# Data Storage and Query Answering
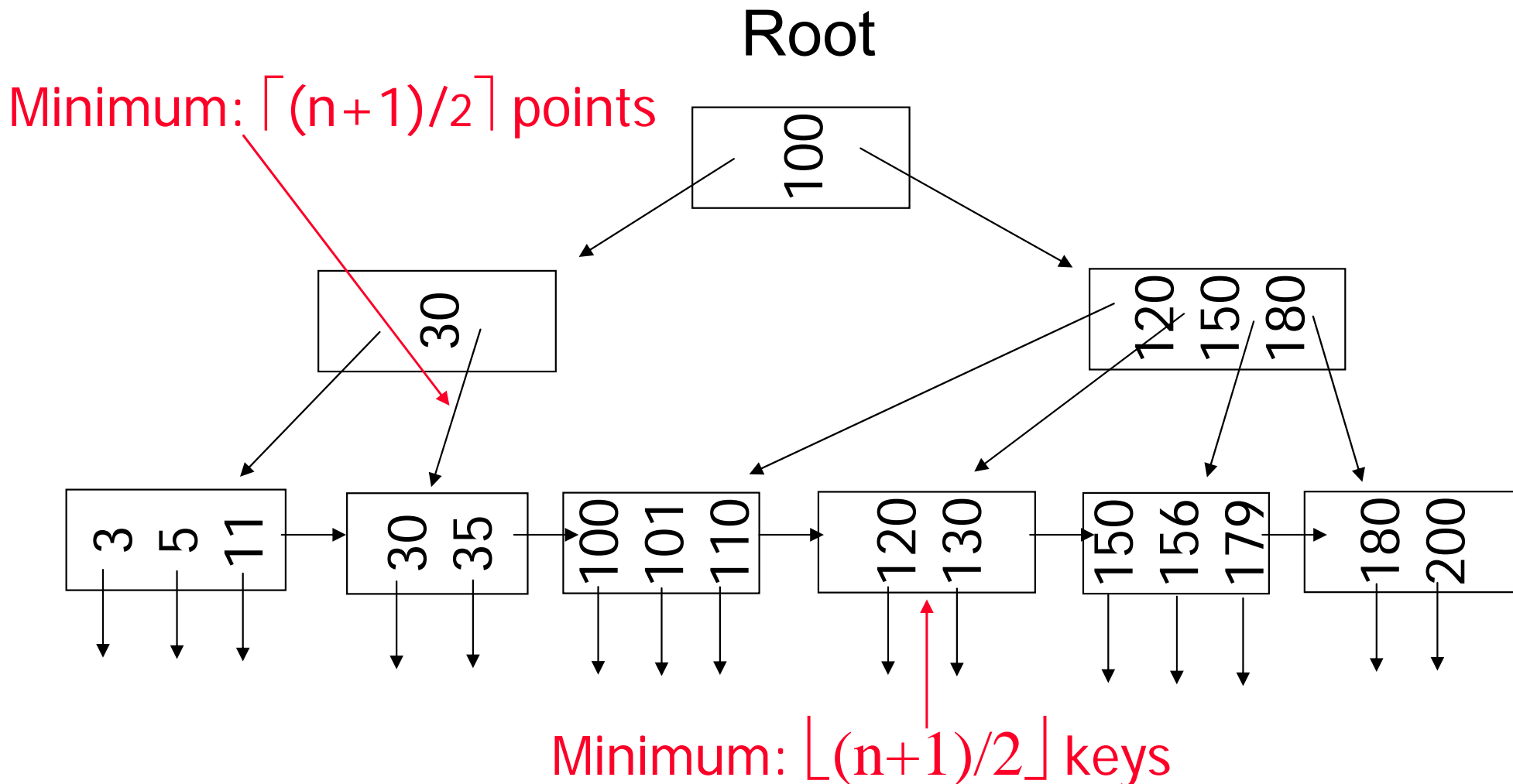
## Indexing and Hashing (3)

# B+ Tree Rules

- (1) Each node is a disk block
  - An I/O access retrieves the whole block into main memory

- (2) All leaves at same lowest level (balanced tree)

- (3) Pointers in leaves point to records, except for "sequence pointer"

- (4) Number of pointers/keys for B+ tree (order n)

|  | 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 |

## B+ Tree Example                    n=3

Root

Minimum: $\lceil (n+1)/2 \rceil$ points

100

30

120 150 180

3 5 11

30 35

100 101 110

120 130

150 156 179

180 200

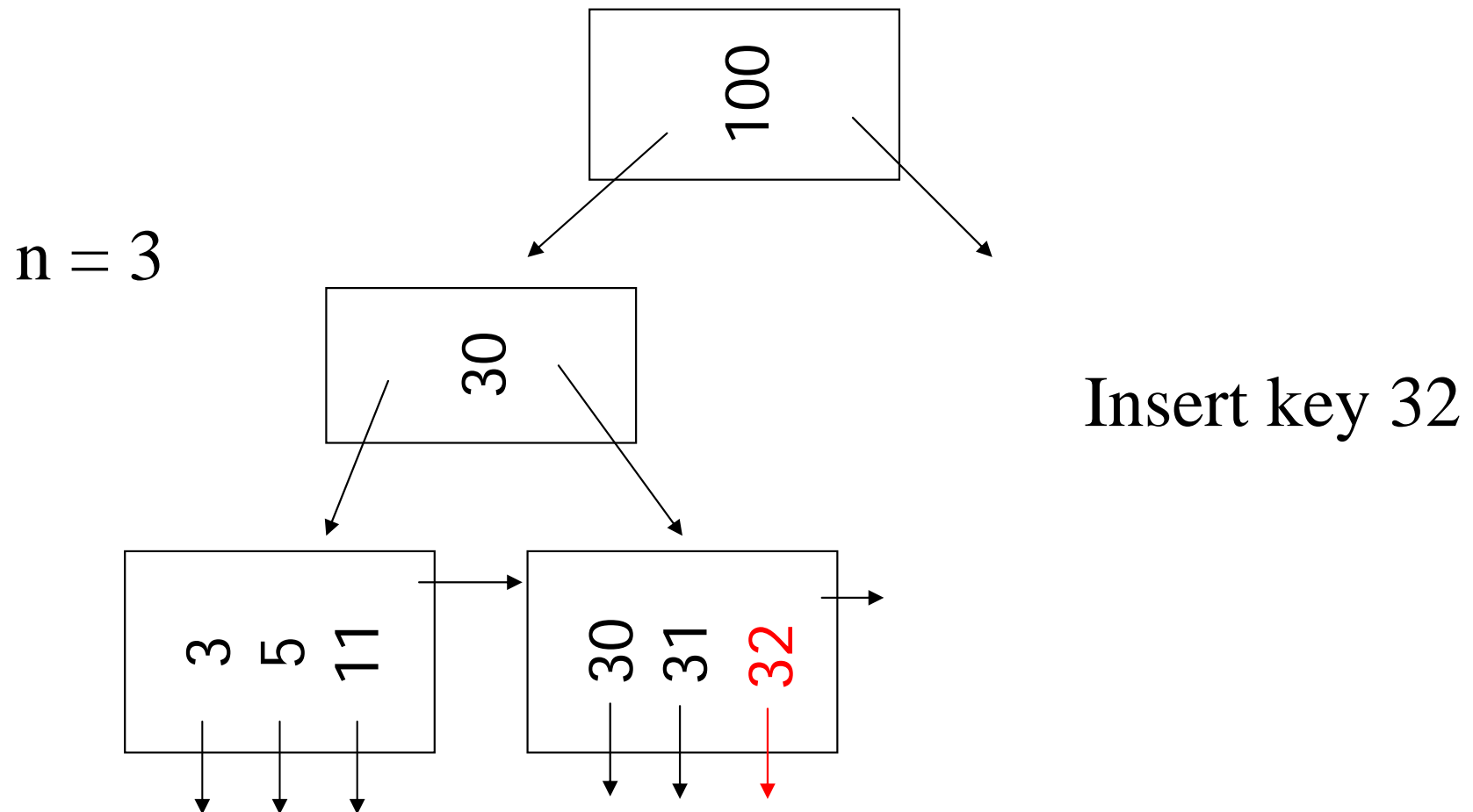Minimum: $\lfloor (n+1)/2 \rfloor$ keys

# B+ Trees

*Insertions*

- Always insert in corresponding leaf.

- Tree grows bottom-up.

- Four different cases:

  - Space available in leaf,

  - Leaf overflow,
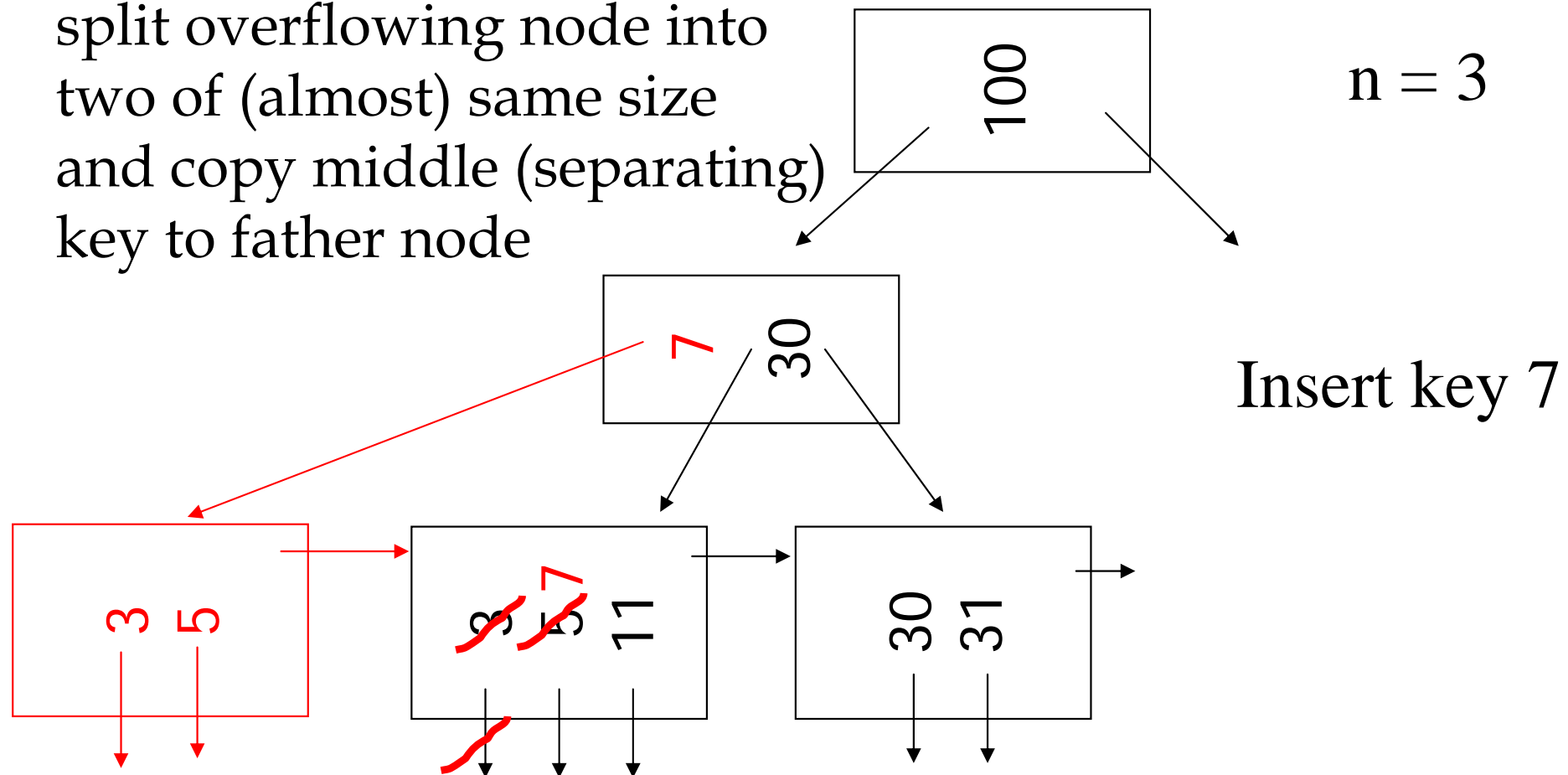
  - Non-leaf overflow,

  - New root.

*Insertions*

- Space available in leaf

n = 3

100

30

3
5
11

30
31
32

Insert key 32

## Insertions
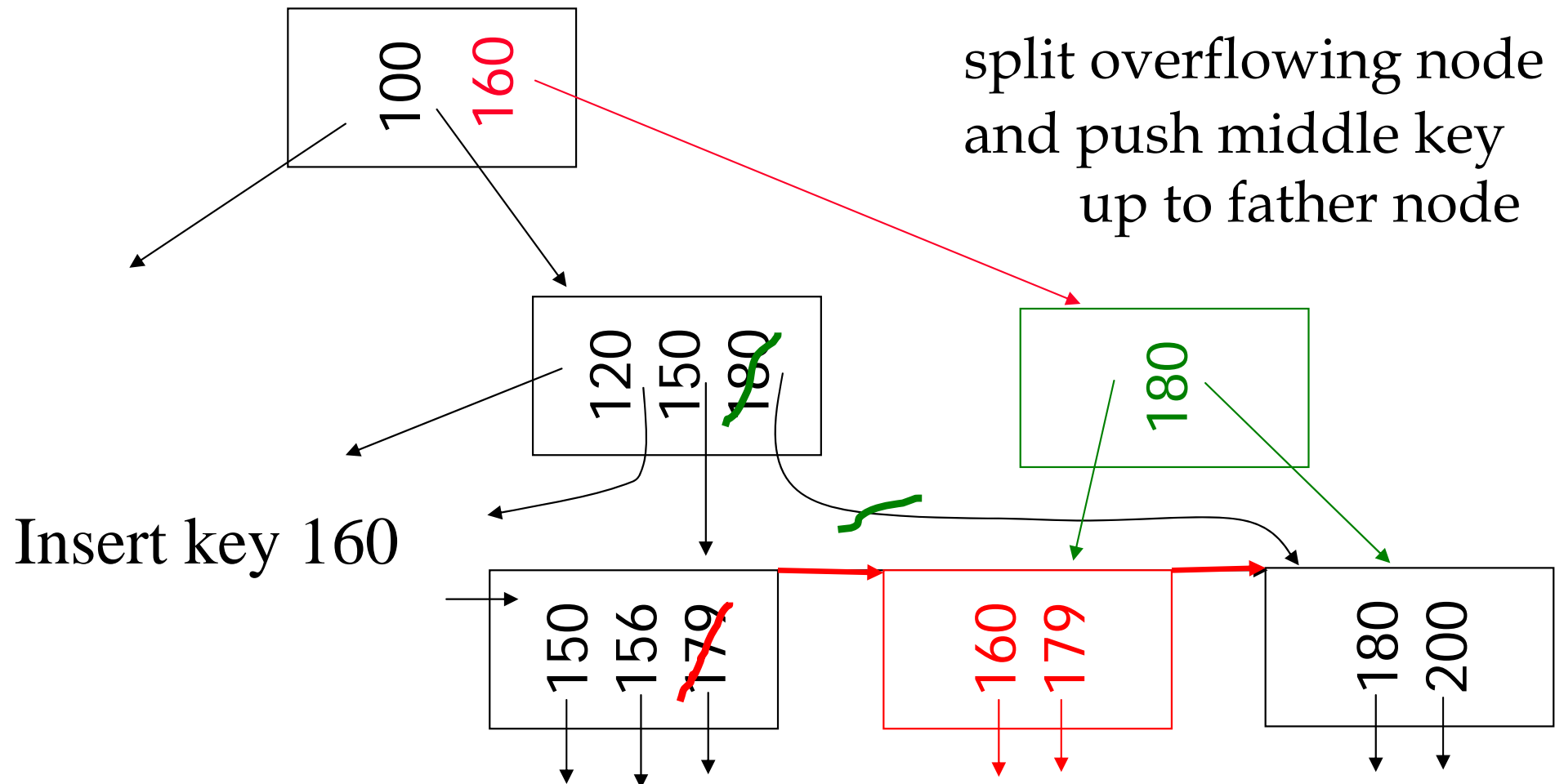
- **Leaf overflow**
  split overflowing node into two of (almost) same size and copy middle (separating) key to father node

100

n = 3

30   7

Insert key 7

3   5

8   5   7   11

30   31

*Insentions*

- Non-leaf overflow

split overflowing node
and push middle key
up to father node

Insert key 160

100 **160**

120 150 180

180

150 156 179

160 179

180 200

*Insertions*

■ New root

Insert key 45

new root

30

split can propagate
up to the root and
result in new root

10 20 30

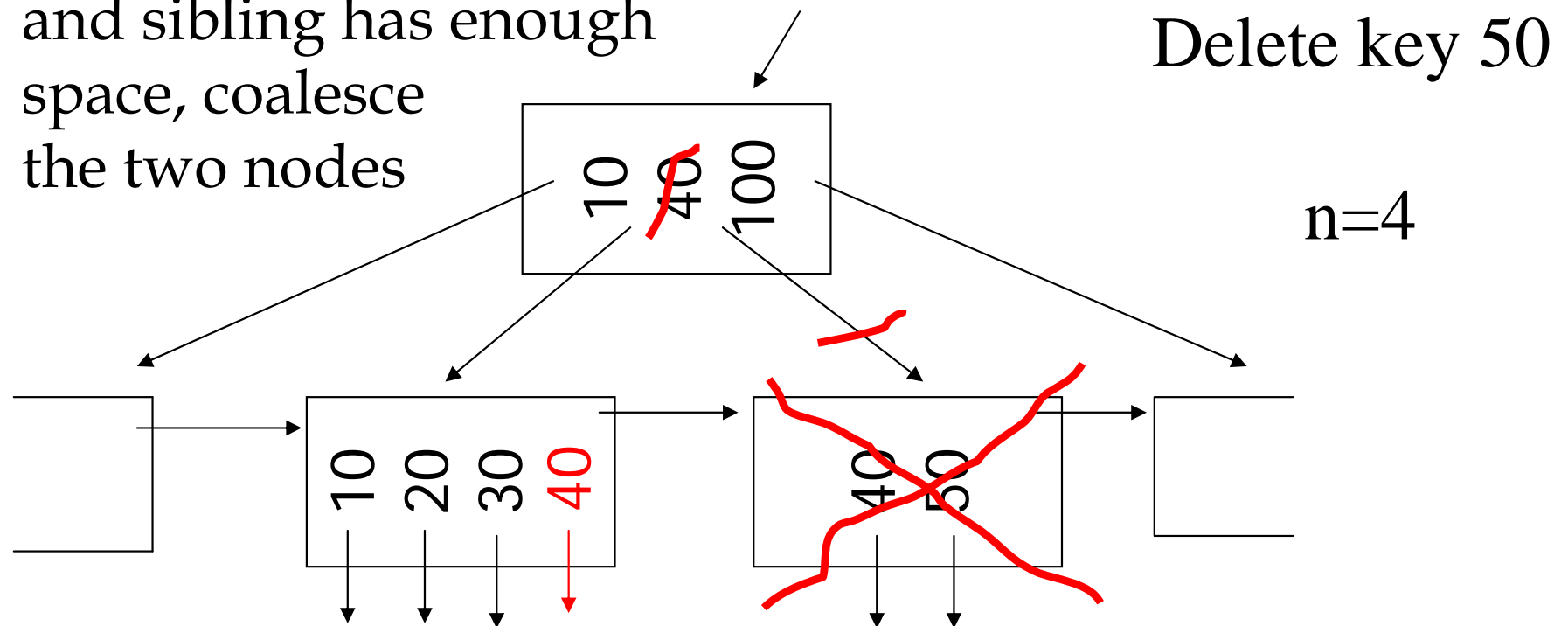40

1 2 3

10 12

20 25

30 32 40

40 45

# B+ Trees

*Deletions*

- Locate corresponding leaf node.

- Delete specified entry.

- Four different cases:

  - Leaf node has still enough entries,

  - Coalesce with neighbor (sibling),

  - Re-distribute keys,

  - Coalesce or re-distribute at non-leaf.

*Deletions*

- Coalesce with neighbor (sibling)

  if node underflows
  and sibling has enough
  space, coalesce
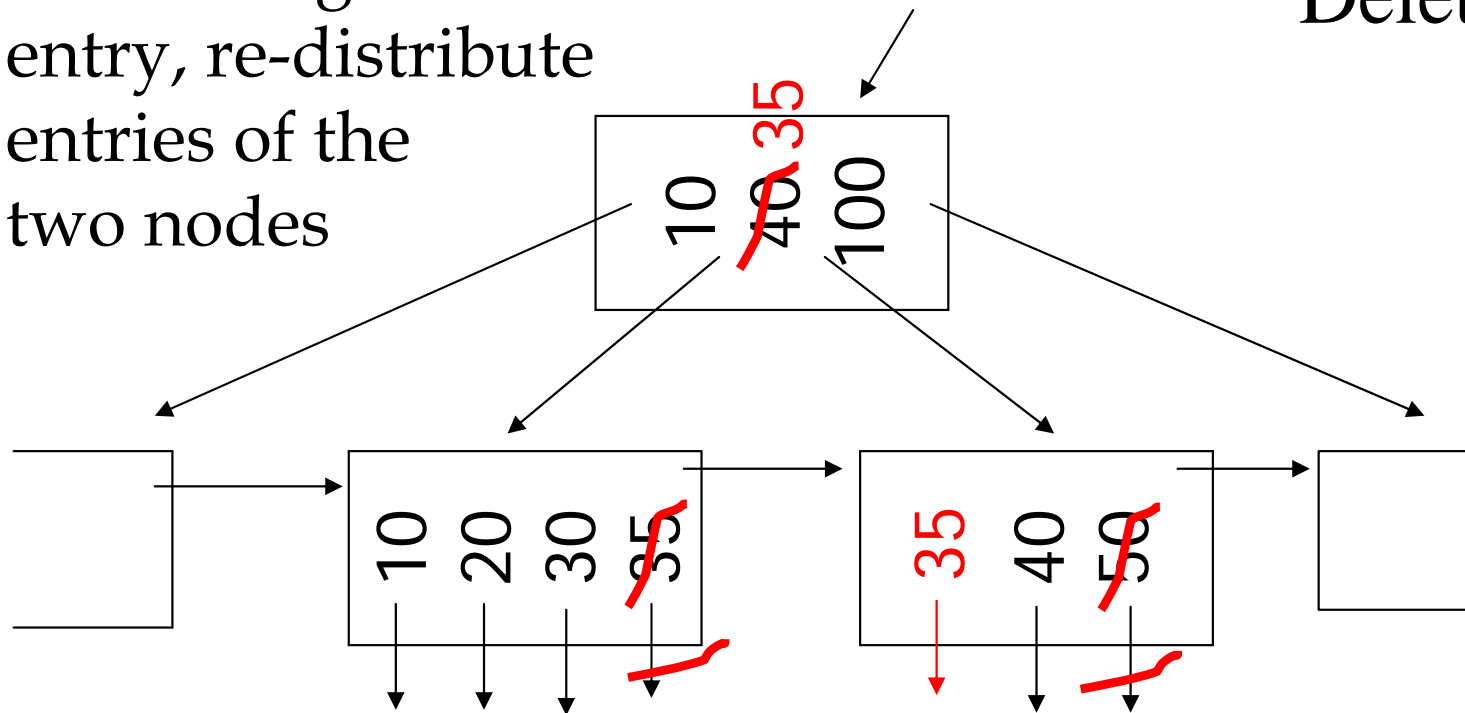  the two nodes

Delete key 50

n=4

# B+ Trees

*Deletions*

- Redistribute keys
  if node underflows
  and sibling has extra
  entry, re-distribute
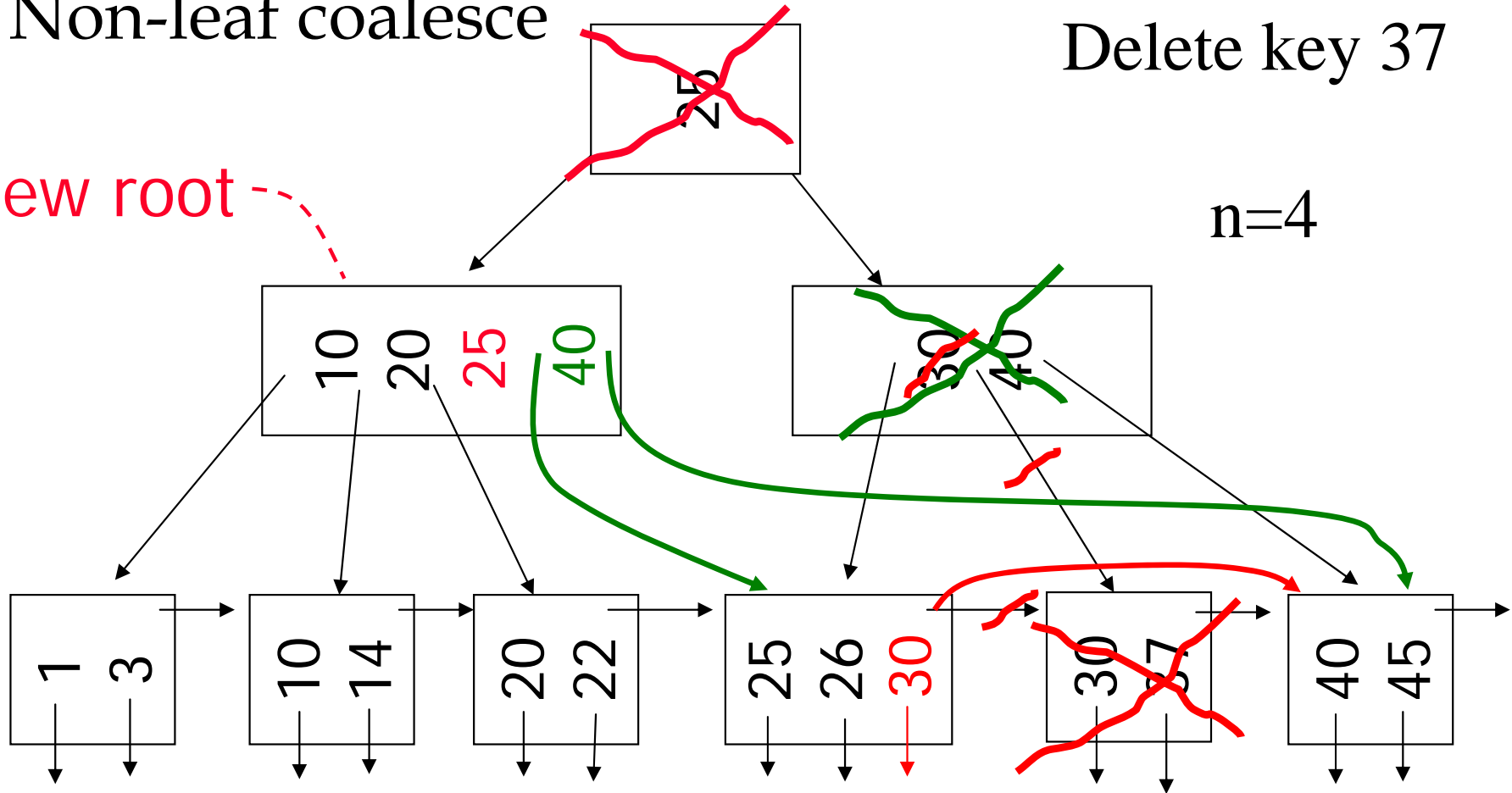  entries of the
  two nodes

Delete key 50

n=4

*Deletions*

- Non-leaf coalesce

Delete key 37

new root

n=4

# B+ Trees

*B+ Trees in Practice*

- Often, coalescing is not implemented.
- It is too hard and typically does not gain a lot of performance.

# B+ Trees

*B+ Trees in Practice*

- Typical order: 200, typical space utilization: 67%, i.e., average fanout = 133.

- Typical capacities:
  - Height 4: $133^4$ = 312,900,700 records,
  - Height 3: $133^3$ =    2,352,637 records.

<div style="border:1px solid red; color:red;">
Fanout: the number of pointers in a node
</div>

- Can often hold top levels in buffer pool:
  - Level 1 =        1 blocks  =    8 Kbytes,
  - Level 2 =     133 blocks =    1 Mbyte,
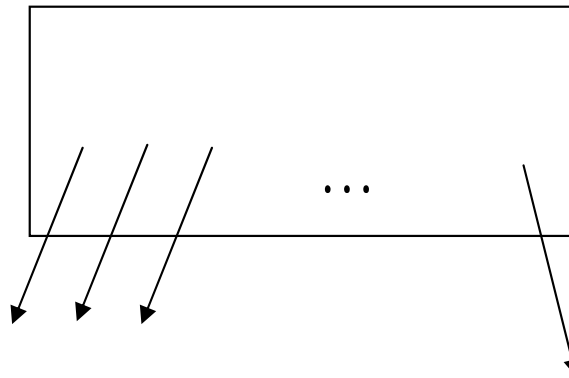  - Level 3 = 17,689 blocks = 133 Mbytes.

# B+ Trees

*B+ Trees in Practice*

- Order (n) concept replaced by physical space criterion in practice ('*at least half-full*').

- Inner nodes can typically hold many more entries than leaf nodes.

- Variable sized records and search keys mean different nodes will contain different numbers of entries.

- Even with fixed length fields, multiple records with the same search key value (duplicates) can lead to variable-sized data entries.

# Interesting Problem

For B+ tree, how large should $n$ be?
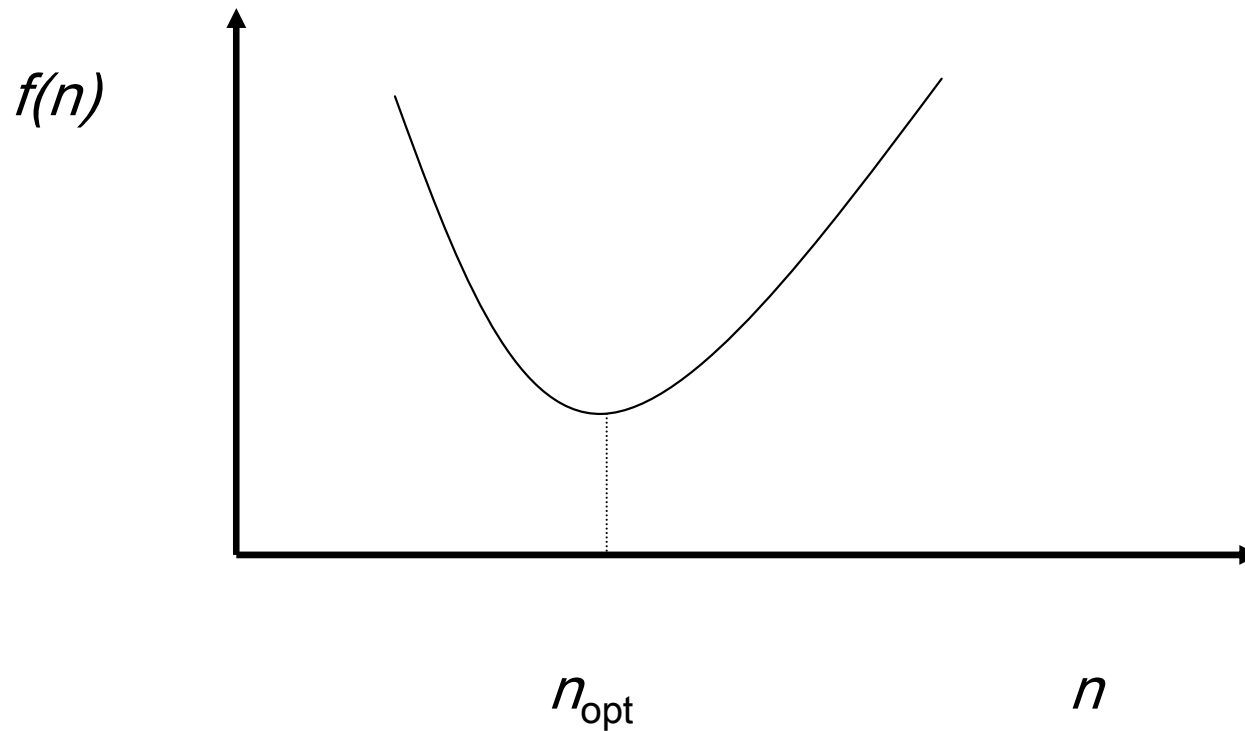


$n$ is number of keys / node

# Sample Assumptions

(1) Time to read node from disk is
$(S+Tn)$ msec.

(2) Once block in memory, use binary
search to locate key:
$(a + b \, LOG_2 \, n)$ msec.

For some constants $a,b$;   Assume a << S

(3) Assume B+ tree is full, i.e.,
# nodes to examine is $LOG_n \, N$
where $N$ = # records

# FIND $n_{opt}$ by $f'(n) = 0$

Answer is $n_{opt}$ = "few hundred"

↘ What happens to $n_{opt}$ as

- Disk gets faster?
- CPU get faster?