

CS 764: Topics in Database Management Systems Lecture 14: Blink Tree

Xiangyao Yu 10/21/2025

Today's Paper: B-tree Locking

Efficient Locking for Concurrent Operations on B-Trees

PHILIP L. LEHMAN
Carnegie-Mellon University
and
S. BING YAO
Purdue University

The B-tree and its variants have been found to be highly useful (both theoretically and in practice) for storing large amounts of information, especially on secondary storage devices. We examine the problem of overcoming the inherent difficulty of concurrent operations on such structures, using a practical storage model. A single additional "link" pointer in each node allows a process to easily recover from tree modifications performed by other concurrent processes. Our solution compares favorably with earlier solutions in that the locking scheme is simpler (no read-locks are used) and only a (small) constant number of nodes are locked by any update process at any given time. An informal correctness proof for our system is given.

Key Words and Phrases: database, data structures, B-tree, index organizations, concurrent algorithms, concurrency controls, locking protocols, correctness, consistency, multiway search trees CR Categories: 3.73, 3.74, 4.32, 4.33, 4.34, 5.24

1. INTRODUCTION

The B-tree [2] and its variants have been widely used in recent years as a data structure for storing large files of information, especially on secondary storage devices [7]. The guaranteed small (average) search, insertion, and deletion time for these structures makes them quite appealing for database applications.

A topic of current interest in database design is the construction of databases that can be manipulated concurrently and correctly by several processes. In this

Agenda

B-Tree Index

Lock coupling

Blink-tree

- Search
- Insert

Optimistic lock coupling (OLC)

Agenda

B-Tree Index

Lock coupling

Blink-tree

- Search
- Insert

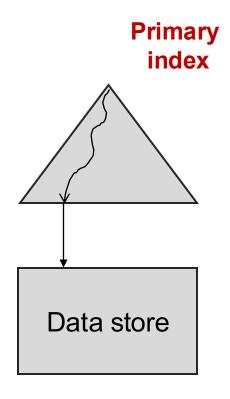
Optimistic lock coupling (OLC)

Index: Accelerate data retrieval operations in a database table

- E.g., random lookup, range scan

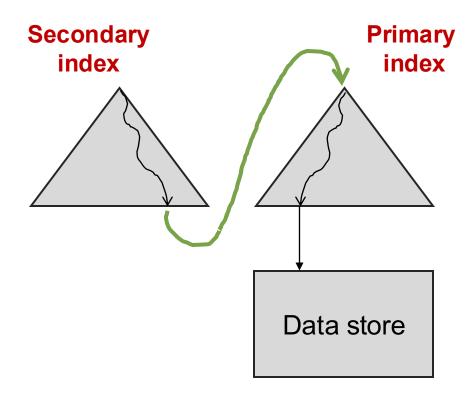
Index: Accelerate data retrieval operations in a database table

- E.g., random lookup, range scan



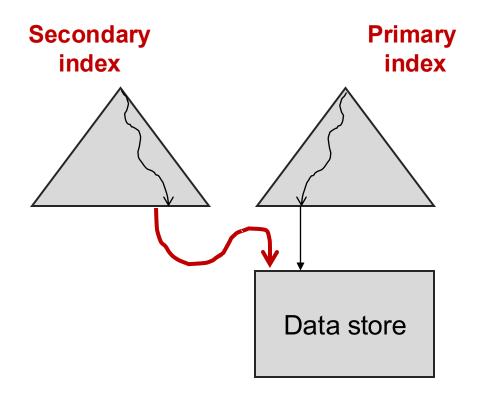
Index: Accelerate data retrieval operations in a database table

- E.g., random lookup, range scan
- Secondary index usually stores the primary key



Index: Accelerate data retrieval operations in a database table

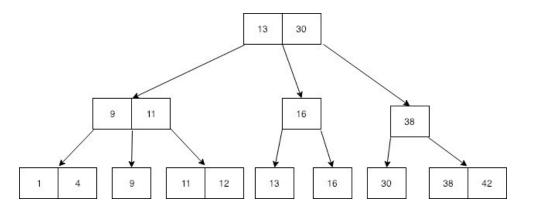
- E.g., random lookup, range scan
- Secondary index can also point to the record directly



B-tree

Balanced tree data structure

- Data is sorted
- Supports: search, sequential scan, inserts, and deletes



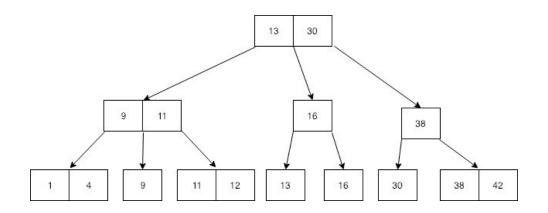
B-tree

Balanced tree data structure

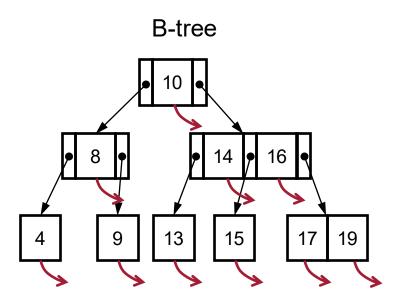
- Data is sorted
- Supports: search, sequential scan, inserts, and deletes

Properties

- Every node contains *k* to 2*k* keys (except root)
- All leaf nodes are at the same level
- k is typically large; a lookup traverses a small number of levels

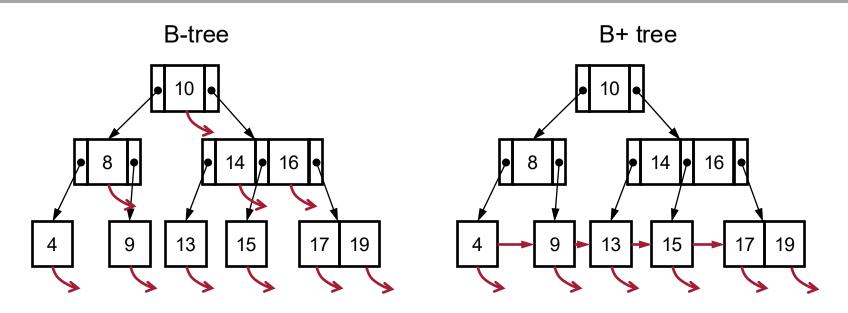


B-tree vs. B+ Tree vs. B* Tree



B-tree: data pointers stored in all nodes

B-tree vs. B+ Tree vs. B* Tree

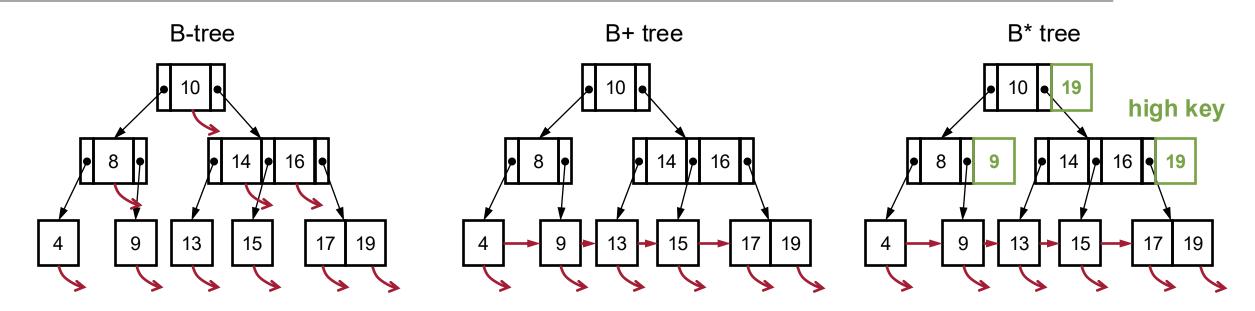


B-tree: data pointers stored in all nodes

B+ tree:

- Data pointers stored only in leaf nodes
- The leaf nodes are linked

B-tree vs. B+ Tree vs. B* Tree



B-tree: data pointers stored in all nodes

B+ tree:

- Data pointers stored only in leaf nodes
- The leaf nodes are linked

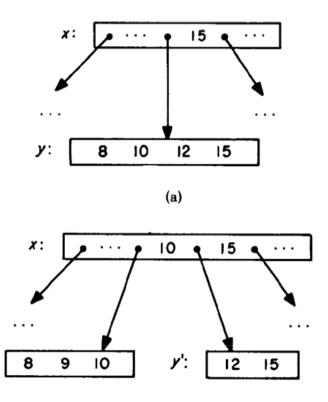
B* tree is a misused term in B-tree literature

- Typically means a variant of B+ tree in which each node is least 2/3 full
- In this paper: B+ tree with high key appended to non-leaf nodes (upper bound on values)

Insert Example

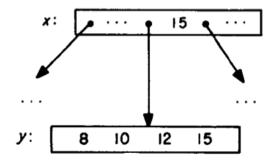
Assume k = 2 (at most 4 keys per node)

insert(9) $A \leftarrow \operatorname{read}(x)$ examine A; get ptr to y $A \leftarrow \operatorname{read}(y)$ insert 9 into A; must split into A, B $\operatorname{put}(B, y')$ $\operatorname{put}(A, y)$ Add to node x a pointer to node y'.



Search Example

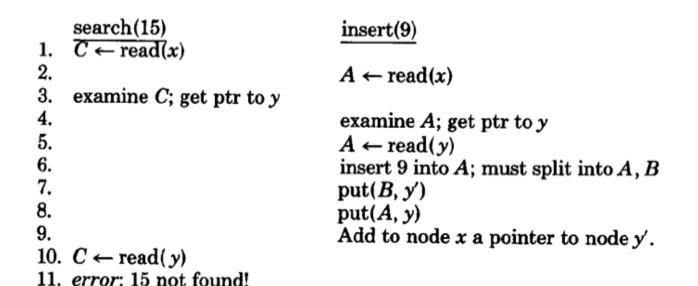
Assume k = 2 (at most 4 keys per node)

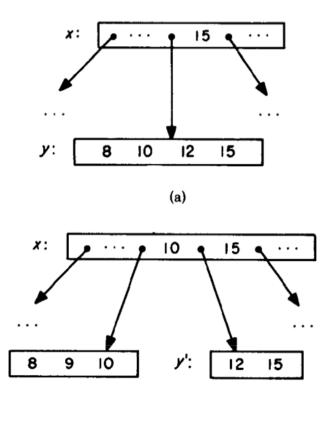


```
search(15)
1. C \leftarrow \operatorname{read}(x)
2.
3. examine C; get ptr to y
4.
5.
6.
7.
8.
9.
10. C \leftarrow \operatorname{read}(y)
```

Concurrency Challenge

Assume k = 2 (at most 4 keys per node) Concurrent search and insert can cause problems





Agenda

B-Tree Index

Lock coupling

Blink-tree

- Search
- Insert

Optimistic lock coupling (OLC)

Lock Coupling

A node is **unsafe** (wrt. insertion) if it is full (i.e., contains 2k keys)

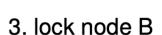
Lock Coupling

A node is **unsafe** (wrt. insertion) if it is full (i.e., contains 2k keys)

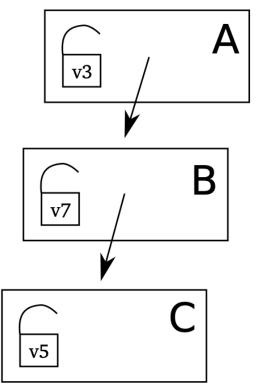
Lock coupling (aka. lock crabbing)

- Lock parent
- Access parent
- Lock child
- Release parent if child is safe

- 1. lock node A
- 2. access node A



- 4. unlock node A
- 5. access node B
- 6. lock node C
- 7. unlock node B
- 8. access node C
- 9. unlock node C



Lock Coupling

A node is **unsafe** (wrt. insertion) if it is full (i.e., contains 2k keys)

Lock coupling (aka. lock crabbing)

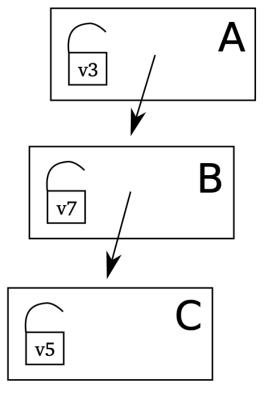
- Lock parent
- Access parent
- Lock child
- Release parent if child is safe
- What if the child is unsafe?

 One solution: split immediately if child is unsafe

- 1. lock node A
- 2. access node A



- 4. unlock node A
- 5. access node B
- 6. lock node C
- 7. unlock node B
- 8. access node C
- 9. unlock node C



Limitation of Lock Coupling

The root is locked for every index access and becomes a scalability bottleneck

Observation: root and upper levels are rarely changed; lock coupling is too conservative

Limitation of Lock Coupling

The root is locked for every index access and becomes a scalability bottleneck

Observation: root and upper levels are rarely changed; lock coupling is too conservative

Concurrency challenge: search may read wrong node due to split

- Lock coupling solution: guard split using a lock
- Blink tree solution: allow search to find the right node

Agenda

B-Tree Index

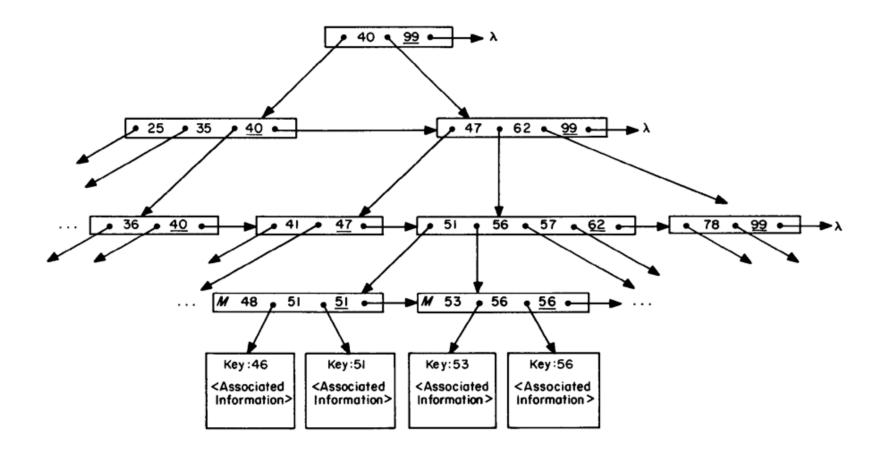
Lock coupling

Blink-tree

- Search
- Insert

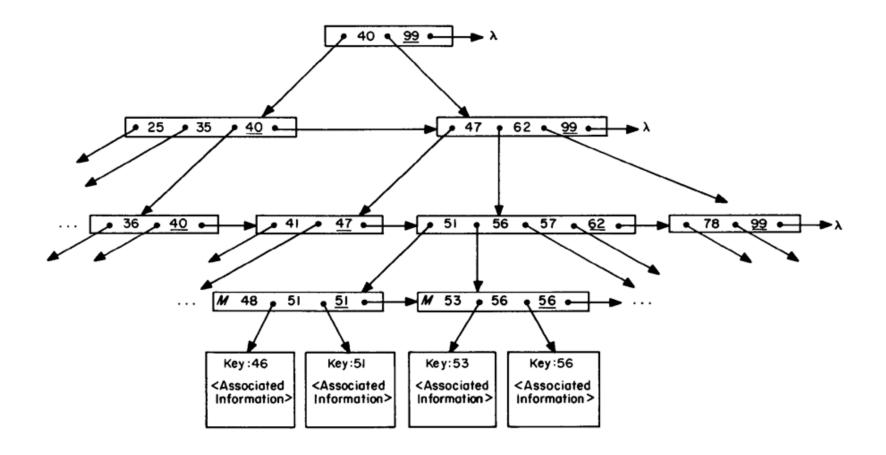
Optimistic lock coupling (OLC)

Blink-Tree



Feature 1: link pointer to next node at each level key idea

Blink-Tree



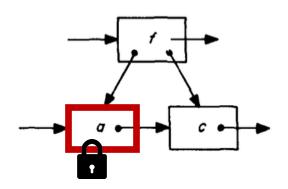
Feature 1: link pointer to next node at each level key idea

Feature 2: high key for each node

Blink-Tree: Locks

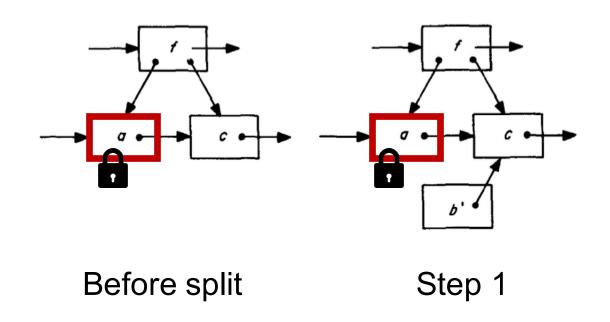
- Hold only a small number of locks at any time
- Reads are not blocked by locks
- Locks prevent only multiple updates

Insert to leaf if the leaf node if not full Illustration of node split (node *a* is split into *a'* and *b'*)

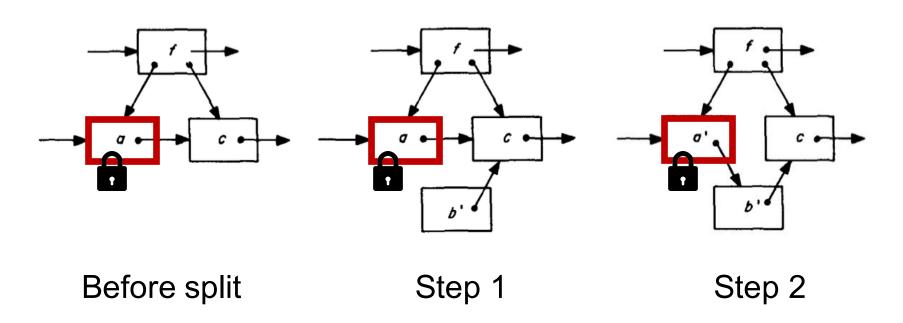


Before split

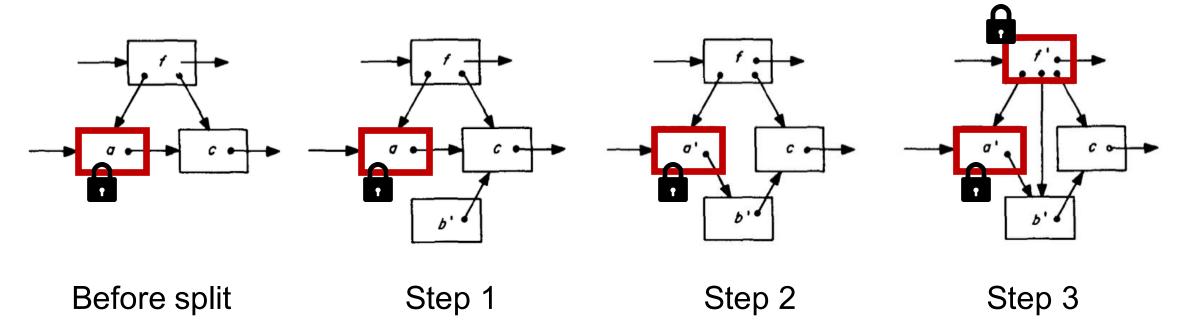
Insert to leaf if the leaf node if not full Illustration of node split (node a is split into a' and b')



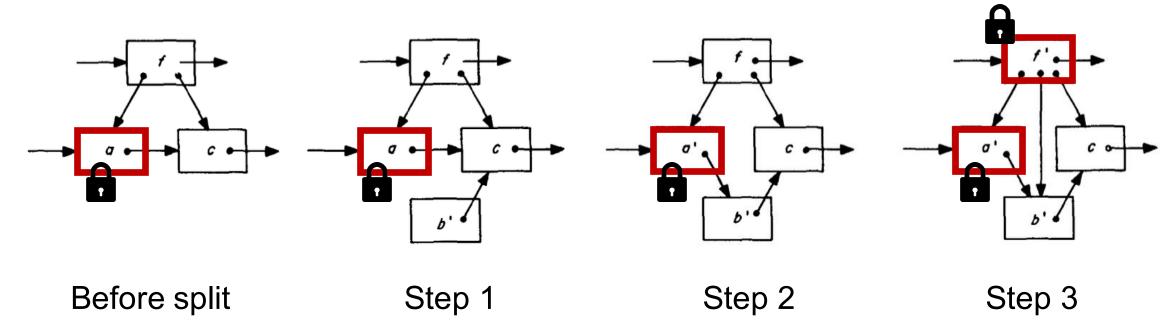
Insert to leaf if the leaf node if not full Illustration of node split (node a is split into a' and b')



Insert to leaf if the leaf node if not full Illustration of node split (node *a* is split into *a'* and *b'*)



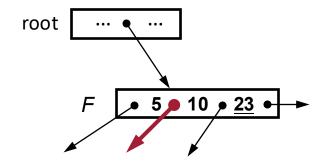
Insert to leaf if the leaf node if not full Illustration of node split (node *a* is split into *a'* and *b'*)



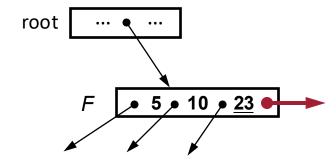
Q: What if another txn searches a key in b' before step 3 finishes?

Blink-Tree: Search Algorithm

May follow the link pointer to find a key



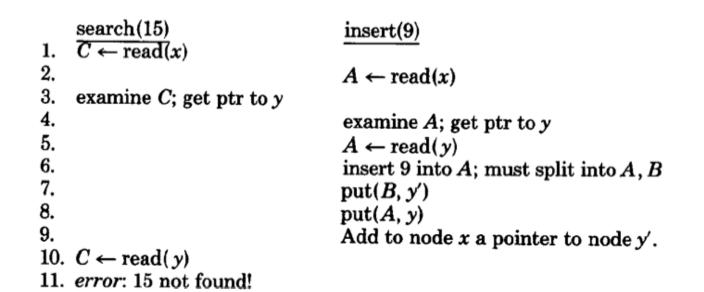
If search for Key=8

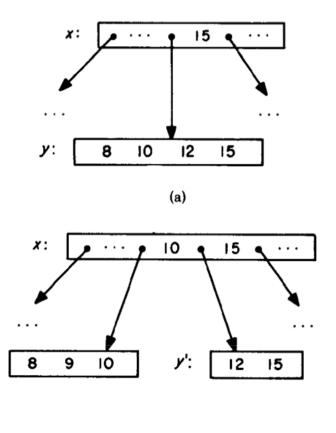


If search for Key=24

Concurrent Search & Insert

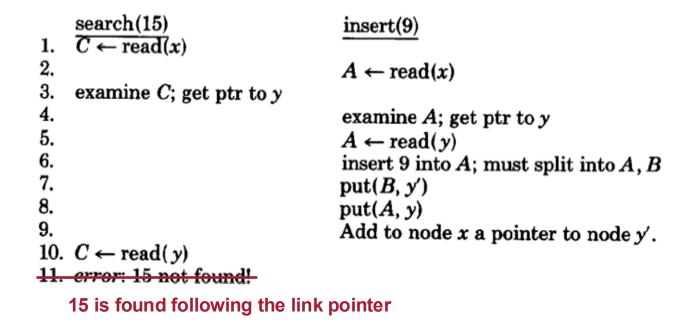
Assume k = 2 (at most 4 keys per node) Concurrency problem is solved in B^{link} tree

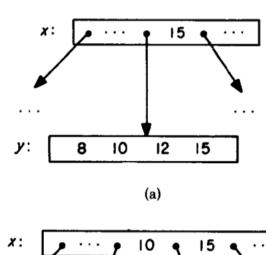


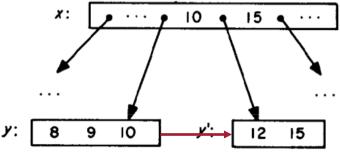


Concurrent Search & Insert

Assume k = 2 (at most 4 keys per node)
Concurrency problem is solved in B^{link} tree
High key indicates when to follow link pointer





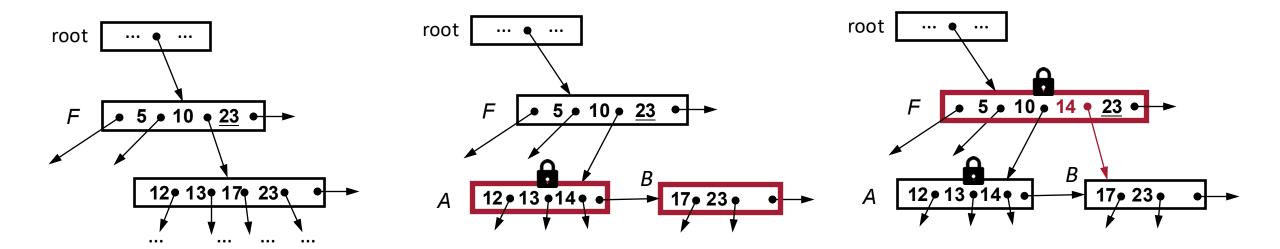


Concurrent Insert & Insert

Before insert 14

Leaf node split

Insert to parent node



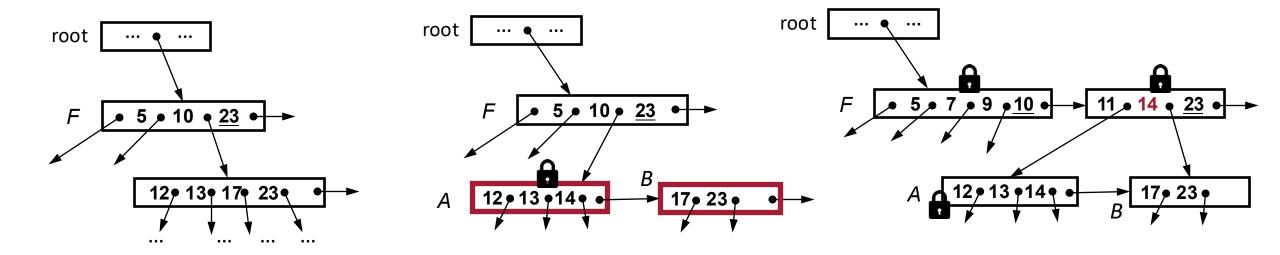
Regular insert process

Concurrent Insert & Insert

Before insert 14

Leaf node split

Insert to parent node



During an insert, the parent node is split by another transaction

- Follow the link point to find the real parent node
- The transaction holds 3 locks in this scenario

Agenda

B-Tree Index

Lock coupling

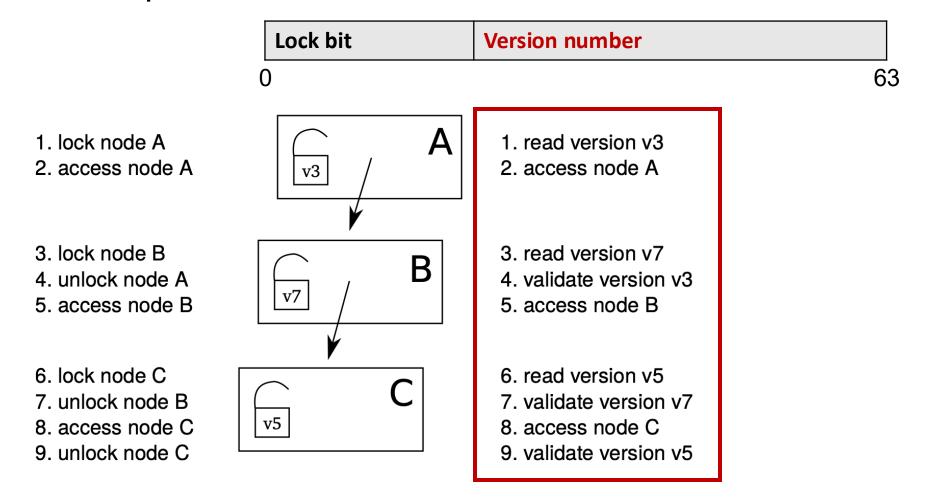
Blink-tree

- Search
- Insert

Optimistic lock coupling (OLC)

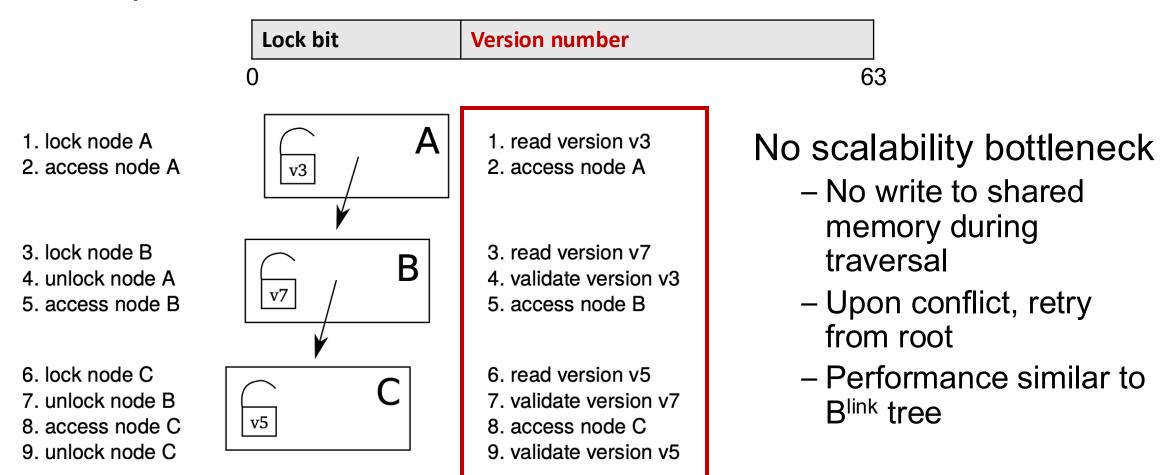
Optimistic Lock Coupling (OLC)

Each tuple contains a 64-bit version counter



Optimistic Lock Coupling (OLC)

Each tuple contains a 64-bit version counter



Evaluation

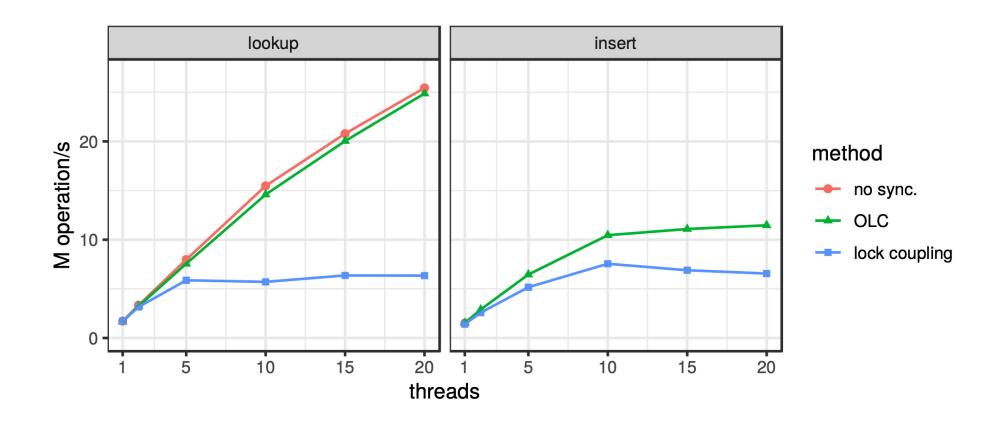


Figure 3: Scalability on 10-core system for B-tree operations (100M values).

Q/A – Blink Tree

Recovery protocol?

Scalability under heavy update workloads?

Repeated splits near the same region?

For high concurrency, contention on upper-level node still negligible?

Excessive link pointer chain-following in the worst case?

Extend to a distributed system?

Next Lecture

Viktor Leis, et al., <u>The Adaptive Radix Tree: ARTful Indexing for Main-Memory Databases</u>. ICDE, 2013