# Indexing

Hector Garcia-Molina

Mahmoud Sakr

# Indexing

Query → Index → Block holding records → candidate records

*we have a user query that valt to access*

*contains data*

*an other way to access the block*

*it's organizing the data of the block and quickly searching for the records*

*returning the candidate records of the query*

# Topics

- Conventional indexes
- B-trees
- Hashing schemes

*Building blocks*

# Sequential File

| 10 | |
| --- | --- |
| 20 | |

| 30 | |
| --- | --- |
| 40 | |

| 50 | |
| --- | --- |
| 60 | |

| 70 | |
| --- | --- |
| 80 | |

| 90 | |
| --- | --- |
| 100 | |

## Dense Index

*sorted by the key*

## Sequential File

Dense Index = a pointer per key

↪ *sequential scan of the index*

How to search for a key= 30 ?

How to search for a key= 25 ?

Can we use a dense index on a non-sequential file ? *Yes*

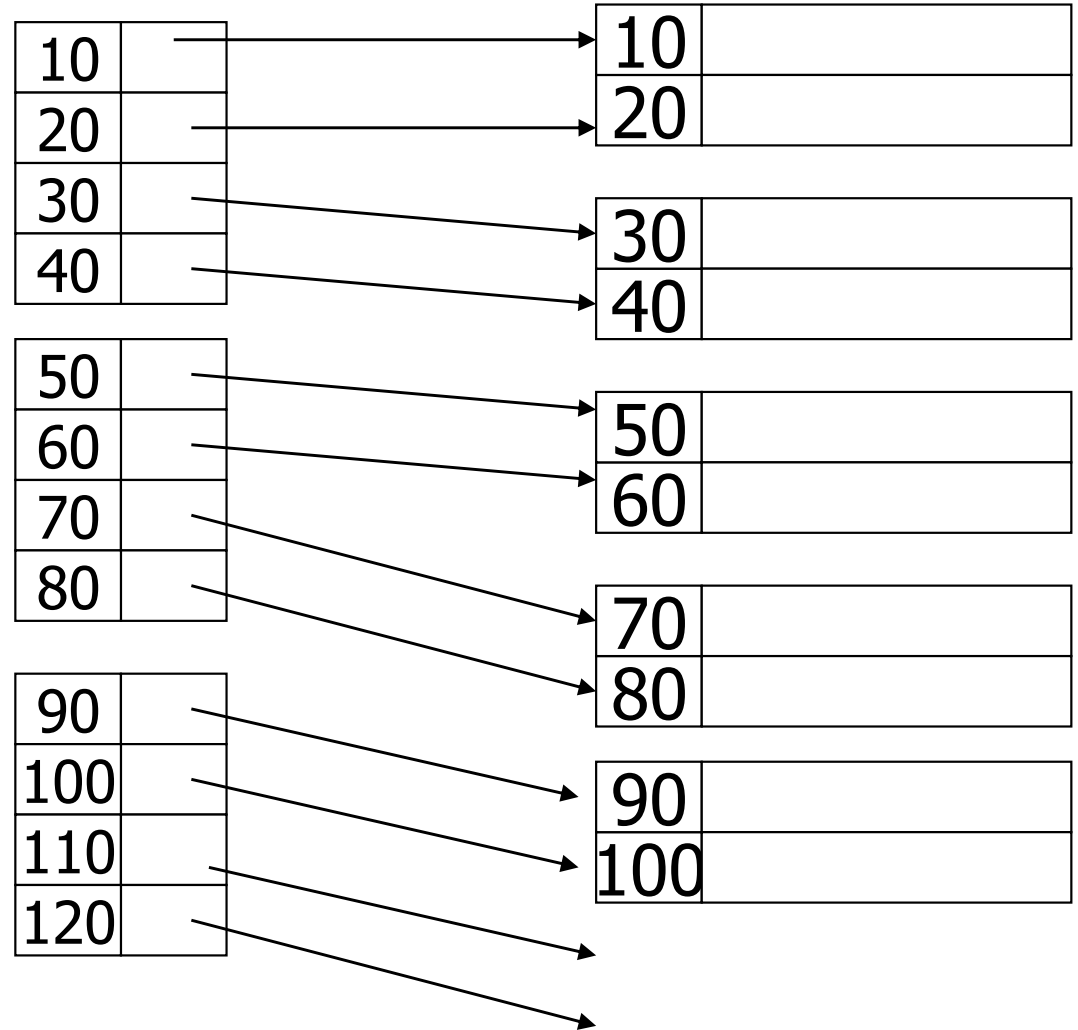↳ *we have a pointer that lead to the records*

Why querying a dense index is more efficient than querying the sequential file ?

→ *allow us to have a different order than the storage*

⟹ *we can sort the data in the index based on what we are looking for*

↳ *search through a sorted list is way faster*

*Binary search*

| 10 | |
| 20 | |
| 30 | |
| 40 | |

| 50 | |
| 60 | |
| 70 | |
| 80 | |

| 90 | |
| 100 | |
| 110 | |
| 120 | |

| 10 | |
| 20 | |

| 30 | |
| 40 | |

| 50 | |
| 60 | |

| 70 | |
| 80 | |

| 90 | |
| 100 | |

## Sparse Index

*we are pointing blocks, not keys*

## Sequential File

Sparse index = a pointer per block
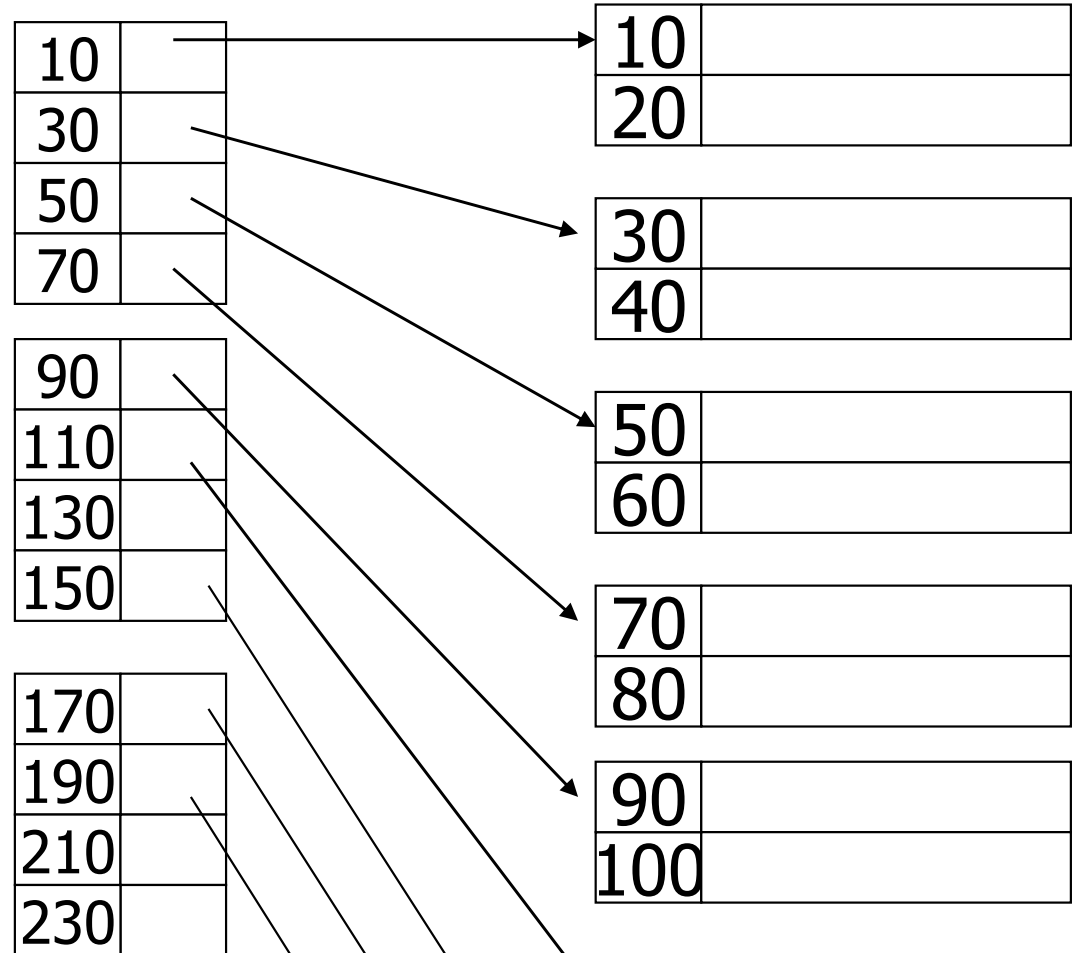↳ *less index entries*

How to search for a key= 30 ?

How to search for a key= 25 ?

Can we use a sparse index on a non-sequential file ? *No*

↳ *for the sparse index to work, we need to be sure that the pages are sorted*
↓
*if we are looking for 20 it must be between 10 and 30*
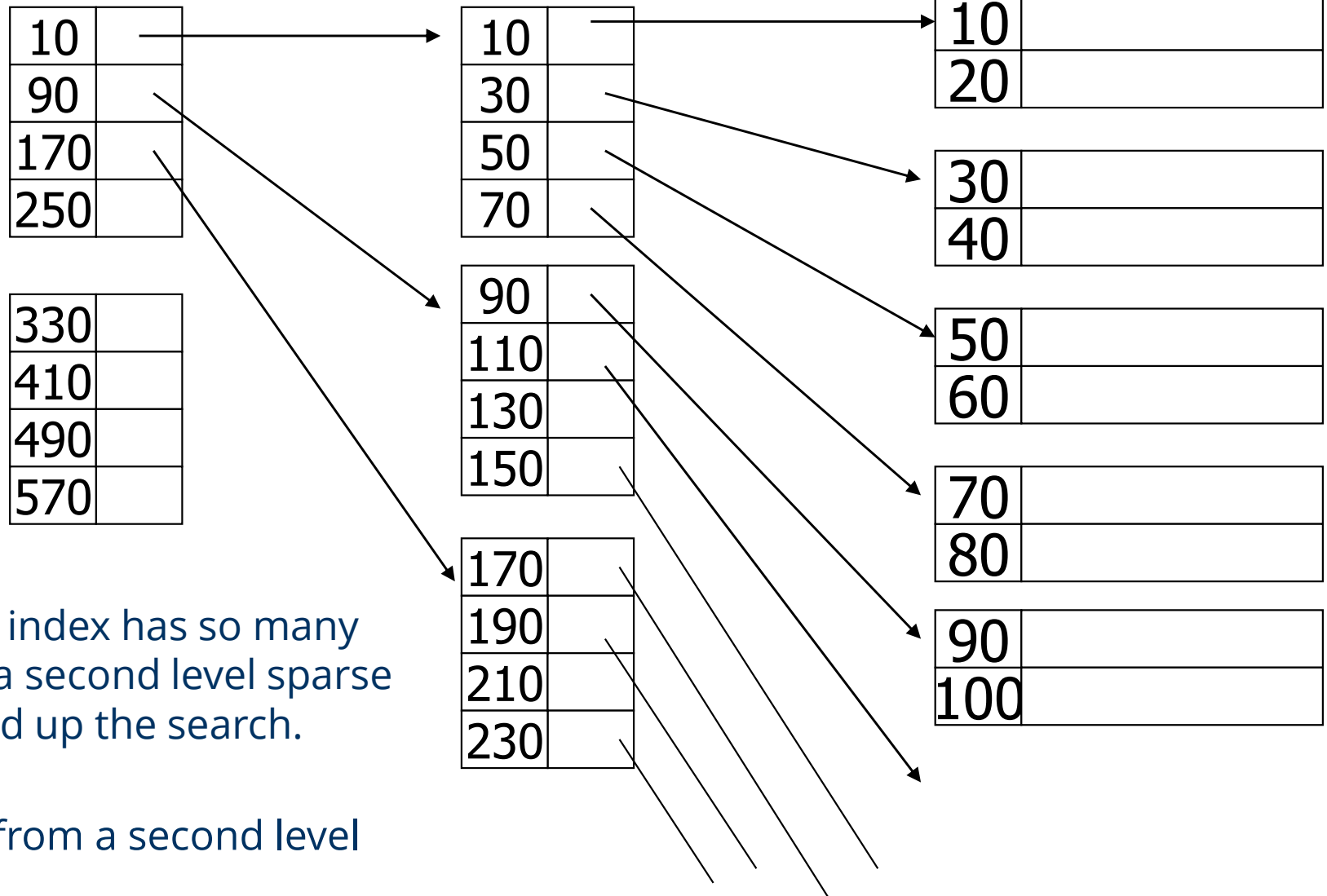⟹ *sorted sequential file on the same key as the index*

*not in the index records ⟹ we look for it in the previous page*

↳ *we do not know if it is there or not, we have to go through the records of the page sequentially*

**Sparse Index table:**

| | |
|---|---|
| 10 | |
| 30 | |
| 50 | |
| 70 | |

| | |
|---|---|
| 90 | |
| 110 | |
| 130 | |
| 150 | |

| | |
|---|---|
| 170 | |
| 190 | |
| 210 | |
| 230 | |

**Sequential File:**

| | |
|---|---|
| 10 | |
| 20 | |

| | |
|---|---|
| 30 | |
| 40 | |

| | |
|---|---|
| 50 | |
| 60 | |

| | |
|---|---|
| 70 | |
| 80 | |

| | |
|---|---|
| 90 | |
| 100 | |

## Sparse 2nd level

## Sequential File

| 10 | |
| 90 | |
| 170 | |
| 250 | |

| 330 | |
| 410 | |
| 490 | |
| 570 | |

| 10 | |
| 30 | |
| 50 | |
| 70 | |

| 90 | |
| 110 | |
| 130 | |
| 150 | |

| 170 | |
| 190 | |
| 210 | |
| 230 | |

| 10 | |
| 20 | |

| 30 | |
| 40 | |

| 50 | |
| 60 | |

| 70 | |
| 80 | |

| 90 | |
| 100 | |

If the first level index has so many pages, adding a second level sparse index can speed up the search.

Do we benefit from a second level dense index ?

⟶ no, we keep the same amount of entries ⟹ useless

# Sparse vs. Dense Tradeoff

- Sparse: Less index space per record
  can keep more of index in memory
- Dense:  Can tell if any record exists
  without accessing file

*but bigger size, we could be unable to store the whole index in the memory*

*reduce number of access (mean cost in the database)*

# Duplicate keys

## Dense index, one way to implement?

Dense index.

Can We do better ?

# Duplicate keys

## Dense index, better way?



require the file to be sorted in the same order as the index

# Duplicate keys

## Sparse index, one way?

Sparse index - place key from block

How to search for 30 ?

1. binary search the index
2. find page of the record
3. scan backward and forward
   ↳
   since store by pages, can be some in the previous and the next page

| 10 | |
|----|--|
| 10 | |
| 20 | |
| 30 | |

| | |
|--|--|
| | |
| | |
| | |
| | |

| 10 | |
|----|--|
| 10 | |

| 10 | |
|----|--|
| 20 | |

| 20 | |
|----|--|
| 30 | |

| 30 | |
|----|--|
| 30 | |

| 40 | |
|----|--|
| 45 | |

# Duplicate keys

## Sparse index, another way?

Sparse index - place
first **new** key from block

How to search for 30 ?

↳ no need to scan back

How to search for 35 ?

↳ go to the previous
pointer and scan
the page for the key

we have to
keep the entry
to be able to
access it even though
there is no new key

| 10 | |
| 20 | |
| 30 | |
| 30 | |

| 10 | |
| 10 | |

| 10 | |
| 20 | |

| 20 | |
| 30 | |

| 30 | |
| 30 | |

| 40 | |
| 45 | |

# Deletion from sparse index

| | |
|---|---|
| 10 | |
| 30 | |
| 50 | |
| 70 | |

| | |
|---|---|
| 90 | |
| 110 | |
| 130 | |
| 150 | |

| | |
|---|---|
| 10 | |
| 20 | |

| | |
|---|---|
| 30 | |
| 40 | |

| | |
|---|---|
| 50 | |
| 60 | |

| | |
|---|---|
| 70 | |
| 80 | |

# Deletion from sparse index

– delete record 40

| | |
|---|---|
| 10 | |
| 30 | |
| 50 | |
| 70 | |

| | |
|---|---|
| 90 | |
| 110 | |
| 130 | |
| 150 | |

| | |
|---|---|
| 10 | |
| 20 | |

| | |
|---|---|
| 30 | |
| 40 | |

| | |
|---|---|
| 50 | |
| 60 | |

| | |
|---|---|
| 70 | |
| 80 | |

# Deletion from sparse index

  – delete record 40



we are lucky
if was not a
key of the index
⇒ just suppress it

+ we could want to
keep the data sequential
in the file
(shifting up)

# Deletion from sparse index

&ndash; delete record 30

# Deletion from sparse index

– delete record 30

# Deletion from sparse index

– delete records 30 & 40

# Deletion from sparse index

– delete records 30 & 40

# Deletion from sparse index

– delete records 30 & 40

# Deletion from dense index

# Deletion from dense index

– delete record 30

# Deletion from dense index

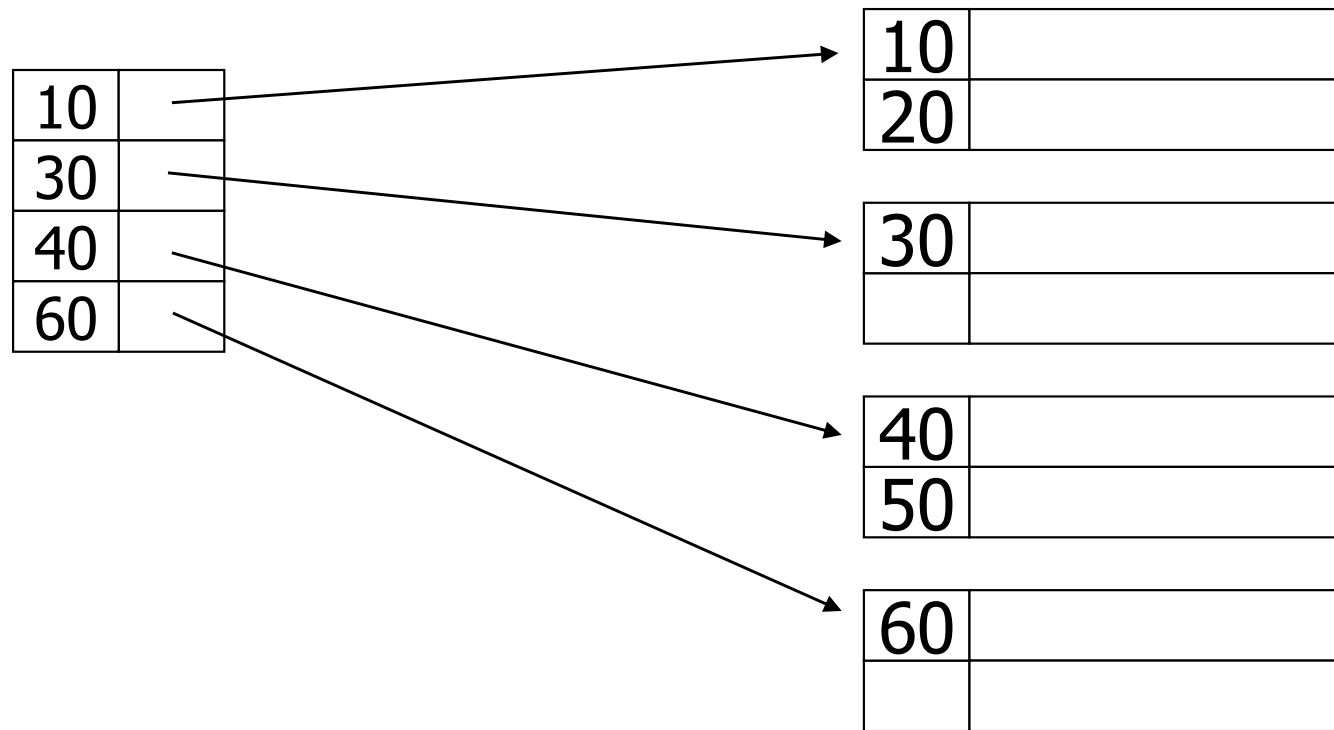– delete record 30

# Deletion from dense index

– delete record 30

# Insertion, sparse index case

# Insertion, sparse index case
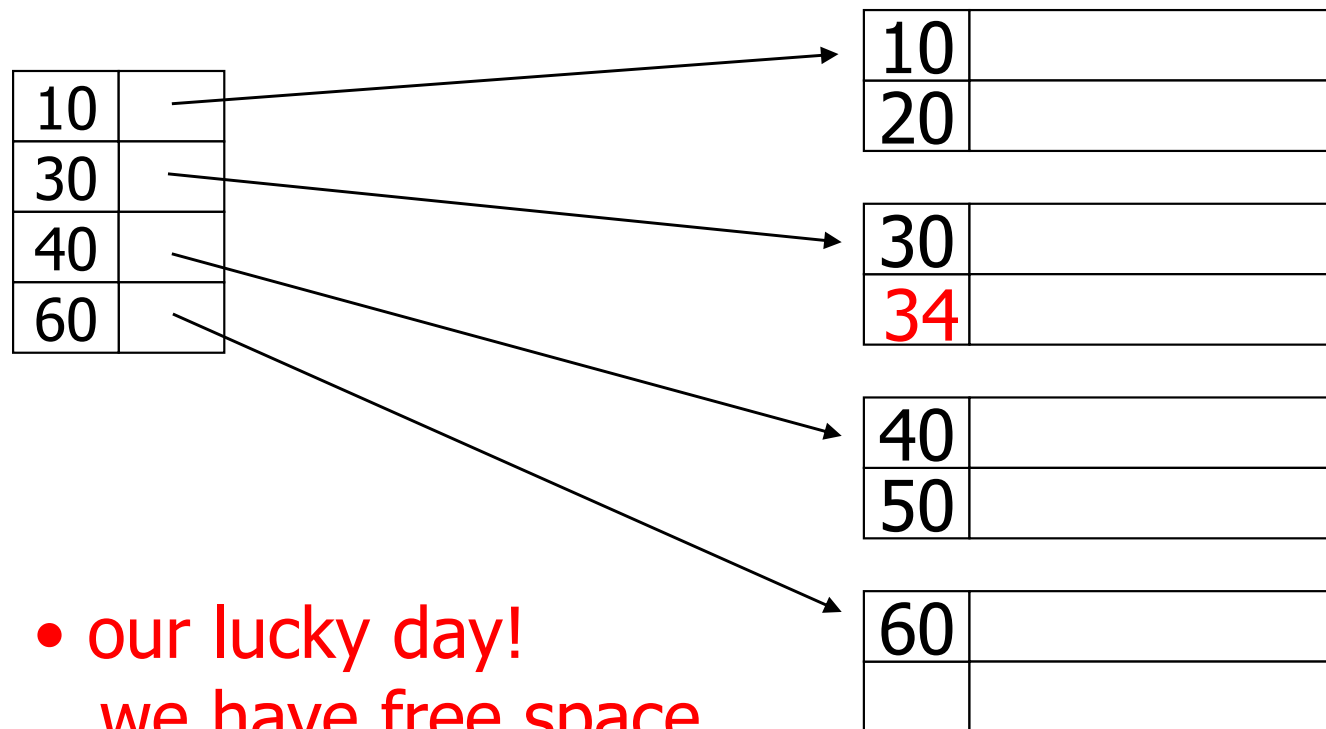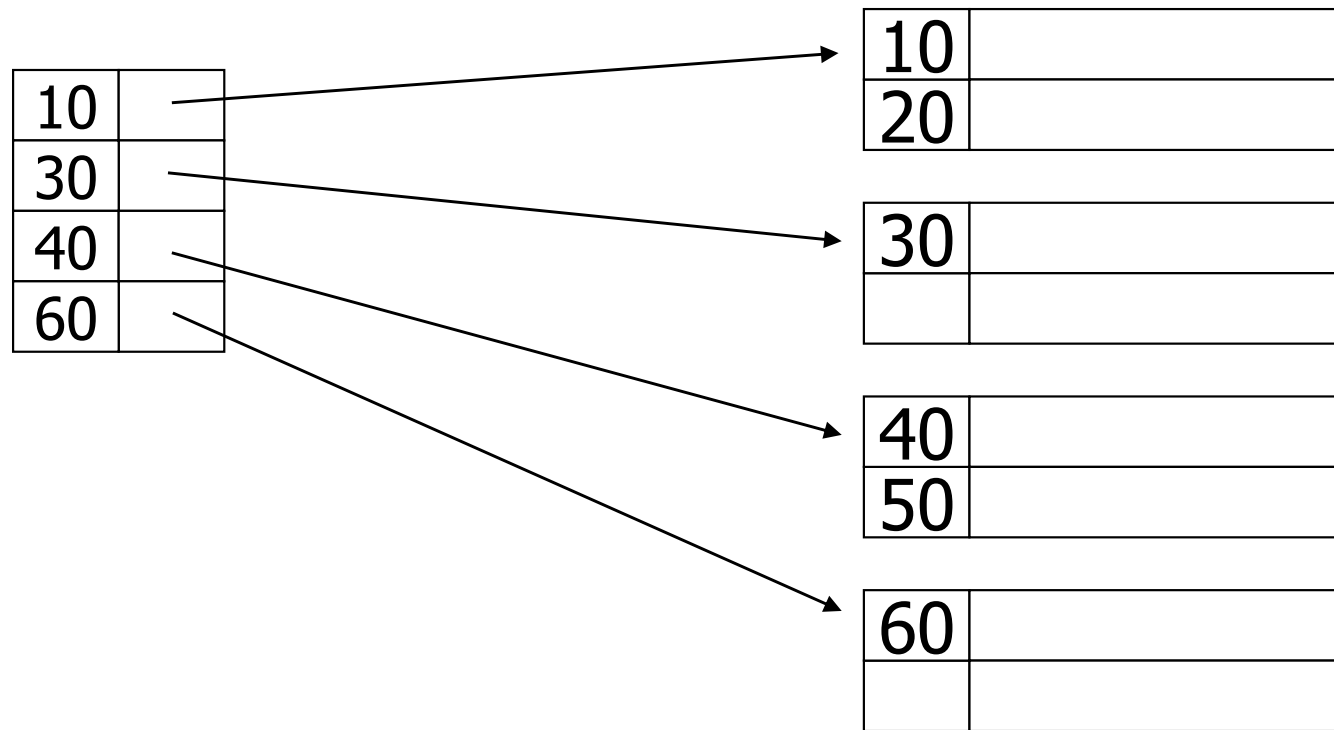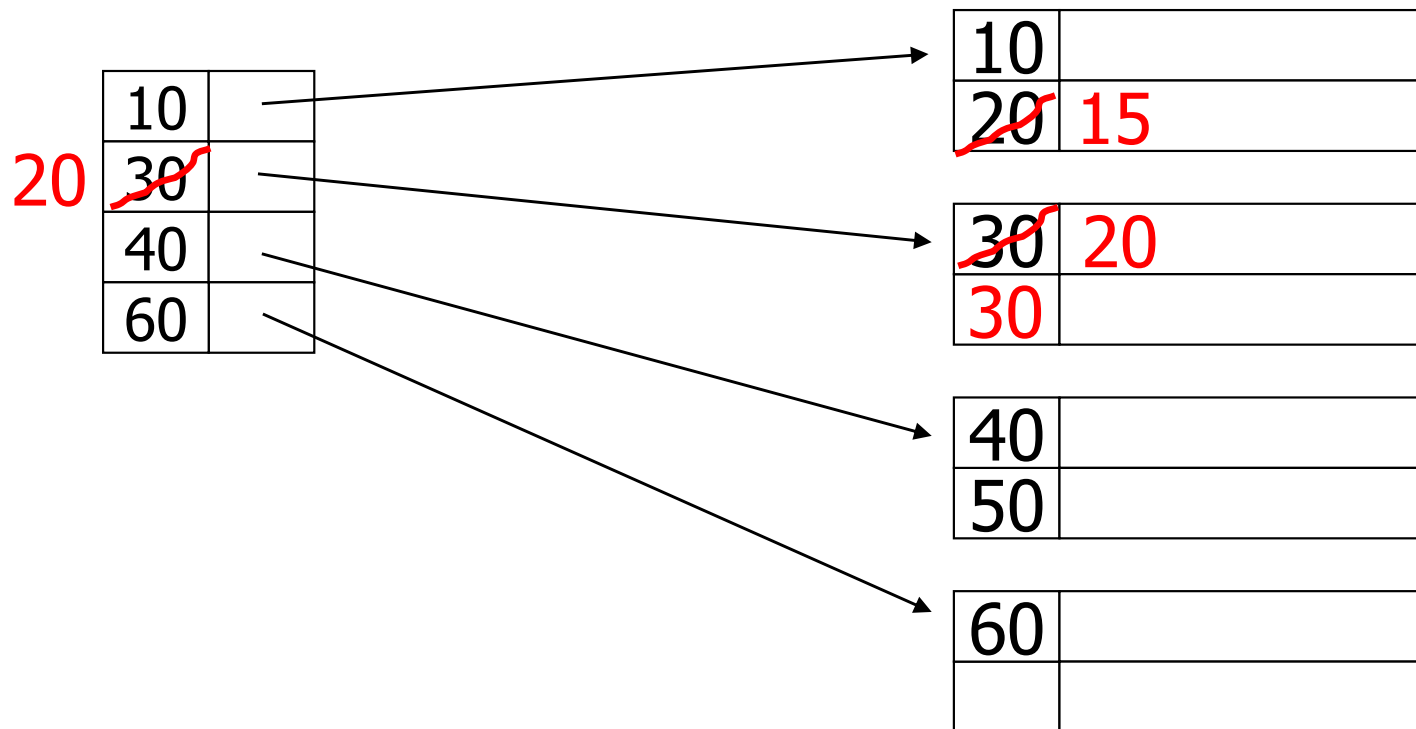
– insert record 34

# Insertion, sparse index case

– insert record 34

| 10 | |
|----|--|
| 30 | |
| 40 | |
| 60 | |

| 10 | |
|----|--|
| 20 | |

| 30 | |
|----|--|
| 34 | |

| 40 | |
|----|--|
| 50 | |

| 60 | |
|----|--|
|    | |

• our lucky day!
we have free space
where we need it!

# Insertion, sparse index case
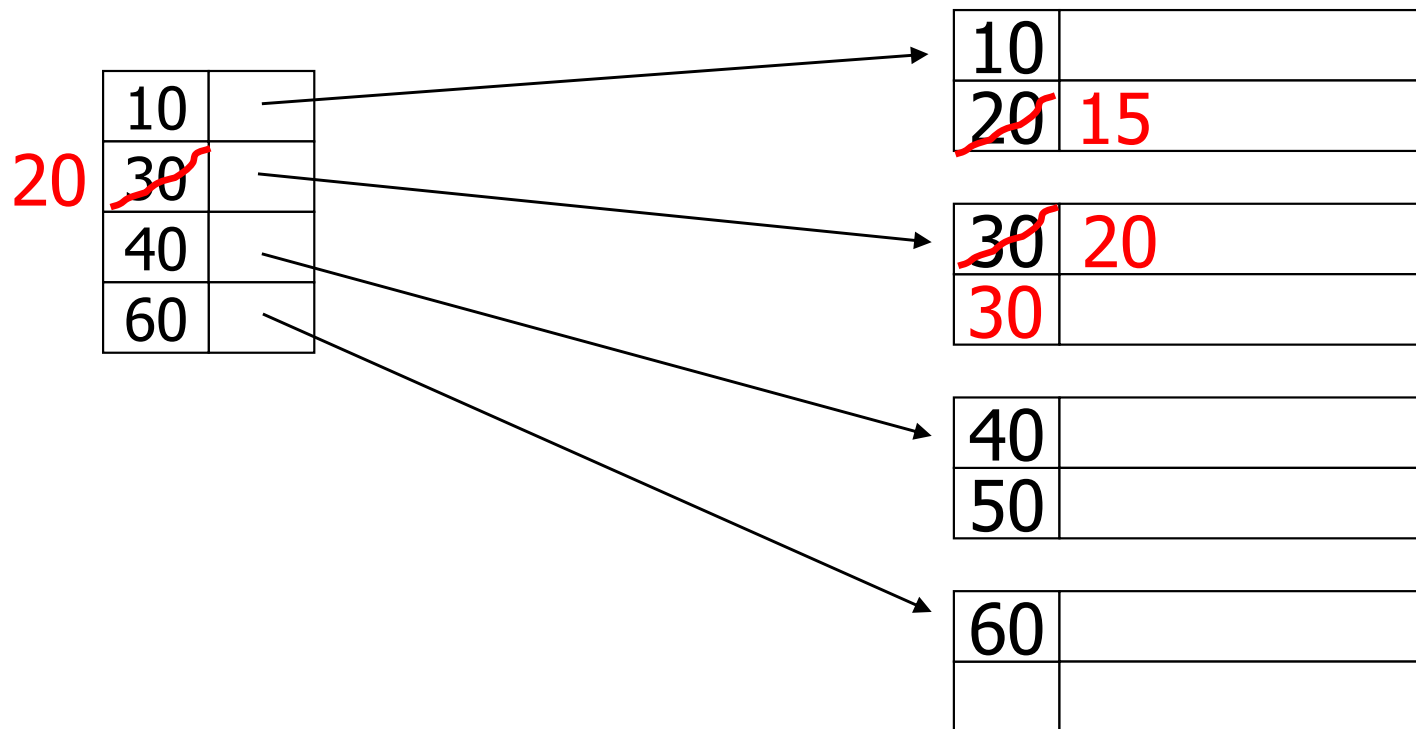
– insert record 15
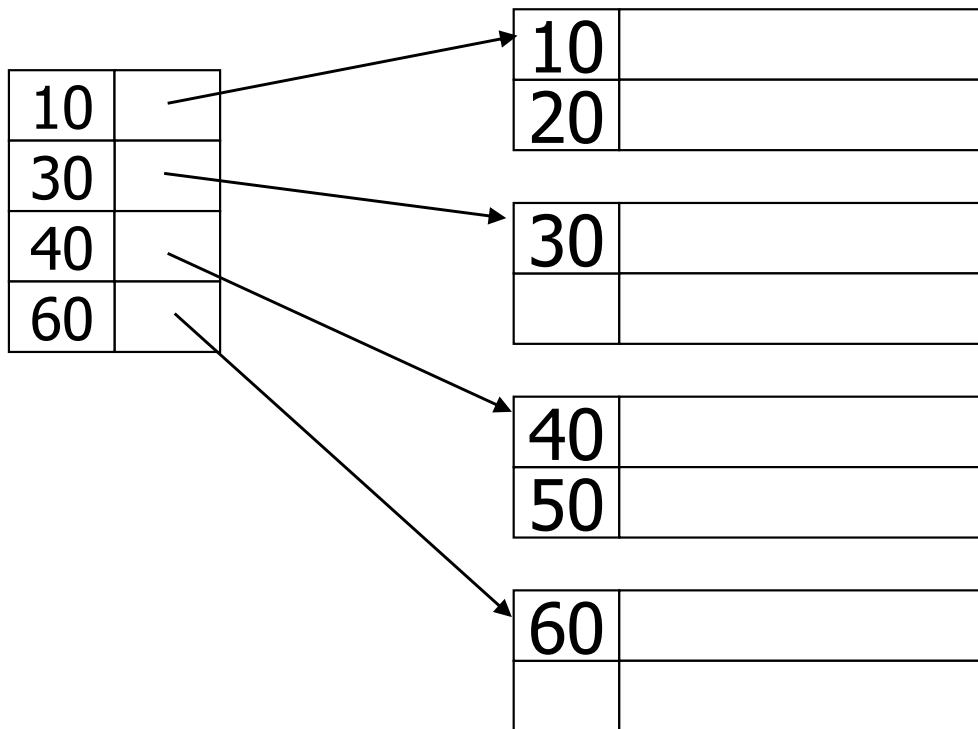
# Insertion, sparse index case

– insert record 15

# Insertion, sparse index case

– insert record 15

# Insertion, sparse index case

– insert record 25

| 10 |  |
|----|--|
| 30 |  |
| 40 |  |
| 60 |  |

| 10 |  |
|----|--|
| 20 |  |

| 30 |  |
|----|--|
|    |  |

| 40 |  |
|----|--|
| 50 |  |

| 60 |  |
|----|--|
|    |  |

1. follow the index
   ↳ have to be after 10
2. we go to the page
3. we scan, we need to go to the next page
4. insert 25 and shift 30
5. replace the index key by 25

# Insertion, dense index case

- Similar

- Often more expensive . . .

# Secondary indexes

Sequence field

| 30 | |
|----|--|
| 50 | |

| 20 | |
|----|--|
| 70 | |

| 80 | |
|----|--|
| 40 | |

| 100 | |
|-----|--|
| 10  | |

| 90 | |
|----|--|
| 60 | |

# Secondary indexes

- Sparse index

Sequence field

| 30 | |
|----|--|
| 20 | |
| 80 | |
| 100 | |

| 90 | |
|----|--|
| ... | |
| | |
| | |

| 30 | |
|----|--|
| 50 | |

| 20 | |
|----|--|
| 70 | |

| 80 | |
|----|--|
| 40 | |

| 100 | |
|-----|--|
| 10 | |

| 90 | |
|----|--|
| 60 | |

34

# Secondary indexes

- Sparse index

Sequence field

| 30 | |
|----|--|
| 20 | |
| 80 | |
| 100 | |

| 90 | |
|----|--|
| ... | |
| | |
| | |

| 30 | |
|-----|--|
| 50 | |

| 20 | |
|-----|--|
| 70 | |

| 80 | |
|-----|--|
| 40 | |

| 100 | |
|-----|--|
| 10 | |

| 90 | |
|-----|--|
| 60 | |

**does not make sense!**

won't work because file
is not sorted

35

# Secondary indexes

- Dense index

Sequence
field

| 30 | |
|----|----|
| 50 | |

| 20 | |
|----|----|
| 70 | |

| 80 | |
|----|----|
| 40 | |

| 100 | |
|----|----|
| 10 | |

| 90 | |
|----|----|
| 60 | |

# Secondary indexes

- Dense index

| 10 | |
|----|---|
| 20 | |
| 30 | |
| 40 | |

| 50 | |
|----|---|
| 60 | |
| 70 | |
| ... | |

| 30 | |
|----|---|
| 50 | |

| 20 | |
|----|---|
| 70 | |

| 80 | |
|----|---|
| 40 | |

| 100 | |
|-----|---|
| 10 | |

| 90 | |
|----|---|
| 60 | |

*now we can do binary search on the index while we couldn't on the sequence field*

37

# Secondary indexes

- Dense index

| | |
|---|---|
| 10 | |
| 50 | |
| 90 | |
| ... | |

sparse
high
level

| | |
|---|---|
| 10 | |
| 20 | |
| 30 | |
| 40 | |

| | |
|---|---|
| 50 | |
| 60 | |
| 70 | |
| ... | |

| | |
|---|---|
| 30 | |
| 50 | |

| | |
|---|---|
| 20 | |
| 70 | |

| | |
|---|---|
| 80 | |
| 40 | |

| | |
|---|---|
| 100 | |
| 10 | |

| | |
|---|---|
| 90 | |
| 60 | |

wouldn't have
sense to be
dense index

wouldn't have
sense to be sparse
index

# With secondary indexes:

- Lowest level is dense
- Other levels are sparse

# Duplicate values & secondary indexes

| | |
|---|---|
| 20 | |
| 10 | |

| | |
|---|---|
| 20 | |
| 40 | |

| | |
|---|---|
| 10 | |
| 40 | |

| | |
|---|---|
| 10 | |
| 40 | |

| | |
|---|---|
| 30 | |
| 40 | |

# Duplicate values & secondary indexes

one option…

| 10 | |
|----|--|
| 10 | |
| 10 | |
| 20 | |

| 20 | |
|----|--|
| 30 | |
| 40 | |
| 40 | |

| 40 | |
|----|--|
| 40 | |
| … | |
| | |

| 20 | |
|----|--|
| 10 | |

| 20 | |
|----|--|
| 40 | |

| 10 | |
|----|--|
| 40 | |

| 10 | |
|----|--|
| 40 | |

| 30 | |
|----|--|
| 40 | |

# Duplicate values & secondary indexes

one option...

Problem:
excess overhead!
- disk space
- search time

| 10 | |
|----|--|
| 10 | |
| 10 | |
| 20 | |

| 20 | |
|----|--|
| 30 | |
| 40 | |
| 40 | |

| 40 | |
|----|--|
| 40 | |
| ... | |
| | |

| 20 | |
|----|--|
| 10 | |

| 20 | |
|----|--|
| 40 | |

| 10 | |
|----|--|
| 40 | |

| 10 | |
|----|--|
| 40 | |

| 30 | |
|----|--|
| 40 | |

# Duplicate values & secondary indexes

another option...

4 integers in place of 6
⇒ less big

# Duplicate values & secondary indexes

another option...

Problem:
variable size
records in
index!

10

20

30
40

20
10

20
40

10
40

10
40

30
40

↳ we need fixed
size records to do
binary search

# Duplicate values & secondary indexes



buckets

↳ still divided in pages

# Why "bucket" idea is useful

## Indexes

Name: primary

Dept: secondary

Floor: secondary

## Records

EMP (name,dept,floor,...)

# Query: Get employees in
##    (Toy Dept) ∧ (2nd floor)

Dept. index                    EMP                    Floor index

| Toy | |

| 2nd |

because we have the buckets
we can compare the pointers
⟹ do not need to access the records

# Query: Get employees in
## (Toy Dept) $\wedge$ (2nd floor)

Dept. index          EMP         Floor index

Toy                          2nd

$\rightarrow$ Intersect toy bucket and 2nd Floor bucket to get set of matching EMP's

# Summary so far

- Conventional index
  - Basic Ideas: sparse, dense, multi-level…
  - Duplicate Keys
  - Deletion/Insertion
  - Secondary indexes

Note: those are building blocks
not used in practice

# Conventional indexes

## Advantage:

- Simple
- Index is sequential file
  good for scans

## Disadvantage:
- Inserts expensive, or
- Lose sequentiality & balance

↳ the idea is to have indexes that are sorted
conventional indexes were sorted lists of key

⟹ maybe using other data structure would be better
  ↳ binary trees $\log(n)$ insertion and search

50

# Example   Index (sequential)

| | |
|---|---|
| 10 | → |
| 20 | → |
| 30 | → |
| | |

continuous

| | |
|---|---|
| 40 | → |
| 50 | → |
| 60 | → |
| | |

free space

| | |
|---|---|
| 70 | → |
| 80 | → |
| 90 | → |
| | |

# Example

## Index (sequential)



continuous

free space

overflow area
(not sequential)

# Outline:

- Conventional indexes
- B-Trees                    $\Rightarrow$ NEXT
- Hashing schemes

- **NEXT: Another type of index**
  - Give up on sequentiality of index
  - Try to get "balance"

unbalanced trees
are far from the
$\log(n)$ insertion
and search

# B+Tree Example

n=3

Root



| | | |
|---|---|---|
| 100 | | |

| | |
|---|---|
| 30 | |

| | | |
|---|---|---|
| 120 | 150 | 180 |

| | | |
|---|---|---|
| 3 | 5 | 11 |

| | |
|---|---|
| 30 | 35 |

| | | |
|---|---|---|
| 100 | 101 | 110 |

| | |
|---|---|
| 120 | 130 |

| | | |
|---|---|---|
| 150 | 156 | 179 |

| | |
|---|---|
| 180 | 200 |

# Sample non-leaf



to keys  to keys       to keys    to keys

< 57   57≤ k<81      81≤k<95    ≥95

# Sample leaf node:

From non-leaf node

to next leaf
in sequence

57   81   95

To record with key 57
To record with key 81
To record with key 85

Size of nodes:  $\left.\begin{array}{l}\text{n+1 pointers}\\\text{n keys}\end{array}\right\}$ (fixed)

# Don't want nodes to be too empty

- Use at least

Non-leaf: $\lceil (n+1)/2 \rceil$ pointers

Leaf: $\lfloor (n+1)/2 \rfloor$ pointers to data

n=3

Full node     min. node

Non-leaf

120 150 180

30

Leaf

3 5 11

30 35

counts even if null

# B+tree rules            tree of order *n*

(1) All leaves at same lowest level
       (balanced tree)

(2) Pointers in leaves point to records
       except for "sequence pointer"

↳ point to the next leaf

# (3) Number of pointers/keys for B+tree

| | Max ptrs | Max keys | Min ptrs→data | Min keys |
|---|---|---|---|---|
| Non-leaf (non-root) | $n+1$ | $n$ | $\lceil (n+1)/2 \rceil$ | $\lceil (n+1)/2 \rceil - 1$ |
| Leaf (non-root) | $n+1$ | $n$ | $\lfloor (n+1)/2 \rfloor$ | $\lfloor (n+1)/2 \rfloor$ |
| Root | $n+1$ | $n$ | $1$ | $1$ |

rule broken
only for the root

# Insert into B+tree

(a) simple case

    – space available in leaf

(b) leaf overflow

(c) non-leaf overflow

(d) new root

whole idea of using
a binary tree
↳ for cost insertion

# (a) Insert key = 32

$\boxed{n=3}$



There is space we use it

# (a) Insert key = 32

n=3

# (a) Insert key = 7

# (a) Insert key = 7

n=3



we split the leaf

100

30

3
5

8 5 7 11

30 31

# (a) Insert key = 7

then
we look at
the parent node
to see if there
is room for
a new branch

100

30   7

3
5

8  7  11

30
31

68

# (c) Insert key = 160

n=3

100

120 150 180

150 156 179

180 200

69

# (c) Insert key = 160

n=3

# (c) Insert key = 160

n=3

(c) Insert key = 160

n=3

72

# (d) New root, insert 45

$n=3$

→ root is full

```
[ 10 20 30 ]
```

Tree structure:
- Root: 10 20 30
- Leaf 1: 1 2 3
- Leaf 2: 10 12
- Leaf 3: 20 25
- Leaf 4: 30 32 40

# (d) New root, insert 45

Root node: 10 20 30

Leaf nodes: 1 2 3 | 10 12 | 20 25 | 30 32 40 | 40 45

# (d) New root, insert 45

$n=3$

# (d) New root, insert 45

n=3

new root

30

*only one key but its okay because the root can break the minimum key constraint*

10 20 30

40

1 2 3

10 12

20 25

30 32 40

40 45

# Deletion from B+tree

(a) Simple case - no example

(b) Coalesce with neighbor (sibling)
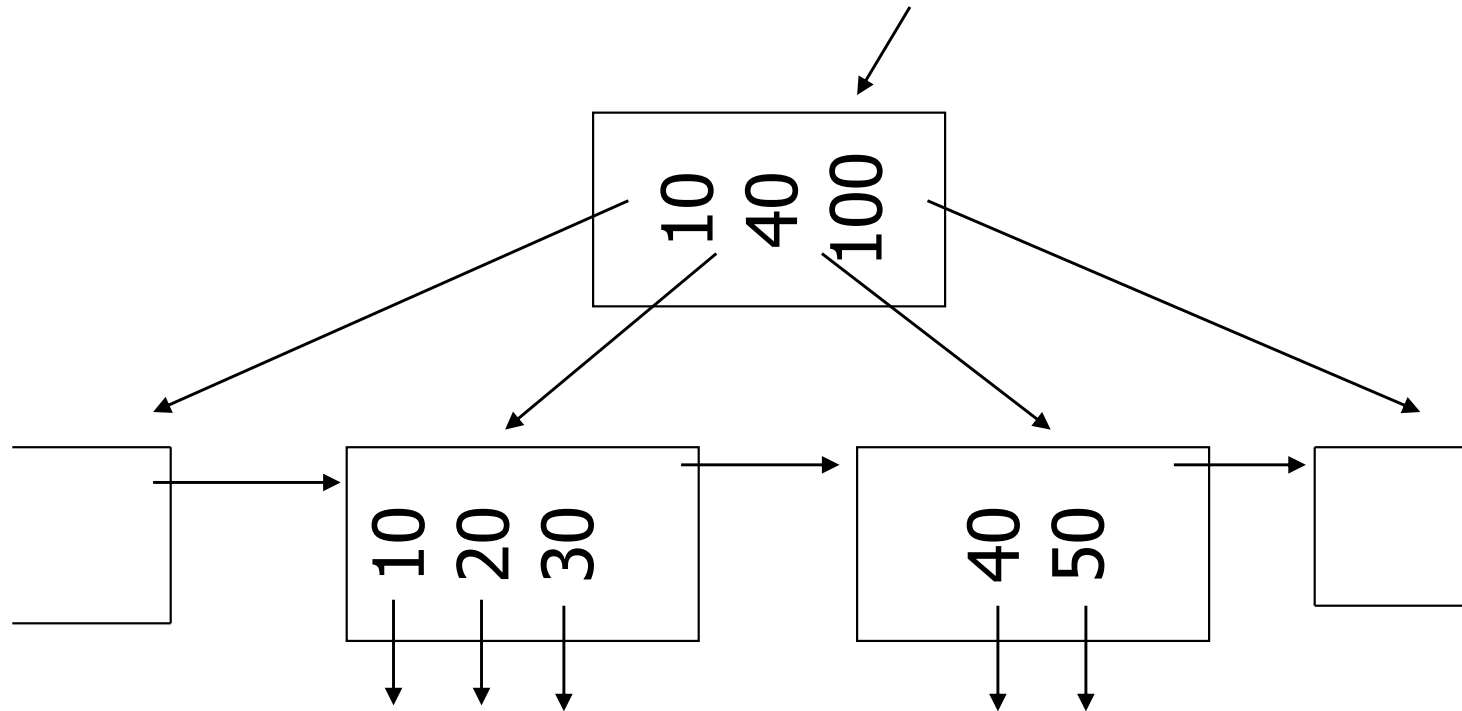
(c) Re-distribute keys

(d) Cases (b) or (c) at non-leaf

go first through it

↳ not implemented
in reality because
too hard and too
expensive

when there are
to many fragments
we rebuild the tree

we only mark
the values as
deleted and go
next

# (b) Coalesce with sibling
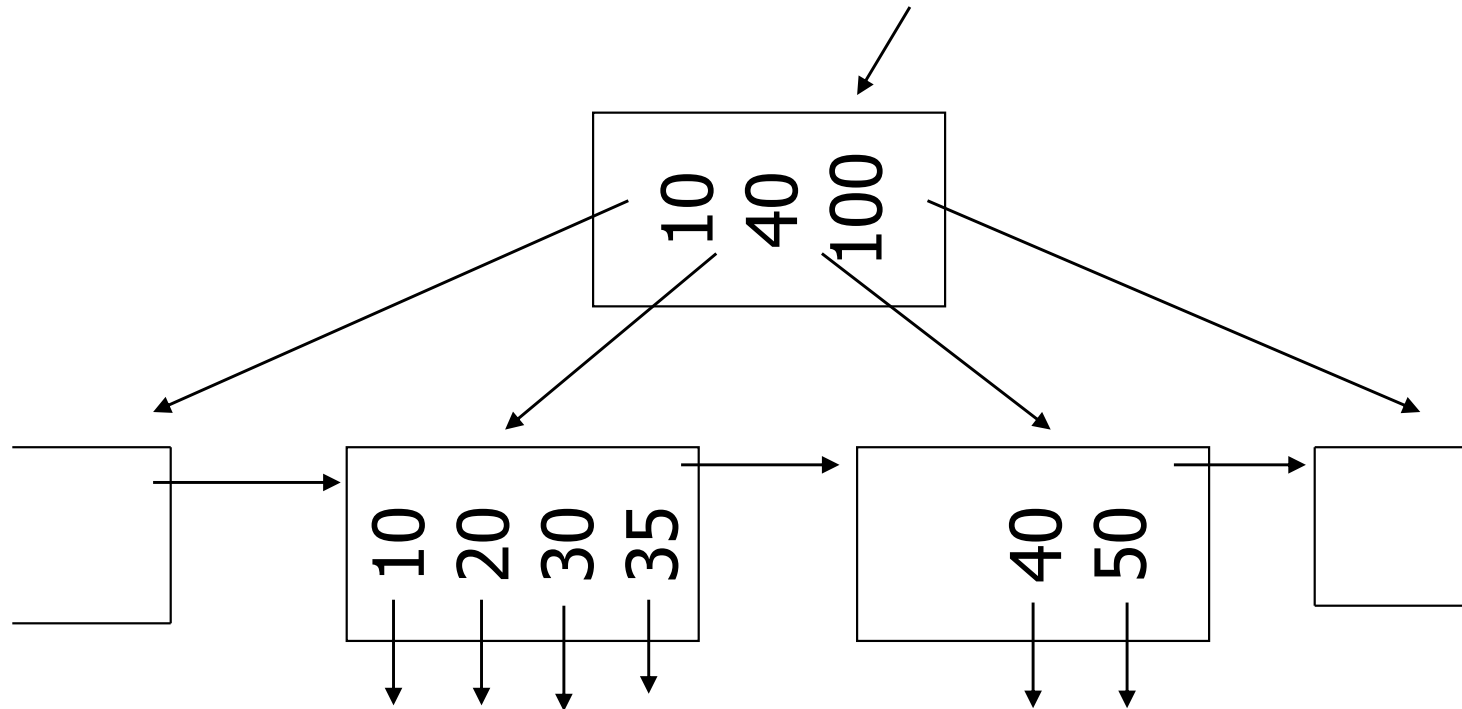
- Delete 50

# (b) Coalesce with sibling
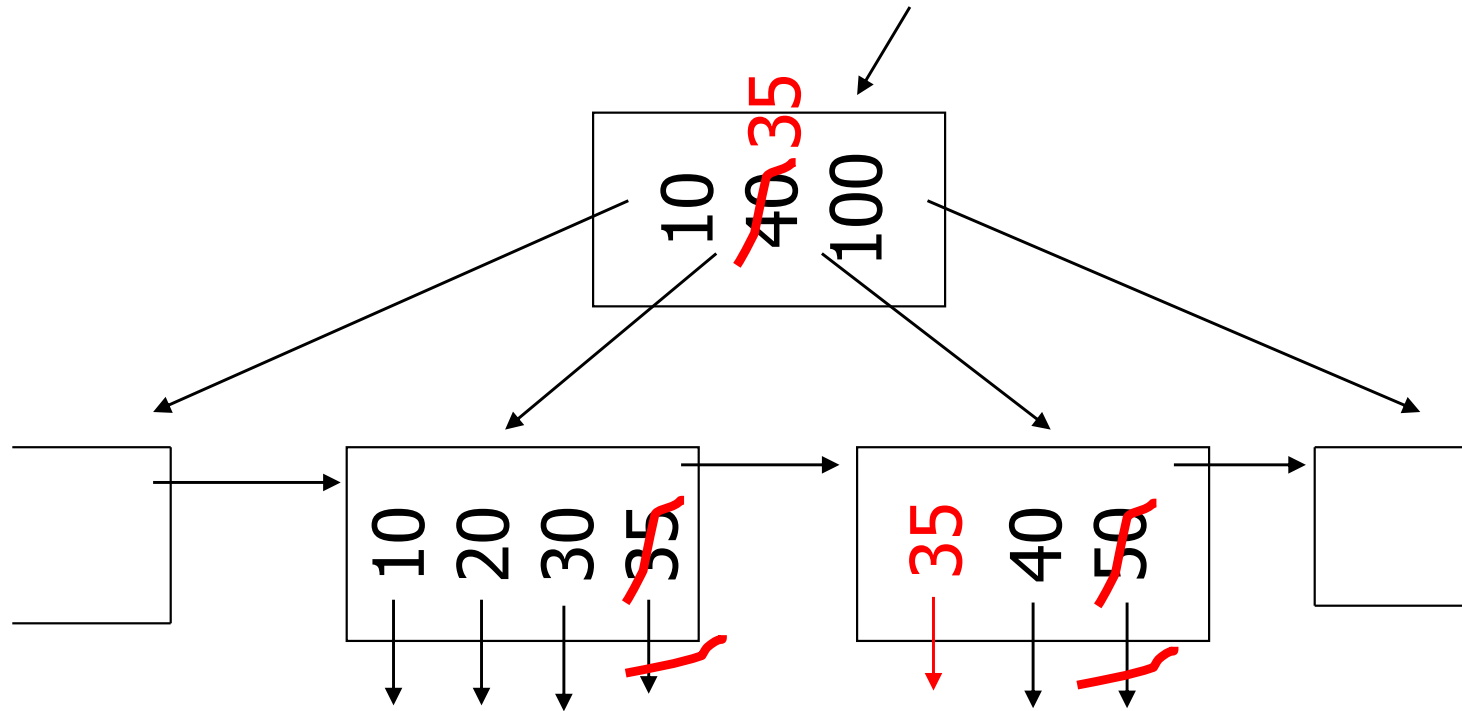 − Delete 50

n=4

# (c) Redistribute keys

– Delete 50

Root node: 10 40 100

Leaf 1: 10 20 30 35

Leaf 2: 40 50

# (c) Redistribute keys
– Delete 50

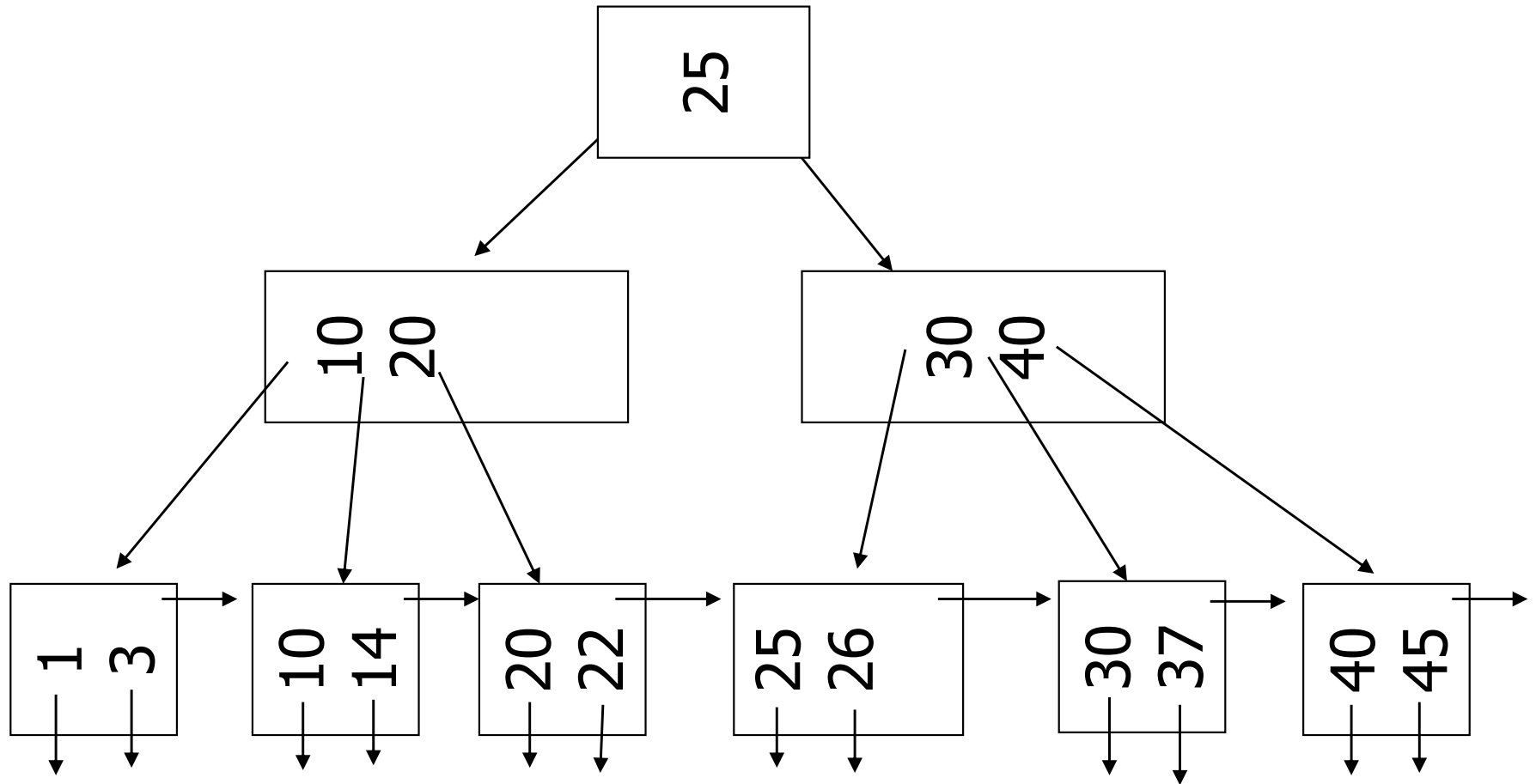n=4

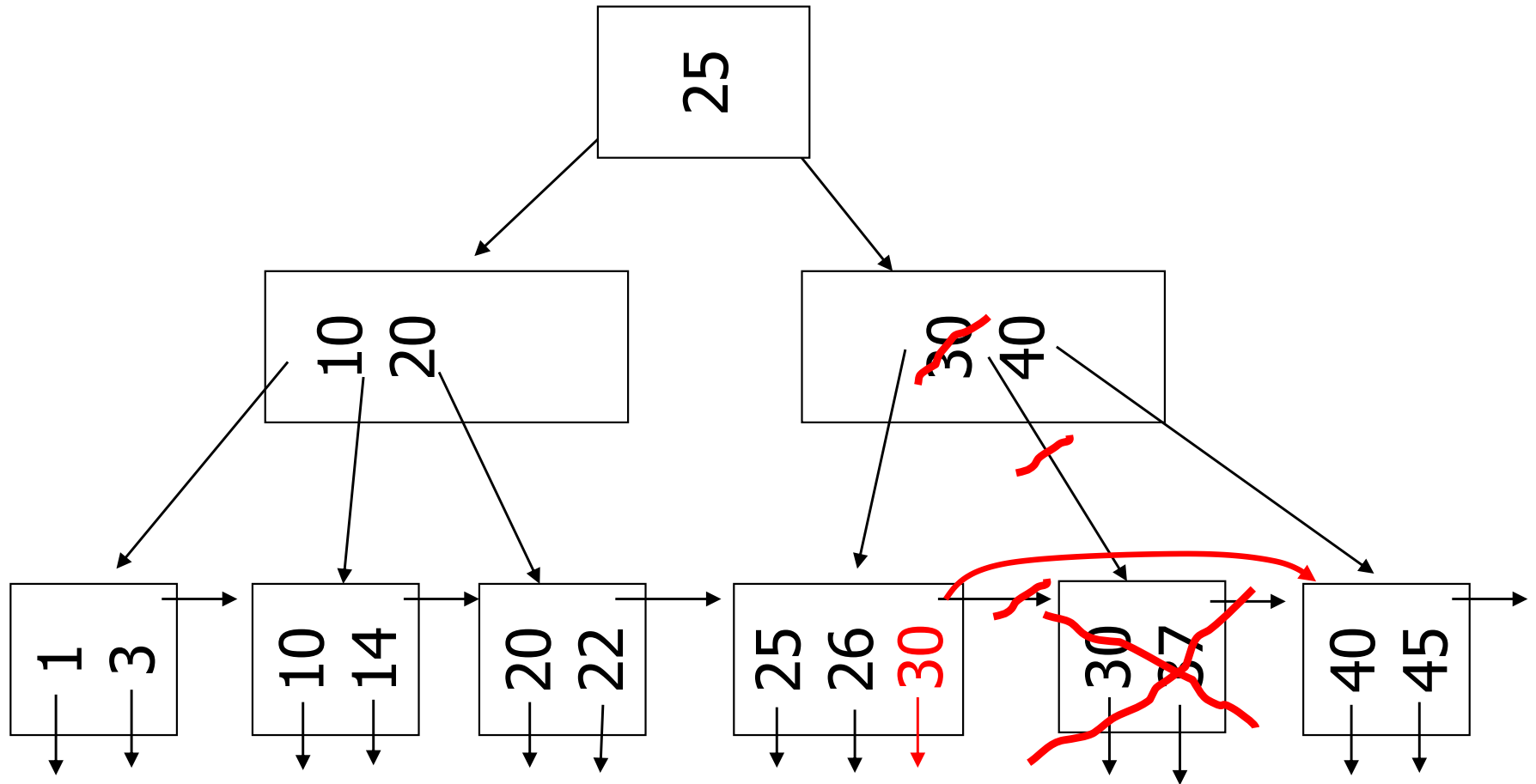# (d) Non-leaf coalesce
 – Delete 37

# (d) Non-leaf coalesce
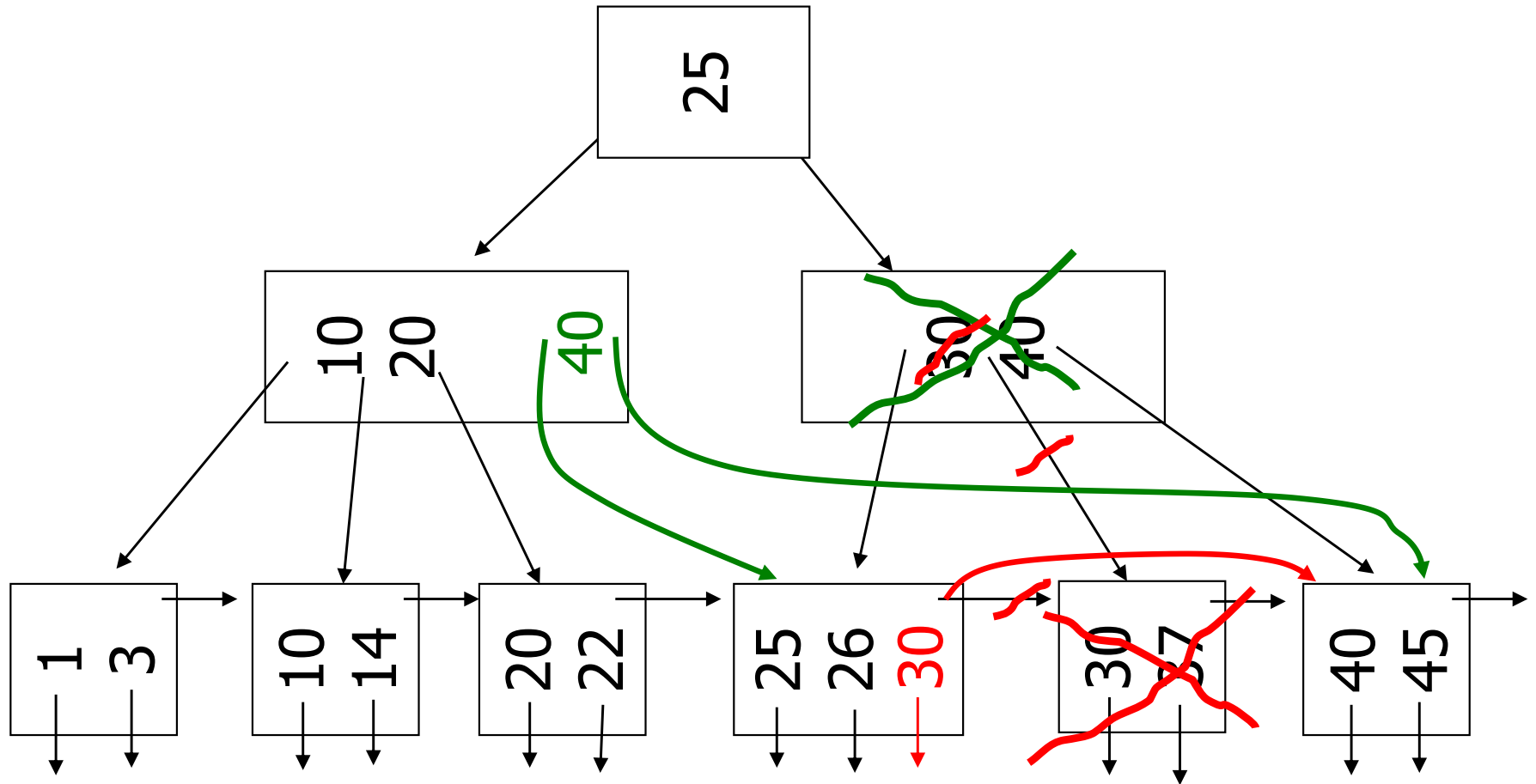  – Delete 37

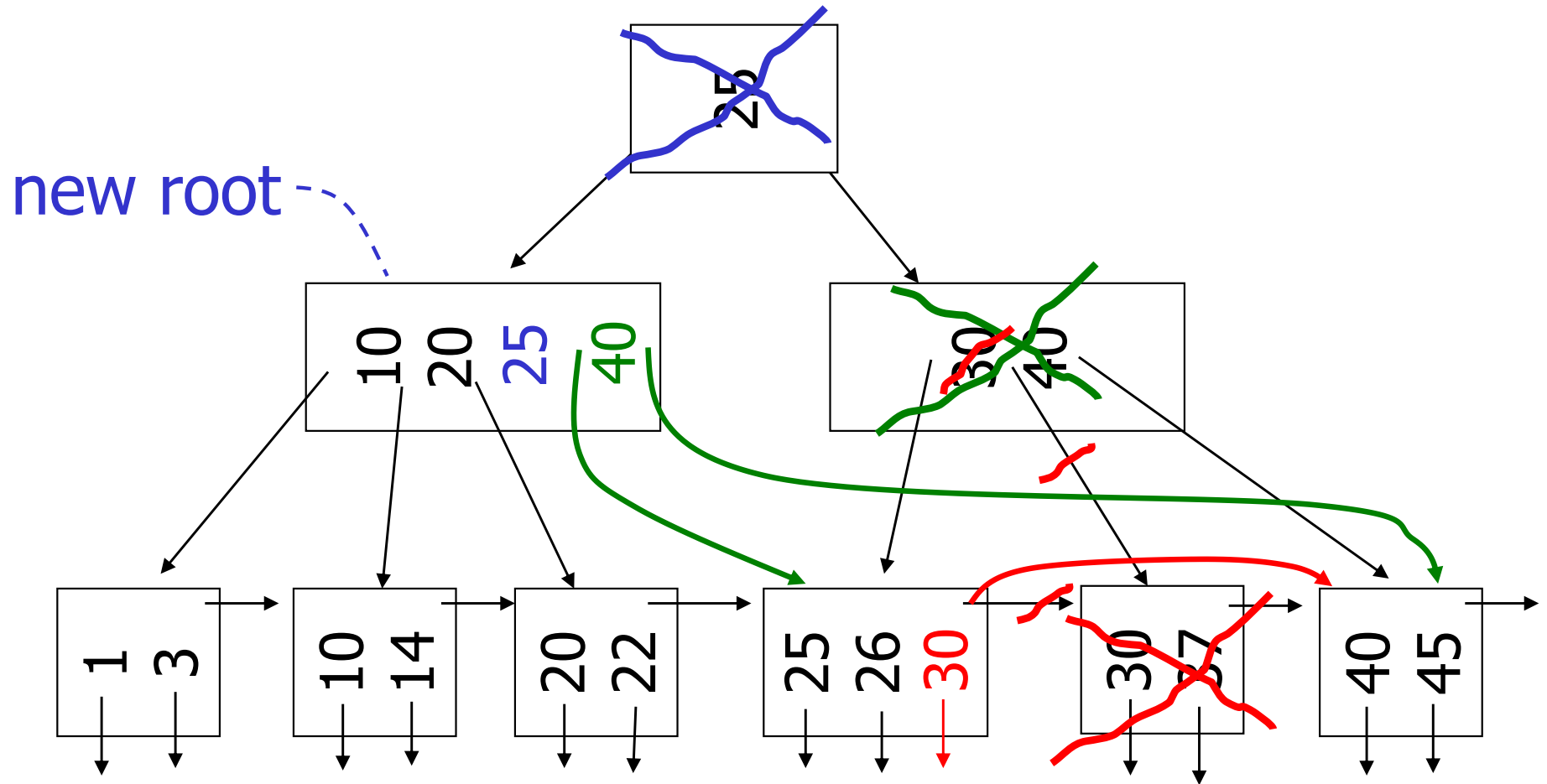# (d) Non-leaf coalesce
– Delete 37

n=4



84

# (d) Non-leaf coalesce
## – Delete 37

n=4

new root

# B+tree deletions in practice

– Often, coalescing is <u>not</u> implemented
   – Too hard and not worth it!

# Outline/summary

- Conventional Indexes
  - Sparse vs. dense
  - Primary vs. secondary

- B trees

- Hashing schemes (recommended reading, not mandatory)

The slides in this lecture are taken from:

- Hector Garcia-Molina, CS 245: Database System Principles, Notes 4: Indexing.

Reading

- Héctor García-Molina, Jeffrey Ullman, and Jennifer Widom. Database Systems: The Complete Book