# Hash-Based Indexes

## Chapter 11
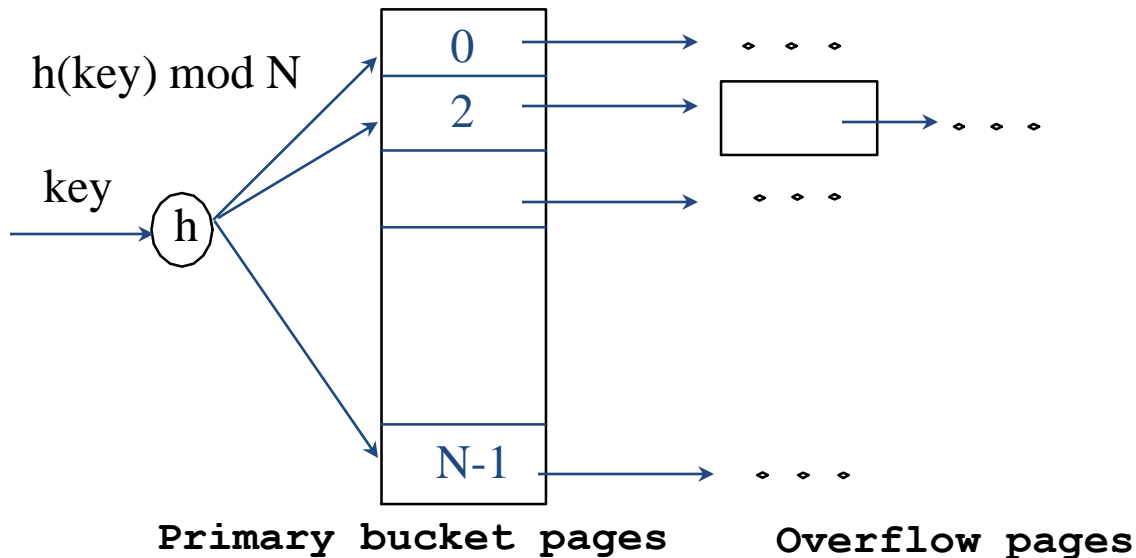
## Sina Meraji

# Introduction

- *As for any index, 3 alternatives for data entries* **k\***:
    - Data record with key value **k**
    - <**k**, rid of data record with search key value **k**>
    - <**k**, list of rids of data records with search key **k**>
    - Choice orthogonal to the *indexing technique*
- *Hash-based* indexes are best for *equality selections*. ***Cannot*** support range searches.
- Static and dynamic hashing techniques exist;

# Static Hashing

- # primary pages fixed, allocated sequentially, never de-allocated; overflow pages if needed.
- **h**(*k*) mod N = bucket to which data entry with key *k* belongs. (N = # of buckets)

h(key) mod N

key → h →

| 0 |
| 2 |
|   |
|   |
| N-1 |

**Primary bucket pages**    **Overflow pages**
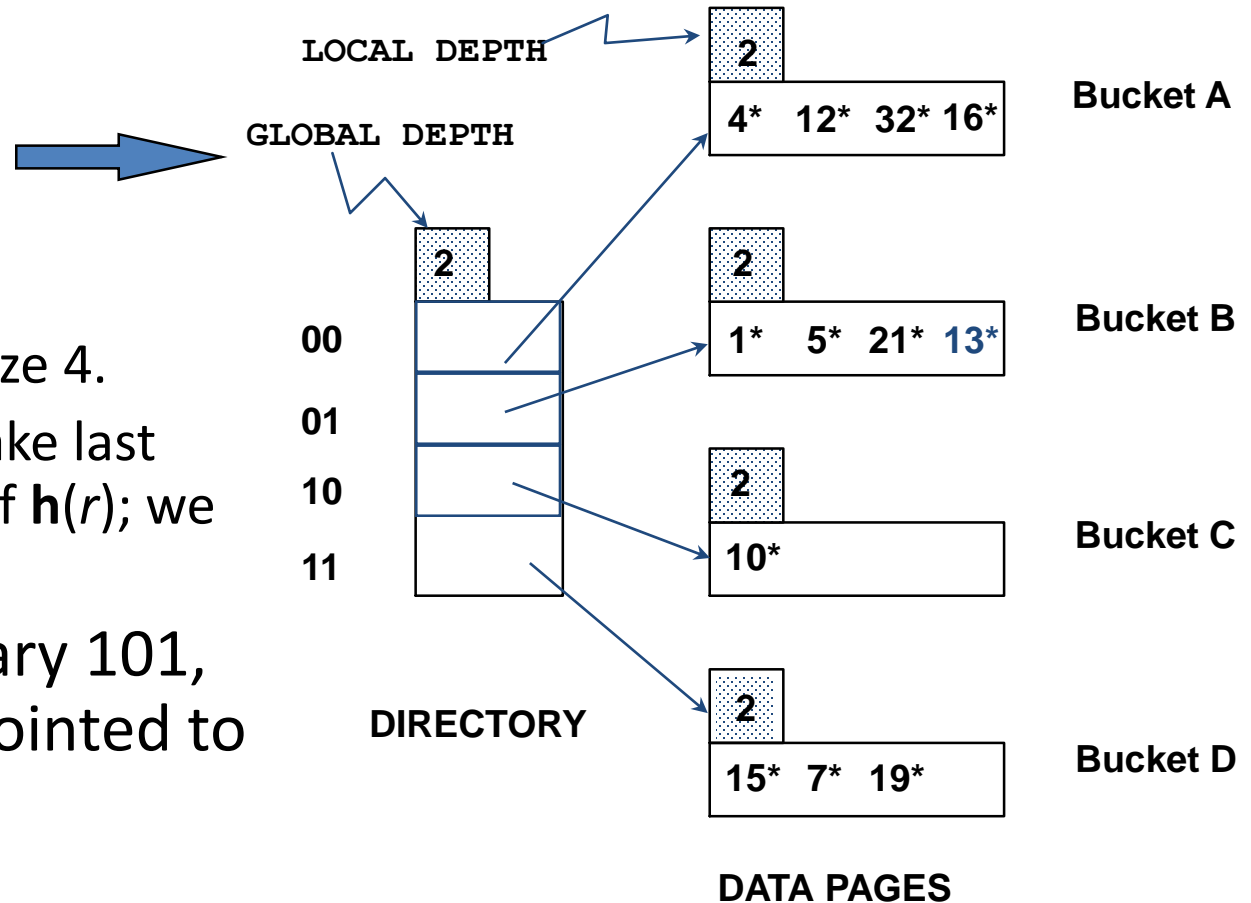
# Static Hashing (Contd.)

- Buckets contain *data entries*.
- Hash fn works on *search key* field of record *r.*  Must distribute values over range 0 … M-1.
  - **h**(*key*) = (a * *key* + b) usually works well.
  - a and b are constants;  lots known about how to tune **h**.
- Long overflow chains can develop and degrade performance.
  - *Extendible* and *Linear Hashing*: Dynamic techniques to fix this problem.

# Extendible Hashing

- Situation: Bucket (primary page) becomes full. Why not re-organize file by *doubling* # of buckets?
  - Reading and writing all pages is expensive!
  - *Idea*:  Use *directory of pointers to buckets*, double # of buckets by *doubling the directory,* splitting just the bucket that overflowed!
  - Directory much smaller than file, so doubling it is much cheaper.  Only one page of data entries is split.  *No overflow page*!
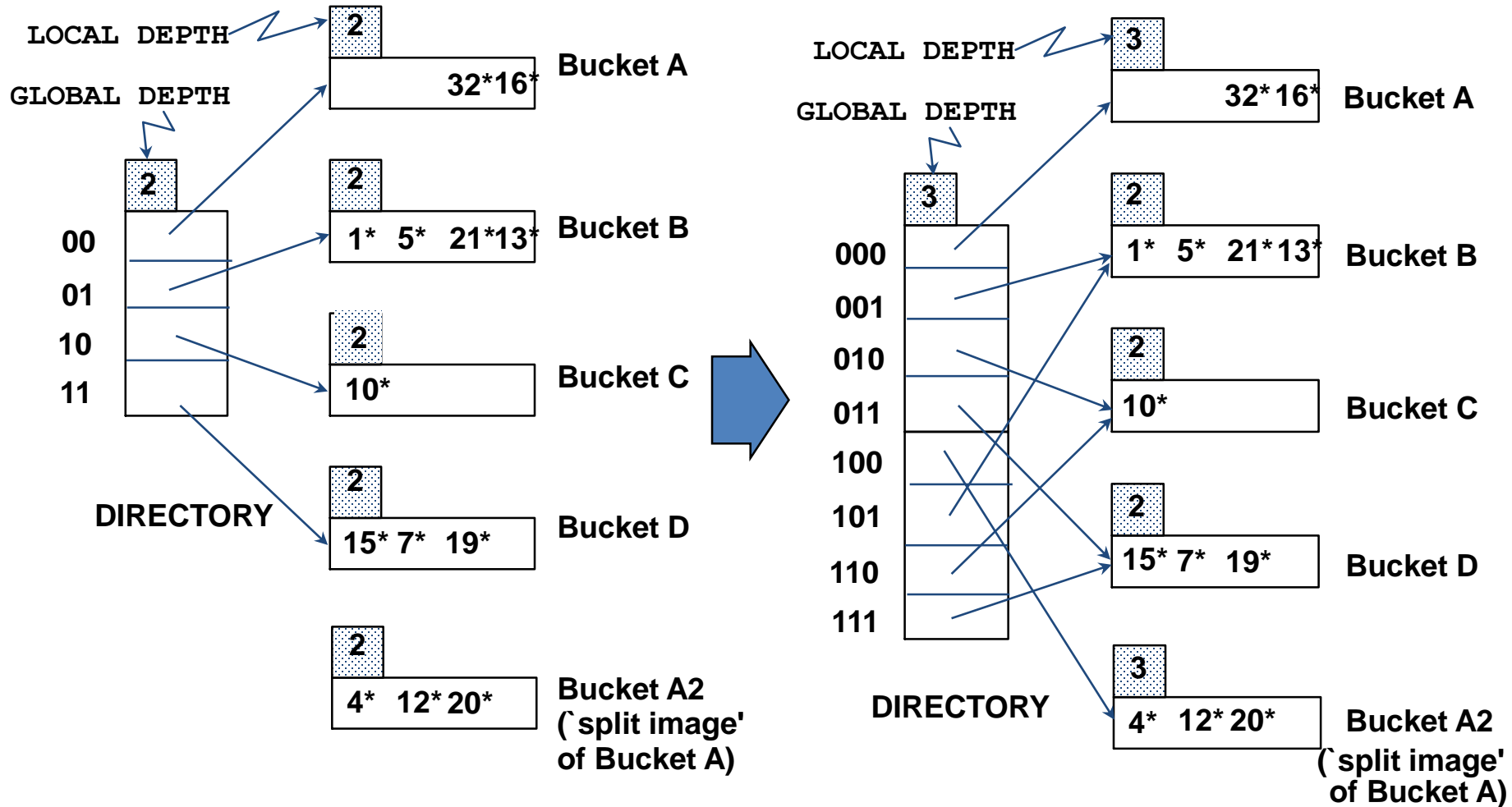  - Trick lies in how hash function is adjusted!

# Example



LOCAL DEPTH

GLOBAL DEPTH

- Directory is array of size 4.
- To find bucket for *r*, take last `*global depth*' # bits of **h**(*r*); we denote *r* by **h**(*r*).

    - If **h**(*r*) = 5 = binary 101, it is in bucket pointed to by 01.

DIRECTORY

DATA PAGES

**2**

**4\*  12\*  32\* 16\***  Bucket A

**2**

**1\*   5\*  21\* 13\***  Bucket B

**2**

**10\***  Bucket C

**2**

**15\*  7\*  19\***  Bucket D

00
01
10
11

❖ **Insert**:  If bucket is full, *split* it (*allocate new page, re-distribute*).

❖ *If necessary*, double the directory.  (As we will see, splitting a bucket does not always require doubling; we can tell by comparing *global depth* with *local depth* for the split bucket.)

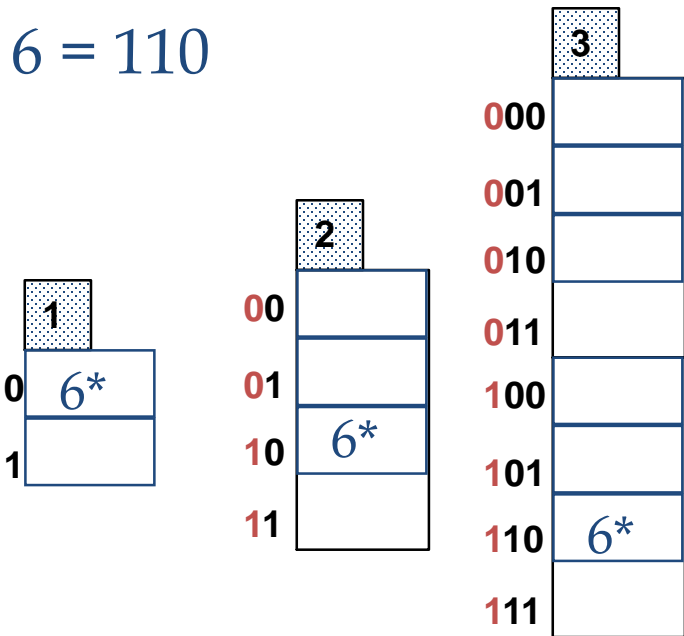# Insert **h**(r)=20 (Causes Doubling)
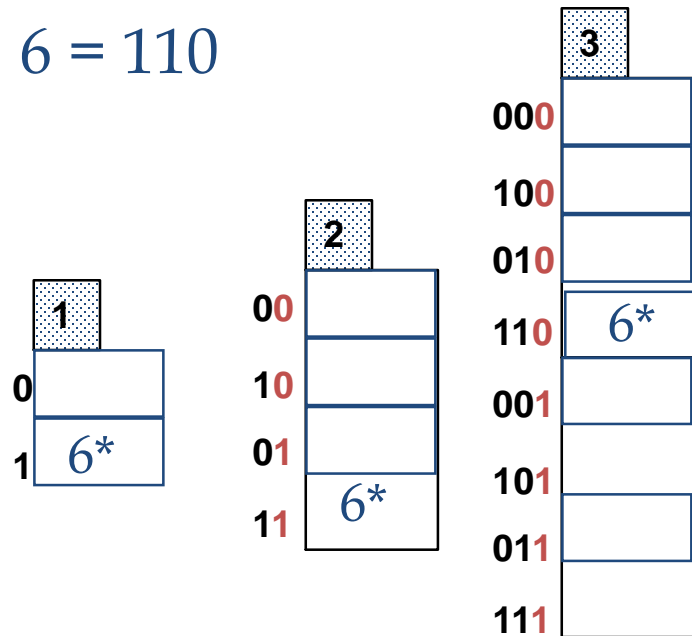
# Points to Note

- 20 = binary 10100.  Last **2** bits (00) tell us *r* belongs in A or A2.  Last **<u>3</u>** bits needed to tell which.
  - *Global depth of directory*:  Max # of  bits needed to tell which bucket an entry belongs to.
  - *Local depth of a bucket*: # of bits used to determine if an entry belongs to this bucket.
- When does bucket split cause directory doubling?
  - Before insert, *local depth* of bucket = *global depth*.  Insert causes *local depth* to become > *global depth*; directory is doubled by *copying it over* and `fixing' pointer to split image page.  (Use of least significant bits enables efficient doubling via copying of directory!)

# Directory Doubling

Why use least significant bits in directory?
⇔ Allows for doubling via copying!

6 = 110

| | |
|---|---|
| **1** | |
| 0 | 6* |
| 1 | |

| | |
|---|---|
| **2** | |
| 00 | |
| 01 | |
| 10 | 6* |
| 11 | |

| | |
|---|---|
| **3** | |
| 000 | |
| 001 | |
| 010 | |
| 011 | |
| 100 | |
| 101 | |
| 110 | 6* |
| 111 | |

6 = 110

| | |
|---|---|
| **1** | |
| 0 | |
| 1 | 6* |

| | |
|---|---|
| **2** | |
| 00 | |
| 10 | |
| 01 | |
| 11 | 6* |

| | |
|---|---|
| **3** | |
| 000 | |
| 100 | |
| 010 | |
| 110 | 6* |
| 001 | |
| 101 | |
| 011 | |
| 111 | |

Least Significant      vs.      Most Significant

# Comments on Extendible Hashing

- If directory fits in memory, equality search answered with one disk access; else two.
    - 100MB file, 100 bytes/rec, 4K pages contains 1,000,000 records (as data entries) and 25,000 directory elements; chances are high that directory will fit in memory.
    - Directory grows in spurts, and, if the distribution *of hash values* is skewed, directory can grow large.
    - Multiple entries with same hash value cause problems!
- **<u>Delete</u>**:  If removal of data entry makes bucket empty, can be merged with `split image'.  If each directory element points to same bucket as its split image, can halve directory.
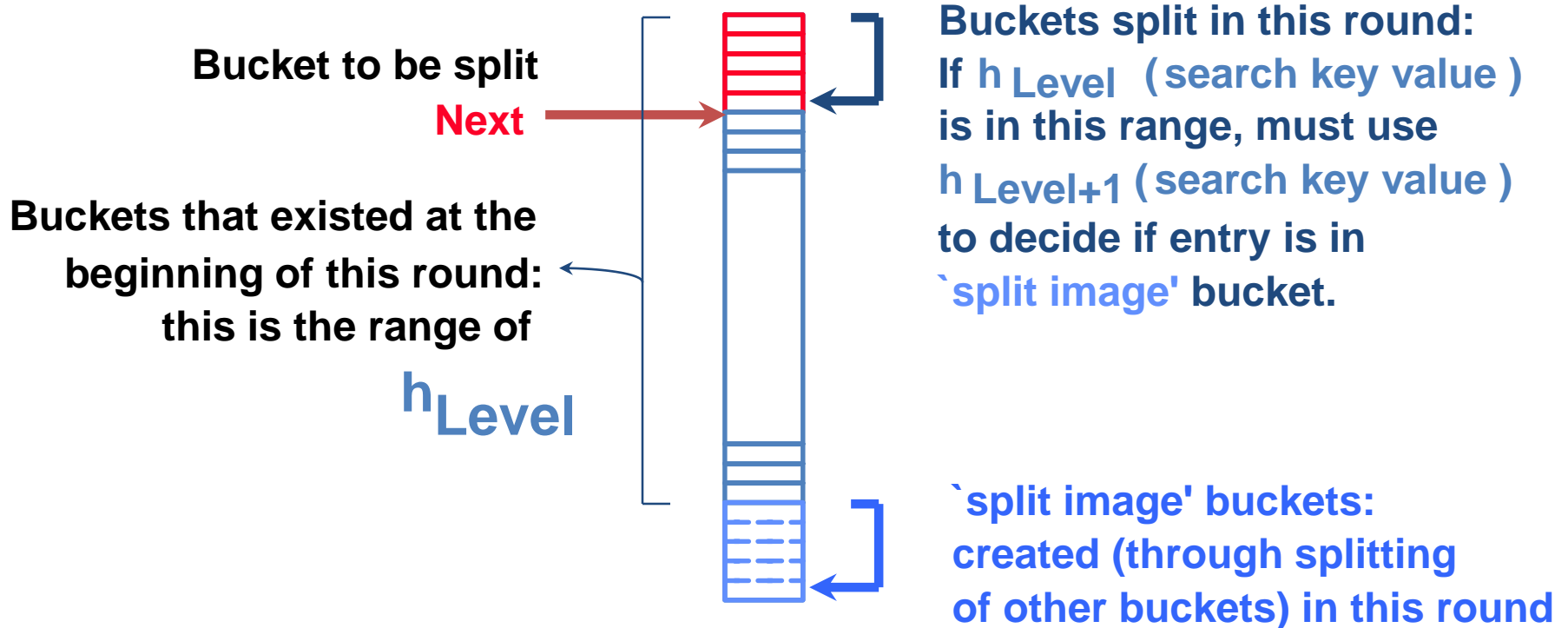
# Linear Hashing

- This is another dynamic hashing scheme, an alternative to Extendible Hashing.

- LH handles the problem of long overflow chains without using a directory, and handles duplicates.

- _Idea_:  Use a family of hash functions $\mathbf{h}_0$, $\mathbf{h}_1$, $\mathbf{h}_2$, …
    - $\mathbf{h}_i(key) = \mathbf{h}(key) \bmod(2^i N)$;  N = initial # buckets
    - $\mathbf{h}$ is some hash function (range is _not_ 0 to N-1)
    - If N = $2^{d0}$, for some _d0_, $\mathbf{h}_i$ consists of applying $\mathbf{h}$ and looking at the last _di_ bits, where _di = d0 + i._
    - $\mathbf{h}_{i+1}$ doubles the range of $\mathbf{h}_i$ (similar to directory doubling)

# Linear Hashing (Contd.)

- Directory avoided in LH by using overflow pages, and choosing bucket to split round-robin.

  - Splitting proceeds in `rounds'.  Round ends when all $N_R$ initial (for round $R$) buckets are split.  Buckets 0 to *Next-1* have been split;  *Next* to $N_R$ yet to be split.

  - Current round number is *Level*.

  - **Search:** To find bucket for data entry *r,* find $\mathbf{h}_{Level}(r)$*:*
    - If $\mathbf{h}_{Level}(r)$ in range `*Next* to $N_R$' , *r* belongs here.
    - Else, r could belong to bucket $\mathbf{h}_{Level}(r)$ or bucket $\mathbf{h}_{Level}(r) + N_R$; must apply $\mathbf{h}_{Level+1}(r)$ to find out.

# Overview of LH File

- In the middle of a round.

**Bucket to be split**
**Next**

**Buckets that existed at the beginning of this round: this is the range of**
$$h_{Level}$$

**Buckets split in this round:**
**If** $h_{Level}$ **( search key value )**
**is in this range, must use**
$h_{Level+1}$ **( search key value )**
**to decide if entry is in**
**`split image' bucket.**

**`split image' buckets:**
**created (through splitting**
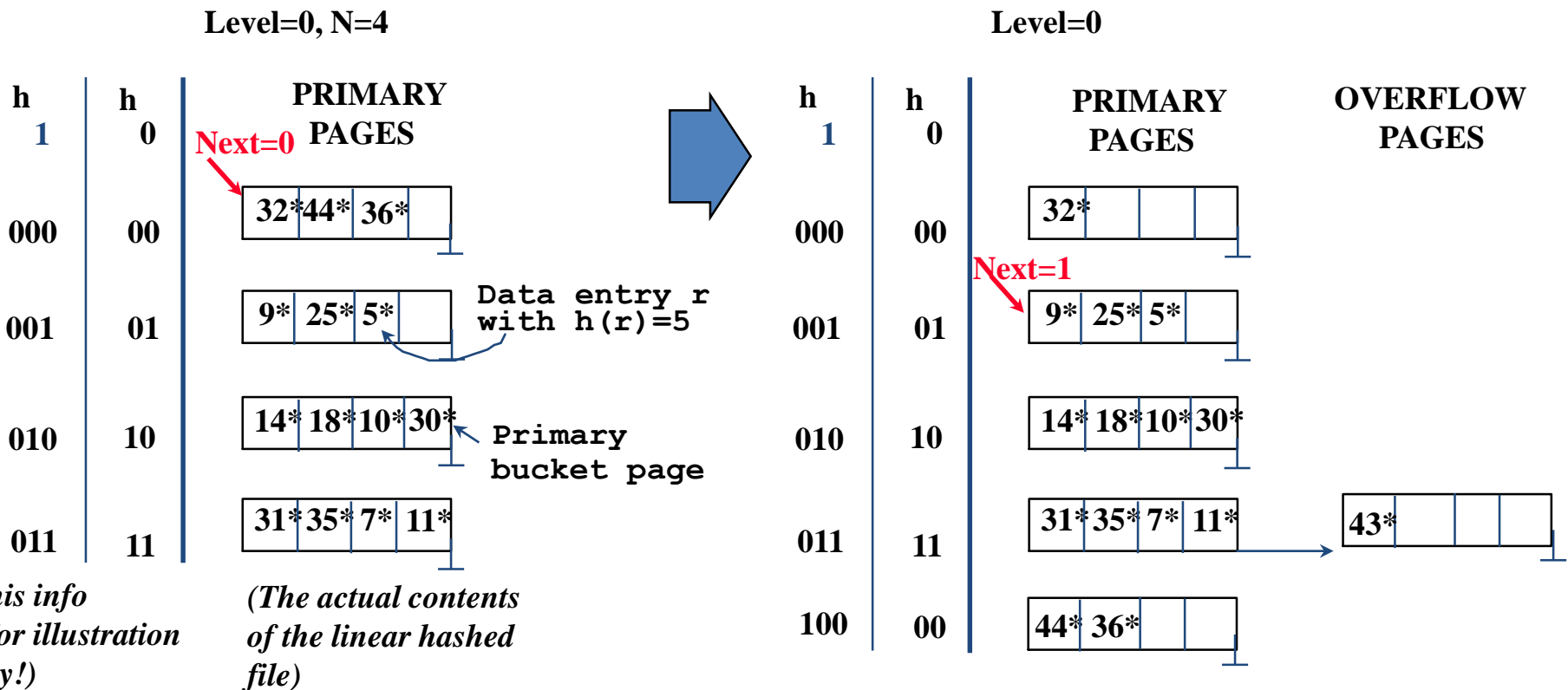**of other buckets) in this round**
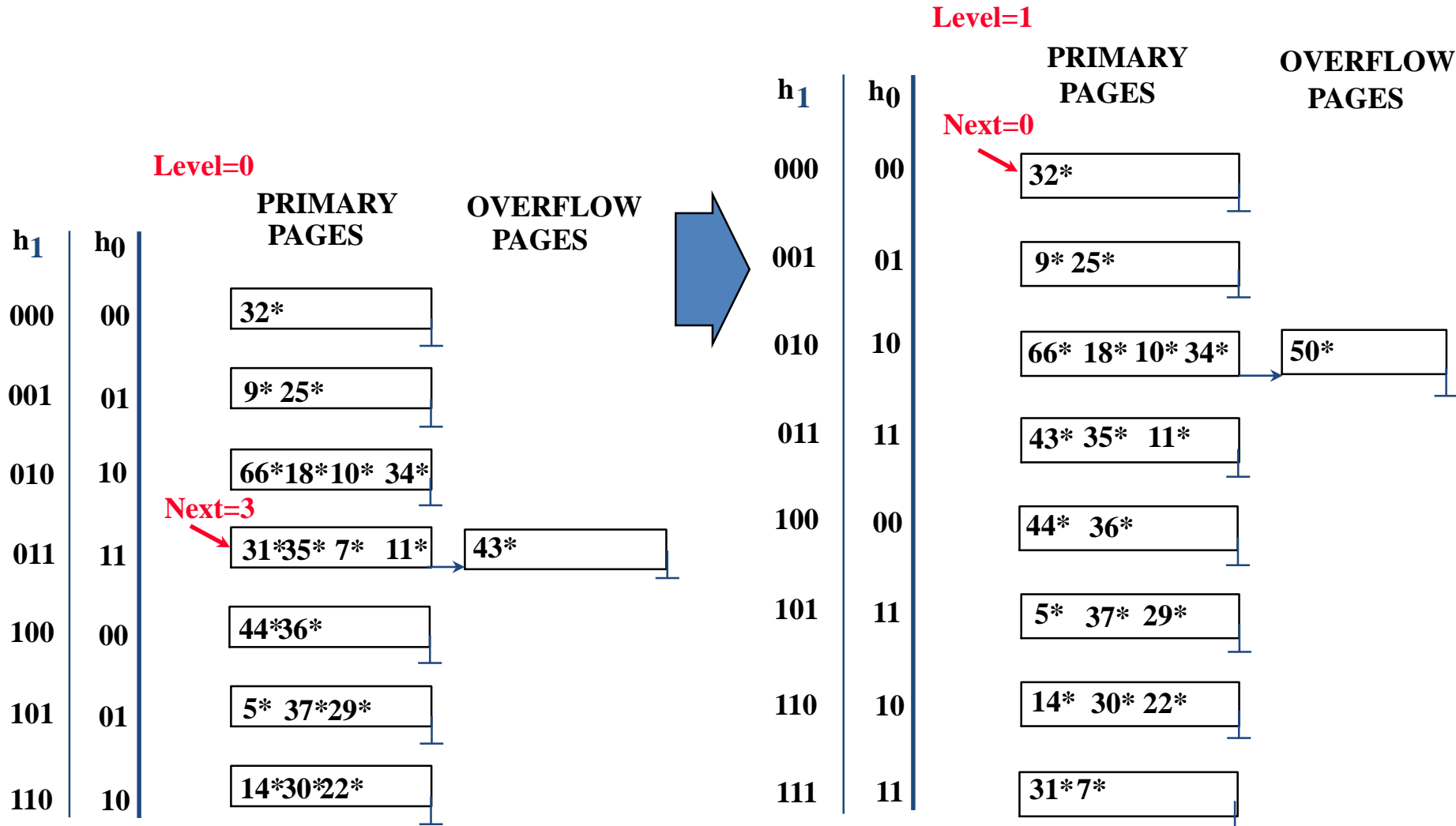
# Linear Hashing (Contd.)

- **<u>Insert</u>**:  Find bucket by applying $h_{Level}$ / $h_{Level+1}$:
  - If bucket to insert into is full:
    - Add overflow page and insert data entry.
    - (*Maybe*) Split *Next* bucket and increment *Next*.
- Can choose any criterion to `trigger' split.
- Since buckets are split round-robin, long overflow chains don't develop!
- Doubling of directory in Extendible Hashing is similar; switching of hash functions is *implicit* in how the # of bits examined is increased.

# Example of Linear Hashing

- On split, $h_{Level+1}$ is used to re-distribute entries.



**Level=0, N=4**

| h 1 | h 0 | | PRIMARY PAGES |
|-----|-----|---|---|
| | | Next=0 | |
| 000 | 00 | | 32* 44* 36* |
| 001 | 01 | | 9* 25* 5* |
| 010 | 10 | | 14* 18* 10* 30* |
| 011 | 11 | | 31* 35* 7* 11* |

Data entry r with h(r)=5

Primary bucket page

*(This info is for illustration only!)*

*(The actual contents of the linear hashed file)*

**Level=0**

| h 1 | h 0 | | PRIMARY PAGES | OVERFLOW PAGES |
|-----|-----|---|---|---|
| 000 | 00 | | 32* | |
| | | Next=1 | | |
| 001 | 01 | | 9* 25* 5* | |
| 010 | 10 | | 14* 18* 10* 30* | |
| 011 | 11 | | 31* 35* 7* 11* | 43* |
| 100 | 00 | | 44* 36* | |

# Example:  End of a Round



Level=0

PRIMARY PAGES

OVERFLOW PAGES

| $h_1$ | $h_0$ | |
|---|---|---|
| 000 | 00 | 32* |
| 001 | 01 | 9* 25* |
| 010 | 10 | 66*18*10* 34* |
| 011 | 11 | 31*35* 7*   11* |
| 100 | 00 | 44*36* |
| 101 | 01 | 5* 37*29* |
| 110 | 10 | 14*30*22* |

Next=3

43*

Level=1

PRIMARY PAGES

OVERFLOW PAGES

| $h_1$ | $h_0$ | |
|---|---|---|
| 000 | 00 | 32* |
| 001 | 01 | 9* 25* |
| 010 | 10 | 66* 18* 10* 34* |
| 011 | 11 | 43* 35*   11* |
| 100 | 00 | 44*   36* |
| 101 | 11 | 5*   37* 29* |
| 110 | 10 | 14*   30* 22* |
| 111 | 11 | 31* 7* |

Next=0

50*

# Summary

- Hash-based indexes: best for equality searches, cannot support range searches.

- Static Hashing can lead to long overflow chains.

- Extendible Hashing avoids overflow pages by splitting a full bucket when a new data entry is to be added to it. *(Duplicates may require overflow pages.)*
  - Directory to keep track of buckets, doubles periodically.
  - Can get large with skewed data; additional I/O if this does not fit in main memory.

# Summary (Contd.)

- Linear Hashing avoids directory by splitting buckets round-robin, and using overflow pages.
  - Overflow pages not likely to be long.
  - Duplicates handled easily.
  - Space utilization could be lower than Extendible Hashing, since splits not concentrated on `dense' data areas.
    - Can tune criterion for triggering splits to trade-off slightly longer chains for better space utilization.
- For hash-based indexes, a *skewed* data distribution is one in which the *hash values* of data entries are not uniformly distributed!

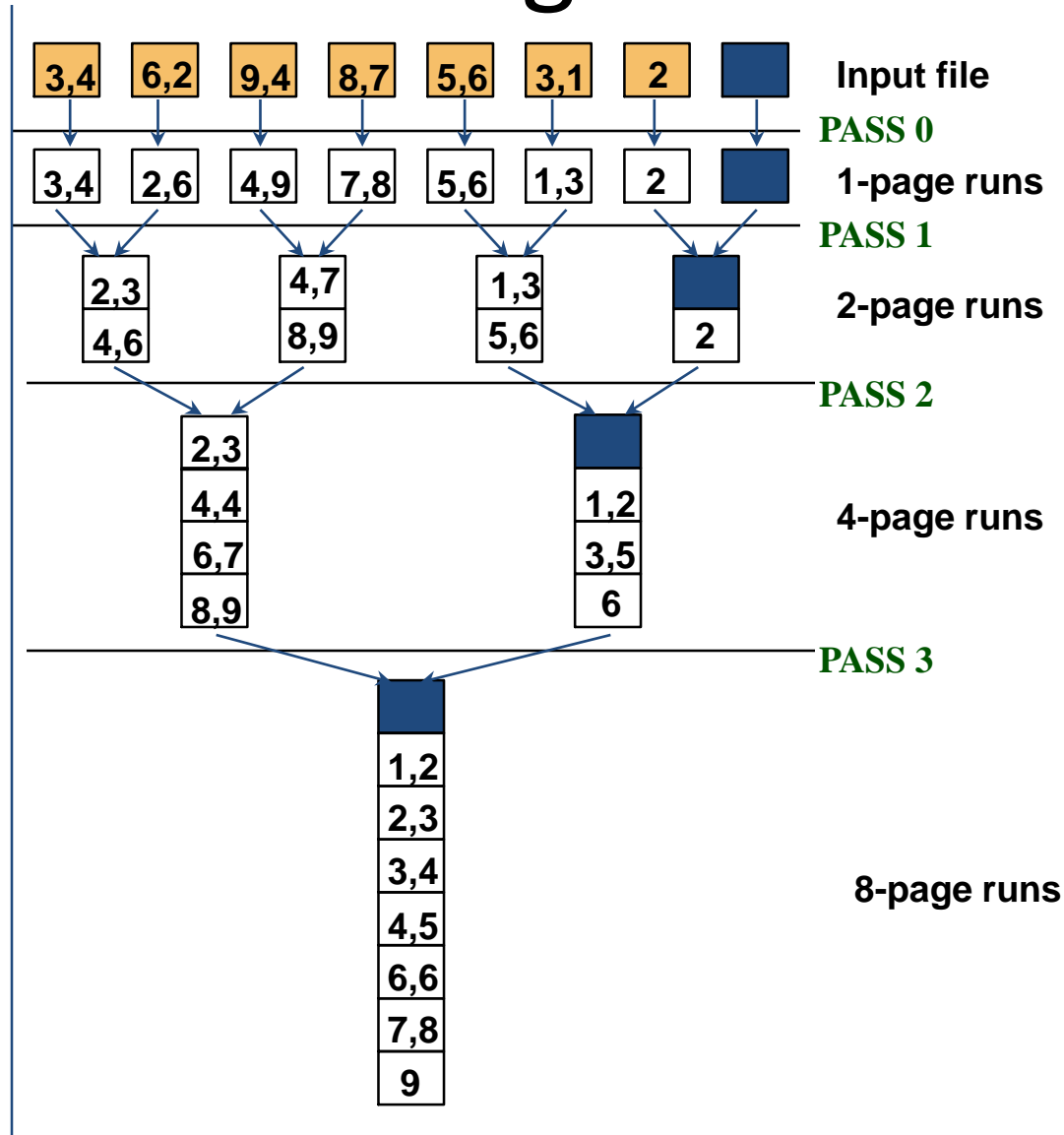# External Sorting

Chapter 13

# Why Sort?

- A classic problem in computer science!
- Data requested in sorted order
  - e.g., find students in increasing *gpa* order
- Sorting is first step in *bulk loading* B+ tree index.
- Sorting useful for eliminating *duplicate copies* in a collection of records (Why?)
- *Sort-merge* join algorithm involves sorting.
- Problem: sort 1Gb of data with 1Mb of RAM.
  - why not virtual memory?

# 2-Way Sort: Requires 3 Buffers

- Pass 1: Read a page, sort it, write it.
  - only one buffer page is used
- Pass 2, 3, …, etc.:
  - three buffer pages used.

# Two-Way External Merge Sort

- Each pass we read + write each page in file.

- N pages in the file => the number of passes
$$= \lceil \log_2 N \rceil + 1$$

- So toal cost is:

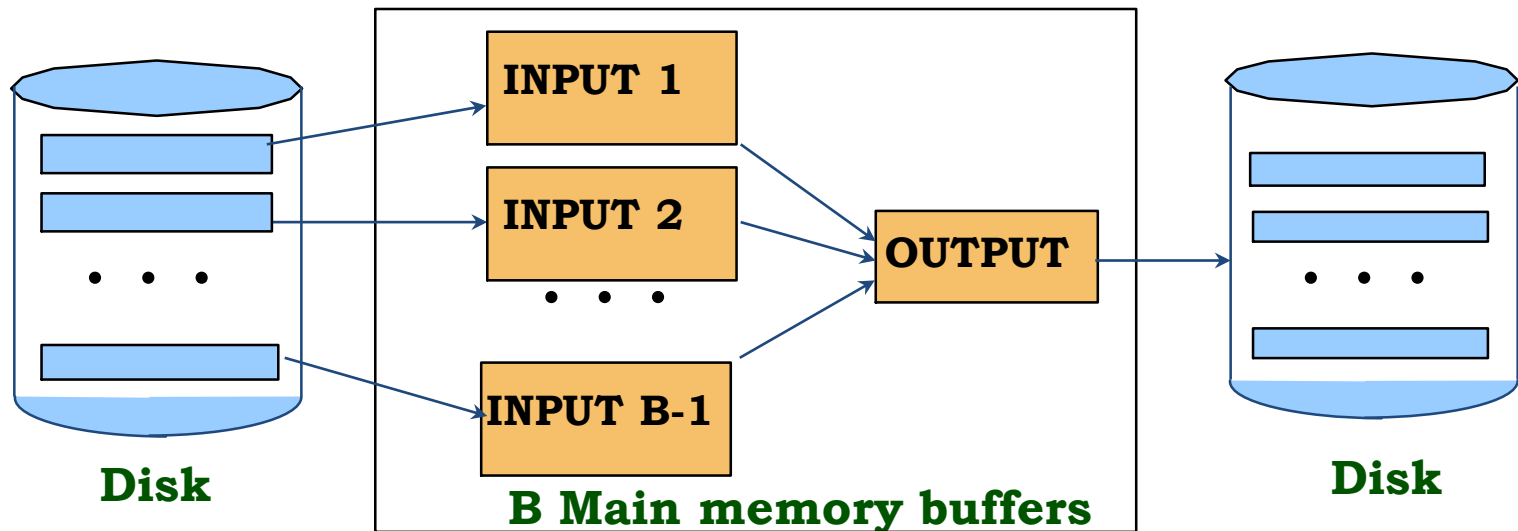$$2N \left( \lceil \log_2 N \rceil + 1 \right)$$

- *Idea:* **Divide and conquer:** sort subfiles and merge

# General External Merge Sort

☐ *More than 3 buffer pages. How can we utilize them?*

- To sort a file with *N* pages using *B* buffer pages:
  - Pass 0: use *B* buffer pages. Produce $\lceil N / B \rceil$ sorted runs of *B* pages each.
  - Pass 2, …, etc.: merge *B-1* runs.



Disk     B Main memory buffers     Disk

# Cost of External Merge Sort

- Number of passes:
- Cost = 2N * (# of passes)  $1 + \lceil \log_{B-1} \lceil N/B \rceil \rceil$
- E.g., with 5 buffer pages, to sort 108 page file:
  - Pass 0: $\lceil 108/5 \rceil$ = 22 sorted runs of 5 pages each (last run is only 3 pages)
  - Pass 1: $\lceil 22/4 \rceil$ = 6 sorted runs of 20 pages each (last run is only 8 pages)
  - Pass 2: 2 sorted runs, 80 pages and 28 pages
  - Pass 3: Sorted file of 108 pages
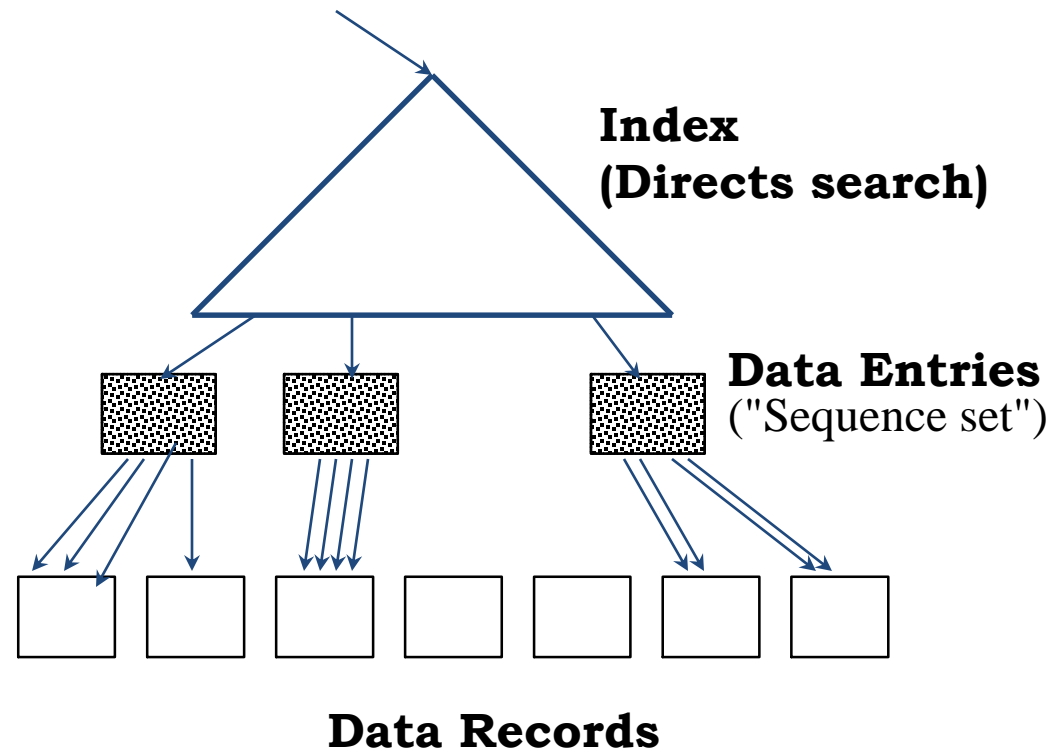
# Number of Passes of External Sort

| N | B=3 | B=5 | B=9 | B=17 | B=129 | B=257 |
|---|---|---|---|---|---|---|
| 100 | 7 | 4 | 3 | 2 | 1 | 1 |
| 1,000 | 10 | 5 | 4 | 3 | 2 | 2 |
| 10,000 | 13 | 7 | 5 | 4 | 2 | 2 |
| 100,000 | 17 | 9 | 6 | 5 | 3 | 3 |
| 1,000,000 | 20 | 10 | 7 | 5 | 3 | 3 |
| 10,000,000 | 23 | 12 | 8 | 6 | 4 | 3 |
| 100,000,000 | 26 | 14 | 9 | 7 | 4 | 4 |
| 1,000,000,000 | 30 | 15 | 10 | 8 | 5 | 4 |

# Using B+ Trees for Sorting

- Scenario: Table to be sorted has B+ tree index on sorting column(s).

- Idea: Can retrieve records in order by traversing leaf pages.

- ***Is this a good idea?***

- Cases to consider:
  - B+ tree is clustered          ***Good idea!***
  - B+ tree is not clustered      ***Could be a very bad idea!***

# Clustered B+ Tree Used for Sorting
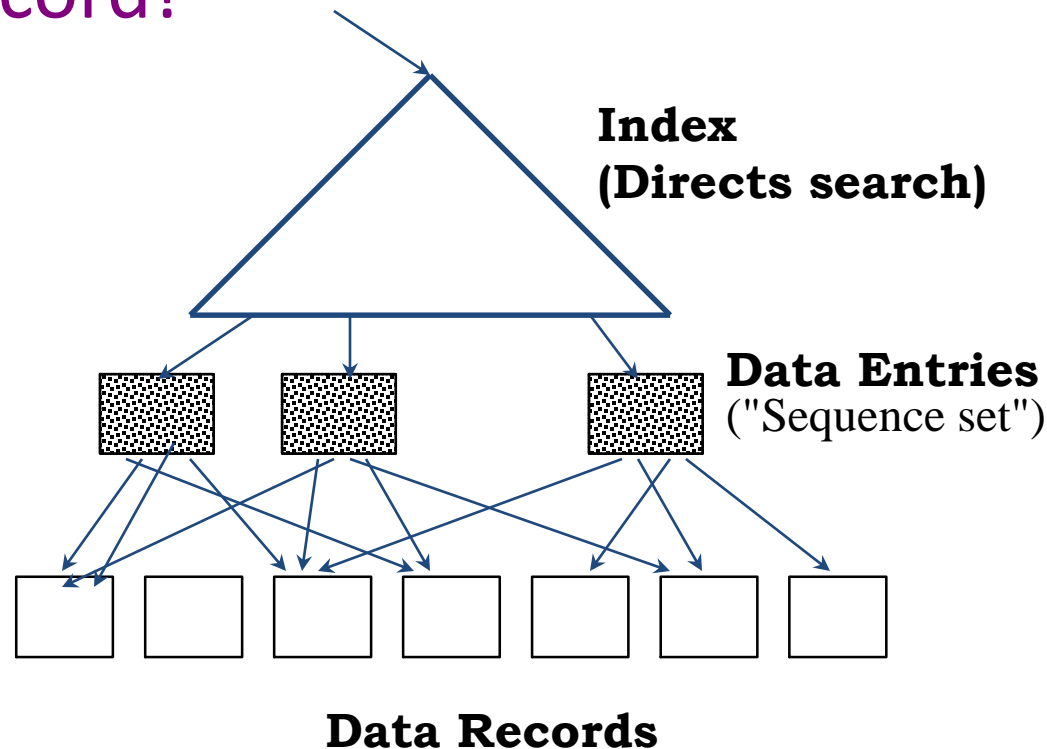
- Cost: root to the left-most leaf, then retrieve all leaf pages

**Index**
**(Directs search)**

**Data Entries**
("Sequence set")

**Data Records**

☐ *Always better than external sorting!*

# Unclustered B+ Tree Used for Sorting

- For data entries; each data entry contains *rid* of a data record.  In general, one I/O per data record!



Index
(Directs search)

Data Entries
("Sequence set")

Data Records

# External Sorting vs. Unclustered Index

| N | Sorting | p=1 | p=10 | p=100 |
|---|---|---|---|---|
| 100 | 200 | 100 | 1,000 | 10,000 |
| 1,000 | 2,000 | 1,000 | 10,000 | 100,000 |
| 10,000 | 40,000 | 10,000 | 100,000 | 1,000,000 |
| 100,000 | 600,000 | 100,000 | 1,000,000 | 10,000,000 |
| 1,000,000 | 8,000,000 | 1,000,000 | 10,000,000 | 100,000,000 |
| 10,000,000 | 80,000,000 | 10,000,000 | 100,000,000 | 1,000,000,000 |

- $p$: # of records per page
- B=1,000 and block size=32 for sorting
- p=100 is the more realistic value.

# Summary

- External sorting is important; DBMS may dedicate part of buffer pool for sorting!
- External merge sort minimizes disk I/O cost:
  - Pass 0: Produces sorted *runs* of size *B* (# buffer pages). Later passes: *merge* runs.
  - # of runs merged at a time depends on *B,* and *block size.*
  - Larger block size means less I/O cost per page.
  - Larger block size means smaller # runs merged.
  - In practice, # of runs rarely more than 2 or 3.

# Summary, cont.

- Choice of internal sort algorithm may matter:
  - Quicksort: Quick!
  - Heap/tournament sort: slower (2x), longer runs
- The best sorts are wildly fast:
  - Despite 40+ years of research, we're still improving!
- Clustered B+ tree is good for sorting; unclustered tree is usually very bad.