# Database Systems I

CMPT 354 Summer 2024

Zhengjie Miao

July 10 2024

# Announcements (Wed. July 10)

- Update on Exam I solution
  - Grades will also be updated
- Makeup OH slots
  - Thu. July 11, 12:45pm - 1:30pm @TASC I 9407
  - Fri. July 12, 11:15am - 12:00pm @TASC I 9407

# Examples of using indexes

- `SELECT * FROM User WHERE name = 'Bart';`

- How is the query processed?

- Without an index on *User.name:* must scan the entire table if we store *User* as a flat file of unordered rows

- With index: go "directly" to rows with name=`'Bart'`

# Examples of using indexes

- ```
  SELECT * FROM User, Member
  WHERE User.uid = Member.uid AND Member.gid = 'cks';
  ```
- How to find relevant *Member* rows directly?
  - With an index on *Member.gid* or (*gid, uid*):
- For each relevant *Member* row, how to directly look up *User* rows with matching *uid*?
  - With an index on *User.uid*
  - Without it: for each *Member* row, scan the entire *User* table for matching *uid*
    - Sorting could help

# Indexes

- An index is an auxiliary persistent data structure
  - Search tree (e.g., B$^+$-tree), lookup table (e.g., hash table), etc.
- Creating and dropping indexes in SQL:
  - CREATE [UNIQUE] INDEX $indexname$ ON $tablename(columnname_1, \ldots, columnname_n)$;
    - With UNIQUE, the DBMS will also enforce that $\{columnname_1, \ldots, columnname_n\}$ is a key of $tablename$
  - DROP INDEX $indexname$;

  You will not have to create indices on these columns after using the unique keyword.

  - Typically, the DBMS will automatically create indexes for PRIMARY KEY and UNIQUE constraint declarations
- Can have many indexes for one table

# Indexes

- An index on $R.A$ can speed up accesses of the form
  - $R.A = value$
  - $R.A > value$ (sometimes; depending on the index type)
- An index on $(R.A_1, \dots, R.A_n)$ can speed up
  - $R.A_1 = value_1 \wedge \cdots \wedge R.A_n = value_n$
  - $(R.A_1, \dots, R.A_n) > (value_1, \dots, value_n)$ (again depends)
- Ordering of index columns is important—is an index on $(R.A, R.B)$ equivalent to one on $(R.B, R.A)$?
- How about an index on $R.A$ plus another on $R.B$?

For the first case, if you access the A column it can be obtained easily. However, if you access the B column, accessing the data will not be as fast.

Ask professor to go over the difference of using separate vs joined indices

If you are checking that R.A equals some value and R.B equals another value, then you can have separate indices.

# Choosing indexes to create

More indexes = better performance?

- Indexes take space

- Indexes need to be maintained when data is updated

- Indexes have one more level of indirection


- Optimal index selection depends on both query and update workload and the size of tables

  - Automatic index selection is now featured in some commercial DBMS

Automatic index problem is common

# Choosing indexes to create

- Make some attribute K a search key if the WHERE clause contains:
  - An exact match on K
  - A range predicate on K
  - A join on K

# The index selection problem 1

- Your workload is

Since name will have less duplication, the indices will serve more for optimizing the query

Age will have at most 100 values. So having a index of that would not be as beneficial compared to name

100000 queries

```
SELECT uid
FROM User
WHERE name = ?
```

100000 queries

```
SELECT uid
FROM User
WHERE age = ?
```

**Which one is better?**

A. Index on name

B. Index on age

# The index selection problem 2

- Your workload is

100000 queries

```
SELECT uid
FROM User
WHERE name = ?
```

100000 queries

```
SELECT uid
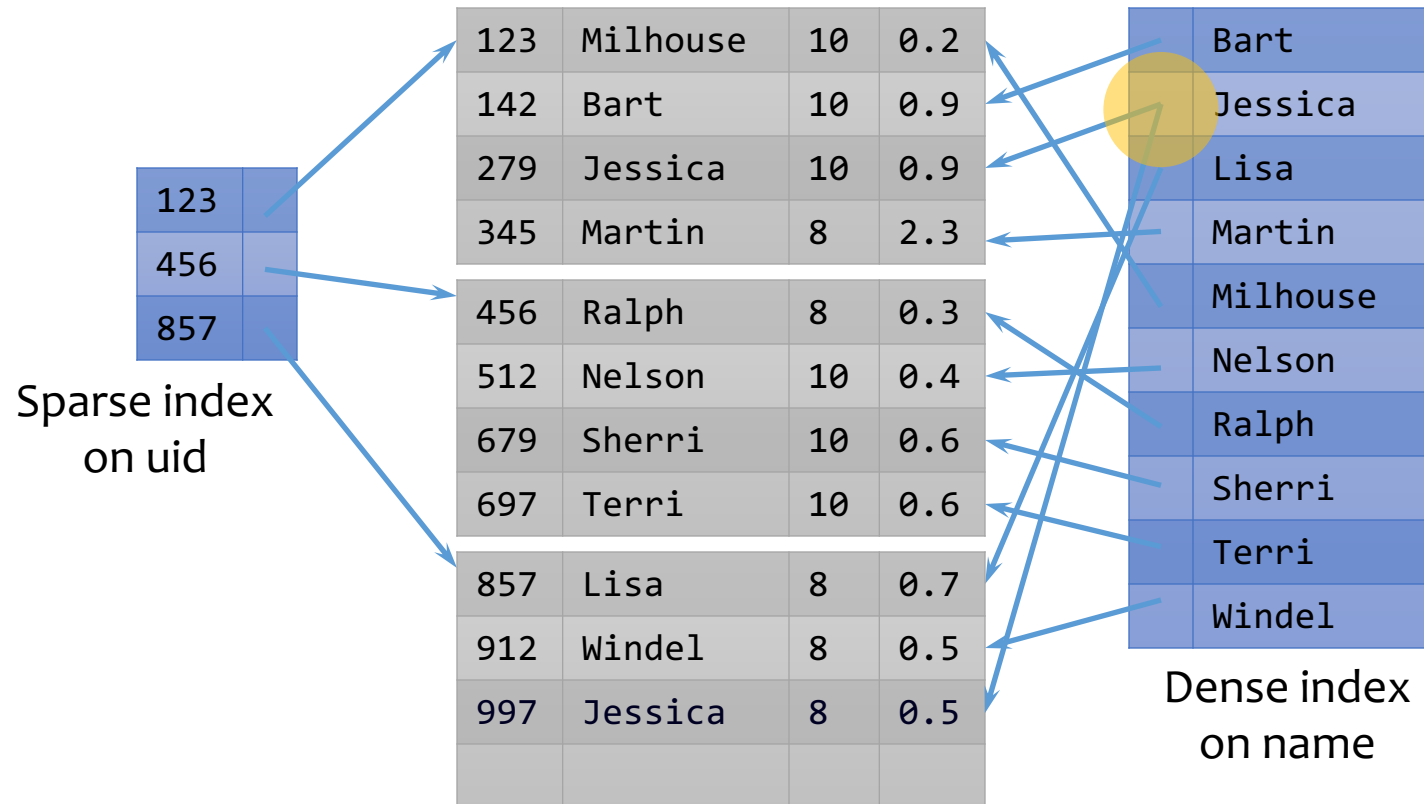FROM User
WHERE name = ? AND age > ?
```

For the first query, the first index is better to allow you quick access to the names

**Which one is better?**

A. Index on (name, age)

B. Index on (age, name)

# Dense and sparse indexes

- **Dense**: one index entry for each search key value
  - One entry may "point" to multiple records (e.g., two users named Jessica)
- **Sparse**: one index entry for each block
  - Records must be **clustered** according to the search key

| | | | |
|---|---|---|---|
| 123 | Milhouse | 10 | 0.2 |
| 142 | Bart | 10 | 0.9 |
| 279 | Jessica | 10 | 0.9 |
| 345 | Martin | 8 | 2.3 |

| | | | |
|---|---|---|---|
| 456 | Ralph | 8 | 0.3 |
| 512 | Nelson | 10 | 0.4 |
| 679 | Sherri | 10 | 0.6 |
| 697 | Terri | 10 | 0.6 |

| | | | |
|---|---|---|---|
| 857 | Lisa | 8 | 0.7 |
| 912 | Windel | 8 | 0.5 |
| 997 | Jessica | 8 | 0.5 |
| | | | |

Sparse index on uid:

| |
|---|
| 123 |
| 456 |
| 857 |

Sparse index on uid

Dense index on name:

| |
|---|
| Bart |
| Jessica |
| Lisa |
| Martin |
| Milhouse |
| Nelson |
| Ralph |
| Sherri |
| Terri |
| Windel |

Dense index on name

# Dense versus sparse indexes

- Index size
  - Sparse index is smaller    One entry corresponds to one block

- Requirement on records
  - Records must be clustered for sparse index

- Lookup
  - Sparse index is smaller and may fit in memory
  - Dense index can directly tell if a record exists

- Update
  - Easier for sparse index

# Primary and secondary indexes

- ## Primary index
  - Created for the primary key of a table
  - Records are usually clustered by the primary key
  - Can be sparse

- ## Secondary index
  - Usually dense

- In SQL
  - PRIMARY KEY declaration automatically creates a primary index, UNIQUE key automatically creates a secondary index
  - Additional secondary index can be created on non-key attribute(s) too
    - `CREATE INDEX UserPopIndex ON User(pop);`

# What if the index is too big as well?

- Put a another (sparse) index on top of that!

| 123 | Milhouse | 10 | 0.2 |
|-----|----------|----|-----|
| 142 | Bart     | 10 | 0.9 |
| 279 | Jessica  | 10 | 0.9 |
| 345 | Martin   | 8  | 2.3 |

| 456 | Ralph  | 8  | 0.3 |
|-----|--------|----|-----|
| 512 | Nelson | 10 | 0.4 |
| 679 | Sherri | 10 | 0.6 |
| 697 | Terri  | 10 | 0.6 |

| 857 | Lisa    | 8 | 0.7 |
|-----|---------|---|-----|
| 912 | Windel  | 8 | 0.5 |
| 997 | Jessica | 8 | 0.5 |
|     |         |   |     |

| 123 |
|-----|
| 456 |
| 857 |

Sparse index
on uid

| Bart     |
|----------|
| Jessica  |
| Lisa     |
| Martin   |
| Milhouse |
| Nelson   |
| Ralph    |
| Sherri   |
| Terri    |
| Windel   |

Dense index
on name

14

# ISAM

- What if an index is still too big?
  - Put a another (sparse) index on top of that!
  - ISAM (Index Sequential Access Method), more or less

Example: look up 197



Index blocks

```
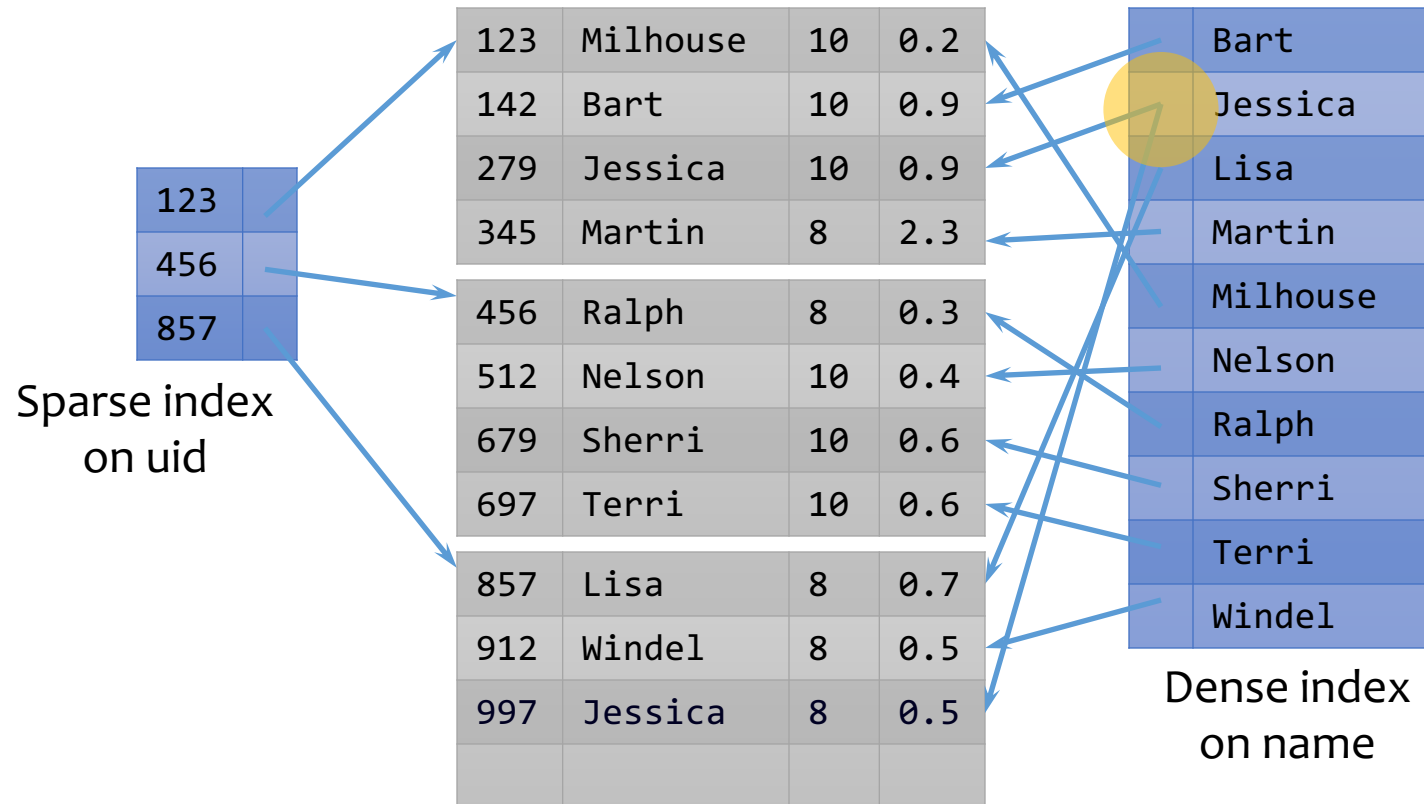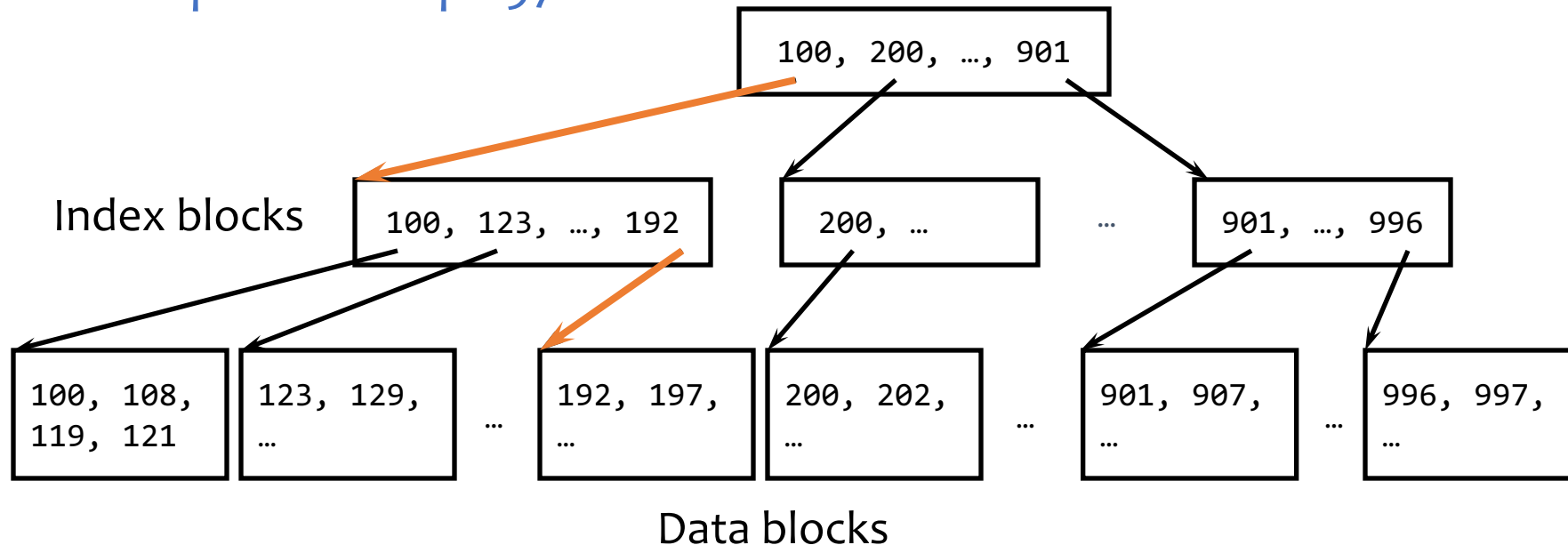100, 200, …, 901
```

```
100, 123, …, 192        200, …        …        901, …, 996
```

```
100, 108,    123, 129,   …   192, 197,   200, 202,   …   901, 907,   …   996, 997,
119, 121     …                …          …               …               …
```

Data blocks

# Updates with ISAM

Example: insert 107
Example: delete 129

100, 200, …, 901

Index blocks

100, 123, …, 192

200, …

…

901, …, 996

100, 108, 119, 121

123, 129, …

…

192, 197, …

200, 202, …

…

901, 907, …

…

996, 997, …

107   Overflow block

Data blocks

- Overflow chains and empty data blocks degrade performance
  - Worst case: most records go into one long chain, so lookups require scanning all data!

# B⁺-tree

- A hierarchy of nodes with intervals
- Balanced (more or less): good performance guarantee
- Disk-based: one node per block; large fan-out



Max fan-out: 4

# Sample B⁺-tree nodes

to keys
$100 \leq k$

Max fan-out: 4

Index Nodes Containing
Index entries

Non-leaf

| 120 | 150 | 180 |

to keys
$100 \leq k < 120$

to keys
$120 \leq k < 150$

to keys
$150 \leq k < 180$

to keys
$180 \leq k$

Leaves are linked

Leaf

| 120 | 130 |

Pointers to next leaf node in sequence

Pointers to records with these $k$ values;
or, store records directly in leaves

# B⁺-tree balancing properties

- Height constraint: all leaves at the same lowest level
- Fan-out constraint: all nodes at least half full  (except root)

|  | Max # pointers | Max # keys | Min # active pointers | Min # keys |
|---|---|---|---|---|
| Non-leaf | $f$ | $f - 1$ | $\lceil f/2 \rceil$ | $\lceil f/2 \rceil - 1$ |
| Root | $f$ | $f - 1$ | 2 | 1 |
| Leaf | $f$ | $f - 1$ | $\lfloor f/2 \rfloor$ | $\lfloor f/2 \rfloor$ |

# Lookups

- SELECT * FROM R WHERE k = 179;
- SELECT * FROM R WHERE k = 32;



Max fan-out: 4

Not found

# Range query

- SELECT * FROM R WHERE k > 32 AND k < 179;



Max fan-out: 4

100

30

Look up 32...

120 150 180

3 5 11

30 35

100 101 110

120 130

150 156 179

180 200

And follow next-leaf pointers until you hit upper bound

# Insertion

- Insert a record with search key value 32



Max fan-out: 4

Look up where the inserted key should go...

And insert it right there

# Another insertion example

- Insert a record with search key value 152



Max fan-out: 4

Oops, node is already full!

# Node splitting



Max fan-out: 4

Oops, that node becomes full!

Need to add to parent node a pointer to the newly created node

1. we "copy up" while splitting leaves – Insertion both at leaf and parent
2. the value inserted at parent may *not* be the new value we are inserting

# More node splitting

Max fan-out: 4

Need to add to parent node a pointer to the newly created node

- In the worst case, node splitting can "propagate" all the way up to the root of the tree (not illustrated here)
  - Splitting the root introduces a new root of fan-out 2 and causes the tree to grow "up" by one level

25

# Performance analysis

- How many I/O's are required for each operation?
  - $h$, the <span style="color:red">height of the tree</span> (more or less)
  - Plus one or two to manipulate actual records
  - Plus $O(h)$ for reorganization (rare if $f$ is large)
  - Minus one if we cache the root in memory
- How big is $h$?
  - Roughly <span style="color:red">$\log_{\text{fanout}} N$</span>, where $N$ is the number of records
  - B[+]-tree properties guarantee that fan-out is least $f/2$ for all non-root nodes
  - Fan-out is typically large (in hundreds)—many keys and pointers can fit into one block
  - A 4-level B[+]-tree is enough for "typical" tables

# B⁺-tree in practice

- Complex reorganization for deletion often is not implemented (e.g., Oracle)
  - Leave nodes less than half full and periodically reorganize
- Most commercial DBMS use B⁺-tree instead of hashing-based indexes because B⁺-tree handles range queries
  - A key difference between hash and tree indexes!

# The Halloween Problem

- Story from the early days of System R…

```
UPDATE Payroll
SET salary = salary * 1.1
WHERE salary >= 100000;
```

  - There is a B$^+$-tree index on *Payroll*(*salary*)
  - The update never stopped (why?)

- Solutions?
  - Scan index in reverse, or
  - Before update, scan index to create a "to-do" list, or
  - During update, maintain a "done" list, or
  - Tag every row with transaction/statement id

# B$^+$-tree versus ISAM

- ISAM is more <span style="color:red">static</span>; B$^+$-tree is more <span style="color:red">dynamic</span>
- ISAM can be more compact (at least initially)
  - Fewer levels and I/O's than B$^+$-tree
- Overtime, ISAM may not be balanced
  - Cannot provide guaranteed performance as B$^+$-tree does

# B$^+$-tree versus B-tree

- B-tree: why not store records (or record pointers) in non-leaf nodes?
  - These records can be accessed with fewer I/O's
- Problems?
  - Storing more data in a node decreases fan-out and increases $h$
  - Records in leaves require more I/O's to access
  - Vast majority of the records live in leaves!

# Beyond ISAM, B-, and B$^+$-trees (skip)

- Other tree-based indexes: R-trees, GiST, etc.
  - How about binary tree?



vs.

- Hashing-based indexes: extensible hashing, linear hashing, etc.
- Text indexes: inverted-list index, suffix arrays, etc.
- Other tricks: bitmap index, bit-sliced index, vector database index, etc.

# Clustered vs. Unclustered Index

Clustered

Unclustered



How does it affect # of page accesses?
Recall that for a disk with block access, sequential IO is much faster than random IO

Data is sorted on search key        Data can be in any order

# Clustered vs. Unclustered Index

- For range search over n values:
  - 1 random IO + n sequential IO vs. n random IO

- SELECT * FROM USER WHERE age = 50
  - Assume 12 users with age = 50
  - Assume one data page can hold 4 User tuples
  - Suppose searching for a data entry requires 3 IOs in a B+-tree, which contain pointers to the data records (assume all matching pointers = data entries are in the same node of B+-tree)
  - What happens if the index is unclustered? (cost = 3+12)
  - What happens if the index is clustered? (cost <= 3 +(3 +1))

# Deletion (skip)

- Delete a record with search key value 130

Max fan-out: 4

Look up the key to be deleted…

If a sibling has more than enough keys, steal one!

And delete it

Oops, node is too empty!

# Stealing from a sibling (skip)



Max fan-out: 4

Remember to fix the key
in the least common ancestor
of the affected nodes

# Another deletion example (skip)

- Delete a record with search key value 179

Max fan-out: 4



Cannot steal from siblings

Then coalesce (merge) with a sibling!

# Coalescing (skip)



Max fan-out: 4

Remember to delete the appropriate key from parent

- Deletion can "propagate" all the way up to the root of the tree (not illustrated here)
  - When the root becomes empty, the tree "shrinks" by one level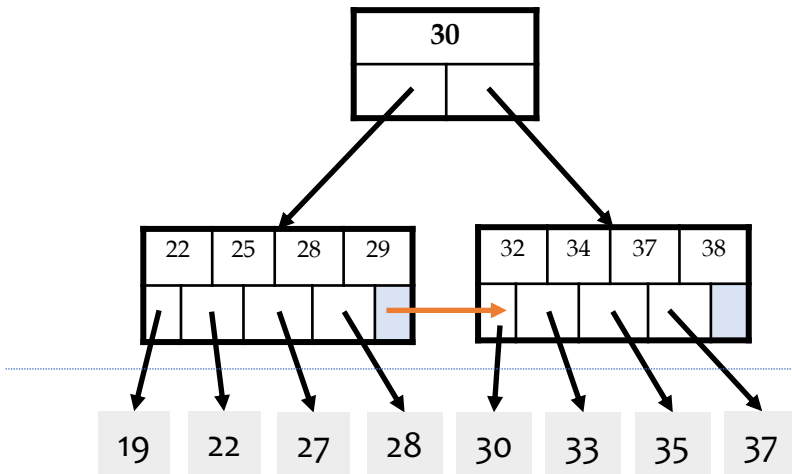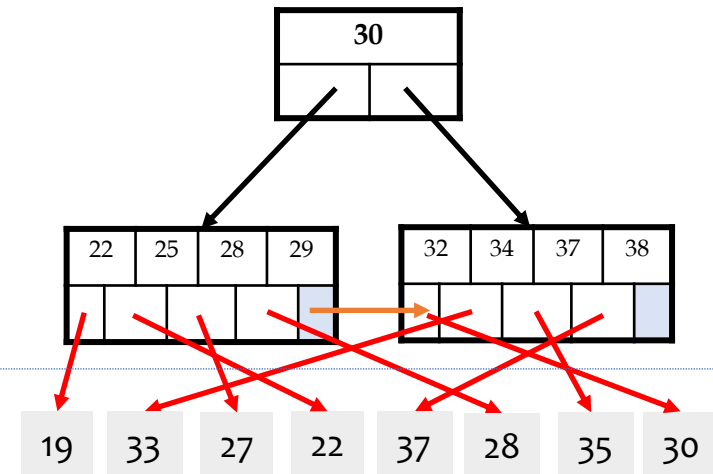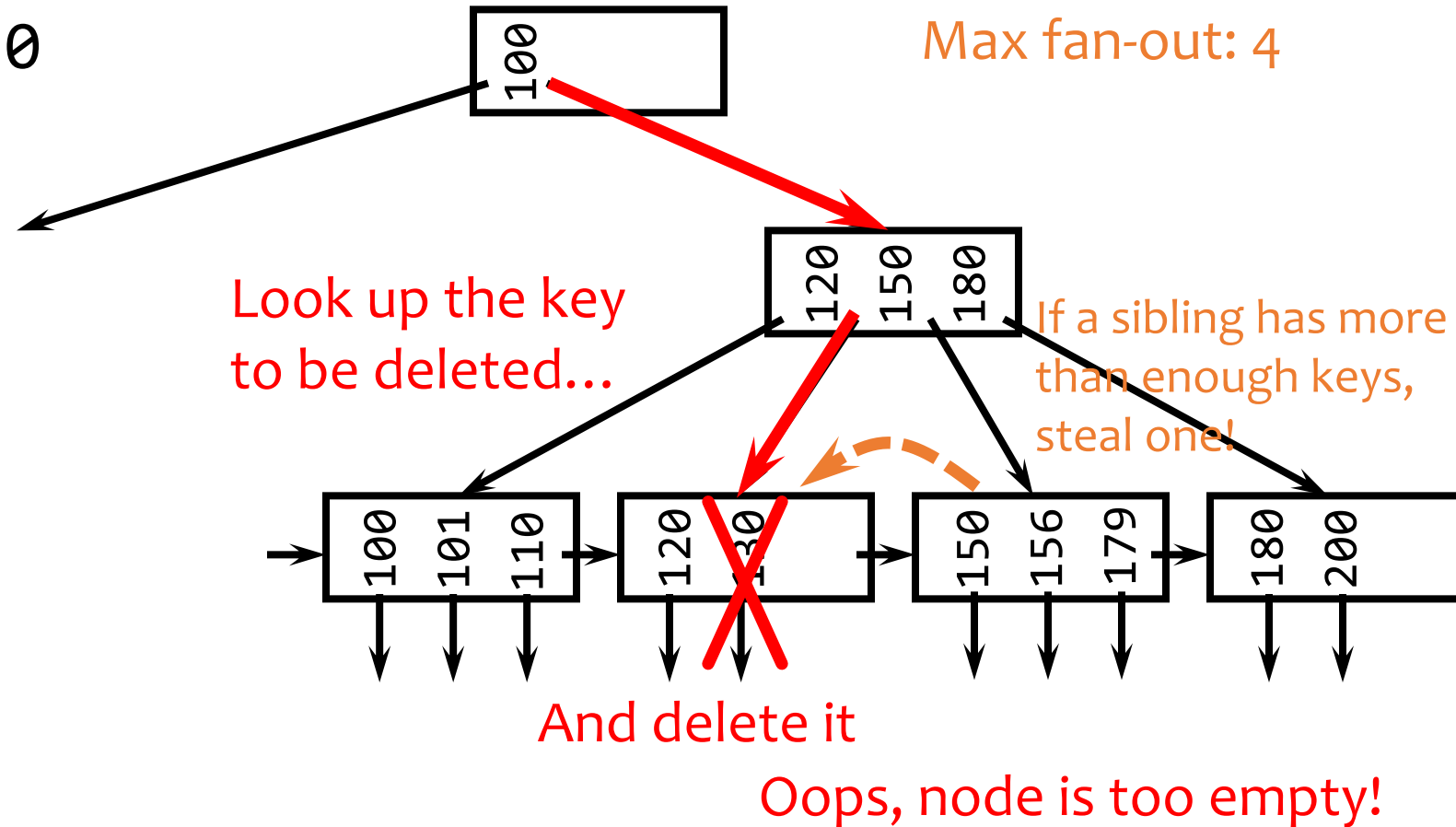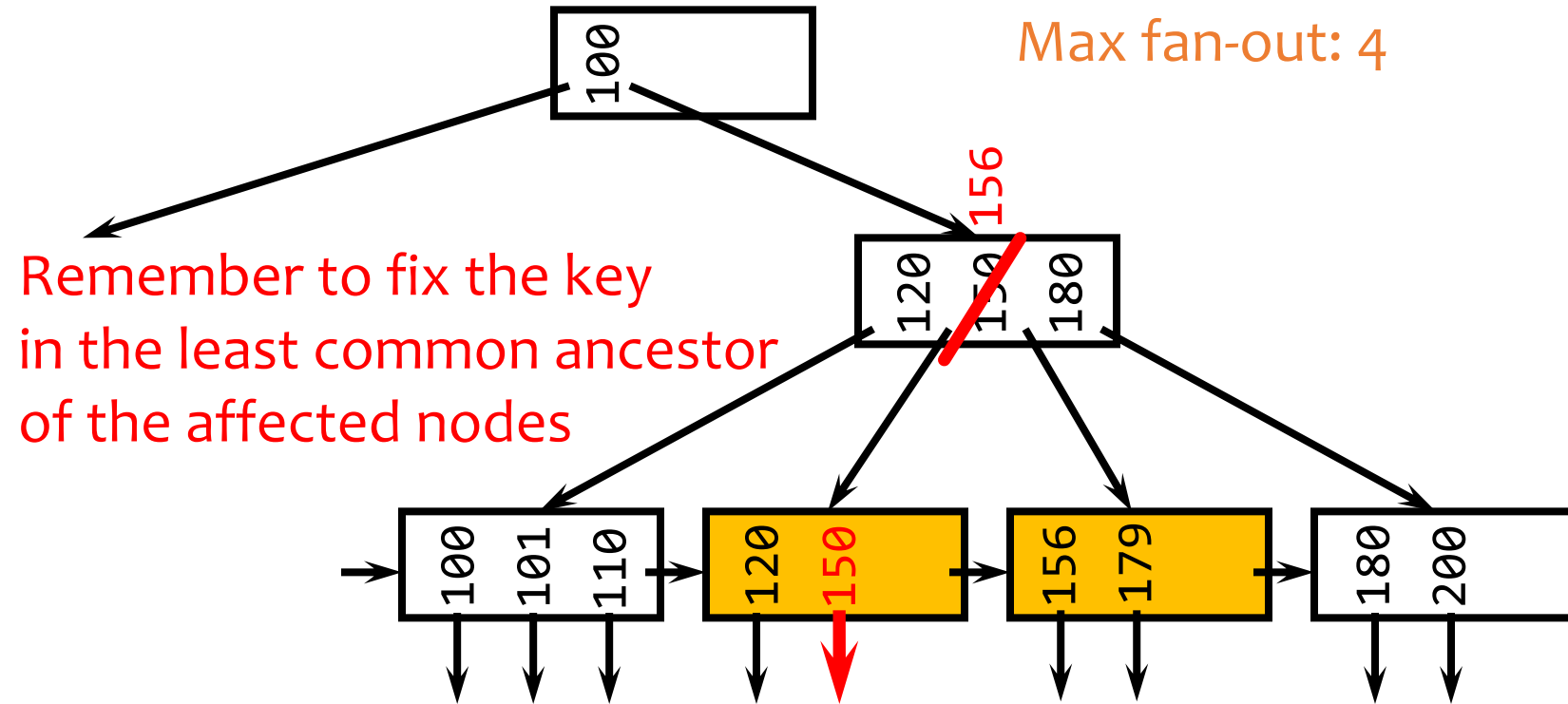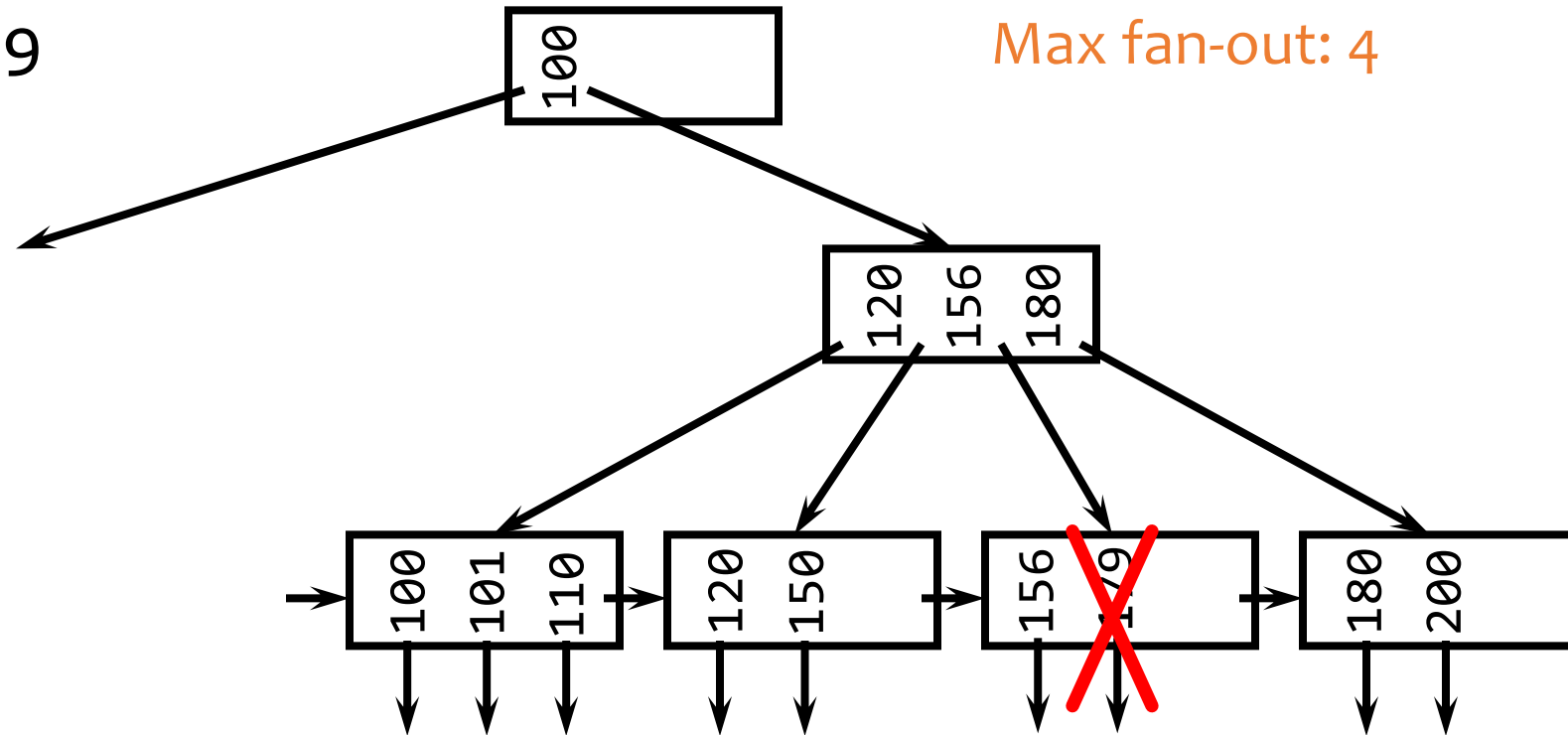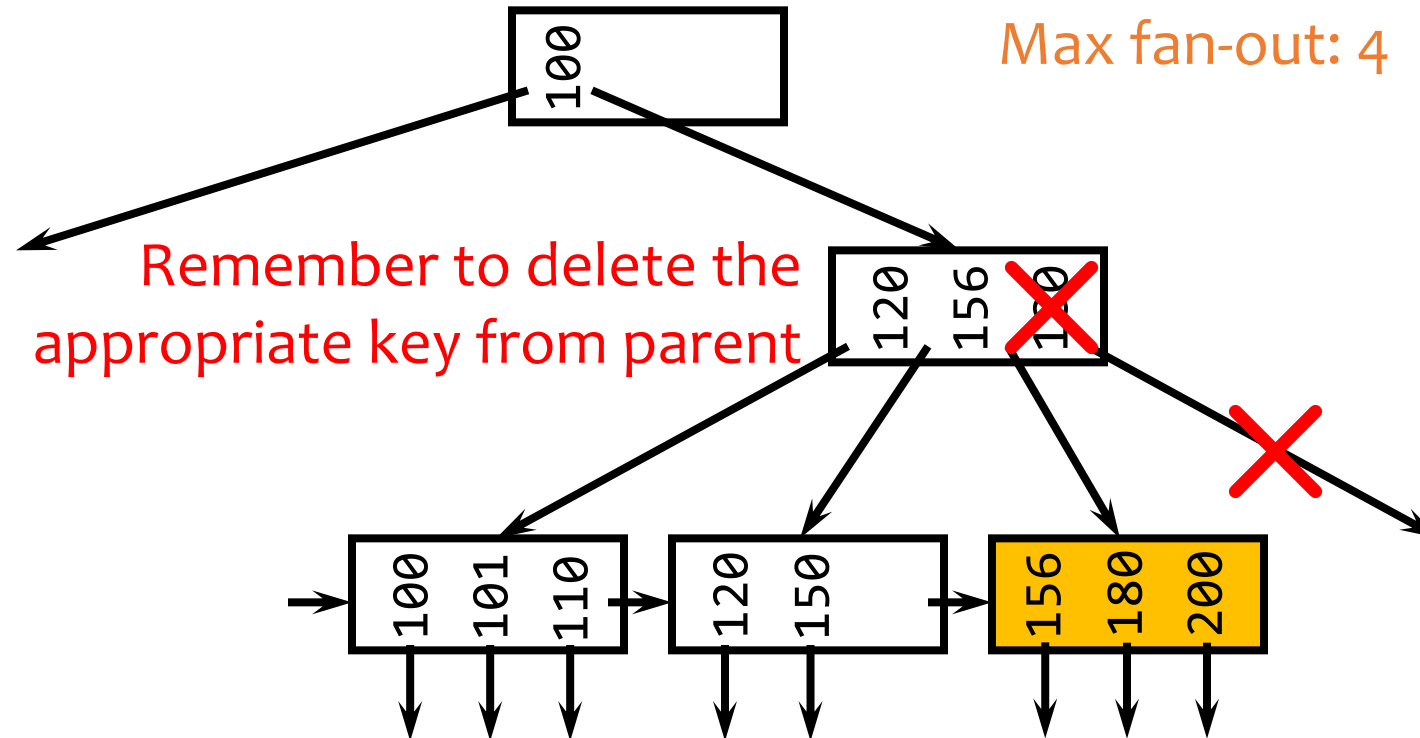