A-Tree: A Bounded Approximate Index Structure

Alex Galakatos¹ Michael Markovitch¹ Carsten Binnig² Rodrigo Fonseca¹ Tim Kraska³ ¹Brown University ²TU Darmstadt ³MIT

ABSTRACT

Index structures are one of the most important tools that DBAs leverage in order to improve the performance of analytics and transactional workloads. However, with the explosion of data that is constantly being generated in a wide variety of domains including autonomous vehicles, Internet of Things (IoT) devices, and E-commerce sites, building several indexes can often become prohibitive and consume valuable system resources. In fact, a recent study has shown that indexes created as part of the TPC-C benchmark can account for 55% of the total memory available in a stateof-the-art in-memory DBMS. This overhead consumes valuable and expensive main memory, and limits the amount of space that a database has available to store new data or process existing data.

In this paper, we present a novel approximate index structure called A-Tree. At the core of our index is a tunable error parameter that allows a DBA to balance lookup performance and space consumption. To navigate this tradeoff, we provide a cost model that helps the DBA choose an appropriate error parameter given either (1) a lookup latency requirement (e.g., 500ns) or (2) a storage budget (e.g., 100MB). Using a variety of real-world datasets, we show that our index structure is able to provide performance that is comparable to full index structures while reducing the storage footprint by orders of magnitude.

ACM Reference format:

Alex Galakatos¹ Michael Markovitch¹ Carsten Binnig² Rodrigo Fonseca¹ Tim Kraska³

¹Brown University

 $^2\mathrm{TU}$ Darmstadt

. 2018. A-Tree: A Bounded Approximate Index Structure. In Proceedings of SIGMOD, Houston, Texas USA, June 2018 (SIGMOD'18), 14 pages.

DOI: 10.1145/nnnnnnn.nnnnnnn

INTRODUCTION

Tree-based index structures such as B+ trees are one of the most important tools that DBAs leverage in order to improve the performance of analytics and transactional workloads. However, for main-memory databases, tree-based indexes can often consume a significant amount of main memory. In fact, a recent study [38] shows that the indexes created for typical OLTP workloads can

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.

SIGMOD'18, Houston, Texas USA

DOI: 10.1145/nnnnnnn.nnnnnnn

© 2018 ACM. 978-x-xxxx-xxxx-x/YY/MM...\$15.00

consume up to 55% of the total memory available in a state-of-theart in-memory DBMS. This overhead not only limits the amount of space that a database has available to store new data but it also reduces space for intermediates that can be helpful when processing existing data.

In an attempt to reduce the storage overhead of B+ trees, various compression schemes have been developed [2, 4, 13, 39]. The main idea of all of these techniques is to remove the redundancy that exists among keys and/or to reduce the size of each key inside each node of the index. For example, prefix and suffix truncation can be used to store common prefixes/suffixes of keys only once per index node, reducing the total size of the tree. Additionally, more expensive compression techniques like Huffmann coding can be applied within each node. However, these techniques come at a higher runtime cost since pages must be decompressed in order search for an item.

Although all of the previously mentioned compression schemes reduce the storage footprint of an index, these indexes still contain one entry (key and pointer) for each distinct value in the dataset, resulting in indexes that still consume a significant amount of available memory if the dataset contains many distinct keys. This observation especially holds true for data types such as timestamps or sensor readings that are generated in a wide variety of modern applications including autonomous vehicles, Internet of Things (IoT), and e-commerce sites, but also other data types such as geo-coordinates or string data that have similar properties. Even worse, the number of unique values for such data types typically grow over time, resulting in indexes that are constantly growing. Consequently, a DBA has no efficient way to restrict the memory consumption other than dropping an index completely.

In this paper, we present A-TREE, a novel approximate index structure that can be bounded in space independent of the number of distinct values that need to be indexed. The main idea is that our index does not contain every unique key found in the data set. Instead, A-Tree only indexes a fraction of these keys and uses interpolation search between two keys. In its original form, A-Tree assumes that the data to be indexed is sorted. To that end, our index is similar to a clustered index that uses fixed-sized pages (segments) whereby only the first key of each page is inserted into the index. However, instead of splitting the data into fixed-size segments, we dynamically partition the data into variable-sized segments whose data is roughly linear. By taking advantage of trends that exist in the data, an A-Tree can help reduce the memory consumption of an index by orders of magnitude compared to a traditional B+ tree, while being able to bound the error associated with the approximation. Interestingly, this approach amounts to approximating the data using piece-wise linear functions, as we discuss later.

At the core of our index structure is a parameter that specifies the amount of acceptable error (i.e., the maximum distance between the predicted and actual position of a key when using interpolation search inside a segment). This tunable parameter allows us to balance the lookup performance and the space consumption of our index, providing DBAs indexes that can be highly performant (i.e., low-latency lookups) or have a low space footprint. To navigate this tradeoff, we also present a cost model that helps a DBA choose an appropriate error term given either (1) a lookup latency requirement (e.g., 1000ns) or (2) a storage budget (e.g., 100MB). Using a variety of real-world datasets, we show that our index structure is able to provide performance that is comparable to full and fixed-paged indexes while reducing the storage footprint by orders of magnitude.

Another interesting observation is that our variable-sized paging approach is orthogonal to node-level compression techniques such as the previously mentioned prefix- or suffix truncation. In other words, since our index internally uses a tree structure to find the page that a given key belongs to, we can still apply these well-known compression techniques to further reduce an index's size. Other techniques such as adaptive indexing [16] and database cracking [15] also reduce the number of keys in an index by building fine-grained indexes for regions that are queried more often. To that end, these techniques are also orthogonal since our technique leverages the data distribution, instead of knowledge about the workload, to compress the index size.

In summary, we make the following contributions:

- We propose A-TREE, a novel index structure that leverages properties about the underlying data distribution in order reduce the size of an index.
- We discuss an efficient one-pass bulk loading technique that creates a compressed A-TREE. Our bulk loading algorithm incorporates a tunable error parameter, which balances the lookup and space footprint of our index.
- We present a cost model that helps a DBA determine an appropriate error threshold given either a latency or storage requirement.
- Using several real-world datasets, we show that our index is able to provide comparable (or in some cases even better) performance compared to existing index structures while consuming orders of magnitude less space.

The remainder of the paper is organized as follows. In Section 2, we first present an overview of our new index structure called A-Tree. Afterwards, we discuss the main index operations: bulk loading (Section 3), lookups (Section 4) and insertion (Section 5). Section 6 then elaborates on the cost model, that allows a DBA to navigate the tradeoff of lookup performance and the space consumption of A-Tree indexes. Finally, in Section 7, we discuss the results of our evaluation on real and synthetic datasets, we discuss related work in Section 8, and finally conclude in Section 9.

2 OVERVIEW

At a high level, indexes (and B+ trees over sorted attributes in particular) can be represented by a function that map key values (e.g., a timestamp) to storage location. Using this representation, A-Tree partitions the key space into a series of disjoint linear segments that approximate the true function. Instead of storing all values in

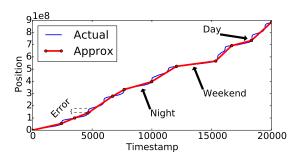


Figure 1: Timestamp to Position Mapping for IoT Dataset

the key space, A-Tree stores only (1) the keys at the boundaries and (2) the slope of the linear function in order to compute a key's approximate position using linear interpolation. At the core of this process is a tunable error threshold, which represents the maximum distance that the predicted location of any key inside a segment is from its actual location.

In the following, we first discuss how the idea of using functions to map key values to storage locations intuitively works. Next, we discuss how we leverage this function representation to efficiently implement our approximate index structure on top of a B+ tree for clustered indexes (over a sorted attribute). Finally, we then present how our ideas can also be applied to compress indexes over non-clustered attributes (i.e., secondary indexes).

2.1 Function Representation

One key insight to our approach is that we can abstractly model an index as a monotonically increasing function that maps keys (i.e., values of the indexed attribute) to storage location (i.e., its page and the offset within that page). In order to explain this intuition, we first assume that all keys to be indexed are stored in a sorted array and the storage location can therefore be abstracted as a index into the sorted array. In the following section, we will discuss how this translates into a tree-index where keys are stored in pages and the storage location is given by a page identifier and the offset of the key within the page.

For example, consider an Internet of Things (IoT) dataset, which contains events from various devices (e.g., door sensors, motion sensors, power monitors) installed throughout a university building. In this dataset, we store the data sorted by the timestamp of an event (e.g., door opened, motion detected) and construct a function that maps each timestamp (i.e., the key) to its position in the dataset (i.e., the index in a sorted array), as shown in Figure 1. Unsurprisingly, since the installed IoT devices monitor human activity, the timestamps of the recorded actions follow a pattern. As shown, during the weekend and night hours, there is little activity, corresponding to large timestamp intervals spanning relatively few positions in the sorted list. On the other hand, due to the increased level of activity during the daytime hours, a small range of timestamps spans a large number of positions.

Since a function that represents an index can be arbitrarily complex and data-dependent, maintaining the precise function that maps keys to positions may not be possible to learn and is expensive to build and update. Therefore, our goal is to approximate the function that represents the mapping of a key to a position. By

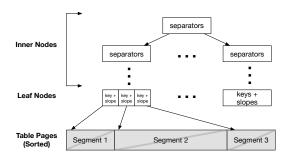


Figure 2: A clustered A-Tree index

approximating this function, the predicted location for a given key does not necessarily point to the true location of the key.

To be able to compactly capture trends that exist in the data while being able to efficiently build a new index and handle updates, we use a series of piece-wise linear functions to approximate an arbitrary function. More formally, to approximate any function that represents an index structure, we partition the function's domain into variable length consecutive disjoint partitions and represent each partition using a linear function. As shown in Figure 1, for example, our segmentation algorithm partitions the timestamp values into several linear segments that are able to accurately reflect the various trends that exist in the data (e.g., less activity during the weekend). In Section 3, we describe our segmentation algorithm in more detail.

Although piece-wise linear approximation can help us model any arbitrary function, the resulting function is not precise (i.e., a key's predicted location is not guaranteed to be its true position). We therefore define the *error* associated with our approximation as the maximum distance between the actual and predicted location of any key, as shown below, where $pred_pos(k)$ and $true_pos(k)$ return the predicated and actual position of an element k in sorted array respectively.

2.1.1 Clustered Indexes.

$$error = \max(|pred_pos(k) - true_pos(k)|) \ \forall \ k \in keys$$
 (1)

This formulation allows us to concretely define the core building block that A-Tree uses, a *segment*. A segment is a contiguous region of a sorted array for which any point is no more than a specified error threshold from its interpolated position.

Depending on the data distribution and the error threshold, the segmentation process will yield either a lower or higher number of segments that approximate the true function. Therefore, importantly, the error threshold enables us to balance main memory consumption and performance of the index. After the segmentation process, A-Tree stores the boundaries and slope of each segment in a B+ tree (instead of each individual key), reducing the overall memory footprint of the index.

2.2 A-Tree Design

As previously mentioned, our segmentation process partitions the key space of an attribute into disjoint linear segments such that the predicted position of any key inside a segment is no more than a bounded distance from the key's true position. A-Tree then

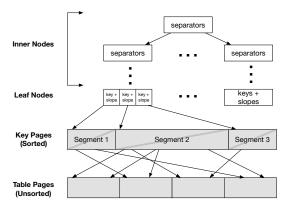


Figure 3: A non-clustered A-Tree index

organizes these segments in a tree to efficiently support insert and lookup operations.

In the following, we discuss clustered indexes, where records are already sorted by the key that is being indexed. Afterwards, we show how our technique can be extended to provide similar benefits for secondary indexes.

In a traditional clustered B+ tree (e.g., an index-organized table), the table data is stored in fixed-sized pages and the leaf level of the index contains only the first key of each of these tables pages. Unlike a clustered B+ tree, in a clustered A-Tree, the table data is partitioned into variable sized segments that satisfy the given error threshold.

Figure 2 shows the general structure of a clustered A-Tree index. Depending on the error parameter and the data distribution, several consecutive keys can thus be summarized into a single segment. As we show in our experiments that use several real-world data sets, our index is able to effectively reduce the storage footprint of an index by orders of magnitude without sacrificing performance. Details of the segmentation algorithm that divides the table data into variable-sized segments (pages) are discussed in Section 3.

Another important difference to a traditional B+ tree is the data that A-Tree needs to store in each leaf level entry. More specifically, to enable interpolation search on top of the variable-sized table pages, we store the segment's slope, starting key, and a pointer to the table page. This allows us to implement interpolation search on each table page since the data within this page is approximated by a linear function that is given by the slope.

The inner nodes of the index are the same as a normal B+ tree (i.e., lookup and insert operations are identical to a normal B+ tree). However, once a lookup or an insert reaches the leaf level, A-Tree needs to perform additional work. For lookups, we need to use the slope and the distance to the starting key in order to calculate the key's approximate position (offset in the array). Since the resulting position is approximate, A-Tree must perform a local search (e.g., binary, linear) to find the item. Details about looking up values in an A-Tree are discussed in Section 4.

Insert operations require some additional work as well upon reaching the leaf level of an A-Tree. First, each table page contains additional space to accommodate new keys. However, unlike a normal B+ tree, we must ensure that the error threshold is satisfied, even after adding new keys. Second, once the table page is full, it

needs to be split. Instead of dividing the node into two equal sized pages, we need to run our segmentation algorithm over data in the segment to find the optimal splitting point(s) that satisfy the error constraint.

Since the table pages may be quite large (due to their variable size), it is possible that inserts can very expensive. In the ideal case when data is completely uniform (in terms of compression), we only need one large segment to index the data. To tackle this problem, we insert new tuples into a fixed-size buffer and merge the buffer frequently into the sorted segments. A more through discussion of how A-Tree handles inserts is presented in Section 5.

Finally, instead of internally using a standard B+ tree to be able to efficiently find the segment that a given key belong to, A-Tree could instead use any other tree-based index structure. For example, if the workload is read-only, other index structures such as the FAST tree [19] could be used.

2.2.1 Non-clustered Indexes. Secondary indexes can dramatically improve the performance of queries involving selections over a non-primary key attribute. Without secondary indexes, these queries would need to examine all tuples in the table, which can often be prohibitive. However, unlike a clustered index, a non-primary key attribute is not sorted and might contain duplicates.

The main difference between a non-clustered (secondary) A-Tree index and a clustered A-Tree index is that a non-cluster index introduces an additional level that stores all indexed values of the non-primary column in sorted order, as shown in Figure 3 (called key pages). This level is then divided into variable-sized pages based on the error parameter using the same segmentation strategy for a clustered index. Each key page stores the keys of the attribute together with a pointer to the attribute in the table page.

The sorted level of key pages in a non-clustered A-Tree index introduces additional overhead compared to a clustered A-Tree index. However, this overhead occurs in any non-clustered (secondary) index. More importantly, as we will show in our experiments, a non-clustered A-Tree yields an index that has significantly fewer leaf and internal nodes when compared to a classical non-clustered B+ tree with fixed-size pages, reducing the storage footprint of a secondary index. Furthermore, by using different error parameters, the DBA can trade memory consumption for runtime as our cost model (Section 6) shows.

3 BULK LOADING

When building an index over a new attribute, A-Tree begins by partitioning the key space into disjoint segments, each of which can be represented by a linear function. However, to bound the time required to locate an element, we add the additional restriction (i.e., an error threshold) that no element's position is more than a constant distance from its predicated location in the segment (determined through linear interpolation).

In the following, we describe how A-Tree partitions the key space of an attribute into segments that satisfy the specified error threshold. As previously mentioned, after this process, each segment is inserted into a standard B+ tree with the starting key of each segment as the key and a pointer to the segment along with the segment's slope as the value. Later, in Section 5, we show how to insert new items into an existing A-Tree.

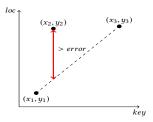


Figure 4: A segment from (x_1, y_1) to (x_3, y_3) is not valid if the point (x_2, y_2) further than *error* from the interpolated line.

3.1 Design Choices

Recall that we can abstractly model an index as a function that maps keys to positions. Therefore, our goal is to approximate this function using a series of piece-wise linear approximations. After using a tree structure to find the segment for a given key, we can use interpolation search (since segments are approximately linear) to find the key's estimated position inside the segment. In the following, we address how to efficiently segment a dataset into a series of linear approximations for a given error threshold.

A common objective function when fitting a function to a dataset is to minimize the least square error (minimizing the second error norm E_2). However, such an objective does not suit our needs: it does not provide a guarantee as to the maximal error and therefore does not provide a bound on the number of locations which must be scanned after interpolating a key's position. This realization leads to an important conclusion: the objective is to satisfy a maximal error (E_∞) , demonstrated in Figure 4, which is not necessarily compatible with minimizing the least square error.

As previously mentioned, to approximate the function that represents an index, we employ piece-wise linear approximation. In this setting, a segment is a region of the key space that can be represented by a single linear function whereby all keys inside are within a bounded distance from their lineally interpolated position. The simplest way to define a segment is to select the first point (first key) and the last point (the last key) in the segment. Using this simple segment definition, we can fit a linear function to the locations of keys in the segment (using the start and the end of the segment as well as the number of positions). We can then interpolate the approximate location of the key (using the segment's slope), and scan from predicted location of the key.

Our objective when segmenting the data is to satisfy a maximal error (i.e., the starting point predicated by the interpolation search is at most *error* number of elements away from the real position). This leads to an important property of a maximal segment (a segment is maximal when the addition of a key will violate the specified error):

THEOREM 3.1. The minimal number of locations covered by a maximal linear segment is error + 1.

PROOF. Consider 3 arbitrary points in a 2-dimensional space (x_1, y_1) , (x_2, y_2) , (x_3, y_3) , where $x_1 < x_2 < x_3$ and $y_1 < y_2 < y_3$ (since we approximate a monotonic increasing function where the x axis is the key and the y axis is the location, these 3 points serve as a generalization). Note that values in the y axis are integers as they represent locations, and since the function is monotonic increasing any y value can appear only once. By definition, the

Algorithm 1 Optimal Segmentation

```
1 for every k \in \text{keys} (in increasing order)

2 do find the set S of feasible segments ending at k

3 find segment [j, k] \in S for which T[j-1] is minimal

4 T[k] \leftarrow 1 + T[j-1]
```

linear function starts at the first point in a segment, and ends at the last point in the segment. The linear segment is not feasible if the distance on the y axis (loc) is larger than the defined error, demonstrated in Figure 4.

Therefore, given the 3 points, a linear segment starting at (x_1, y_1) and ending at x_3, y_3 is not feasible if:

$$y_2 - err > \frac{y_3 - y_1}{x_3 - x_1} (x_2 - x_1) + y_1$$
 (1)

By rearranging the inequality we get:

$$err < y_2 - y_1 - \frac{y_3 - y_1}{x_3 - x_1} (x_2 - x_1)$$
 (2)

$$= (y_3 - y_1) \cdot \left(1 - \frac{x_2 - x_1}{x_3 - x_1}\right) - (y_3 - y_2) \tag{3}$$

$$\leq (y_3 - y_1) \cdot \left(1 - \frac{x_2 - x_1}{x_3 - x_1}\right) - 1 \tag{4}$$

In (3) y_3 was added and subtracted, and in (4) we use the fact that y_2 and y_3 are integers (thus $y_3 - y_2 \ge 1$). This provides a lower bound for the distance on the y axis between the first point in a segment and the first point in the following segment:

$$y_3 - y_1 > \frac{err + 1}{1 - \frac{x_2 - x_1}{x_3 - x_1}} = (err + 1) \cdot \frac{x_3 - x_1}{x_3 - x_2} > err + 1$$
 (5)

$$\Rightarrow y_3 - y_1 > err + 1 \tag{6}$$

Since (x_3, y_3) is the first point outside of the segment, the total number of locations in the segment is $y_3 - 1 - y_1 \ge err + 1$.

3.2 Optimal Segmentation

Since the number of segments (the number of linear piece-wise functions) is equal to the number of leafs in the index tree (i.e., each segment is stored in a variable-sized page), our optimization goal is to approximate a dataset using the smallest number of piecewise linear functions, such that the maximal error *error* for each function is satisfied.

Given our design choices, we present Algorithm 1, a dynamic programming algorithm to optimally divide a dataset to non overlapping segments. The algorithm fills a table T[i], where each entry in the table holds the minimal number of segments required for segmenting up to the ith key (keys are sorted in increasing order).

This optimal algorithm has a runtime $O(n^3)$: finding the set of feasible segments ending at some key k takes $O(n^2)$ since the verification of every segment takes O(n) time. Additionally, the memory cost for this algorithm is O(n).

A run time of $O(n^3)$ is not practical for any moderate to large dataset. The runtime can be reduced to $O(n^2)$, however such a reduction comes at the cost of requiring $O(n^2)$ memory, which is still not practical.

Algorithm 2 ShrinkingCone Segmentation

```
sl_{high} \leftarrow \infty
 2
     sl_{low} \leftarrow 0
     the first key is the segment origin
     for every k \in \text{keys} (in increasing order)
 5
            do if k is in the cone:
                    then update sl_{high}
 6
                           update sl_{low}
 7
 8
                    else key k is the origin of a new segment
 9
                           sl_{high} \leftarrow \infty
                           sl_{low} \leftarrow 0
10
```

3.3 Online Segmentation Algorithm

Since using the optimal segmentation algorithm is not practical, we need an algorithm with linear runtime (i.e., O(n)). Therefore, we present a greedy algorithm ShrinkingCone (Algorithm 2), which, given a starting point (key) of a segment, attempts to maximize the length of a segment while satisfying a given error threshold. The main principle of ShrinkingCone is that a new key can be added to a segment if and only if it does not violate the error constraint of any previous key in the segment, assuming that it is the last key in the segment.

We define a cone using the triple: origin point (the key and its location), high slope, and low slope. The combination of the starting point and the low slope gives the lower bound of the cone, and the combination of the starting point and the high slope gives the upper bound of the cone. The origin point is the starting point of a segment.

The cone represents the family of feasible linear functions for a segment starting at the origin of the cone (the high and low slopes represent the range of valid slopes). If a new key to be added to the segment is not inside the cone, there must exist at least one previous key in the segment for which the error constraint will not be satisfied given a segment starting at the origin and ending at the new key. Therefore, a new key that is not inside the cone cannot be included in the segment, and will instead become the origin point of the new segment.

When a new key is added to the segment the cone either narrows (the high slope decreases and/or the low slope increases), or stays the same. The cone cannot widen because that would violate a constraint of a previous key. Additionally, the angle is asymptotically decreasing to 0, since any new key is further away from the origin than any previous key (and as a result the distance from the origin key may become much larger than *error*).

Figure 5 provides an illustration of how the cone is updated: point 1 is the origin of the cone. Point 2 updates both the high slope and the low slope, so the difference between the point and the cone bounds is at most *error* (the arrows up and down are of size *error*). Point 3 is inside the cone, however it only updates the upper bound of the cone (point 3 is less than *error* above the lower bound). Point 4 is outside of the updated cone, and therefore will be the first point of a new segment.

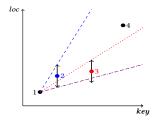


Figure 5: Shrinking cone illustration. Point 1 is the origin of the cone. Point 2 is then added, resulting in the dashed cone. Point 3 is added next, yielding in the dotted cone. Point 4 is outside the dotted cone and therefore starts a new segment.

3.4 Greedy Algorithm Analysis

While the ShrinkingCone algorithm has a run-time of O(n) and uses a constant amount memory (to keep track of the cone), the number of segments that it produces can be arbitrarily worse than the optimal algorithm (proved in Appendix A.3). However, in the following, we show that the algorithm yields results that are close to the optimal algorithm using a variety of real-world datasets.

While ShrinkingCone is not competitive, it does provide a guarantee about the number of segments: the maximal number of segments is at most min $\left(\frac{|keys|}{2}, \frac{|D|}{error+1}\right)$, where |D| is the total number of elements in the dataset (with duplicates). This guarantee stems from Theorem 3.1: no input with less than 3 keys spanning at least error+2 memory locations will cause ShrinkingCone to create a new segment.

In order to evaluate ShrinkingCone on real world datasets, we implemented the optimal algorithm with a runtime of $O(n^2)$ and memory consumption of $O(n^2)$, using a sparse matrix to store whether two keys can be endpoints of a valid segment. We limited ourselves to samples containing 10^6 elements, since even samples of this size can require more than a TB of memory. We performed experiments using contiguous samples from datasets with different distributions and show the number of segments generated by the optimal algorithm and by ShrinkingCone in Table 1.

First, we use three different attributes from the NYC Taxi Dataset [24] (pickup time, drop longitude and drop latitude), each of which has a different underlying distribution. We also used a sample of longitude data from OpenStreetMap data [25], timestamps of requests to a webserver (Weblogs), and timestamps of recorded events (e.g., door opening, motion detected) from IoT devices (IoT). The results show that the number of segments that ShrinkingCone yields is comparable to the number of segments in the optimal case for several error thresholds. For error thresholds for which we did not include results, the matrix used to store whether two keys can be endpoints of a valid segment exceeded the amount of available memory on a server with 768GB of RAM.

4 INDEX LOOKUPS

The most basic operation that an index must support is to be able to lookup the values for either a single key or a range of keys. However, since each entry in the leaf level of A-Tree points to a segment rather than an individual value for a given key, performing a lookup requires first locating the segment that a key belongs to and then performing local search inside the segment in order to

Dataset	error	ShrinkingCone	Optimal	Ratio
Taxi drop lat	10	5358	4996	1.07
Taxi drop lat	100	351	271	1.29
Taxi drop lat	1000	51	48	1.06
Taxi drop lon	10	1198	1138	1.05
Taxi drop lon	100	371	325	1.14
Taxi drop lon	1000	40	37	1.08
Taxi pick time	10	6238	4359	1.43
Taxi pick time	100	165	137	1.2
OSM lon	10	7727	6027	1.28
OSM lon	100	101	63	1.6
Weblogs	10	16961	14179	1.2
Weblogs	100	909	642	1.42
IoT	10	8605	6945	1.24
IoT	100	723	572	1.26

Table 1: ShrinkingCone compared to optimal

find the given element. In the following, we first describe in detail how A-Tree performs lookup operations a single key and then show how we can extend this technique to range predicates.

4.1 Point Queries

The process for searching an A-Tree for a single element involves two steps: (1) searching the tree to find the segment that the element belongs to, and (2) finding the element within a segment. These steps are outlined in Algorithm 3.

4.1.1 Tree Search. The first step in performing a lookup for a given key is to find the segment that the key belongs to. Since, as previously described, each segment is stored in a standard B+ tree (with its starting position as the key and the segment's slope and a pointer to the table page as its value), we must first search the tree in order to find the segment that a given key belongs to.

To find the segment that a given element belongs to, we begin traversing the B+ tree from the root to the leaf, using the standard tree traversal algorithms. More specifically, starting from the root of the tree, we recursively follow the pointer to the child node whose range includes the key to be located. These steps, outlined in the SEARCHTREE function of Algorithm 3 terminate when reaching a leaf node where the corresponding entry in the leaf node which points to a the table page that contains the key.

Since the B+ tree is used to index the segments rather than individual points, the runtime for searching for the segment that a given key belongs in is $O(log_b(p))$, where b is the fanout of the tree (number of separators/pointers inside an inner node) and p is the number of segments created during the segmentation process.

4.1.2 Segment Search. Once the segment for a key has been located, A-Tree must find element's position inside the segment. Recall that segments are created such that an element is no more than a constant distance (error) from the element's position determined through linear interpolation. Other techniques for interpolation search inside a fixed-sized index page are discussed in [12].

To compute the approximate location of a key *key* within a given segment *s*, we subtract the key from the first key that appears in the segment *s.start*. Then, we multiply the difference by the segment's slope *s.slope*, as shown below in the following equation.

$$pred_pos = (key - s.start) \times s.slope$$
 (1)

Algorithm 3 Lookup Algorithm

```
LOOKUP(tree, key)
1 seg \leftarrow SearchTree(tree.root, key)
  val \leftarrow SearchSegment(seg, key)
   return val
SEARCHTREE(node, key)
  i \leftarrow 0
   while key < node.keys[i]
2
3
         do i \leftarrow i + 1
   if node.value[i].isLeaf()
       then j \leftarrow 0
6
            while key < node.values[j]
                 do i \leftarrow i + 1
7
8
            return node.values[j]
9
   return SearchTree(node.values[i], key)
SEARCHSEGMENT(seg, key)
   pos \leftarrow (key - seg.start) \times seg.slope
  return BinarySearch(seg. data, pos – error, pos + error, key)
```

After interpolating the element's position, the true position of the element is guaranteed to be within the error threshold. Therefore, the next step is to locally search the following region.

$$true_pos \in [pred_pos - error, pred_pos + error]$$
 (2)

To search this region, it is possible to utilize any well-known search algorithm, including linear search, binary search, or exponential search. For simplicity, A-Tree uses binary search to find the element inside the specified region, as shown in the Search-Segment function of Algorithm 3. However, for very small error thresholds, linear search can be faster. By guaranteeing that segments satisfy the specified error condition, the cost of searching for an element inside a segment is bounded. More specifically, the runtime for searching for an element inside a segment is O(error) where error is constant.

4.2 Range Queries

In addition to point queries, range queries are common in many database workloads. Range queries, unlike point queries, have the additional requirement that they need to examine every data item in the specified range. Therefore, for range queries, the selectivity of the query (i.e., the number of tuples that satisfy the query's predicate) has a large influence on the total runtime of the query.

However, like point queries, range queries must also find a single tuple: either the start or the end of the range. Therefore, A-Tree uses the previously described point query lookup techniques in order to find the beginning of the specified range. Then, since segments either store keys contiguously (clustered index) or have an indirection layer with pointers that is sorted by the key (nonclustered index), A-Tree can simply scan from the starting location until a key is found that lies outside of the specified range. For a clustered index, scanning the relevant range performs only sequential access, while for a non-clustered index, range queries requires random memory accesses (which is true for any non-clustered index).

5 INDEX INSERTS

Along with determining up the location of a given key, an index needs to be able to handle insert operations. In a typical B+ tree that uses paging, pages are left partially filled and new values are inserted into an empty slot. When a given page is full, the node is split into two nodes, and the changes are propagated up the tree (i.e., the inner nodes in the tree are updated to reflect the change).

Although similar, inserting into an A-Tree requires additional consideration due to the precise error guarantees that we provide. More specifically, a new item cannot necessarily be inserted into a new segment, since inserting it may violate the specified error condition. Moreover, as discussed before, segments are stored in a variable-sized pages that can be large in size and thus make inserts more expensive (sort order within a segment must be kept).

Therefore, in addition to the data items for a given segment, each segment in an A-Tree contains an additional fixed-sized buffer. As shown in Algorithm 4, new keys to be inserted into an A-Tree are added to the buffer portion of the segment for which the key belongs to (line 2). This buffer is kept sorted in order to enable efficient search and merge operations. Once the buffer reaches its predetermined size (buffer_size), the buffer is combined with the data in the segment (creating one sorted array of data) (line 5).

Then, to ensure that the error condition is always satisfied when merging the buffer with the segment, we use the previously described segmentation algorithm (Algorithm 2) in order to create a series of valid segments that satisfy the error threshold (line 6). Note that depending on the data, the number of segments generated after this process can be one (i.e., the data inserted into the buffer does not violate the error threshold) or several. Finally, for each of the new segments that are generated after the segmentation process, each is inserted into the tree (line 7-8) and the old page is removed (line 9).

Importantly, however, storing additional data inside a segment impacts how to locate a given item, as well as how the error is defined. More specifically, since searching for an item requires searching the segment's data (as described in the SearchSegment function in Algorithm 3 of Section 4) as well as searching for the item in the buffer, we need to make sure that the error threshold that the user specifies is satisfied.

Since adding a buffer for each segment can violate the error guarantees that A-Tree provides, we transparently incorporate the buffer's size into the error parameter for the segmentation process. More formally, given a specified error of *error*, we transparently set the error threshold for the segmentation process to (*error* – *buffer_size*). This ensures that, for any specified error, a lookup operation will satisfy the specified error even if the element is located in the buffer.

The overall runtime for inserting a new element into an A-Tree is the time required to locate the segment that the element belongs to, as well as add the element to the sorted array. With p pages stored in an A-Tree, inserting a new key has the following runtime.

$$insert\ runtime: O(log_b p) + O(buffer_size)$$
 (1)

Note that when the buffer is full and the segment needs to be merged and segmented, the runtime has an additional cost of O(d), where d is the sum of a segment's data and buffer size.

Algorithm 4 Insert Algorithm

```
InsertKey(tree, seg, key)
1 seg \leftarrow SearchTree(tree.root, key)
    seg.buffer.insert(key)
 3
    if seg.buffer.isFull()
 4
        then
 5
              d = merge(seg.data, seg.buffer)
              segs = \texttt{segmentation}(d)
 6
 7
              for s \in segs
 8
                   do tree.insert(s)
 9
              tree.remove(seg)
10
    return
```

Note that if the write-rate is very high, we could also support merging algorithms that use a second buffer similar to how column stores merge a write-optimized delta to the main compressed column. However, this is an orthogonal consideration that heavily depends on the ratio of reads and writes in a workload and is outside the scope of this paper.

6 COST MODEL

Since the specified error threshold influences both the performance of lookup and insert operations as well as the index's size, the natural question follows: how should a DBA pick the error threshold for a given workload? Therefore, as part of this work, we provide a cost model that helps a DBA to pick a "good" error threshold when creating a new A-Tree. At a high level, there are two main objectives that a DBA can optimize: performance (i.e., lookup latency) and space consumption. Therefore, we present two ways to apply our cost models that help a DBA choose the corresponding error threshold given either a performance (lookup/insert latency) or a space requirement.

6.1 Latency Guarantee

For a given workload, it is valuable to be able to provide latency guarantees to an application. For example, an application may require that lookups to a database take no more than a specified time threshold (e.g., 1000ns) due to SLAs in cloud scenarios or application-specific requirements (e.g., for interactive applications). Since A-Tree incorporates an error term that in turn affects performance, we can model the index's latency in order to pick an error threshold that satisfies the specified latency requirement.

As previously mentioned, lookups require first finding the relevant segment and then searching the segment (data and buffer) for the element. Since the error threshold influences the number of segments that are created (i.e., a smaller error threshold yields more segments), we can use a function that takes an error value and returns the number of segments that are needed to satisfy the specified error. This function can either be learned for a specific dataset (i.e., segment the data using different error threshold and record the number of segments created) or a general function can be used (e.g., make the simplifying assumption that the number of segments decreases linearly as the error increases). We use S_e to represent the total number of resulting segments for given dataset using an error threshold of e.

Therefore, the total estimated latency for a lookup query for an error threshold of e can be modeled by the following expression, where b is the tree's fanout, buff is a segment's maximum buffer size, and c is a constant representing the latency (in ns) of a cache miss on the given hardware (e.g., 100ns). Moreover, the cost function assumes binary search for the area that needs to be searched within a segment bounded by e as well as searching the complete buffer.

LATENCY(e) =
$$c \underbrace{\left[log_b(S_e) + log_2(e) + log_2(buff)\right]}_{\text{Tree Search}}$$
 Segment Search Buffer Search (1)

Setting c to a constant value implies that all random memory accesses have a constant penalty. However, depending on the size of the data, caching can often change the penalty for a random access since, for smaller dataset sizes, a higher percent of the data can fit in the cache. In theory, instead of being a constant, c could be a function that returns the penalty of a random access for a given dataset size. However, for our estimates, we make the simplifying that c is a constant value.

Using this cost estimate, the index with the smallest storage footprint that will satisfy the given latency requirement L_{req} (given in nanoseconds) is given by the following expression, where E represents a set of all possible error values (e.g., $E = \{10, 100, 1000\}$) and INDEX_SIZE is a function that returns the estimated size of an index (defined in the next section).

$$e = \underset{\{e \in E \mid \text{LATENCY}(e) \le L_{req}\}}{\arg \min} \left(\text{INDEX_SIZE}(e) \right)$$
 (2)

In addition to modeling the latency for lookup operations, we can similarly model the latency for insert operations. However, there are a few important differences. First, inserts do not have to probe the segment. Also, instead of searching a segment's buffer, inserts require adding the item to the buffer in sorted order. Finally, we must also consider the cost associated with splitting a full segment.

6.2 Space Budget

Instead of specifying a bound on the latency for looking up a key/inserting a key in the index, a DBA can give A-Tree a storage budget to use. In this case, the goal becomes to provide the highest performance (i.e., lowest latency for lookups and inserts) while not exceeding the specified storage budget.

More formally, we can estimate the size of the a clustered index (in bytes) for a given error threshold of e using the following function, where again S_e is the number of segments that are created for an error threshold of e, f is the fill ratio of the tree (e.g., 0.5), and b is the fanout of the tree.

$$SIZE(e) = \underbrace{(f \cdot S_e \cdot log_b(S_e) \cdot 16B)}_{Tree} + \underbrace{(S_e \cdot 24B)}_{Segment}$$
 (1)

The first term is a pessimistic bound on the storage cost of the tree (leaf + internal nodes using 8 byte keys/pointers), while the second term represents the added metadata about each segment (i.e., each segment has a starting key, slope, and pointer to the underlying data, each 8 bytes).

Therefore, similar to the cost model presented for lookups, the smallest error threshold that satisfies a given storage budget S_{req}

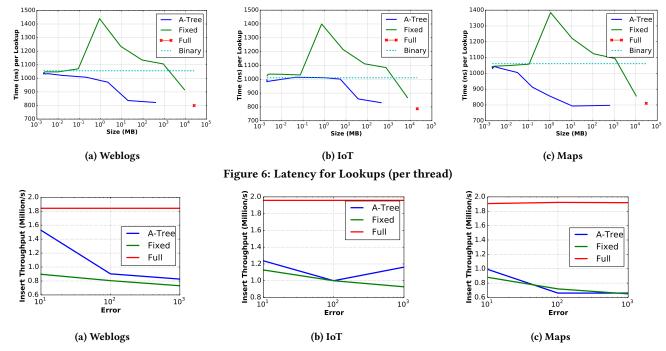


Figure 7: Throughput for Inserts (per thread)

(given in bytes) is given by the following expression where again E represents a set of all possible error values (e.g., $E = \{10, 100, 1000\}$).

$$e = \underset{\{e \in E \mid \text{SIZE}(e) \leq S_{req}\}}{\arg \min} \left(\text{LATENCY}(e) \right)$$
(2)

where $\text{constants} = \text{constants} = \text{$

As we show in our experiments, our cost model is accurately able to estimate the size of an A-Tree over real-world datasets, providing DBAs with a valuable way to balance performance (i.e., latency) with the storage footprint of an A-Tree.

7 EVALUATION

This section evaluates A-Tree and the techniques presented in this paper. First, in Section 7.1 we compare the overall performance of A-Tree, measuring its lookup and insert performance for a variety of real-world datasets. Overall, we see that A-TREE achieves comparable performance to both a full index as well as fixed-sized indexes while using orders of magnitude less space. Then, in Section 7.2 and in Section 7.3, we present several microbenchmarks, including results that show how A-TREE performs for an adversarial synthetically generated dataset (i.e., a dataset with a worst-case distribution of values) as well as the scalability of our index over different dataset sizes. Finally, in Section 7.4 we show that our cost model is able to accurately model the behavior of an A-Tree and can pick an error threshold that satisfies a specified lookup latency or space consumption requirement. In the Appendix, we include a breakdown of the time spent searching for an element as well as the effect of a segment's fill factor on performance.

All experiments were conducted on a single server with an Intel E5-2660 CPU (2.2GHz, 10 cores, 25MB L3 cache) and 256GB RAM. Furthermore, all index and table data was held in memory for all experiments.

7.1 Exp. 1: Overall Performance

In the following, we evaluate the overall lookup and insert performance of A-Tree. For these comparisons, we benchmark A-Tree against both a full index (i.e, a dense index) as well as an index that uses fixed size pages and indexes only the first key per page (i.e., a sparse index). More specifically, since at a high level A-Tree uses paging to implement segments, we compare our approximate and variable sized paging techniques to commonly used fixed-sized pages. However, typically a full index can be seen as best case baseline for the lookup performance and thus gives us an interesting reference point.

For the two baselines (full and fixed-sized paging), we use a popular B+ tree implementation (STX-tree [34] version 0.9). Our A-Tree prototype also uses the STX-tree to index the segments that we generate, with lookup and insert operations implemented as previously described. Note that, as previously mentioned, any other tree implementation could also be used to serve as the organization layer under A-Tree. For example, a tree that internally performs compression for each tree node can be used to index our generated segments and will still provide benefits (i.e., compressing the data that A-Tree stores in its underlying tree can improve performance), but this benefit is orthogonal to the benefits that A-Tree achieves through its approximate and variables sized segmentation of the underlying data. Therefore, in order to ensure a fair comparisons, it is important that we keep the underlying tree implementation the same for all baselines as well as for our A-Tree prototype.

7.1.1 Datasets. Since performance of our index depends on the distributions of elements in a given dataset, we evaluate A-Tree on real-world datasets with different distributions. For our evaluation, we have chosen the following three different real-world

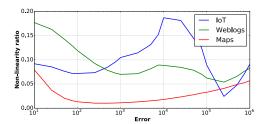


Figure 8: Non-linearity of tested datasets

datasets, each with very different underlying data distributions: (1) Weblogs [35], (2) IoT [17], and (3) Maps [25].

The Weblogs dataset contains $\approx 715M$ log entries for every request to the computer science department at a university over the past 14 years. This dataset contains subtle trends, such as the fact that more requests occur during certain times (e.g., school year vs summer, daytime vs night time). On the other hand, the IoT dataset contains $\approx 5M$ readings from approximately 100 different IoT sensors (e.g., door, motion, power) installed throughout an academic building at a university. Since these sensors generally reflect human activity, this dataset has interesting patterns, such as the fact there is more activity during certain hours because classes are in session. Finally, the Maps dataset contains the longitude of $\approx 2B$ user-maintained features (e.g., roads, museums, coffee shops) across the world. Unsurprisingly, the longitude of locations is relatively linear and does not contain many periodic trends.

For our approach, the most important aspect of a dataset that impacts A-Tree's performance is the periodicity. For now, think of the periodicity as the distance between two "humps" in a stepwise function that maps keys to storage locations as shown as in Figure 9a (blue line). If the error of our index is larger than the periodicity (green line), the segmentation results in a single segment. However, if the error is smaller than the periodicity (red line), we need multiple segments to approximate the data distribution.

Therefore, we introduce a new measure that can reflect the periodicity of a dataset. The measure is defined as follows: First, we compute the number of segments required to cover the dataset for a given error threshold. We then normalize this result by the number of segments required for a dataset of the same size with periodicity equal to the error (which is the worst case, or the most "non-linear" in that scale). We refer to this metric as the *non-linearity ratio*.

To show that all these datasets contain a distinct periodicity pattern, we report the non-linearity ratio of each of dataset as shown in Figure 8. The IoT dataset has a very significant bump, signifying that there is very strong periodicity the scale of 10⁴, likely due to patterns that following human behavior (e.g., day/night hours). Weblogs, on the other hand, has multiple bumps which are likely correlated to different periodic patterns (e.g., more requests during the school year). The Maps dataset, unlike the other two, is very linear at small scales (but has stronger periodicity at scales larger than presented).

7.1.2 Lookups. The first series of benchmarks that we present show how A-Tree compares to (1) a full index, (2) an index that uses fixed-sized paging, and (3) binary search on one large segment. We included binary search since it represents the most extreme case where the error is equal to the data size (i.e., our segmentation algorithm would return one large segment). Moreover, for the

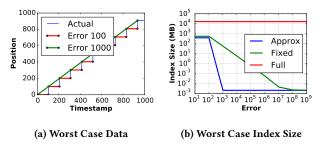


Figure 9: Worst Case Analysis

Weblog and IoT dataset, we created a clustered index using the timestamp attribute which is the primary key of these datasets. For the Maps dataset, we created a non-clustered index over the longitude attribute, which is not unique.

The results are shown in Figure 6 and plot the lookup latency for various sizes of the index for the Weblogs (scaled to 1.5B records), IoT (scaled to 1.5B records), and Maps (not scaled, 2B records) datasets. More specifically, since the size of both A-Tree and the fixed-size paging baseline can be varied (i.e., A-Tree's error term and the page size influence the number of indexed keys), we show how the performance of each of these approaches scales with the size of the index. Note that the size of a full index cannot be varied and is therefore a single point in the plot. Additionally, since typical binary search does not have any additional storage requirement, its size is zero but is visualized as a dotted line.

In general, the results show that our A-Tree always has better performance than an index that uses fixed-sized paging. Most importantly, however, the results show that an A-Tree offers significant space savings compared to both fixed-sized paging as well as a full index. For example, in the Maps dataset, A-Tree is able to match the performance of a full index using only 609MB of memory, while a full index consumes over 30GB of space. Moreover, compared to a fixed sized tree, an A-Tree which consumes only 1MB of memory is able to match the performance of a fixed-sized index which consumes over 10GB of memory, offering a space savings of four orders of magnitude.

Furthermore, as expected for very small index sizes (i.e., very large page sizes in fixed-sized paging and a high error threshold in A-Tree), both A-Tree and fixed-sized paging mimic the performance of binary search. This is because there are only a small number (in some cases one) of pages, and the vast majority of the time is spent performing binary search to locate the item inside the large page. On the other hand, as the index grows, the performance of both fixed-sized paging as well as A-Tree converge to that of a full index due to the fact that pages contain very few elements and almost all of the time is spent traversing the tree to find the relevant page. Note that the spike in the graph for the fixed-sized index is due to the fact that the index begins to fall out of the CPU's L2 cache.

Finally, as expected, the data distribution impacts the performance of A-Tree. More specifically, we can see that A-Tree is able to more quickly match the performance of a full tree with the Maps dataset, compared to the Weblogs and IoT datasets. This is due to the fact that the Maps dataset is relatively linear, when compared to the Weblogs and IoT datasets as shown in Figure 8.

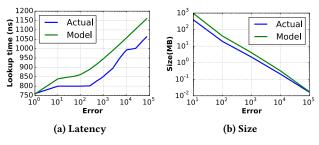


Figure 10: Cost Model Accuracy

7.1.3 Inserts. Next, we compare the insert performance of A-Tree to both a full index, as well as an index that leverages fixed-sized paging as previously described. As an insertion strategy, both indexes use the previously described techniques to insert new keys into the buffer of the segment that they key belongs in. Specifically, in order to ensure a fair comparison and that the A-Tree is not unfairly write-optimized, we set the size of the buffer to half of the specified error (i.e, for an error threshold of 100, the underlying data for each segment is has a maximum error of 50 and each segment's buffer has a size of has room for 50 elements). For the index with fixed-sized pages, the page size is given by the error threshold we used for the A-Tree and half of the page size is used as the buffer size. As usual, once the buffer is full, the page is split into two pages.

The results, shown in Figure 7 compare each of the index's insert throughput for various error thresholds. As shown, A-Tree is able to achieve insert performance that is, in general, comparable to an index that uses fixed-sized paging. Unsurprisingly, a full index is able to handle a higher write load than either A-Tree or an index that uses fixed-sized paging. This is because both A-Tree and an index that uses fixed-sized paging need to periodically split pages that become full which is not required for the full index.

Unlike a full B+ tree, for A-Tree and an index that uses fixedsized paging, the pages that the index points to at the leaf level need to be split once the page becomes full. This results in additional overhead that a full B+ tree does not have. Even though both A-Tree and fixed-sized paging need to split pages, A-Tree needs to also execute the segmentation algorithm, explaining the performance discrepancy between A-Tree and fixed-sized paging.

Interestingly, however, A-Tree is faster than fixed-sized paging in some cases. This is due to the fact that the error threshold determines the number of segments in the tree. For a small error, there typically are more segments generated, which reduces the number of times that A-Tree needs to merge the buffer with a node and then re-segment the node.

7.2 Exp. 2: Worst Case Analysis

Since the data distribution influences the performance of A-Tree, we synthetically generated data to illustrate how our index performs with data that represents a worst-case. More specifically, we generate data using a step function with a fixed step size of 100, as shown in Figure 9a. Since the step size is fixed, an error threshold less than the step size results in a single segment per step, representing the worst-case distribution of data. However, given an error threshold larger than the step size, our segmentation algorithm will be able to use a single segment to represent the entire dataset.

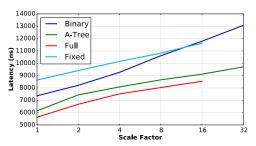


Figure 11: Data Scalability

Figure 9b shows the size of each index built over this worst case dataset. As shown, for error thresholds of less than 100, the size of an A-Tree is the same as a fixed-sized index but still smaller than a full index. This is due to the fact that for the error thresholds less than the step size, A-Tree creates segments of size *error*, resulting in a large number of nodes in the tree. On the other hand, for an error threshold of larger than 100, A-Tree is able to represent the step dataset with only a single segment, dramatically reducing the index's size.

7.3 Exp. 3: Data Size Scalability

To evaluate how A-Tree performs for various dataset sizes, we measure the lookup latency for the Weblogs dataset using various scale factors where both the error threshold and fixed-page size are set to 100, which is optimal for this dataset for our A-Tree index. Since the performance of our index depends on the underlying data distribution, we scale the dataset while maintaining the underlying trends. Furthermore, we omit the result for all other datasets here (IoT and Maps) since they follow similar trends.

The data scalability results, shown in Figure 11 show that the tree-based approaches (i.e., A-Tree, a full index, and an index that uses fixed-sized paging) scale better than binary search due to the better theoretical asymptotic runtime $(log_b(n) \text{ vs } log_2(n))$. Additionally, A-Tree's performance over various dataset sizes closely follows that of a full index which offers the best performance, demonstrating that our techniques offer valuable space savings without sacrificing much of the optimal performance. More importantly, neither a full index nor an index that uses fixed-sized paging was able to scale to a scale factor of 32 since the size of the index exceeded the amount of available memory, which again shows that A-Tree is able to offer valuable space savings and thus allows databases to scale to much larger datasets with comparable performance to full index.

7.4 Exp. 4: Accuracy of Cost Model

Since, as previously described, the error threshold influences both the latency as well as space consumption of our index, our cost model presented in Section 6 aims to guide a DBA when determining what error threshold to use for an A-Tree. More specifically, given a latency requirement (e.g., 1000ns) or a space budget (e.g., 2GB), our cost model automatically determines an appropriate error threshold that satisfies the given constraint.

Figure 10a shows the estimated and actual lookup latency for various error threshold on the Weblogs dataset using a value of 50ns for c (the cost of a random memory access) determined through a memory benchmarking tool on the given hardware. As shown, our model for the estimated latency predicts an upper bound for

the actual latency of a lookup operation allowing a DBA to set the lookup latency as a constraint. This can be attributed to the fact that our cost models does not reflect effects of CPU caches and thus overestimates the lookup costs.

In addition to our cost model that estimates the latency for lookup operations on an A-Tree index for a given error threshold, we also provide DBAs with the ability to specify a storage budget. Therefore, to evaluate our model that predicts the size of an index for a given error threshold, we show the predicted and actual size of an A-Tree for various error thresholds in Figure 10b. As shown, our model is able to accurately predict the size of an index for a given error threshold while ensuring that our estimates are pessimistic (i.e., the estimated cost higher than the true cost).

8 RELATED WORK

The presented techniques in this paper have overlap with work in different areas ranging from (1) index compression, over (2) partial/adaptive indexes, to (3) function approximations.

Index Compression: Since B+ trees can often consume a large amount of storage space, it is not surprising that several index compression techniques have been proposed. These approaches primarily aim to reduce the size of keys in internal nodes by applying techniques such as prefix/suffix truncation, dictionary compression, and key normalization [11, 13, 23]. This allows these techniques to reduce the overall number of nodes used to store an index. Importantly, these techniques can also be applied within A-Tree in order to further reduce the size of the underlying tree structure.

Similar to B+ tree compression, several methods have been proposed in order to more compactly represent bitmap indexes [3, 6, 27, 31, 36]. Many of these techniques are specific to bitmap indexes, which are primarily only useful for attributes with few distinct values and not the general workloads that A-Tree targets.

Moreover, other techniques have been developed to better align the node layout with the requirements of modern hardware by also utilizing compression. For example, CSB+ trees [28] remove the need to store pointers by using offsets. However, the main target is read-heavy workloads. FAST [19] is another more recent tree structure that organizes tree elements in order to store then index in a more compact representation and exploit modern hardware features (e.g., SIMD, cache line size) for read-heavy workloads. Similarly, an Adpative Radix Tree (ART) [20] also tried to leverage CPU caches for in-memory indexing. Another idea that was discussed in [38] are hybrid indexes. The idea is to separate the index in hot and cold data where cold data is stored in a compressed format since it is updated less often. As previously mentioned, all these techniques are orthogonal to A-Tree and we could use these techniques to store the underlying tree structure to efficiently locate the segment that a given key belongs to and optimize for read-heavy workloads or hot/cold data.

The closest work to A-Tree is the BF-Tree [2]. Similar to A-Tree, the BF-Tree uses a B+ tree to store information about a region of the dataset, instead of the indexing individual keys. However, leaf nodes in a BF-Tree are bloom filters instead of linear segments. Additionally, unlike A-Tree, BF-tree does not exploit properties about the data's distribution when segmenting a dataset.

Sparse indexes like Hippo [37], Block Range Indexes [33], and Small Materialized Aggregates (SMAs) [22] all store information about value ranges similar to the idea of segments that we store in A-Tree. However, these techniques do not consider underlying data distribution or bound the latency for lookups/inserts.

Finally, several approximation techniques have been proposed in order to improve the performance of similarity search [9, 14, 21, 29] (for string or multimedia data), unlike A-Tree which uses approximation for compressing indexes optimized for traditional point and range queries.

Partial and Adaptive Indexes: Partial indexes [32] aim to reduce the storage footprint of an index since they built the index only for a subset of data that is of interest to the user. The main idea is that they only take the workload into account to index the "important data". Another example are tail indexes [10] that index only rare data items (i.e., outliers) in aggregate queries to reduce the storage footprint of the overall index. A-Tree, on the other hand, supports queries over all attribute values but could be extended to index only "important" data ranges as well. Furthermore, database cracking [15] is a technique that physically adaptively reorders values in a column store in order to more efficiently support selection queries without needing to store secondary indexes. Since database cracking reorganizes values based on past queries, it does not efficiently support for ad-hoc queries, like A-Tree can.

Function Approximation: The main idea of an A-Tree index is that it approximates the data distribution using a piece-wise linear function. The general problem of approximating curves using piece-wise functions is not new [5, 8]. The error metrics E_2 (integral square error) and E_∞ (maximal error) for these approximations have been discussed as well as different segmentation algorithms [26]. However, unlike the prior work, we consider only monotonic increasing functions and E_∞ . Moreover, none of these techniques have been used to build approximate indexes so far.

More recent work [7, 18, 30] for time series datasets also leverages piece-wise linear approximations to store patterns in a dataset and for similarity search, instead of for indexing. While these works also trade-off the number of segments with the accuracy of the approximated representation, they do not aim to provide the lookup and space consumption guarantees that A-Tree does.

Finally, other work [1] leverages piece-wise linear functions to compress inverted lists by storing functions and the distances of elements from the extrapolated functions. However, these approximations use linear regression (which minimizes E_2 per segment), and there are no bounds on the error (neither E_2 nor E_∞).

9 CONCLUSION

In this paper, we present A-Tree a new approximate index structure that incorporates tunable error parameter to allow a DBA to balance lookup performance and space consumption of an index. To navigate this tradeoff, we presented a cost model that determines an appropriate error parameter given either (1) a lookup latency requirement (e.g., 500ns) or (2) a storage budget (e.g., 100MB). We evaluated A-Tree using several real-world datasets and showed that we can provide performance that is comparable to a full index structure while reducing the storage footprint by orders of magnitude.

REFERENCES

- [1] N. Ao, F. Zhang, D. Wu, D. S. Stones, G. Wang, X. Liu, J. Liu, and S. Lin. Efficient Parallel Lists Intersection and Index Compression Algorithms Using Graphics Processing Units. VLDB, pages 470-481, 2011.
- [2] M. Athanassoulis and A. Ailamaki. BF-tree: Approximate Tree Indexing. In VLDB, pages 1881-1892, 2014.
- [3] M. Athanassoulis, Z. Yan, and S. Idreos. UpBit: Scalable In-Memory Updatable Bitmap Indexing. In SIGMOD, pages 1319-1332, 2016.
- [4] R. Bayer and K. Unterauer. Prefix B-trees. ACM Trans. Database Syst., pages
- [5] D. Braess. Chebyshev Approximation by Spline Functions with Free Knots. Numerische Mathematik, pages 357-366, 1971.
- [6] C.-Y. Chan and Y. E. Ioannidis. Bitmap Index Design and Evaluation. In SIGMOD, pages 355-366, 1998.
- [7] T. chung Fu. A Review on Time Series Data Mining. Engineering Applications of Artificial Intelligence, pages 164 - 181, 2011.
- R. Esch and W. Eastman. Computational Methods for Best Spline Function Approximation. Journal of Approximation Theory, pages 85 - 96, 1969.
- [9] A. Esuli. Use of Permutation Prefixes for Efficient and Scalable Approximate Similarity Search. *Inf. Process. Manage.*, pages 889–902, 2012.
 [10] A. Galakatos, A. Crotty, E. Zgraggen, C. Binnig, and T. Kraska. Revisiting Reuse
- for Approximate Query Processing. In VLDB, pages 1142–1153, 2017.

 [11] J. Goldstein, R. Ramakrishnan, and U. Shaft. Compressing Relations and Indexes.
- In ICDE, pages 370-379, 1998.
- [12] G. Graefe, B-tree Indexes, Interpolation Search, and Skew, In DaMon, 2006.
- [13] G. Graefe and P. Larson. B-Tree Indexes and CPU Caches. In ICDE, pages 349-358,
- [14] M. E. Houle and J. Sakuma. Fast Approximate Similarity Search in Extremely High-Dimensional Data Sets. In *ICDE*, pages 619–630, 2005.
- [15] S. Idreos, M. L. Kersten, and S. Manegold. Database Cracking. In CIDR, pages 68-78, 2007.
- [16] S. Idreos, S. Manegold, H. Kuno, and G. Graefe. Merging What's Cracked, Cracking What's Merged: Adaptive Indexing in Main-memory Column-stores. VLDB, pages 586-597, 2011.
- [17] Omitted due to double blind requirement.
- [18] E. Keogh, S. Chu, D. Hart, and M. Pazzani. An Online Algorithm for Segmenting Time Series. In ICDM, pages 289–296. IEEE, 2001.
- [19] C. Kim, J. Chhugani, N. Satish, E. Sedlar, A. D. Nguyen, T. Kaldewey, V. W. Lee, S. A. Brandt, and P. Dubey. FAST: Fast Architecture Sensitive Tree Search on Modern CPUs and GPUs. In SIGMOD, pages 339-350, 2010.
- V. Leis, A. Kemper, and T. Neumann. The Adaptive Radix Tree: ARTful Indexing for Main-memory Databases. In ICDE, pages 38-49, 2013.
- [21] C. Li, E. Chang, H. Garcia-Molina, and G. Wiederhold. Clustering for Approximate Similarity Search in High-Dimensional Spaces. IEEE Trans. on Knowl. and Data Eng., pages 792-808, 2002.
- [22] G. Moerkotte. Small Materialized Aggregates: A Light Weight Index Structure for Data Warehousing. In VLDB, pages 476-487, 1998
- [23] T. Neumann and G. Weikum. RDF-3X: A RISC-style Engine for RDF. Proc. VLDB Endow., pages 647-659, 2008.
- [24] NYC Taxi & Limousine Commission Trip Record Data. http://www.nyc.gov/ html/tlc/html/about/trip_record_data.shtml.
- [25] OpenStreetMap database @OpenStreetMap contributors. https://aws.amazon. com/public-datasets/osm.
- [26] T. Pavlidis and S. L. Horowitz. Segmentation of Plane Curves. IEEE Trans. Comput., pages 860-870, 1974.
- [27] A. Pinar, T. Tao, and H. Ferhatosmanoglu. Compressing Bitmap Indices by Data Reorganization. In ICDE, pages 310-321, 2005.
- [28] J. Rao and K. A. Ross. Making b+-trees cache conscious in main memory. In SIGMOD, pages 475-486, 2000
- [29] K. V. Ravi Kanth, D. Agrawal, and A. Singh. Dimensionality Reduction for Similarity Searching in Dynamic Databases. In SIGMOD, pages 166-176, 1998.
- [30] H. Shatkay and S. B. Zdonik. Approximate Queries and Representations for
- Large Data Sequences. In *ICDE*, pages 536–545. IEEE, 1996.
 [31] M. Stabno and R. Wrembel. RLH: Bitmap Compression Technique Based on Run-length and Huffman Encoding. Inf. Syst., pages 400-414, 2009.
- M. Stonebraker. The Case for Partial Indexes. SIGMOD Record, pages 4-11, 1989.
- M. Stonebraker and L. A. Rowe. The Design of POSTGRES. In SIGMOD, pages 340-355, 1986,
- STX B+ Tree. https://panthema.net/2007/stx-btree/.
- [35] Omitted due to double blind requirement.
- K. Wu, E. J. Otoo, and A. Shoshani. Optimizing Bitmap Indices with Efficient Compression. ACM Trans. Database Syst., pages 1-38, 2006.
- [37] J. Yu and M. Sarwat. Two Birds, One Stone: A Fast, Yet Lightweight, Indexing Scheme for Modern Database Systems. In VLDB, pages 385–396, 2016.
- H. Zhang, D. G. Andersen, A. Pavlo, M. Kaminsky, L. Ma, and R. Shen. Reducing the Storage Overhead of Main-Memory OLTP Databases with Hybrid Indexes

In SIGMOD, pages 1567–1581, 2016. M. Zukowski, S. Heman, N. Nes, and P. Boncz. Super-Scalar RAM-CPU Cache Compression. In ICDE, pages 59-, 2006.

A APPENDIX

A.1 Lookup Breakdown

A lookup operation, as described in Section 4 involves two steps (i.e., locating the segment where a key belong to and the searching the segment's data in order to find the item within the segment). Therefore, we examine the amount of time spent in each of these two steps for A-Tree as well as an index that uses fixed-sized paging for various error thresholds.

The results, presented in Figure 13, show that in both cases the majority of time is spent searching the tree to find the page where the data item belongs for smaller error thresholds (and page sizes),. Since A-Tree is able to leverage properties of the underlying data distribution in order to create variable sized segments, the resulting tree is significantly smaller. Therefore, A-Tree spends less time searching the tree to find the corresponding segment for a given key.

Varying Fill Factor

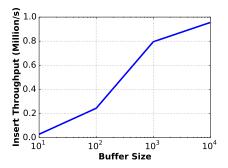
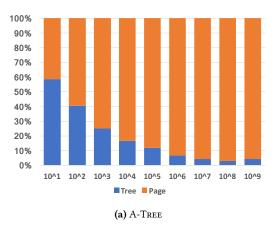


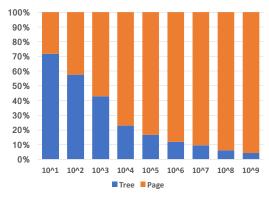
Figure 12: Insert Throughput / Varying Buffer Size

As previously mentioned, the buffer size of a segment determines the amount of space that a segment reserves to hold new data items. Once the segment's insert buffer reaches this threshold, the data from the segment and the segment's buffer is merged, and A-Tree executes the previously described segmentation algorithm to generate new segments that satisfy the specified error threshold.

Therefore, in Figure 12, we vary the buffer size and measure the total throughput that A-TREE can achieve using the Weblogs dataset with an error threshold of e = 20,000. As shown, the size of the buffer can dramatically impact the write throughput of an A-Tree. More specifically, larger buffers result in fewer splitting operations, improving performance. However, a buffer that is too large will result in longer lookup latencies (modeled in the cost model in Section 6).

Therefore, the fill factor of an A-TREE can be effectively used by a DBA to tune an A-Tree to be more read or write optimized, depending on the workload.





(b) Fixed-sized Index

Figure 13: Lookup Breakdown

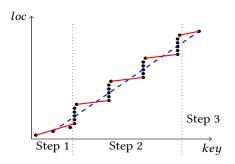


Figure 14: Competitive analysis sketch: the dots are the input, the dashed line is the optimal segmentation (the first dot is a segment), and the solid lines are the segments created by ShrinkingCone

A.3 SHRINKINGCONE Competitive Analysis

We prove that ShrinkingCone can be arbitrarily worse than the optimal solution (not competitive)

PROOF. Given the error threshold E = 100, consider the following input to ShrinkingCone:

- (1) 3 keys (x_1, y_1) , (x_2, y_2) , (x_3, y_3) where $y_1 = 1$, $y_2 = 2$, $y_3 = 3$ and $x_3 x_2 = x_2 x_1 = \frac{E}{2}$ (this is step 1 in Figure 14).
- (2) The key $x_4 = x_3 + \frac{1}{E}$ repeated E + 1 times (using E + 1 consecutive locations), and the key $x_5 = x_4 + \frac{1}{E}$ without repetitions (using 1 location).

After that repeat for $i \in [1, N]$ the following pattern: the key $x_{2(i+2)} = x_{2(i+2)-1} + E$ repeated E + 1 times, and a single appearance of the key $x_{2(i+2)+1} = x_{2(i+2)} + \frac{1}{E}$ (this is step 2 in Figure 14).

- (3) The key $x_{2(N+1+2)} = x_{2(N+1+2)-1} + \frac{E}{2}$ (step 3 in Figure 14). The algorithm will then create the following segments (an illustration is shown in Figure 14):
 - $[x_1, x_4]$ (with slope $\frac{3}{E+\frac{1}{E}}$): adding the key x_5 will result in the slope $\frac{3+E+1}{E+\frac{2}{E}}$ which will not satisfy the error requirement for x_4 , $y_1 + \frac{3+E+1}{E+\frac{2}{E}} \cdot (x_4 x_1) y_4 = 1 + \frac{3+E+1}{E+\frac{2}{E}} \cdot (E + \frac{1}{E}) 4 = 100.98 > E$.
 - Each of the next segments will contain exactly two keys (where the first key appears once, and the second key appears E+1 times), since otherwise the error for the second key will be $\frac{1+E+1}{E+\frac{1}{E}} \cdot E 1 = 100.98 > E$. Just like before, the E+1 repetitions of a single key will cause a violation of the error (due to the spacing between subsequent keys).

Therefore, the algorithm will create N+2 segments given this input.

On the other hand, the optimal algorithm will need only 2 segments: the first segment is the first key, and the second segment covers the rest of the input since the line starting at the second key and ending at the last key is never further away then E from any key, due to the construction of the input. The slope of the second segment will be $\frac{3+(N+1)\cdot(E+2)}{E+(N+1)\cdot(E+\frac{1}{E})}$, and the first key in the segment is

about $\frac{E}{2}$ away on the x axis from the first repeated key. Since the repeated keys are spaced evenly (distance on the x axis of $E + \frac{1}{E}$), the linear function will not violate the error threshold for any key. An illustration is shown in Figure 14.

Since N can be arbitrarily large, the algorithm is not competitive.

П