



CS 764: Topics in Database Management Systems

Lecture 15: Blink Tree

Xiangyao Yu

10/26/2022

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

B^{link}-tree

- Search

- Insert

Optimistic lock coupling (OLC)

Agenda

B-Tree Index

Lock coupling

B^{link}-tree

- Search
- Insert

Optimistic lock coupling (OLC)

Index

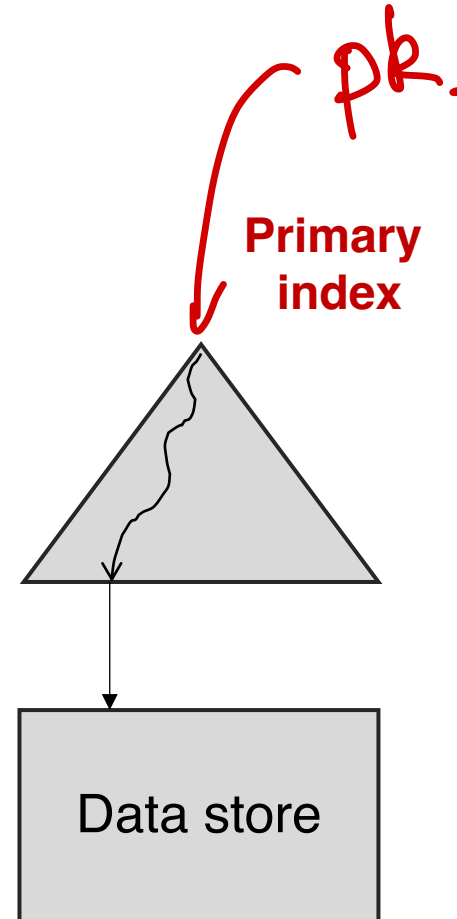
Index: Accelerate data retrieval operations in a database table

– E.g., random lookup, range scan

Index

Index: Accelerate data retrieval operations in a database table

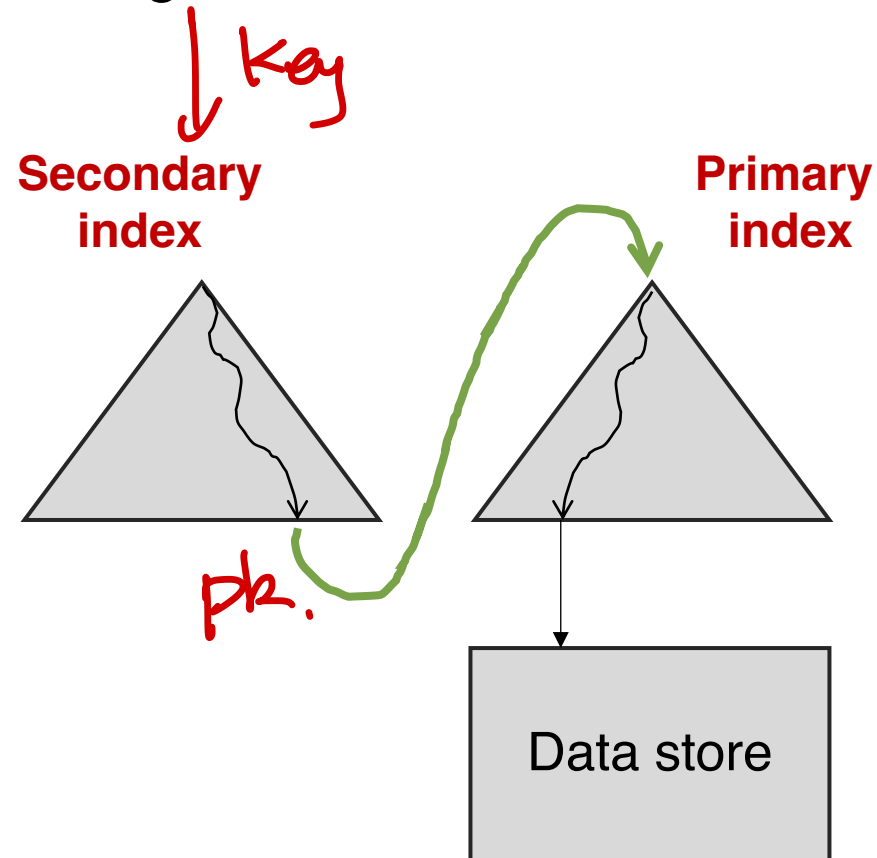
- E.g., random lookup, range scan



Index

Index: Accelerate data retrieval operations in a database table

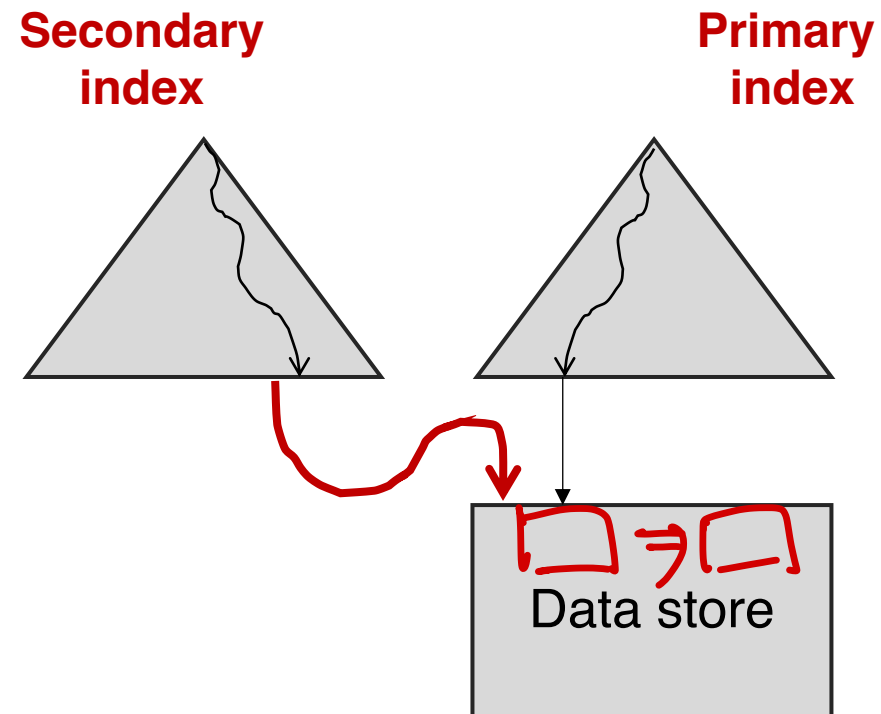
– E.g., random lookup, range scan



Index

Index: Accelerate data retrieval operations in a database table

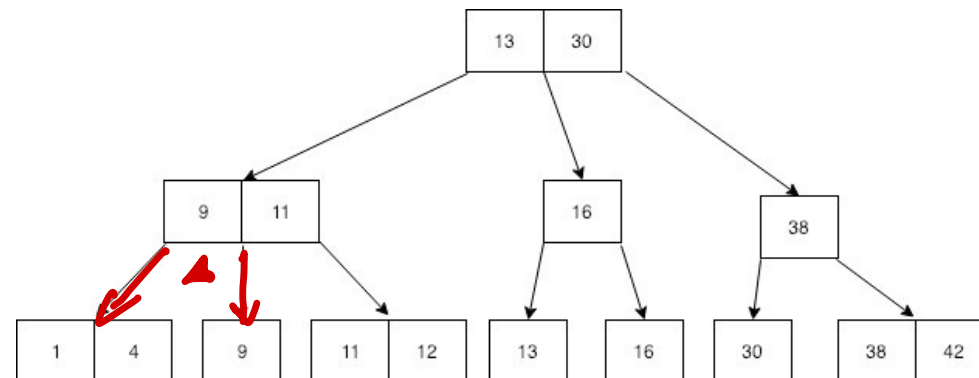
– E.g., random lookup, range scan



B-tree

Balanced tree data structure

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



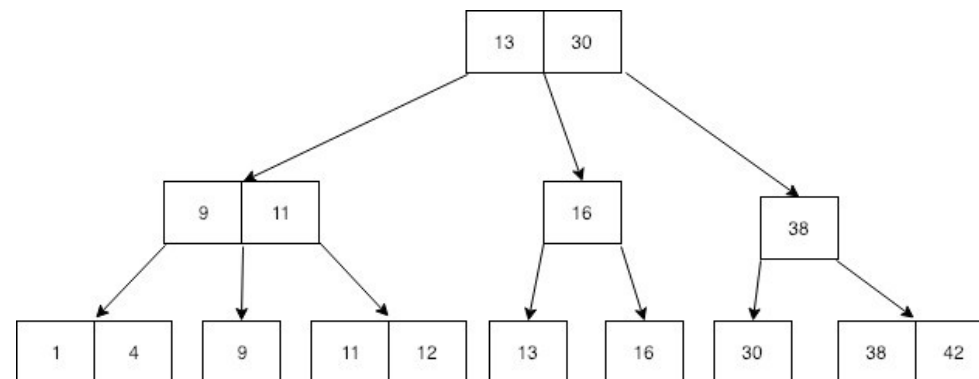
B-tree

Balanced tree data structure

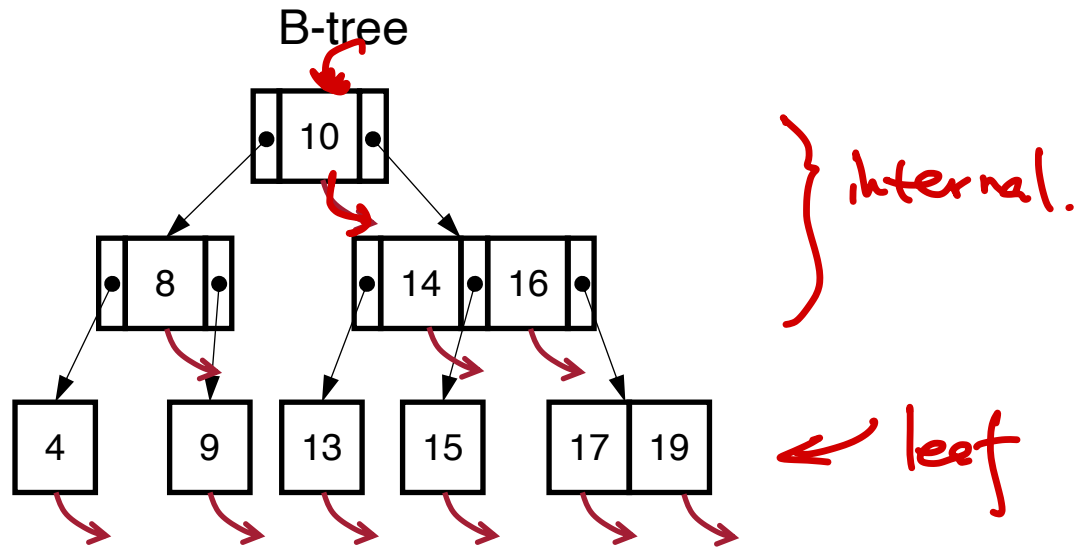
- Data is sorted
- Supports: search, sequential scan, inserts, and deletes

Properties

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

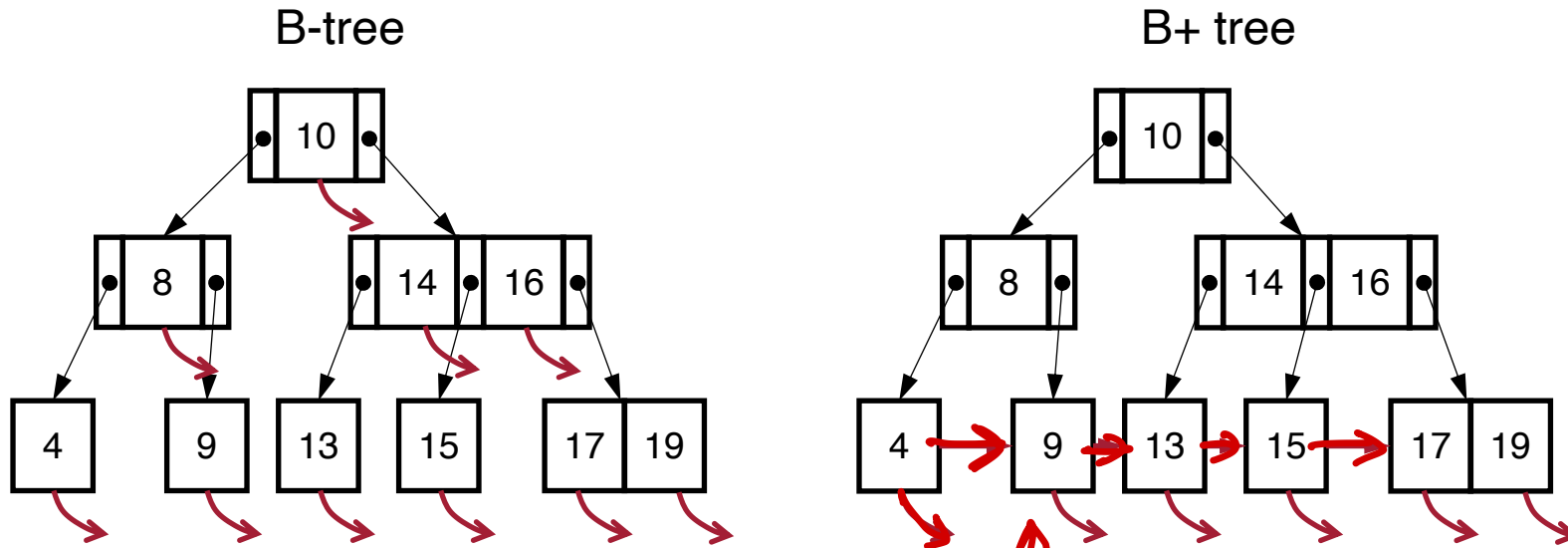


B-tree vs. B+ Tree vs. B* Tree



B-tree: data pointers stored in all nodes

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

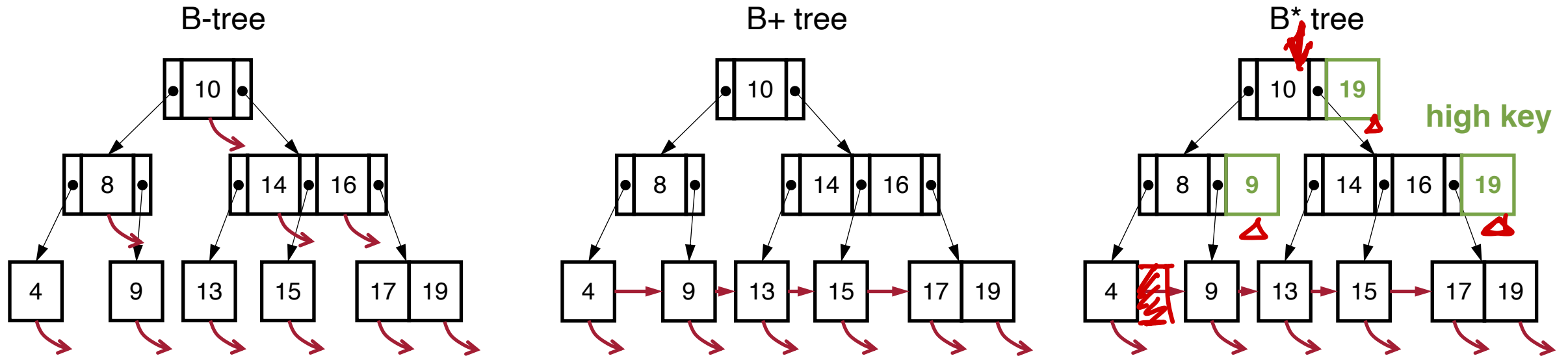


B-tree: data pointers stored in all nodes

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)

$A \leftarrow \text{read}(x)$

examine A ; get ptr to y

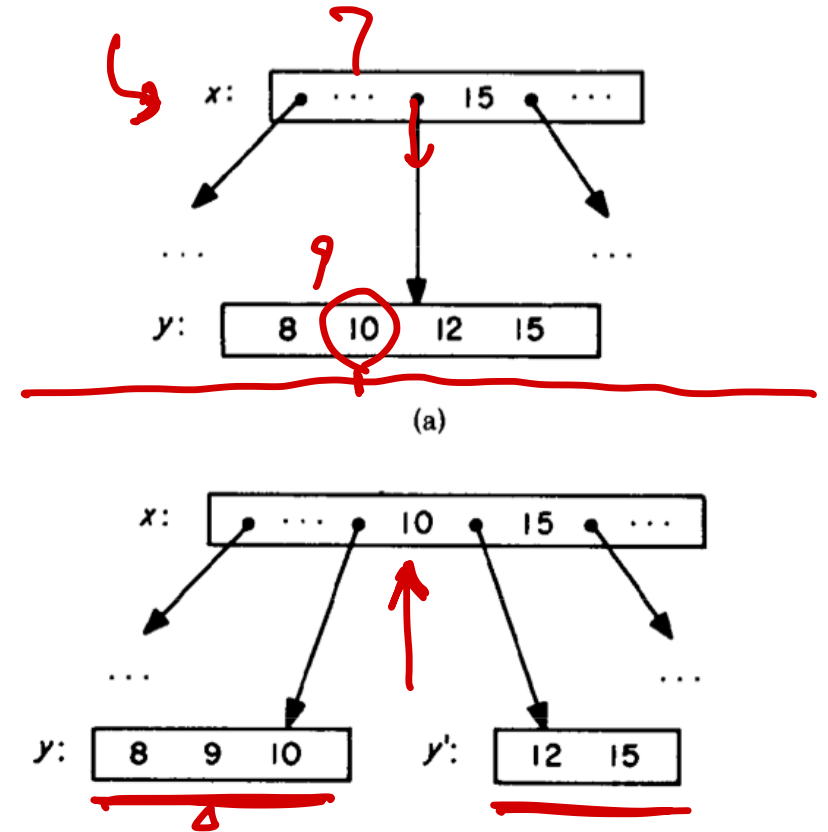
$A \leftarrow \text{read}(y)$

insert 9 into A ; must split into A, B

put(B, y')

put(A, y)

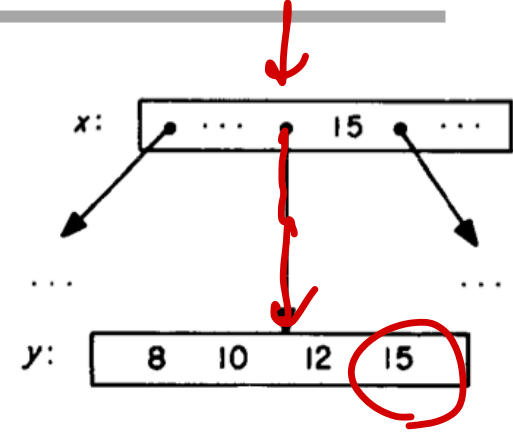
Add to node x a pointer to node y' .



Search Example

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

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

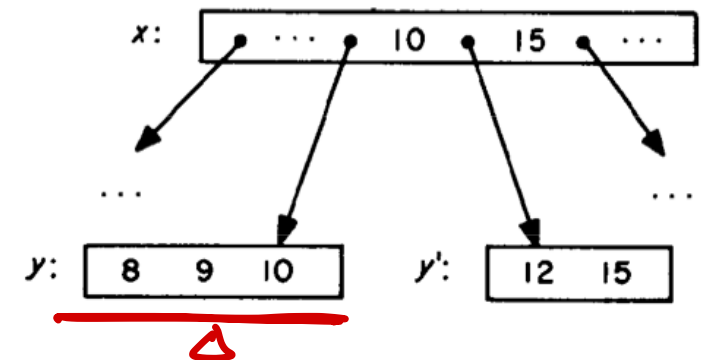
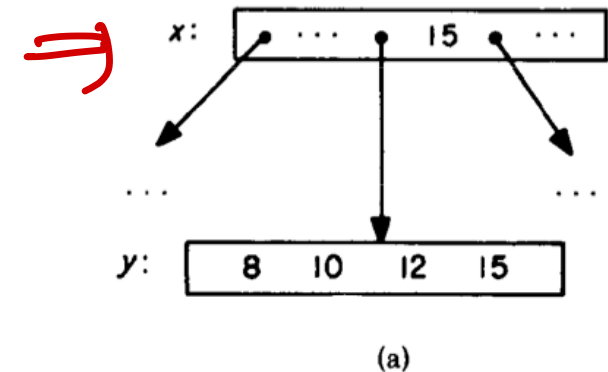


Concurrency Challenge

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

1. search(15)
 $C \leftarrow \text{read}(x)$
- 2.
3. examine C : get ptr to y
- 4.
- 5.
- 6.
- 7.
- 8.
- 9.
10. $C \leftarrow \text{read}(y)$
11. error: 15 not found!

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



Agenda

B-Tree Index

Lock coupling

B^{link}-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

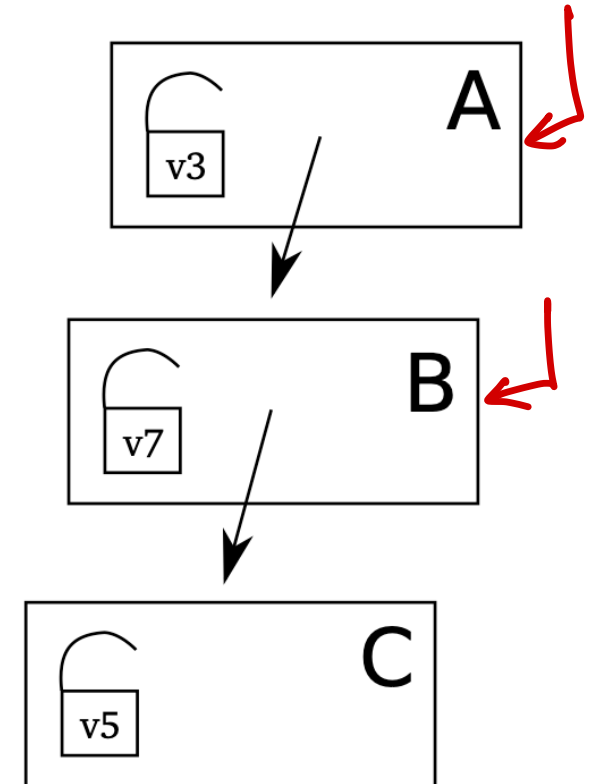
A node is **unsafe** (wrt. insertion) if it is full (i.e., contains $2k$ keys)

Lock coupling (aka. lock crabbing)

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

insert.

1. lock node A
2. access node A
3. lock node B
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

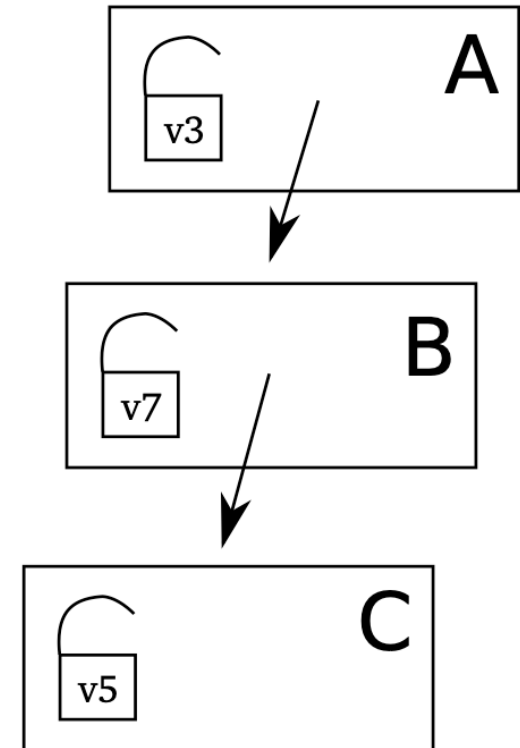
What if the child is unsafe?

- One solution: split immediately if child is unsafe

1. lock node A
2. access node A

3. lock node B
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
- **B^{link} tree solution:** allow search to find the right node

Agenda

B-Tree Index

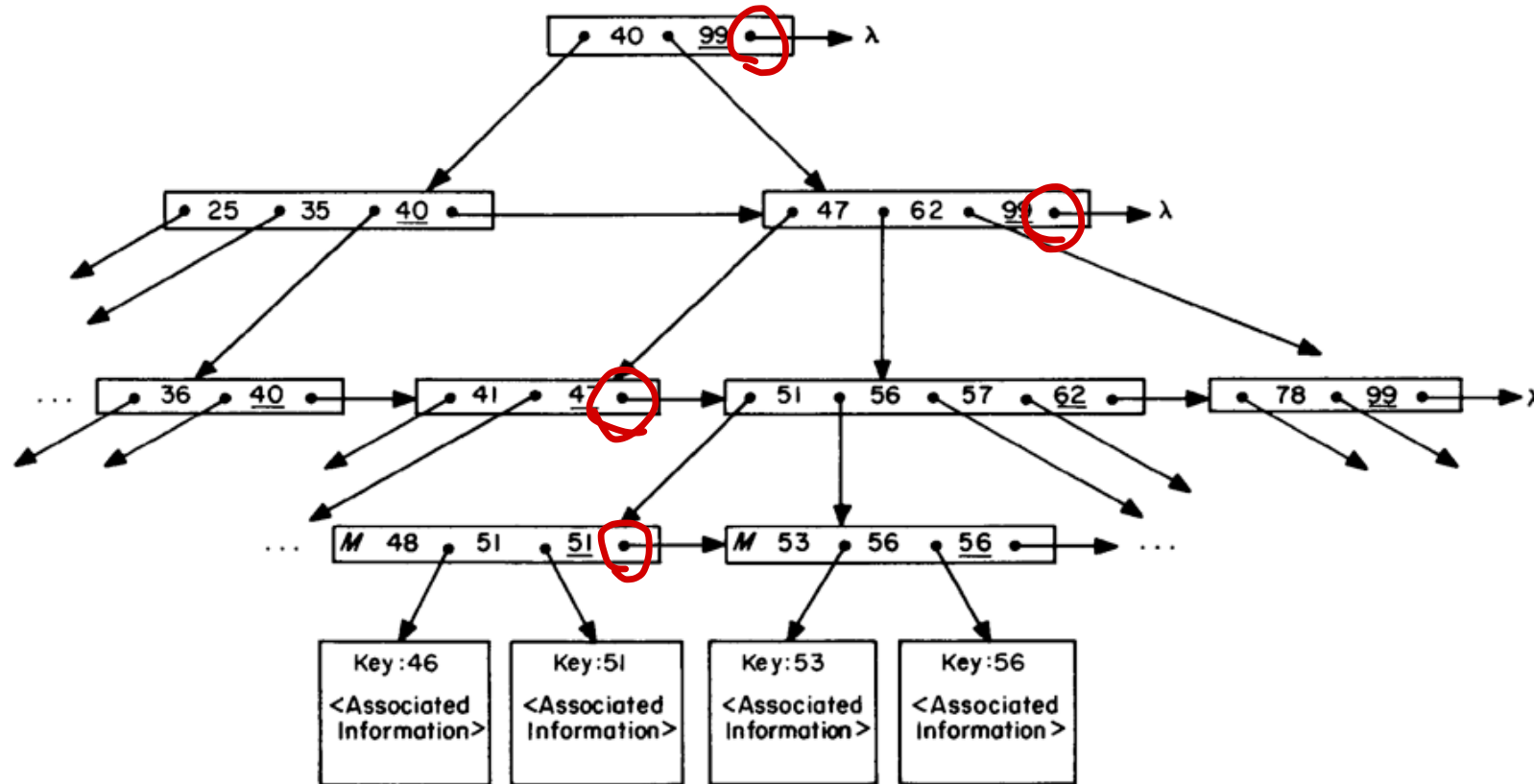
Lock coupling

Blink-tree

- Search
- Insert

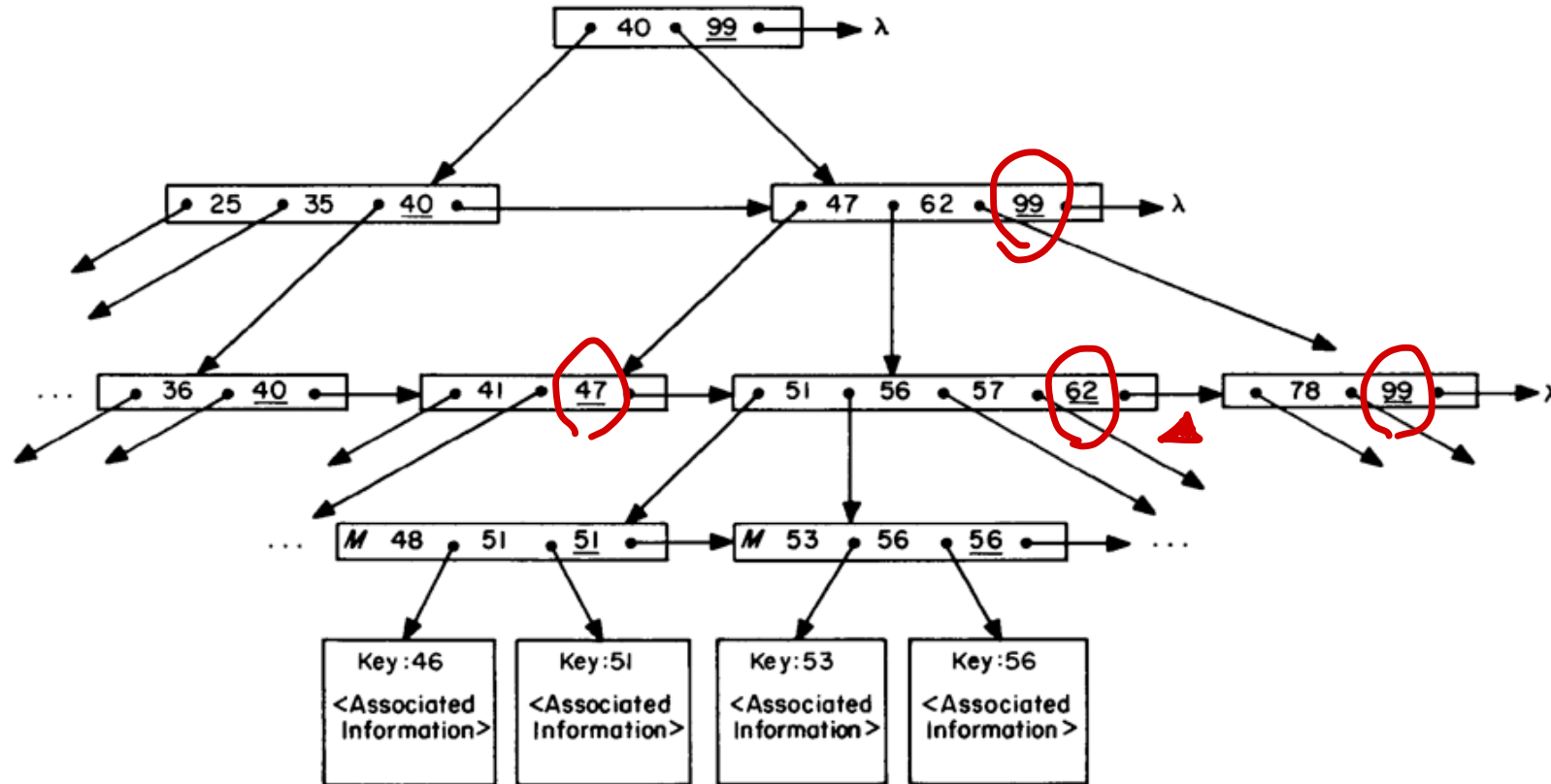
Optimistic lock coupling (OLC)

Blink-Tree



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

Blink-Tree



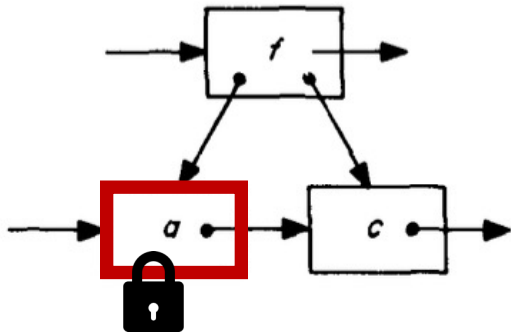
Feature 1: **link pointer** to next node at each level \longrightarrow **key idea**

Feature 2: **high key** for each node

B^{link}-Tree: Insert Algorithm

Insert to leaf if the leaf node is not full

Illustration of node split (node a is split into a' and b')

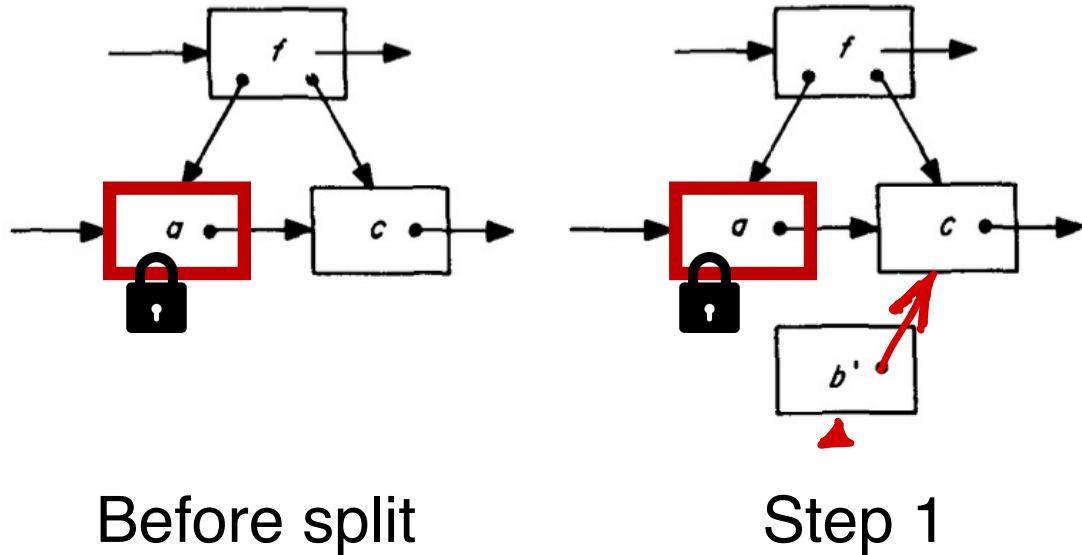


Before split

B^{link}-Tree: Insert Algorithm

Insert to leaf if the leaf node is not full

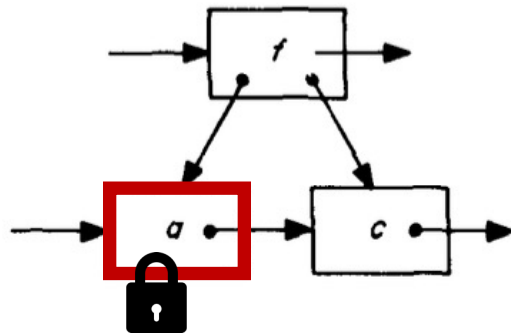
Illustration of node split (node a is split into a' and b')



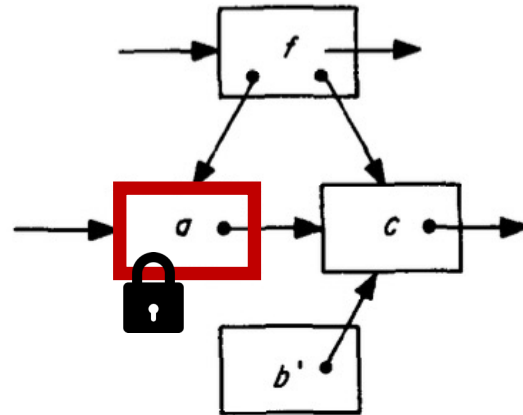
B^{link}-Tree: Insert Algorithm

Insert to leaf if the leaf node is not full

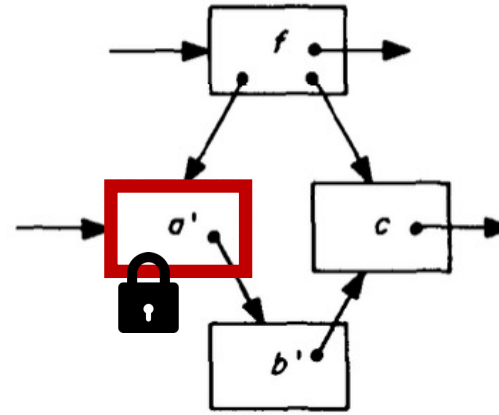
Illustration of node split (node a is split into a' and b')



Before split



Step 1

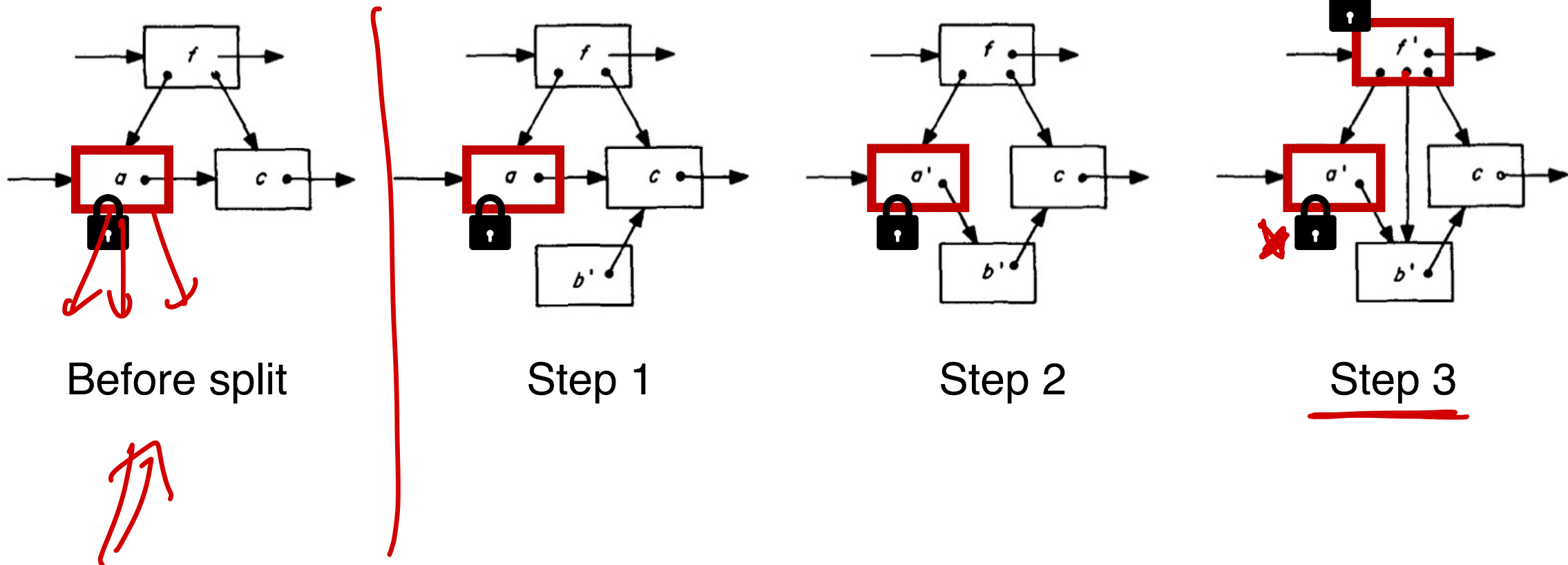


Step 2

B^{link}-Tree: Insert Algorithm

Insert to leaf if the leaf node is not full

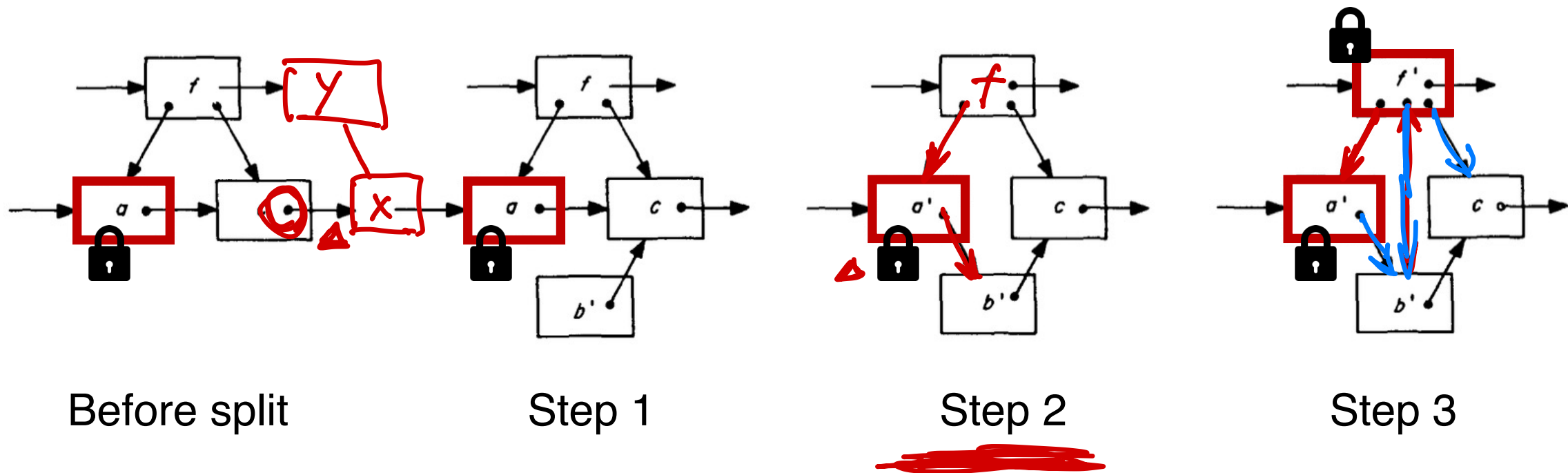
Illustration of node split (node a is split into a' and b')



Blink-Tree: Insert Algorithm

Insert to leaf if the leaf node is not full

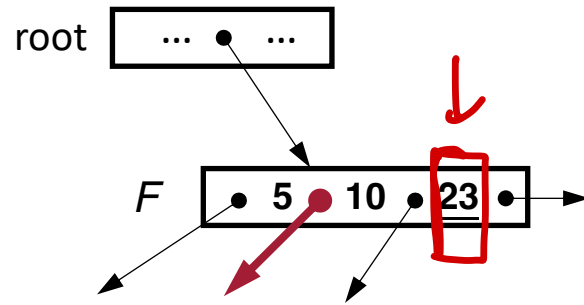
Illustration of node split (node a is split into a' and b')



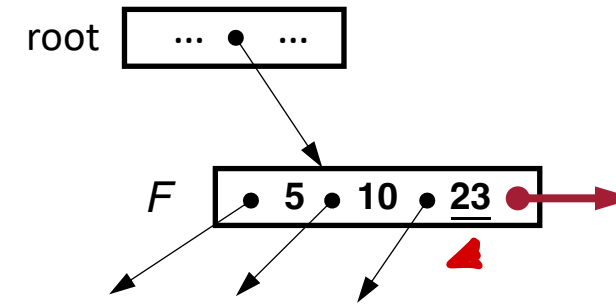
Q: What if another txn searches a key in b' before step 3 finishes?

B^{link}-Tree: Search Algorithm

May follow the link pointer to find a key



If search for Key=8



If search for Key=24

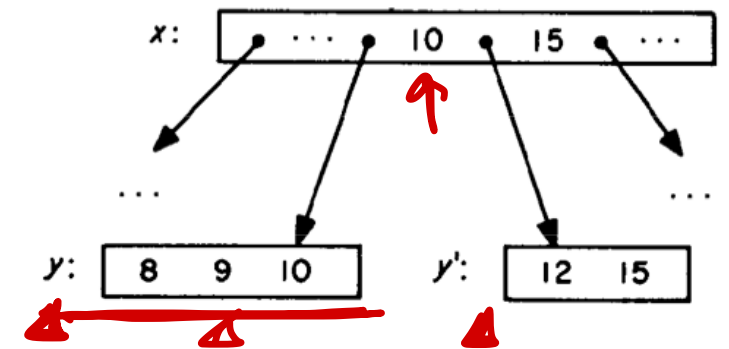
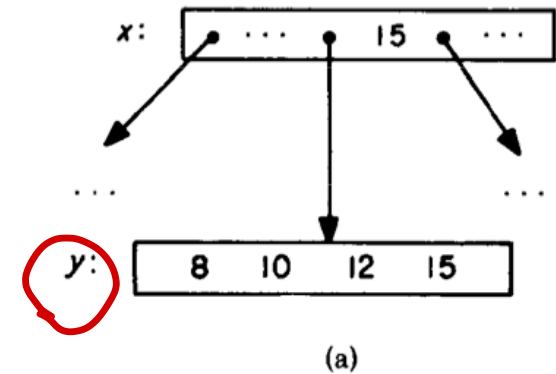
Concurrent Search & Insert

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

Concurrency problem is solved in B^{link} tree

1. search(15)
2. $C \leftarrow \text{read}(x)$
3. examine C ; get ptr to y
4. examine C ; get ptr to y
- 5.
- 6.
- 7.
- 8.
- 9.
10. $C \leftarrow \text{read}(y)$
11. error: 15 not found!

1. insert(9)
2. $A \leftarrow \text{read}(x)$
3. examine A ; get ptr to y
4. $A \leftarrow \text{read}(y)$
5. insert 9 into A ; must split into A, B
6. put(B, y')
7. put(A, y)
8. Add to node x a pointer to node y' .



Concurrent Search & Insert

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

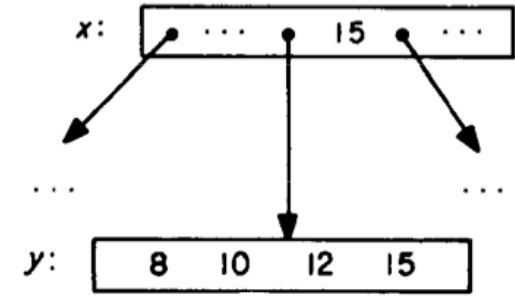
Concurrency problem is solved in B^{link} tree

High key indicates when to follow link pointer

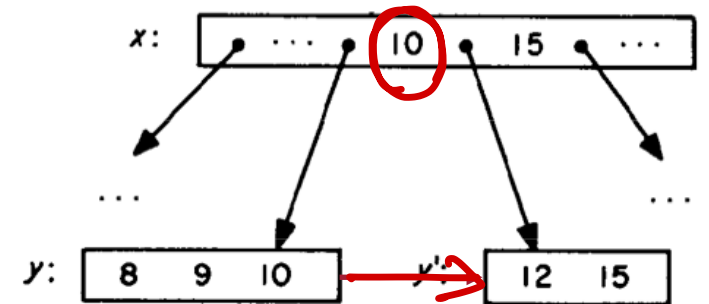
1. $\text{search}(15)$
 $C \leftarrow \text{read}(x)$
- 2.
3. examine C ; get ptr to y
- 4.
- 5.
- 6.
- 7.
- 8.
- 9.
10. $C \leftarrow \text{read}(y)$
11. ~~error: 15 not found!~~

1. $\text{insert}(9)$
 $A \leftarrow \text{read}(x)$
- 2.
3. examine A ; get ptr to y
4. $A \leftarrow \text{read}(y)$
5. insert 9 into A ; must split into A, B
6. $\text{put}(B, y')$
7. $\text{put}(A, y)$
- 8.
9. Add to node x a pointer to node y' .

15 is found following the link pointer

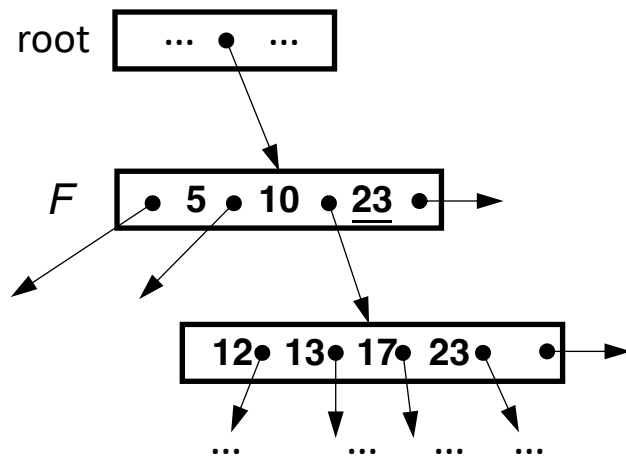


(a)

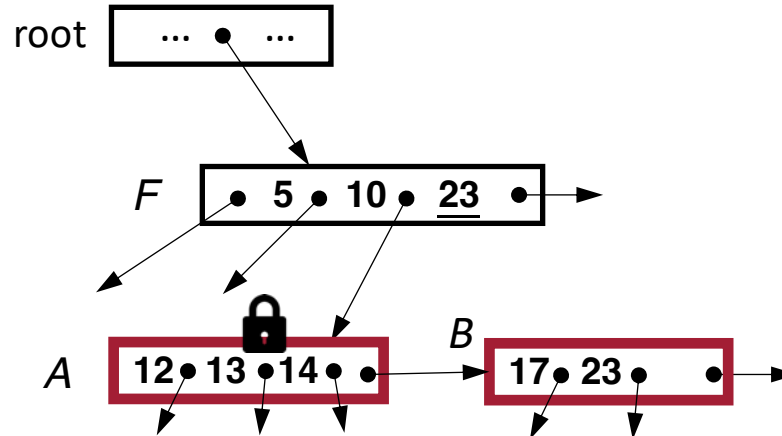


Concurrent Insert & Insert

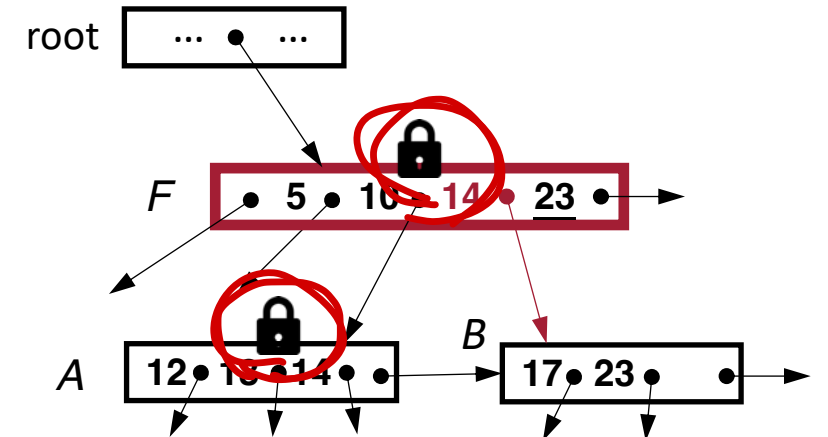
Before insert 14



Leaf node split



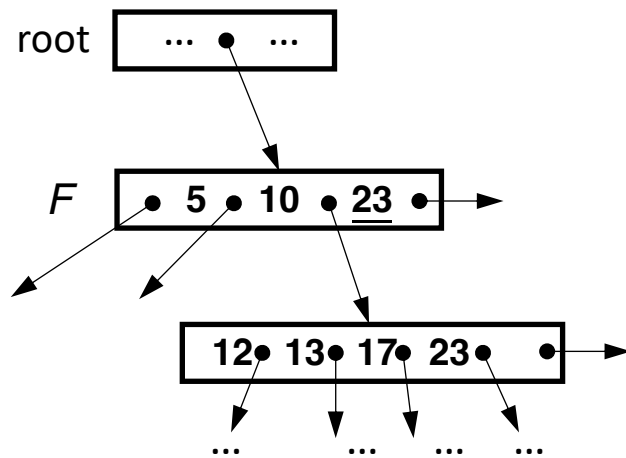
Insert to parent node



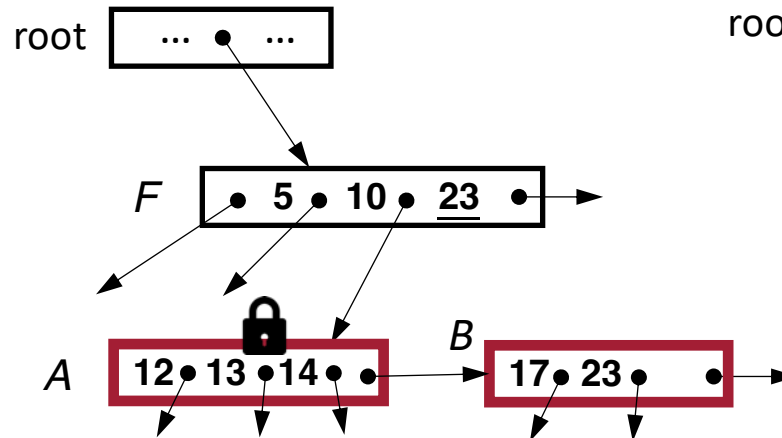
Regular insert process

Concurrent Insert & Insert

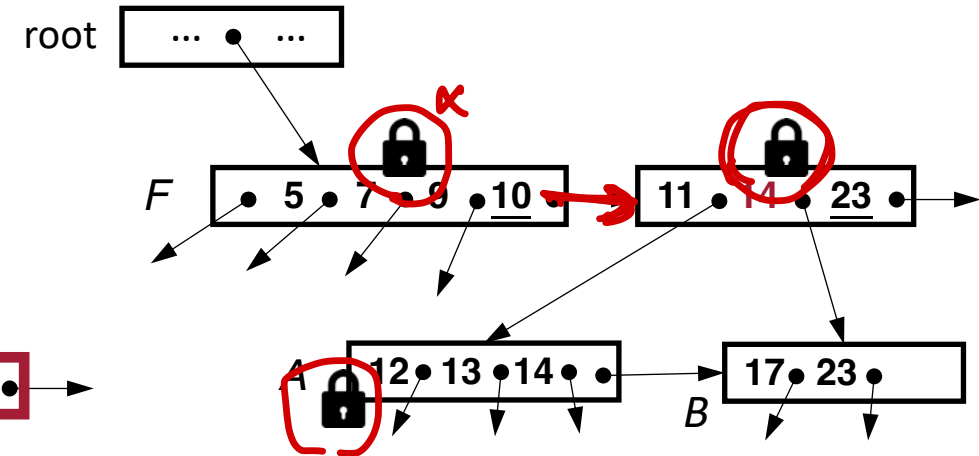
Before insert 14



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

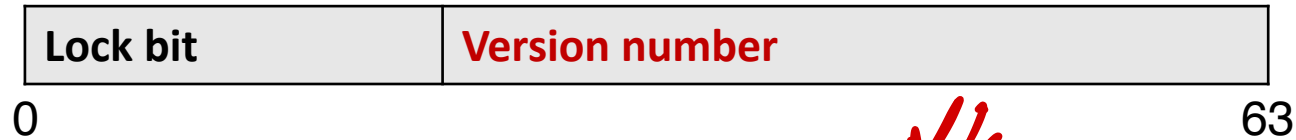
B^{link}-tree

- Search
- Insert

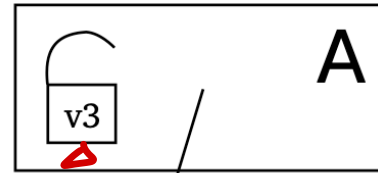
Optimistic lock coupling (OLC)

Optimistic Lock Coupling (OLC)

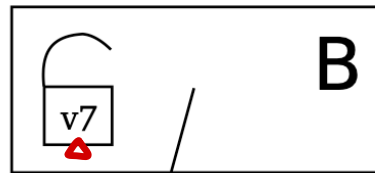
Each tuple contains a 64-bit version counter



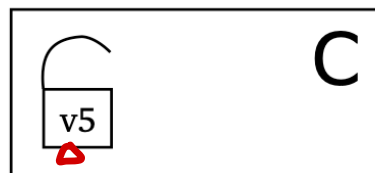
1. lock node A
2. access node A



3. lock node B
4. unlock node A
5. access node B



6. lock node C
7. unlock node B
8. access node C
9. unlock node C



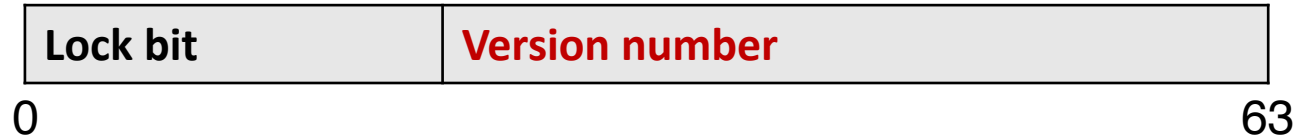
1. read version v3
2. access node A

3. read version v7
4. validate version v3
5. access node B

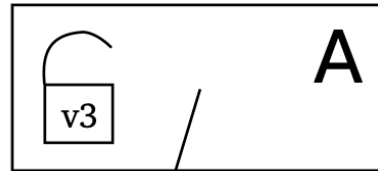
6. read version v5
7. validate version v7
8. access node C
9. validate version v5

Optimistic Lock Coupling (OLC)

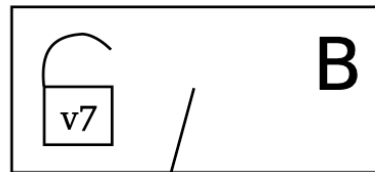
Each tuple contains a 64-bit version counter



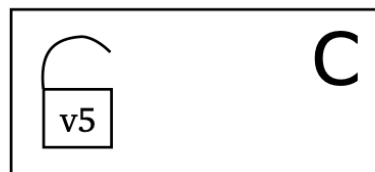
1. lock node A
2. access node A



3. lock node B
4. unlock node A
5. access node B



6. lock node C
7. unlock node B
8. access node C
9. unlock node C



1. read version v3
2. access node A
3. read version v7
4. validate version v3
5. access node B
6. read version v5
7. validate version v7
8. access node C
9. validate version v5

No scalability bottleneck

- No write to shared memory during traversal
- Upon conflict, retry from root
- Performance similar to B^{link} tree

Evaluation

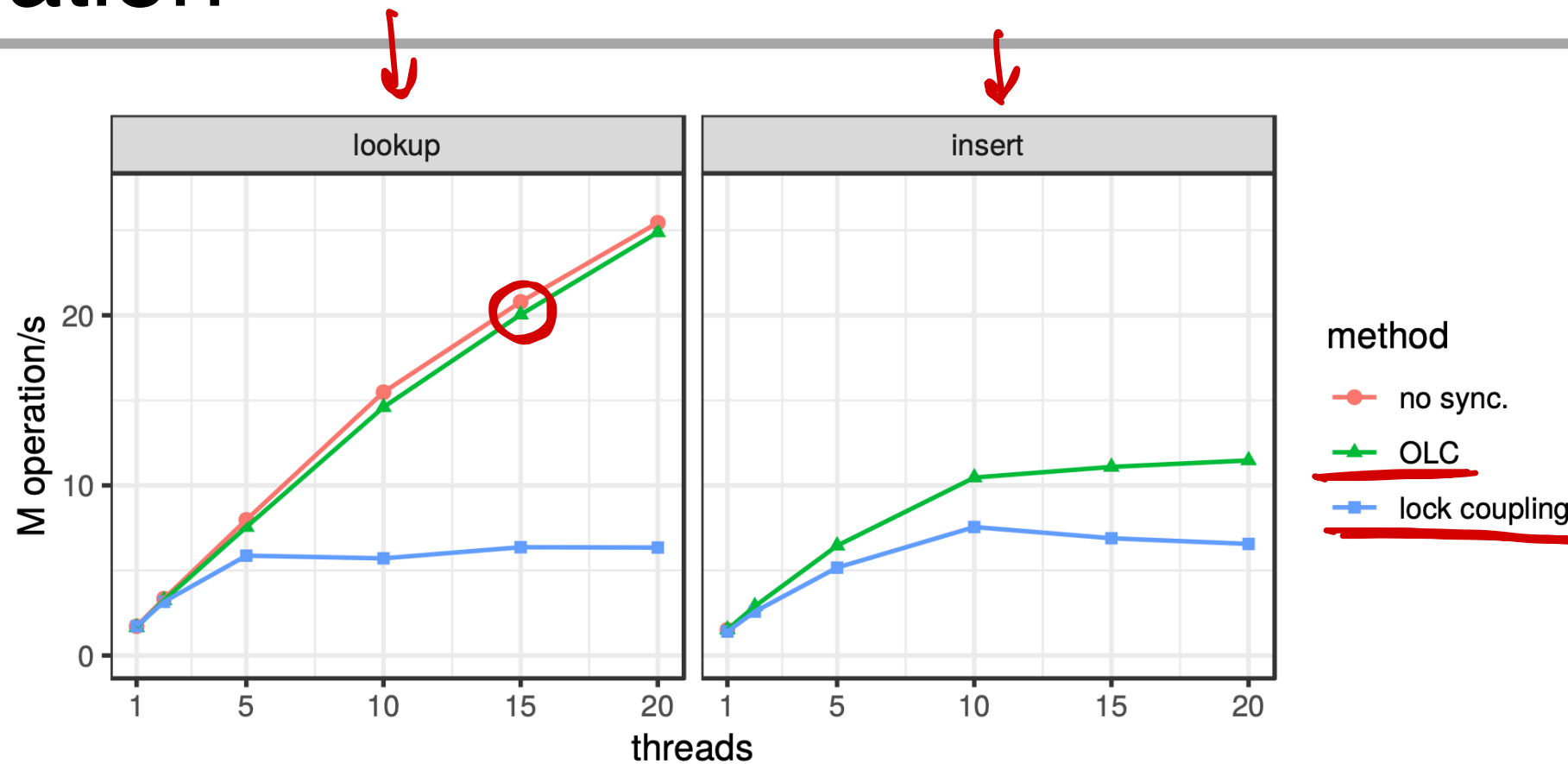


Figure 3: Scalability on 10-core system for B-tree operations (100M values).

Q/A – Blink Tree

Is Blink-tree optimization used in practice?

Is LSM tree more performance-friendly to concurrent operations?

Actual performance benchmark?

If lock blocks only writers, does a reader see inconsistent data while a writer is modifying the data?

Why only three locks are needed?

Can the same method apply to main-memory database?

Before Next Wednesday

Submit review for

- Viktor Leis, et al., [The Adaptive Radix Tree: ARTful Indexing for Main-Memory Databases](#). ICDE, 2013