, 19, 25, 72] has been a conventional choice for many data systems such as PostgreSQL [56], MySQL [8], ZLog [7], BTRFS [58], etc. More recently, it has been shown that by using *learned models* (i.e., regression functions offering approximate pointers) in place of the (exact) child pointers

SIGMOD '24, June 11-16, 2024, Santiago, Chile (Accepted 23 May 2023)

© 2024 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

https://doi.org/XXXXXXXXXXXXXXXX



Figure 1: Expected optimal structures by I/O characteristics. AIRINDEX finds the optimal structure in a large design space.

inside a B-tree's internal nodes, we can reduce the amount of data we need to fetch for each layer, lowering the overall lookup latency [31, 33, 41, 43, 50]. In general, the amount of data indicated by those exact/approximate branching pointers—along with the number of layers—determine the overall lookup latency of an index.

Despite their effectiveness, (most) existing indexes have a common limitation, making them often suboptimal for novel system environments with different I/O characteristics. That is, due to their nearly fixed index structures (e.g., the number of layers, the types of layers), existing indexes cannot deliver as fast performance as the ones specifically designed in consideration of the data access cost of a target system environment. For instance, in Figure 1, if an index is maintained directly on remote storage with relatively high I/O latency (e.g., in cloud systems like Amazon Aurora [67] and Delos [16]), we might be able to achieve faster lookup speeds by creating a wider/shallower index (e.g., larger fanout in B-trees, more accurate models in learned indexes), mitigating the impact of high I/O latency with fewer data fetches required for index traversal. In contrast, if an index is maintained on a local SSD (e.g., SQLite [63]) with fast I/O latency (or relatively smaller bandwidth), we can create a narrower/taller index (e.g., smaller fanout in B-trees, less accurate models in learned indexes) for faster lookup. In other words, depending on a different system environment, an optimal index structure may vary, causing significant performance gaps between fixed index structures and the ones adapting to target environments (§2). This observation closely resembles the one made by Gray and Graefe [36] for determining an optimal B-tree page size in relation to page access cost. However, its utility-based approach is specific to conventional indexes with exact child pointers, meaning it cannot be generalized to a wider class of indexes with learned models.

In this work, we tackle this limitation with *AIRINDEX*, a general index-building framework that can efficiently find an (optimal) lowlatency index structure in consideration of profiled I/O characteristics as well as data distribution. The core difference of AIRINDEX is that it can balance the properties of all the layers simultaneously (including their total count) to optimize the expected (cache-aware) end-to-end latency by solving a mathematical optimization problem. During the optimization, AIRINDEX considers many different design choices (e.g., the number of layers, fanout, model types, model

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Method	Approx. Pointer?	I/O Aware?	Novelty	Weakness/Difference
B-trees [17, 25]	X	×	First B-tree proposal	No search for optimal fanout
RMI [45]	 ✓ 	X	First approximate branching	No hyperparameter tuning, fixed two-layer structure
ALEX/APEX [30]	 ✓ 	X	Updatable with learned models	No optimization for end-to-end latency
5minRules [36]	X	~	Heuristic for B-tree page size	Restricted tuning, fixed branching across layers
CDFShop [53]	 ✓ 	X	Searches Pareto efficient RMIs	Inconclusive tuning, not optimize for latency
DATA CALCULATOR [38]	A	~	Evaluates end-to-end performance	Inefficient tuning, restricted branching functions
Ours (AIRINDEX)	~	~	I/O-aware exact/approx. layers	Optimization focuses on lookup than updates

Table 1: Summary of existing work. ▲/X indicates limited/no support.

size/accuracy), often producing a heterogeneous index consisting of different types of layers: an index may have a B-tree-like design for one layer and a model-based design for another layer. While the primary contribution of AIRINDEX is demonstrating significant improvements in lookup speeds enabled by its principled search technique, its approach can easily integrate with existing orthogonal work (e.g., ALEX [31]) to allow insert/delete without complete reconstruction. To our knowledge, *no previous work has formally studied an efficient search technique for optimizing the entire index structure in consideration of I/O performance and has evaluated its practical performance benefits with comprehensive experiments.*

Challenge. The core challenge in finding an optimal index structure is the lack of an efficient mechanism that can compare the quality of exponentially many candidate structures. For example, we need to consider indexes with different heights (i.e., 1, 2, ...); for an index of height *L*, there are *L* layers to build; and for each layer, we need to compare various branching functions (e.g., B-tree-like child pointers, regression-based learned models). Moreover, the design decision we make for layers are dependent on one another; that is, if we alter the design of a layer, it affects the other components that depend on this layer (i.e., all the other layers on top of it) due to the inherent nature of hierarchical indexes, complicating the search process. Finally, we must be able to evaluate the *goodness* of each index for (arbitrary) target storage media.

Our Approach. AIRINDEX can be understood as a search process (AIRTUNE) inside a high-dimensional design space consisting of exponentially many feasible candidate structures. Each (entire) candidate structure is a model: We mathematically represent each candidate structure using a unified index model (called AIRINDEX-MODEL), a high-level abstraction that expresses diverse hierarchical indexes parametrically $\mathcal{L}(\Theta)$, where $\Theta = (\Theta_1, \dots, \Theta_L)$ are layer-specific design parameters (such as fanout, model accuracy/size), and $\mathcal{L}(\Theta)$ is the (cache-aware) end-to-end lookup latency (i.e., cumulative time for traversing the index). That is, given a parameter set Θ , there exists an associated physical index; and the index takes $\mathcal{L}(\Theta)$ time for finding a relevant item via traversal (thus, the lower the better). Search: Our key contribution is an efficient search algorithm for finding an optimal design Θ^* at which $\mathcal{L}(\Theta)$ has the smallest value. Since the design includes both numeric and categorical values, gradient-descent-like optimization algorithms are inapplicable. We devise a novel algorithm by translating our problem into a graph traversal; that is, each design represents a vertex (or state) in the search space and an edge is drawn from Θ_A to Θ_B if Θ_B is immediately reachable from Θ_A by

stacking another layer on top. While our algorithm involves actual index construction—because for some types of layers, their quality depends on data distribution requiring actual construction—our algorithm can quickly find Θ^* with bounded time by avoiding visits to the (low-quality) states unlikely to reach an optimal one.

Orthogonal Work. Database research has a long history in index design such as (1) leveraging the properties of emerging hardware, (2) employing machine learning for compact layers, and (3) targetting for different data layouts. We summarize more closely related work in Table 1 while providing comprehensive discussions in §8. Our work is largely orthogonal to those existing work because we do not propose a new type of index per se; instead, we study how to combine existing techniques under a unified framework to build an index specifically tuned for target I/O characteristics.

Contribution. In this work, we make the following contributions:

- We illustrate the significance of optimizing indexes with both data and storage (§2).
- (2) We formulate the index optimization as a search for the optimal design parameters of a *unified index model* (§4).
- (3) We design an efficient graph-based search method (§5).
- (4) We empirically study AIRINDEX with various datasets and storage and show that AIRINDEX can offer up to 2.0×-46.3× faster lookup speed compared to state-of-the-art methods (§7).

Finally, §8 discusses related work, and §9 concludes this work.

2 MOTIVATION

We present why we need a new index builder that considers the end-to-end lookup latency (§§ 2.1 to 2.3).

2.1 Need for I/O-Aware Optimization

In this section, we motivate the need for environment-specific index optimization with concrete examples showing there is no single dominant index structure (e.g., B-trees with fixed fanout) that can offer superior performance for all system environments.

Example. We have two candidate B-tree structures: B200 and B5000 (Note: special cases of AIRINDEX-MODEL). B200 consists of 4 KB nodes, each with 200 fanout. B5000 consists of 100 KB nodes, each node with 5,000 fanout. Both B200 and B5000 index the same dataset with one million distinct keys, stored in 4 KB pages.

To index the dataset, B200 needs three layers (because the third layer can hold up to $200^3 = 8M$ pointers). Likewise, B5000 needs two layers (because the second layer can hold up to $5000^2 = 25M$ pointers). Note that while B5000 is shallower than B200, fetching

SIGMOD '24, June 11-16, 2024, Santiago, Chile (Accepted 23 May 2023)



Figure 2: Need for I/O-aware Optimization. Depending on system environments (SSD and CloudStorage), employing different B-tree designs (B200 and B5000) achieve higher expected performance in end-to-end lookup latency.



Figure 3: Indexes built for gmm dataset stored on SSD ($250\mu s$ latency, 175MB/s bandwidth) and their estimated latency. Bold numbers show the average read sizes of components in a query: root layer, partial first index layer, and partial data layer.

each node of B5000 takes longer because its page size is 25× larger. Figures 2a and 2b depict these structures.

Interestingly, neither of these two indexes (i.e., B200 and B5000) is superior to the other, if we compare their lookup latencies based on a widely used formula: (data transfer time) = (latency) + (datasize) / (bandwidth); that is, (1) B200 offers higher performance than B5000 if we store data on SSD having 100 μ s latency and 1 GB/s bandwidth, and that (2) B5000 offers higher performance than B200 if we store data on CloudStorage having 100 ms latency and 100 MB/s bandwidth. Specifically, in SSD, B5000 is 21% slower than B200 because B200 needs to retrieve 3 nodes and 1 data page, taking 416 μ s (= 3 × (100 μ s + 4KB / (1GB/s)) + (100 μ s + 4KB / (1GB/s))) while B5000 needs to retrieve 2 nodes and 1 data page, taking 504 $\mu s (= 2 \times (100 \mu s + 100 \text{KB} / (1 \text{GB/s})) + (100 \mu s + 4 \text{KB} / (1 \text{GB/s})))$. In contrast, in CloudStorage, B200 is 32% slower than B5000 with B200 taking 400.16 ms and B5000 taking 302.04 ms. Figure 2c summarizes this relative performance strength. For each environment, the figure reports the relative difference in end-to-end lookup time. This shows that different index structures offer higher lookup performance, depending on the storage device.

2.2 Need for Layer-Wise Optimization

In this section, we explain why we need to consider different types of branching functions for different layers. That is, a homogeneous index—consisting of the same type of layers—may offer poorer performance even with careful tuning in comparison to a tuned heterogeneous index—consisting of different types of layers.

General Lookup Process. In general, we can consider a lookup process with a hierarchical index as a series of data fetch operations that proceed as follows: the root (or the *i*-th) layer is fetched, based on which we determine what data we need to fetch in the next (or the (i - 1)-th) layer. This iteration repeats until we reach the data layer. This means that in each layer, we can use any type of monotonically increasing function (with respect to search keys)

that can tell us what data we need to fetch in the subsequent layer. For example, a B-tree layer has a property such that for all the keys between two adjacent separators, we get the same child pointer (or the same range of data we need to fetch), which can be expressed as a step function that *jumps* at separators. This gives rise to our *unified index model*, AIRINDEX-MODEL, allowing different types of branching functions in different layers (§4).

Formally, we express such a design space with Θ_l for *l*-th layer, and $\Theta = (\Theta_1, \dots, \Theta_L)$ describes an entire index design. The following examples demonstrate that by allowing $\Theta_i \neq \Theta_j$ for $i \neq j$, we can construct an index with lower end-to-end latency.

Concrete Example. Suppose two types of branching functions step functions (Step) appearing in B-trees and piece-wise linear functions (PWL) employed in RMI [45]. Formally, a step function is a *p*-piece constant function that $\hat{y}_{step}(x) = [b_i, b_{i+1})$, while a PWL is a *p*-piece band function (widen linear functions) $\hat{y}_{PWL}(x) =$ $[m_i x + c_i - \delta_i, m_i x + c_i + \delta_i)$ for $x \in [a_i, a_{i+1})$. We present a case where by combining Step and PWL, we can build a higher-performing index than the ones solely comprising Step or PWL, respectively.

First, we construct an optimal B-tree index (Figure 3a) by comparing the latencies of multiple B-trees with different node sizes (100 bytes–10 MB): for the dataset we use, 16KB nodes offer the fastest lookup. In Figure 3d, we show the costs of index traversal to fetch each layer from the root to the data layer to find the relevant key-value; For the tuned B-tree, the end-to-end latency takes slightly longer than $1,000\mu$ s. Likewise, we construct an optimal index solely consisting of PWL layers (Figure 3b). The dataset we consider has its keys distributed favorably for PWL; thus, we need smaller layers than the optimal B-tree, needing to fetch 96B for the root, 8KB for Layer 1, and 10KB for the data layer. Accordingly, the overall latency shown in Figure 3d is also lower than the B-tree.

In contrast, by combing Step and PWL, we can discover an index (denoted by AIRINDEX) with significantly more compact layers



Figure 4: AIRINDEX Architecture. AIRINDEX manages indexes on Storage Layer (e.g., SSD, NFS, cloud storage APIs) via its Storage Layer Interface that abstracts read/write operations.

overall (Figure 3c). Specifically, while AIRINDEX's root layer is bigger in size compared to the tuned PWL index, it allows fetching significantly smaller amounts of data for the other layers (i.e., Layer 1 and the data layer), lowering the end-to-end latency. In theory, an optimal index with heterogeneous layers is guaranteed to deliver performance not worse than the optimal index with homogeneous layers. Nevertheless, the amount of improvement we can offer with layer-wise optimization well compensates for the increase in search effort, as we demonstrate with more empirical results in §7.

2.3 Need for Novel Index Tuning

New Search Space. Because of the new connection between machine learning and indexing [45], the search space has grown exponentially, including models types (e.g., linear models to neural networks) as well as their parameters (e.g., regression coefficients, weights), hyperparameters (e.g., fitting methods, regularization weights), and error correction algorithm (i.e., last mile search). The search space expands exponentially as each index sub-structure, like layer or node, can choose its model design individually.

Lack of Predictability. Unlike traditional indexes, learned indexes have an unpredictable performance *before* fitting them to the data, which consequently increases tuning costs. While we can rigorously analyze some traditional indexes like B-trees [17, 25] and skiplists [57] to predict their worst-case or average-case performance given a set of hyperparameters (e.g., fanout, level fraction, key prefix size), analyzing learned indexes and learned models is much more challenging because: (1) learning performance depends on nontrivial statistics of dataset, in which the dependency may not yet be explainable, (2) learned models often comprise of many learning components whose existing theories may not combine together, and (3) learning may unreliably produce poorly fitted models due to randomization or numerical instability. As a result, learned index tuning needs to pay the price to fit models or restrict the search space to a limited set of well-understood learned models.

Existing Tuning's Limitations. Existing tuning methods are either restricted, inconclusive, or inefficient. Traditional tuners like [36] only apply to restricted sets of indexes prior to learned indexes, missing opportunities to fit better to data patterns and so save lookup latency. On the other hand, many existing learned index tuners inconclusively recommend many index designs, for example, CDFSHOP [53] finds many Pareto efficient RMIs while PGM-INDEX [33] and PLEX [64] reduce their parameters into fewer hyperparameter(s), delegating selection or hyperparameter tuning to their users. Lastly, naïve methods such as brute-force search, grid search, or binary search incur large overheads due to excessive fitting or multiple sequential passes on the full dataset. Later, our

Supawit Chockchowwat, Wenjie Liu, Yongjoo Park

Table 2: Notations and their meaning.

Input to AIRINDEX		
$x = [y^-, y^+)$	Key Position (range) on storage Key-position dataset	
$T(\Delta)$	Storage profile, time to read Δ bytes from storage	
Design Parameters (§4)		
$ \begin{array}{c} L \\ \boldsymbol{\Theta} \\ \boldsymbol{\Theta}_{l} \\ \boldsymbol{s}(\boldsymbol{\Theta}_{l}) \\ \boldsymbol{\Delta}(\boldsymbol{x};\boldsymbol{\Theta}_{l}) \\ \boldsymbol{\Delta}(\boldsymbol{D};\boldsymbol{\Theta}_{l}) \\ \boldsymbol{\hat{y}} = [\hat{y}^{-}, \hat{y}^{+}) \\ \mathcal{L}_{SM} \end{array} $	Number of index layers Parameters across all index layers Parameters describing the <i>l</i> -th index layer Size of the <i>l</i> -th index layer Read size of <i>x</i> predicted by <i>l</i> -th index layer Average read size over keys in dataset <i>D</i> A predicted position on storage Index lookup cost from storage model	
Search Proces	s (§5)	
$F \in \mathcal{F}$ k $1 + \epsilon$	Node builder mapping a D into a Θ Number of top candidates to branch out exponential base for granularity exponentiation Bound for granularity exponentiation	
$\tau(D)$ $\tau(D)$ $\hat{\tau}(D)$	Index complexity, the ideal index lookup cost Step index complexity, an upper bound to $\tau(D)$	

experiments (§7.5) verify that these existing methods find suboptimal indexes and incur large tuning overheads. Therefore, a novel index tuning is needed to identify the optimal learned index from a much larger search space by efficiently balancing learning costs.

3 SYSTEM OVERVIEW

We describe AIRINDEX's core components (§3.1) and a storage profile representing I/O performance (§3.2).

3.1 Architecture

AIRINDEX is an index library for sorted key-value data stored on a Storage Layer (e.g., SSD, NFS). AIRINDEX can work with various storage devices by extending its storage interface. AIRINDEX stores indexes together with the data on Storage Layer while part of them may be cached. Internally, AIRINDEX consists of three components-Lookup, Builder, and Storage Layer Interface-as depicted in Figure 4. First, Storage Layer Interface provides the consistent abstraction over different storage interfaces/devices (e.g., virtual file systems (VFS), over SSD or object storage, cloud storage services) such as creating files and reading/writing serialized objects. Based on this consistent interface, we can profile the time needed to read Δ bytes (§3.2). Second, Lookup provides a querying interface. Given a key, the module traverses an index as it fetches necessary data (e.g., ranges within index layers) if not cached, caches the fetched data if space is available (§4), and makes inferences on what data it needs to fetch next (§4.2), until it finds the value associated with the key (or if finds there are no such keys indexed). Finally, Builder finds optimal index designs and builds/stores actual indexes on storage (§§ 5.1 to 5.3). How Builder finds high-quality index designs efficiently is the key contribution of this work.

3.2 Storage Model

AIRINDEX relies on *storage performance profile* $T(\Delta)$, which represents the time taken for the storage layer to read consecutive



Figure 5: A hierarchical index of L layers. The l-th layer looks up the data in the (l-1)-th layer. The close-up diagram layouts key-value pairs (ellipses) stored in pages. It shows layer's size $s(\Theta_l)$ and precision $\Delta(x;\Theta_l)$ of a key whose node on l-th layer and relevant key-value on (l-1)-th are colored.

data of size Δ bytes. Considering that the read time \mathcal{T} is probabilistic (e.g., due to variability, system loads, access paths, address alignment, lower-level optimizations), AIRINDEX is interested in the conditional expectation $T(\Delta) = \mathbb{E}[\mathcal{T}|\Delta]$. For example, $T(\Delta)$ can be an affine function $T_{\text{aff}}(\Delta) = \ell + \Delta/B$ with latency ℓ and bandwidth *B*. If the latency and bandwidth uniformly varies in $\lfloor \ell_0, \ell_1 \rfloor$ and $\lfloor B_0, B_1 \rfloor$, after calculating the expectation, the storage profile becomes $T_{\text{aff}}(\Delta) = \frac{\ell_1 + \ell_0}{2} + \frac{\Delta(\ln B_1 - \ln B_0)}{B_1 - B_0}$. While we implement the affine storage model $T_{\text{aff}}(\Delta)$ parameterized by ℓ and *B*, our optimization works with any monotonically increasing $T(\Delta)$.

Currently, we consider a deterministic and time-independent $T(\Delta)$, summarizing over other variables that are unnecessary for index tuning. Future directions can extend $T(\Delta)$ from a deterministic function to a distribution conditioned on the read size Δ to incorporate the randomness. Such storage models would enrich AIRINDEX to tune on more complex goals such as fastest tail latency (e.g., lowest p99) or highest reliability (e.g., lowest latency variance).

3.3 Points of Applications

AIRINDEX applies to any index building in a system life cycle when lookup operations from storage need to be performant. AIRINDEX primarily aims for two types of applications. (1) Immutable indexes' bulk loading: AIRINDEX naturally builds high-speed indexes that do not change; nonetheless, the application can still support changing data with existing techniques such as gapped arrays and LSM-tree compaction. (2) Updatable indexes' initial design and maintenance: users can utilize AIRINDEX to build the initial index structure, then follow any updatable index protocol compatible with the structure. After the index has mutated significantly, users can re-build the index using AIRINDEX as a part of vacuum processes.

AIRINDEX incurs computation overhead but is designed to minimize its real-time overhead. Although not required, a system with high parallelism is preferable. Alternatively, users can trade off tuning accuracy for an overhead reduction through configurations.

4 AIRINDEX-MODEL: UNIFIED INDEX MODEL

In this section, we mathematically model hierarchical indexes by representing them with design parameters Θ , which are also used to express their lookup latencies. These mathematical models are the foundations of our optimization process described in §5.

(*a_i*, *b_i*) (7, [25, 35)) (5, [20, 25))

Figure 6: Two node types with key-position examples. Left: 5-piece step. Right: band with points and width.

 (x_1, y_1)

4.1 Hierarchical Indexes

(4, [10, 20)

(1, [0, 10))

A hierarchical index is a data structure that maintains data locations in a layered internal structure, consisting of explicit *L* index layers— $I_1, I_2, \ldots I_L$. Here, I_L refers to the root and I_1 refers to the bottommost one, next to the data layer denoted I_0 . A data system can query a hierarchical index for the location of a desired item, then retrieve the item from the location. Accordingly, a high-quality index must quickly find the location in a database, namely, *keyposition collection* $D = \{(x_i, y_i)\}_{i=1}^n$, where x_i is the *i*-th key with data on a position $y_i = y(x_i) = [y^-(x_i), y^+(x_i))$ (i.e. range of data).

Index Layer. Given a key, an index layer points to a range of data (in the next layer) that contains the information associated with the key. The system then can read the range of data from the next layer, and continues the index traversal until it reaches the data layer. To support partial reads, an index layer consists of one or more *nodes*; for a specific key, there is only one node associated with it. Conceptually, a node is a function from a key to its position (as defined shortly), and an index layer is a piecewise function comprising its node functions. This composite structure enables partial data reads and lower data access costs.

Node. A node is the smallest data structure that maps a key to its position in the next layer. Given a key x, a node predicts a position $\hat{y} : X \to \mathcal{Y}$ in range that encompasses the actual position y(x).

Node \hat{y} : $\hat{y}(x) = [\hat{y}^-(x), \hat{y}^+(x)) \supseteq [y^-(x), y^+(x)) = y(x)$ (1)

While our optimization framework can support any such function \hat{y} that satisfies Eq (1), AIRINDEX currently implements two types of nodes sketched in Figure 6. First, a step node (step) is a step function, or a p-piece constant function parameterized by partition keys $\mathbf{a} = (a_1, \dots, a_p)$ and partition positions $\mathbf{b} = (b_1, \dots, b_p)$. In other words, $\hat{y}_{step}(x; \mathbf{a}, \mathbf{b}) = [b_i, b_{i+1})$ for $x \in [a_i, a_{i+1})$. Second, a linear band node (band) is a thick linear function traveling through two key-position points (x_1, y_1) and (x_2, y_2) with a width δ : $\hat{y}_{\text{band}}(x; x_1, y_1, x_2, y_2, \delta) = [mx + c - \delta, mx + c + \delta)$ where $m = (y_2 - y_1)/(x_2 - x_1)$ and $c = y_1 - mx_1$. A step node is 16p bytes in size while a linear band node is 40 bytes. For example, consider a dataset with 4 key-position pairs: $D = \{(x, y)\} =$ {(1, [0, 10)), (4, [10, 20)), (5, [20, 25)), (7, [25, 35))} as in Figure 6. A step node can represent this key-position collection with p = 5, $a = (1, 4, 5, 7, \infty)$, and b = (0, 10, 20, 25, 35), guaranteeing the exact prediction $\hat{y}_{step}(x; \mathbf{a}, \mathbf{b}) = y$ for all $(x, y) \in D$. Note that a step node having fewer pieces can be valid but would have a higher error. Alternatively, a band node can approximately represent D with $(x_1, y_1) = (1, 0), (x_2, y_2) = (4, 10), \delta = 15$. We encourage readers to verify that $\hat{y}_{\mathsf{band}}(x; 4, 0, 7, 25, 35) \supseteq y$ for all $(x, y) \in D$.

Al	Algorithm 1: AIRINDEX Query Process, Lookup $(x; \hat{y}_L, L)$		
Input: Query key <i>x</i> , root position \hat{y}_L , number of layers <i>L</i> Output: Relevant key-value (x, v)			
1 for l from L to 0 do			
2	$\{(x[i], v[i])\}_i \leftarrow Read(\hat{y}_l(x))$	<pre>// Storage access</pre>	
3	$v_l \leftarrow \text{Search}(x, \{(x[i], v[i])\}_i)$		
4	if In index layer, $l \ge 1$ then		
5	$\hat{y}_{l-1} \leftarrow \text{ReconstructNode}(v_l)$		
$6 \text{return} \ (x, v) = (x, v_0)$			

Currently, these two node types are sufficient. Together, they fit sorted key-position collections accurately. step adapts to discontinuities while band regresses well with regularly sized key-value pairs [32]. They also possess many efficient fitting methods that pass over the entire key-position collection only once. Furthermore, with these types, AIRINDEX-MODEL can already incorporate many existing learned indexes such as PGM-Index [33], ALEX [31], RadixSpline [43], SIndex [69], XModel [70].

Parameters (to be tuned). The collection of all parameters Θ is a nested tuple of parameters that is sufficient to represent a hierarchical index instance. In particular, Θ specifies the number of layers *L* followed by *L* layer-wise parameters Θ_i which in turn specifies the node type, number of nodes n_l , and all node parameters.

$$\begin{split} &\Theta = (L, (\Theta_1, \ldots, \Theta_L)), \quad \Theta_l = (\mathsf{NodeType}, n_l, (\theta_1, \ldots, \theta_{n_l})) \quad (2) \\ &\mathsf{Node-specific parameters} \ (\theta_1, \ldots, \theta_{n_l}) \ \text{for} \ n_l \ \text{nodes depend on} \\ &\mathsf{the corresponding NodeType}. \ \mathsf{For example, if NodeType} = \mathsf{step}(p), \\ &\theta_i = (\mathbf{a}_i, \mathbf{b}_i) \ \text{where} \ |\mathbf{a}_i| = |\mathbf{b}_i| = p. \ \text{If NodeType} = \mathsf{band}, \ \theta_i = (x_1^{(i)}, y_1^{(i)}, x_2^{(i)}, y_2^{(i)}, \delta_i). \ \text{We choose to specify a single node type} \\ &\mathsf{per layer over multiple node types per layer because (1) \ it reduces \\ &\mathsf{the serialization overhead per node which reduces read volume and \\ &\mathsf{so overall lookup latency, and} \ (2) \ it \ \text{simplifies our tuning (§5.2).} \end{split}$$

Examples. As a general class of indexes, hierarchical index examples include traditional indexes like B-tree and learned indexes like RMI, PGM-INDEX, balanced ALEX/APEX, and PLEX. For example, 2-layer balanced B-tree with fanout 3 is a hierarchical index as in Eq (3) where $\theta_{i,j}$ and θ_k are appropriate step function parameters to encode keys and pointers in the leaf and root nodes respectively.

$$\Theta_{B-Tree} = (2, (\Theta_1, \Theta_2)),$$

$$\Theta_2 = (\text{step}, 3, (\theta_1, \theta_2, \theta_3)),$$

$$\Theta_1 = (\text{step}, 3^2, (\theta_{i,j})_{i,j \in \{1,2,3\}})$$
(3)

RMI with a cubic root node $ax^3 + bx^2 + cx + d$ and 8 linear intermediate nodes $(\alpha_i x + \beta_i) \pm \varepsilon_i$ is a hierarchical index as in Eq (4). PGM-INDEX, balanced ALEX/APEX, and PLEX can be similarly instantiated as hierarchical indexes by extending the node types.

$$\begin{split} \Theta_{\text{RMI}} &= (2, (\Theta_1, \Theta_2)), \\ \Theta_2 &= (\text{cubic}, 1, ((a, b, c, d))), \\ \Theta_1 &= (\text{linear}, 8, ((\alpha_i, \beta_i, \varepsilon_i))_{i \in \{1, \dots, 8\}}) \end{split}$$
(4)

4.2 Query Process

Functionally, a query process takes in a search key x and outputs its relevant key-value (x, v). It internally consults index layers inside a hierarchical index, starting from its root layer L, traversing down

Table 1	3: 0	otimization	summary
---------	------	-------------	---------

Design variables Θ:			
L	Number of layers		
NodeType _l	Node type in layer $l \in \{1, \ldots, L\}$		
n_l	Number of nodes in layer $l \in \{1,, L\}$		
$\theta_{l,i}$	Parameters of the <i>i</i> -th node in layer $l, i \in \{1,, n_l\}$		
Fixed variables:			
Т	Storage profile		
X	Query key distribution		
D	Key-position collection $D = \{(x_i, y_i)\}_{i=1}^n$		
Constraints	Valid index $\hat{y}(x) \supseteq y$ for all $(x, y) \in D$		
Objective	Minimize expected latency $\mathcal{L}_{SM}(X; \Theta, T)$, Eq (7)		
Algorithm	AirTune (§5)		

one index layer at a time to look up the relevant position of the query key, and finally retrieving the target value.

Alg. 1 formalizes AIRINDEX's overall query process, traversing through the index hierarchy of relevant node(s) $\hat{y}_L, \ldots, \hat{y}_0$ to retrieve the relevant key-value (x, v). There are mainly three steps in each iteration in an index layer: (1) AIRINDEX reads potentially relevant raw bytes $\{(x[i], v[i])\}_i$, (2) it then searches for the relevant raw bytes v_l based on the tagged key $\{x[i]\}_i$, and (3) AIRINDEX deserializes v_l to reconstruct the next relevant node \hat{y}_{l-1} and predicts the next position $\hat{y}_{l-1}(x)$. At the end when l = 0, AIRINDEX returns the value with the query key $(x, v_0) = (x, v)$.

Such a clearly defined query process allows us to estimate a hierarchical index performance for tuning. Given appropriate statistics on the hierarchical index and a storage model, we can translate Alg. 1 into the latency formula incorporating different lookup costs.

4.3 Latency Under Storage Model

For index traversals, the dominant costs are storage accesses (Alg. 1, line 2), compared to other internal computations including relevant value searching, data deserialization, and node prediction. Following the iterations in Alg. 1, there are exactly L + 1 sequential storage accesses corresponding to reading L index layers and the data layer. Specifically, AIRINDEX first reads the entire root layer of size $s(\Theta_L)$ bytes, then partially reads $\Delta(x; \Theta_{l+1}) = |\hat{y}_l(x)|$ bytes from the next index layer, and so on until it reaches the data layer to read $\Delta(x; \Theta_1) = |\hat{y}_0(x)|$ bytes. Using a storage profile T (§3.2), we express $\mathcal{L}_{SM}(x; \Theta, T)$, the query latency to find a value for a key x under the storage model, as follows:

$$\mathcal{L}_{SM}(x; \Theta, T) = T(s(\Theta_L)) + \sum_{l=1}^{L} T(\Delta(x; \Theta_l))$$
(5)

To obtain expected latency over multiple keys, we aggregate query latencies, each specific to key x, over a query key distribution X. In this work, we set X to be a uniform distribution over existing keys in the key-position collection D.

$$\mathcal{L}_{SM}(\mathcal{X};\Theta,T) = \mathop{\mathbb{E}}_{x \sim \mathcal{X}} \left[T(s(\Theta_L)) + \sum_{l=1}^{L} T(\Delta(x;\Theta_l)) \right]$$
(6)

Given a key-position collection *D*, a storage profile *T*, and a query key distribution X, our objective (Eq (7)) is to minimize this expected query latency where Θ represents a hierarchical index



Figure 7: An instance of AIRTUNE execution, starting from the bottom (data layer) to the top. The label "root" indicates that the candidate should be the root layer with no further index layer. In this illustration, the algorithm has 7 layer builders and selects k = 3 top candidates to branch out.

valid over D. Table 3 summarizes this optimization problem.

$$\Theta^* = \arg\min_{\Theta} \mathop{\mathbb{E}}_{x \sim X} \left[T(s(\Theta_L)) + \sum_{l=1}^{L} T(\Delta(x;\Theta_l)) \right]$$
(7)

5 AIRTUNE: SEARCH WITH BOUNDED VISITS

This section describes AIRINDEX's optimization algorithm, AIR-TUNE, which solves the optimization defined earlier (Table 3). We first describe the overall process (§5.1). Then we explain important components in detail: index layer builders (§5.2), pruning technique for computational efficiency (§5.3), and parallelization techniques (§5.4). Finally, we analyze its time complexity (§5.5).

5.1 Guided Graph Search

AIRTUNE is a guided graph search; it starts from an origin vertex, walks over edges to different vertices, and stops when a stopping criterion is satisfied. Each vertex represents a layer in a hierarchical index. As the special case, the origin vertex u_0 represents the dataset being indexed. Each edge from a vertex u to another vertex u'represents a *layer builder* building an index layer u' based on u. A path from the origin represents a particular hierarchical index design. For example, a path (u_0, u_1, u_2) is roughly equivalent to data and index layers (I_0, I_1, I_2) where the index layers have parameters $\Theta = (2, (\Theta_1, \Theta_2))$. Two paths with common vertices represent two hierarchical index designs that share lower layers; thus, AIRTUNE can reuse layer-building results in those common lower layers.

On a vertex, AIRTUNE first explores all outgoing edges, i.e. possible candidate index layers from all available layer builders (§5.2, Alg. 2, lines 3–6). Layer builder explorations are the most expensive step, but they are independent of one another and embarrassingly parallelizable. AIRTUNE leverages this observation to reduce the tuning time overhead (§5.4). Upon receiving all candidates, it consults heuristic guidance to select only the top-k candidates (§5.3, Alg. 2 line 7) to continue the search recursively (Alg. 2 lines 8–12). Selecting a few promising candidates is an important mechanism to limit the branching factor, avoiding exponential time complexity (see analysis in §5.5). At the end of a recursive search, AIRTUNE

SIGMOD '24, June 11-16, 2024, Santiago, Chile (Accepted 23 May 2023)

Algorithm 2: AIRINDEX Index Tuning, AIRTUNE $(D; T, \mathcal{F})$			
Input: Key-position collection $D = \{$ layer builders \mathcal{F} Output: Index structure Θ^*	(x_i, y_i) ^{<i>n</i>} _{<i>i</i>=1} , storage profile <i>T</i> ,		
// Check stopping criterion 1 if $\mathcal{L}_{SM}(D; (), T) < \text{IdealLatencyWithIndex}(T)$ then 2 \lfloor return () // Cannot improve with additional layer			
<pre>// Build multiple next layer candi</pre>	dates		
3 for F in \mathcal{F} do			
4 $\Theta_{\text{next}} \leftarrow F(D)$	// Build next layer (§5.2)		
$_{5}$ $D_{\text{next}} \leftarrow \text{Outline}(\Theta_{\text{next}})$	<pre>// Turn into key-positions</pre>		
	<pre>// Append candidate</pre>		
// Select top-k candidates (§5.3) 7 $C \leftarrow \text{Select}(C, k)$			
// Build indexes on top- k candidat	tes		
s for $(\Theta_{next}, D_{next})$ in C do			
9 $\Theta_{\text{next}}^* \leftarrow \text{AirTune}(D_{\text{next}}; T, \mathcal{F})$	<pre>// Call recursively</pre>		
10 $\Theta_{\text{new}} \leftarrow (\Theta_{\text{next}}) \oplus \Theta_{\text{next}}^*$	// Prepend layer		
11 if $\mathcal{L}_{SM}(D; \Theta_{new}, T) < \mathcal{L}_{SM}(D)$	$\mathbf{if} \ \mathcal{L}_{SM}(D; \mathbf{\Theta}_{ngw}, T) < \mathcal{L}_{SM}(D; \mathbf{\Theta}^*, T) \mathbf{then}$		
$12 \qquad \qquad$	Select the better structure		
13 return Θ^*			

compares all the options using the latency formula under a storage model (Eq (6)) and returns the best hierarchical index design.

To decide when to stop searching, AIRTUNE determines whether an additional index layer will be beneficial. That is, if the *ideal index layer*—the best (possibly impossible) layer we can build on top—does not reduce the query latency, AIRTUNE stops further exploration and declare the current vertex as the root index layer (Alg. 2, lines 1–2). Specifically, an ideal index has the minimal size of $s(\Theta) = 1$ byte and the finest precision $\Delta(x; \Theta) = 1$ byte.

Example. If we have 7 layer builders and set k = 3, an execution of AIRTUNE could look like Figure 7. That is, AIRTUNE starts from the origin data-layer vertex (the long rectangle at the bottom) and explores all 7 layer builders, resulting in 7 vertices (rectangles of varying sizes in the middle). The guidance then tells AIRTUNE to search deeper into k = 3 candidates of them (red and blue highlighted rectangles). AIRTUNE stops at one of the candidates because it is too small to gain any benefit from an ideal index layer. For the rest, AIRTUNE continues exploring, selecting top-3 candidates, and finally stops at 2nd layer candidates. After comparing all options, AIRTUNE reports the best path (all red highlighted rectangles), representing the fastest hierarchical index design.

5.2 Layer Builders

A layer builder is a method to produce a valid index layer on top of existing index layer(s). In other words, it is a mapping $F(D) = \Theta$ such that Θ satisfies a valid index layer $\hat{y}(x) \supseteq y$ for all $(x, y) \in D$.

In theory, there is a large number of ways of building index layers. For example, a method A_1 can find the smallest collection of band that covers D with error at most λ bytes using $O(n^2)$ for n key-position pairs. In O(n), another method A_2 could quickly and mindlessly connect every other m key-position points into a collection of band. To avoid exploring every possible method, we



Figure 8: Index complexity $\hat{\tau}(D;T)$ as function of data size s_D and affine storage profile T (parameterized by latency ℓ , and bandwidth B). Different line represent different variations in bandwidth (left) and latency (right).

choose a set of good layer builders that (1) run quickly (say, time complexity O(n)), (2) build small and accurate index layers, and (3) synergetically cover different data patterns together. In the earlier examples, method A_1 builds the optimal band index layer but is too slow, while method A_2 is fast but builds a suboptimal band index layer. In addition, $\{A_1, A_2\}$ is also not a good set of methods, because they only cover the band node type.

To cover the two types of nodes (step and band) on different data patterns, AIRINDEX currently deploys three types of layer builders, each generating many layer builders by varying hyperparameters. (1) **Greedy Step** (GStep(p, λ_{GS})) builds p-piece step nodes with precision at most λ_{GS} bytes by greedily packing key-position pairs. (2) **Greedy Band** (GBand(λ_{GB})) builds band nodes by greedily fitting as many key-position pairs as possible using the monotone chain convex hull [10]. (3) **Equal Band** (EBand(λ_{EB})) builds band nodes by grouping key-position pairs in equal-size position ranges. Please see our extended script [66] for further details.

AIRINDEX generates the set of candidate layer builders \mathcal{F} by sampling the granularity exponentially: λ_{low} , $\lambda_{low}(1 + \epsilon)$, $\lambda_{low}(1 + \epsilon)^2$, ..., λ_{high} where (λ_{low} , λ_{high}) are the bounds and $\epsilon > 0$ controls the exponentiation base. For example, if $\lambda_{low} = 2^8$, $\lambda_{high} = 2^{20}$, and $1 + \epsilon = 2$ with p = 16, then \mathcal{F} contains 39 builders in total:

$$\mathcal{F} = \{ \mathsf{GStep}(16, 2^8), \mathsf{GStep}(16, 2^9), \dots, \mathsf{GStep}(16, 2^{20}) \} \\ \cup \{ \mathsf{GBand}(2^8), \mathsf{GBand}(2^9), \dots, \mathsf{GBand}(2^{20}) \}$$
(8)
$$\cup \{ \mathsf{EBand}(2^8), \mathsf{EBand}(2^9), \dots, \mathsf{EBand}(2^{20}) \}$$

5.3 Top-k Candidates by Index Complexity

Before branching out, AIRTUNE selects only top-*k* candidates with the highest potential to be in the optimal design. For each candidate (Θ_i , D_i), AIRTUNE evaluates its quality as a summation of a "remaining work" heuristic function $\hat{\tau}(D; T)$ and its layer-specific lookup latency. Then, it selects top-*k* candidates with *k* lowest estimated costs (arg min^{*k*} denotes top-*k* arguments of the minima).

$$\{(\Theta_i, D_i)\}_{i=1}^k = \underset{(\Theta_i, D_i) \in \mathcal{C}}{\arg\min}^k \quad \hat{\tau}(D_i) + \underset{x \sim \mathcal{X}}{\mathbb{E}} \left[T(\Delta(x; \Theta_i)) \right] \quad (9)$$

Choice of Heuristic Function. Ideally, if there exists an oracle that reveals the optimal search latency, we could simply select the best candidate and avoid branching out entirely. This optimal search latency of dataset *D* under storage profile *T* is called *index*

complexity $\tau(D;T)$. Unfortunately, $\tau(D;T)$ is unknown for a large class of indexes supported by AIRINDEX.

Instead, AIRTUNE estimates candidates' quality using an upper bound to the index complexity: *step index complexity* $\hat{\tau}(D;T) \ge \tau(D;T)$. $\hat{\tau}(D;T)$ is the optimal search latency considering only step index layers (i.e. B-tree layers). Since the quality of stepbased layers can be analytically computed, we obtain an efficient algorithm that depends only on the collection size s_D and storage profile *T*. Figure 8 shows the general shape of $\hat{\tau}(D;T)$ solved with the algorithm. Please see our extended report [66] for more details.

5.4 Parallel Tuning

AIRINDEX is highly parallelizable from three sources of parallelisms as described below, ordered from finest to coarsest layers. It is worth noting that, with a proper branching (Theorem 5.1), the node-building step is the primary target for parallelization.

From Data Partitioning. AIRINDEX partitions the key-position collections, uses a layer builder to build sub-range candidates, and merges them into a candidate. This is possible because AIRINDEX's existing layer builders generate a piecewise function that can be merged together across different key ranges. By default, AIRINDEX breaks a key-position collection into partitions, each containing 1 million key-position pairs. Thus, AIRINDEX can scale with growing data sizes by increasing the level of parallelisms accordingly.

From Across Layer Builders. Although Alg. 2 calls layer builders in loops, the invocations are independent of one another. AIRINDEX conveniently turns the for-loop into a parallel mapping to produce $(\Theta_{next}, D_{next})$ and collecting into the candidate set *C*. With more parallelisms, AIRINDEX can then scale along with the diversity of layer builders to capture wider key-position patterns.

From Branching. AIRINDEX can recursively call Alg. 2 in parallel and select the best index structure with the minimum storage model cost. This source of parallelism allows AIRINDEX to explore more candidate branches and especially higher structures when favored by the storage profile (e.g., low bandwidth and low latency).

5.5 Analysis

Branching recursive optimization requires a balancing between the branching factor and depth. If AIRTUNE branches out to too many candidates relative to the index depth, it would become uncontrollably slow. On the other extreme, it could miss the optimal candidate branch. We first analyze the tuning time complexity with respect to a choice of hyperparameter. Then, we attempt to provide an approximation factor of the automatic tuning.

THEOREM 5.1. Time complexity: Let there be n key-position pairs, L layers to be explored at most according to the storage profile T, and $|\mathcal{F}|$ layer builders. If AIRTUNE selects at most $k \leq {}^{L+}\sqrt{n}$ candidates then the time complexity is at most $O((L+1)|\mathcal{F}|n)$.

PROOF. Suppose the data *D* consist of *n* key-position pairs with total data size s_D . Let *L* be the maximum number of layers expected to be explored, which has the upper bound $L_{\max} \ge L$: the number of layers chosen by the step index complexity.

Because our layer builders all process a O(n) key-position collection in O(n) time, the time to build all $|\mathcal{F}|$ candidates is $O(|\mathcal{F}|n)$.

Next, considering the worst case compression ratio $O(n^{\frac{L}{L+1}})$ and the number of branches k, we expect $O(n^{\frac{L}{L+1}})$ key-position pairs and so the time complexity in the next layer is $O(|\mathcal{F}|kn^{\frac{L}{L+1}})$. The branching and compression of key-position pairs continue until the root layer. In summary, the total time complexity in *L*-layer branching recursion is as followed whose last step is the closed-form formula to the geometric series.

$$O\left(|\mathcal{F}|\sum_{l=0}^{L}k^{l}n^{\frac{L+1-l}{L+1}}\right) = O\left(|\mathcal{F}|n\frac{1-k^{L+1}/n}{1-k/n^{\frac{1}{L+1}}}\right)$$
(10)

Under $k^{L+1}/n \le 1$ (i.e. $k \le {}^{L+1}\sqrt{n}$), this reduces to $O(|\mathcal{F}|(L+1)n)$ when $n \to \infty$.

5.6 Other Implementation Details

We implement AIRINDEX in Rust [1]. Like many data systems, it has two explicit levels in its memory hierarchy: the underlying storage and its internal read-through cache. AIRINDEX interacts with storage through an abstract interface, in which concrete implementations serve partial range reads at their best effort. Currently, AIRINDEX's internal zero-copy read-through cache employs a first-in-first-out (FIFO) eviction policy due to its admission simplicity. Apart from customized node type format, AIRINDEX serializes the metadata together with root layer as a byte array in the Postcard [55] format via Serde [60].

6 EXTENSIONS

Supporting Updates. Although AIRINDEX does not optimize for write workloads, it can tune and build an index that supports write operations. For example, we can augment the data layer into a gapped array, allowing insertion into gaps and deletion without changing index layers. When gaps are filled or expected to be filled, we can enlarge data layer's gaps and build a new index with AIRINDEX. AIRINDEX can also serves as the initial bulk loading of an updatable index (e.g. ALEX/APEX). However, the updatable index may evolve suboptimally. To reduce the frequency of structural index updates, we can enlarge the position granularity from bytes to pages, or we can buffer writes similarly to LSM-trees [26, 54].

Cache-aware Optimization. AIRINDEX can find an optimal *cache-aware* index design by additionally considering the distribution of cache hit C_l and the cost of cache access \mathcal{L}_{cache} . This modification then only affects the index complexity $\hat{\tau}$ and candidate selection in AIRTUNE. While this work does not include explicit results for cache-aware optimization, Figure 10 indicates that the indexes built with cache-pessimistic optimization already offer high performance than other existing methods across a wide range of cache warmness.

Pre-Search Assessment. Upon significant data or workload change, we can first assess the potential performance gain through the step index complexity (§5.3). Based on the assessment, users can better decide whether to tune the index. Step index complexity is a loose upper bound of the gain, however. Future works combining more accurate index complexity, what-if index design techniques, and data/workload trackers would help avoid unnecessary search costs for marginal performance gain.

7 EXPERIMENT

We empirically study AIRINDEX to demonstrate its faster search (§7.2), benefits of automatic index designs (§7.3), adaptability under wide ranges of I/O profile (§7.4), and quick build time (§7.5).

7.1 Setup

The experiment locates on two physical components: compute and storage. The former hosts benchmark scripts to execute queries against systems and measure their performance. These scripts together with required binaries are stored on local storage (i.e. Azure OS Disk). Meanwhile, the latter stores both datasets and indexes. Our benchmark consists of 40 runs in each setting: the *i*-th prepares the environment (e.g., clear cache¹), loads/executes the *i*-th list of one million query keys sequentially, and measures the elapsed time. We summarise those runtimes with average and standard deviation.

System Environment. We use Azure cloud platform [12], specifically D8s_v3 (8 vCPUs, 32 GiB RAM) with Ubuntu 20.04. The VM connects to two types of storage. (1) *NFS*: Azure network file system [15] hosted on Azure Blob Storage [13] (StorageV2, standard performance, zone-redundant storage, hot access tier). (2) *SSD*: Azure Premium SSD [14] with P20 performance tier (256 GiB, 2300 IOPS, 150 MBps, read/write host caching). 3) HDD: Azure Standard HDD [14] (1024 GiB, 500 IOPS, 60 MBps, no host caching). All resources are allocated within the same East US region.

Baselines. We compare AIRINDEX to a traditional database index, learned indexes, and our manual configuration counterpart. Although many of them are in-memory indexes, we integrate them onto external storages whose implementations are in their respective forks [2-6]. We manually tune each of the baselines through microbenchmarks. (1) LMDB: LMDB [48] is a B-tree database that accesses its data on storage through mmap. (2) RMI: RMI [41, 53] is a top-down learned index with a compact two-layer structure where the top one contains only one perfectly accurate node partitioning key space to the bottom nodes. We utilize CDFSHOP [53] to recommend function types and select the most accurate RMI across all datasets. (3) PGM-INDEX: PGM-INDEX [33] is a learned index with bounded precision across all layers. PGM-INDEX partitions the keyposition collection to build the next bottom-most layer towards the top. (4) ALEX/APEX: ALEX/APEX [31] is an updatable learned index built top-down like RMI but further arrange key-value pairs in its layout (notably, "gapped array" to buffer structural changes). (5) PLEX: PLEX [64] is a learned index with compact Hist-Tree (CHT) layered on top of RadixSpline [43] (RS). Although PLEX optimizes most parameters, its user need to specify the maximum prediction error ϵ . We select $\epsilon = 2048$ based on a benchmark on a setting (Figure 12d). (6) DATA CALCULATOR: DATA CALCULATOR [38, 39] is a data layout design engine that calculates the performance of a data structure. We follow its auto-completion and build its recommended data layout within AIRINDEX's framework. (7) B-TREE: A B-tree-like structure implemented using AIRINDEX (4KB pages and 255 fanout), which serves as the most controlled baseline.

¹For both NFS and SSD, we execute sysct1 vm.drop_caches=3 on a VM, clearing Linux-related caches such as page cache, entries, and inodes. For NFS, we also unmount and re-mount the Azure Blob Storage NFS to reload its client.



Figure 9: Average first-query latency comparison on NFS, SSD, and HDD storages across different datasets. In each storage and dataset setting, the bars represent the average latencies over random keys. LMDB, RMI, PGM-INDEX, ALEX/APEX, PLEX, DATA CALCULATOR, B-TREE, and AIRINDEX, from left to right. The shorter a bar, the faster the corresponding method serves queries.



Figure 10: Average latency curves on NFS and SSD with the books and osm datasets. A combination of a line color and a marker style represents each method: LMDB, RMI, PGM-INDEX, ALEX/APEX, PLEX, DATA CALCULATOR, B-TREE, AIRINDEX, respectively. Note the logarithm scales on both axes. Latency curves to the bottom-left corner represent faster methods.

Datasets. First, we use the SOSD benchmark [42, 52], including books (800M), fb (200M), osm (800M), and wiki (200M). Each of these contains 200-800 million 64-bit integer keys stored consecutively in an array. Given a query integer key, the task is to find its offset position in the array. Equivalently, it asks the systems for the rank of the query integer. As an unusual dataset, wiki contains many duplicated keys in which the task is to find the smallest offset of the key. Second, for more diverse data patterns, we also use a synthetic dataset, gmm, generated from a Gaussian mixture model (GMM) of 100 normal distribution clusters over 800 million keys.

7.2 Faster End-to-end Lookup Speed

We study cold-state and warm-state latencies separately. Cold-state latency is useful for understanding the performance under shortlived executions (e.g., serverless, ad-hoc workloads) and very large data (e.g., many tables, large indexes). Afterward, we study warmstate latency curves over different warmnesses. §7.4 later discusses index structures discovered by AIRINDEX.

Cold-state Latency. AIRINDEX is consistently one of the fastest methods at searching the first query, across datasets and storage (Figure 9). Compared to LMDB (B-tree), AIRINDEX is $2.4\times-2.7\times$ faster on NFS and $2.6\times-4.1\times$ faster on HDD. AIRINDEX is on par with LMDB on SSD. Similarly, AIRINDEX is $2.0\times-2.4\times$ faster than B-TREE on NFS but performs equally well on SSD on HDD. This difference across storages suggests that LMDB and B-TREE structures are tuned to disk-scale storage profiles like SSD's, so they underperform on storage on a different scale like NFS.

Compared to learned index baselines, AIRINDEX is reliably faster without unexpectedly long latency arising from tuning difficulties. While RMI performs reasonably well, it has two limitations.

First, its fixed two-layer structure makes it rigid to the underlying storage. Second, because RMI's top-down building assigns a disproportionate amount of data to intermediate nodes, the second-layer precision varies substantially. This later effect is more pronounced in gmm. Overall, AIRINDEX delivers lower latencies more reliably than RMI, being $1.2 \times -2.6 \times$ faster on NFS, $1.0 \times -2.2 \times$ faster on SSD, and 1.4×-5.9× faster on HDD. PGM-INDEX suffers on books, fb, and osm, but is competitive on wiki and gmm. Upon closer inspection, PGM-INDEX fits poorly with the former three, creating larger indexes than those in the latter two. AIRINDEX is on par with PGM-INDEX on wiki and gmm, but outperforms on other datasets with 3.0×-7.2×, 5.8×-11.7×, and 9.0×-15.6× speedup on NFS, SSD, and HDD respectively. Similarly, AIRINDEX is faster than ALEX/APEX in all settings except in gmm SSD, with 1.7×-10.1×, 1.3×-46.3×, and 3.0×-22.2× speedup on NFS, SSD, and HDD. ALEX/APEX performs poorly in osm because its root node holds 2M child pointers with more than 15MB of data. Lastly, AIRINDEX is faster than PLEX on NFS and HDD with 1.5×-2.2×, and 1.7×-3.3× speedup. Both methods perform equally well on SSD: AIRINDEX is 1.7× faster at searching osm, but 1.4× slower at gmm. Upon closer inspection, PLEX's compact histogram tree fit osm poorly (762KB in size) but fit gmm exceptionally better (0.6KB in size).

Compared to DATA CALCULATOR, AIRINDEX searches in a richer set of indexes and so it is consistently faster: $1.4 \times -2.0 \times$, $1.2 \times -1.5 \times$, and $1.0 \times -1.4 \times$ speedup on NFS, SSD, and HDD. On the bright side, DATA CALCULATOR generally has a faster lookup than B-TREE because of its storage-aware tuning and B-TREE's fixed structure.

Warm-state Latency Curve. As we continue querying (and caching more), the queries become progressively faster shown as per-query latency curves in Figure 10. Here, a point (x, y) in the latency curve

SIGMOD '24, June 11-16, 2024, Santiago, Chile (Accepted 23 May 2023)



Figure 11: Comparison of average first-query latencies between AIRINDEX-tuned designs (in red) and manual alternatives across NFS and SSD (fb dataset) varying numbers of layers L and granularities λ . The error bars display standard deviations.



Figure 12: Comparison of average first-query latencies on the books dataset in NFS between AIRINDEX-tuned design (left) and other methods with varying knobs. RMI's settings are recommended by CDFSHOP with the least accurate (fewer models) on the left and the most accurate (more models) on the right. The error bars display standard deviations of the latency.

implies that a system completes x queries in $x \times y$ seconds. The differences in accelerations across methods can be explained by their index structures. In hierarchical indexes, shorter and narrower indexes accelerate more aggressively than taller and wider indexes because fewer (random) queries are needed to touch all nodes in a narrower index layer. For example, in osm dataset, ALEX/APEX accelerates faster than LMDB because of its shorter index with a one-node root as opposed to a full B-tree. Even though AIRINDEX optimizes for cold-state latency, the tuned structure is still faster for warm-state latency, ranging from 100 to 100K queries. Such a range of warmness is useful for a large collection of datasets, short-lived search sessions, or limited memory environments.

7.3 Layer-wise Optimization Helps

Speedup over Hierarchical Indexes. To empirically verify that AIRINDEX tunes accurately and finds a fast index, we compare AIRINDEX's tuned index designs against manually configured ones in terms of their first-query latencies. Figure 11 presents the comparison on numbers of layers *L* and granularity hyperparameter λ across the two storage within the same dataset fb. In all settings and variable dimensions, AIRINDEX consistently finds the fastest index designs. Inspecting the trends, we observe that λ forgivingly admits a larger optimal region, even in the logarithmic scale, than the number of layers *L*. This allows AIRINDEX to select a coarse granularity exponentiation base $1 + \varepsilon$ without risking suboptimality.

We have also experimented with other dimensions such as the granularity exponentiation base $1 + \varepsilon$ and the set of node estimators to discover any trade-off. Lower bases $1 + \varepsilon$ result in faster indexes but with only insignificant gain for a higher cost in a longer tuning time. A wider set of node estimators and types provides some fitness advantages. The impact is clear when the data pattern is exclusive to a node type, for example, band nodes fit perfectly on a uniformly random key set (uden64 from [42]) while step nodes do not.

Speedup over Well-tuned Baselines. AIRINDEX is faster than all baseline configurations, and so is faster than optimal baselines. Figure 12 varies all permissible page sizes of LMDB, all 10 RMI settings recommended by its optimizer, and 8 chosen ε in PLEX to cover the optimal valley. We observe that AIRINDEX is 2.7×, 1.5×, and 1.7× faster that the optimal LMDB, RMI, and PLEX, respectively. Our similar experiment with B-TREE by varying λ granularity shows 1.3× speed up. These gaps from optimal baselines suggest the benefit to consider a larger class of indexes (i.e. AIRINDEX-MODEL).

7.4 Adaptive to I/O Performance

On Testing Storages. Indeed, AIRINDEX discovers different optimal index structures for the NFS and SSD/HDD storages in previous experiments. NFS indexes have L = 1 index layer with only band node types, while SSD/HDD indexes have L = 2 layers with a mix of band-band, band-step, and step-band node types. The sizes of root layers $s(\Theta_L)$ and precisions Δ range from 36KB to 328KB in NFS and 864B to 16KB in SSD/HDD, depending on dataset size and complexity. Among 5 datasets, osm is the most challenging one, reflecting the same observation from [52]. Please see our extended manuscript [66] for specific index structures.

On Latency-Bandwidth Spectrum. If we have a storage with latency ℓ and bandwidth *B*, what would the fastest index look like? We answer this question as a whole, on a wide spectrum of latency $\ell \in [1\mu s, 1000s]$ and bandwidth $B \in [1\text{KB/s}, 1\text{TB/s}]$. Figure 13 shows AIRINDEX adapting its index to the diverse range of storage profiles. Higher bandwidth or latency promotes shallower indexes with coarser precision (larger total read volume). In the extreme, AIRINDEX proposes no index at all, i.e. fetching the entire data layer to search locally. On the other hand, lower bandwidth or latency promotes taller indexes with finer precision. Although this trend is similar to a well-known tuning rule of thumb for B-tree, AIRINDEX offers a more complete tuning on a much larger class of indexes, for any data pattern, and storage profile.



Figure 13: Impact of storage latency/bandwidth on AIRINDEX's index design. The fb dataset is used. Note the logarithm scales covering 1KB/s – 1TB/s bandwidth and 1 μ s – 1000s latency. The number of layers *L*, total read volume $s(\Theta_L) + \sum_{l=1}^{L} E_{x \sim X} \Delta(x; \Theta_l)$, and the optimal costs are displayed in color annotated in the sidebar. NFS, SSD, and HDD performances are marked accordingly.



Figure 14: Extreme errors (± 3 and ± 2 magnitude differences in latencies/bandwidths) expectedly make tuned indexes suboptimal compared to correctly tuned ones. Left: the index is tuned for NFS (50 ms, 12 MB/s). Right: the index is tuned for SSD (250 μ s, 175 MB/s). fb is the underlying dataset.

With Storage Variability. Storage performance may vary. If it varies within a magnitude, AIRINDEX's tuned index mostly stays optimal. In fact, it does so if the performance remains within the same band (i.e. same color in Figure 13a and Figure 13b) as the profiled performance. However, if the actual performance T' varies across many magnitudes, the index tuned with inaccurate storage profile T can be suboptimal as shown in Figure 14 through the relative slowdown of the index Θ tuned with the inaccurate profile T based on the index Θ' tuned with the actual profile T' in retrospect: $\mathcal{L}_{SM}(X;\Theta,T)/\mathcal{L}_{SM}(X;\Theta',T')$. For example, if the actual NFS latency is 0.001 times smaller than the profiled NFS latency (i.e., 50μ s instead of 5ms), the originally tuned index would be 35 times slower than the accurately tuned index.

7.5 Competitive Build Time

Total Build Time. Figure 15a measures index build times on a machine with 2 AMD EPYC 7552 48-Core Processors (192 CPUs in total). Build times in LMDB, RMI, PGM-INDEX, and ALEX/APEX only account for data loading, inserting into the system, and writing to files, excluding their manual hyperparameter tuning. DATA CALCULATOR's build time includes parallelized autocompletion and index building. All methods use all available cores.

Thanks to its parallel tuning (§5.4), AIRINDEX's total build time is competitive with other baselines. AIRINDEX is 3.8×, 5.9×, and 4.3× faster than LMDB, RMI/CDFSHOP, and ALEX/APEX. while AIRINDEX both tunes and builds as fast as RMI (excluding CDFSHOP time), PGM-INDEX, and PLEX build their indexes. Compared to DATA CALCULATOR's autocompletion and building, AIRINDEX is 2.3× faster despite exploring a larger class of indexes. Nonetheless, AIRINDEX tuning and building are 2.7× slower compared to its fixed structure counterpart B-TREE. If needed to be faster, AIRINDEX can relax some hyperparameters (e.g. number of candidates, granularity base and bounds) to trade its speed with its tuning accuracy.

Search Overhead. Figure 15b measures search overheads—the differences between total build time and build time given a known configuration (i.e., CDFSHOP's search procedure on RMI structures, DATA CALCULATOR's autocompletion, AIRINDEX's AIRTUNE excluding the time for building the optimal index). AIRINDEX incurs non-negligible search overhead (around 50 ns/key on the 192-core machine, or 9.6 μ s/key on single-core machines); however, this search overhead is lower compared to other methods. Because of its parallelization and top-*k* candidate selection limiting branching, AIRINDEX finds its index structure 7.8× and 3.3× faster than learned index tuning method CDFSHOP and traditional index tuning method DATA CALCULATOR, respectively. In contrast, DATA CALCULATOR's parallelized autocompletion is slow because it tries all design combinations similarly to a grid search. We note that CDFSHOP outputs multiple index structures on a Pareto front.

7.6 Applicable to Read-Write Workloads

We implement a proof-of-concept updatable AIRINDEX based on the gapped array [31] that allocates empty gaps on data layer for AIRINDEX to insert a key-value at any available gap within the predicted position $\hat{y}(x)$. Our read-write benchmark follows that of [31]. It initially inserts 100M keys sampled from osm and measures the time to complete 10K queries consisting of cycles of r read and w write operations. Four workloads vary the read/write proportion: (1) Read-Only: (r, w) = (1, 0), (2) Read-Write: r = 19, w = 1, (3)Write-Heavy: r = 1, w = 1, (4) Write-Only: r = 0, w = 1. The benchmark samples read keys uniformly from the inserted key set and samples write keys from the non-inserted key set.

In Figure 16, our prototype remains the fastest compared to updatable baselines (LMDB and ALEX/APEX) across all workloads, confirming that tuning for lookup speed is relevant to both read and write performances. Apart from the direct relation to read operations, lookup speed is relevant to write operations because



Figure 15: Index build and search overhead for different data sizes (200, 400, 600, and 800 million keys) from the gmm dataset. Build times for LMDB, RMI, PGM-INDEX, ALEX/APEX, and PLEX do *not* include their manual tuning time.



Figure 16: Average throughputs of LMDB, ALEX/APEX, and AIRINDEX across read-write workloads on osm dataset in SSD.

all methods need to first look up the insertion position from their indexes before writing the target key-value pair.

8 RELATED WORK

Our work is built on top of the vast amount of existing research on index design and optimization as summarized below; however, our unified model (AIRINDEX-MODEL) and efficient search (AIRTUNE) enables a high-quality data and I/O-aware hierarchical index.

Storage-aware Indexes. Besides the original B-trees [17, 25], many works have studied unique storage properties to design indexes specifically optimized for certain storage, such as CPU cache [35], SRAM cache [23], disk [24], NVMe [68], and distributed cloud [19, 71, 72]. As the most generic of all, [18] designs B-Trees that perform well for any I/O page size *P*; however, its storage profile (cache-oblivious model) only limits to $T(\Delta) = O(\lfloor \Delta/P \rfloor)$. Also, skip list has been adapted to various settings, such as multi-core [29], cache-sensitive for range queries [62], non-uniform access [27], and distributed nodes [37]. In contrast, AIRINDEX takes a general approach by composing an optimization problem in consideration of storage profiles, which makes it possible to adapt its structure without re-evaluating the parameters when the transfer size changes.

Index Tuning. AIRINDEX automates index designs by improving on long-standing heuristics such as "use larger pages for larger bandwidth" [11, 36, 49]. Rather than deciding *what* index to build, other index tuning techniques determine *when and where* to build index, as a well-known index selection problem (ISP) [21, 28, 34, 40, 51, 59, 65], which are orthogonal to our work.

DATA CALCULATOR [38, 39] helps designing efficient data structures by evaluate the cost of a structure in a what-if fashion in relation to workload and hardware. However, its auto-completion search—recursively trying all possible designs ($|\mathcal{E}| = 10^{16}$ discretized designs)— scales poorly, which worsens if we extend DATA CALCULATOR's periodic table and cost synthesis flowchart with learned indexes. AIRINDEX solves this challenge for index design.

Learned Indexes. Previous works largely focus on tuning in-memory indexes. [53] demonstrates an interactive model tool that allows

users to modify RMI configuration (per-layer node type and branching factor) and observe the resulting model fitness. Although the tool provides automatic tuning, it measures the lookup latency by benchmarking each configuration, which can be expensive on a larger scale. [64] formulates a lookup cost function and tunes a single maximum error hyperparameter to build a combination of RadixSpline [43] and a compact histogram tree. In contrast, AIRINDEX encompasses a larger index design space, and more importantly, targets a different cost setting where I/O cost is dominant.

Many manually tuned learned indexes have shown success stories in the context of external storage. [26] studies favorable conditions to learn data patterns and integrates learned indexes as an optimization into an LSM-based storage system on disk. For NVM storages, [22] adapts ALEX [31] to cooperate with the preferred access pattern. Instead of predicting locations of written records, [9] uses learned indexes to distribute data into blocks in BigTable [20]. Similarly, [47] deploys learned indexes to organize on-disk spatial data into shards and pages, reducing I/O costs over other spatial trees. With the learning paradigm at the core, [44] re-designs a database system that supports data persistence on disk but is only evaluated in the in-memory mode without disk accesses. Many of these works have already involved parameter tuning but as a manual tuning step for each setting. While we share common tuning principles, AIRINDEX can be seen as a fine-grained automatic optimization of learned index structures.

9 CONCLUSION

This work presents a novel index-building technique, AIRINDEX, that can build high-speed hierarchical indexes by learning from both data and I/O characteristics-the first of its kind. To achieve its goal, AIRINDEX formulates an optimization problem consisting of a large hierarchical index search space and a lookup latency objective function (AIRINDEX-MODEL). To overcome the computational challenges rising from the inter-dependency between index layers and exponentially many candidate designs, AIRINDEX explores the search space using a purpose-built graph-search method (AIRTUNE). Our experiments verify that AIRINDEX accurately finds the optimal configuration and provides performance gains over conventional indexes as well as state-of-the-art learned indexes. In many applications, the decisions on data placements-local disks, cloud storage, network file system-are relatively fixed, which makes AIRINDEX's data-and-I/O-aware optimization appealing to achieve significantly faster lookup compared to the ones not specifically optimized.

ACKNOWLEDGMENTS

This work is supported in part by Microsoft Azure.

Supawit Chockchowwat, Wenjie Liu, Yongjoo Park

REFERENCES

- [1] [n.d.]. https://github.com/illinoisdata/airindex-public.
- [2] [n.d.]. https://github.com/illinoisdata/lmdb.
- [3] [n.d.]. https://github.com/illinoisdata/RMI.
- [4] [n.d.]. https://github.com/illinoisdata/PGM-index.
- [5] [n.d.]. https://github.com/illinoisdata/ALEX_ext.
- [6] [n.d.]. https://github.com/illinoisdata/airindex-public/tree/main/src/bin/data_ calculator.rs.
- [7] [n.d.]. A high-performance distributed shared-log for Ceph. https://github.com/ cruzdb/zlog. [Online; accessed December-27-2022].
- [8] [n.d.]. MySQL. https://www.mysql.com/. [Online; accessed December-27-2022].
- [9] Hussam Abu-Libdeh, Deniz Altinbüken, Alex Beutel, Ed H. Chi, Lyric Doshi, Tim Kraska, Xiaozhou Li, Andy Ly, and Christopher Olston. 2020. Learned Indexes for a Google-scale Disk-based Database. *CoRR* abs/2012.12501 (2020).
- [10] A. M. Andrew. 1979. Another Efficient Algorithm for Convex Hulls in Two Dimensions. Inf. Process. Lett. 9, 5 (1979), 216–219.
- [11] Raja Appuswamy, Goetz Graefe, Renata Borovica-Gajic, and Anastasia Ailamaki. 2019. The five-minute rule 30 years later and its impact on the storage hierarchy. *Commun. ACM* 62, 11 (2019), 114–120.
- [12] Microsoft Azure. [n.d.]. Azure. https://azure.microsoft.com. [Online; accessed Jul-17-2022].
- [13] Microsoft Azure. [n.d.]. Azure Blob Storage. https://azure.microsoft.com/enus/services/storage/blobs/. [Online; accessed Jul-17-2022].
- [14] Microsoft Azure. [n.d.]. Azure managed disk types. https://docs.microsoft.com/ en-us/azure/virtual-machines/disks-types. [Online; accessed Jul-17-2022].
- [15] Microsoft Azure. [n.d.]. Network File System (NFS) 3.0 protocol support for Azure Blob Storage. https://docs.microsoft.com/en-us/azure/storage/blobs/networkfile-system-protocol-support. [Online; accessed Jul-17-2022].
- [16] Mahesh Balakrishnan, Jason Flinn, Chen Shen, Mihir Dharamshi, Ahmed Jafri, Xiao Shi, Santosh Ghosh, Hazem Hassan, Aaryaman Sagar, Rhed Shi, et al. 2020. Virtual consensus in delos. In 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20). 617–632.
- [17] R. Bayer and E. M. McCreight. 1972. Organization and maintenance of large ordered indexes. Acta Informatica 1 (1972), 173–189. Issue 3. https://doi.org/10. 1007/BF00288683
- [18] Michael A. Bender, Erik D. Demaine, and Martin Farach-Colton. 2000. Cache-Oblivious B-Trees. In FOCS. IEEE Computer Society, 399–409.
- [19] Huang Bin and Peng Yuxing. 2014. An efficient distributed B-tree index method in cloud computing. Open Cybernetics and Systemics Journal 8 (2014). Issue 1. https://doi.org/10.2174/1874110x01408010302
- [20] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Michael Burrows, Tushar Chandra, Andrew Fikes, and Robert Gruber. 2006. Bigtable: A Distributed Storage System for Structured Data (Awarded Best Paper!). In OSDI. USENIX Association, 205–218.
- [21] Surajit Chaudhuri and Vivek R. Narasayya. 1998. AutoAdmin 'What-if' Index Analysis Utility. In SIGMOD Conference. ACM Press, 367–378.
- [22] Leying Chen and Shimin Chen. 2021. How Does Updatable Learned Index Perform on Non-Volatile Main Memory?. In *ICDE Workshops*. IEEE, 66–71.
- [23] Shimin Chen, Phillip B. Gibbons, and Todd C. Mowry. 2001. Improving index performance through prefetching. SIGMOD Record (ACM Special Interest Group on Management of Data) 30 (2001). Issue 2. https://doi.org/10.1145/376284.375688
- [24] Shimin Chen, Phillip B. Gibbons, Todd C. Mowry, and Gary Valentin. 2002. Fractal prefetching B+-Trees: Optimizing both cache and disk performance. Proceedings of the ACM SIGMOD International Conference on Management of Data.
- [25] Douglas Comer. 1979. UBIQUITOUS B-TREE. Comput Surv 11 (1979). Issue 2.
- [26] Yifan Dai, Yien Xu, Aishwarya Ganesan, Ramnatthan Alagappan, Brian Kroth, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. 2020. From WiscKey to Bourbon: A Learned Index for Log-Structured Merge Trees. In OSDI. USENIX Association, 155–171.
- [27] Henry Daly, Ahmed Hassan, Michael F. Spear, and Roberto Palmieri. 2018. Nu-Mask: High performance scalable skip list for NUMA. *Leibniz International Proceedings in Informatics, LIPIcs* 121. https://doi.org/10.4230/LIPIcs.DISC.2018.18
- [28] Debabrata Dash, Neoklis Polyzotis, and Anastasia Ailamaki. 2011. CoPhy: A Scalable, Portable, and Interactive Index Advisor for Large Workloads. Proc. VLDB Endow. 4, 6 (2011), 362–372.
- [29] Ian Dick, Alan Fekete, and Vincent Gramoli. 2017. A skip list for multicore. Concurrency Computation 29 (2017). Issue 4. https://doi.org/10.1002/cpe.3876
- [30] Jialin Ding, Umar Farooq Minhas, Jia Yu, Chi Wang, Jaeyoung Do, Yinan Li, Hantian Zhang, Badrish Chandramouli, Johannes Gehrke, Donald Kossmann, David Lomet, and Tim Kraska. 2020. ALEX: An Updatable Adaptive Learned Index. Proceedings of the ACM SIGMOD International Conference on Management of Data. https://doi.org/10.1145/3318464.3389711
- [31] Jialin Ding, Umar Farooq Minhas, Jia Yu, Chi Wang, Jaeyoung Do, Yinan Li, Hantian Zhang, Badrish Chandramouli, Johannes Gehrke, Donald Kossmann, David B. Lomet, and Tim Kraska. 2020. ALEX: An Updatable Adaptive Learned Index. In SIGMOD Conference. ACM, 969–984.

- [32] Paolo Ferragina, Fabrizio Lillo, and Giorgio Vinciguerra. 2020. Why Are Learned Indexes So Effective?. In Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research), Hal Daumé III and Aarti Singh (Eds.), Vol. 119. PMLR, 3123–3132. https://proceedings.mlr.press/ v119/ferragina20a.html
- [33] Paolo Ferragina and Giorgio Vinciguerra. 2020. The PGM-index: a fully-dynamic compressed learned index with provable worst-case bounds. *Proc. VLDB Endow.* 13, 8 (2020), 1162–1175.
- [34] Martin R. Frank, Edward Omiecinski, and Shamkant B. Navathe. 1992. Adaptive and Automated Index Selection in RDBMS. In EDBT (Lecture Notes in Computer Science), Vol. 580. Springer, 277–292.
- [35] Goetz Graefe and Per Åke Larson. 2001. B-tree indexes and CPU caches. Proceedings - International Conference on Data Engineering (2001). https://doi.org/ 10.1109/ICDE.2001.914847
- [36] Jim Gray and Goetz Graefe. 1997. The Five-Minute Rule Ten Years Later, and Other Computer Storage Rules of Thumb. SIGMOD Rec. 26, 4 (1997), 63–68.
- [37] Jing He, Shao wen Yao, Li Cai, and Wei Zhou. 2018. SLC-index: A scalable skip list-based index for cloud data processing. *Journal of Central South University* 25 (2018). Issue 10. https://doi.org/10.1007/s11771-018-3927-0
- [38] Stratos Idreos, Kostas Zoumpatianos, Brian Hentschel, Michael S. Kester, and Demi Guo. 2018. The Data Calculator: Data Structure Design and Cost Synthesis from First Principles and Learned Cost Models. In SIGMOD Conference. ACM, 535–550.
- [39] Stratos Idreos, Kostas Zoumpatianos, Brian Hentschel, Michael S. Kester, and Demi Guo. 2018. The Internals of the Data Calculator. *CoRR* abs/1808.02066 (2018).
- [40] Ivo Jimenez, Huascar Sanchez, Quoc Trung Tran, and Neoklis Polyzotis. 2012. Kaizen: a semi-automatic index advisor. In SIGMOD Conference. ACM, 685–688.
- [41] Andreas Kipf, Thomas Kipf, Bernhard Radke, Viktor Leis, Peter Boncz, and Alfons Kemper. 2018. Learned cardinalities: Estimating correlated joins with deep learning. arXiv preprint arXiv:1809.00677 (2018).
- [42] Andreas Kipf, Ryan Marcus, Alexander van Renen, Mihail Stoian, Alfons Kemper, Tim Kraska, and Thomas Neumann. 2019. SOSD: A Benchmark for Learned Indexes. *NeurIPS Workshop on Machine Learning for Systems* (2019).
- [43] Andreas Kipf, Ryan Marcus, Alexander van Renen, Mihail Stoian, Alfons Kemper, Tim Kraska, and Thomas Neumann. 2020. RadixSpline: a single-pass learned index. In *aiDM@SIGMOD*. ACM, 5:1–5:5.
- [44] Tim Kraska, Mohammad Alizadeh, Alex Beutel, Ed H. Chi, Ani Kristo, Guillaume Leclerc, Samuel Madden, Hongzi Mao, and Vikram Nathan. 2019. SageDB: A Learned Database System. In CIDR. www.cidrdb.org.
- [45] Tim Kraska, Alex Beutel, Ed H. Chi, Jeffrey Dean, and Neoklis Polyzotis. 2018. The case for learned index structures. Proceedings of the ACM SIGMOD International Conference on Management of Data. https://doi.org/10.1145/3183713.3196909
- [46] Kenneth Lange. 2016. MM optimization algorithms. SIAM.
- [47] Pengfei Li, Hua Lu, Qian Zheng, Long Yang, and Gang Pan. 2020. LISA: A Learned Index Structure for Spatial Data. In SIGMOD Conference. ACM, 2119–2133.
- [48] LMDB. [n.d.]. Lightning Memory-Mapped Database Manager. http://www.lmdb. tech/doc/ Online; accessed Jul-17-2022.
- [49] David B. Lomet. 1998. B-tree Page Size When Caching is Considered. SIGMOD Rec. 27, 3 (1998), 28–32.
- [50] Baotong Lu, Jialin Ding, Eric Lo, Umar Farooq Minhas, and Tianzheng Wang. 2021. APEX: A High-Performance Learned Index on Persistent Memory. arXiv preprint arXiv:2105.00683 (2021).
- [51] Vincent Y. Lum and Huei Ling. 1971. An optimization problem on the selection of secondary keys. In ACM Annual Conference. ACM, 349–356.
- [52] Ryan Marcus, Andreas Kipf, Alexander van Renen, Mihail Stoian, Sanchit Misra, Alfons Kemper, Thomas Neumann, and Tim Kraska. 2020. Benchmarking Learned Indexes. Proc. VLDB Endow. 14, 1 (2020), 1–13.
- [53] Ryan Marcus, Emily Zhang, and Tim Kraska. 2020. CDFShop: Exploring and Optimizing Learned Index Structures. In SIGMOD Conference. ACM, 2789–2792.
- [54] Patrick E. O'Neil, Edward Cheng, Dieter Gawlick, and Elizabeth J. O'Neil. 1996. The Log-Structured Merge-Tree (LSM-Tree). Acta Informatica 33, 4 (1996), 351– 385.
- [55] Postcard. [n.d.]. Postcard: A no_std + serde compatible message library for Rust. https://github.com/jamesmunns/postcard. [Online; accessed July-17-2022].
- [56] PostgreSQL. [n.d.]. PostgreSQL: The World's Most Advanced Open Source Relational Database. https://www.postgresql.org. [Online; accessed July-17-2022].
- [57] William Pugh. 1990. Skip lists: a probabilistic alternative to balanced trees. Commun. ACM 33, 6 (1990), 668–676.
- [58] Ohad Rodeh, Josef Bacik, and Chris Mason. 2013. BTRFS: The Linux B-tree filesystem. ACM Transactions on Storage (TOS) 9, 3 (2013), 1–32.
- [59] Mario Schkolnick. 1975. The Optimal Selection of Secondary Indices for Files. Inf. Syst. 1, 4 (1975), 141–146.
- [60] Serde. [n.d.]. Serde. https://serde.rs. [Online; accessed July-17-2022].
- [61] Facebook Open Source. [n.d.]. RocksDB: A persistent key-value store. https: //rocksdb.org/. [Online; accessed July-17-2022].

- [62] Stefan Sprenger, Steffen Zeuch, and Ulf Leser. 2017. Cache-sensitive skip list: Efficient range queries on modern CPUs. Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) 10195 LNCS. https://doi.org/10.1007/978-3-319-56111-0_1
- [63] SQLite. [n.d.]. SQLite. https://www.sqlite.org. [Online; accessed April-24-2021].
 [64] Mihail Stoian, Andreas Kipf, Ryan Marcus, and Tim Kraska. 2021. PLEX: Towards Practical Learned Indexing. *CoRR* abs/2108.05117 (2021).
- [65] Michael Stonebraker. 1974. The choice of partial inversions and combined indices. Int. J. Parallel Program. 3, 2 (1974), 167–188.
- [66] Yongjoo Park Supawit Chockchowwat, Wenjie Liu. 2023. AirIndex: Versatile Index Tuning Through Data and Storage (Extended Version). arXiv preprint (2023).
- [67] Alexandre Verbitski, Anurag Gupta, Debanjan Saha, Murali Brahmadesam, Kamal Gupta, Raman Mittal, Sailesh Krishnamurthy, Sandor Maurice, Tengiz Kharatishvili, and Xiaofeng Bao. 2017. Amazon aurora: Design considerations for high throughput cloud-native relational databases. In Proceedings of the 2017 ACM International Conference on Management of Data. 1041–1052.
- [68] Li Wang, Zining Zhang, Bingsheng He, and Zhenjie Zhang. 2020. PA-Tree: Polledmode asynchronous B+ tree for NVMe. Proceedings - International Conference on Data Engineering 2020-April. https://doi.org/10.1109/ICDE48307.2020.00054
- [69] Youyun Wang, Chuzhe Tang, Zhaoguo Wang, and Haibo Chen. 2020. SIndex: A Scalable Learned Index for String Keys. In Proceedings of the 11th ACM SIGOPS Asia-Pacific Workshop on Systems (Tsukuba, Japan) (APSys '20). Association for Computing Machinery, New York, NY, USA, 17–24. https://doi.org/10.1145/ 3409963.3410496
- [70] Xingda Wei, Rong Chen, and Haibo Chen. 2020. Fast RDMA-based Ordered Key-Value Store using Remote Learned Cache. In 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20). USENIX Association, 117–135. https://www.usenix.org/conference/osdi20/presentation/wei
- [71] Sai Wu, Dawei Jiang, Beng Chin Ooi, and Kunlung Wu. 2010. Efficient btree based indexing for cloud data processing. *Proceedings of the VLDB Endowment* 3 (2010). Issue 1. https://doi.org/10.14778/1920841.1920991
- [72] Wei Zhou, Jin Lu, Zhongzhi Luan, Shipu Wang, Gang Xue, and Shaowen Yao. 2014. SNB-index: A SkipNet and B+ tree based auxiliary Cloud index. *Cluster Computing* 17 (2014). Issue 2. https://doi.org/10.1007/s10586-013-0246-y

A AIRINDEX IN DETAIL

A.1 Layer Builders

A layer builder turns the target key-position collection *D* into an index layer Θ . In other words, it is a function $F(D) = \Theta$ where $D = \{(x_i, y_i)\}_{i=1}^n$ is a collection of keys x_i together with their positions y_i , and $\Theta = \{(z_j, \theta_j)\}_{j=1}^{n^+}$ is the index layer containing the per-node parameters θ_j with per-node key range $[z_j, z_{j+1})$. To build two types of nodes (step and band, sketched in Figure 6), AIRINDEX currently deploys three types of layer builders below.

- (1) Greedy Step (GStep(p, λ_{GS})) builds a *p*-piece step function with precision at most λ_{GS} bytes. It iterates over each keyposition pair (x_i, y_i) and greedily determines whether to create the next constant function if y_i⁺ b_k > λ_{GS} where (a_k, b_k) represents the current constant function. If so, the next constant function has a partition key a_{k+1} = x_i and partition position b_{k+1} = y_i. Once it reaches *p* pieces of constant functions, GStep(p, λ_{GS}) generates a *p*-piece step node.
- (2) **Greedy Band** (GBand(λ_{GB})) builds linear band nodes by greedily fitting as many key-position pairs as possible, using a method called monotone chain convex hull [10], until we have $\Delta(x_i) > \lambda_{GB}$. Then, it generates a band node with $(x_1, y_1, x_2, y_2, \delta)$. For *n* key-position pairs, the convex hull inserts a key-position pair in O(n) time (O(1) amortized insertion) and answers a feasibility query in $O(n \log m)$ time ($O(\log m)$ amortized query) where *m* is the average number of key-position pairs included in one linear band. Typically, *m* is small when λ_{GB} is small.
- (3) Equal Band (EBand(λ_{EB})) builds linear band nodes by grouping key-position pairs in equal-size position ranges. That is, each group {(x_l, y_l), (x_{l+1}, y_{l+1}), ..., (x_r, y_r)} has a bounded

position range $|y_l^- - y_r^+| \le \lambda_{EB}$. It then fits a linear band function to each group and creates a band node. Note that the precision $\Delta(x_i)$ can vary depending on how "linear" the group is. EBand(λ_{EB}) groups by position ranges rather than key ranges so that the worst-case precision is controlled by λ_{EB} .

First, $GStep(p, \lambda_{GS})$ is equivalent to bulk indexing in a sparse B-tree with a fanout p and page size λ_{GS} bytes. Second, $GBand(\lambda_{GB})$ is a generalization from step functions to linear functions with precision $\Delta(x_i) \leq \lambda_{GB}$. Third, $EBand(\lambda_{EB})$ is another generalization focusing on the key-position group size $|y_I^- - y_r^+| \leq \lambda_{EB}$.

Granularity Exponentiation. λ_{GS} , λ_{GB} , and λ_{EB} are called *granularity*, which roughly control layer builders' tendency to split key-position pairs apart. Because they correlate with the resulting node's precision $\Delta(x; \Theta)$, AIRINDEX needs to determine the appropriate granularity for each node type.

AIRINDEX creates many candidate layer builders \mathcal{F} by sampling the granularity on an exponential grid: λ_{low} , $\lambda_{low}(1 + \epsilon)$, $\lambda_{low}(1 + \epsilon)^2$, ..., λ_{high} where (λ_{low} , λ_{high}) are the bounds and $\epsilon > 0$ controls the exponentiation base. Smaller ϵ implies a finer search which improves optimization accuracy but increases tuning time.

A.2 Cache

AIRINDEX uses a read-through cache in local memory. That is, if a read range is present in the cache, AIRINDEX gets the layer view directly from it. Otherwise, AIRINDEX reads from storage and fills in the corresponding cache page(s). Filling a cache page and getting a layer view are both zero-copy operations: they do not copy the raw bytes but only their references. As the cache fills up, AIRINDEX avoids more expensive storage reads, offering faster lookup speed; nevertheless, to offer consistent performance, our optimization minimizes the worst-case latency when there is no cached data.

A.3 Index Complexity

Before branching out, AIRINDEX selects only top-*k* candidates with the highest potential to be in the optimal design. Suppose there exists an oracle that minimizes the objective function on a keyposition collection *D* under the storage profile. We call the optimal search cost under storage profile *T*: *index complexity* $\tau(D;T)$. If $\tau(D;T)$ is known, selection of candidates $C = \{(\Theta_i, D_i)\}_{i=1}^{|\mathcal{F}|}$ would be as straightforward as Eq (11). That is, we could simply select the top-1 candidate and avoid branching out.

$$(\Theta^*, D^*) = \underset{(\Theta_i, D_i) \in \mathcal{C}}{\arg\min} \tau(D_i) + \underset{x \sim \mathcal{X}}{\operatorname{E}} \left[T(\Delta(x; \Theta_i)) \right]$$
(11)

Unfortunately, $\tau(D;T)$ is unknown for a large class of indexes supported by AIRINDEX. Instead, AIRINDEX uses a surrogate upper bound to the index complexity: *step index complexity* $\hat{\tau}(D;T) \ge$ $\tau(D;T)$. Specifically, $\hat{\tau}(D;T)$ imitates a step-function hierarchical index that partitions the key-position collection D into groups with arbitrary position ranges. Let $s_D = y_{|D|}^+ - y_1^-$ be the total size of the candidate layer, $\hat{\tau}(D;T)$ tries to build ideal step indexes with different numbers of layers $L \in \{0, 1, \ldots, O(\log s_D)\}$. For each L, it perfectly balances both root layer size and subsequent precisions so that $s(\Theta_L) = \Delta(x;\Theta_l) = {}^{L+1}\sqrt{s_D s_{step}^L}$. This considers the total size s_D and the ideal size for each 1-piece step node s_{step} (e.g., 16 bytes for 8-byte key and position types). Lastly, $\hat{\tau}(D;T)$ then outputs the

Supawit Chockchowwat, Wenjie Liu, Yongjoo Park



Figure 17: Latency breakdown of B-TREE and AIRINDEX by the time spent reading different layers. The books dataset is used. We vary the warmness by the number of queries from one to 1 million from left to right in logarithmic scale. Layer-wise latencies are stacked from bottom up in their retrieval order (root index layer to data layer). For visual purposes, the plots then interpolate linearly and fill areas in between with different colors indicating different layers.



Figure 18: Latency breakdown of AIRINDEX's non-I/O operations over different levels of warmness (indicated by the numbers of queries from one to 1 million). Left: the fraction of non-I/O operations—the rest is for I/O operations. Middle and right: a deeper breakdown of latency spent on different types of I/O operations on NFS and SSD, respectively.

lowest storage model cost as the step index complexity (Eq (12)).

$$\hat{\tau}(D;T) = \min_{L \in \{0,1,\dots,O(\log s_D)\}} (L+1) \times T\left(\sqrt[L+1]{s_D s_{step}^L}\right)$$
(12)

Note that $\hat{\tau}(D; T)$ is only interested in the integer s_D (not the distribution of D); thus, $\hat{\tau}(D; T)$ can be arithmetically computed (hence cheap). Figure 8 shows the general shape of $\hat{\tau}(D; T)$ with respect to the collection size s_D , under an affine storage profile T. Notice the sudden index complexity cliffs (such as those around $s_D = 1$ MB in Figure 8b), marking the boundaries between different chosen numbers of layers. The technique to minimize an objective based on its cheaper upper bound is related to majorize-minimization algorithms [46].

B BASELINES

We compare AIRINDEX to a traditional database index, learned indexes, and our manual configuration counterpart, hosted at their respective forks ².

LMDB. LMDB [48] is a B-tree database that accesses its data on storage through mmap. We have also tested PostgreSQL [56] and RocksDB [61] but decided to present LMDB due to its competitive performance in our setting.

RMI. RMI [41] is a top-down learned index with a compact twolayer structure where the top one contains only one perfectly accurate node partitioning key space to the bottom nodes. To build one, we execute its provided optimizer for each dataset and select the most accurate configuration describing index size and model types. We then integrate RMI onto external memory setting by mmap-ing its parameter arrays so RMI can access its parameters through the OS buffer cache.

PGM-INDEX. PGM-INDEX [33] is a learned index with bounded precision across all layers. PGM-INDEX partitions the key-position collection to build the next bottom-most layer towards the top. To benchmark in our settings, we use the MappedPGMIndex variant that operates on file systems. Although its tuner does not satisfy our target (fastest index, regardless of size), we adjust its error level $\varepsilon \in \{16, 32, ..., 1024\}$ to microbenchmark on wiki (chosen one arbitrarily) and finally pick $\varepsilon = 128$.

ALEX/APEX. ALEX/APEX [31] is an updatable learned index built top-down like RMI but further arrange key-value pairs in its layout (notably, "gapped array" to buffer structural changes). We integrate ALEX/APEX onto external memory setting by mmap-ing its serialized node objects, key arrays, and value arrays.

PLEX. PLEX [64] is a learned index with compact Hist-Tree (CHT) layered on top of RadixSpline [43] (RS). We integrate PLEX onto external memory setting through mmap similarly to ALEX/APEX. Although PLEX optimizes most parameters, its user need to specify the maximum prediction error ϵ . We select $\epsilon = 2048$ based on a benchmark on a setting (Figure 12d).

DATA CALCULATOR. DATA CALCULATOR [38, 39] is a data layout design engine that calculates the performance of a data structure. For index learning, we extend DATA CALCULATOR to autocomplete recursion allowed (number of layers), fanout, key partitioning. By following its cost synthesis flowcharts (Fig 5 in [38] and Fig 29 in [39]), we identify data access primitives, profile them on SSD and NFS for cost models, and execute a parallelized auto-completion flow. Later, the selected data layout is built within AIRINDEX's framework.

²https://github.com/illinoisdata/lmdb, https://github.com/illinoisdata/RMI, https://github.com/illinoisdata/PGM-index, https://github.com/illinoisdata/ALEX_ext, https://github.com/illinoisdata/arindex-public/tree/main/src/bin/data_calculator.rs



Figure 19: Effects of skewed query key distribution latency across different methods on books dataset. Left: first query latency over skewnesses 0.5, 1.0, and 2.0. Right: 100th latency over the same skewnesses.

B-TREE. A B-tree-like structure implemented using AIRINDEX's framework. It has 255-piece step nodes built with GStep($p = 255, \lambda_{GS} = 4096$), which is equivalent to a B-tree with 4KB node pages and 255 fanout factors. This is the most controlled baseline where the only difference is AIRINDEX's storage- and data-aware tuning.

C EXTENDED EXPERIMENTS

C.1 Latency Breakdown

We investigate the effect of our optimization by studying latency breakdowns, by layers and by operations, respectively.

By Layer. We first decompose the end-to-end latency into the times spent in retrieving each layer which includes I/O, cache read, and all internal computation. Apart from some small exceptions, latency measurements across layers roughly follow the storage profiles: larger root size $s(\Theta_L)$ and coarser precision $\Delta_l(x; \Theta_l)$ reflects in a longer latency spent in the corresponding layer. Over numbers of queries, we also observe the alternating acceleration phenomenon in more details. That is, the fast acceleration region indicate that the topmost partially cached index layer is becoming fully cached. For example, in Figure 17c between 10^3 and 10^4 queries, B-TREE searches faster because the speedup in its layer-1 index.

This breakdown also reveals factors unaccounted for in AIRINDEX. For a prominent example, in the first query under SSD (Figures 17c and 17d), both AIRINDEX and B-TREE spend more time reading their root and data layers as opposed to reading other index layers. This is because index layers and data layer are stored in separate directories, forcing the file system (ext4) to slowly walk through different paths of directory entries (dentry) to fetch index-layer and datalayer inodes. Consequently, reading subsequent index layers stored in the same directory is then significantly faster than expected. Although these missing characteristics are crucial for future works towards maximally fast indexes, we believe that AIRINDEX's storage profile is at an appropriate level of abstraction to adapt to diverse types of storage.

By Operation. We categorize AIRINDEX's operations into I/O and non-I/O groups. The latter group contains in-memory caching, data structure deserialization, node prediction, relevant key-value finding, and other negligible steps. Figure 18a shows that these non-I/O operations only account for up to 1.0% on NFS and 9.0% on SSD, even after 1 million queries. These small non-I/O fractions matches with our expectation. Because of its 4KB–5KB average





Figure 20: Effects of hyperparameter k in selecting top-k candidates to overall AIRINDEX's building time and cost \mathcal{L} . We use books dataset and SSD profile.

precisions of tuned structures on NFS and SSD, AIRINDEX needs more than a million query to completely cache the 6.4GB data layer. As an improvement on its warm-state performance, AIRINDEX can prefetch to warm up its cache more aggressively.

Figures 18b and 18c delve deeper into non-I/O operations to put their significance into perspective. We found that caching occupies non-I/O latencies increasingly over warmness. In the first query, deserialization dominates because it needs to deserialize the rootmetadata file, send the root layer to cache, and reconstruct a nested data structure for query processing. These two observations explain the U-shape fraction of non-I/O operations on SSD. Aside from those, finding operations (i.e. the last mile binary search) take a considerable portion within non-I/O latency, but incomparably small compared to I/O. If they were to grow larger, AIRINDEX would have to consider smoothly transitioning to in-memory index rather than a binary search on cache. Lastly, as expected, node prediction is insignificant given that our current node types are simple.

C.2 Skewed Workload

Although AIRINDEX's objective (Eq (7)) considers the query distribution X, future query distributions may change unexpectedly. This experiment (Figure 19) builds AIRINDEX on the uniform distribution X but requests keys sampled from a Zipf distribution with parameters 0.5 (least skewed), 1.0, and 2.0 (most skewed). The more skewed the query is, the faster all methods can respond at warm state (Figure 19b). However, the skew does not affect first-query latency as much across all methods (Figure 19a). Because AIRINDEX is tuned for cold-state latency, higher skewness results in a quicker takeover. For example, it takes 12k uniform queries for any methods (PGM-INDEX) to take over AIRINDEX, but only 70, 41, and 725 queries in 0.5, 1.0, and 2.0 Zipf query.

C.3 Top-k Candidate Parameter Sweep

Across all experiments, we set k = 5 as an arbitrary constant greater than one and less than the number of node builders ($|\mathcal{F}| = 45$). In this experiment, we vary this hyperparameter k to verify our understanding: as k increases, build time should increase in L-degree polynomial (L = 2 in this setting) while the optimized cost should monotonically decrease. Figure 20 reaffirms this hypothesis, but also shows that the available parallelism (192 CPUs) is able to hide the polynomial build time more than we had expected (taking around 50 seconds up until k = 20), implying that we could have selected a higher k to get a faster index at no additional build time cost.

Supawit Chockchowwat, Wenjie Liu, Yongjoo Park

Received 15 January 2023; revised 20 April 2023; accepted 23 May 2023