

Carnegie Mellon University

Database Systems

15-445/645 SPRING 2026

ANDY PAVLO

JIGNESH PATEL

Lecture #08

Indexes + Filters:
B+Trees (aka The GOAT
Data Structure)



ADMINISTRIVIA



Project #1 is due Sunday Feb 15th @ 11:59pm

→ Recitation Video + Slides ([@64](#))

→ Perf Recitation Video + Slides ([@79](#))

→ Special OH on Saturday Feb 14th @ 3:00-5:00pm (GHC 5207)

Homework #3 is due Sunday Feb 22nd @ 11:59pm

Mid-Term Exam is on Wednesday Feb 25th

→ Lectures #01–11 (inclusive)

→ Study guide will be released early next week.

UPCOMING DATABASE TALKS

Amazon Aurora DSQL (DB Seminar)

→ Monday Feb 9th @ 4:30pm ET

→ Zoom



Amazon
Aurora DSQL

TopK (DB Seminar)

→ Monday Feb 16th @ 4:30pm ET

→ Zoom



Microsoft HorizonDB (DB Seminar)

→ Monday Feb 23rd @ 4:30pm ET

→ Zoom



Azure HorizonDB

LAST CLASS



Hash tables are important data structures that are used all throughout a DBMS.

→ Space Complexity: **$O(n)$**

→ Average Time Complexity: **$O(1)$**

Static vs. Dynamic Hashing schemes

DBMSs use mostly hash tables for their internal data structures.

INDEXES VS. FILTERS

An **index** data structure of a subset of a table's attributes that are organized and/or sorted to the location of specific tuples using those attributes.

→ Example: B+Tree

A **filter** is a data structure that answers set membership queries; it tells you whether a key (likely) exists in a set but not where it is located.

→ Example: Bloom Filter

TODAY'S AGENDA

B+Tree Overview

Design Choices

Optimizations



B-TREE FAMILY

There is a specific data structure called a **B-Tree**.

People also use the term to generally refer to a class of balanced tree data structures:

- **B-Tree** (1970)
- **B+Tree** (1973)
- **B*Tree** (1977?)
- **B^{link}-Tree** (1981)
- **B ϵ -Tree** (2003)
- **Bw-Tree** (2013)
- **Bf-Tree** (2024)

B-TREE FA

There is a specific data structure

People also use the term to gener-
alized balanced tree data structures:

- **B-Tree** (1970)
- **B+Tree** (1973)
- **B*Tree** (1977?)
- **Blink-Tree** (1981)
- **Bε-Tree** (2003)
- **Bw-Tree** (2013)
- **Bf-Tree** (2024)

107

D1-82-0989

ORGANIZATION AND MAINTENANCE OF LARGE

ORDERED INDICES

by

R. Bayer

and

E. McCreight

Mathematical and Information Sciences Report No. 20

Mathematical and Information Sciences Laboratory

BOEING SCIENTIFIC RESEARCH LABORATORIES

July 1970

B-TREE FA

There is a specific data structure

People also use the term to gener-
balanced tree data structures:

- **B-Tree** (1970)
- **B+Tree** (1973)
- **B*Tree** (1977?)
- **Blink-Tree** (1981)
- **B ϵ -Tree** (2003)
- **Bw-Tree** (2013)
- **Bf-Tree** (2024)

The Ubiquitous B-Tree

DOUGLAS COMER

Computer Science Department, Purdue University, West Lafayette, Indiana 47907

B-trees have become, de facto, a standard for file organization. File indexes of users, dedicated database systems, and general-purpose access methods have all been proposed and implemented using B-trees. This paper reviews B-trees and shows why they have been so successful. It discusses the major variations of the B-tree, especially the B*-tree, contrasting the relative merits and costs of each implementation. It illustrates a general purpose access method which uses a B-tree.

Keywords and Phrases: B-tree, B*-tree, B ϵ -tree, file organization, index

CR Categories: 3.73 3.74 4.33 4.34

INTRODUCTION

The secondary storage facilities available on large computer systems allow users to store, update, and recall data from large collections of information called files. A computer must retrieve an item and place it in main memory before it can be processed. In order to make good use of the computer resources, one must organize files intelligently, making the retrieval process efficient.

The choice of a good file organization depends on the kinds of retrieval to be performed. There are two broad classes of retrieval commands which can be illustrated by the following examples:

Sequential: "From our employee file, prepare a list of all employees' names and addresses," and

Random: "From our employee file, extract the information about employee J. Smith".

We can imagine a filing cabinet with three drawers of folders, one folder for each employee. The drawers might be labeled "A-G," "H-R," and "S-Z," while the folders

might be labeled with the employees' last names. A sequential request requires the searcher to examine the entire file, one folder at a time. On the other hand, a random request implies that the searcher, guided by the labels on the drawers and folders, need only extract one folder.

Associated with a large, randomly accessed file in a computer system is an *index* which, like the labels on the drawers and folders of the file cabinet, speeds retrieval by directing the searcher to the small part of the file containing the desired item. Figure 1 depicts a file and its index. An index may be physically integrated with the file, like the labels on employee folders, or physically separate, like the labels on the drawers. Usually the index itself is a file. If the index file is large, another index may be built on top of it to speed retrieval further, and so on. The resulting hierarchy is similar to the employee file, where the topmost index consists of labels on drawers, and the next level of index consists of labels on folders.

Natural hierarchies, like the one formed by considering last names as index entries, do not always produce the best perform-

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its copy otherwise, or to republish, requires a fee and/or specific permission.
© 1979 ACM 0010-4892/79/0600-0121 \$00.75

B-TREE FAMILY

There is a specific data structure called a

People also use the term to generally refer to balanced tree data structures:

- **B-Tree** (1970)
- **B+Tree** (1973)
- **B*Tree** (1977?)
- **Blink-Tree** (1981)
- **B ϵ -Tree** (2003)
- **Bw-Tree** (2013)
- **Bf-Tree** (2024)

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 paper, we consider a simple variant of the B-tree (actually of the B*-tree, system proposed by Wedekind [15]) especially well suited for use in a concurrent database system.

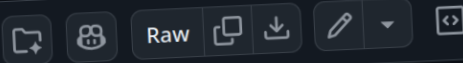
Methods for concurrent operations on B*-trees have been discussed by Bayer and Schkolnick [3] and others [6, 12, 13]. The solution given in the current paper

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

This research was supported by the National Science Foundation under Grant MCS76-16694. Authors' present addresses: P. L. Lehman, Department of Computer Science, Carnegie-Mellon University, Pittsburgh, PA 15213; S. B. Yao, Department of Computer Science and College of Business and Management, University of Maryland, College Park, MD 20742.

© 1981 ACM 0362-5915/81/1200-0650 \$00.75
ACM Transactions on Database Systems, Vol. 6, No. 4, December 1981, Pages 650-670.

Code Blame



1 src/backend/access/nbtree/README

2
3 Btree Indexing

4 =====

5
6 This directory contains a correct implementation of **Lehman and Yao's**
7 high-concurrency B-tree management algorithm (P. Lehman and S. Yao,
8 Efficient Locking for Concurrent Operations on B-Trees, ACM Transactions
9 on Database Systems, Vol 6, No. 4, December 1981, pp 650-670). We also
10 use a simplified version of the deletion logic described in Lanin and
11 Shasha (V. Lanin and D. Shasha, A Symmetric Concurrent B-Tree Algorithm,
12 Proceedings of 1986 Fall Joint Computer Conference, pp 380-389).

13
14 The basic Lehman & Yao Algorithm

15 -----

16
17 Compared to a classic B-tree, L&Y adds a right-link pointer to each page,
18 to the page's right sibling. It also adds a "high key" to each page, which
19 is an upper bound on the keys that are allowed on that page. These two
20 additions make it possible to detect a concurrent page split, which allows
21 the tree to be searched without holding any read locks (except to keep a
22 single page from being modified while reading it).

Efficient Locking for Concurrent Operations on B-Trees

LIP L. LEHMAN
Carnegie-Mellon University

SHANG YAO
Carnegie Mellon University

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 inherent difficulty of concurrent operations on such structures, using a storage model. A single additional "link" pointer in each node allows a process to easily perform 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 the number of nodes are locked by any update process at any given time. An correctness proof for our system is given.
Key Words and Phrases: database, data structures, B-tree, index organizations, concurrent algorithms, synchronization, control, locking protocols, correctness, consistency, multiway search trees
Categories: 3.73, 3.74, 4.32, 4.33, 4.34, 5.24

INTRODUCTION

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. The guaranteed small (average) search, insertion, and deletion time of B-trees makes them quite appealing for database applications. One of the areas of current interest in database design is the construction of databases that can be manipulated concurrently and correctly by several processes. In this paper we consider a simple variant of the B-tree (actually of the B*-tree, or B+ tree, as Wedekind [15] especially well suited for use in a concurrent database system). For concurrent operations on B*-trees have been discussed by Bayer and Rieckert [3] and others [6, 12, 13]. The solution given in the current paper

This work was supported by the National Science Foundation under Grant MCS76-16694. The authors would like to thank the ACM for permission to publish this work. Any copying without fee all or part of this material is granted provided that the copies are not made for direct commercial advantage, the ACM copyright notice and the title of the publication appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

Address: P. L. Lehman, Department of Computer Science, Carnegie-Mellon University, Pittsburgh, PA 15213; S. B. Yao, Department of Computer Science and College of Business Administration, University of Maryland, College Park, MD 20742.

ACM Transactions on Database Systems, Vol. 6, No. 4, December 1981, Pages 650-670.

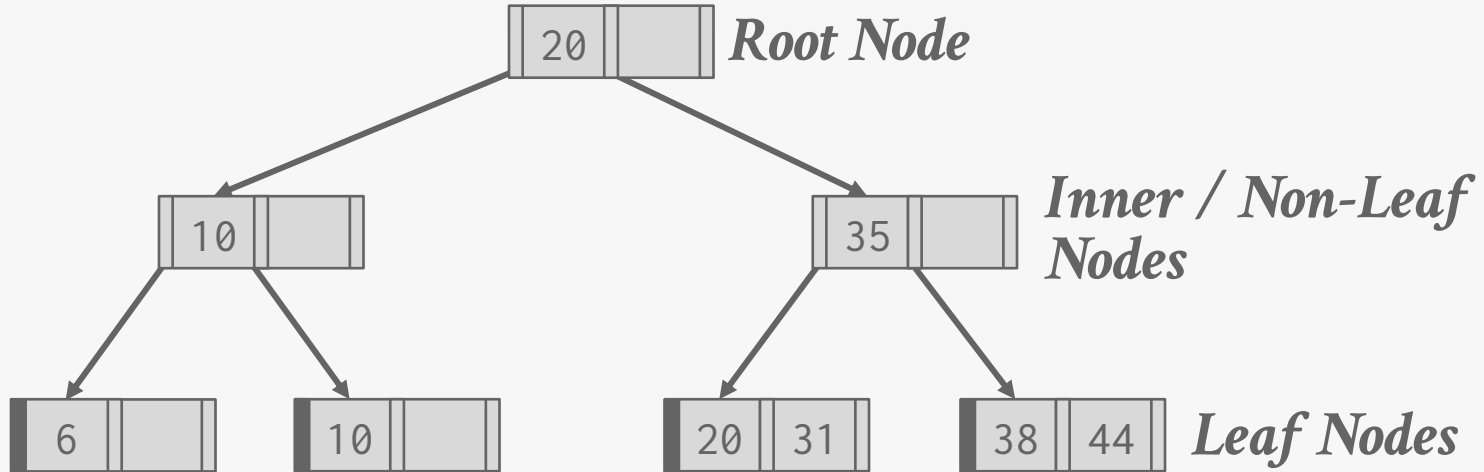
B+TREE

A **B+Tree** is a self-balancing, ordered m -way tree for searches, sequential access, insertions, and deletions in $O(\log_m n)$ where m is the tree fanout.

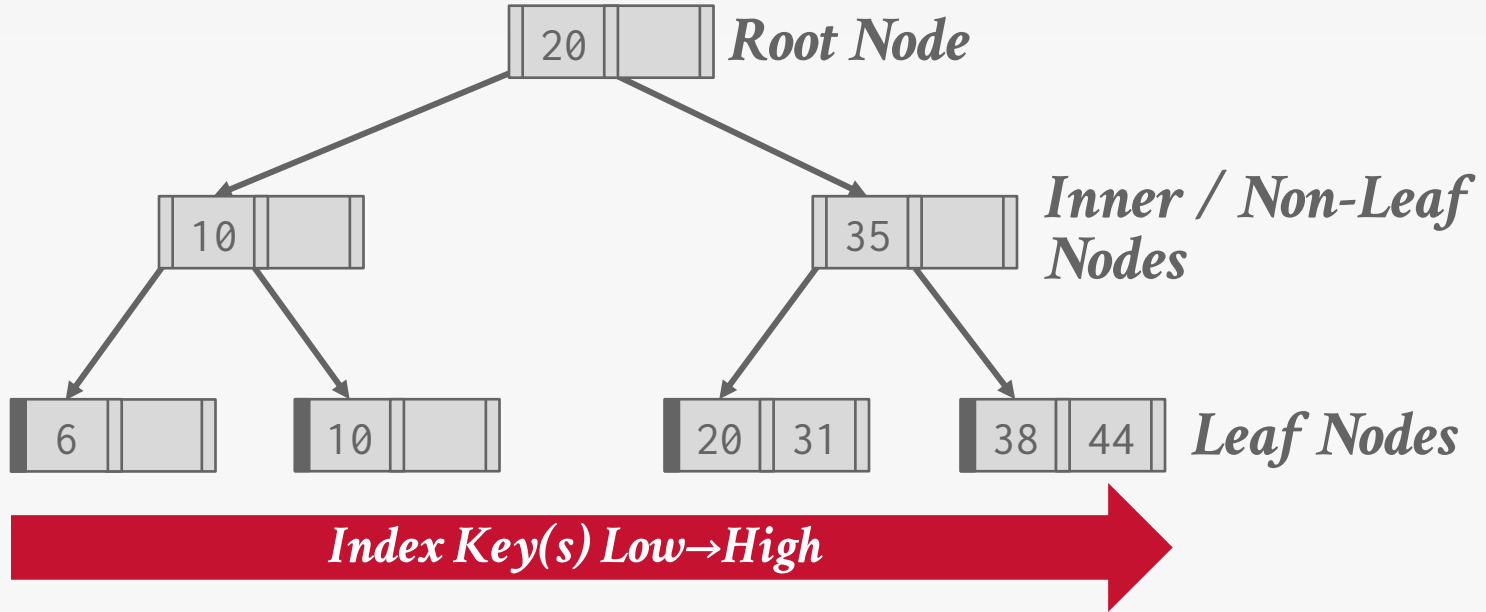
- It is perfectly balanced (i.e., every leaf node is at the same depth in the tree)
- Every node other than the root is at least half-full
 $m/2 - 1 \leq \#keys \leq m - 1$
- Every inner node with k keys has $k + 1$ non-null children.
- Optimized for reading/writing large data blocks.

Some real-world implementations relax these properties, but we will ignore that for now...

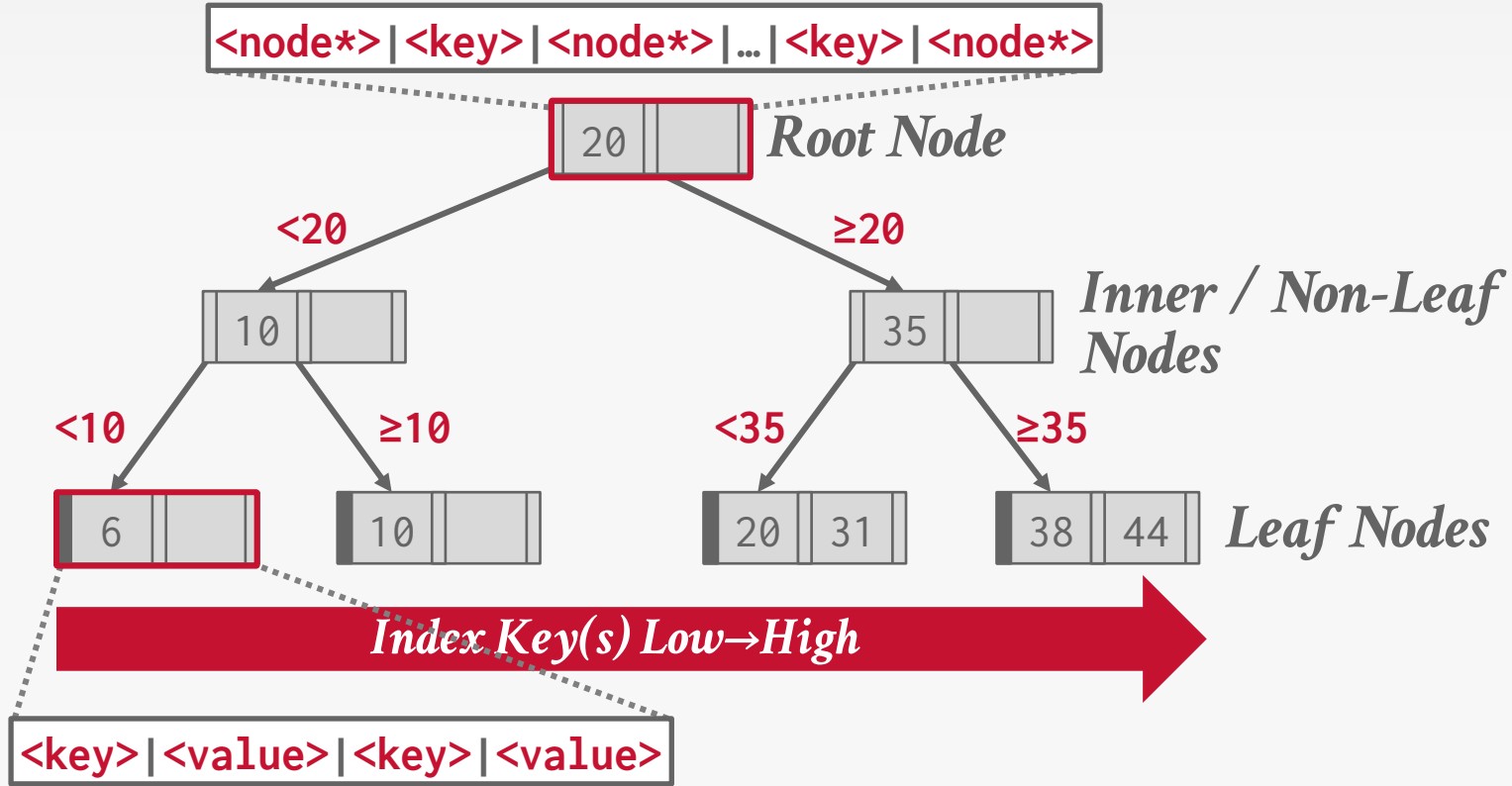
B+TREE EXAMPLE



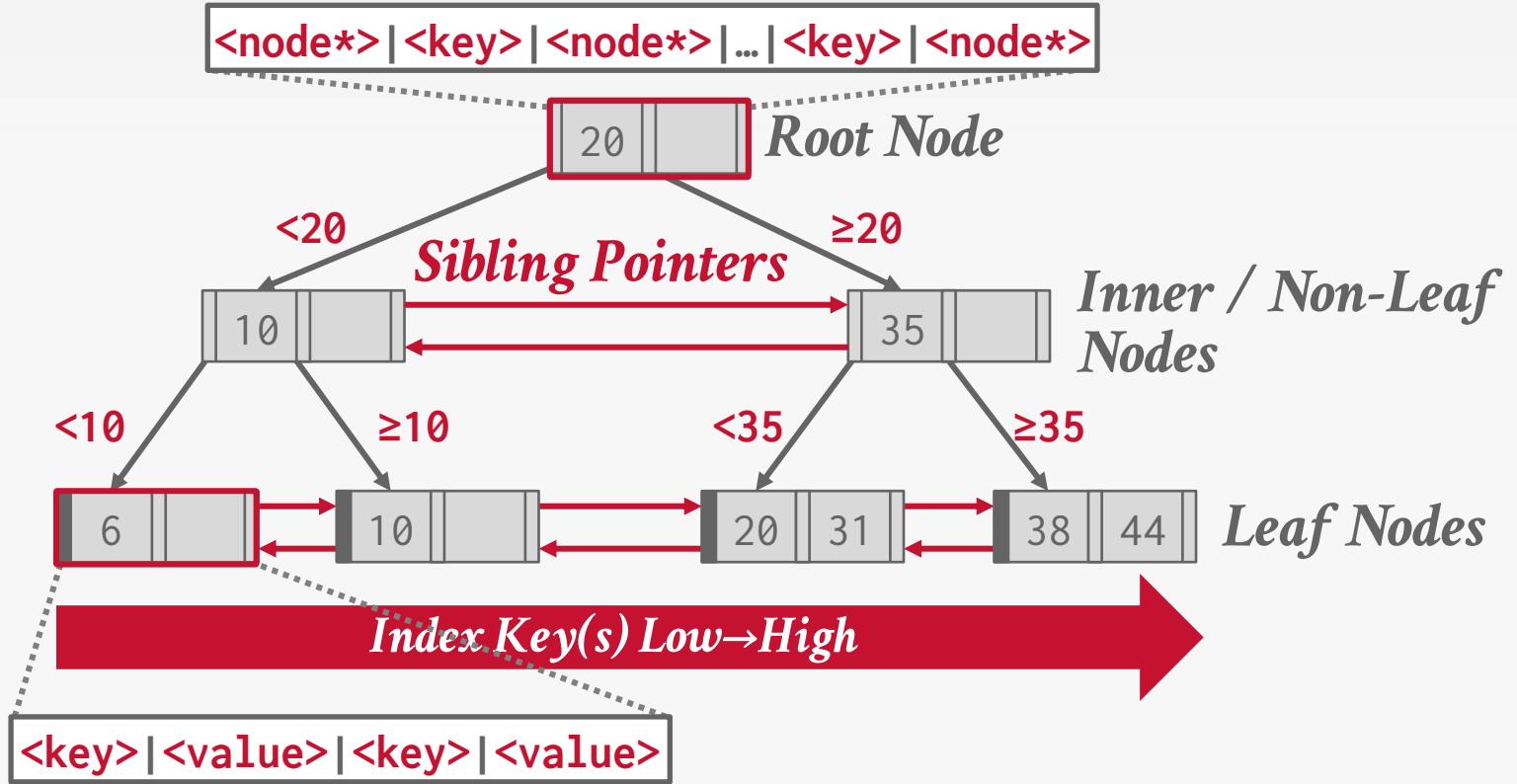
B+TREE EXAMPLE



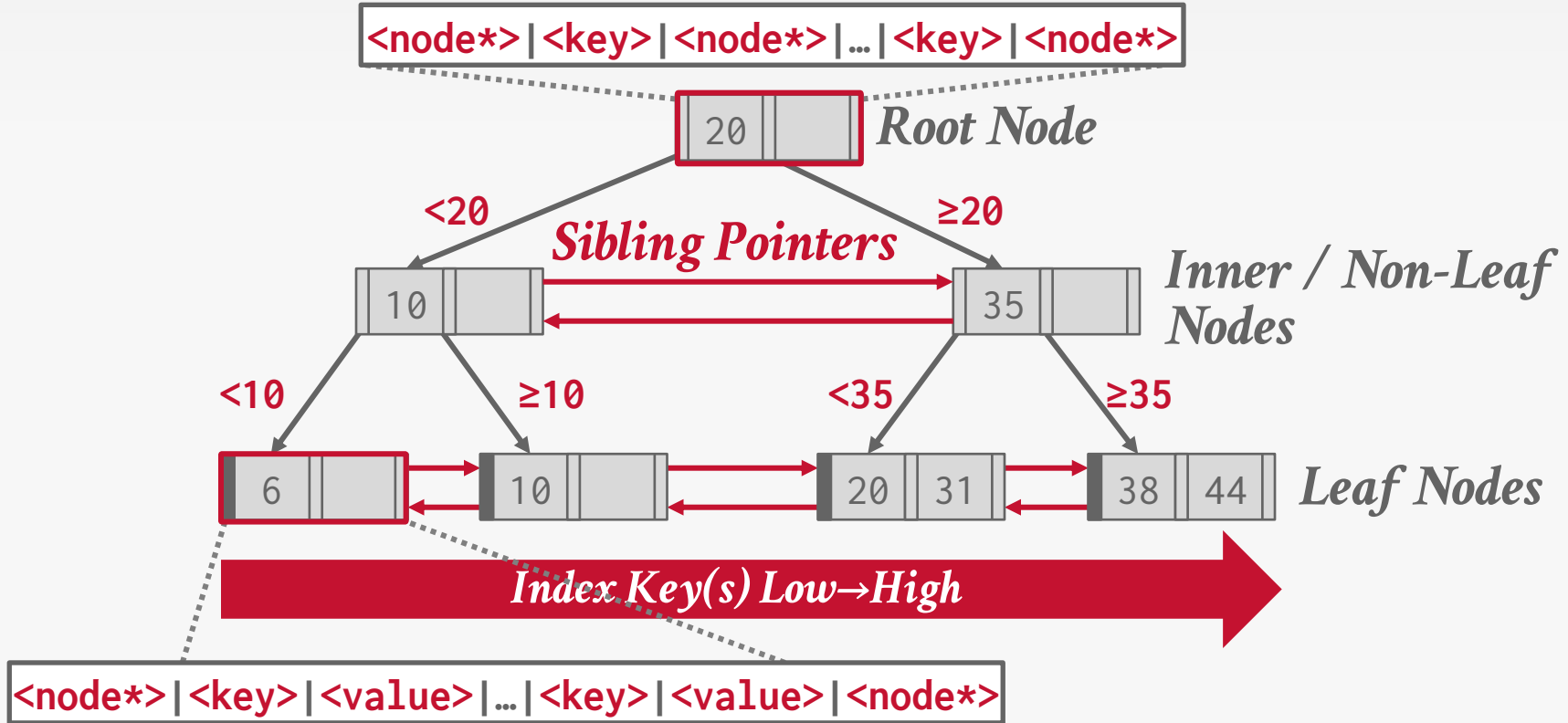
B+TREE EXAMPLE



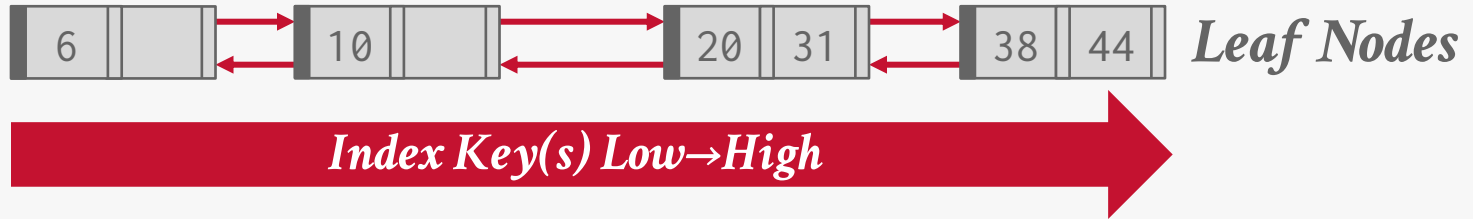
B+TREE EXAMPLE



B+TREE EXAMPLE



B+TREE EXAMPLE



NODES

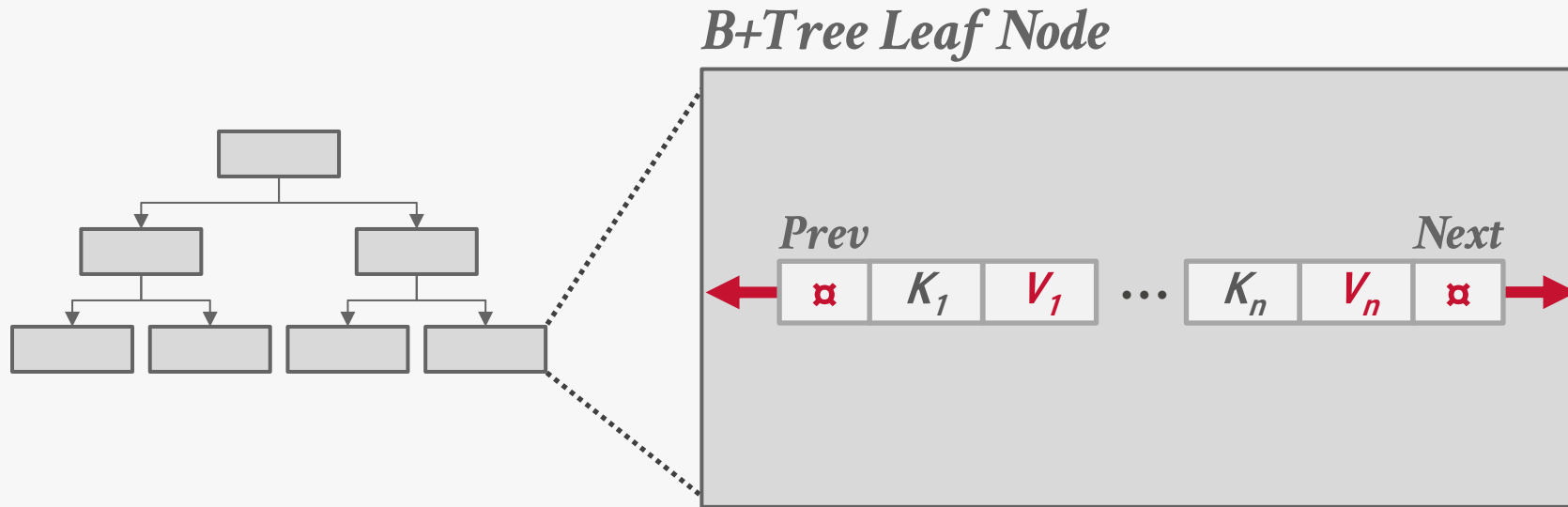
Every B+Tree node is comprised of an array of key/value pairs.

- The keys are derived from the index's target attribute(s).
- The values will differ based on whether the node is classified as an inner node or a leaf node.

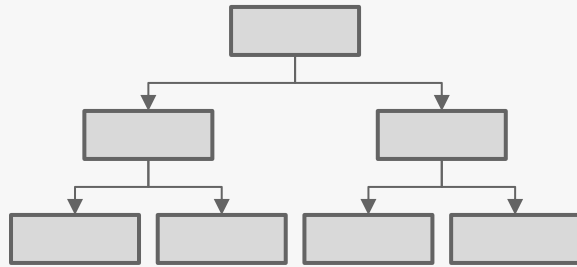
The arrays are (usually) kept in sorted key order.

Store all **NULL** keys at either first or last leaf nodes.

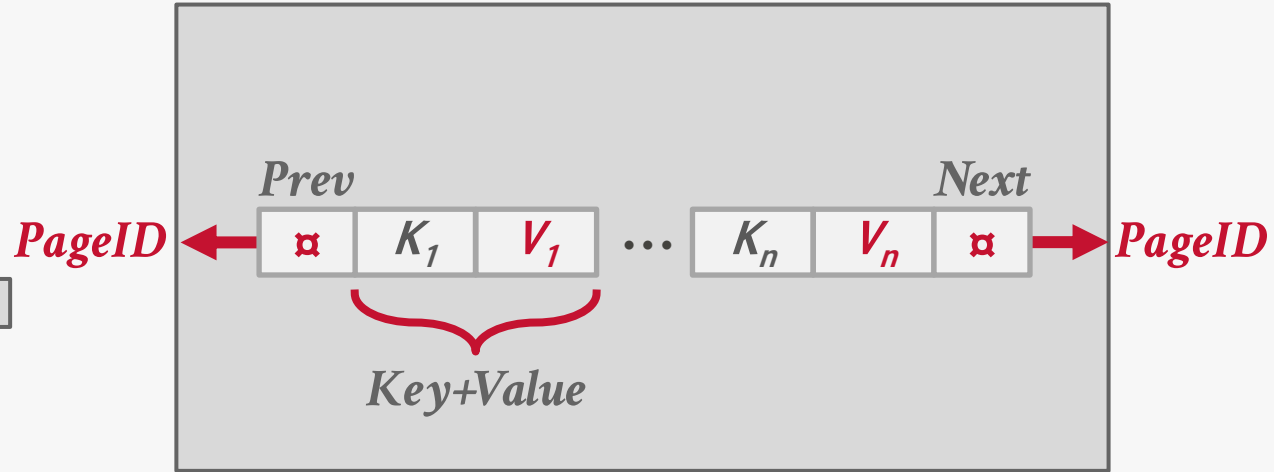
B+TREE LEAF NODES



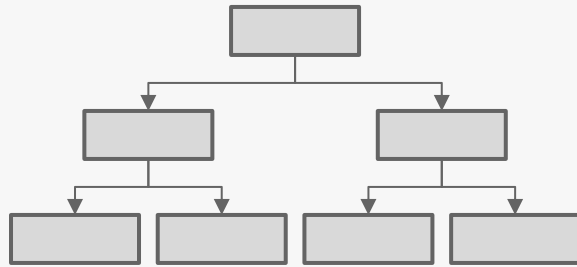
B+TREE LEAF NODES



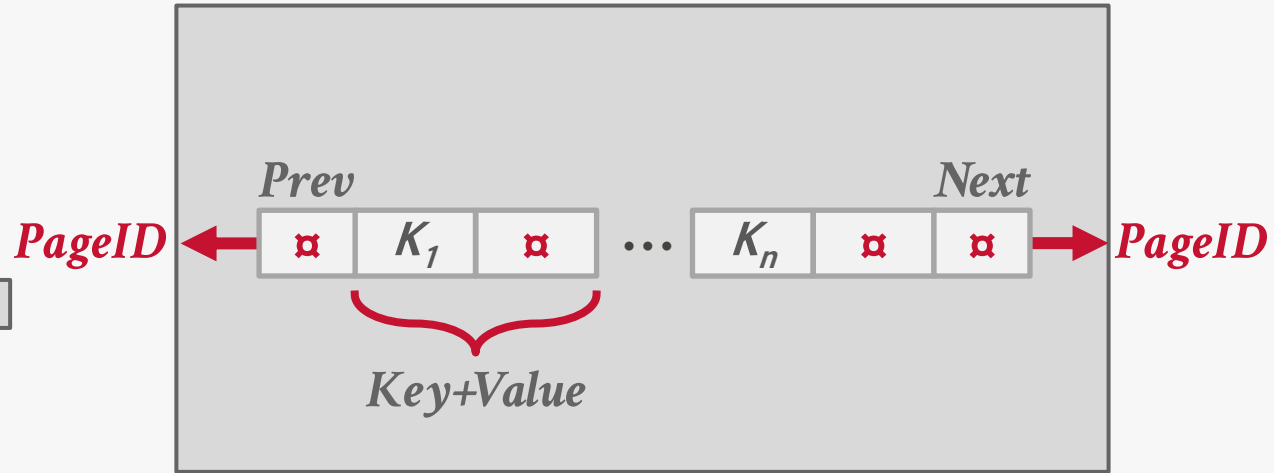
B+Tree Leaf Node



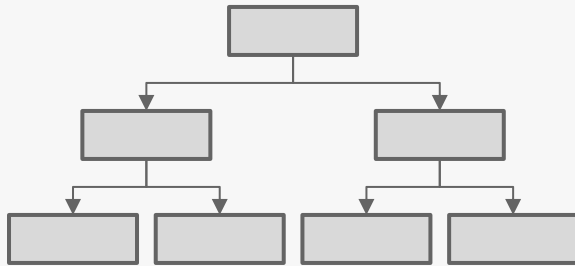
B+TREE LEAF NODES



B+Tree Leaf Node



B+TREE LEAF NODES



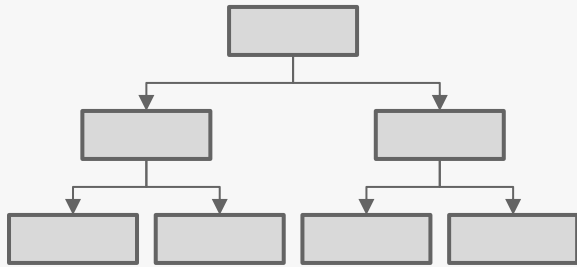
B+Tree Leaf Node

<i>Level</i>	<i>Slots</i>	<i>Prev</i>	<i>Next</i>	<i>High Key</i>
#	#	☐	☐	K_x

Sorted Key/Value Pairs

K_1	☐	K_2	☐	K_3	☐
K_4	☐	K_5	☐	...	

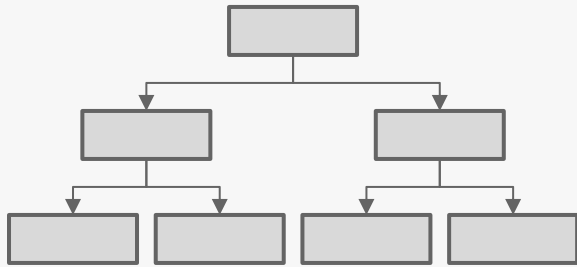
B+TREE LEAF NODES



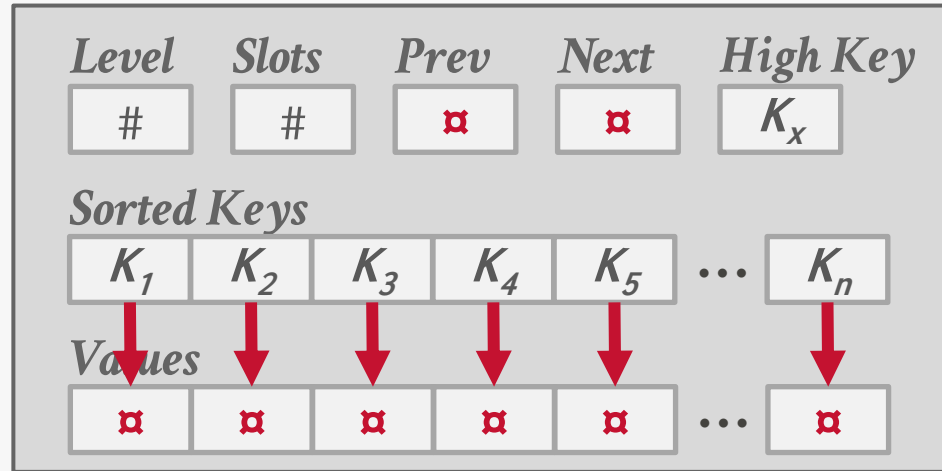
B+Tree Leaf Node



B+TREE LEAF NODES



B+Tree Leaf Node



LEAF NODE VALUES

Approach #1: Record IDs

- A pointer to the location of the tuple to which the index entry corresponds.
- Most common implementation.



Approach #2: Tuple Data

- Index-Organized Storage (Lecture #05)
- Primary Key Index: Leaf nodes store the contents of the tuple.
- Secondary Indexes: Leaf nodes store tuples' primary key as their values.



B-TREE VS. B+TREE

The original **B-Tree** from 1970 stored keys and values in all nodes in the tree.

- More space-efficient, since each key only appears once in the tree.
- Harder to support concurrent access.

The **B+Tree** only stores values in leaf nodes. Inner nodes only guide the search process.

- A deleted key must be removed from the leaf node, but it may still exist in inner nodes.

B+TREE: INSERT

Find correct leaf node **L**.

Insert data entry into **L** in sorted order.

If **L** has enough space, done!

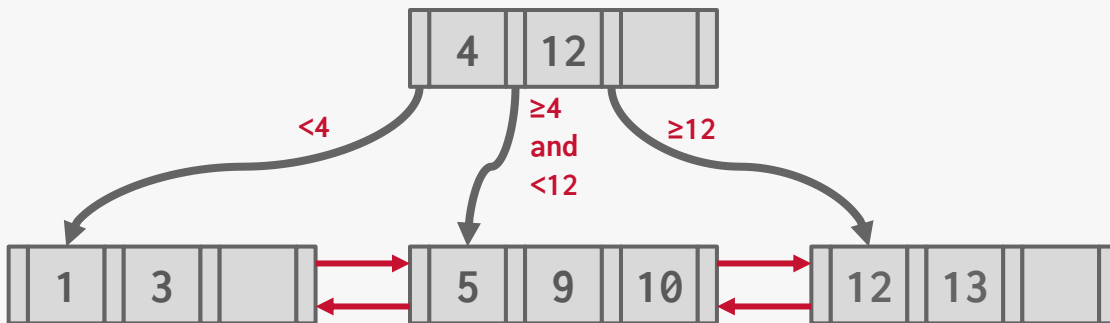
Otherwise, split **L** keys into **L** and a new node **L₂**

→ Redistribute entries evenly, copy up middle key.

→ Insert index entry pointing to **L₂** into parent of **L**.

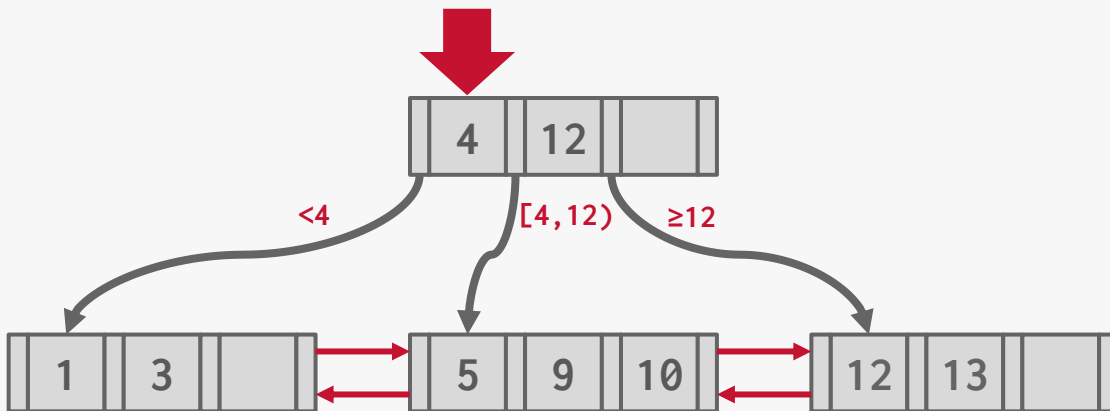
To split inner node, redistribute entries evenly, but push up middle key.

B+TREE: INSERT EXAMPLE (1)



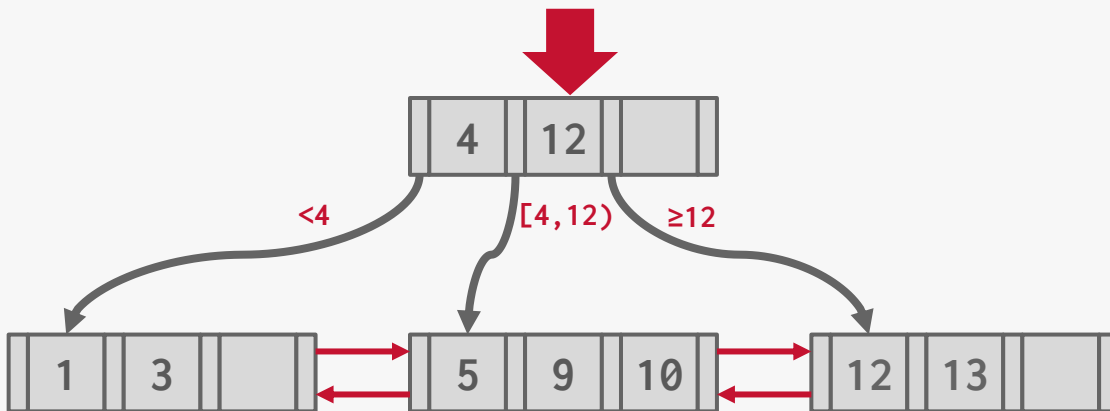
B+TREE: INSERT EXAMPLE (1)

Insert 6



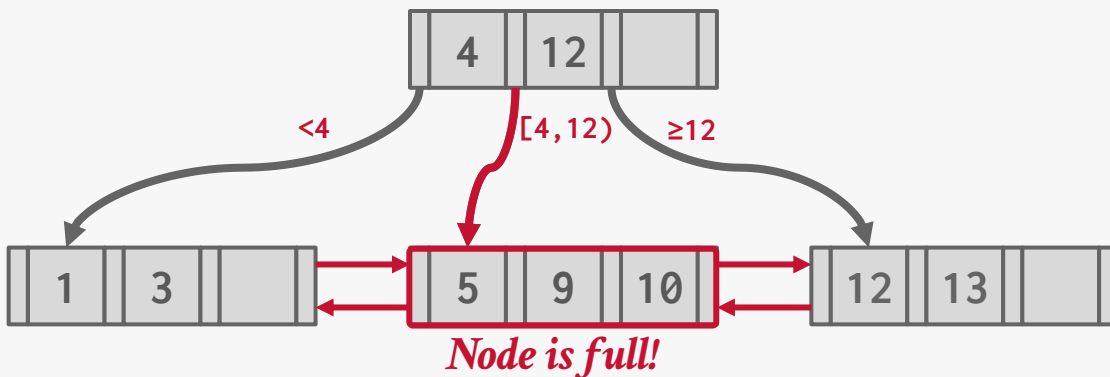
B+TREE: INSERT EXAMPLE (1)

Insert 6



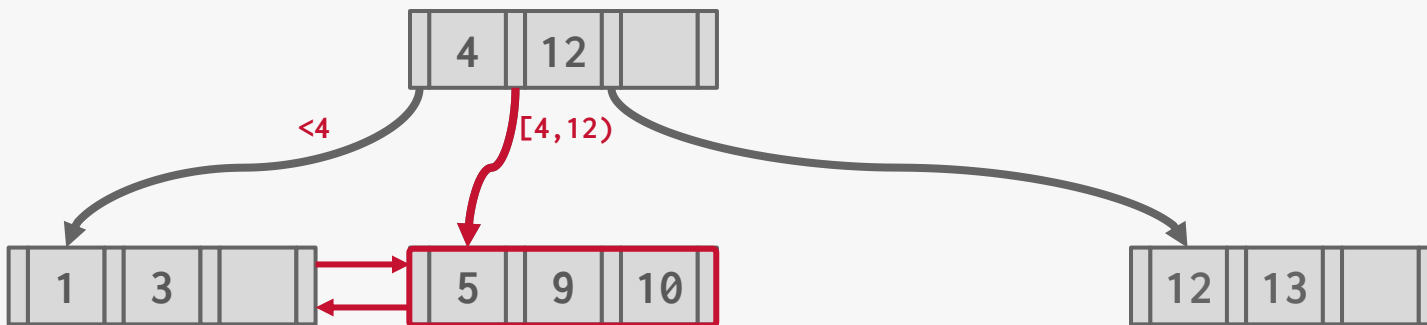
B+TREE: INSERT EXAMPLE (1)

Insert 6



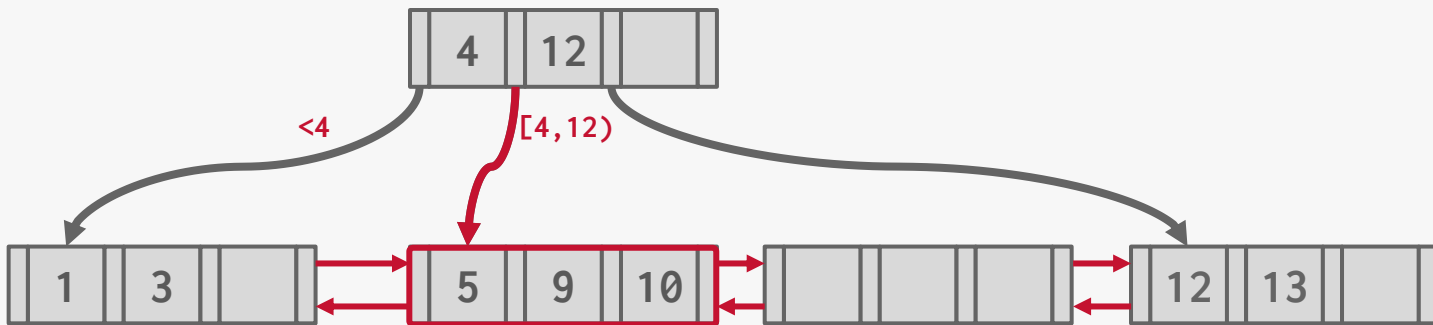
B+TREE: INSERT EXAMPLE (1)

Insert 6



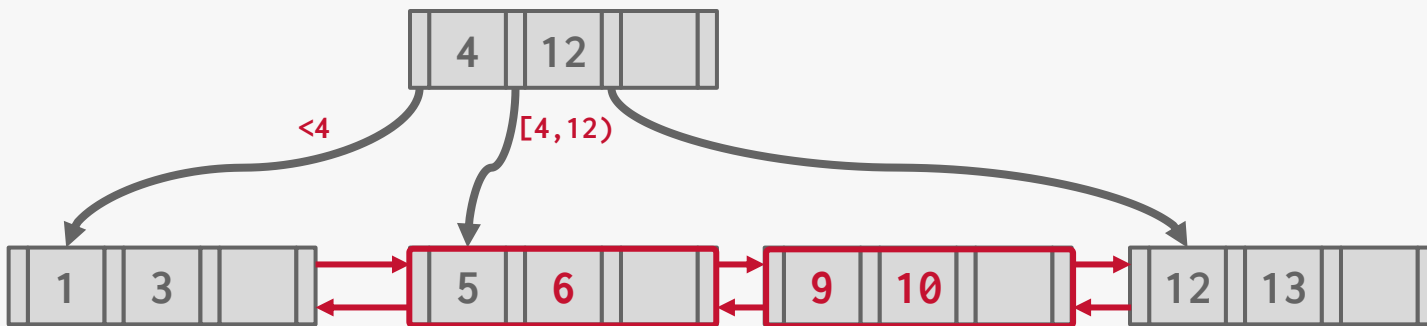
B+TREE: INSERT EXAMPLE (1)

Insert 6



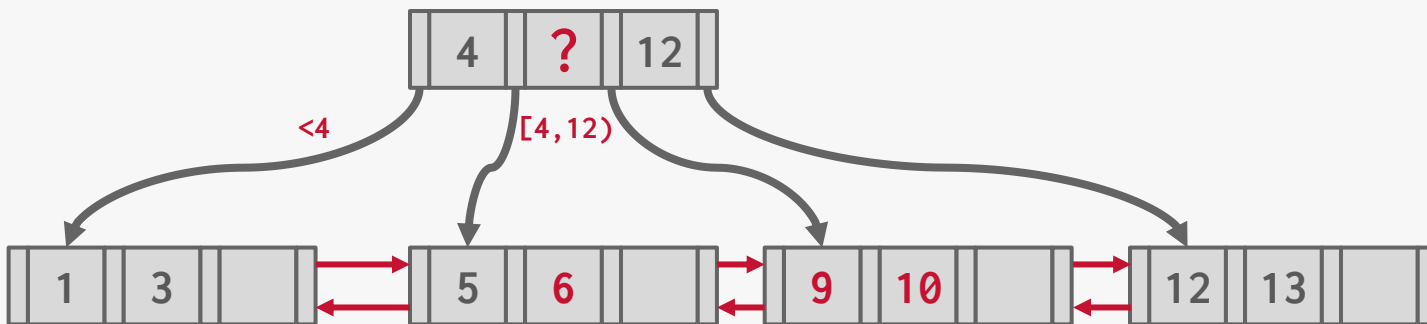
B+TREE: INSERT EXAMPLE (1)

Insert 6



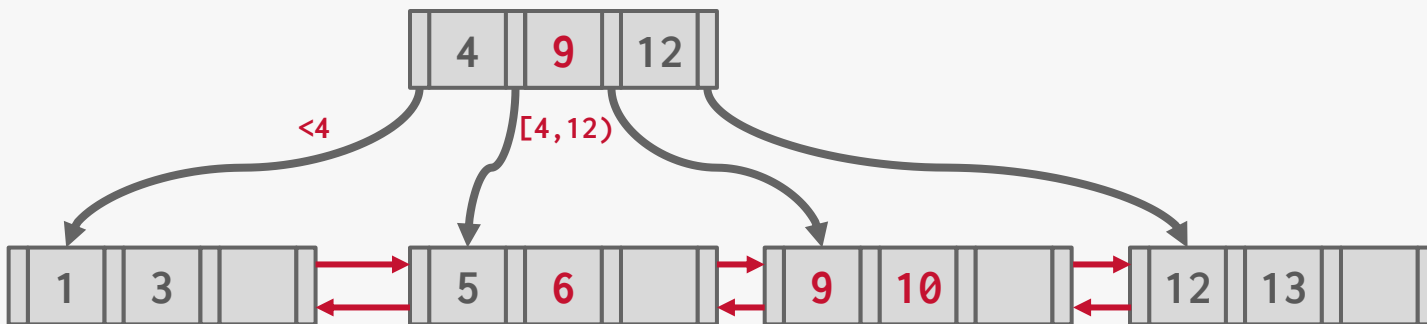
B+TREE: INSERT EXAMPLE (1)

Insert 6



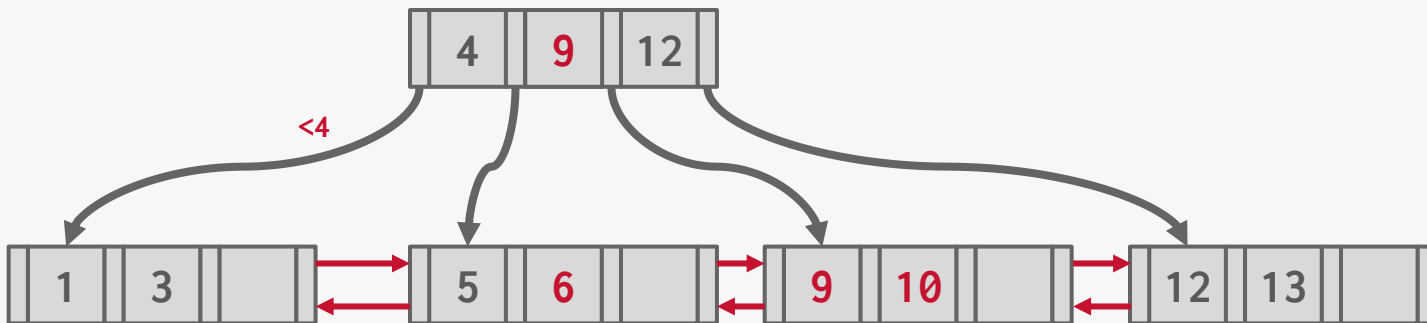
B+TREE: INSERT EXAMPLE (1)

Insert 6



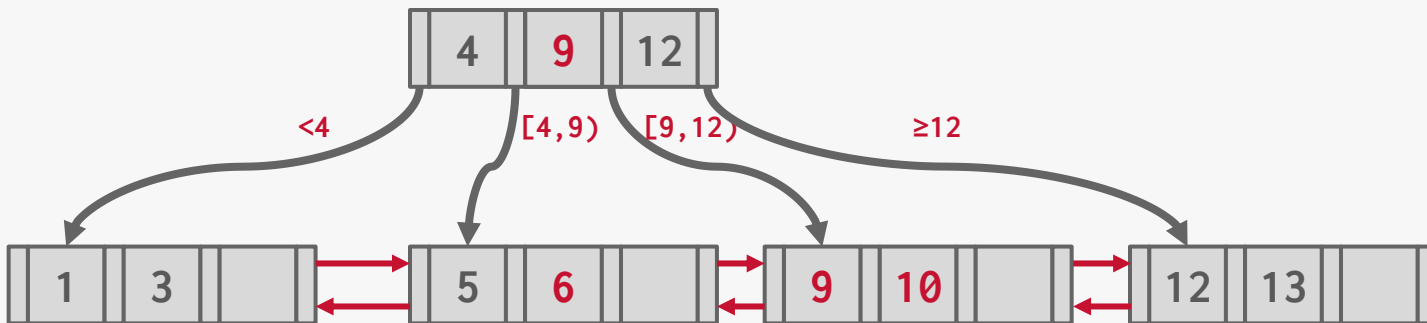
B+TREE: INSERT EXAMPLE (1)

Insert 6



B+TREE: INSERT EXAMPLE (1)

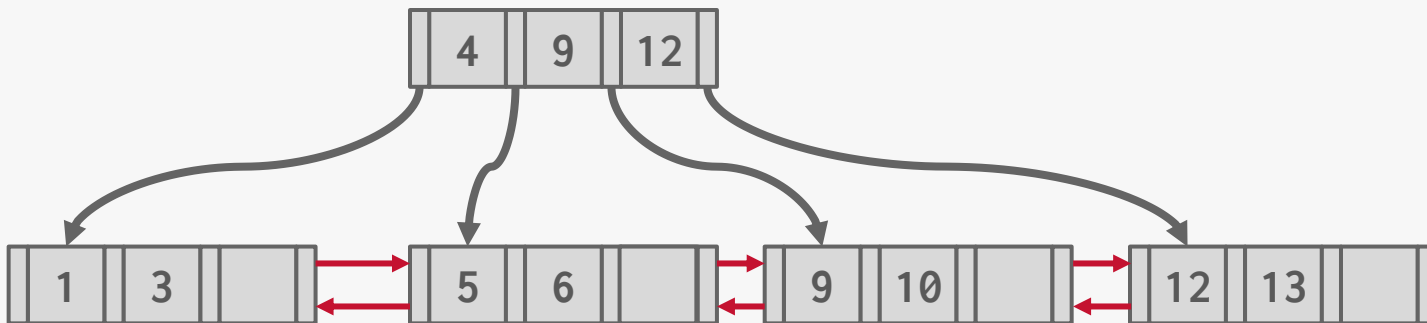
Insert 6



B+TREE: INSERT EXAMPLE (1)

Insert 6

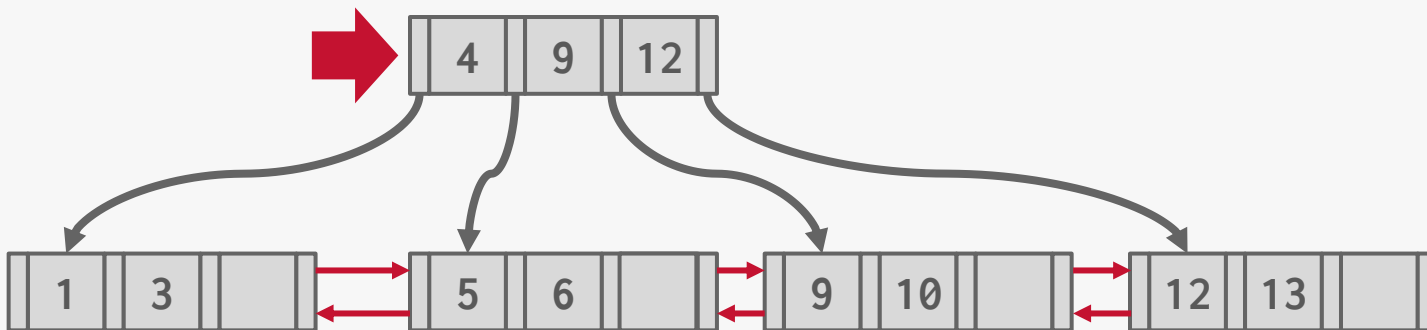
Insert 8



B+TREE: INSERT EXAMPLE (1)

Insert 6

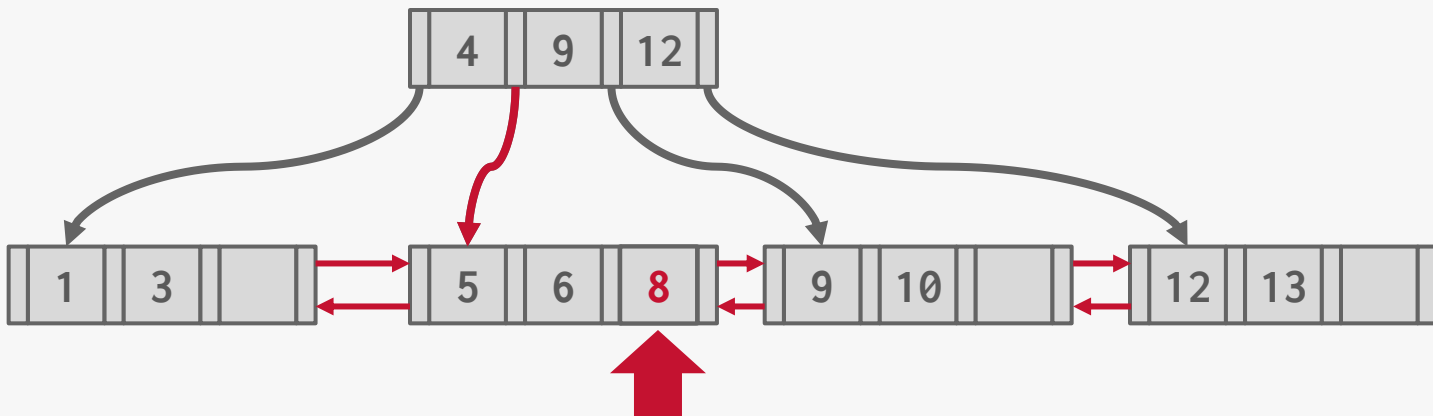
Insert 8



B+TREE: INSERT EXAMPLE (1)

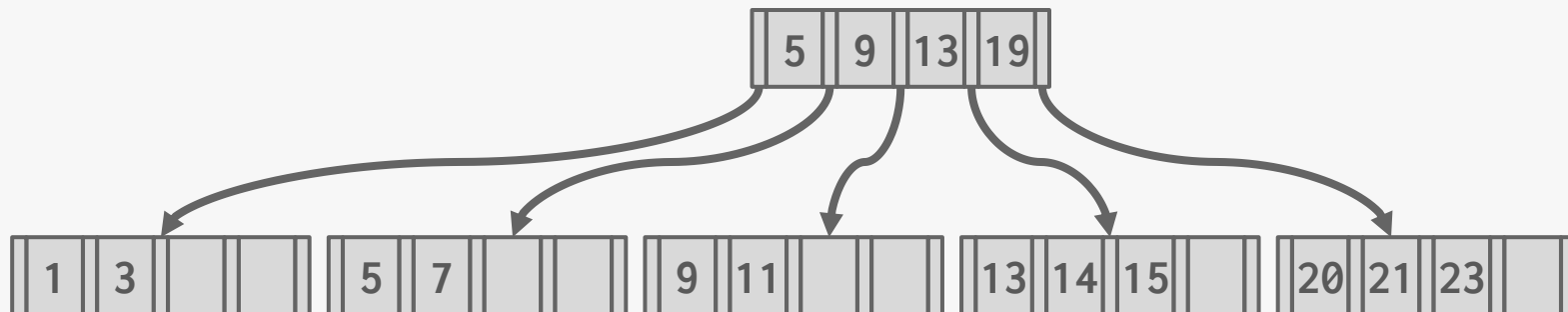
Insert 6

Insert 8



B+TREE: INSERT EXAMPLE (2)

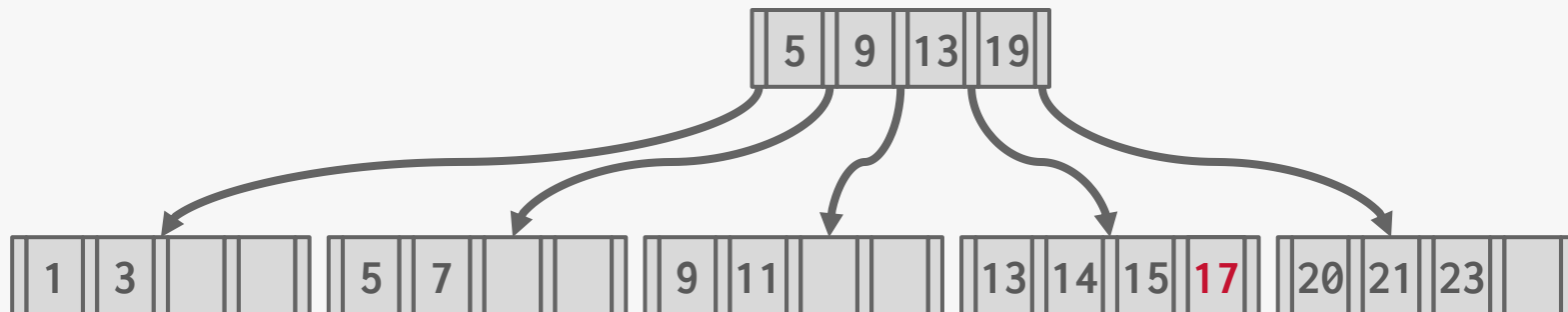
Insert 17



Note: New Example/Tree.

B+TREE: INSERT EXAMPLE (2)

Insert 17

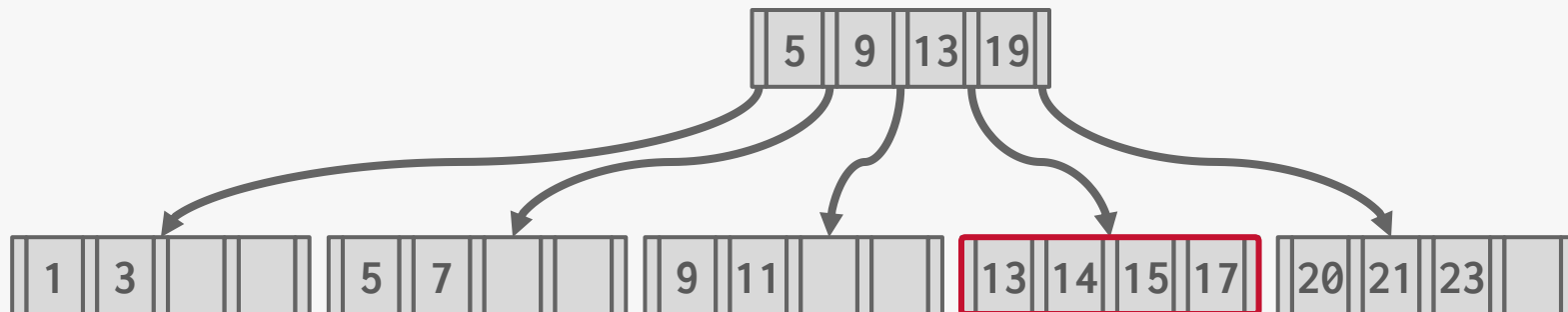


Note: New Example/Tree.

B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16

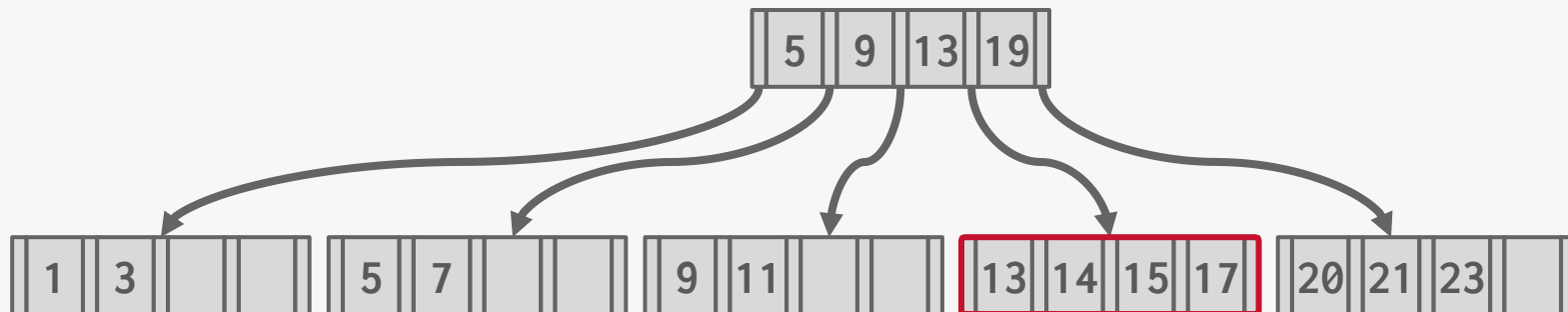


*No space in the node where
the new key "belongs".*

B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16

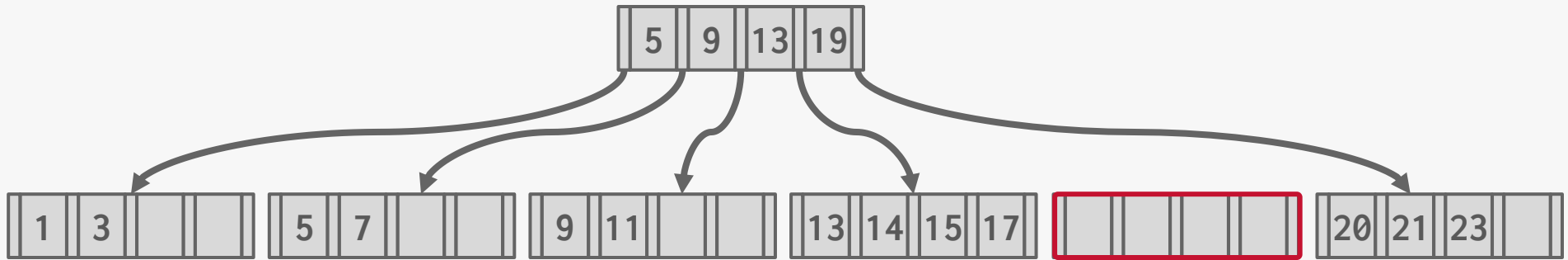


*Split the node!
Copy the middle key.
Push the key up.*

B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16

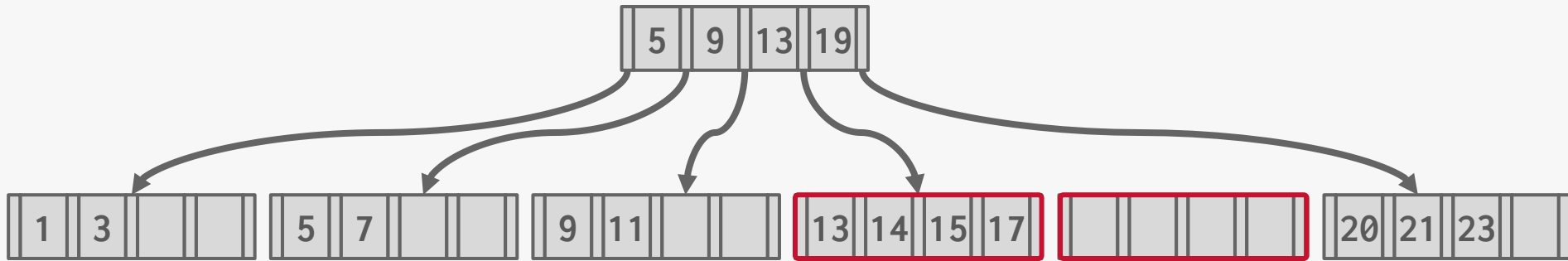


New Node!
Shuffle keys from the node
that triggered the split.

B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16

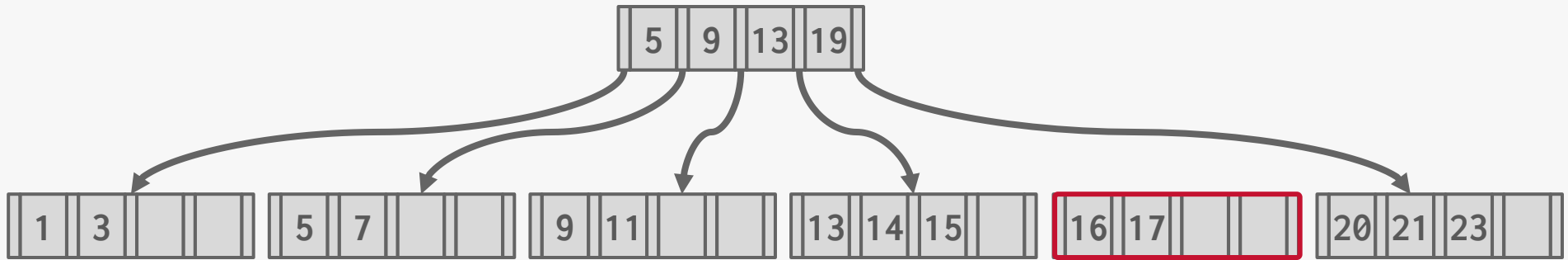


New Node!
Shuffle keys from the node that triggered the split.

B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16



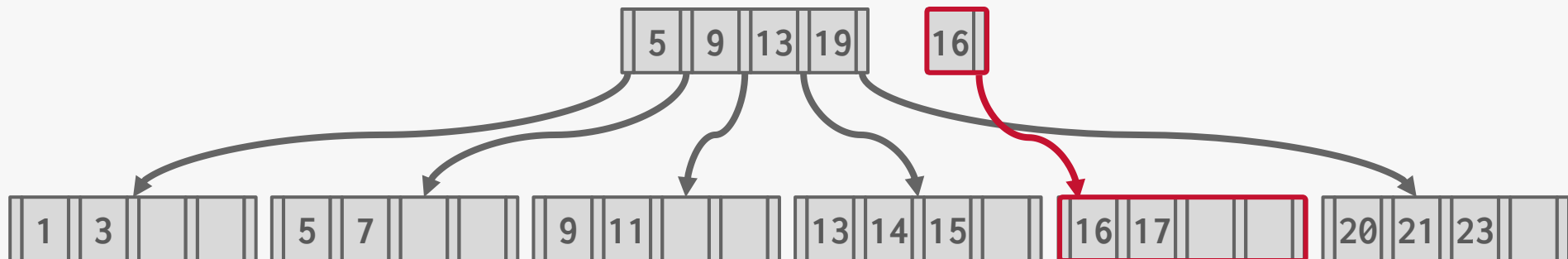
*But this is an "orphan" node!
No parent node points to it.*

B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16

Want to create a key, pointer pair like this. But cannot insert it in the root node, which is full.



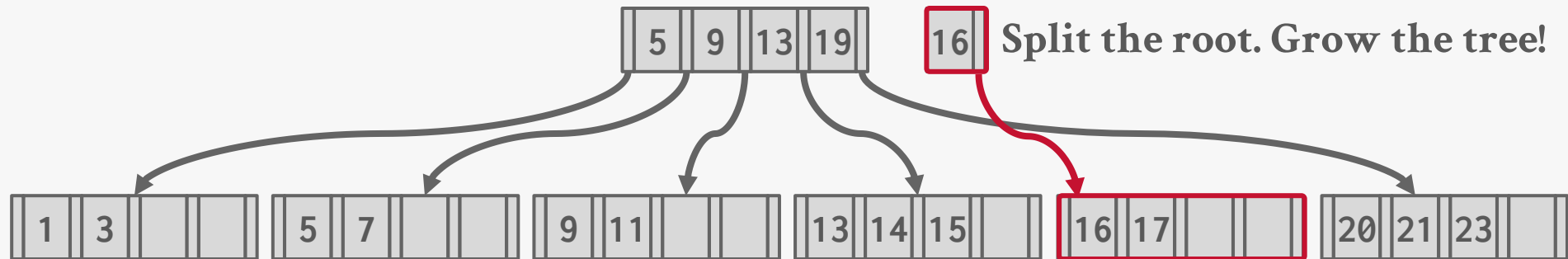
*But this is an "orphan" node!
No parent node points to it.*

B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16

Want to create a key, pointer pair like this. But cannot insert it in the root node, which is full.

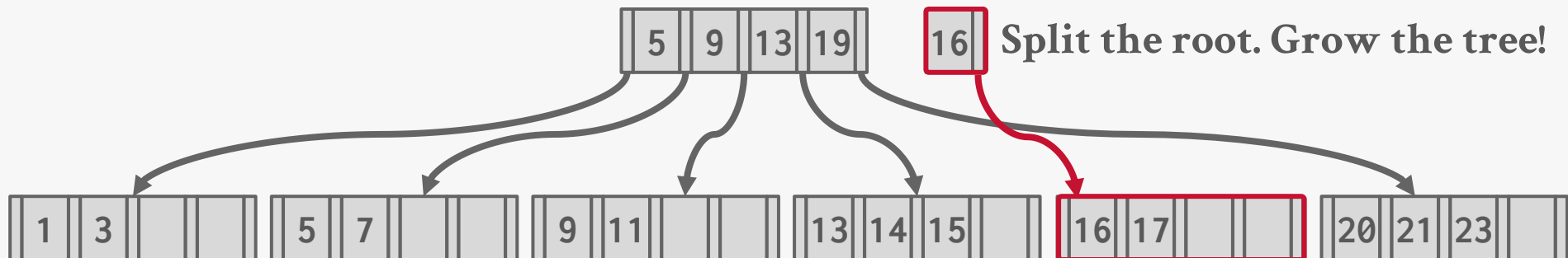


*But this is an "orphan" node!
No parent node points to it.*

B+TREE: INSERT EXAMPLE (3)

Insert 17

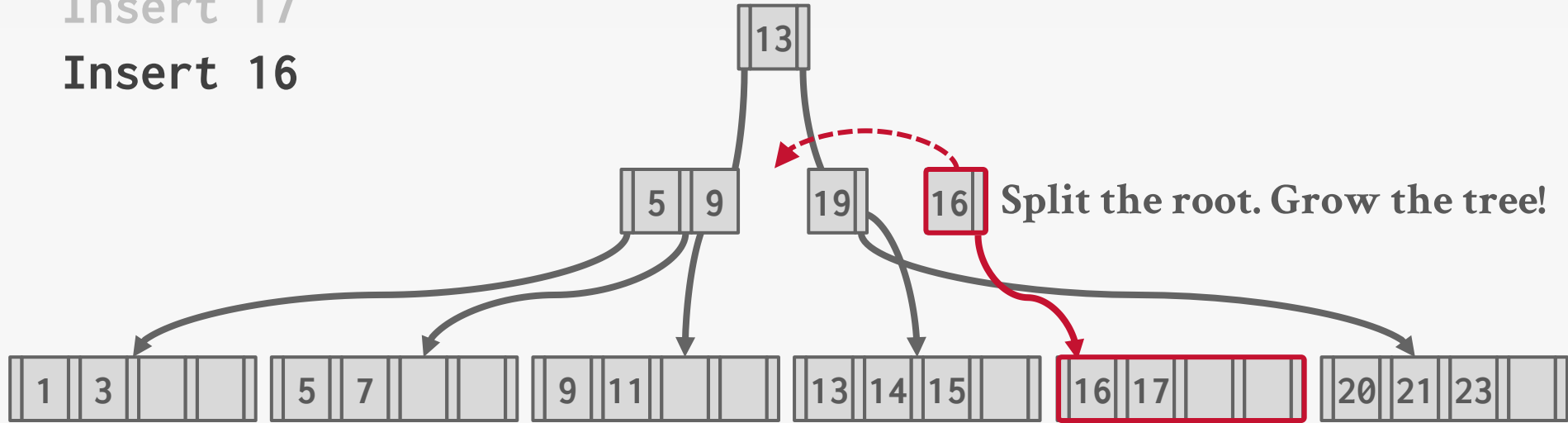
Insert 16



B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16



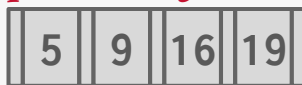
B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16



Next, need to split the "old" root, then point to the split nodes from the new root.



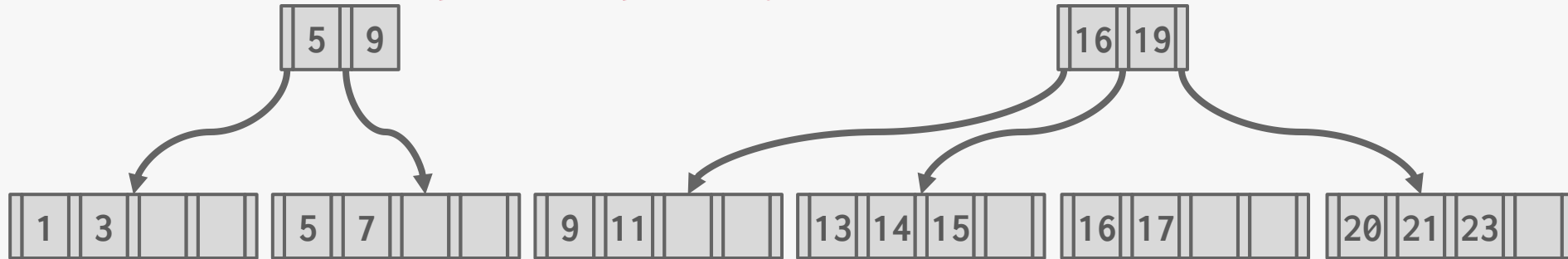
B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16



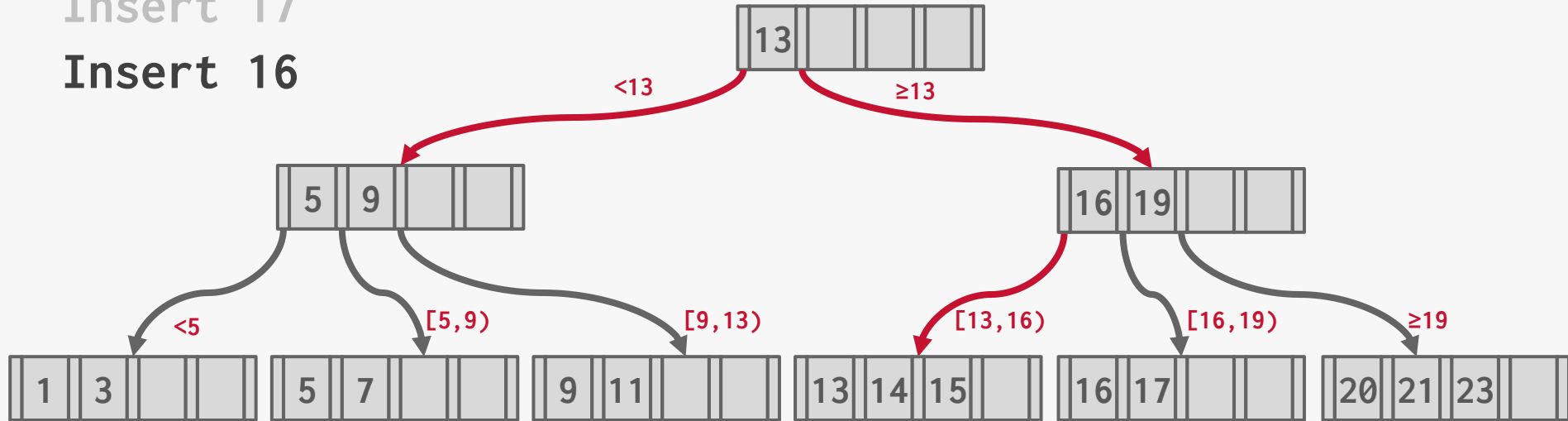
Next, need to split the "old" root, then point to the split nodes from the new root.



B+TREE: INSERT EXAMPLE (3)

Insert 17

Insert 16



B+TREE: DELETE

Start at root, find leaf **L** where entry belongs.

Remove the entry.

If **L** is at least half-full, done!

If **L** has only $m/2-1$ entries,

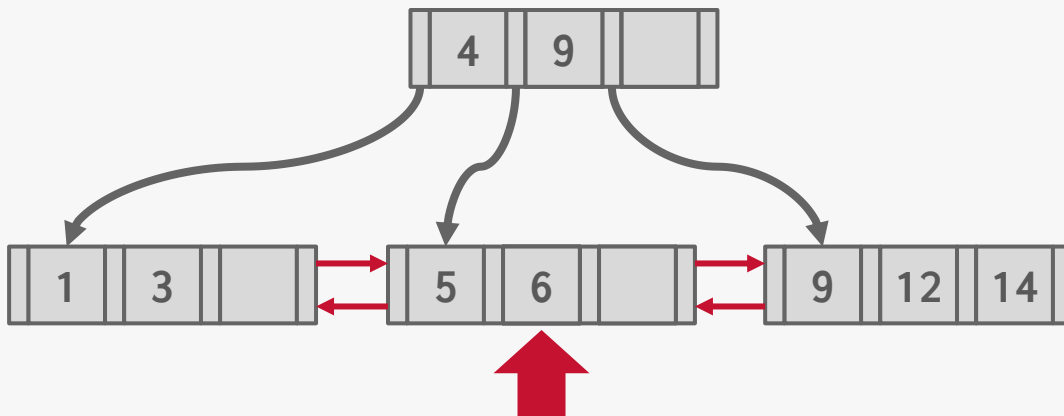
→ Try to re-distribute, borrowing from sibling (adjacent node with same parent as **L**).

→ If re-distribution fails, merge **L** and sibling.

If merge occurred, must delete entry (pointing to **L** or sibling) from parent of **L**.

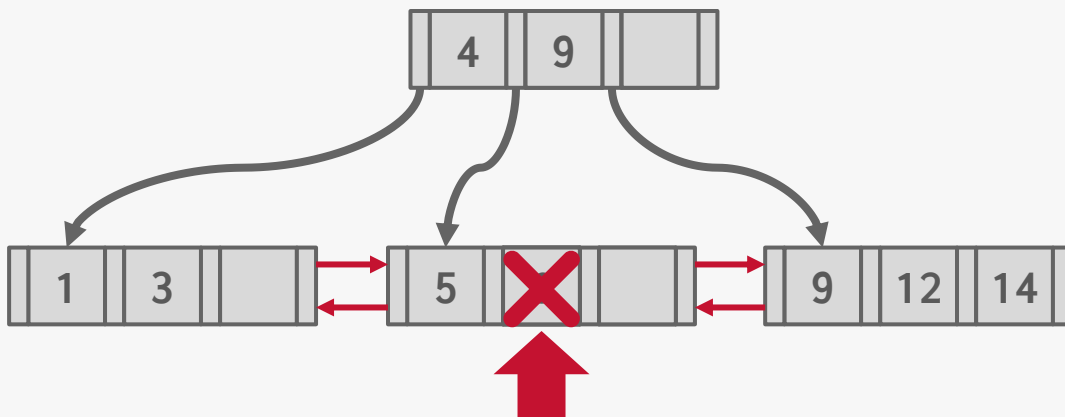
B+TREE: DELETE EXAMPLE (1)

Delete 6



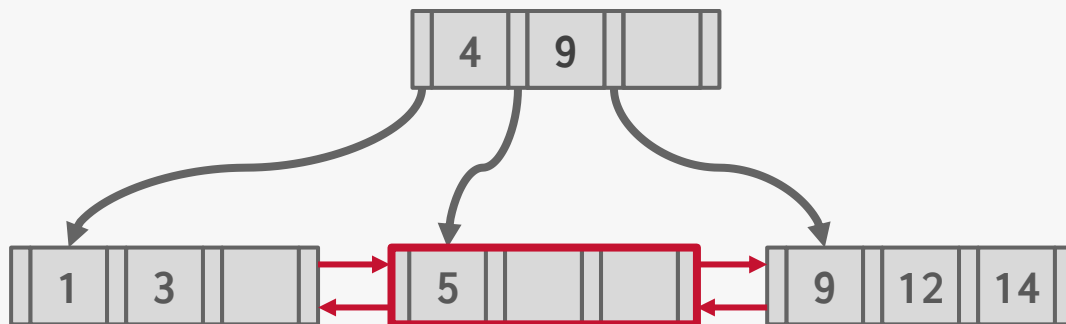
B+TREE: DELETE EXAMPLE (1)

Delete 6



B+TREE: DELETE EXAMPLE (1)

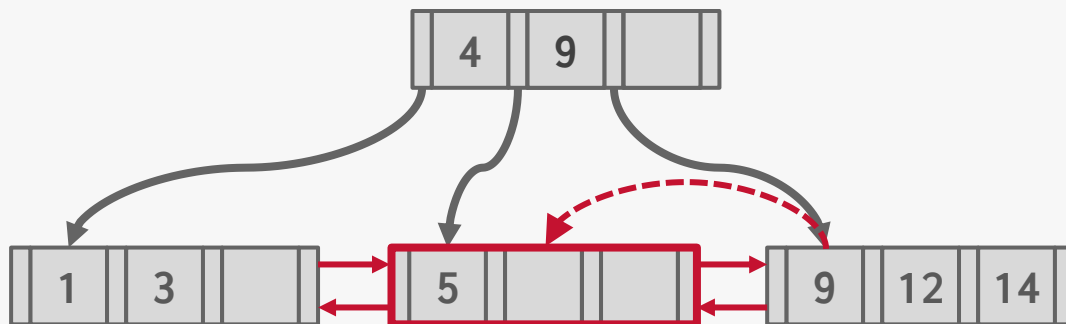
Delete 6



*Borrow from a "rich" sibling node.
Could borrow from either sibling.*

B+TREE: DELETE EXAMPLE (1)

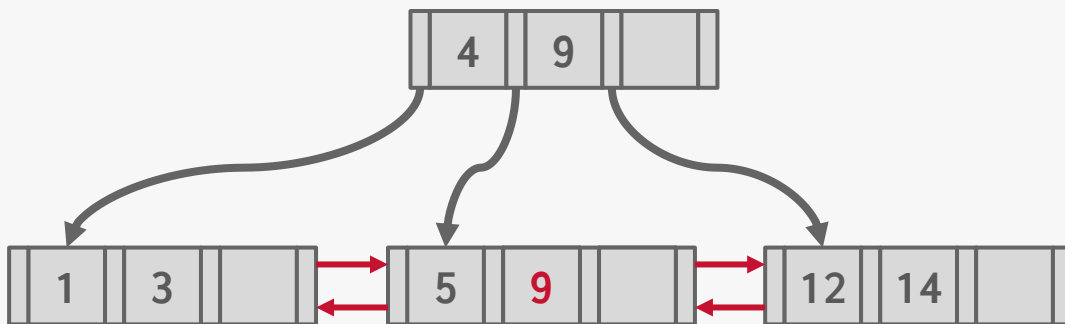
Delete 6



*Borrow from a "rich" sibling node.
Could borrow from either sibling.*

B+TREE: DELETE EXAMPLE (1)

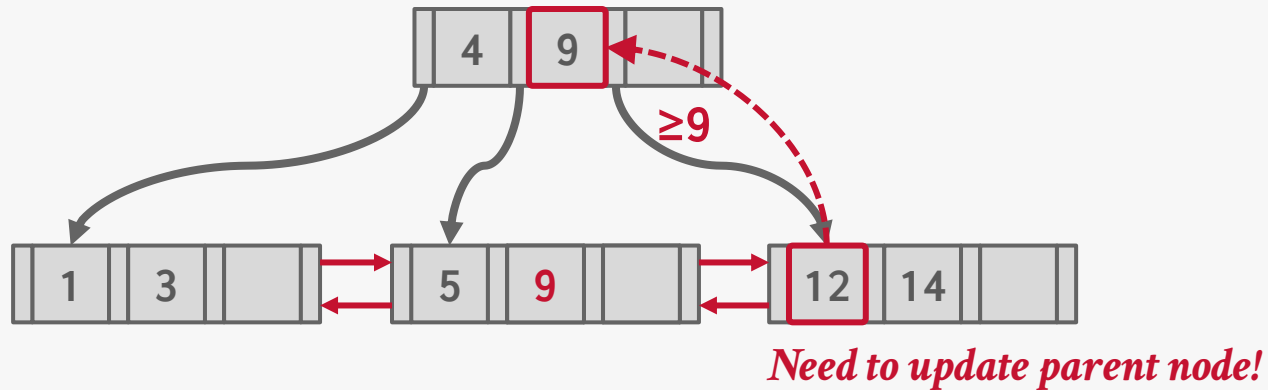
Delete 6



Need to update parent node!

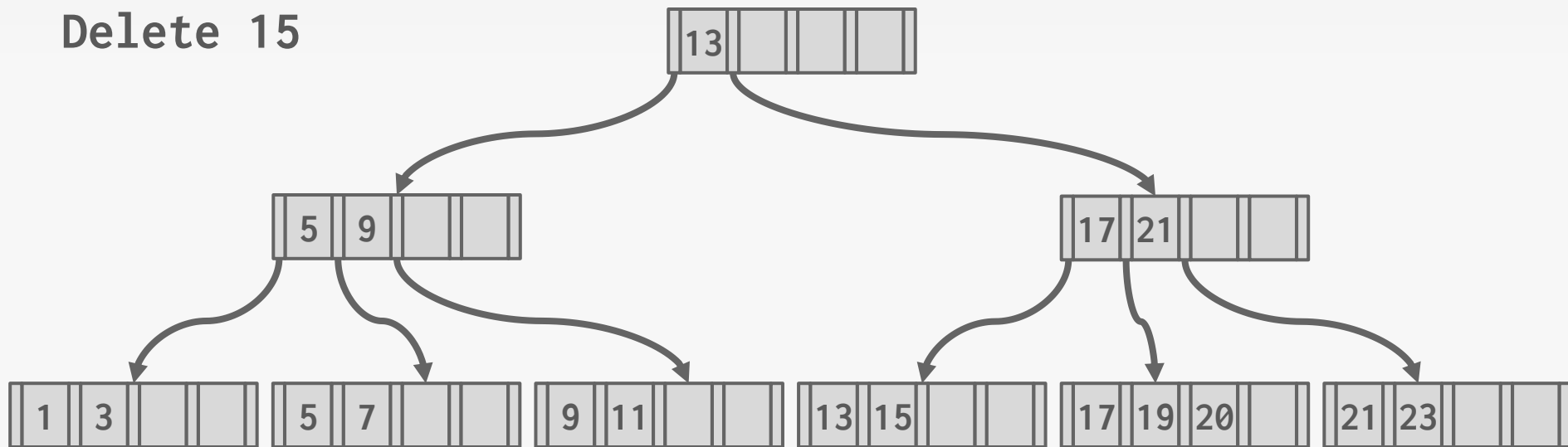
B+TREE: DELETE EXAMPLE (1)

Delete 6



B+TREE: DELETE EXAMPLE (2)

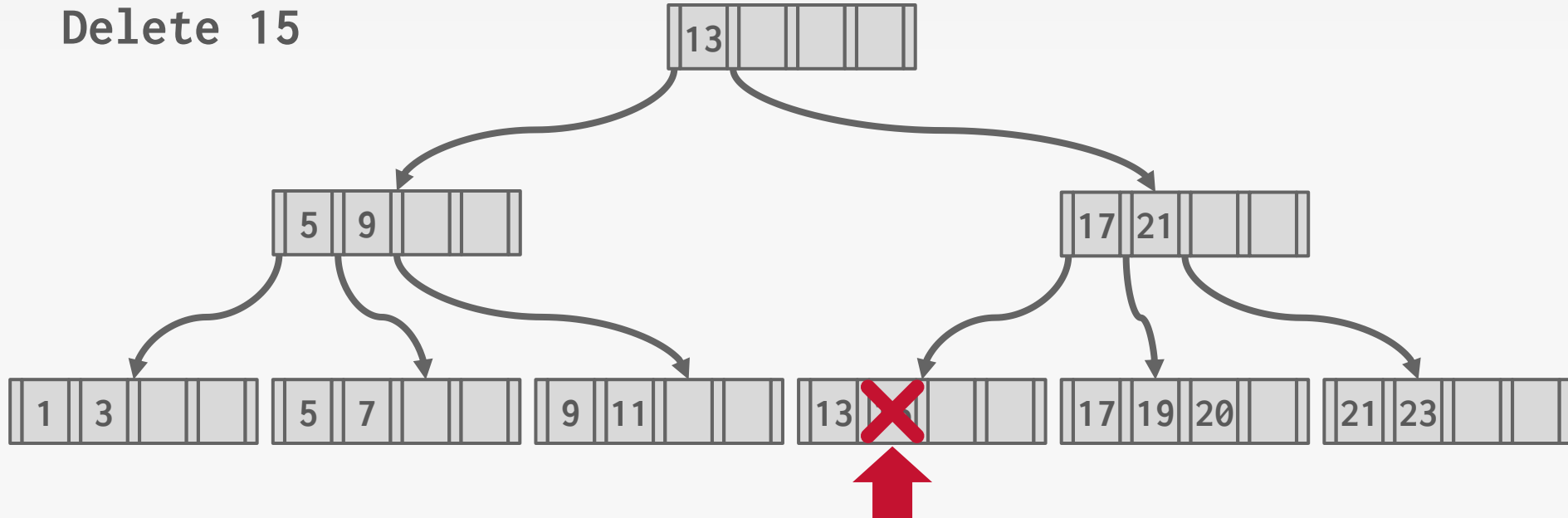
Delete 15



Note: New Example/Tree.

B+TREE: DELETE EXAMPLE (2)

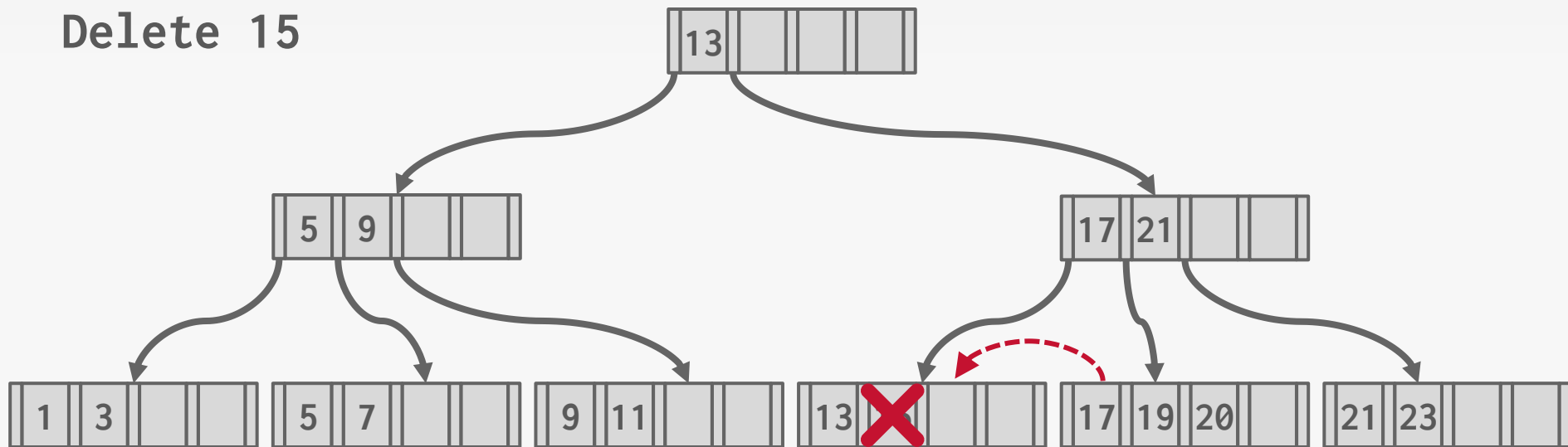
Delete 15



Note: New Example/Tree.

B+TREE: DELETE EXAMPLE (2)

Delete 15

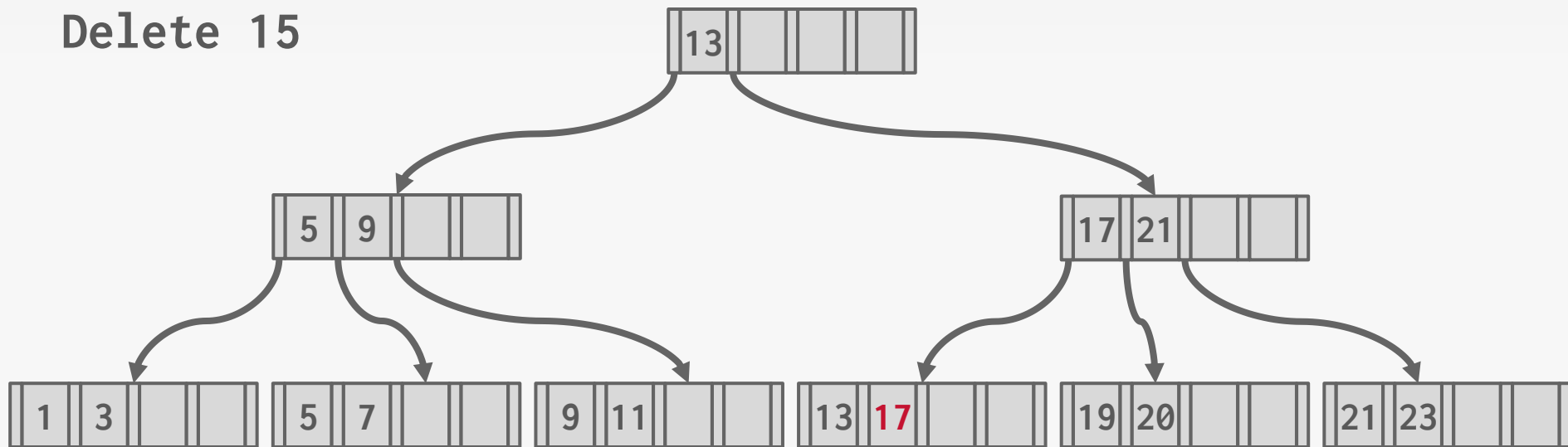


Borrow from a "rich" sibling node.

Note: New Example/Tree.

B+TREE: DELETE EXAMPLE (2)

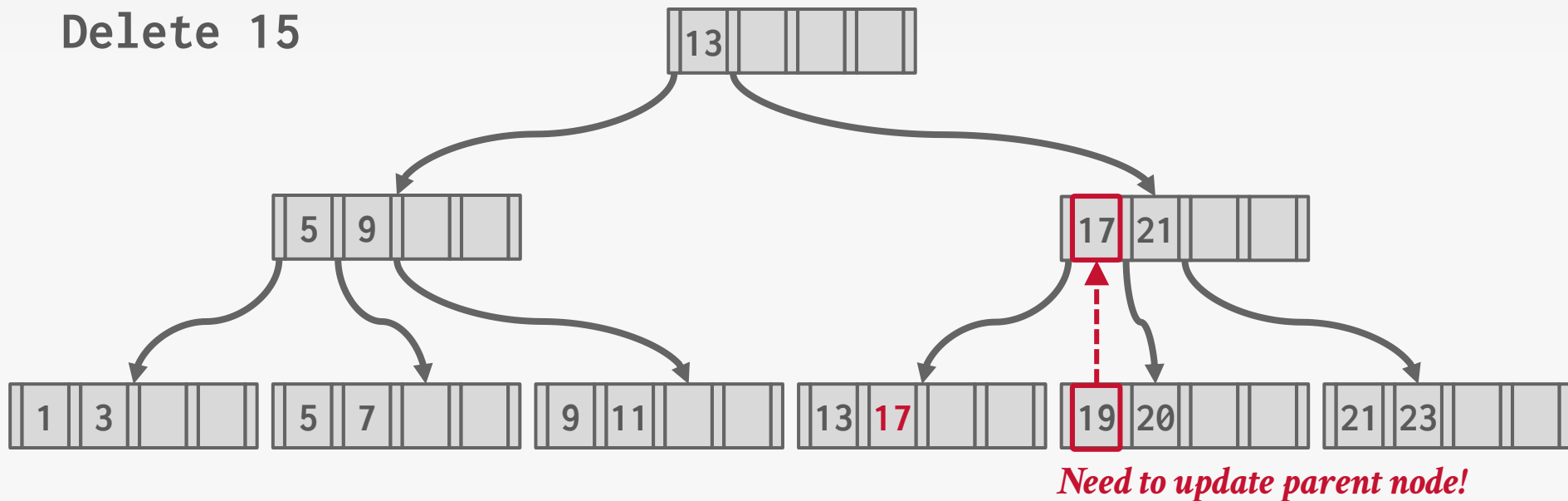
Delete 15



Need to update parent node!

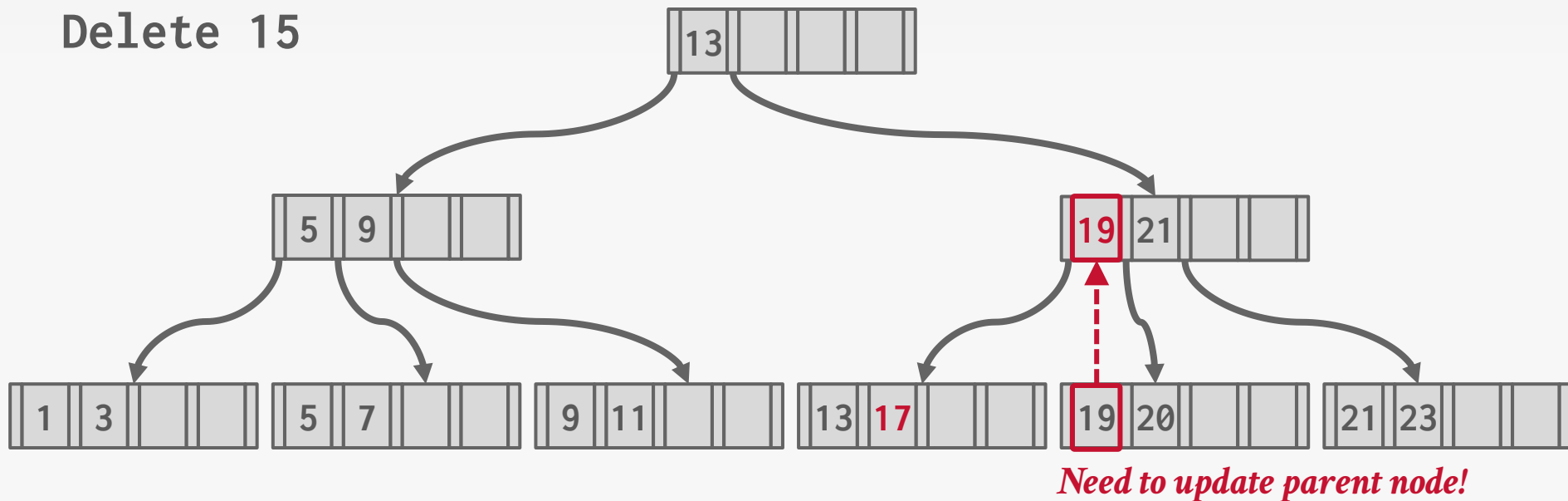
B+TREE: DELETE EXAMPLE (2)

Delete 15



B+TREE: DELETE EXAMPLE (2)

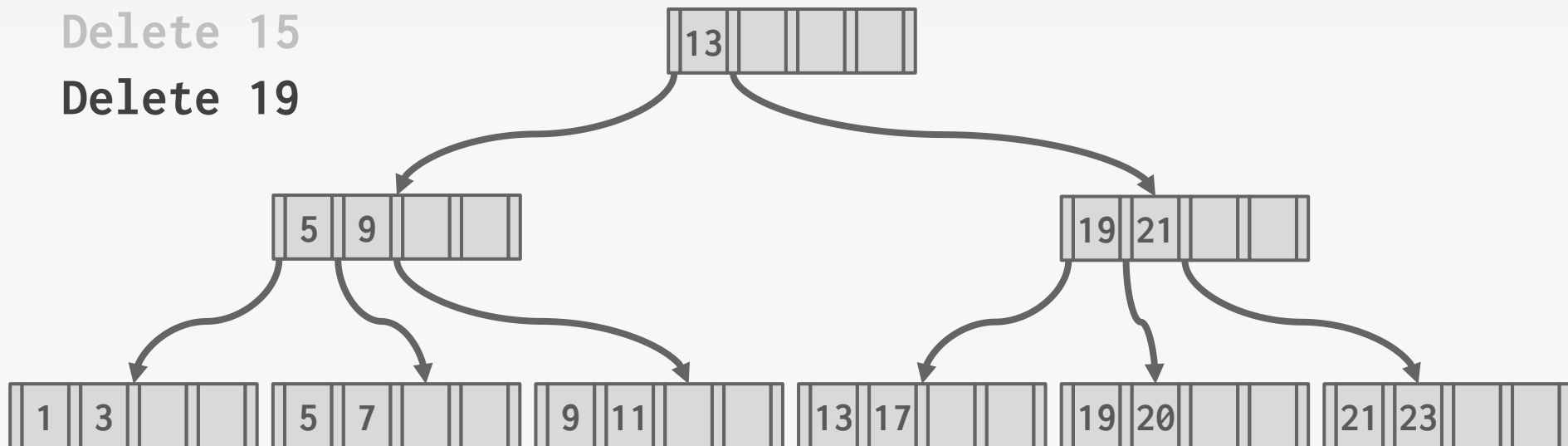
Delete 15



B+TREE: DELETE EXAMPLE (3)

Delete 15

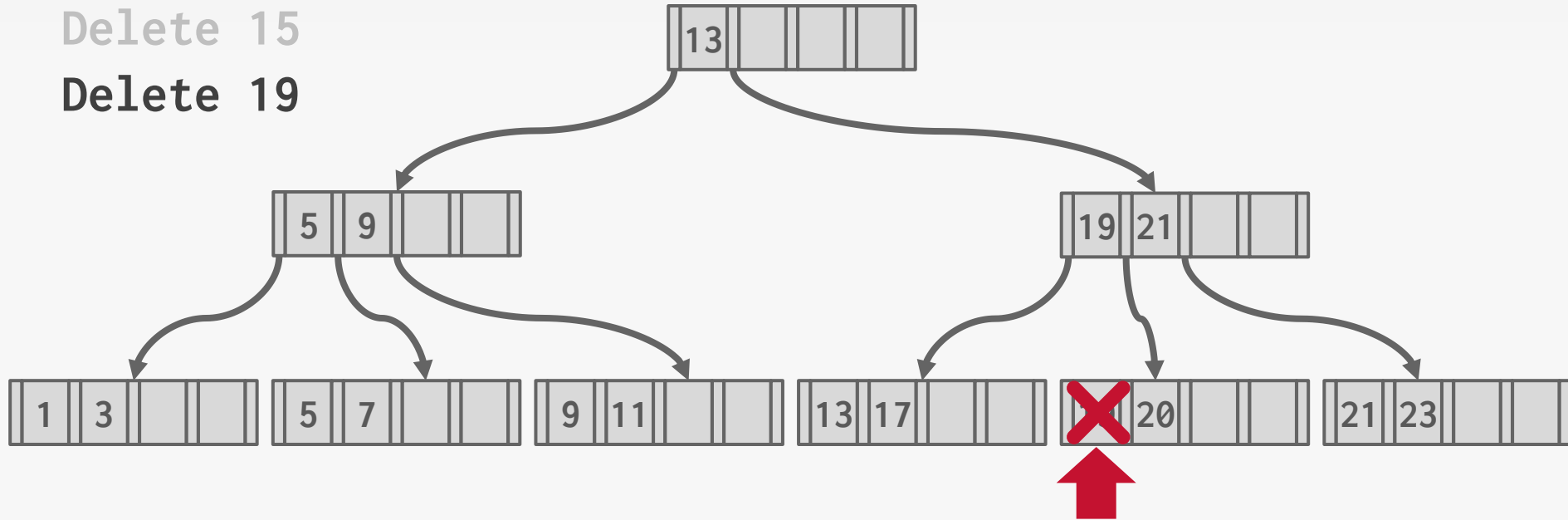
Delete 19



B+TREE: DELETE EXAMPLE (3)

Delete 15

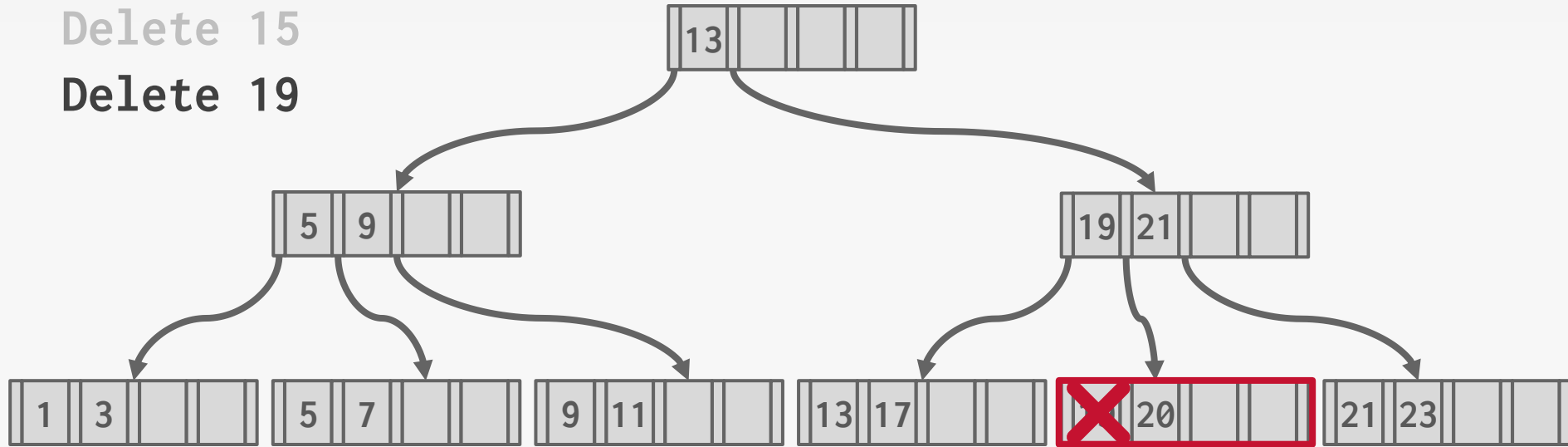
Delete 19



B+TREE: DELETE EXAMPLE (3)

Delete 15

Delete 19

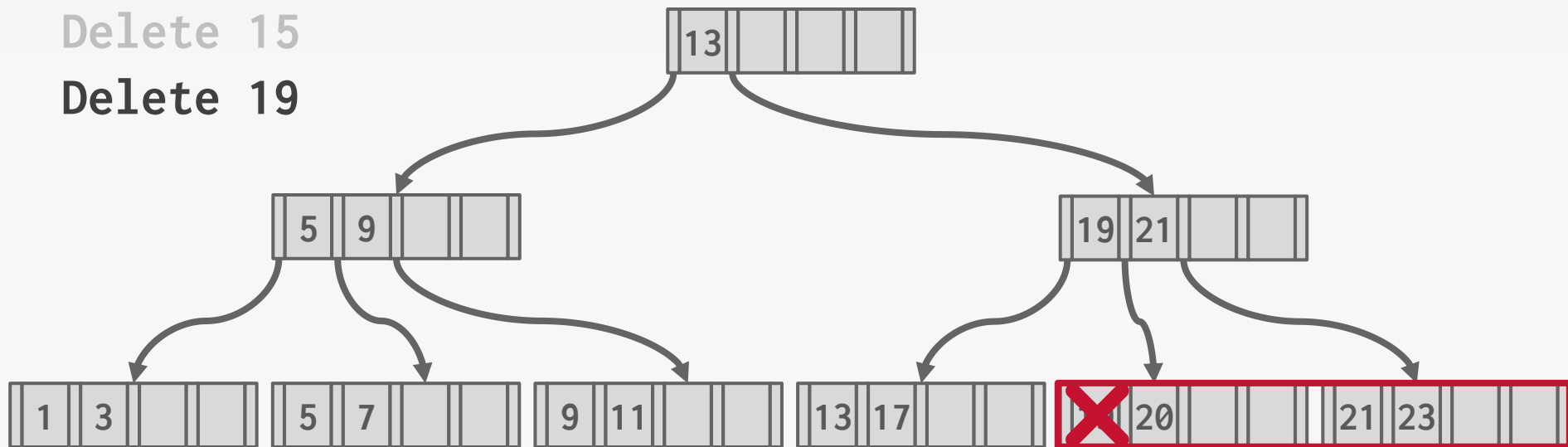


Under-filled!
No "rich" sibling nodes to borrow.
Merge with a sibling

B+TREE: DELETE EXAMPLE (3)

Delete 15

Delete 19

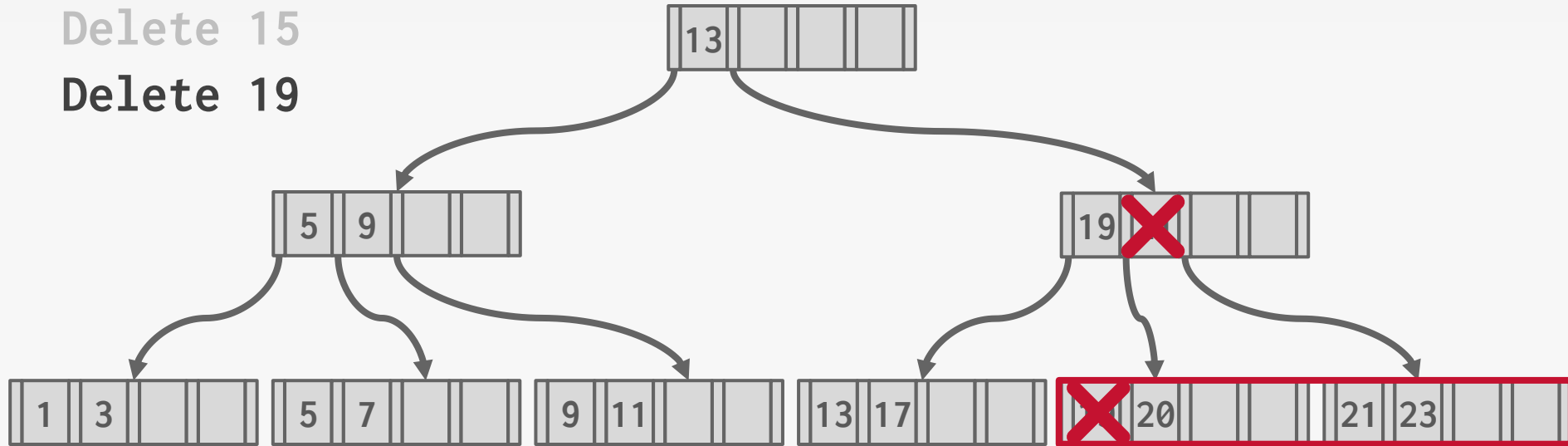


Under-filled!
No "rich" sibling nodes to borrow.
Merge with a sibling

B+TREE: DELETE EXAMPLE (3)

Delete 15

Delete 19



Under-filled!

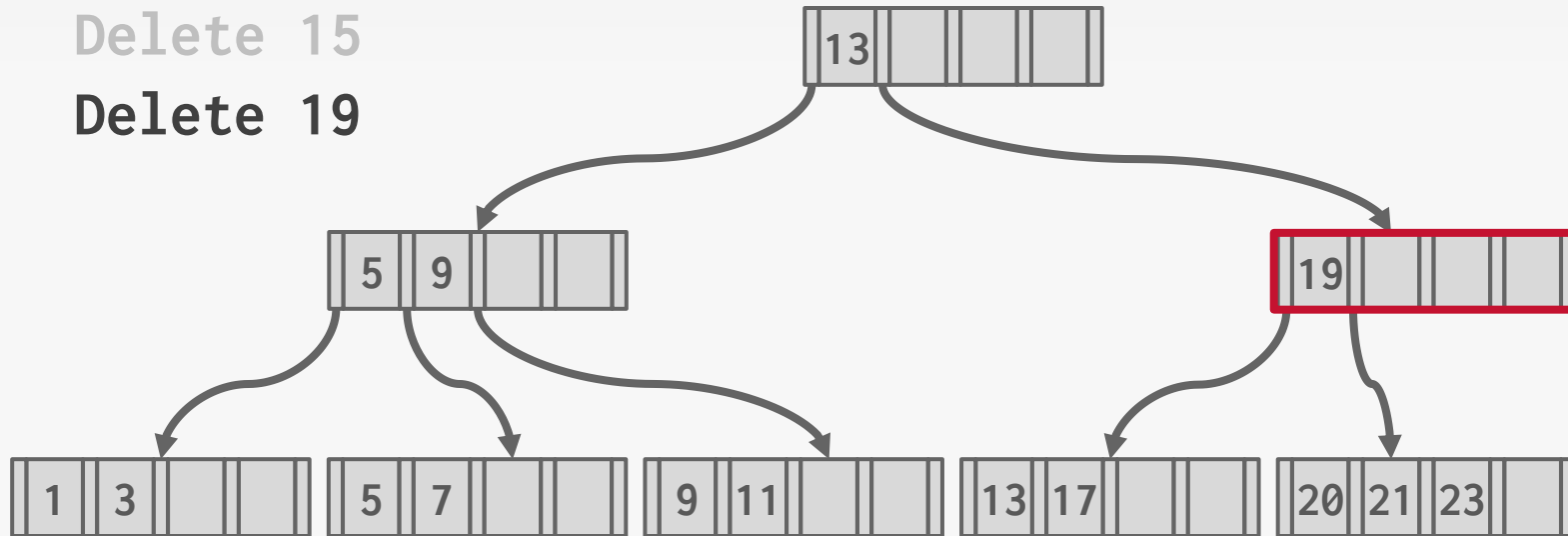
No "rich" sibling nodes to borrow.

Merge with a sibling

B+TREE: DELETE EXAMPLE (3)

Delete 15

Delete 19



*This node is under-filled!
Pull-down.*

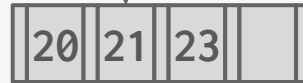
B+TREE: DELETE EXAMPLE (3)

Delete 15

Delete 19



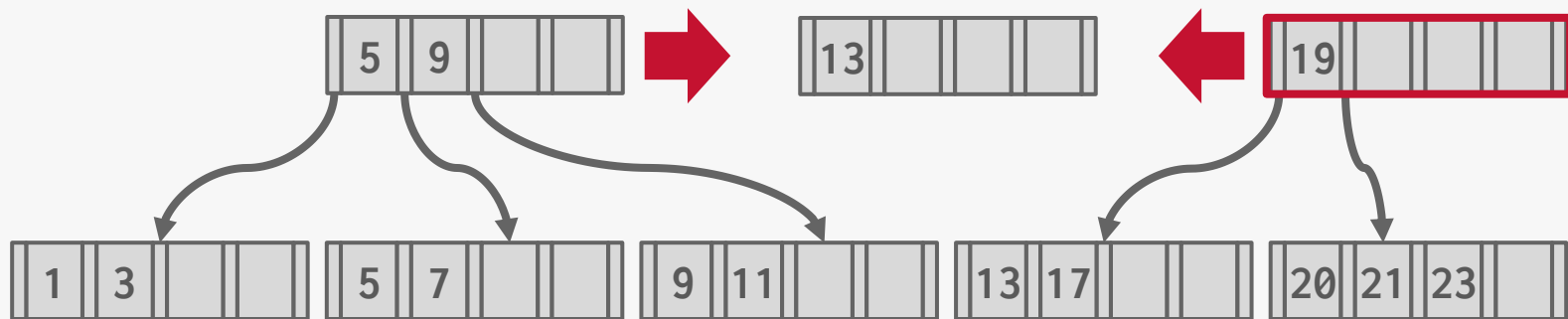
*This node is
under-filled!
Pull-down.*



B+TREE: DELETE EXAMPLE (3)

Delete 15

Delete 19



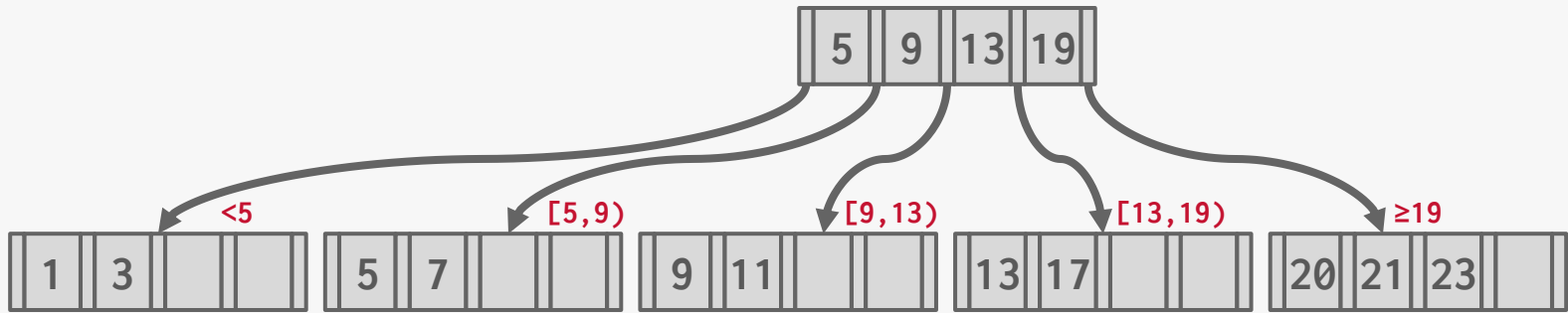
*This node is under-filled!
Pull-down.*

B+TREE: DELETE EXAMPLE (3)

Delete 15

Delete 19

The tree has shrunk in height.



COMPOSITE INDEX

A composite index is when the key is comprised of two or more attributes.

→ Example: Index on **<a, b, c>**

```
CREATE INDEX my_idx ON xxx (a, b DESC, c NULLS FIRST);
```

DBMS can use B+Tree index if the query provides a “prefix” of composite key.

→ Supported: **(a=1 AND b=2 AND c=3)**

→ Supported: **(a=1 AND b=2)**

→ Rarely Supported: **(b=2), (c=3)**

COMPOSITE INDEX

A composite index is when the key is comprised of two or more attributes.

→ Example: Index on **<a, b, c>**

```
CREATE INDEX my_idx ON xxx (a, b DESC, c NULLS FIRST);
```



DBMS can use B+Tree index if the query provides a “prefix” of composite key.

→ Supported: **(a=1 AND b=2 AND c=3)**

→ Supported: **(a=1 AND b=2)**

→ Rarely Supported: **(b=2), (c=3)**

COMPOSITE INDEX

A composite index is when the key is comprised of two or more attributes.

→ Example: Index on **<a, b, c>**

```
CREATE INDEX my_idx ON xxx (a, b DESC, c NULLS FIRST);
```

Sort Order

Null Handling

DBMS can use B+Tree index if the query provides a “prefix” of composite key.

→ Supported: **(a=1 AND b=2 AND c=3)**

→ Supported: **(a=1 AND b=2)**

→ Rarely Supported: **(b=2), (c=3)** ← *Skip Scans*

ORACLE®

PostgreSQL

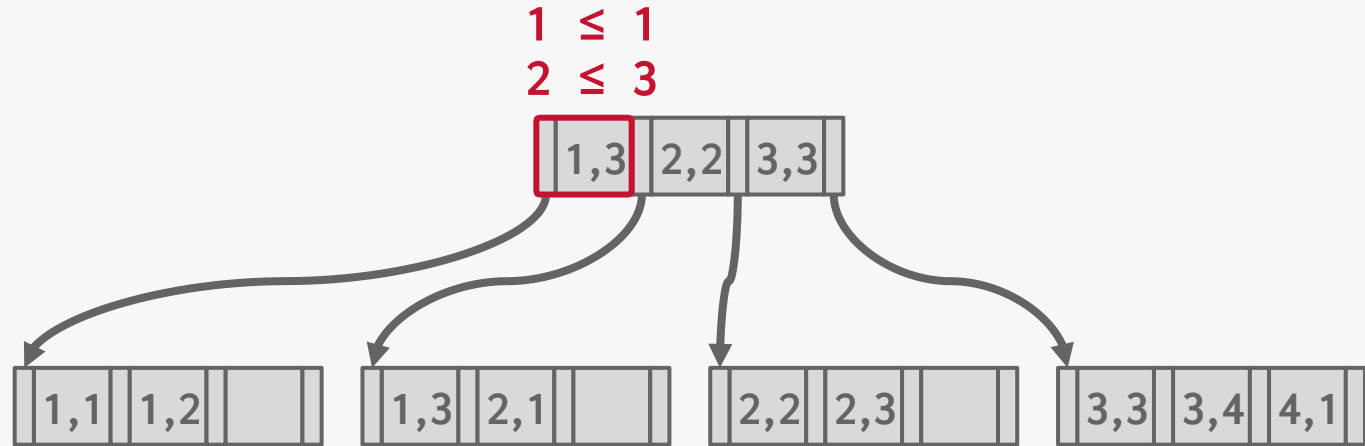
MySQL™

MariaDB

yugabyteDB

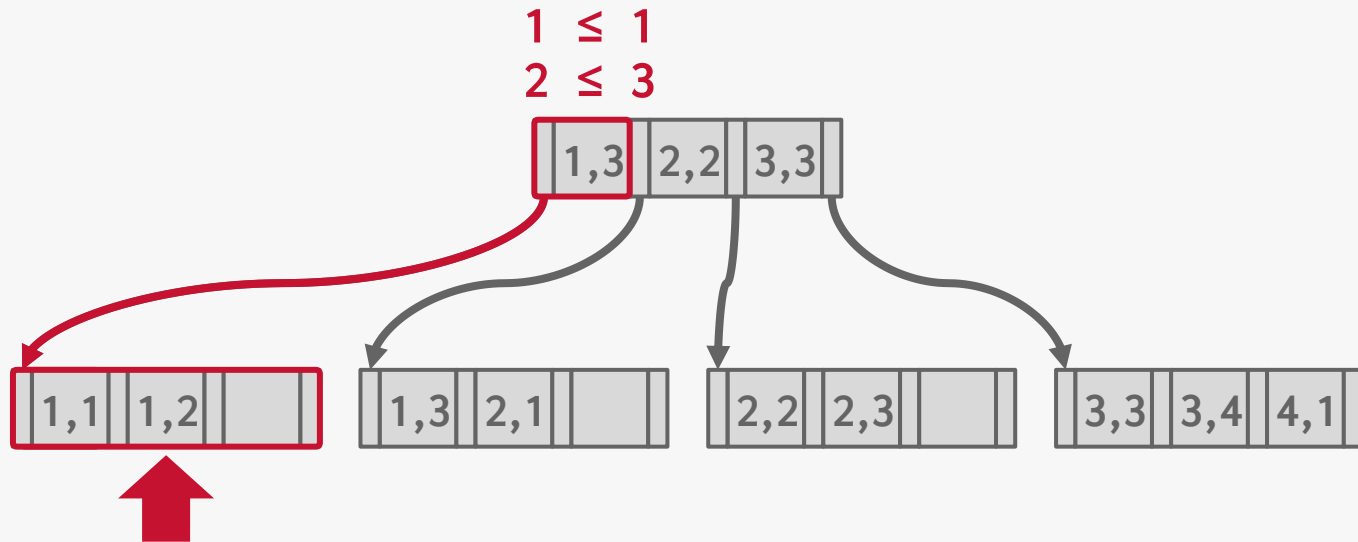
SELECTION CONDITIONS

Find Key=(1,2)



SELECTION CONDITIONS

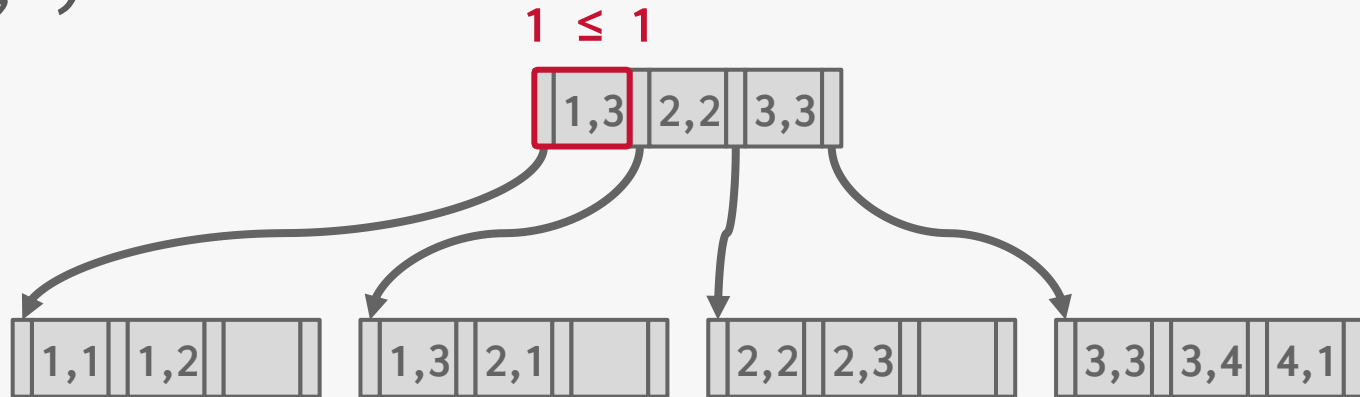
Find Key=(1,2)



SELECTION CONDITIONS

Find Key=(1,2)

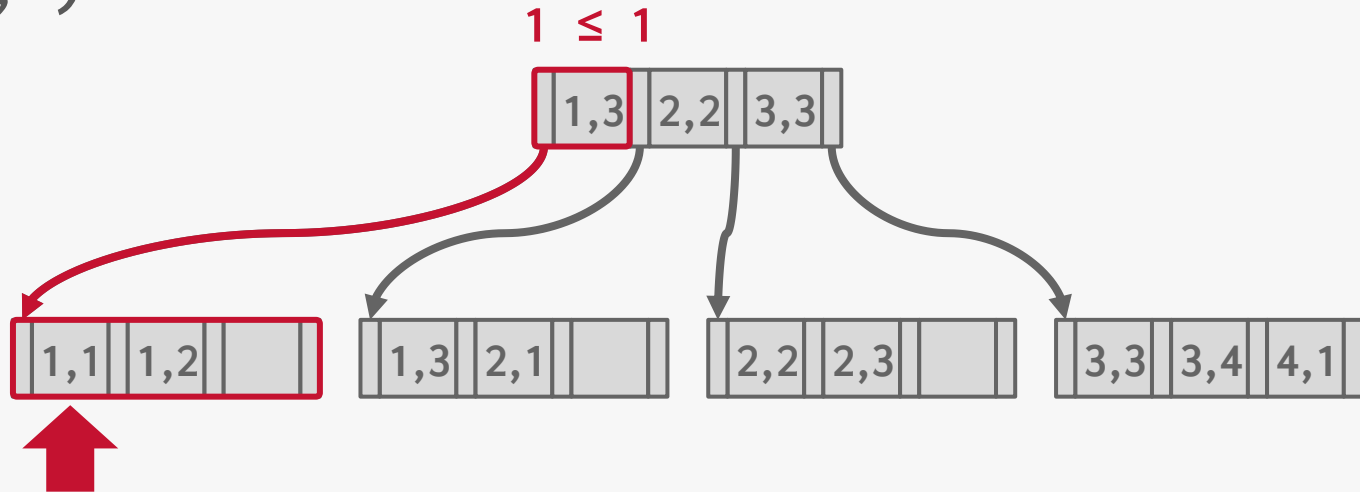
Find Key=(1,*)



SELECTION CONDITIONS

Find Key=(1,2)

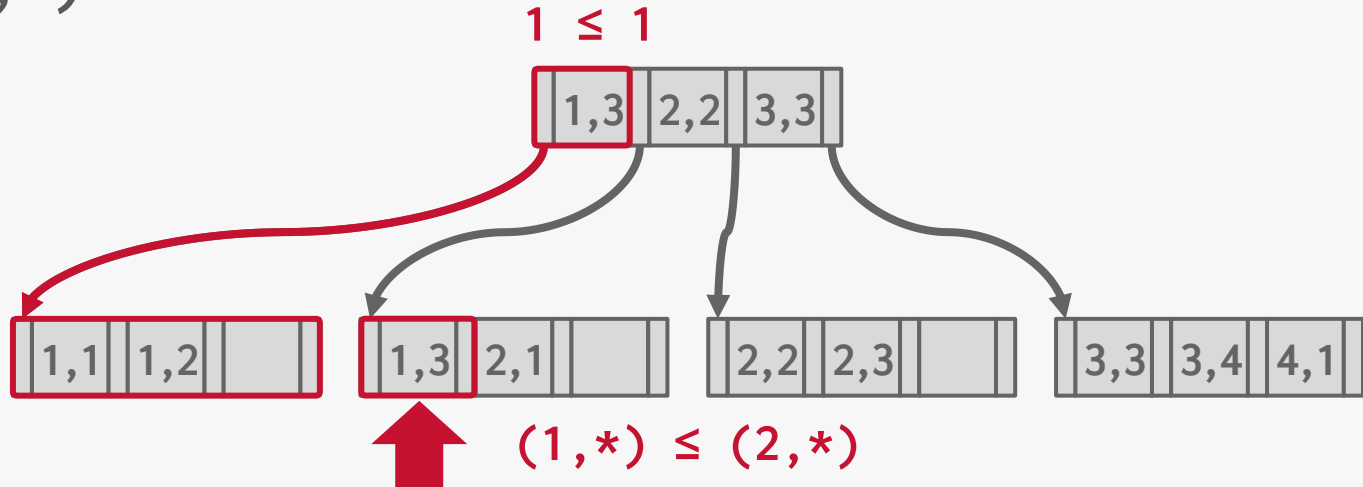
Find Key=(1,*)



SELECTION CONDITIONS

Find Key=(1,2)

Find Key=(1,*)

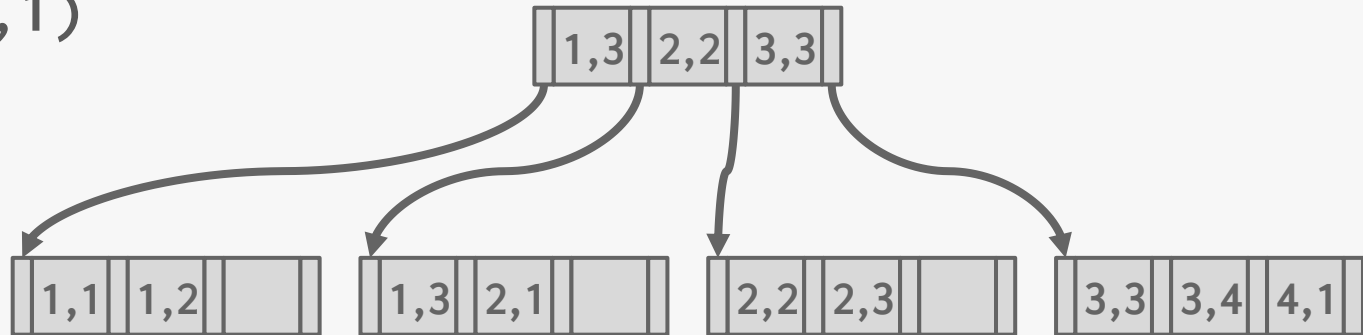


SELECTION CONDITIONS

Find Key=(1,2)

Find Key=(1,*)

Find Key=(*,1)

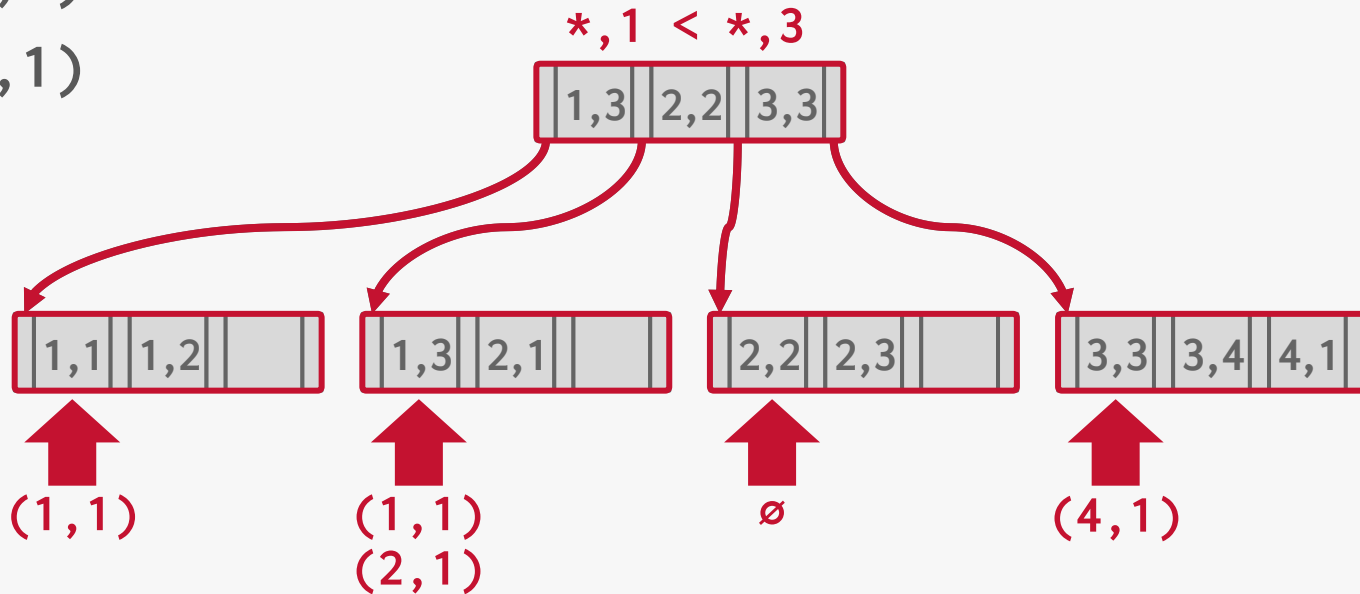


SELECTION CONDITIONS

Find Key=(1,2)

Find Key=(1,*)

Find Key=(*,1)



B+TREE: DUPLICATE KEYS

Approach #1: Append Record ID

- Add the tuple's unique Record ID as part of the key to ensure that all keys are unique.
- The DBMS can still use partial keys to find tuples.

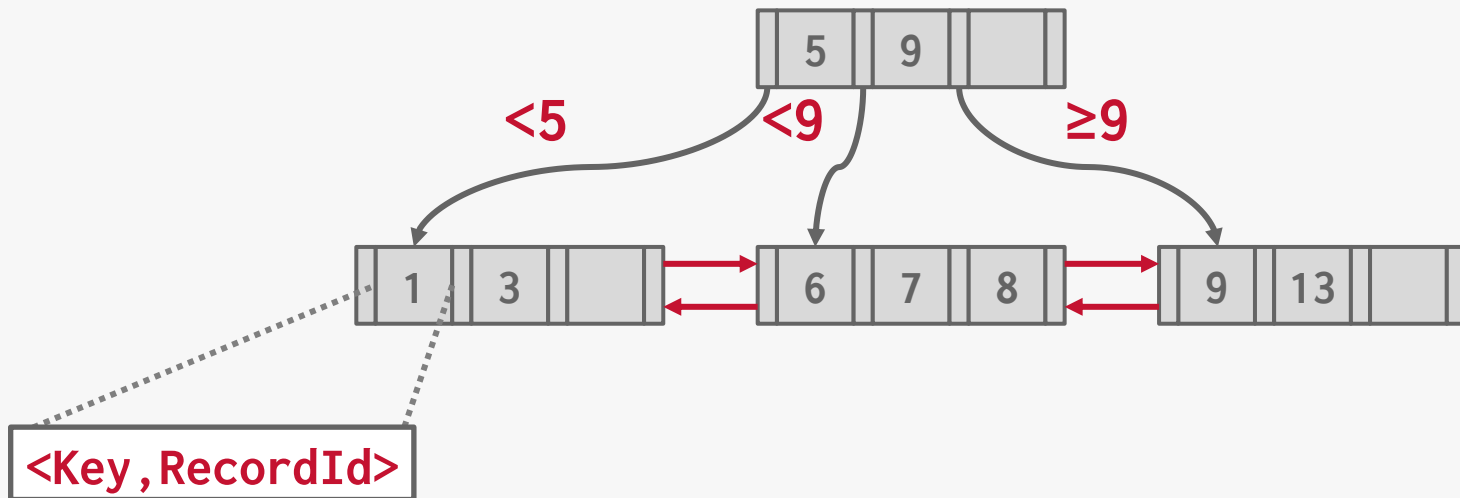


*Don't
Do This!*

Approach #2: Overflow Leaf Nodes

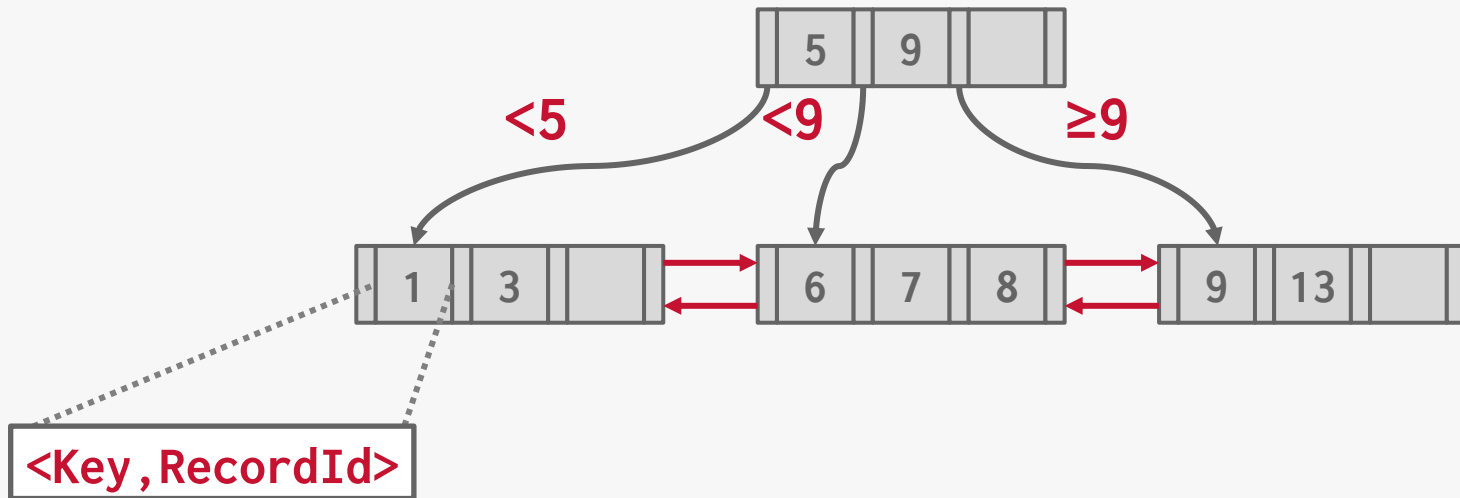
- Allow leaf nodes to spill into overflow nodes that contain the duplicate keys.
- This is more complex to maintain and modify.

B+TREE: APPEND RECORD ID



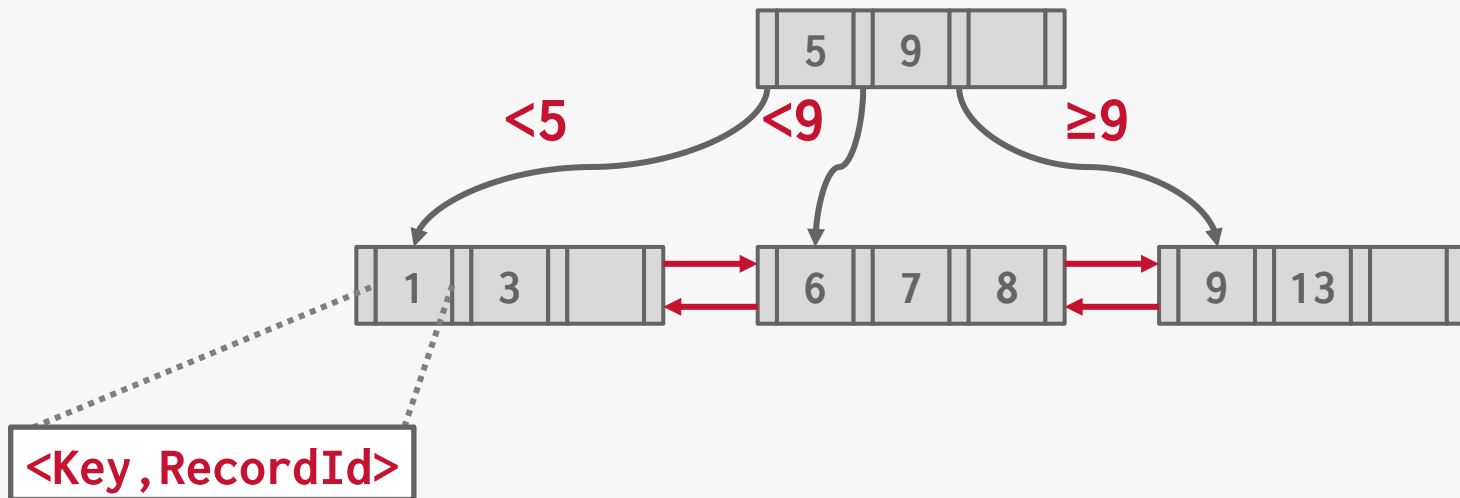
B+TREE: APPEND RECORD ID

Insert 6



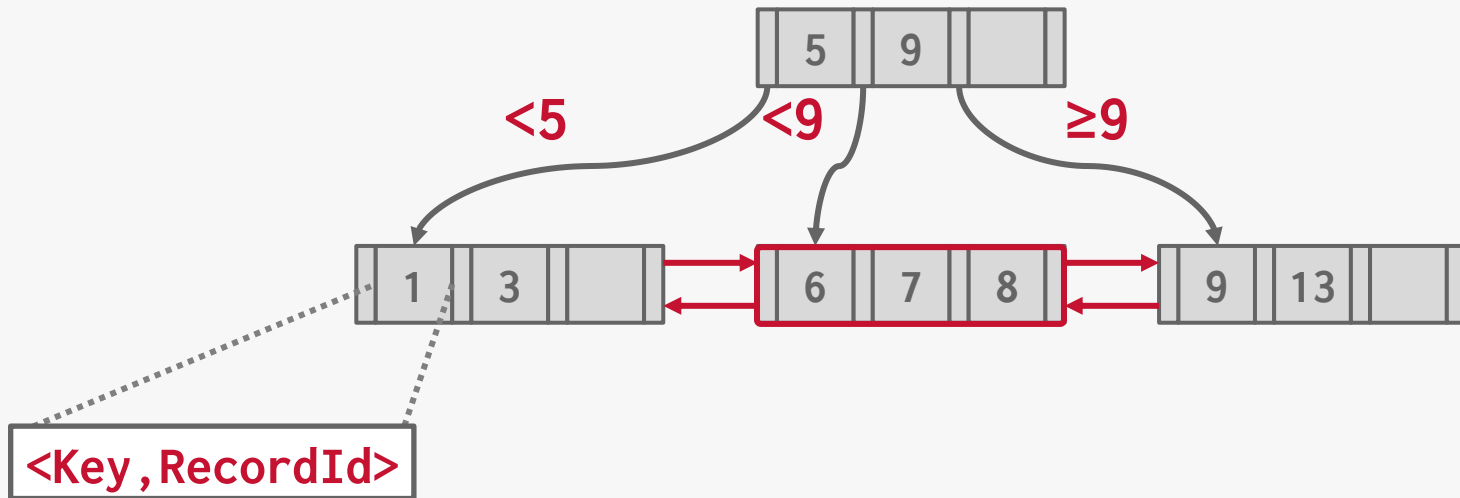
B+TREE: APPEND RECORD ID

Insert $\langle 6, (\text{Page}, \text{Slot}) \rangle$



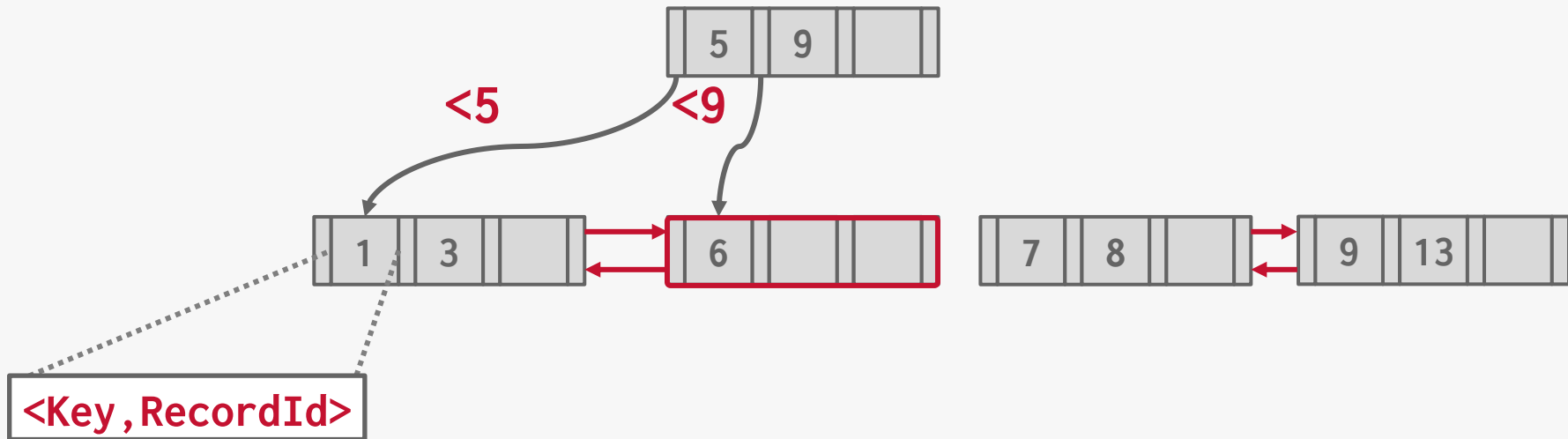
B+TREE: APPEND RECORD ID

Insert $\langle 6, (\text{Page}, \text{Slot}) \rangle$



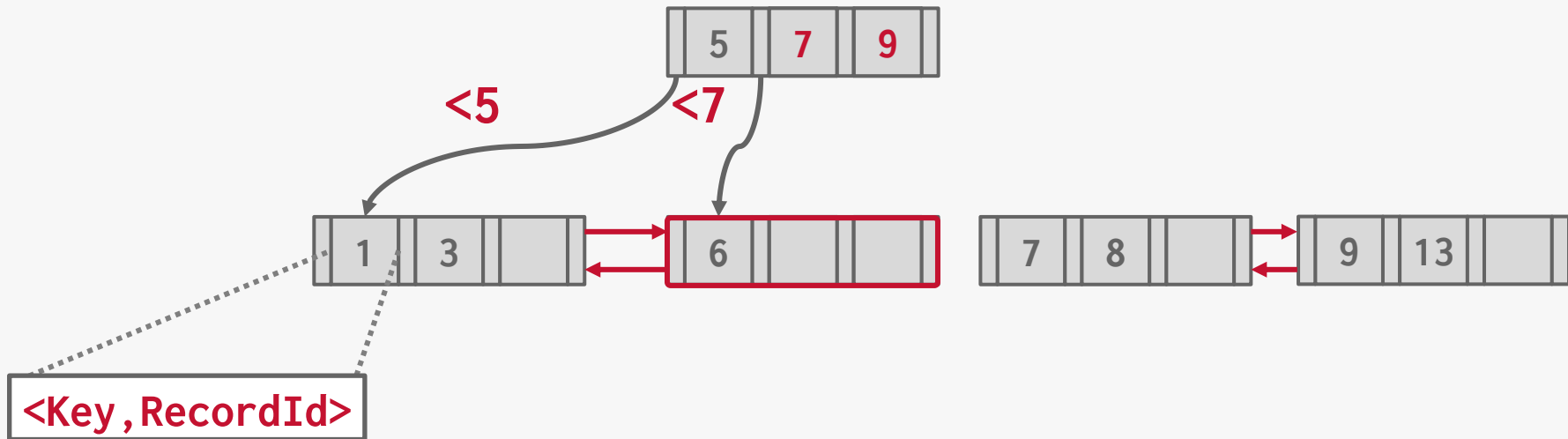
B+TREE: APPEND RECORD ID

Insert $\langle 6, (\text{Page}, \text{Slot}) \rangle$



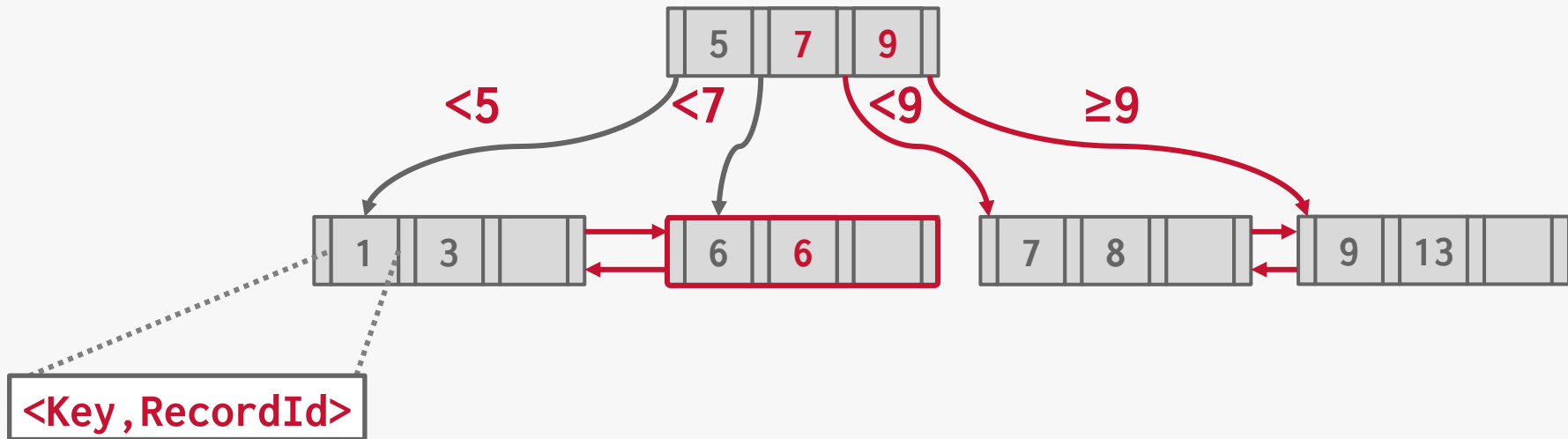
B+TREE: APPEND RECORD ID

Insert $\langle 6, (\text{Page}, \text{Slot}) \rangle$



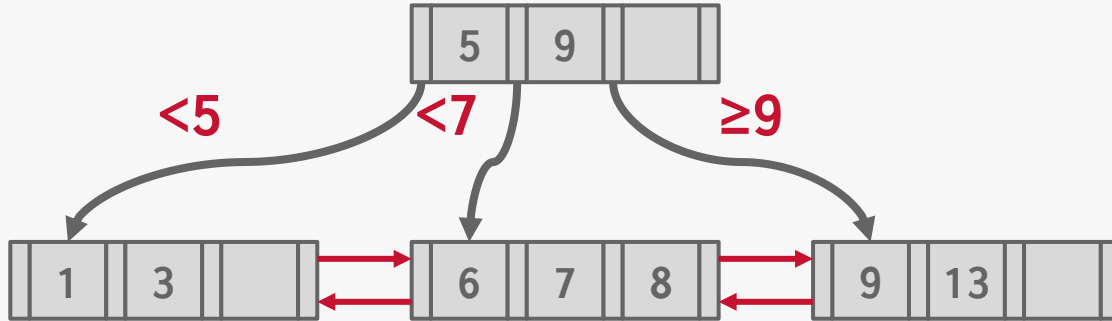
B+TREE: APPEND RECORD ID

Insert $\langle 6, (\text{Page}, \text{Slot}) \rangle$



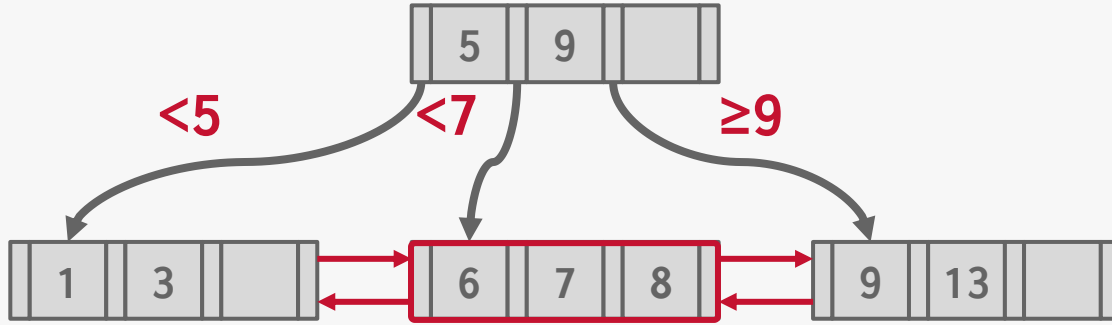
B+TREE: OVERFLOW LEAF NODES

Insert 6



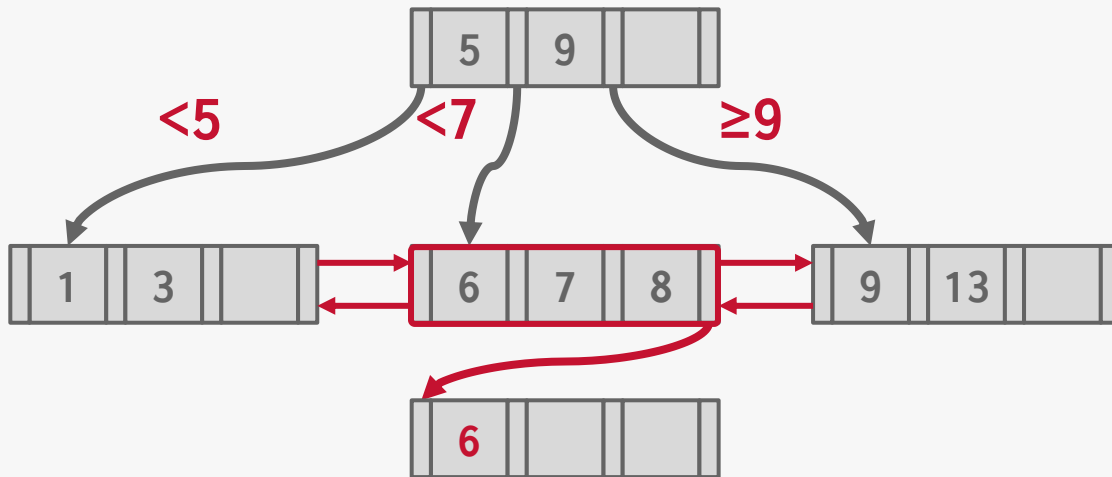
B+TREE: OVERFLOW LEAF NODES

Insert 6



B+TREE: OVERFLOW LEAF NODES

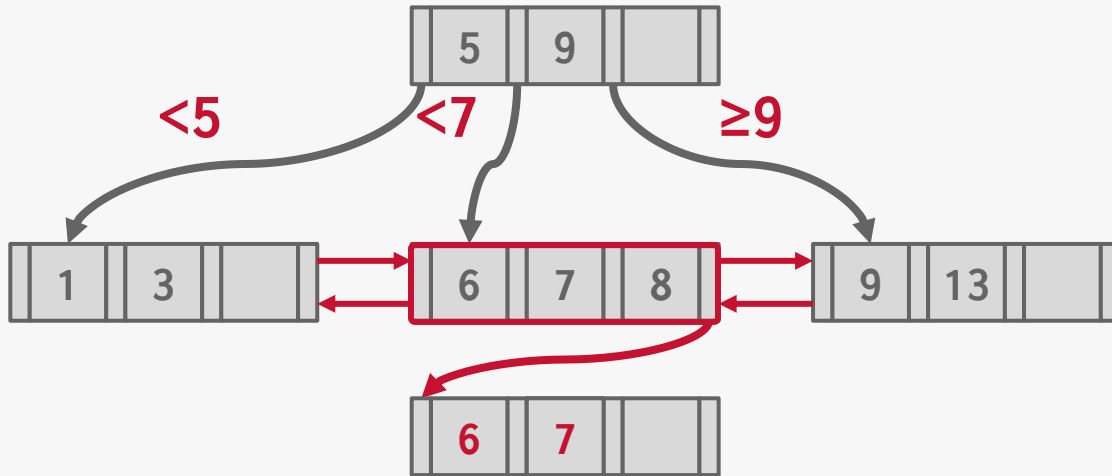
Insert 6



B+TREE: OVERFLOW LEAF NODES

Insert 6

Insert 7

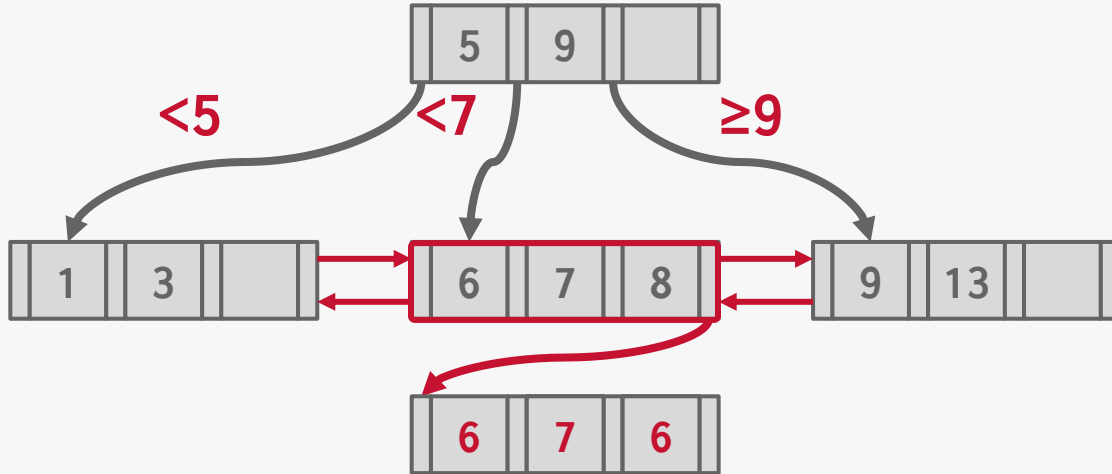


B+TREE: OVERFLOW LEAF NODES

Insert 6

Insert 7

Insert 6



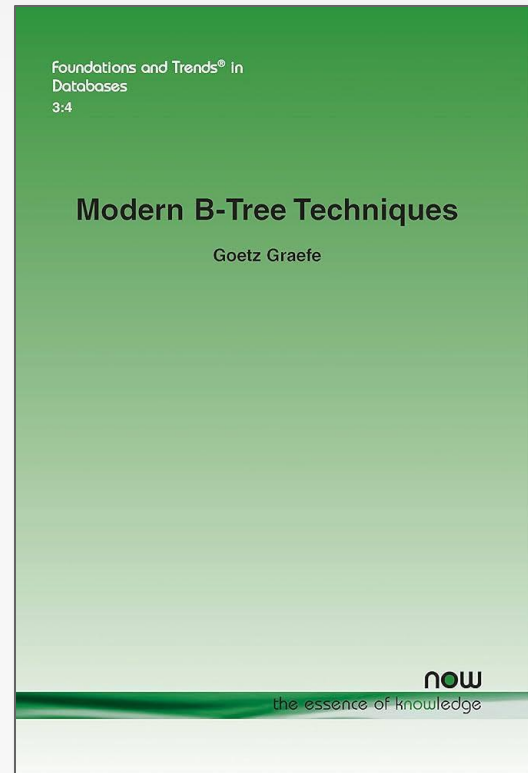
B+TREE DESIGN CHOICES

Node Size

Merge Threshold

Variable-Length Keys

Intra-Node Search



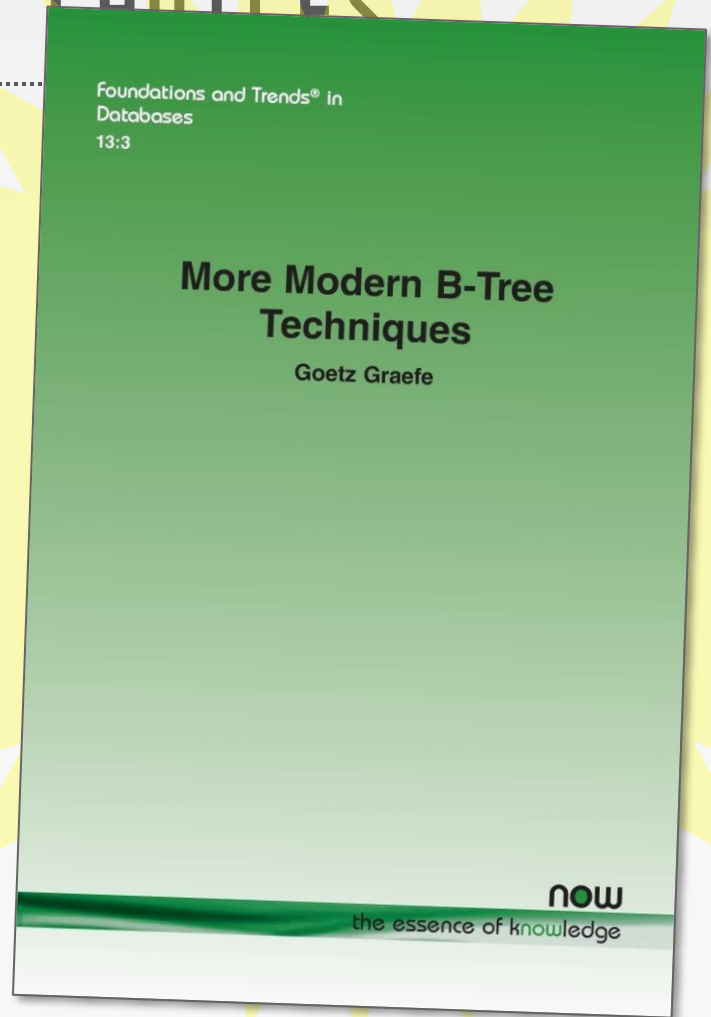
B+TREE DESIGN CHOICES

Node Size

Merge Threshold

Variable-Length Keys

Intra-Node Search



NODE SIZE

The slower the storage device, the larger the optimal node size for a B+Tree.

→ HDD: ~1MB

→ SSD: ~10KB

→ In-Memory: ~512B

Optimal sizes can vary depending on the workload

→ Leaf Node Scans vs. Root-to-Leaf Traversals

MERGE THRESHOLD

Some DBMSs do not always merge nodes when they are half full.

→ Average occupancy rate for B+Tree nodes is 69%.

Delaying a merge operation may reduce the amount of reorganization.

It may also be better to let underfilled nodes exist and then periodically rebuild entire tree.

This is why PostgreSQL calls their B+Tree a "non-balanced" B+Tree ([nbtree](#)).

VARIABLE-LENGTH KEYS

Approach #1: Pointers

- Store the keys as pointers to the tuple's attribute.
- Also called **T-Trees** (in-memory DBMSs)

Approach #2: Variable-Length Nodes

- The size of each node in the index can vary.
- Requires careful memory management.

Approach #3: Padding

- Always pad the key to be max length of the key type.

Approach #4: Key Map / Indirection

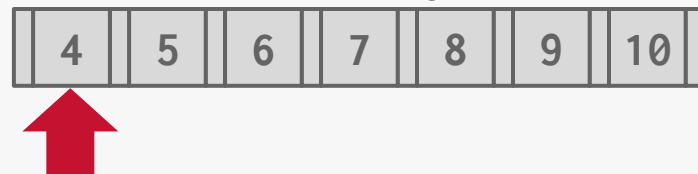
- Embed an array of pointers that map to the key + value list within the node.

INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Find Key=8

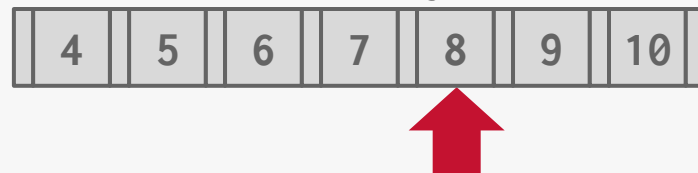


INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Find Key=8



INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Find Key=8

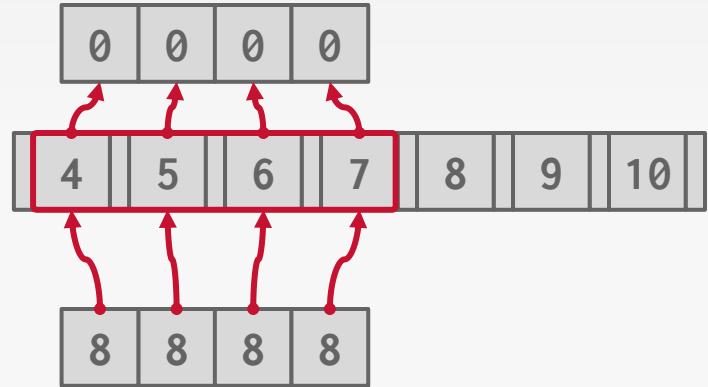


```
_mm_cmpeq_epi32_mask(a, b)
```

INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

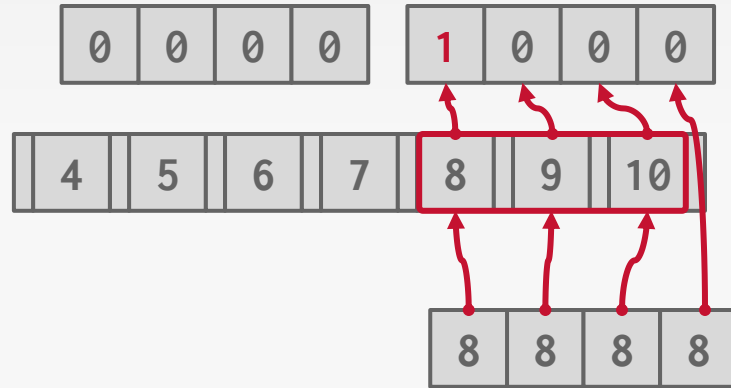


```
_mm_cmpeq_epi32_mask(a, b)
```

INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.



```
_mm_cmpeq_epi32_mask(a, b)
```

INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Approach #2: Binary

- Jump to middle key, pivot left/right depending on comparison.



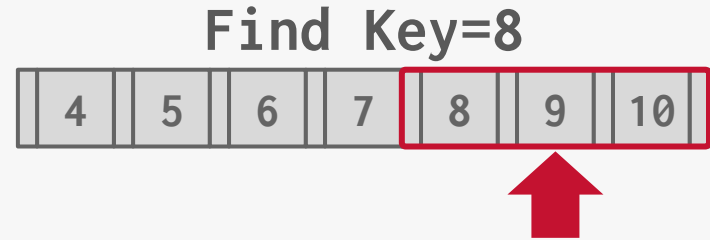
INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Approach #2: Binary

- Jump to middle key, pivot left/right depending on comparison.



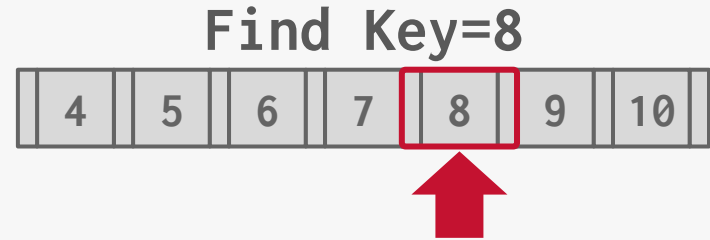
INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Approach #2: Binary

- Jump to middle key, pivot left/right depending on comparison.



INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Approach #2: Binary

- Jump to middle key, pivot left/right depending on comparison.

Approach #3: Interpolation

- Approximate location of desired key based on known distribution of keys.



INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Approach #2: Binary

- Jump to middle key, pivot left/right depending on comparison.

Approach #3: Interpolation

- Approximate location of desired key based on known distribution of keys.

$$\text{Offset: } (8-4) \times 7 / (10-4) = 4$$

4	5	6	7	8	9	10
---	---	---	---	---	---	----

INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Approach #2: Binary

- Jump to middle key, pivot left/right depending on comparison.

Approach #3: Interpolation

- Approximate location of desired key based on known distribution of keys.

$$\text{Offset: } (8-4) \times 7 / (10-4) = 4$$



INTRA-NODE SEARCH

Approach #1: Linear

- Scan node keys from beginning to end.
- Use SIMD to vectorize comparisons.

Approach #2: Binary

- Jump to middle key, pivot left/right depending on comparison.

Approach #3: Interpolation

- Approximate location of desired key based on known distribution of keys.

$$\text{Offset: } (8-4) \times 7 / (10-4) = 4$$



OPTIMIZATIONS

Pointer Swizzling

Buffered Updates (*B ϵ -trees*)

Partial Indexes

Include Columns

Prefix Compression

Deduplication

Suffix Truncation

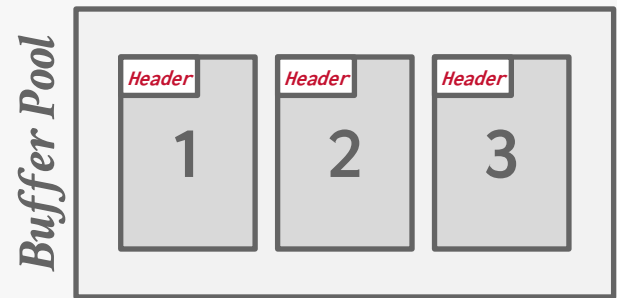
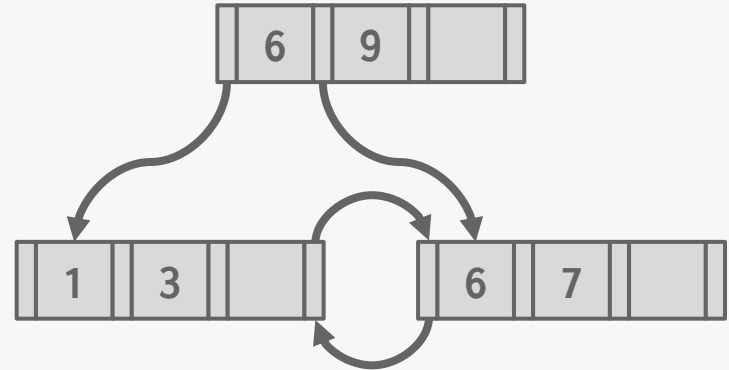
Bulk Insert

Many more...

POINTER SWIZZLING

Nodes use page ids to reference other nodes in the index. The DBMS must get the memory location from the page table during traversal.

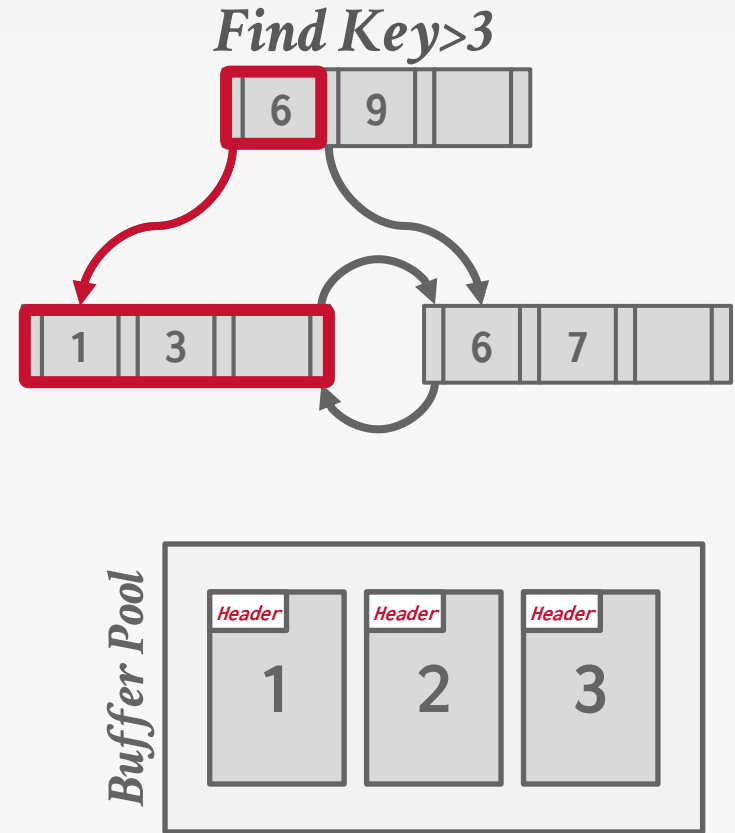
If a page is pinned in the buffer pool, then DBMS can store raw pointers instead of page ids. This avoids address lookups from the page table.



POINTER SWIZZLING

Nodes use page ids to reference other nodes in the index. The DBMS must get the memory location from the page table during traversal.

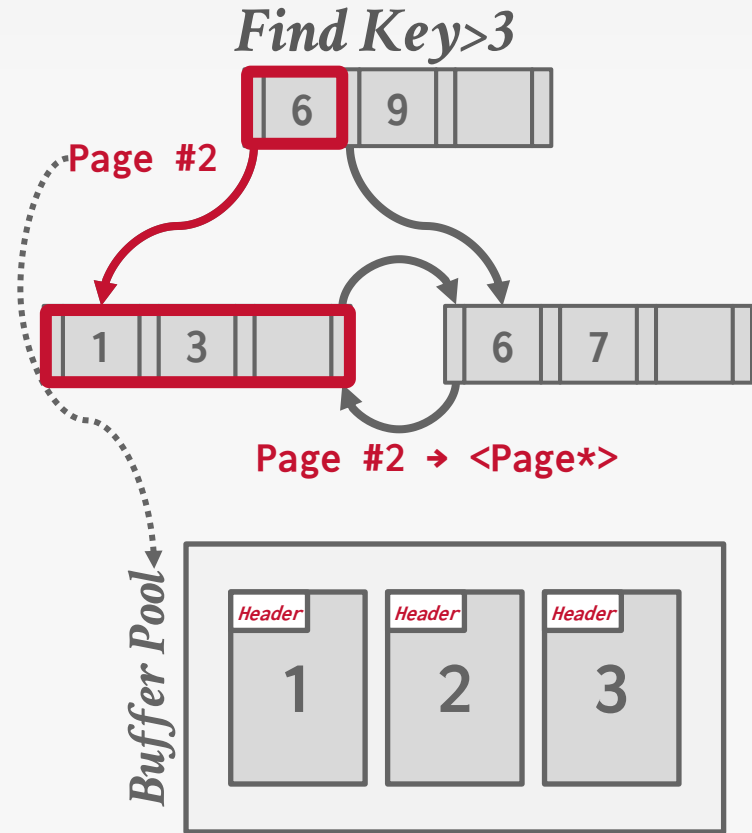
If a page is pinned in the buffer pool, then DBMS can store raw pointers instead of page ids. This avoids address lookups from the page table.



POINTER SWIZZLING

Nodes use page ids to reference other nodes in the index. The DBMS must get the memory location from the page table during traversal.

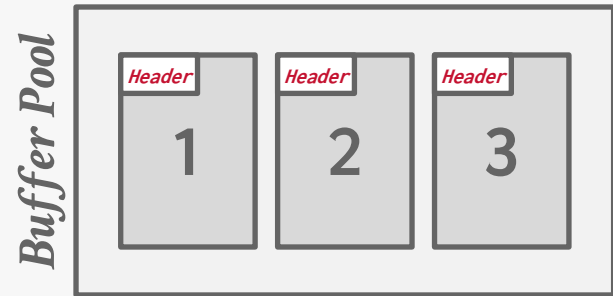
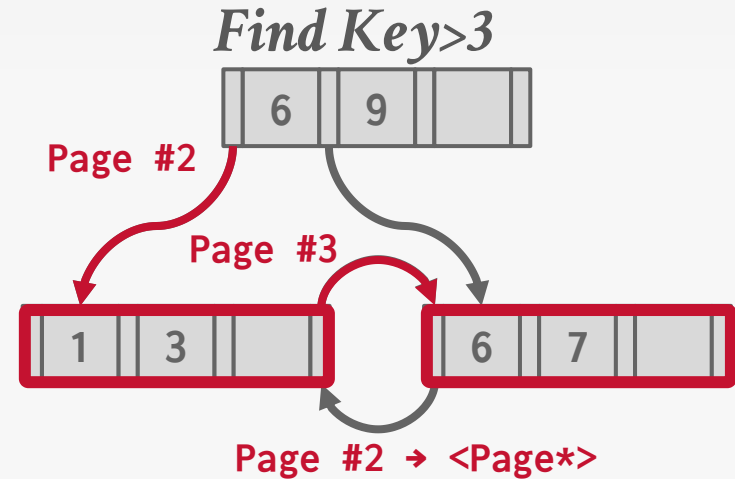
If a page is pinned in the buffer pool, then DBMS can store raw pointers instead of page ids. This avoids address lookups from the page table.



POINTER SWIZZLING

Nodes use page ids to reference other nodes in the index. The DBMS must get the memory location from the page table during traversal.

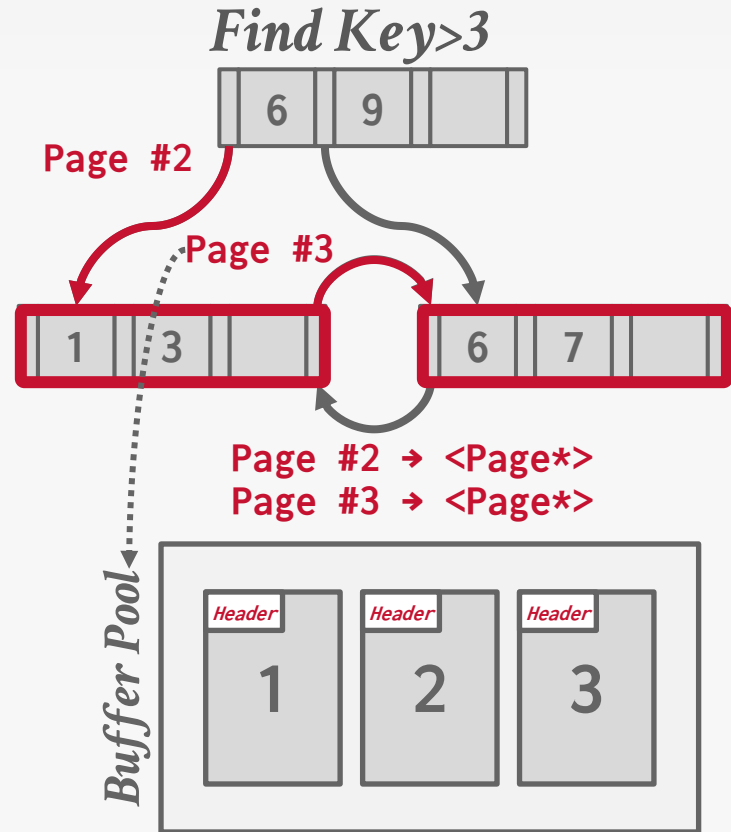
If a page is pinned in the buffer pool, then DBMS can store raw pointers instead of page ids. This avoids address lookups from the page table.



POINTER SWIZZLING

Nodes use page ids to reference other nodes in the index. The DBMS must get the memory location from the page table during traversal.

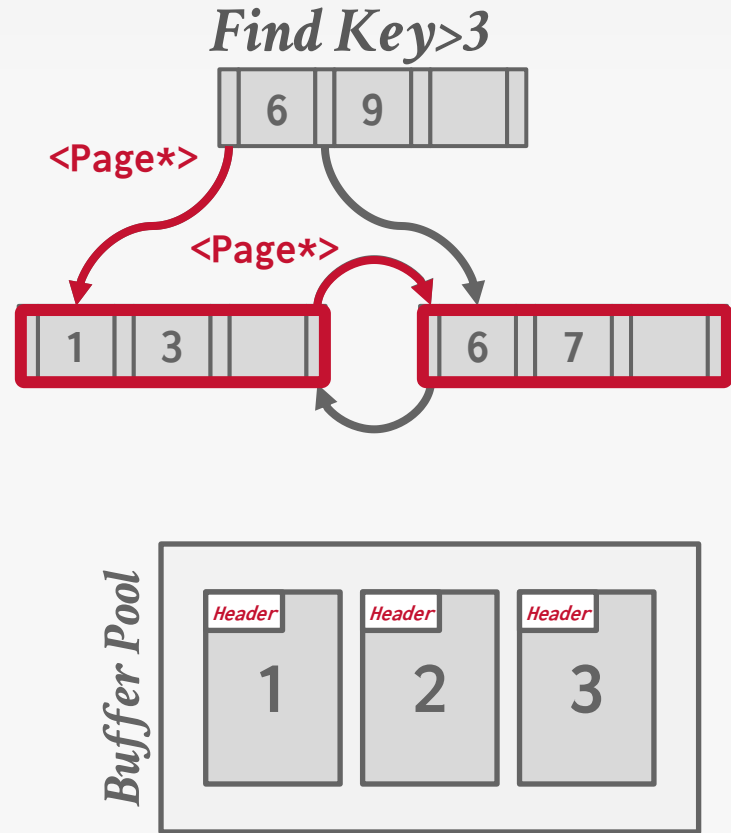
If a page is pinned in the buffer pool, then DBMS can store raw pointers instead of page ids. This avoids address lookups from the page table.



POINTER SWIZZLING

Nodes use page ids to reference other nodes in the index. The DBMS must get the memory location from the page table during traversal.

If a page is pinned in the buffer pool, then DBMS can store raw pointers instead of page ids. This avoids address lookups from the page table.



OBSERVATION

Modifying a B+tree is expensive when the DBMS has to split/merge nodes.

- Worst case is when DBMS reorganizes the entire tree.
- The worker that causes a split/merge is responsible for doing the work.

What if the DBMS can delay updates to the data structure's organization and then apply multiple changes together in a batch?

WRITE-OPTIMIZED B+TREE

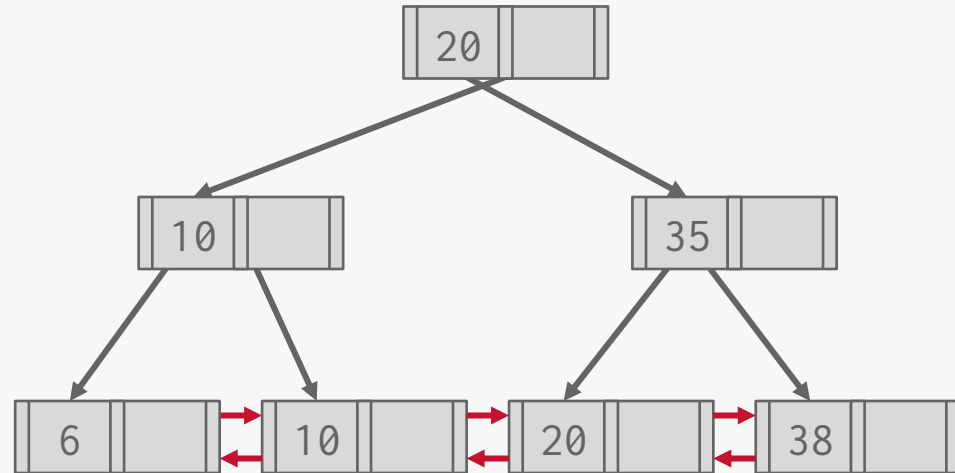
Instead of immediately applying updates, store changes to key/value entries in log buffers at inner nodes.
 → aka **Fractal Trees** / **Bε-trees**.

Updates cascade down to lower nodes incrementally when buffers get full.

Tokutek  SPLINTERDB

 **RelationalAI**  **ChromoDB**

STSDb

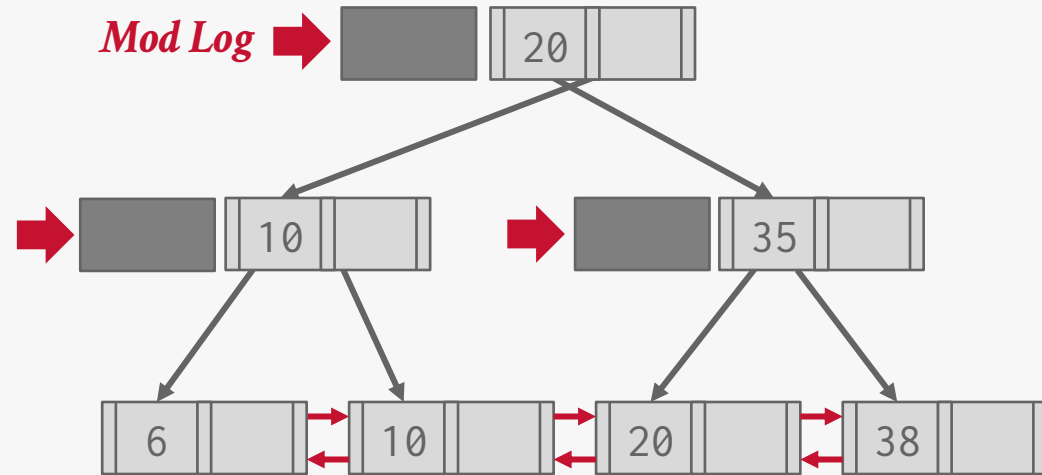


WRITE-OPTIMIZED B+TREE

Instead of immediately applying updates, store changes to key/value entries in log buffers at inner nodes.

→ aka **Fractal Trees** / **Bε-trees**.

Updates cascade down to lower nodes incrementally when buffers get full.



Tokutek  SPLINTERDB

 **RelationalAI**  **ChromoDB**

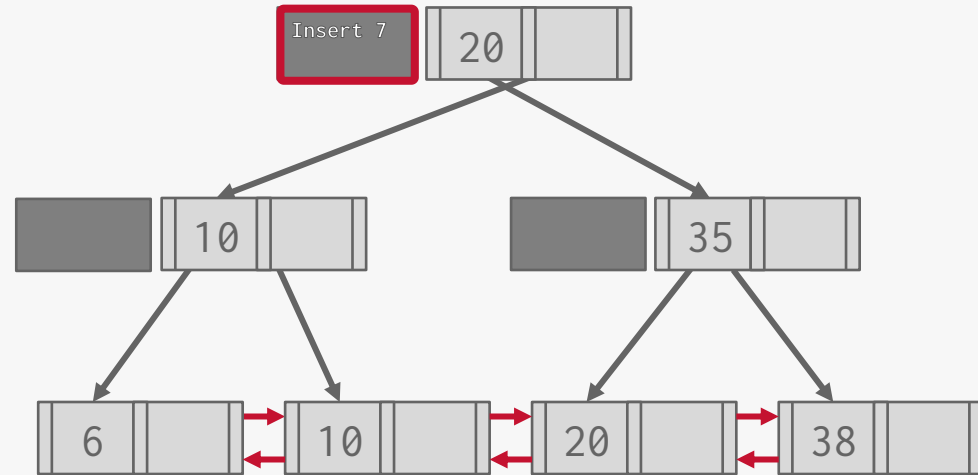
STS DB

WRITE-OPTIMIZED B+TREE

Instead of immediately applying updates, store changes to key/value entries in log buffers at inner nodes.
 → aka **Fractal Trees** / **Bε-trees**.

Updates cascade down to lower nodes incrementally when buffers get full.

Insert 7



Tokutek  SPLINTERDB

 **RelationalAI**  **ChromoDB**

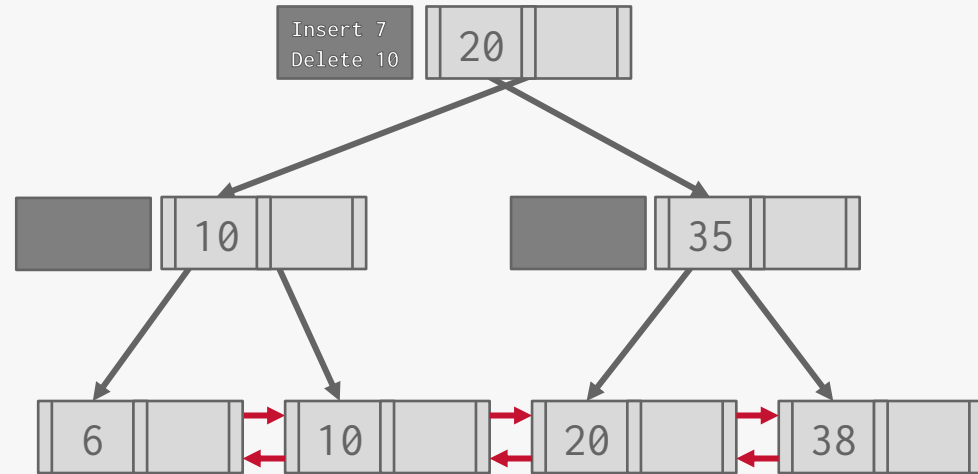
STSDb

WRITE-OPTIMIZED B+TREE

Instead of immediately applying updates, store changes to key/value entries in log buffers at inner nodes.
 → aka **Fractal Trees** / **Bε-trees**.

Updates cascade down to lower nodes incrementally when buffers get full.

Insert 7
Delete 10



Tokutek  SPLINTERDB

 **RelationalAI**  **ChromoDB**

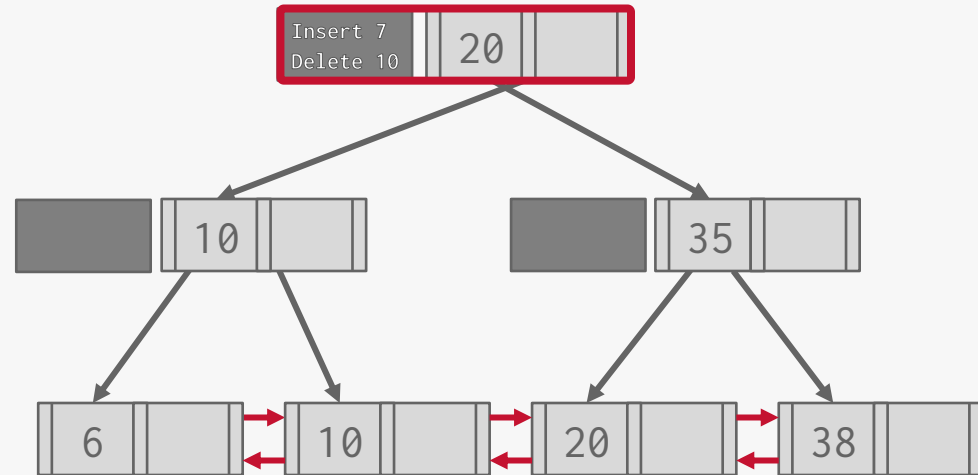
STSDb

WRITE-OPTIMIZED B+TREE

Instead of immediately applying updates, store changes to key/value entries in log buffers at inner nodes.
 → aka **Fractal Trees** / **Bε-trees**.

Updates cascade down to lower nodes incrementally when buffers get full.

Find 10



Tokutek  SPLINTERDB

 **RelationalAI**  **ChromoDB**

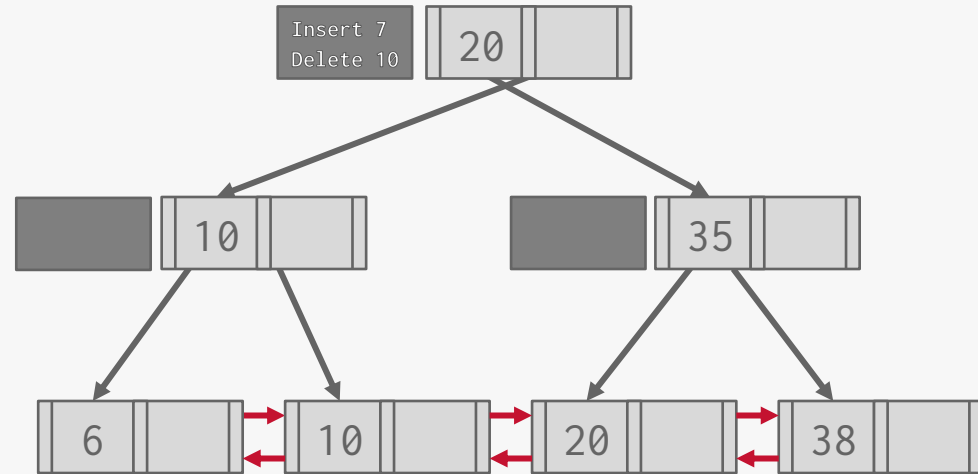
STSDb

WRITE-OPTIMIZED B+TREE

Instead of immediately applying updates, store changes to key/value entries in log buffers at inner nodes.
 → aka **Fractal Trees** / **Bε-trees**.

Updates cascade down to lower nodes incrementally when buffers get full.

Insert 40



Tokutek  SPLINTERDB

 **RelationalAI**  **ChromoDB**

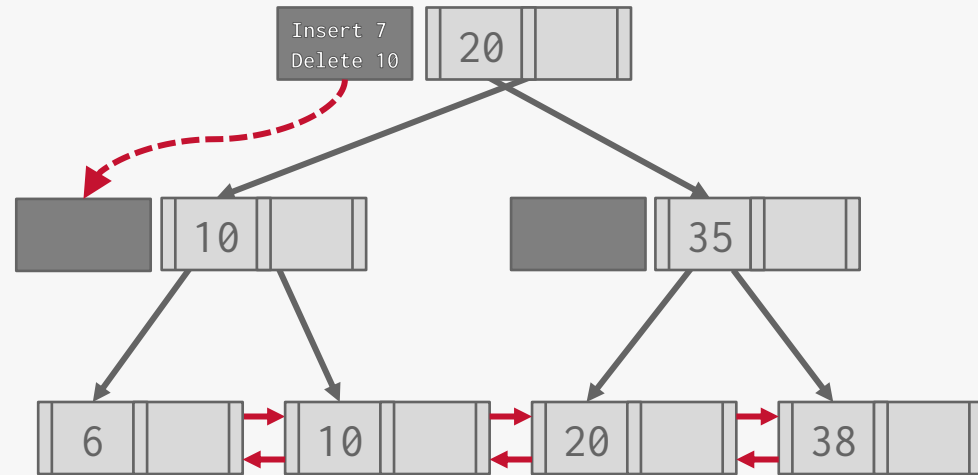
STS DB

WRITE-OPTIMIZED B+TREE

Instead of immediately applying updates, store changes to key/value entries in log buffers at inner nodes.
 → aka **Fractal Trees** / **Bε-trees**.

Updates cascade down to lower nodes incrementally when buffers get full.

Insert 40



Tokutek  SPLINTERDB

 **RelationalAI**  **ChromoDB**

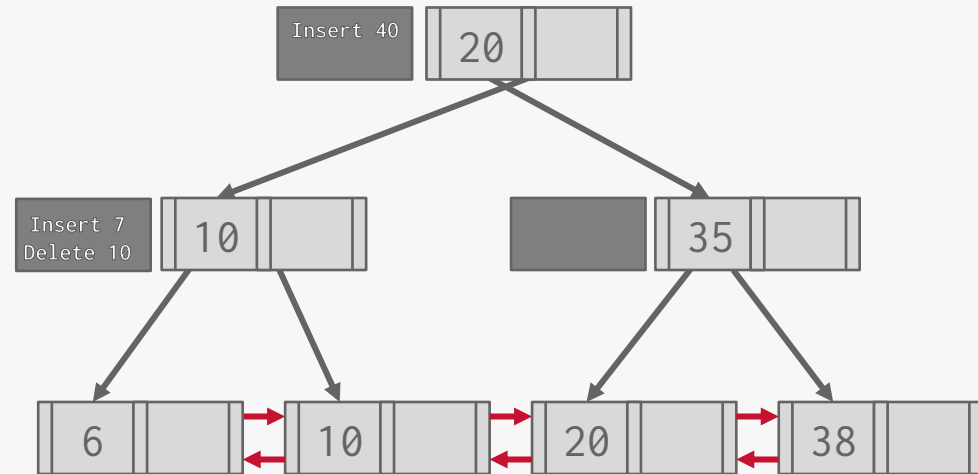
STSDb

WRITE-OPTIMIZED B+TREE

Instead of immediately applying updates, store changes to key/value entries in log buffers at inner nodes.
 → aka **Fractal Trees** / **Bε-trees**.

Updates cascade down to lower nodes incrementally when buffers get full.

Insert 40



Tokutek  SPLINTERDB

 **RelationalAI**  **ChromoDB**

STS DB

PARTIAL INDEXES

Create an index on a subset of the entire table. This potentially reduces its size and the amount of overhead to maintain it.

One common use case is to partition indexes by date ranges.
→ Create a separate index per month, year.

```
CREATE INDEX idx_foo
      ON foo (a, b)
      WHERE c = 'WuTang';
```

PARTIAL INDEXES

Create an index on a subset of the entire table. This potentially reduces its size and the amount of overhead to maintain it.

One common use case is to partition indexes by date ranges.
→ Create a separate index per month, year.

```
CREATE INDEX idx_foo
      ON foo (a, b)
      WHERE c = 'WuTang';
```

```
SELECT b FROM foo
      WHERE a = 123
      AND c = 'WuTang';
```

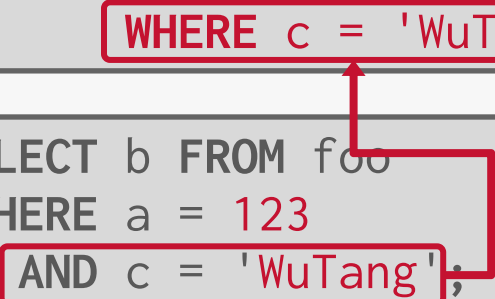
PARTIAL INDEXES

Create an index on a subset of the entire table. This potentially reduces its size and the amount of overhead to maintain it.

One common use case is to partition indexes by date ranges.
→ Create a separate index per month, year.

```
CREATE INDEX idx_foo
      ON foo (a, b)
      WHERE c = 'WuTang';
```

```
SELECT b FROM foo
      WHERE a = 123
      AND c = 'WuTang';
```




PARTIAL INDEXES


Create an index on a subset of the entire table. This potentially reduces its size and the amount of overhead to maintain it.

One common use case is to partition indexes by date ranges.
→ Create a separate index per month, year.

```
CREATE INDEX idx_foo
      ON foo (a, b)
      WHERE c = 'WuTang';
```



```
SELECT b FROM foo
      WHERE a = 123
```



INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries.

These extra columns are only stored in the leaf nodes and are not part of the search key.

```
CREATE INDEX idx_foo
      ON foo (a, b)
      INCLUDE (c);
```

INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries.

These extra columns are only stored in the leaf nodes and are not part of the search key.

```
CREATE INDEX idx_foo
      ON foo (a, b)
      INCLUDE (c)
```

INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries.

These extra columns are only stored in the leaf nodes and are not part of the search key.

```
CREATE INDEX idx_foo
      ON foo (a, b)
      INCLUDE (c);
```


```
SELECT b FROM foo
WHERE a = 123
      AND c = 'WuTang';
```

INDEX INCLUDE COLUMNS


Embed additional columns in indexes to support index-only queries.

These extra columns are only stored in the leaf nodes and are not part of the search key.

```
CREATE INDEX idx_foo  
ON foo (a, b)  
INCLUDE (c);
```



```
SELECT b FROM foo  
WHERE a = 123  
AND c = 'WuTang';
```



INDEX INCLUDE COLUMNS

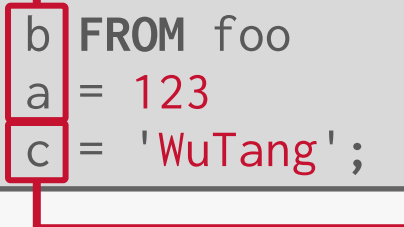
Embed additional columns in indexes to support index-only queries.

These extra columns are only stored in the leaf nodes and are not part of the search key.

```
CREATE INDEX idx_foo  
ON foo (a, b)  
INCLUDE (c);
```



```
SELECT b FROM foo  
WHERE a = 123  
AND c = 'WuTang';
```

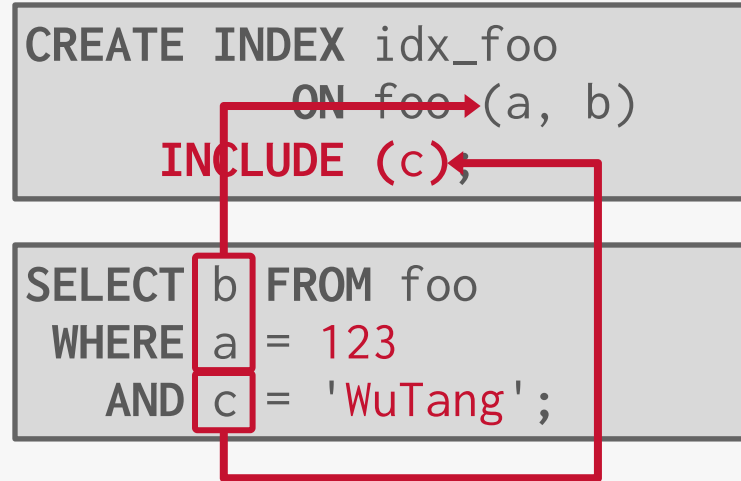


INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries.

These extra columns are only stored in the leaf nodes and are not part of the search key.

If all attributes a query needs are available in an index, then the DBMS does not need to retrieve the tuple.
→ AKA covering index or index-only scans.



PREFIX COMPRESSION

Sorted keys in the same leaf node are likely to have the same prefix.

Instead of storing the entire key each time, extract common prefix and store only unique suffix for each key.

→ Many variations.

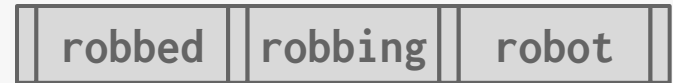
robbed	robbing	robot
--------	---------	-------

PREFIX COMPRESSION

Sorted keys in the same leaf node are likely to have the same prefix.

Instead of storing the entire key each time, extract common prefix and store only unique suffix for each key.

→ Many variations.



DEDUPLICATION

Non-unique indexes can end up storing multiple copies of the same key in leaf nodes.

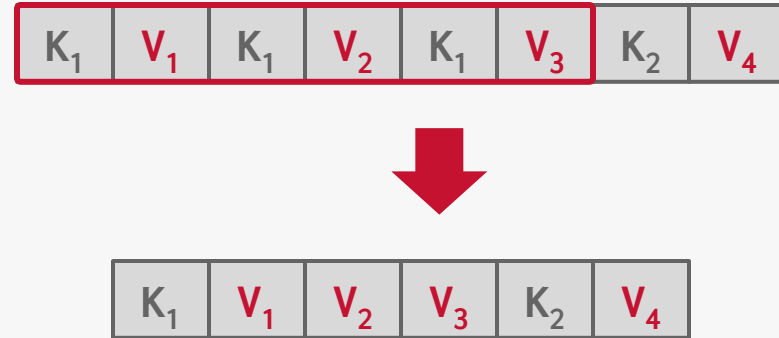
The leaf node can store the key once and then maintain a "posting list" of tuples with that key (similar to what we discussed for hash tables).

K_1	V_1	K_1	V_2	K_1	V_3	K_2	V_4
-------	-------	-------	-------	-------	-------	-------	-------

DEDUPLICATION

Non-unique indexes can end up storing multiple copies of the same key in leaf nodes.

The leaf node can store the key once and then maintain a "posting list" of tuples with that key (similar to what we discussed for hash tables).

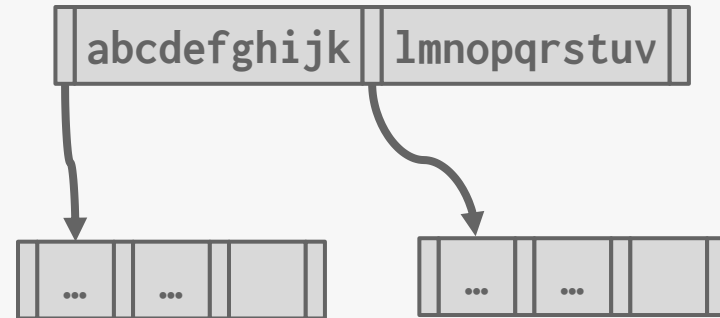


SUFFIX TRUNCATION

The keys in the inner nodes are only used to "direct traffic".

→ We don't need the entire key.

Store a minimum prefix that is needed to correctly route probes into the index.

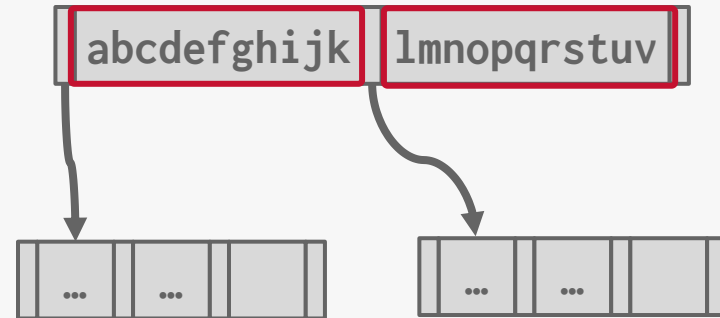


SUFFIX TRUNCATION

The keys in the inner nodes are only used to "direct traffic".

→ We don't need the entire key.

Store a minimum prefix that is needed to correctly route probes into the index.

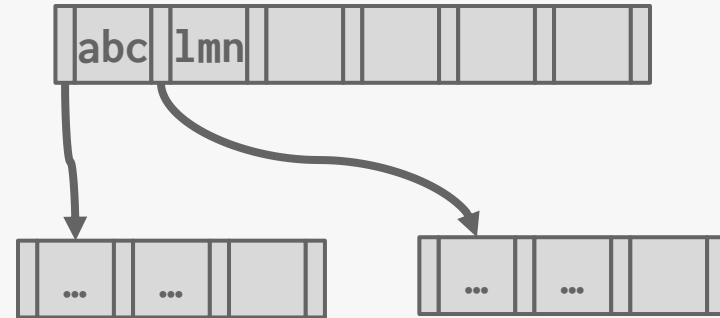


SUFFIX TRUNCATION

The keys in the inner nodes are only used to "direct traffic".

→ We don't need the entire key.

Store a minimum prefix that is needed to correctly route probes into the index.



BULK INSERT

The fastest way to build a new B+Tree for an existing table is to first sort the keys and then build the index from the bottom up.

BULK INSERT

The fastest way to build a new B+Tree for an existing table is to first sort the keys and then build the index from the bottom up.

Keys: 3, 7, 9, 13, 6, 1

Sorted Keys: 1, 3, 6, 7, 9, 13

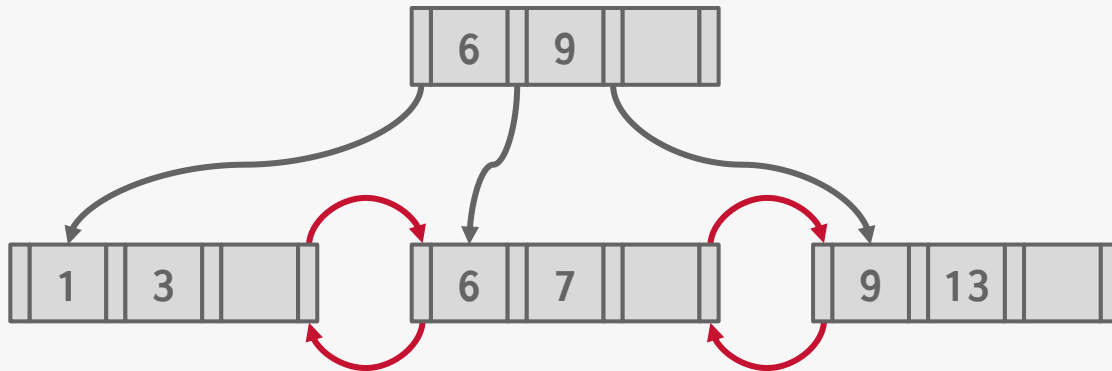


BULK INSERT

The fastest way to build a new B+Tree for an existing table is to first sort the keys and then build the index from the bottom up.

Keys: 3, 7, 9, 13, 6, 1

Sorted Keys: 1, 3, 6, 7, 9, 13



B+Tree vs. Hash Indexes

Table Clustering

CONCLUSION

The venerable B+Tree is (almost) always a good choice for your DBMS.

Supporting concurrent access is tricky. We will cover that next week.

NEXT CLASS

Bloom Filters

Tries / Radix Trees / Patricia Trees

Skip Lists

Inverted Indexes

Vector Indexes