

An Introduction to B^ϵ -trees and Write-Optimization

MICHAEL A. BENDER, MARTIN FARACH-COLTON, WILLIAM JANNEN, ROB JOHNSON, BRADLEY C. KUSZMAUL, DONALD E. PORTER, JUN YUAN, AND YANG ZHAN



Michael A. Bender is a Professor of Computer Science at Stony Brook University in Stony Brook, New York. His research focuses on algorithms, particularly on out-of-core algorithms. Bender co-founded the database company Tokutek, which was recently acquired by Percona. He has won several awards, including an R&D 100 award, a Test of Time award, a Best Paper award, a Best Newcomer award, and five teaching awards. bender@cs.stonybrook.edu



Martin Farach-Colton is a Professor of Computer Science at Rutgers University, New Brunswick, New Jersey. His research focuses on both the theory and practice of external memory and storage systems. He was a pioneer in the theory of cache oblivious analysis. His current research focuses on the use of write optimization to improve performance in both read- and write-intensive big data systems. He has also worked on the algorithmics of strings and metric spaces, with applications to bioinformatics. In addition to his academic work, Professor Farach-Colton has extensive industrial experience. He is CTO and co-founder of Tokutek, a database company that was founded to commercialize his research. During 2000-2002, he was a Senior Research Scientist at Google. farach@cs.rutgers.edu



William Jannen is a PhD student at Stony Brook University, where he attempts to design systems that accommodate the physical characteristics of their underlying media. He is also an artist and a player of games. wjannen@cs.stonybrook.edu

A B^ϵ -tree is an example of a *write-optimized* data structure and can be used to organize on-disk storage for an application, such as a database or file system. A B^ϵ -tree provides a key-value API, similar to a B-tree, but with better performance, particularly for inserts, range queries, and key-value updates. This article describes the B^ϵ -tree, compares its asymptotic performance to B-trees and Log-Structured Merge trees (LSM-trees), and presents real-world performance measurements. After finishing this article, a reader should have a basic understanding of how a B^ϵ -tree works, its performance characteristics, how it compares to other key-value stores, and how to design applications to gain the most performance from a B^ϵ -tree.

B^ϵ -trees

B^ϵ -trees were proposed by Brodal and Fagerberg [1] as a way to demonstrate an asymptotic performance tradeoff curve between B-trees [2] and buffered repository trees [3]. Both data structures support the same operations, but a B-tree favors queries, whereas a buffered repository tree favors inserts.

Researchers, including the authors of this article, have recognized the practical utility of a B^ϵ -tree when configured to occupy the “middle ground” of this curve—realizing query performance comparable to a B-tree but insert performance orders of magnitude faster than a B-tree. The B^ϵ -tree has since been used by both the high-performance, commercial TokuDB database [4] and the BetrFS research file system [5]. For the interested reader, we have created a simple, reference implementation of a B^ϵ -tree, available at <https://github.com/oscarlab/Be-Tree>.

We first explain how the basic operations on a B^ϵ -tree work. We then give the motivation behind these design choices and illustrate how these choices affect the asymptotic analysis.

API and basic structure. A B^ϵ -tree is a B-tree-like search tree for organizing on-disk data, as illustrated in Figure 1. Both B-trees and B^ϵ -trees export a key-value store API:

- ◆ $\text{insert}(k, v)$
- ◆ $\text{delete}(k)$
- ◆ $v = \text{query}(k)$
- ◆ $[v_1, v_2, \dots] = \text{range-query}(k_1, k_2)$

Like a B-tree, the node size in a B^ϵ -tree is chosen to be a multiple of the underlying storage device’s block size. Typical B^ϵ -tree node sizes range from a few hundred kilobytes to a few megabytes. In both B-trees and B^ϵ -trees, internal nodes store pivot keys and child pointers, and leaves store key-value pairs, sorted by key. For simplicity, one can think of each key-value or pivot-pointer pair as being unit size; both B-trees and B^ϵ -trees can store keys and values of different sizes in practice. Thus, a leaf of size B holds B key-value pairs, which we call items below.

The distinguishing feature of a B^ϵ -tree is that internal nodes also allocate some space for a buffer, as shown in Figure 1. The buffer in each internal node is used to store messages, which encode updates that will eventually be applied to items in leaves under this node. This

An Introduction to B^ϵ -trees and Write-Optimization

Rob Johnson is a Research Professor at Stony Brook University and conducts research on security, big data algorithms, and cryptography.

He does theoretical work with an impact on the real world. rob@cs.stonybrook.edu



Bradley C. Kuzmaul is a Research Scientist in the Computer Science and Artificial Intelligence Laboratory at the Massachusetts Institute of

Technology (MIT CSAIL). His research focuses on performance engineering of multicore software as well as on data structures and algorithms that optimize cache and disk I/O. bradley@mit.edu



Donald E. Porter is an Assistant Professor of Computer Science at Stony Brook University in Stony Brook, New York. His research aims to improve

computer system efficiency and security. In addition to recent work on write optimization in file systems, recent projects have developed lightweight guest operating systems for virtual environments, system security abstractions, and efficient data structures for caching.

porter@cs.stonybrook.edu



Jun Yuan is a PhD student in computer science at Stony Brook University in Stony Brook, New York. Her research interest primarily lies in

compiler and system security. In addition to write-optimized file systems, she has recently studied access control on the Android OS.

junyuan@cs.stonybrook.edu



Yang Zhan is a PhD student in the Department of Computer Science at Stony Brook University. His research interests include file system

and system performance.

yazhan@cs.stonybrook.edu

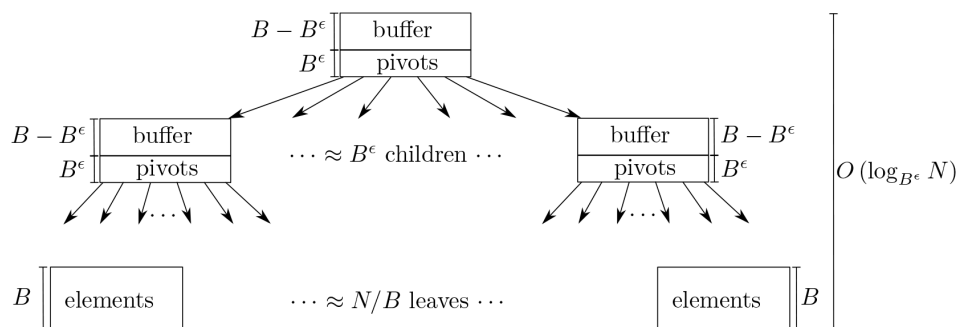


Figure 1: A B^ϵ -tree. Each node is roughly of size B , and ϵ controls how much of an internal node's space is used for pivots (B^ϵ) and how much is used for buffering pending updates ($B - B^\epsilon$). As in a B-tree, items are stored in leaves, and the height of the tree is logarithmic in the total number of items (N), based on the branching factor (here B^ϵ).

buffer is not an in-memory data structure; it is part of the node and is written to disk, evicted from memory, etc., whenever the node is. The value of ϵ , which must be between 0 and 1, is a tuning parameter that selects how much space internal nodes use for pivots ($\approx B^\epsilon$) and how much space is used as a buffer ($\approx B - B^\epsilon$).

Inserts and deletes. Insertions are encoded as “insert messages,” addressed to a particular key and added to the buffer of the root node of the tree. When enough messages have been added to a node to fill the node's buffer, a batch of messages are flushed to one of the node's children. Generally, the child with the most pending messages is selected. Over the course of flushing, each message is ultimately delivered to the appropriate leaf node, and the new key and value are added to the leaf. When a leaf node becomes too full, it splits, just as in a B-tree. Similar to a B-tree, when an interior node gets too many children, it splits and the messages in its buffer are distributed between the two new nodes.

Moving messages down the tree in batches is the key to the B^ϵ -tree's insert performance. By storing newly inserted messages in a buffer near the root, a B^ϵ -tree can avoid seeking all over the disk to put elements in their target locations. The B^ϵ -tree only moves messages to a subtree when enough messages have accumulated for that subtree to amortize the I/O cost. Although this involves rewriting the same data multiple times, this can *improve* performance for smaller, random inserts, as our analysis in the next section shows.

B^ϵ -trees delete items by inserting “tombstone messages” into the tree. These tombstone messages are flushed down the tree until they reach a leaf. When a tombstone message is flushed to a leaf, the B^ϵ -tree discards both the deleted item and the tombstone message. Thus, a deleted item, or even entire leaf node, can continue to exist until a tombstone message reaches the leaf. Because deletes are encoded as messages, deletions are algorithmically very similar to insertions.

A high-performance B^ϵ -tree should detect and optimize the case where a large series of messages all go to one leaf. Suppose a series of keys are inserted that will completely fill one leaf. Rather than write these messages to an internal node only to immediately rewrite them to each node on the path from root to leaf, these messages should flush directly to the leaf, along with any other pending messages for that leaf. The B^ϵ -tree implementation in TokuDB and BetrFS includes some heuristics to avoid writing to intermediate nodes when a batch of messages are all going to a single child.

Point and range queries. Messages addressed to a key k are guaranteed to be applied to k 's leaf or in some buffer along the root-to-leaf path towards key k . This invariant ensures that

An Introduction to B^ϵ -trees and Write-Optimization

point and range queries in a B^ϵ -tree have a similar I/O cost to a B-tree.

In both a B-tree and a B^ϵ -tree, a point query visits each node from the root to the correct leaf. However, in a B^ϵ -tree, answering a query also means checking the buffers in nodes on this path for messages, and applying relevant messages before returning the results of the query. For example, if a query for key k finds an entry (k, v) in a leaf and a tombstone message for k in the buffer of an internal node, then the query will return “NOT FOUND”, since the entry for key k has been logically deleted from the tree. Note that the query need not update the leaf in this case—it will eventually be updated when the tombstone message is flushed to the leaf. A range query is similar to a point query, except that messages for the entire range of keys must be checked and applied as the appropriate subtree is traversed.

In order to make searching and inserting into buffers efficient, the message buffers within each node are typically organized into a balanced binary search tree, such as a red-black tree. Messages in the buffer are sorted by their target key, followed by timestamp. The timestamp ensures that messages are applied in the correct order. Thus, inserting a message into a buffer, searching within a buffer, and flushing from one buffer to another are all fast.

Performance Analysis

We analyze the behavior of B-trees, B^ϵ -trees, and LSM-trees in this article in terms of I/Os. Our primary interest is in data sets too large to fit into RAM. Thus, the first-order performance impact is how many I/O requests must be issued to complete each operation. In the algorithms literature, this is known as the disk-access-machine (DAM) model, the external-memory model, or the I/O model [6].

Performance model. In order to compare B-trees and B^ϵ -trees in a single framework, we make a few simplifying assumptions. We assume that all key-value pairs are the same size and that each node in the tree can hold B key-value pairs. The entire tree stores N key-value pairs. We also assume that each node can be accessed with a single I/O transaction—i.e., on a rotating disk, the node is stored contiguously and requires only one random seek.

This model focuses on the principal performance characteristics of a block storage device, such as a hard drive or SSD. For instance, on a hard drive, this model captures the latency of a random seek to read a node. In the case of an SSD, the model captures the I/O bandwidth costs, i.e., the number of blocks that must be read or written from the device per operation. Regardless of whether the device is bandwidth or latency bound, for a given node size B , minimizing the number of nodes accessed minimizes both bandwidth and latency costs.

B^ϵ -tree I/O performance. Table 1 lists the asymptotic complexities of each operation in a B-tree and B^ϵ -tree. We will explain upserts and epsilon (ϵ), as well as how they affect performance, later in the article. For this discussion, note that ϵ is a tuning parameter between 0 and 1; ϵ is generally set at design time and becomes a constant in the analysis.

The point-query complexities of a B-tree and a B^ϵ -tree are both logarithmic in the number of items ($O(\log_B N)$); a B^ϵ -tree adds a constant overhead of $1/\epsilon$. Compared to a B-tree with the same node size, a B^ϵ -tree reduces the fanout from B to B^ϵ , making the tree taller by a factor of $1/\epsilon$. Thus, for example, querying a B^ϵ -tree where $\epsilon = 1/2$ will require, at most, twice as many I/Os.

Range queries incur a logarithmic search cost for the first key, as well as a cost that is proportional to the size of the range and how many disk blocks the range is distributed across. The scan cost is roughly the number of keys read (k) divided by the block size (B). The total cost of a range query is $O(k/B + \log_B N)$ I/Os.

Compared to a B-tree, batching messages *divides the insertion cost by the batch size* ($B^{1-\epsilon}$). For example, if $B = 1024$ and $\epsilon = 1/2$, a B^ϵ -tree can perform inserts $\approx \epsilon B^{1-\epsilon} = \frac{1}{2} \sqrt{1024} = 16$ times faster than a B-tree.

Write optimization. Batching small, random inserts is an essential feature of write-optimized data structures (WODS), such as a B^ϵ -tree or LSM-tree. Although a WODS may issue a small write multiple times as a message moves down the tree, once the I/O cost is divided among a large batch, the cost per insert or delete is much smaller than one I/O per operation. In contrast, a workload of random inserts to a B-tree requires a *minimum* of one I/O per insert—to write the new element to its target leaf.

The B^ϵ -tree flushing strategy is designed to ensure that it can always move elements in large batches. Messages are only flushed to a child when the buffer of a node is full, containing $B - B^\epsilon \approx B$ messages. When a buffer is flushed, not all messages are necessarily flushed—messages are only flushed to children with enough pending messages to offset the cost of rewriting the parent and child nodes. Specifically, at least $(B - B^\epsilon)/B^\epsilon \approx B^{1-\epsilon}$ messages are moved from the parent’s buffer to the child’s on each flush. Consequently, any node in a B^ϵ -tree is only rewritten if a sufficiently large portion of the node will change.

Caching. Most systems cache a subset of the tree in RAM. With an LRU replacement policy, accesses to the top of the tree are likely to hit in the cache, whereas accesses to leaves and “lower nodes” will more commonly miss. Thus, when the cache is warm, the actual cost of a search may be much less than $O(\log_B N)$ I/Os. For both B-trees and B^ϵ -trees, if only the leaves are out-of-cache, point queries and updates require a single I/O, whereas a range query has an I/O cost that is linear in the number of leaves read.

The Impact of Node Size (B) on Performance

B-trees have small nodes to balance the cost of insertions and range queries. B-tree implementations face a tradeoff between update and range-query performance. A larger node size B favors range queries, and a smaller node size favors inserts and deletes. Larger nodes help range-query performance because the I/O costs, such as seeks, can be amortized over more data. However, larger nodes make updates more expensive because a leaf node and possibly internal nodes must be completely rewritten each time a new item is added to the index, and larger nodes mean more to rewrite.

Thus, many B-tree implementations use small nodes (tens to hundreds of KB), resulting in sub-optimal range-query performance. As free space on disk becomes fragmented, B-tree nodes may also become scattered on disk; this is sometimes called *aging*. Now a range query must seek for each leaf in the scan, resulting in poor bandwidth utilization.

For example, with 4 KB nodes stored on a disk with a 5 ms seek time and 100 MB/s bandwidth, updating a single key only rewrites 4 KB. Range queries, however, must perform a seek for each 4 KB leaf node, resulting in a net bandwidth of 800 KB/s, less than 1% of the disk's potential bandwidth.

B^ϵ -trees have efficient updates and range queries even when nodes are large. In contrast, batching in a B^ϵ -tree allows B to be much larger in a B^ϵ -tree than in a B-tree. In a B^ϵ -tree the bandwidth cost per insert is $O(\frac{\log_B N}{\epsilon B^{1-\epsilon}})$, which grows much more slowly as B increases. As a result, B^ϵ -trees use node sizes of a few hundred kilobytes to a few megabytes.

By using large B , B^ϵ -trees can perform range queries at near disk bandwidth. For example, a B^ϵ -tree using 4 MB nodes need perform only one seek for every 4 MB of data it returns, yielding a net bandwidth of over 88 MB/s on the same disk as above.

In the comparison of insert complexities above, we stated that a B^ϵ -tree with $\epsilon = 1/2$ would be twice as deep as a B-tree, as some

fanout is sacrificed for buffer space. This is only true when the node size is the same. Because a B^ϵ -tree can use larger nodes in practice, a B^ϵ -tree can still have close to the same fanout and height as a B-tree.

The Role of ϵ

The parameter ϵ in a B^ϵ -tree was originally designed to show that there is an optimal tradeoff curve between insert and point query performance. Parameter ϵ ranges between 0 and 1. As we explain in the rest of this section, making ϵ an exponent simplifies the asymptotic analysis and creates an interesting tradeoff curve.

Intuitively, the tradeoff with parameter ϵ is how much space in the node is used for storing pivots and child pointers ($\approx B^\epsilon$) and how much space is used for message buffers ($\approx B - B^\epsilon$). As ϵ increases, so does the branching factor (B^ϵ), causing the depth of the tree to decrease and searches to run faster. As ϵ decreases, the buffers get larger, batching more inserts for every flush and improving overall insert performance.

At one extreme, when $\epsilon = 1$, a B^ϵ -tree is just a B-tree, since interior nodes contain only pivot keys and child pointers. At the other extreme, when $\epsilon = 0$, a B^ϵ -tree is a binary search tree with a large buffer at each node, called a buffered repository tree [3].

The most interesting configurations place ϵ strictly between 0 and 1, such as $\epsilon = 1/2$. For such configurations, a B^ϵ -tree has the same asymptotic point query performance as a B-tree, but asymptotically better insert performance.

For inserts, setting $\epsilon = 1/2$ *divides* the cost by the square root of node size. Formally, the cost then becomes: $O(\frac{\log_B N}{\epsilon B^{1-\epsilon}}) = O(\frac{\log_B N}{\sqrt{B}})$. Since the insert cost is divided by \sqrt{B} , selecting larger nodes (increasing B) can dramatically improve insert performance.

Assuming all other parameters are the same, decreasing ϵ slows down point queries by a constant $1/\epsilon$. To see the query performance for $\epsilon = 1/2$, evaluate the point query cost in Table 1:

$O(\frac{\log_B N}{\epsilon}) = O(\frac{\log_B N}{1/2}) = O(2 \log_B N)$ —doubling the number of I/Os. Changing ϵ from $1/2$ to $1/4$ would make this a factor of 4. This cost can be offset by increasing B , which, as pointed out above, also improves insert performance.

The above analysis assumes all keys have unit size and that nodes can hold B keys; real systems must deal with variable-sized keys, so B , and hence ϵ , are not fixed or known a priori. Nonetheless, the main insight of B^ϵ -trees—that we can speed up insertions by buffering items in internal nodes and flushing them down the tree in batches—still applies in this setting.

Data Structure	Insert	Point Query		Range Query
		no Upserts	w/ Upserts	
B^ϵ -tree	$\frac{\log_B N}{\epsilon B^{1-\epsilon}}$	$\frac{\log_B N}{\epsilon}$	$\frac{\log_B N}{\epsilon}$	$\frac{\log_B N}{\epsilon} + \frac{k}{B}$
B^ϵ -tree ($\epsilon = 1/2$)	$\frac{\log_B N}{\sqrt{B}}$	$\log_B N$	$\log_B N$	$\log_B N + \frac{k}{B}$
B-tree	$\log_B N$	$\log_B N$	$\log_B N$	$\log_B N + \frac{k}{B}$
LSM	$\frac{\log_B N}{\epsilon B^{1-\epsilon}}$	$\frac{\log_B^2 N}{\epsilon}$	$\frac{\log_B^2 N}{\epsilon}$	$\frac{\log_B^2 N}{\epsilon} + \frac{k}{B}$
LSM+BF	$\frac{\log_B N}{\epsilon B^{1-\epsilon}}$	$\log_B N$	$\frac{\log_B^2 N}{\epsilon}$	$\frac{\log_B^2 N}{\epsilon} + \frac{k}{B}$

Table 1: Asymptotic I/O costs of important operations. B^ϵ -trees simultaneously support efficient inserts, point queries (even in the presence of upserts), and range queries. These complexities apply for $0 < \epsilon \leq 1$. Note that ϵ is a design-time constant. We show the complexity for general ϵ and evaluate the complexity when ϵ is set to a typical value of $1/2$. The $1/\epsilon$ factor evaluates to a constant that disappears in the asymptotic analysis.

An Introduction to B^ε-trees and Write-Optimization

In practice, B^ε-tree implementations select a fixed physical node size and fanout (B^ε). For the implementation in TokuDB and BetrFS, nodes are approximately 4 MB, and the branching factor ranges from 4 to 16. As a result, the B^ε-tree can always flush data in batches of at least 256 KB.

How to Speed up Applications by Using a B^ε-tree

A practical consequence of the analysis above is that a B^ε-tree can perform updates orders of magnitude faster than point queries. This search-insert asymmetry has two implications for designing applications on B^ε-trees.

Performance rule. *Avoid query-before-update whenever possible.*

Because of the search-insert asymmetry, the common read-modify-write (or query-modify-insert) pattern will be bound to the speed of a query, which is no faster in a B^ε-tree than in a B-tree.

Upserts. B^ε-trees support a new upsert operation, to help applications bridge this performance asymmetry. An upsert is a type of message that encodes an update with a callback function which does not require first knowing the value of the key.

Upserts can encode any modification that is asynchronous and depends only on the key, the old value, and some auxiliary data that can be stored with the upsert message. Tombstones are a special case of upserts. Upserts can also be used to increment a counter, update the access time on a file, update a user's account balance after a withdrawal, and many other operations.

With upserts, an application can update the value associated with key k in the B^ε-tree by inserting an “upsert message” $(k, (f, \Delta))$ into the tree, where f is a call-back function and Δ is auxiliary data specifying the update to be performed. This upsert message is semantically equivalent to performing a query followed by an insert:

$$v \leftarrow \text{query}(k); \text{insert}(k, f(v, \Delta)).$$

However, the upsert does not perform these operations. Rather, the message $(k, (f, \Delta))$ is inserted into the tree like an insert or tombstone message.

When an upsert message $(k, (f, \Delta))$ is flushed to a leaf, the value v associated with k in the leaf is replaced by $f(v, \Delta)$ and the upsert message is discarded. If the application queries k before the upsert message reaches a leaf, then the upsert message is applied to v before the query returns.

As with any insert or delete message, multiple upserts can be buffered for the same key between the root and leaf. If a key is queried with multiple upserts pending, each upsert must be collected on the path from root to leaf and applied to the key in the order they were inserted into the tree.

The upsert mechanism does not interfere with I/O performance of searches, because the upsert messages for a key k always lie on the search path from the root of the B^ε-tree to the leaf containing k . Thus, the upsert mechanism can accelerate updates by one to two orders of magnitude without slowing down queries.

Performance rule. *Use insert performance to improve query performance by maintaining appropriate indices.*

Secondary indices. In a database, secondary indices can greatly speed up queries. For example, consider a database table with three columns, k_1 , k_2 , and k_3 , and an application that sometimes performs queries using k_1 and sometimes using k_2 . If the table is implemented as a B-tree sorted on k_1 , then queries using k_1 are fast, but queries using k_2 are extremely slow—they may have to scan essentially the entire database. To solve this problem, the table can be configured to maintain two indices—one sorted by k_1 and one sorted by k_2 . Queries can then use the appropriate index based on the type of the query.

When multiple indices are maintained with B-trees, each index update requires an additional insert. Because inserts are as expensive as a point query, maintaining an index on each column is often impractical. Thus, the table designer must carefully analyze factors such as the expected type of queries and distribution of keys in deciding which columns to index, in order to ensure good overall performance.

B^ε-trees turn these issues upside down. Indices are cheap to maintain. Point queries are fundamentally expensive—B^ε-tree point queries are no faster than in a B-tree. Thus, B^ε-tree applications should maintain whatever indices are needed to perform queries efficiently.

There are three rules for designing good B^ε-tree indices.

First, maintain indices sorted by the keys used to query the database. For example, in the above example, the database should maintain two B^ε-trees—one sorted by k_1 and one sorted by k_2 .

Second, ensure that each index has all the information required to answer the intended queries. For example, if the application looks up the k_3 value using key k_2 , then the index sorted by k_2 should store the corresponding k_3 value for each entry. In many databases, the secondary index contains only keys into the primary index. Thus, for example, a query on k_2 would return the primary key value, k_1 . To obtain k_3 , the application has to perform another query in the primary index using the k_1 value obtained from the secondary index. An index that contains all the information relevant to a query is called a covering index for that query.

Finally, design indices to enable applications to perform range queries whenever possible. For example, if the application wants to look up all entries (k_1, k_2, k_3) for which $a \leq k_1 \leq b$, and k_2 satisfies

some predicate, then the application should maintain a secondary index sorted by k_1 that only contains entries for which k_2 matches the predicate.

Log-Structured Merge-Trees

Log-structured merge trees (LSM-trees) [7] are WODS with many variants [8, 9]. An LSM-tree typically consists of a logarithmic number of B-trees of exponentially increasing size. Once an index at one level fills up, it is emptied by merging it into the index at the next level. The factor by which each level grows is a tunable parameter comparable to the branching factor (B^ϵ) in a B^ϵ -tree. For ease of comparison, Table 1 gives the I/O complexities of operations in an LSM-tree with growth factor B^ϵ .

LSM-trees can be tuned to have the same insertion complexity as a B^ϵ -tree, but queries in a naively implemented LSM-tree can require $O(\frac{\log_B^2 N}{\epsilon})$ I/Os because the query must be repeated in $O(\log_B N)$ B-trees. Most LSM-tree implementations use Bloom filters to avoid queries in all but one of the B-trees, improving point query performance to $O(\frac{\log_B N}{\epsilon})$ I/Os.

One problem for LSM-trees is that the benefits of Bloom filters do not extend to range queries. Bloom filters are only designed to improve point queries and do not support range queries. Thus, a range query must be done on every level of the LSM-tree—squaring the search overhead in Table 1 and yielding strictly worse asymptotic performance than a B^ϵ -tree or a B-tree.

A second advantage of a B^ϵ -tree over an LSM-tree is that B^ϵ -trees can effectively use upserts, whereas upserts in an LSM-tree will ruin the performance advantage of adding Bloom filters. As discussed above, upserts address a search-insert asymmetry common to any WODS, including LSM-trees. When an application uses upserts, it is possible for a message for that key to be present in every level of the tree containing a pending message for the key. Thus, a subsequent point query will still have to query every level of the tree, defeating the purpose of adding Bloom filters. Note that querying every level of an LSM-tree also squares the overhead compared to a B^ϵ -tree or B-tree, and is more expensive than walking the path from root-to-leaf in a B^ϵ -tree.

In summary, Bloom-filter-enhanced LSM-trees can match the performance of B^ϵ -trees for some but not all workloads. B^ϵ -trees asymptotically dominate LSM-tree performance. In particular, B^ϵ -trees are asymptotically faster than LSM-trees for small range queries and point queries in upsert-intensive workloads.

Performance Comparison

To give a sense of how B^ϵ -trees perform in practice, we present some data from BetrFS, an in-kernel, research file system based on B^ϵ -trees. We compare BetrFS to other file systems, including

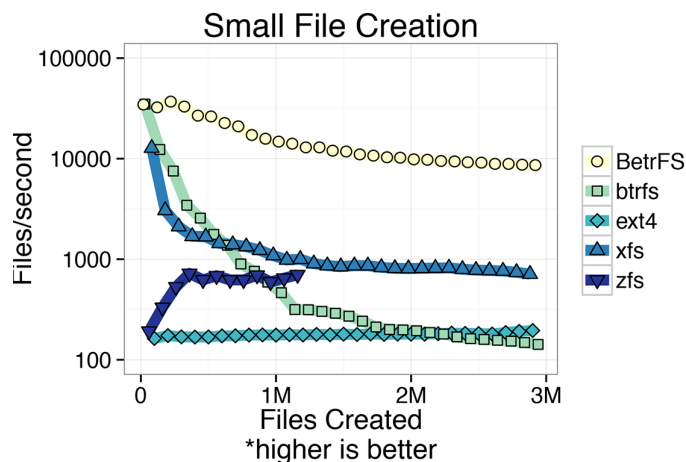


Figure 2: Sustained rate of file creation for 3 million 200-byte files, using four threads. Higher is better.

Btrfs, which is built with B-trees. A more thorough evaluation appears in our recent FAST paper [5].

All experimental results were collected on a Dell Optiplex 790 with a four-core 3.40 GHz Intel Core i7 CPU, 4 GB RAM, and a 250 GB, 7200 RPM ATA disk. Each file system used a 4096-byte block size. The system ran Ubuntu 13.10, 64-bit, with Linux kernel version 3.11.10. Each experiment compared several general-purpose file systems, including Btrfs, ext4, XFS, and ZFS. Error bars and \pm ranges denote 95% confidence intervals. Unless otherwise noted, benchmarks are cold-cache tests.

Small writes. We used the TokuBench benchmark [10] to create 3 million 200-byte files in a balanced directory tree with fanout of 128, using four threads (one per CPU). In BetrFS, file creations are implemented as B^ϵ -tree inserts, and small file writes are implemented using upserts, so this benchmark demonstrates the B^ϵ -tree's performance on these two operations. Figure 2 shows the sustained rate of file creation in each file system (note the log scale). In the case of ZFS, the file system crashed before completing the benchmark, so we reran the experiment five times and used data from the longest-running iteration. BetrFS is initially among the fastest file systems, and continues to perform well for the duration of the experiment. The steady-state performance of BetrFS is an order of magnitude faster than the other file systems.

This performance distinction is attributable to both fewer total writes and fewer seeks per byte written—i.e., better aggregation of small writes. Based on profiling from blktrace, one major distinction is total bytes written: BetrFS writes 4–10x fewer total MB to disk, with an order of magnitude fewer total write requests. Among the other file systems, ext4, XFS, and ZFS wrote roughly the same amount of data, but realized widely varying underlying write throughput.

An Introduction to B^ε-trees and Write-Optimization

FS	find	grep
BetrFS	0.36 ± 0.06	3.95 ± 0.28
Btrfs	3.87 ± 0.94	14.91 ± 1.18
ext4	2.47 ± 0.07	46.73 ± 3.86
XFS	19.07 ± 3.38	66.20 ± 15.99
ZFS	11.60 ± 0.81	41.74 ± 0.64

Table 2: Directory operation benchmarks, measured in seconds. Lower is better.

Locality and directory operations. In BetrFS, fast range queries translate to fast large directory scans. Table 2 reports the time taken to run “find” and “grep -r” on the Linux 3.11.10 source tree, starting from a cold cache. The grep test recursively searches the file contents for the string “cpu_to_be64”, and the find test searches for files named “wait.c”.

Both the find and grep benchmarks do well because file system data and metadata are stored in large nodes and sorted lexicographically by full path. Thus, related files are stored near each other on disk. BetrFS also maintains a second index that contains only metadata, so that metadata scans can be implemented as range queries. As a result, BetrFS can search directory metadata and file data one or two orders of magnitude faster than the other file systems.

Limitations. It is important to note that BetrFS is still a research prototype and currently has three primary cases where it performs considerably worse than other file systems: large directory renames, large deletes, and large sequential writes (more

details in [5]). Renames and deletes are slow because BetrFS does not map them directly onto B^ε-tree operations. Sequential writes are slow largely because the underlying B^ε-tree appends all data to a log before inserting it into the tree, so everything is written to disk at least twice. We believe these issues can be addressed in ongoing research and development efforts; our goal, supported by the asymptotic analysis, is for BetrFS to match or exceed the performance of other file systems on all workloads.

Conclusion

B^ε-tree implementations can match the search performance of B-trees, perform inserts and delete orders of magnitude faster, and execute range queries at near disk bandwidth. The design and implementation of B^ε-trees converts a tradeoff between update and range query costs into a mutually beneficial synergy between batching small updates and large nodes. Our results with BetrFS demonstrate that the asymptotic improvements of B^ε-trees can yield practical performance improvements for applications that are designed for B^ε-tree’s performance characteristics.

Acknowledgments

This work was supported in part by NSF grants CNS-1409238, CNS-1408782, CNS-1408695, CNS-1405641, CNS-1149229, CNS-1161541, CNS-1228839, CNS-1408782, IIS-1247750, CCF-1314547, Sandia National Laboratories, and the Office of the Vice President for Research at Stony Brook University.

References

- [1] G. S. Brodal and R. Fagerberg, “Lower Bounds for External Memory Dictionaries,” in *Proceedings of the 14th Annual ACM-SIAM Symposium on Discrete Algorithms (ACM)*, 2003, pp. 546–554.
- [2] D. Comer, “The Ubiquitous B-tree,” *ACM Computing Surveys*, vol. 11, June 1979, pp. 121–137.
- [3] A. L. Buchsbaum, M. Goldwasser, S. Venkatasubramanian, and J. R. Westbrook, “On External Memory Graph Traversal,” in *Proceedings of the 11th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2000, pp. 859–860.
- [4] Tokutek, Inc., TokuDB: MySQL Performance, MariaDB Performance, 2013: <http://www.tokutek.com/products/tokudb-for-mysql/>.
- [5] W. Jannen, J. Yuan, Y. Zhan, A. Akshintala, J. Esmet, Y. Jiao, A. Mittal, P. Pandey, P. Reddy, L. Walsh, M. Bender, M. Farach-Colton, R. Johnson, B. C. Kuszmaul, and D. E. Porter, “BetrFS: A Right-Optimized Write-Optimized File System,” in *Proceedings of the USENIX Conference on File and Storage Technologies (FAST)*, 2015, pp. 301–315.
- [6] A. Aggarwal and J. S. Vitter, “The Input/Output Complexity of Sorting and Related Problems,” *Communications of the ACM*, vol. 31, Sept. 1988, pp. 1116–1127.
- [7] P. O’Neil, E. Cheng, D. Gawlic, and E. O’Neil, “The Log-Structured Merge-Tree (LSM-tree),” *Acta Informatica*, vol. 33, no. 4, 1996, pp. 351–385.
- [8] R. Sears and R. Ramakrishnan, “bLSM: A General Purpose Log Structured Merge Tree,” in *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, ACM, 2012, pp. 217–228.
- [9] P. Shetty, R. P. Spillane, R. Malpani, B. Andrews, J. Seyster, and E. Zadok, “Building Workload-Independent Storage with VT-trees,” in *Proceedings of the USENIX Conference on File and Storage Technologies (FAST)*, 2013, pp. 17–30.
- [10] J. Esmet, M. A. Bender, M. Farach-Colton, and B. C. Kuszmaul, “The TokuFS Streaming File System,” in *Proceedings of the 4th USENIX Workshop on Hot Topics in Storage (HotStorage)*, June 2012.