# CS222/CS122C: Principles of Data Management

# UCI, Fall 2019
# Notes #09

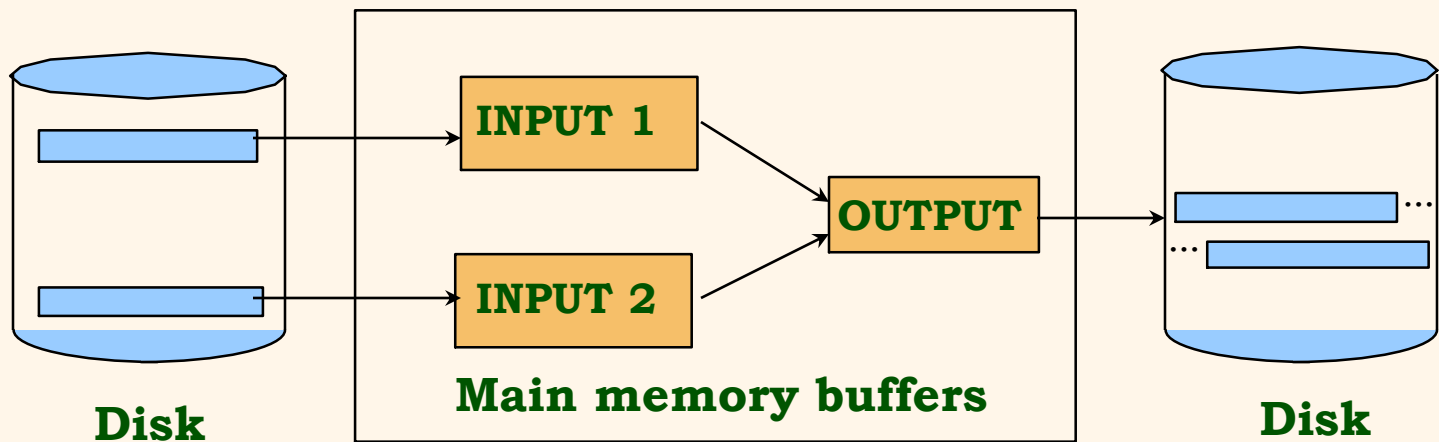## External Sorting

*Disk based*

Instructor: Chen Li

# *Why Sort?*

❖ A classic problem in computer science! (☺)

❖ Data sometimes requested in sorted order.
  - E.g., find students in decreasing *gpa* order.

❖ Sorting is first step in *bulk loading* B+ tree index.

❖ Sorting useful for eliminating *duplicate copies* (i.e., SELECT DISTINCT) in a collection of records. (Why?)

❖ *Sort-merge* join algorithm involves sorting.

❖ Problem: sort 1 TB of data with 1 GB of RAM.
  - *Q*: Why not just use virtual memory?

❖ Access heap file using sorted RIDs from B+ tree.

# *2-Way Sort: Requires 3 Buffers*

❖ Pass 1: Read a page, sort it, write it out
  - only one buffer page is used
❖ Pass 2, 3, …, etc.: Read and merge pairs of runs
  - three buffer pages are used:



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

# *Sorting N=2$^k$ Pages of Data*

❖ Pass 0:

  ▪ Read, sort, write ➔ 2$^k$ 1-page runs (subfiles)

❖ Pass 1:

  ▪ Read+merge 1-page pairs, write ➔ 2$^{k-1}$ 2-page runs

❖ Pass 2:

  ▪ Read+merge 2-page pairs, write ➔ 2$^{k-2}$ 4-page runs

...

❖ Pass *k*-1:

  ▪ Read+merge 2$^{k-2}$-page pairs, write ➔ 2 2$^{k-1}$-page runs

❖ Pass *k*:

  ▪ Read+merge 2$^{k-1}$-page pairs, write ➔ 1 2$^k$-page result

# *Two-Way External Merge Sort*
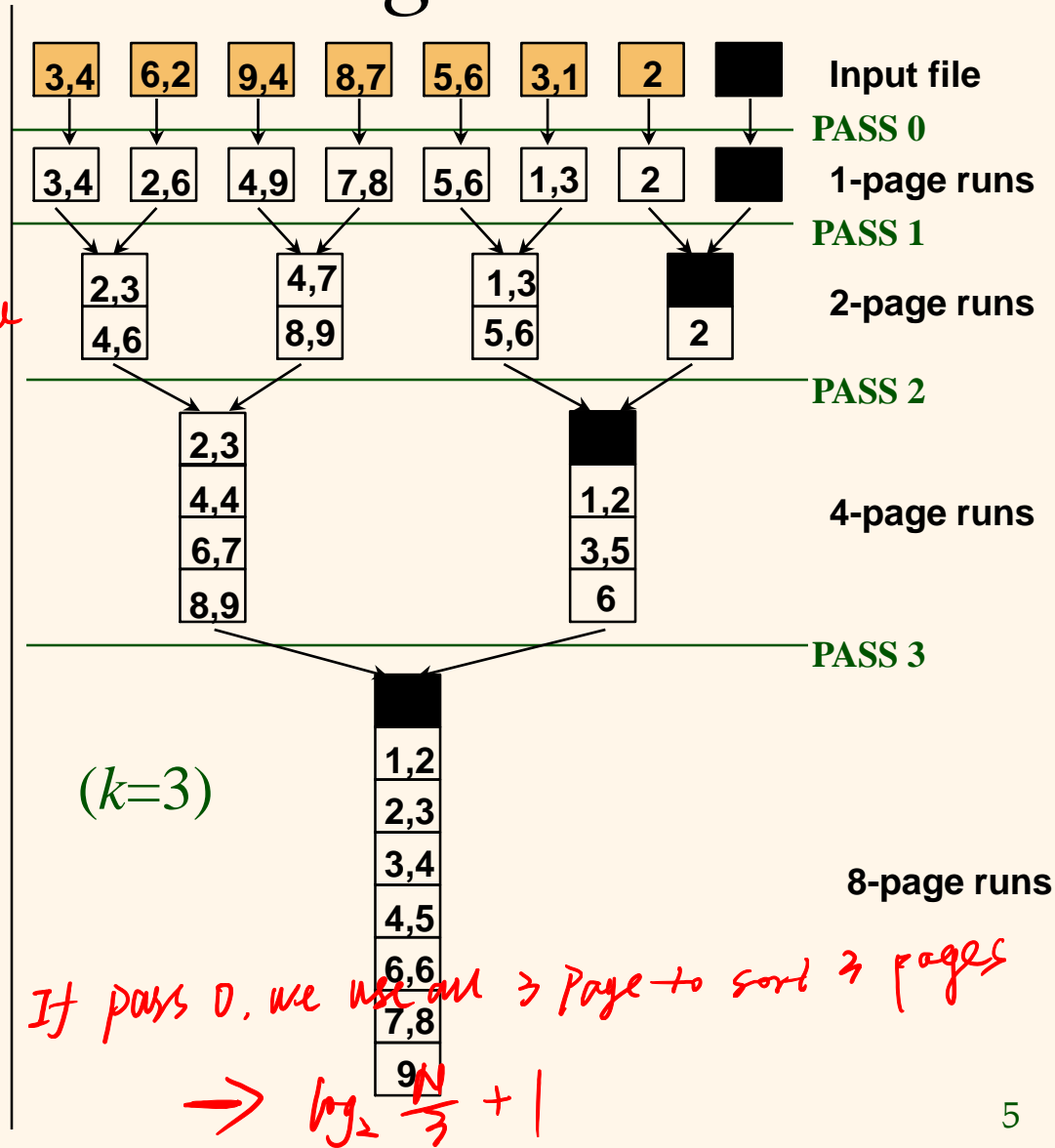
❖ Each pass we read + write each page in file.

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

❖ So total I/O cost is:
$$2N\left(\lceil \log_2 N \rceil + 1\right)$$

❖ *Idea:* **Divide and conquer:** sort subfiles and merge
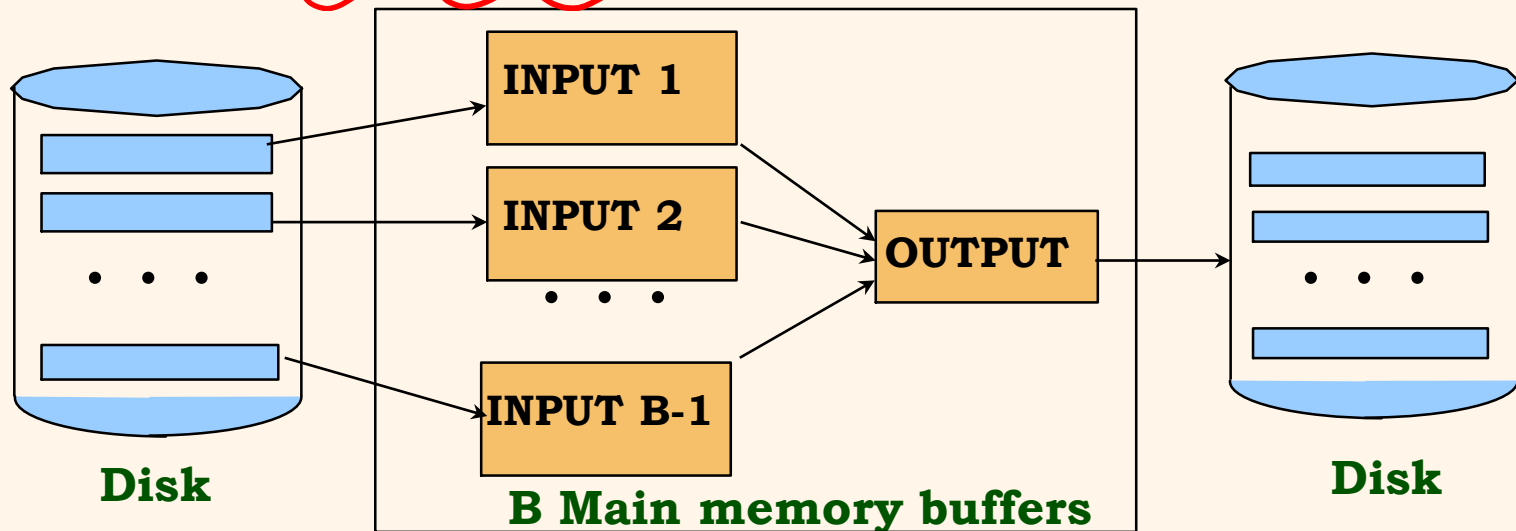
❖ *Q:* See any room to do better, w/just 3 pages?



**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 / 4,6 ; 4,7 / 8,9 ; 1,3 / 5,6 ; ■ / 2

PASS 2

**4-page runs**
2,3 / 4,4 / 6,7 / 8,9 ; 1,2 / 3,5 / 6

PASS 3

**8-page runs**
1,2 / 2,3 / 3,4 / 4,5 / 6,6 / 7,8 / 9

(*k*=3)

*(handwritten, red)* → Sort the page

*(handwritten, red)* If pass 0, we use all 3 Page to sort 3 pages
→ $\log_2 \frac{N}{3} + 1$

5

# *General External Merge Sort*

*(handwritten: { Read B pages use quick sort, and flush back to disk. No output page is lost.)*

## *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. *(Actually B-1 instead of B w/variable-length records.)*

- Pass 2, …,  etc.: merge $B-1$ runs.



**Disk**          **INPUT 1**   **INPUT 2**   **INPUT B-1**   **OUTPUT**          **Disk**

**B Main memory buffers**

# *Cost of External Merge Sort*
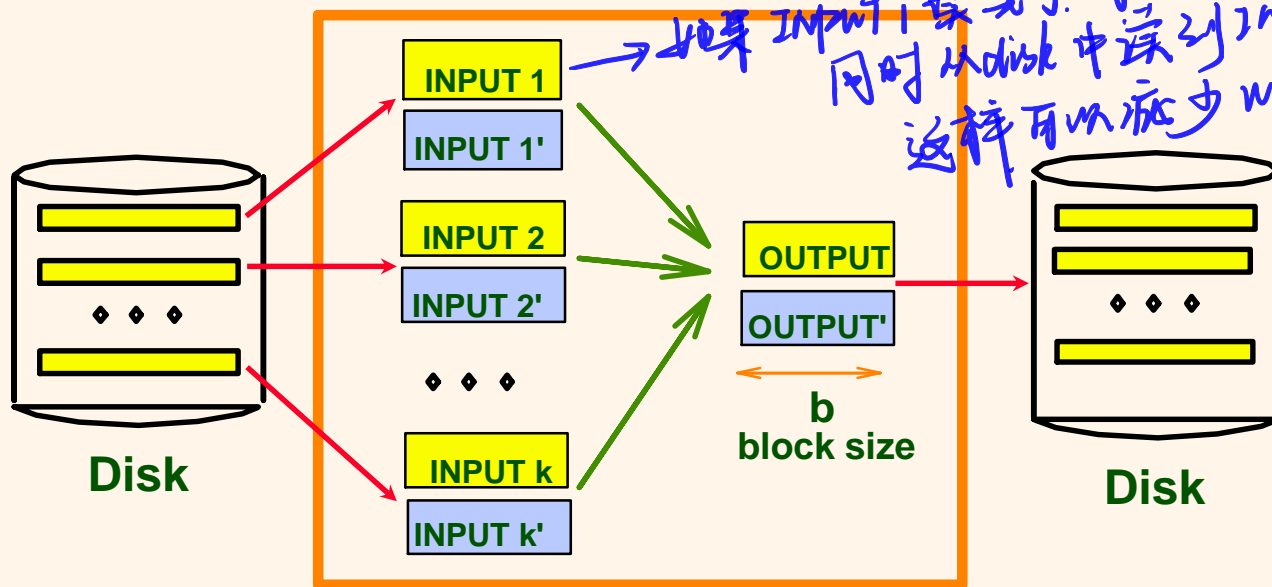
❖ Number of passes:   $1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil$

❖ Cost = 2N * (# of passes)

❖ E.g., with 5 buffer pages, to sort 108 page file:

*In memory sort B pages*

- 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 |

# *Double Buffering*

❖ To reduce wait time for I/O request to complete, can *prefetch* into `shadow block'.

- Potentially, more passes; in practice, most files *still* sorted in 2-3 passes.



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

# *Sorting Summary*

❖ External sorting is important; DBMS may dedicate part of buffer pool for sorting!

❖ External merge sort minimizes disk I/O cost:

  ▪ Pass 0:  Produce sorted *runs* of size *B* (# buffer pages), or of size *B-1* if we are handling variable-length records. Passes > 0:  *Merge* runs (until just one run is produced).

  ▪ # of runs merged at a time depends on *B* and *block size*.

  ▪ Larger block size means less I/O cost per page of data.

  ▪ Larger block size means fewer runs merged per step.

  ▪ In practice, # of passes needed rarely more than 2 or 3.

SQL :-        called  SPJ =      select

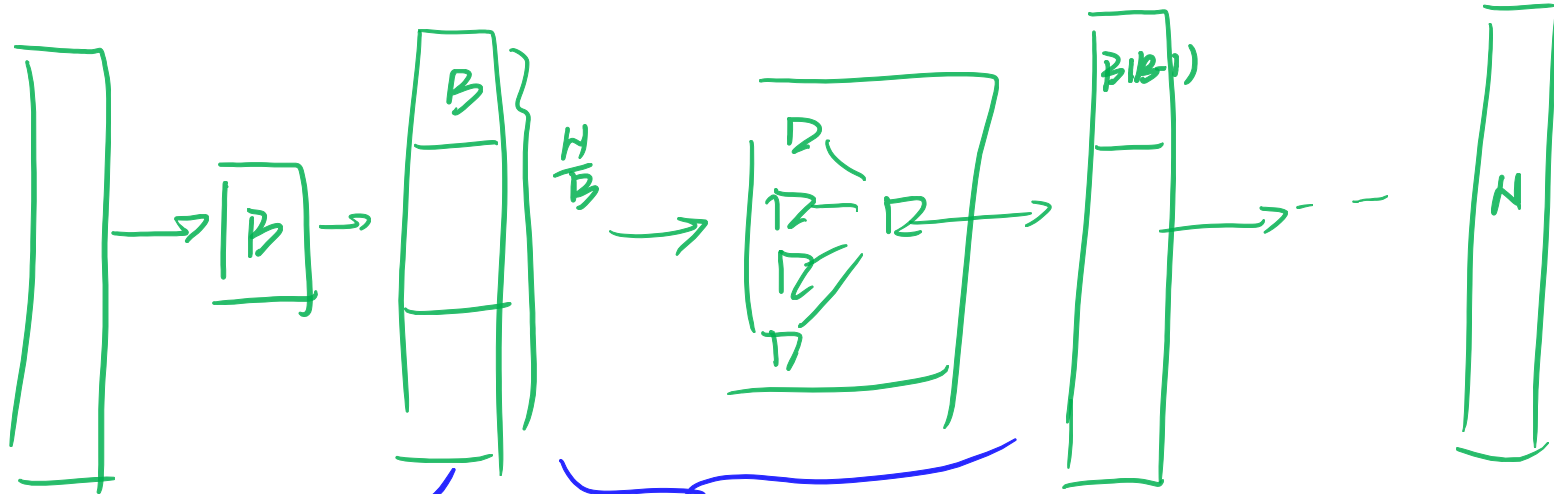                              From

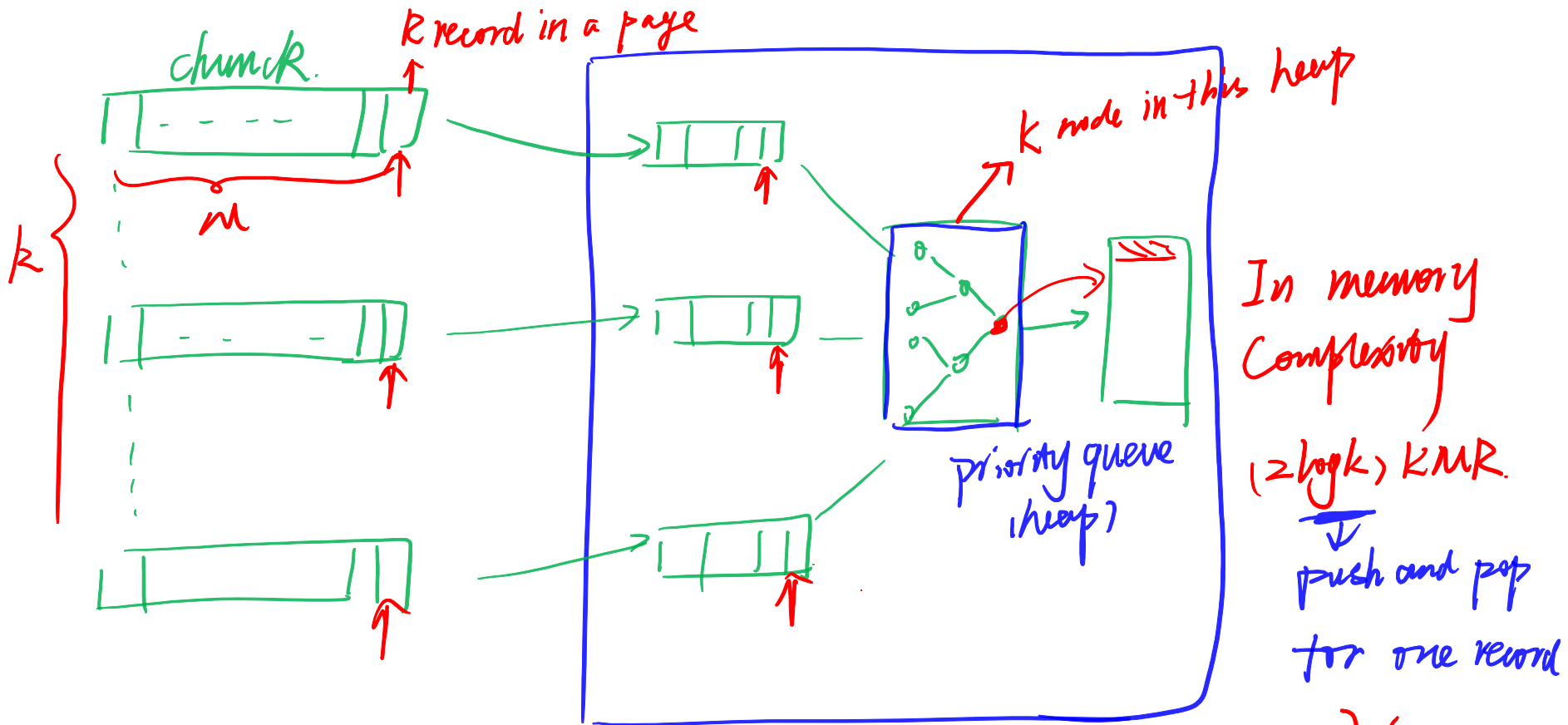                    selection   where.

                    projects    Group by

                    join        having

                                order by

$B$

$\dfrac{H}{B}$

$B(B-1)$

$N$

$BR \log_1(B \cdot R) \cdot \dfrac{N}{B}$

in-memory sort for B pages

$(\geq \log_1(B-1)) \, N R$

push/pop time for heap

chunk.

R record in a page

k

M

k node in this heap

Priority queue (heap)

In memory Complexity

$(\geq \log k)$ KMR.

push and pop for one record

$\geq \log(B-1)$ NR

- Build the initial heap
- keep a cursor with one input page
- use pop/push to go through each page.
- whenever the output page is full, flush it to disk
- if one input page is exhausted, read next page

Q: for N page, what's the minimum # of pages that can sort N pages using (2 iteration)

(pass 0, pass 1)

in-memory sort

merge sort

$$\begin{cases} \dfrac{N}{B} \leq B-1 \\ \text{or} \\ N \leq B(B-1) \end{cases}$$