CS411
Database Systems
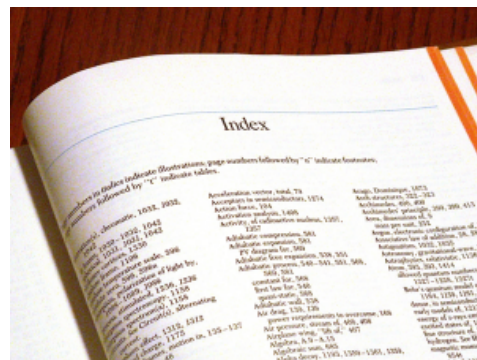
07: Indexing
14.1-14.3

## Why Do We Learn This?

## Q: What is "indexing"?

- What is an index?

- To build an index.

- To maintain an index

## What is an index?

## Indexes in databases

- An *index* on a file speeds up selections on the *search key field(s)*
- Search key = any subset of the fields of a relation
  - *Search key* is not the same as *key* (minimal set of fields that uniquely identify a record in a relation).
- Entries in an index: (k, r), where:
  - k = the key
  - r = the record OR record id OR record ids

## Some terminology

- *Data file*: has the data corresponding to a relation

- *Index file*: has the index

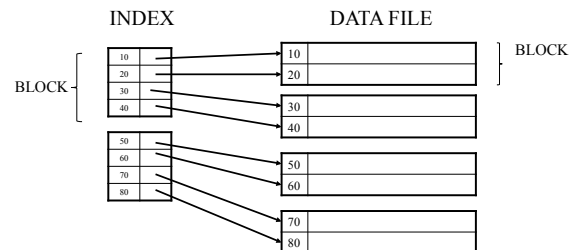- File consists of smaller units called "blocks" (e.g. of size 4 KB or 8 KB)

## Types of Indexes

- Clustered/unclustered
  - Clustered = records sorted in the key order
  - Unclustered = no
- Dense/sparse
  - Dense = each record has an entry in the index
  - Sparse = only some records have
- Primary/secondary
  - Primary = on the primary key
  - Secondary = on any key
  - Some textbooks interpret these differently
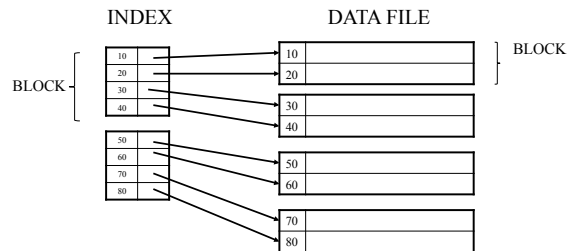- B+ tree / Hash table / …

## Ex: Clustered, Dense Index

- Clustered: File is sorted on the index attribute
- *Dense*: sequence of (key,pointer) pairs

## Ex: Clustered, Dense Index

- Clustered: File is sorted on the index attribute
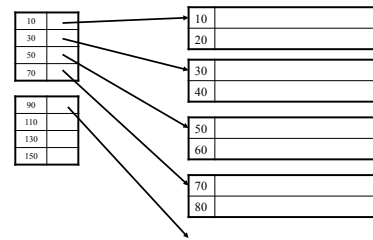- *Dense*: sequence of (key,pointer) pairs

INDEX                    DATA FILE

BLOCK

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

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

| 10 |
| 20 |

| 30 |
| 40 |

| 50 |
| 60 |

| 70 |
| 80 |

BLOCK

- Index blocks are much smaller than data blocks. Index may even fit into main memory

9

---

## Clustered, Sparse Index

- *Sparse* index: one key per data block, corresponding to the lowest search key in that block
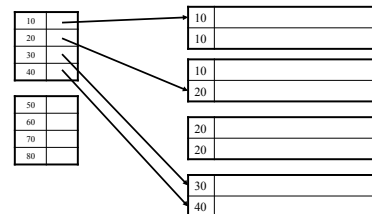
| 10 |
| 30 |
| 50 |
| 70 |

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

| 10 |
| 20 |

| 30 |
| 40 |

| 50 |
| 60 |

| 70 |
| 80 |

10

---

## What if there are duplicate keys?

11

---

## Clustered Index with Duplicate Keys

- Dense index: point to the first record with that key

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

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

| 10 |
| 10 |

| 10 |
| 20 |

| 20 |
| 20 |

| 30 |
| 40 |

12

3

## Clustered Index with Duplicate Keys

- Sparse index: pointer to lowest search key in each block:

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

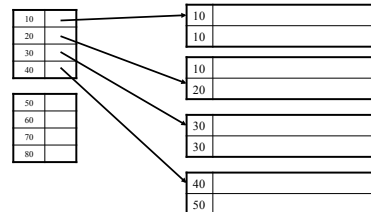| 10 |
| 10 |

| 10 |
| 20 |

| 20 |
| 20 |

| 30 |
| 40 |

- OK?

  Try search for 20

13

## Clustered Index with Duplicate Keys

- Better: pointer to lowest new search key in each block:
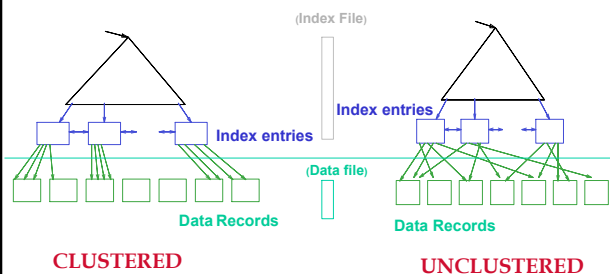- Search for 20

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

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

| 10 |
| 10 |

| 10 |
| 20 |

| 30 |
| 30 |

| 40 |
| 50 |

14

## Unclustered Indexes

- Often for indexing other attributes than primary key
- Always dense (why ?)

| 10 |
| 10 |
| 20 |
| 20 |

| 20 |
| 30 |
| 30 |
| 30 |

| 20 |
| 30 |

| 30 |
| 20 |

| 10 |
| 20 |

| 10 |
| 30 |

15

## Summary Clustered vs. Unclustered Index



**(Index File)**

**Index entries**

**Index entries**

**(Data file)**

**Data Records**

**Data Records**

**CLUSTERED**

**UNCLUSTERED**

16

4

# B+ Trees

---

## B-Trees/B+Trees: B__?__?__ Trees

- Intuition:
  - The index can be very large.
  - Index of index?
  - Index of index of index?
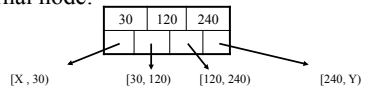  - How best to create such a multi-level index?

- B+trees:
  - Textbook refers to B+trees (a popular variant) as B-trees (as most people do)
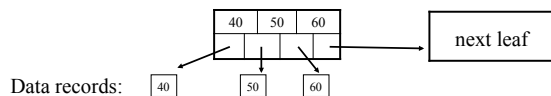  - Distinction will be clear later

---

## B+ Trees Basics

- Parameter d = the *degree*
- Each node has [d, 2d] keys (except root)
  - Internal node:
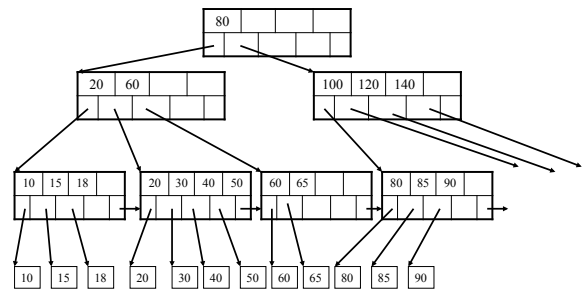
| 30 | 120 | 240 |
|----|-----|-----|

[X , 30)   [30, 120)   [120, 240)   [240, Y)

  - Leaf:

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

→ next leaf

Data records:   40    50    60

---

## B+ Tree Example

d = 2

| 80 |
|----|

| 20 | 60 |
|----|----|

| 100 | 120 | 140 |
|-----|-----|-----|

| 10 | 15 | 18 |    | 20 | 30 | 40 | 50 |    | 60 | 65 |    | 80 | 85 | 90 |

10   15   18    20   30   40   50    60   65    80   85   90

---

## B+ Tree Design

- How large is d ?
- Example:
  - Key size = 4 bytes
  - Pointer size = 8 bytes
  - Block size = 4096 byes
- $2d \times 4 + (2d+1) \times 8 \leq 4096$
- d = 170

## Searching a B+ Tree

- Exact key values:
  - Start at the root
  - Proceed down, to the leaf

Select name
From people
Where age = 25

- Range queries:
  - As above
  - Then sequential traversal
  - This is where the "next leaf"
    pointer is useful

Select name
From people
Where 20 <= age
   and  age <= 30

## Some applications of B+ trees

1. Search key is primary key; index is dense. Data file may or may not be sorted by key.

2. Data file is sorted by primary key; index is sparse.

## B+ Trees in Practice

- Typical d: 100.  Typical fill-factor: 67%.
  - average "fanout" = 133

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

- Can often hold top levels in buffer pool:
  - Level 1 =          1 page  =      8 Kbytes
  - Level 2 =       133 pages =     1 Mbyte
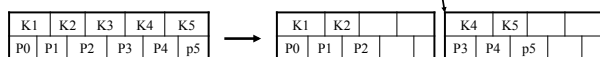  - Level 3 = 17,689 pages = 133 MBytes

## Insertion in a B+ Tree

Assume dense index.

Insert (K, P)
- Find leaf where K belongs, insert
- If no overflow (2d keys or less), halt
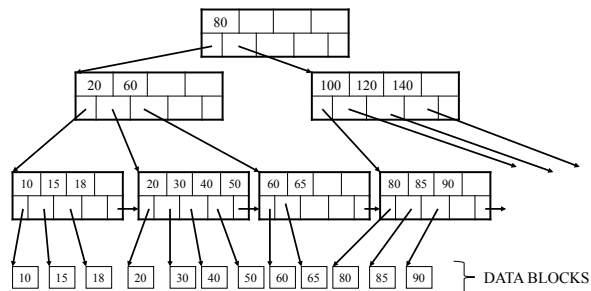- If overflow (2d+1 keys), split node, insert in parent:

(K3,   ) to parent

| K1 | K2 | K3 | K4 | K5 |
|----|----|----|----|----|
| P0 | P1 | P2 | P3 | P4 | p5 |

→

| K1 | K2 |  |  |
|----|----|--|--|
| P0 | P1 | P2 |  |  |

| K4 | K5 |  |
|----|----|--|
| P3 | P4 | p5 |  |

- If leaf, keep K3 too in right node
- When root splits, new root has 1 key only
  - that's why root is special for degree satisfaction

25

---
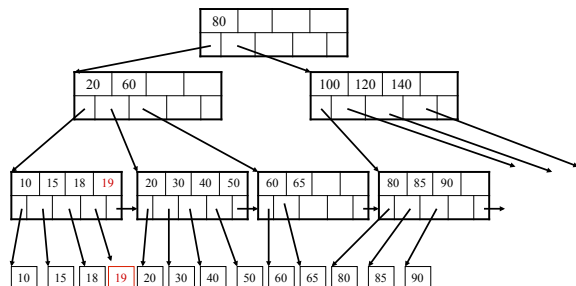
## Insertion in a B+ Tree

Assume d=2.

Insert K=19



DATA BLOCKS

26

---

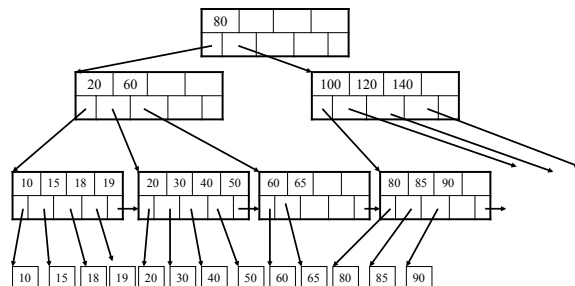## Insertion in a B+ Tree

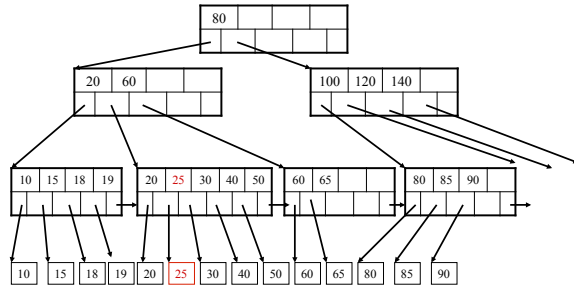After insertion



27

---

## Insertion in a B+ Tree

Now insert 25



28

7

## Insertion in a B+ Tree
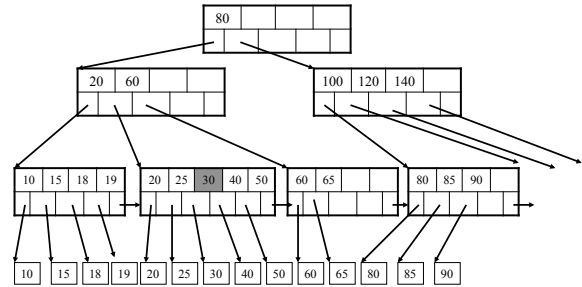
After insertion



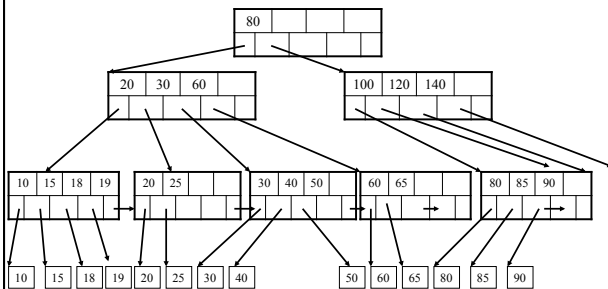## Insertion in a B+ Tree
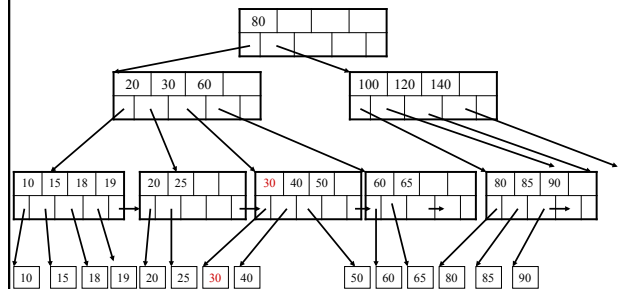
But now have to split !



## Insertion in a B+ Tree

After the split



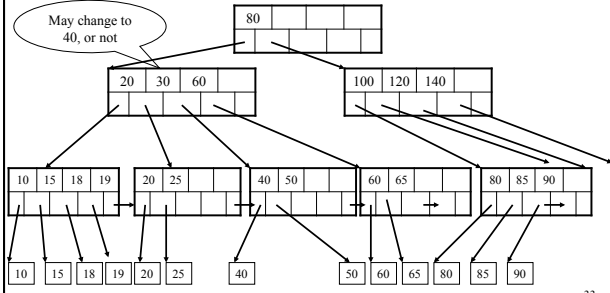## Deletion from a B+ Tree

Delete 30



29

30

31

32

8

## Deletion from a B+ Tree
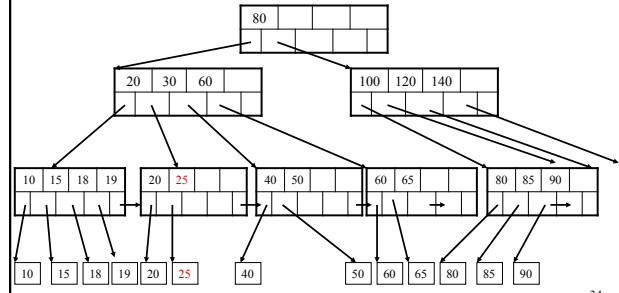
After deleting 30

May change to 40, or not



33

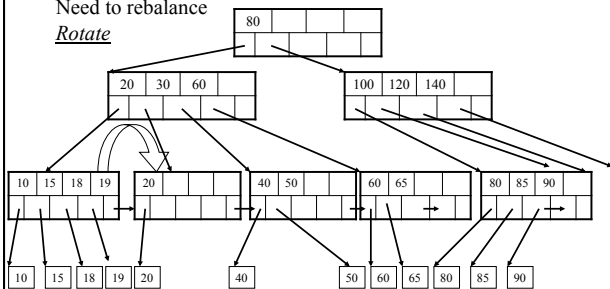## Deletion from a B+ Tree

Now delete 25



34

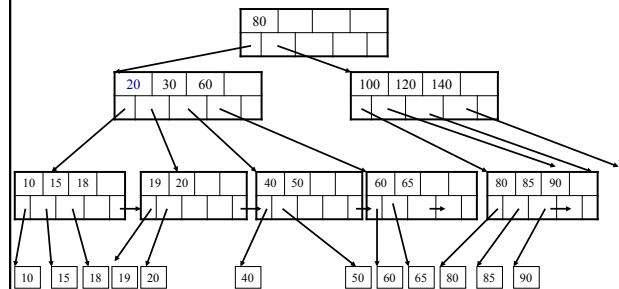## Deletion from a B+ Tree

After deleting 25
Need to rebalance
*Rotate*

Rotation in general can involve either sibling, but here only the left sibling can help
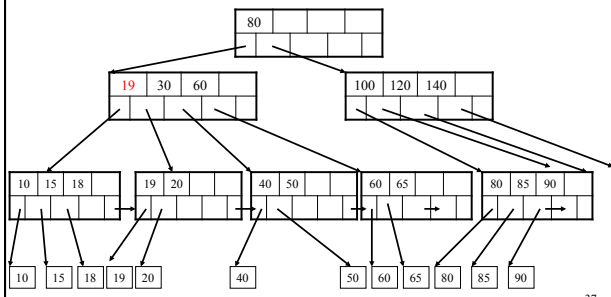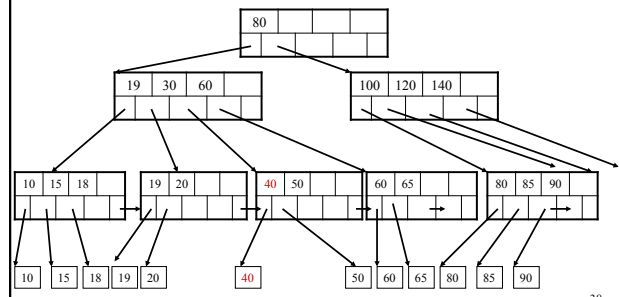


35

## Deletion from a B+ Tree
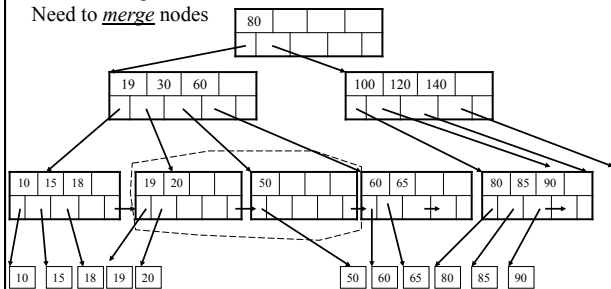


36

Deletion from a B+ Tree



Deletion from a B+ Tree

Now delete 40



Deletion from a B+ Tree
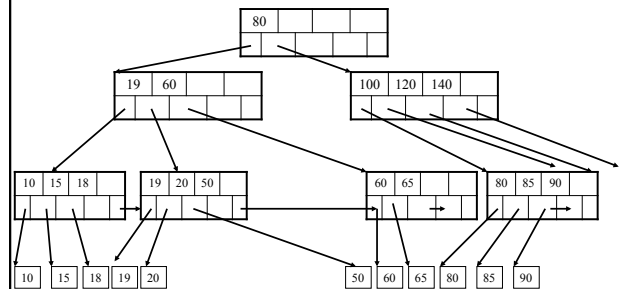
After deleting 40
Rotation not possible
Need to *merge* nodes



Deletion from a B+ Tree

Final tree

# Hash Tables

41

---

## Hash Tables

- Secondary storage hash tables are much like main memory ones
- Recall basics:
  - There are n *buckets*
  - A hash function f(k) maps a key k to {0, 1, …, n-1}
  - Store in bucket f(k) a pointer to record with key k
- Secondary storage: bucket = block
  - Store in bucket f(k) any record with key k
  - use overflow blocks when needed

42

---

## Hash Table Example

- Assume 1 bucket (block) stores 2 records
- $h(e)=0$
- $h(b)=h(f)=1$
- $h(g)=2$
- $h(a)=h(c)=3$

```
   ┌─────────┐
 0 │ e       │
   ├─────────┤
 1 │ b       │
   │ f       │
 2 │ g       │
   │         │
 3 │ a       │
   │ c       │
   └─────────┘
```

43

---

## Searching in a Hash Table

- Search for a:
- Compute $h(a)=3$   HOW ?
- Read bucket 3
- 1 disk access

Main memory may have an array of pointers (to buckets) accessible by bucket number.

```
   ┌─────────┐
 0 │ e       │
   ├─────────┤
 1 │ b       │
   │ f       │
 2 │ g       │
   │         │
 3 │ a       │
   │ c       │
   └─────────┘
```

44

11

## Insertion in Hash Table

- Place in right bucket, if space
- E.g. h(d)=2

| | |
|---|---|
| 0 | e |
| 1 | b<br>f |
| 2 | g<br>d |
| 3 | a<br>c |

45

## Insertion in Hash Table

- Create overflow block, if no space
- E.g. h(k)=1

| | | |
|---|---|---|
| 0 | e | |
| 1 | b<br>f | → k |
| 2 | g<br>d | |
| 3 | a<br>c | |

- More over-flow blocks may be needed

46

## Hash Table Performance

- Excellent, if no overflow blocks
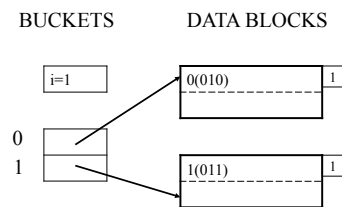
- Degrades considerably when many overflow blocks.

47

## Extensible Hash Table

- Allows hash table to grow, to avoid performance degradation
- Assume a hash function h that returns numbers in $\{0, \ldots, 2^k - 1\}$
- Start with $n = 2^i << 2^k$ , only look at first i most significant bits
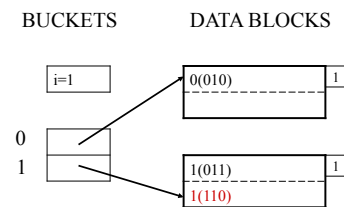
48

## Extensible Hash Table

- E.g. i=1, n=2, k=4

  BUCKETS          DATA BLOCKS

  i=1              0(010)        1

  0
  1                1(011)        1

- Note: we only look at the first bit (0 or 1)

49

## Insertion in Extensible Hash Table

- Insert 1110

  BUCKETS          DATA BLOCKS

  i=1              0(010)        1

  0
  1                1(011)        1
                   1(110)

50

## Insertion in Extensible Hash Table

- Now insert 1010

  BUCKETS          DATA BLOCKS

  i=1              0(010)        1

  0
  1                1(011)        1
                   1(110), 1(010)
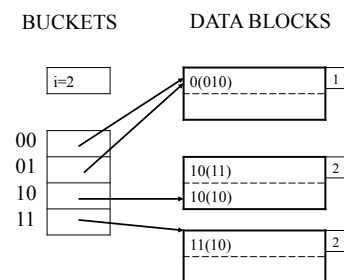
- Need to split block and extend bucket array
- i becomes 2

51

## Insertion in Extensible Hash Table

- Now insert 1010 (cont.)

  BUCKETS          DATA BLOCKS

  i=2              0(010)        1

  00
  01               10(11)        2
  10               10(10)
  11
                   11(10)        2

52

13

## Insertion in Extensible Hash Table

- Now insert 0000, then 0101

BUCKETS           DATA BLOCKS

```
i=2                    0(010)              1
                       0(000), 0(101)

00
01                     10(11)              2
10                     10(10)
11
                       11(10)              2
```
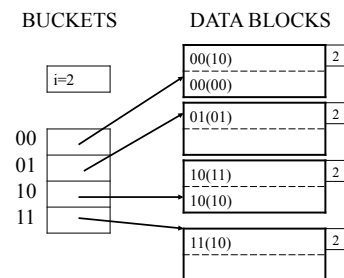
- Need to split block, but not bucket array

53

## Insertion in Extensible Hash Table

- After splitting the block

BUCKETS           DATA BLOCKS

```
                       00(10)              2
i=2                    00(00)

                       01(01)              2
00
01                     10(11)              2
10                     10(10)
11
                       11(10)              2
```

54

## Performance Extensible Hash Table

- No overflow blocks: access always one read
- BUT:
  - Extensions can be costly and disruptive
  - After an extension table may no longer fit in memory
  - Imagine three records whose keys share the first 20 bits. These three records cannot be in same block (assume two records per block). But a block split would require setting i = 20, i.e., accommodating for $2^{20}$ = 1 million buckets, even though there may be only a few hundred records.
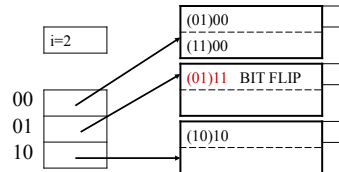
55

## Linear Hash Table

- Idea: add only one bucket at a time
- Problem: n = no longer a power of 2
- Let i be #bits necessary to address n buckets.
  - $2^{i-1} < n <= 2^i$
- After computing h(k), use last i bits:
  - If last i bits represent a number >= n, change msb from 1 to 0 (get a number < n)
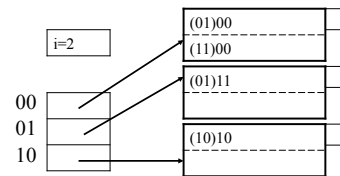- Also, allow overflow blocks
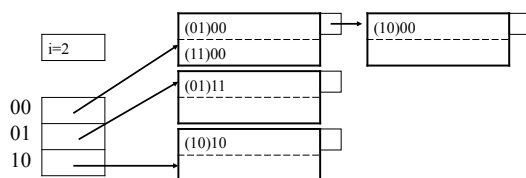
56

14

## Linear Hash Table Example

- N=3

i=2

00
01
10

(01)00
(11)00

(01)11   BIT FLIP

(10)10

57

## Linear Hash Table Example

- Insert 1000:

i=2

00
01
10

(01)00
(11)00

(01)11

(10)10

58

## Linear Hash Table Example

- Insert 1000: overflow blocks…

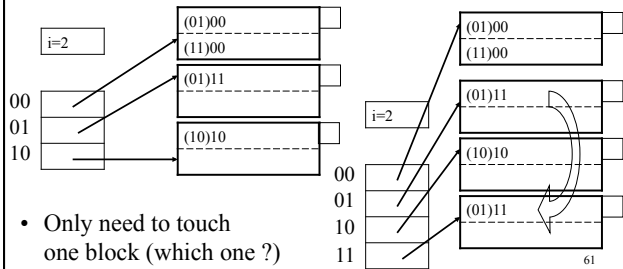i=2

00
01
10

(01)00
(11)00

(10)00

(01)11

(10)10

59

## Linear Hash Tables

- Extend n:=n+1 when average number of records per block exceeds (say) 80%

60

15

## Linear Hash Table Extension

- From n=3 to n=4

i=2

(01)00
(11)00

(01)11

(10)10

00
01
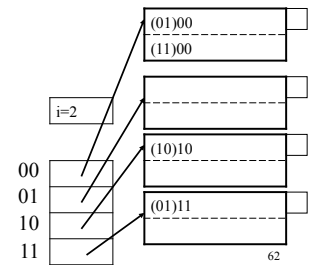10

(01)00
(11)00

(01)11

(10)10

(01)11

i=2

00
01
10
11

- Only need to touch one block (which one ?)

61

## Linear Hash Table Extension

- From n=3 to n=4 finished

- Extension from n=4 to n=5 (new bit)
- Need to touch every single block (why ?)

(01)00
(11)00

(10)10

(01)11

i=2

00
01
10
11

62

- See examples in text.

63