# Database Applications (15-415)

# DBMS Internals- Part V
# Lecture 15, March 15, 2015

Mohammad Hammoud

# Today…

- **Last Session:**
  - DBMS Internals- Part IV
    - Tree-based (i.e., B+ Tree) and Hash-based (i.e., Extendible Hashing) indexes
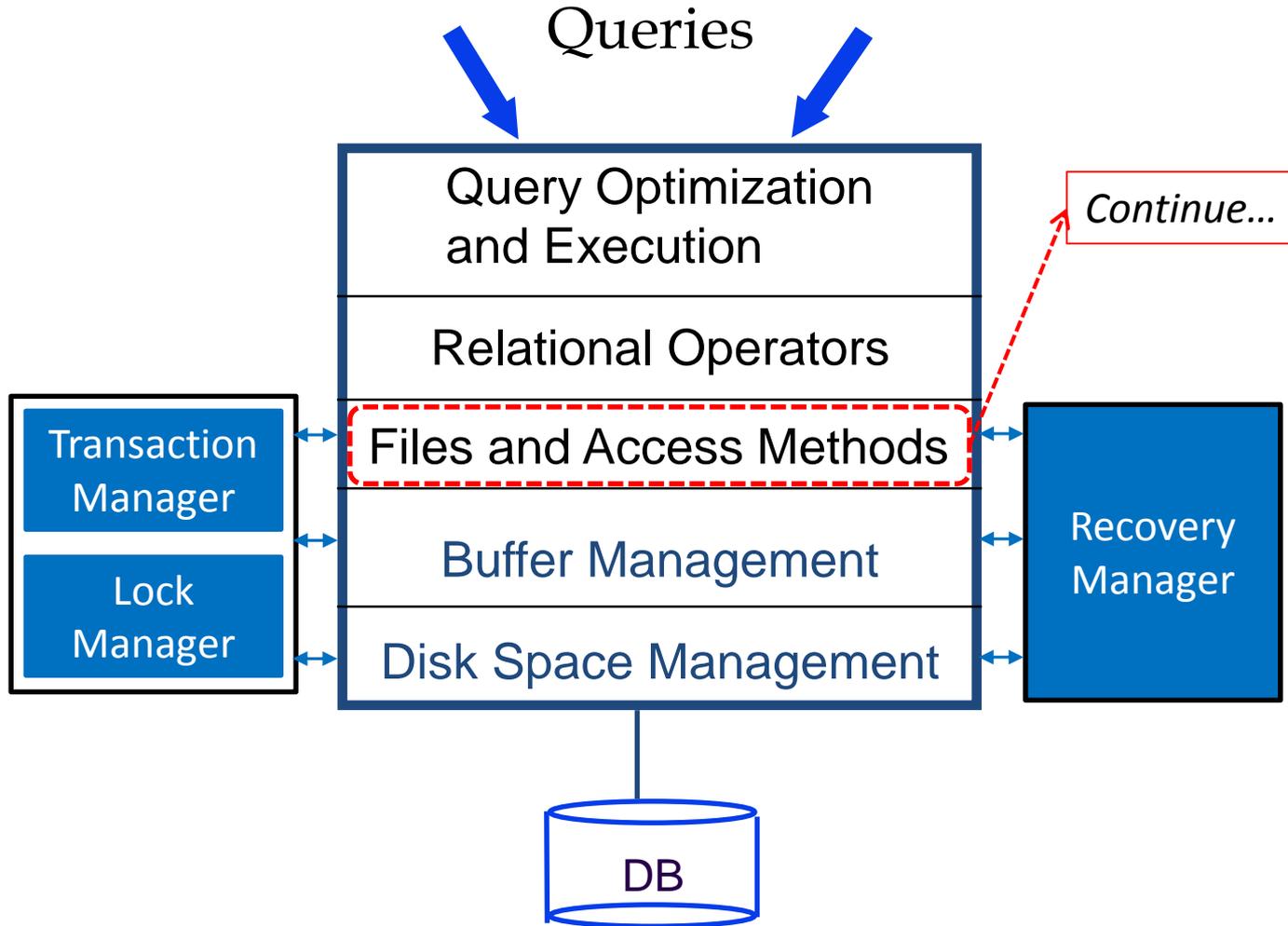
- **Today's Session:**
  - DBMS Internals- Part V
    - Hash-based indexes (Cont'd) and External Sorting
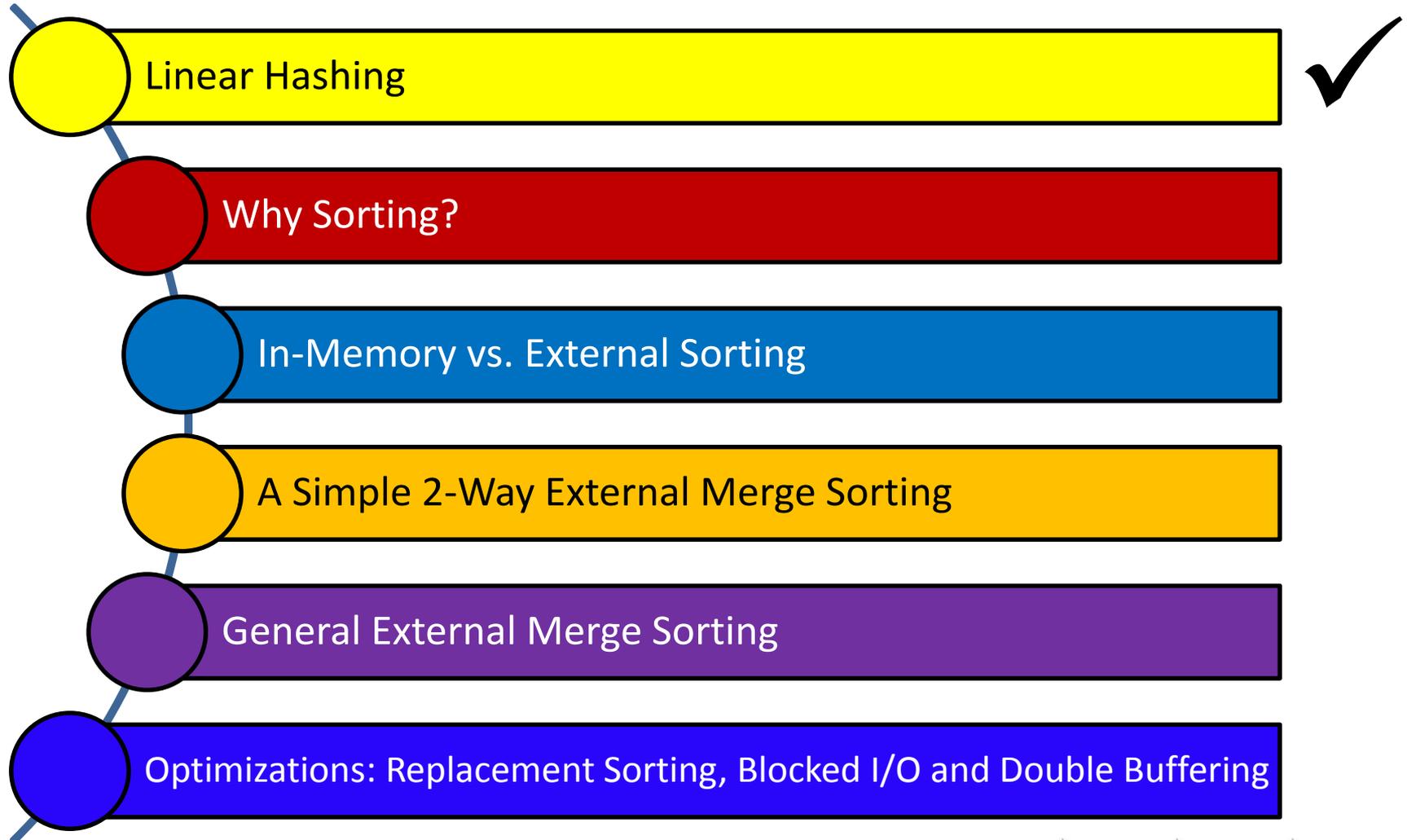
- **Announcements:**
  - Project 2 is due today by midnight. Student demos will be conducted on Tuesday/Thursday
  - PS3 is now posted and it is due on March 26 by midnight
  - Project 3 will be posted by Thursday

# DBMS Layers

# Outline

**Linear Hashing** ✓

**Why Sorting?**

**In-Memory vs. External Sorting**

**A Simple 2-Way External Merge Sorting**

**General External Merge Sorting**

**Optimizations: Replacement Sorting, Blocked I/O and Double Buffering**

# Linear Hashing

- Another way of adapting gracefully to insertions and deletions (i.e., pursuing dynamic hashing) is to use Linear Hashing (LH)

- In contrast to Extendible Hashing, LH
  - Does not require a directory
  - Deals naturally with collisions
  - Offers a lot of flexibility w.r.t the timing of bucket split (allowing trading off greater overflow chains for higher average space utilization)
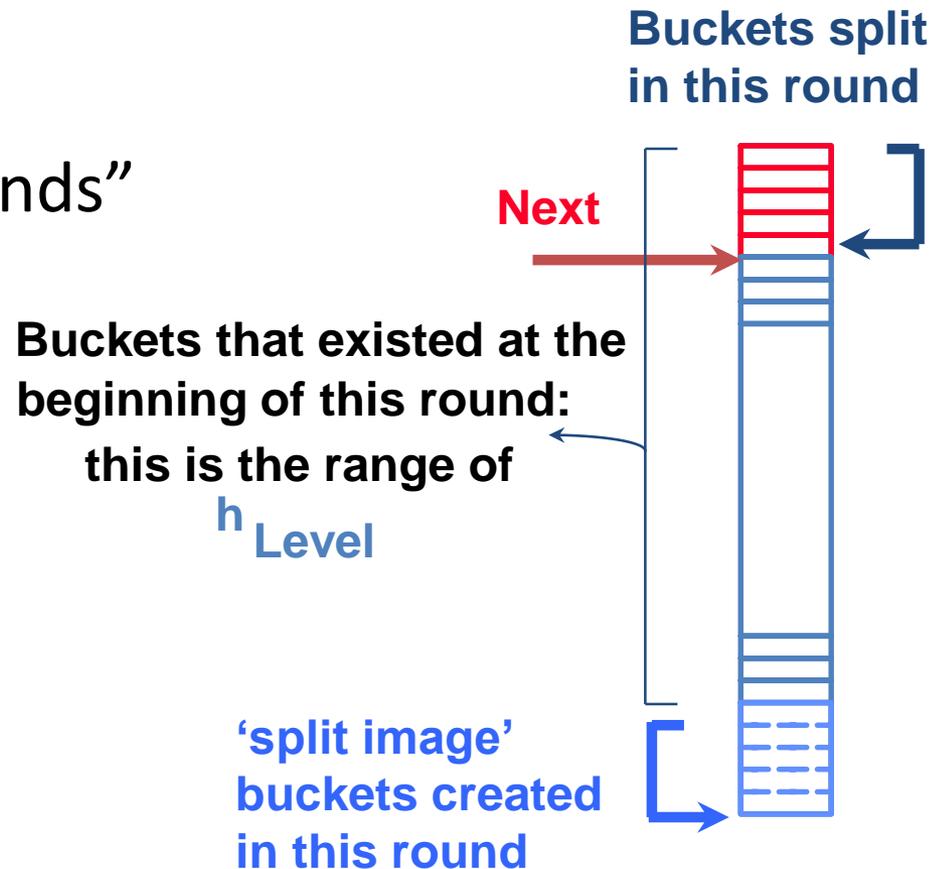
# How Linear Hashing Works?

- LH uses a family of hash functions $h_0$, $h_1$, $h_2$, ...

  - $h_i(key) = h(key) \bmod (2^i N)$;  N = initial # buckets

  - $h$ is some hash function (range is *not* 0 to N-1)

  - If N = $2^{d0}$, for some *d0*, $h_i$ consists of applying $h$ and looking at the last *di* bits, where *di = d0 + i*

  - $h_{i+1}$ doubles the range of $h_i$ (*similar to directory doubling*)

# How Linear Hashing Works? (Cont'd)

- LH uses overflow pages, and chooses buckets to split in a *round-robin* fashion

- Splitting proceeds in "rounds"
  - A round ends when all $N_R$ (for round $R$) initial buckets are split
  - Buckets 0 to *Next-1* have been split; *Next* to $N_R$ yet to be split
  - Current round number is referred to as *Level*

**Buckets split in this round**

**Next**

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

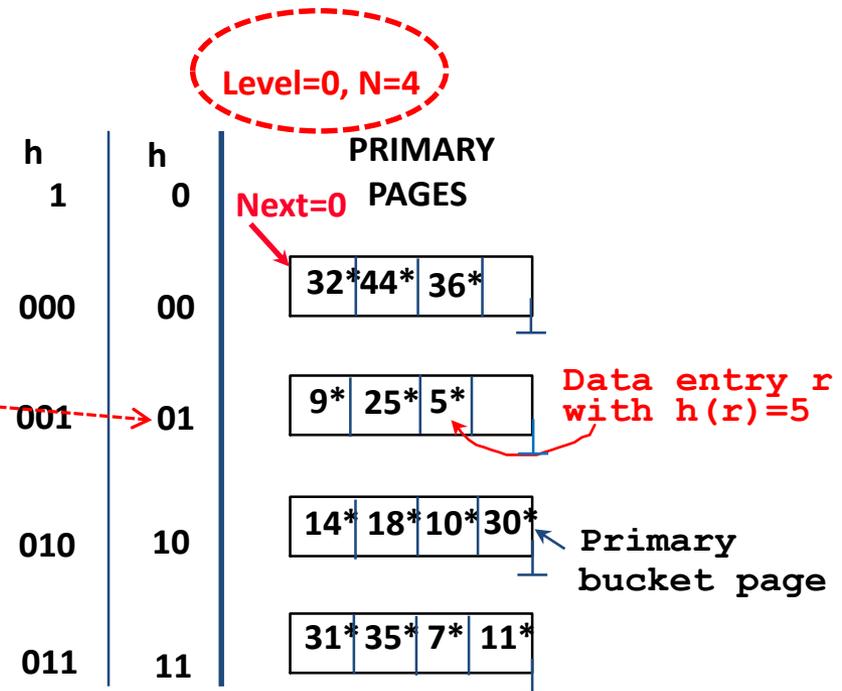**'split image' buckets created in this round**

# Linear Hashing: Searching For Entries

- To find a bucket for data entry $r$, find $h_{Level}(r)$:
  - If $h_{Level}(r)$ in range `Next to $N_R$' , $r$ belongs there
  - Else, $r$ could belong to bucket $h_{Level}(r)$ or bucket $h_{Level}(r) + N_R$; must apply $h_{Level+1}(r)$ to find out

- Example: search for 5*

**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* |

Level = 0 ➜ h0
5* = 101 ➜ 01

Data entry r with h(r)=5

Primary bucket page

# Linear Hashing: Inserting Entries

- Find bucket as in search
  - If the bucket to insert the data entry into is full:
    - Add an overflow page and insert data entry
    - (*Maybe*) Split *Next* bucket and increment *Next*

- Some points to Keep in mind:
  - Unlike Extendible Hashing, when an insert triggers a split, the bucket into which the data entry is inserted is not necessarily the bucket that is split

  - As in Static Hashing, an overflow page is added to store the newly inserted data entry

  - However, since the bucket to split is chosen in a round-robin fashion, eventually *all* buckets will be split
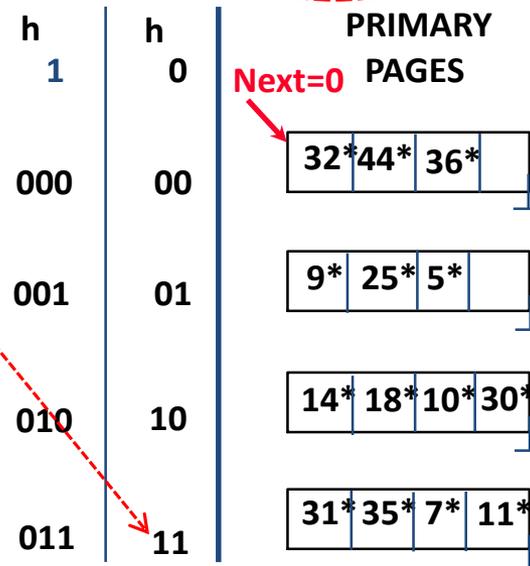
# Linear Hashing: Inserting Entries

- Example: insert 43*

**Level = 0 ➔ h0**

**43* = 101011 ➔ 11**

**Level=0, N=4**

| h1 | h0 | PRIMARY PAGES |
|---|---|---|
| 000 | 00 | 32* 44* 36* |
| 001 | 01 | 9* 25* 5* |
| 010 | 10 | 14* 18* 10* 30* |
| 011 | 11 | 31* 35* 7* 11* |

**Next=0**
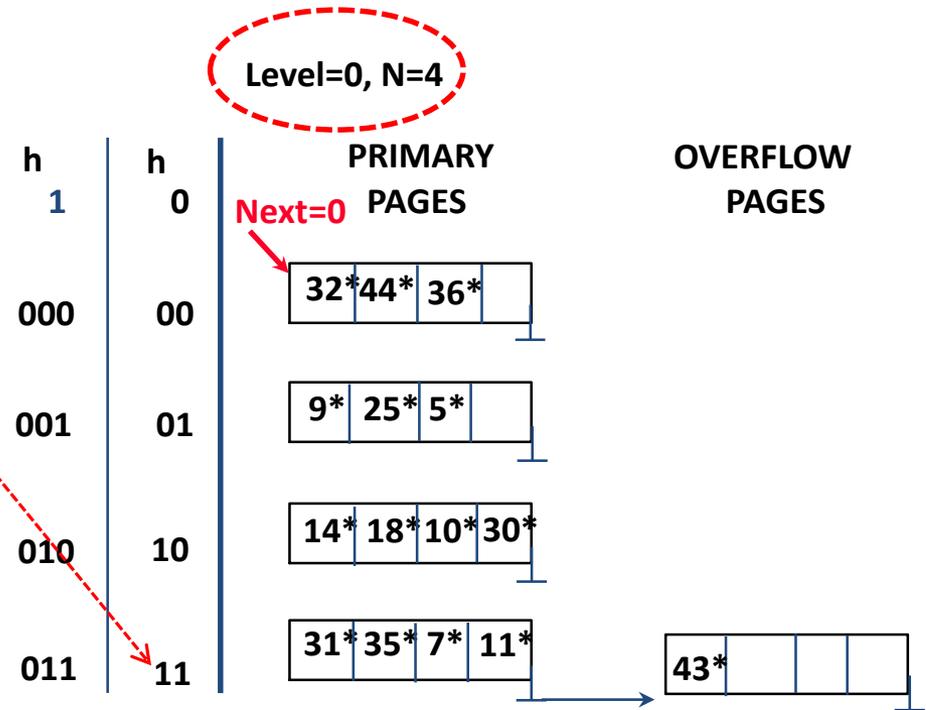
Add an overflow page and insert data entry

# Linear Hashing: Inserting Entries

- Example: insert 43*

**Level = 0 ➜ h0**

**43* = 101011 ➜ 11**

**Level=0, N=4**

| h 1 | h 0 | PRIMARY PAGES | OVERFLOW 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* | 43* |

Split *Next* bucket and increment *Next*

# Linear Hashing: Inserting Entries

- Example: insert 43*

Level = 0 ➜ h0
43* = 101011 ➜ 11

Almost there…

Level=0, N=4

| h 1 | h 0 | PRIMARY PAGES | OVERFLOW PAGES |
|-----|-----|---------------|----------------|

Next=0

000 | 00 | 32* | | |

001 | 01 | 9* | 25* | 5* |

010 | 10 | 14* | 18* | 10* | 30* |

011 | 11 | 31* | 35* | 7* | 11* |  →  43*

100 | 00 | 44* | 36* | |

# Linear Hashing: Inserting Entries

- Example: insert 43*

**Level = 0 → h0**

**43* = 101011 → 11**

**Level=0, N=4**

**FINAL STATE!**

| h$_1$ | h$_0$ | PRIMARY PAGES | | | | OVERFLOW PAGES |
|---|---|---|---|---|---|---|
| 000 | 00 | 32* | | | | |
| 001 | 01 | 9* | 25* | 5* | | |
| 010 | 10 | 14* | 18* | 10* | 30* | |
| 011 | 11 | 31* | 35* | 7* | 11* | 43* |
| 100 | 00 | 44* | 36* | | | |

**Next=1**

# Linear Hashing: Inserting Entries

- Another Example: insert 50*

**Level=0, N= 4**

Level = 0 ➜ h0
50* = 110010 ➜ 10

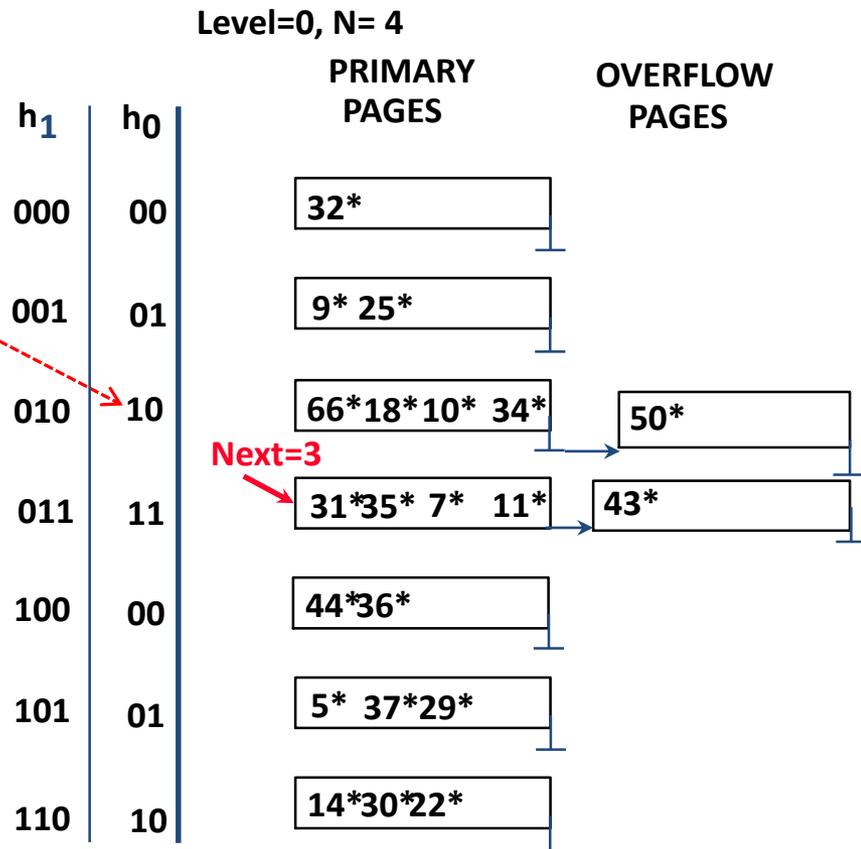|  $h_1$ | $h_0$ | PRIMARY PAGES | OVERFLOW PAGES |
|------|------|---------------|----------------|
| 000 | 00 | 32* | |
| 001 | 01 | 9* 25* | |
| 010 | 10 | 66*18*10* 34* | |
| 011 | 11 | 31*35* 7*  11* | 43* |
| 100 | 00 | 44*36* | |
| 101 | 01 | 5*  37*29* | |
| 110 | 10 | 14*30*22* | |

Next=3

Add an overflow page and insert data entry

# Linear Hashing: Inserting Entries

■ Another Example: insert 50*

Level=0, N= 4

Level = 0 ➜ h0

50* = 110010 ➜ 10

| h1 | h0 | PRIMARY PAGES | OVERFLOW PAGES |
|---|---|---|---|
| 000 | 00 | 32* | |
| 001 | 01 | 9* 25* | |
| 010 | 10 | 66*18*10* 34* | 50* |
| 011 | 11 | 31*35* 7* 11* | 43* |
| 100 | 00 | 44*36* | |
| 101 | 01 | 5* 37*29* | |
| 110 | 10 | 14*30*22* | |

Next=3

Split *Next* bucket and increment *Next*

# Linear Hashing: Inserting Entries

- Another Example: insert 50*

**Level=0**

| $h_1$ | $h_0$ | PRIMARY PAGES | OVERFLOW PAGES |
|-------|-------|---------------|----------------|
| 000 | 00 | 32* | |
| 001 | 01 | 9* 25* | |
| 010 | 10 | 66* 18* 10* 34* | 50* |
| 011 | 11 | 43* 35*  11* | |
| 100 | 00 | 44*  36* | |
| 101 | 11 | 5*  37* 29* | |
| 110 | 10 | 14*  30* 22* | |
| 111 | 11 | 31* 7* | |

Level = 0 ➜ h0
50* = 110010 ➜ 10

Next=3

Almost there…

# Linear Hashing: Inserting Entries

- Another Example: insert 50*

**Level=0**

| $h_1$ | $h_0$ | PRIMARY PAGES | OVERFLOW PAGES |
|-------|-------|---------------|----------------|

Level = 0 ➜ h0
50* = 110010 ➜ 10

Next=0

| | | | |
|---|---|---|---|
| 000 | 00 | 32* | |
| 001 | 01 | 9* 25* | |
| 010 | 10 | 66* 18* 10* 34* | 50* |
| 011 | 11 | 43* 35*  11* | |
| 100 | 00 | 44*  36* | |
| 101 | 11 | 5*  37* 29* | |
| 110 | 10 | 14*  30* 22* | |
| 111 | 11 | 31* 7* | |

Almost there...

# Linear Hashing: Inserting Entries

- Another Example: insert 50*

Level=1

| $h_1$ | $h_0$ | PRIMARY PAGES | OVERFLOW PAGES |
|---|---|---|---|

Next=0

Level = 0 ➜ h0
50* = 110010 ➜ 10

FINAL STATE!

| 000 | 00 | 32* | |
| 001 | 01 | 9* 25* | |
| 010 | 10 | 66* 18* 10* 34* | 50* |
| 011 | 11 | 43* 35*  11* | |
| 100 | 00 | 44*  36* | |
| 101 | 11 | 5*  37* 29* | |
| 110 | 10 | 14*  30* 22* | |
| 111 | 11 | 31* 7* | |

# Linear Hashing: Deleting Entries

- Deletion is essentially the inverse of insertion

- If the last bucket in the file is empty, it can be removed and *Next* can be decremented

- If *Next* is zero and the last bucket becomes empty
  - *Next* is made to point to bucket *M*/2 -1 (where *M* is the current number of buckets)
  - *Level* is decremented
  - The empty bucket is removed

- The insertion examples can be worked out backwards as examples of deletions!

# DBMS Layers

# Outline

Linear Hashing

Why Sorting? ✔

In-Memory vs. External Sorting

A Simple 2-Way External Merge Sorting

General External Merge Sorting

Optimizations: Replacement Sorting, Blocked I/O and Double Buffering

# When Does A DBMS Sort Data?

- Users may want answers in some order
    - **SELECT FROM** student ***ORDER BY*** name
    - **SELECT** S.rating, **MIN** (S.age) **FROM** Sailors S ***GROUP BY*** S.rating

- *Bulk loading* a B+ tree index involves sorting

- Sorting is useful in eliminating duplicates records

- The *Sort-Merge* Join algorithm involves sorting (*next session!*)

# Outline

- **Linear Hashing**
- **Why Sorting?**
- **In-Memory vs. External Sorting** ✔
- **A Simple 2-Way External Merge Sorting**
- **General External Merge Sorting**
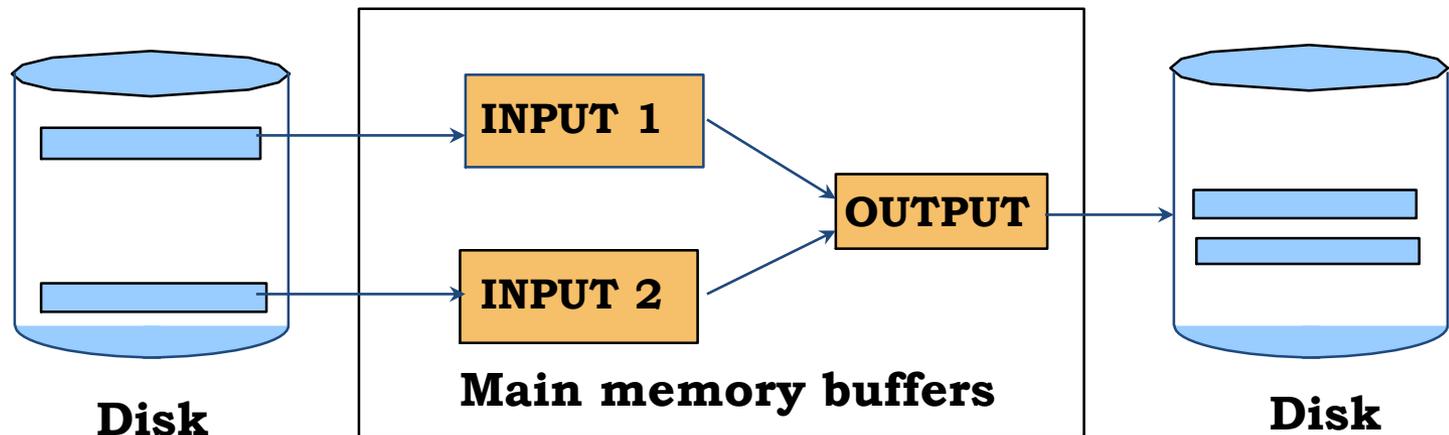- **Optimizations: Replacement Sorting, Blocked I/O and Double Buffering**

Carnegie Mellon University Qatar

# In-Memory vs. External Sorting

- Assume we want to sort 60GB of data on a machine with only 8GB of RAM
  - In-Memory Sort (e.g., Quicksort) ?
    - Yes, but data do not fit in memory
    - What about relying on virtual memory?

  - In this case, external sorting is needed
    - In-memory sorting is *orthogonal* to external sorting!

# Outline

- Linear Hashing
- Why Sorting?
- In-Memory vs. External Sorting
- A Simple 2-Way External Merge Sorting ✔
- General External Merge Sorting
- Optimizations: Replacement Sorting, Blocked I/O and Double Buffering

**Carnegie Mellon University Qatar**

# A Simple Two-Way Merge Sort

- **IDEA**: Sort sub-files that can fit in memory and merge

- Let us refer to each sorted sub-file as a *run*

- Algorithm:
  - Pass 1: Read a page into memory, sort it, write it
    - 1-page runs are produced
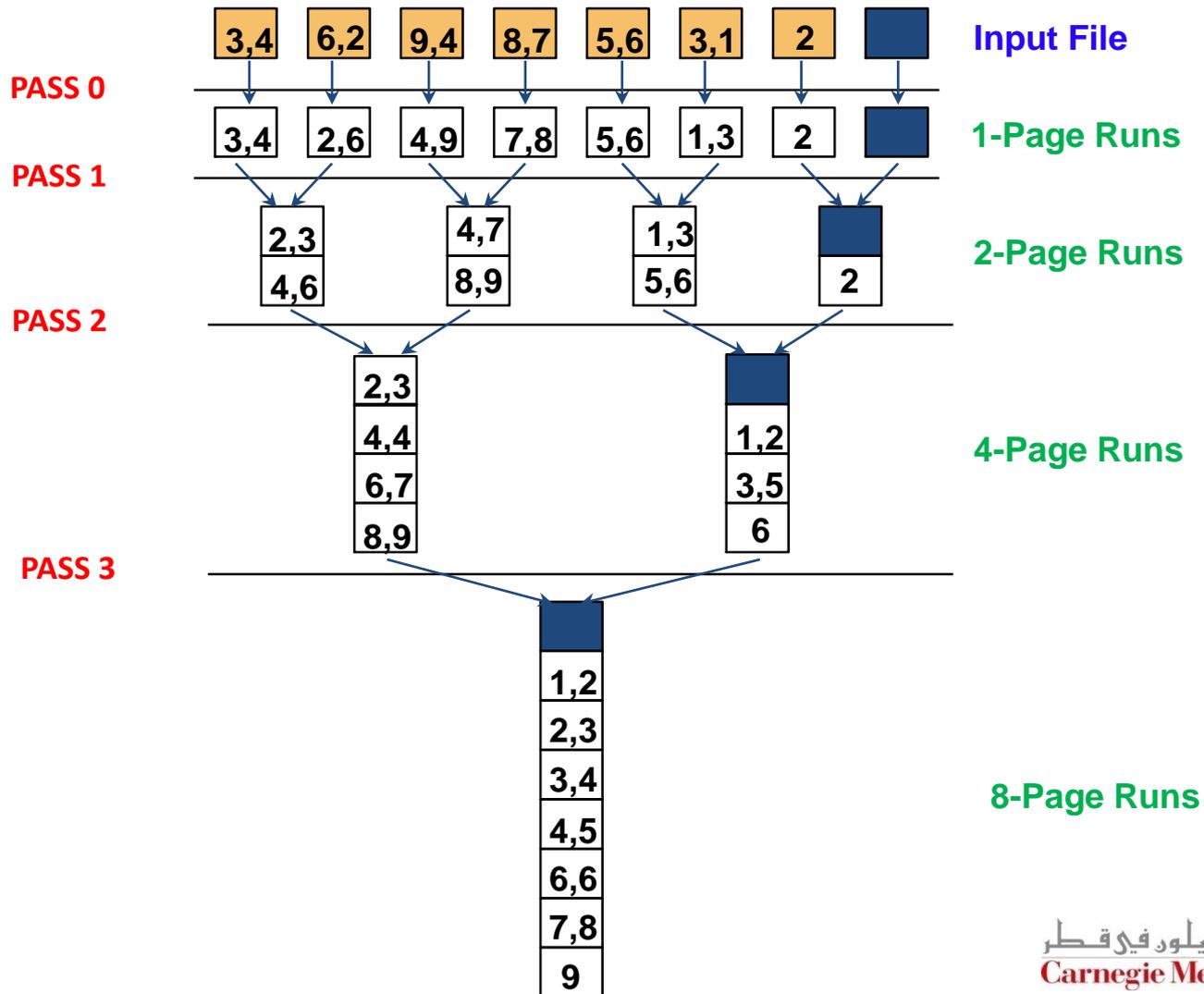  - Passes 2, 3, etc.,: Merge *pairs* (hence, 2-way) of runs to produce longer runs until only one run is left

# A Simple Two-Way Merge Sort

- ## Algorithm:

  - ### Pass 1: Read a page into memory, sort it, write it
    - How many buffer pages are needed? ONE

  - ### Passes 2, 3, etc.,: Merge *pairs* (hence, 2-way) of runs to produce longer runs until only one run is left
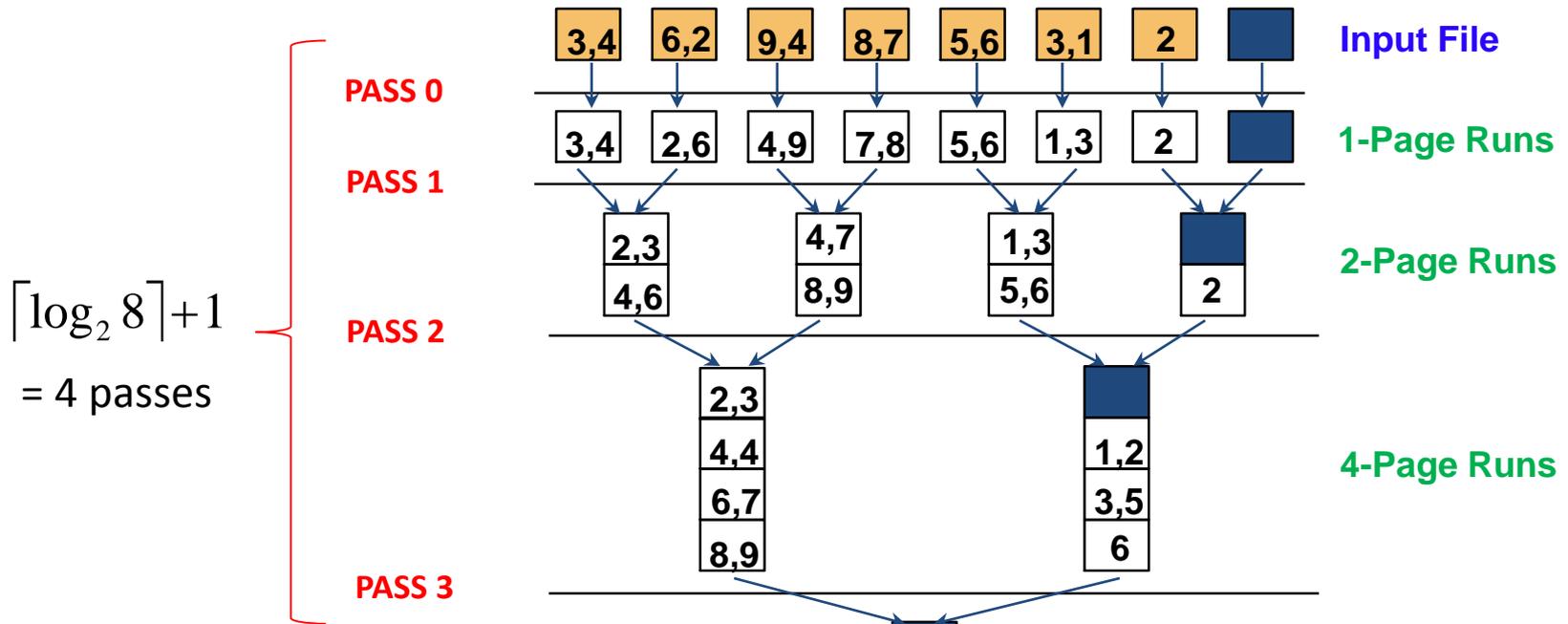    - How many buffer pages are needed? THREE

**Disk** — INPUT 1, INPUT 2, OUTPUT — **Main memory buffers** — **Disk**

# 2-Way Merge Sort: An Example

# 2-Way Merge Sort: I/O Cost Analysis

- If the number of pages in the input file is $2^k$
  - How many runs are produced in pass 0 and of what size?
    - $2^k$ 1-page runs
  - How many runs are produced in pass 1 and of what size?
    - $2^{k-1}$ 2-page runs
  - How many runs are produced in pass 2 and of what size?
    - $2^{k-2}$ 4-page runs
  - How many runs are produced in pass k and of what size?
    - $2^{k-k}$ $2^k$-page runs (or 1 run of size $2^k$)
  - For *N* number of pages, how many passes are incurred?
    - $\lceil \log_2 N \rceil + 1$
  - How many pages do we read and write in each pass?
    - 2*N*
  - *What is the overall cost?*

    $$2N \times (\lceil \log_2 N \rceil + 1)$$

# 2-Way Merge Sort: An Example

**Input File**

| 3,4 | 6,2 | 9,4 | 8,7 | 5,6 | 3,1 | 2 | ▆ |

**PASS 0**

**1-Page Runs**

| 3,4 | 2,6 | 4,9 | 7,8 | 5,6 | 1,3 | 2 | ▆ |

**PASS 1**

**2-Page Runs**

| 2,3 | | 4,7 | | 1,3 | | ▆ |
| 4,6 | | 8,9 | | 5,6 | | 2 |

$$\lceil \log_2 8 \rceil + 1$$

= 4 passes

**PASS 2**

**4-Page Runs**

| 2,3 | | ▆ |
| 4,4 | | 1,2 |
| 6,7 | | 3,5 |
| 8,9 | | 6 |

**PASS 3**

**8-Page Runs**

| ▆ |
| 1,2 |
| 2,3 |
| 3,4 |
| 4,5 |
| 6,6 |
| 7,8 |
| 9 |

**Formula Check:**

$$2N \times (\lceil \log_2 N \rceil + 1)$$

= (2 × 8) × (3 + 1) = 64 I/Os
Correct!

# Outline

Linear Hashing

Why Sorting?

In-Memory vs. External Sorting

A Simple 2-Way External Merge Sorting

General External Merge Sorting ✔

Optimizations: Replacement Sorting, Blocked I/O and Double Buffering

# *B*-Way Merge Sort

- How can we sort a file with *N* pages using **_B_** buffer pages?
  - Pass 0: use *B* buffer pages
    - This will produce $\lceil N / B \rceil$ sorted B-page runs
  - Pass 2, …, etc.: merge *B-1* runs



**Disk**        **B Main memory buffers**        **Disk**

INPUT 1
INPUT 2
INPUT B-1
OUTPUT

# B-Way Merge Sort: I/O Cost Analysis

- I/O cost = 2N × Number of passes

- Number of passes = $1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil$

- Assume the previous example (i.e., 8 pages), *but* using 5 buffer pages (instead of 2)
  - I/O cost = 32 (*as opposed to 64*)

- Therefore, increasing the number of buffer pages minimizes the number of passes and accordingly the I/O cost!

# Number of Passes of B-Way 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 |

High Fan-in during merging is crucial!

How else can we minimize I/O cost?

# Outline

Linear Hashing

Why Sorting?

In-Memory vs. External Sorting

A Simple 2-Way External Merge Sorting

General External Merge Sorting

Optimizations: Replacement Sorting, Blocked I/O and Double Buffering ✔
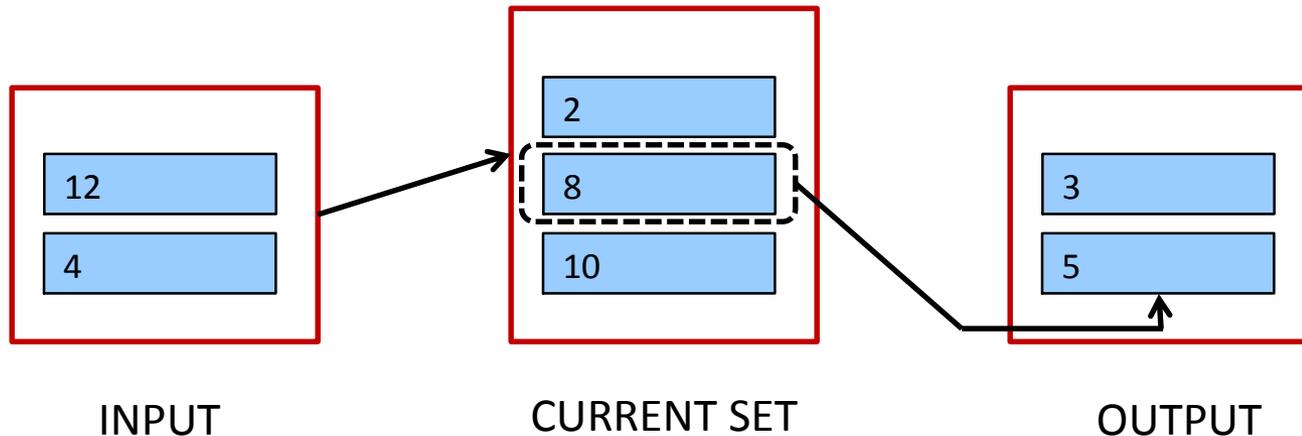
# Replacement Sort

- With a more aggressive implementation of B-way sort, we can write out runs of ~2×B internally sorted pages
  - This is referred to as replacement sort



INPUT          CURRENT SET          OUTPUT

IDEA: Pick the tuple in the *current set* with the smallest value that is greater than the largest value in the *output buffer* and append it to the *output buffer*

# Replacement Sort

- With a more aggressive implementation of B-way sort, we can write out runs of ~2×B internally sorted pages
  - This is referred to as replacement sort


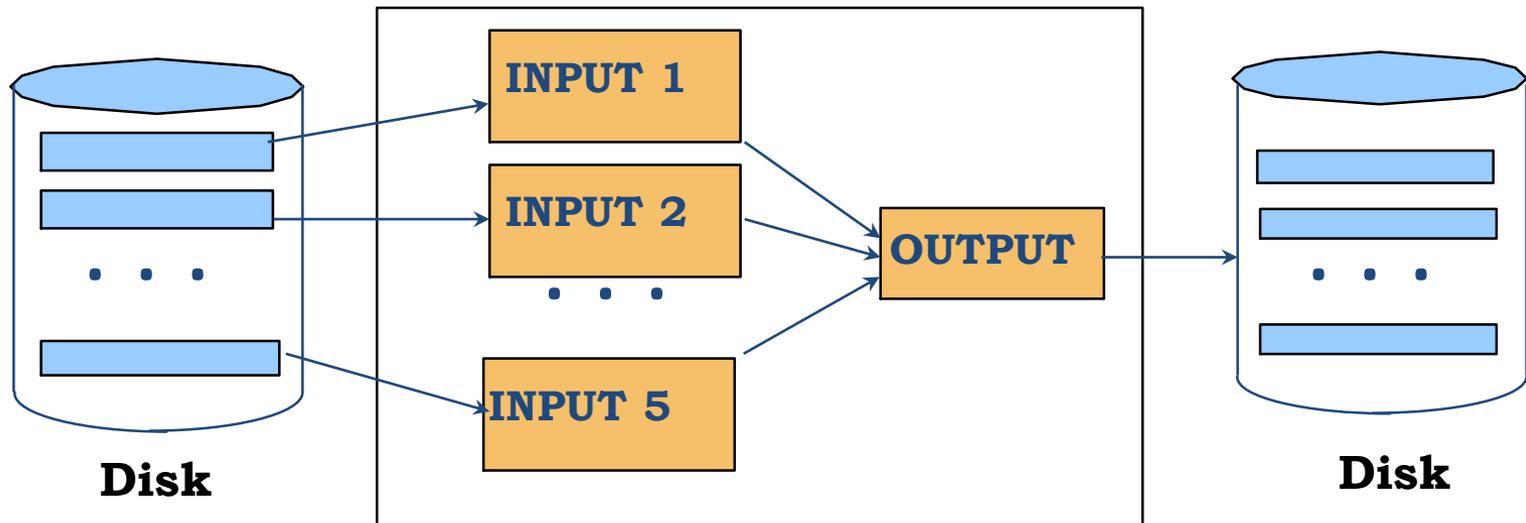
INPUT          CURRENT SET          OUTPUT

When do we terminate the current *run* and start a new one?

# Blocked I/O and Double Buffering

- So far, we assumed random disk accesses

- Would cost change if we assume that reads and writes are done sequentially?
  - Yes

- How can we incorporate this fact into our cost model?
  - Use bigger units (this is referred to as Blocked I/O)
  - Mask I/O delays through pre-fetching (this is referred to as double buffering)

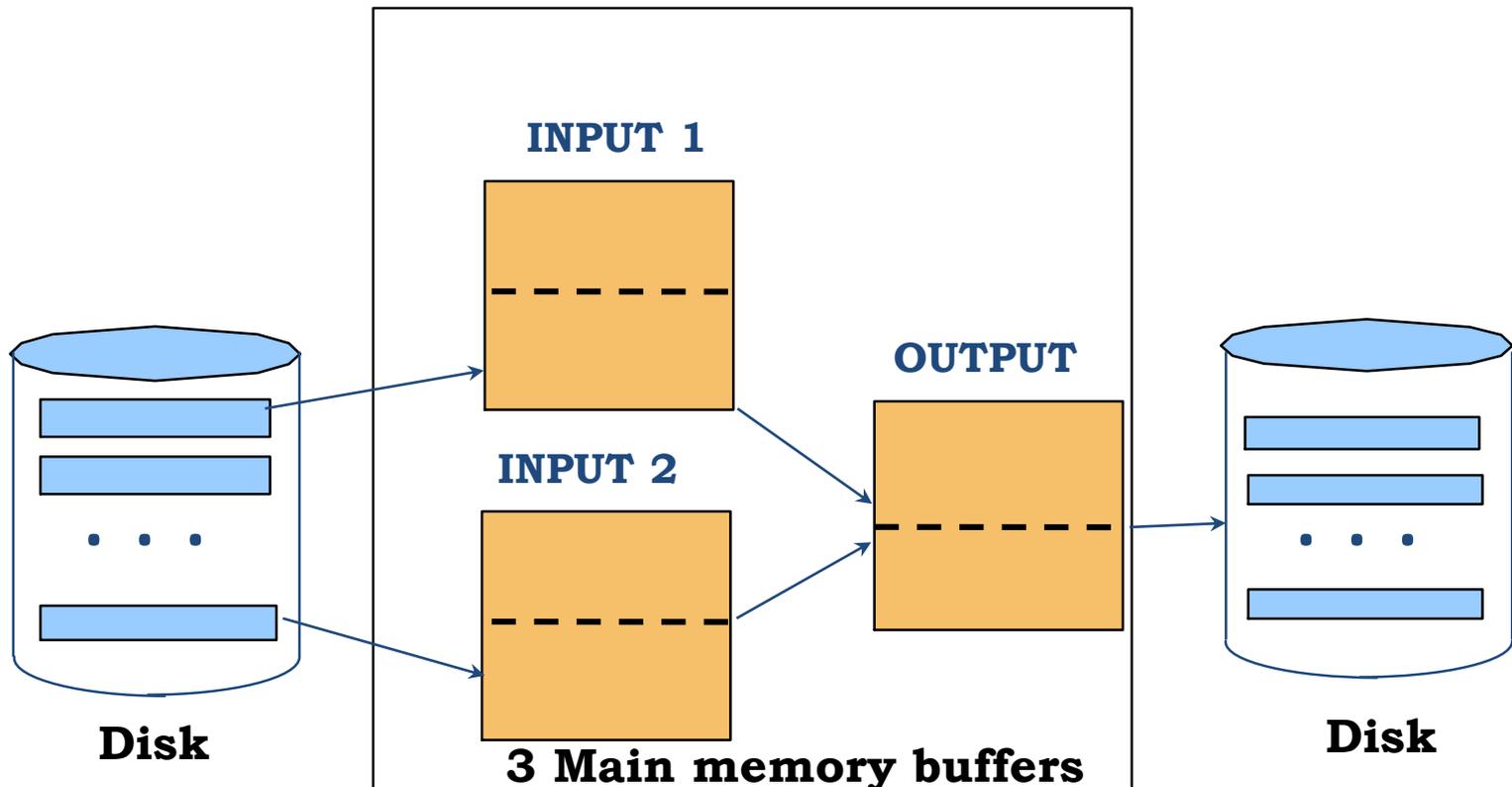# Blocked I/O

- Normally, we go with '$B$' buffers of size (say) 1 page

# Blocked I/O

- Normally, we go with '$B$' buffers of size (say) 1 page
- INSTEAD: let us go with $B/b$ buffers, of size '$b$' pages

**INPUT 1**

**INPUT 2**

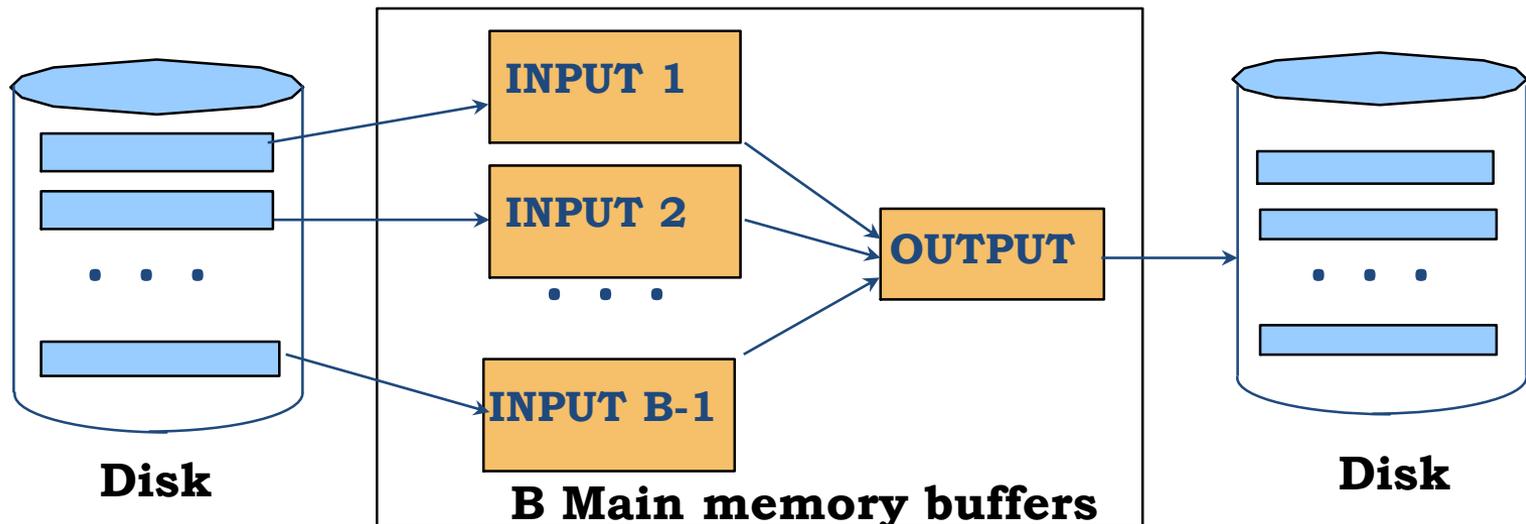**OUTPUT**

**Disk**

**3 Main memory buffers**

**Disk**

# Blocked I/O

- Normally, we go with '$B$' buffers of size (say) 1 page
- INSTEAD: let us go with $B/b$ buffers, of size '$b$' pages

- What is the main advantage?
  - Fewer random accesses (as some of the page will be arranged sequentially!)

- What is the main disadvantage?
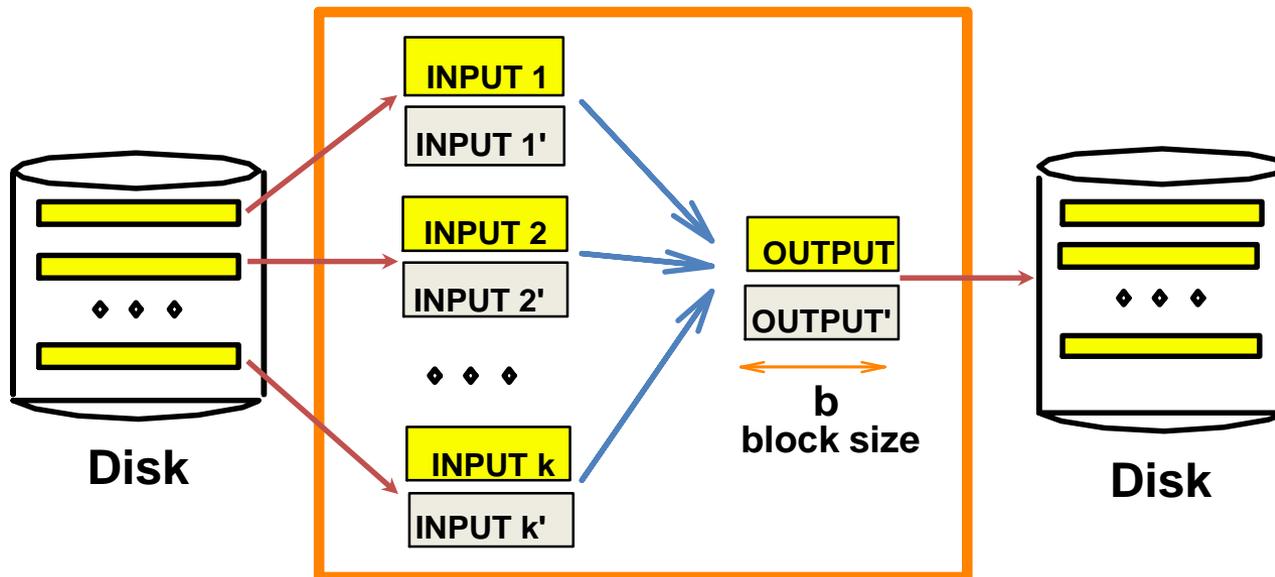  - Smaller fan-in and accordingly larger number of passes!

# Double Buffering

- Normally, when, say 'INPUT1' is exhausted
  - We issue a 'read' request and
  - We wait ...

# Double Buffering

- INSTEAD: *pre-fetch* INPUT1' into a `*shadow block*'
  - When INPUT1 is exhausted, issue a 'read'
  - BUT, also proceed with INPUT1'
  - Thus, the CPU can never go idle!



**B main memory buffers, k-way merge**

# Next Class