



CS 443 External Sorting Chapters 13

Slides adapted from
Ramakrishnan & Gerhke
pages.cs.wisc.edu/~dbbook/

Renée J. Miller

Why Sort?

2

- 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 1 Gb of data with 1 Mb of RAM.
 - Why not virtual memory?

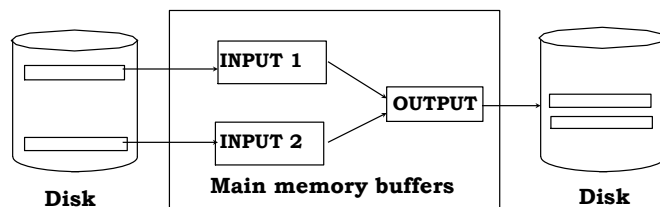
11/16/11



2-Way Sort: Requires 3 Buffers

3

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



11/16/11



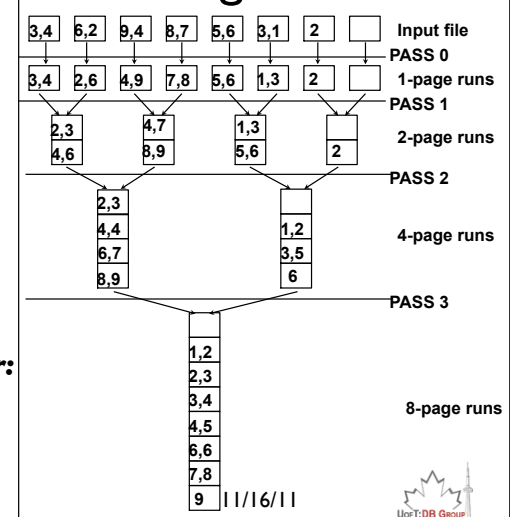
Two-Way External Merge Sort

4

- Each pass we read + write each page in file.
- N pages in the file \Rightarrow the number of passes

$$= \lceil \log_2 N \rceil + 1$$
- So total cost is:

$$2N(\lceil \log_2 N \rceil + 1)$$
- *Idea: Divide and conquer:*
sort subfiles and merge



11/16/11

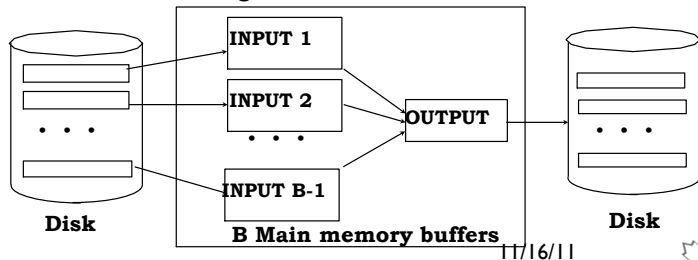


General External Merge Sort

5

➡ *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.



11/16/11



Cost of External Merge Sort

6

- Number of passes: $1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil$
- Cost = $2N * (\text{\# of passes})$
- 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

11/16/11



Number of Passes of External Sort

7

| 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 |

11/16/11



Using B+ Trees for Sorting

8

- 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!***

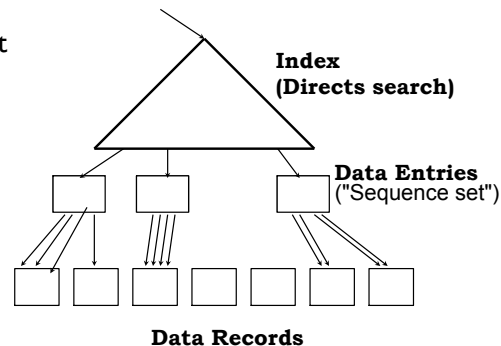
11/16/11



Clustered B+ Tree Used for

9

- Cost: root to the left-most leaf, then retrieve all leaf pages



☛ Always better than external sorting!

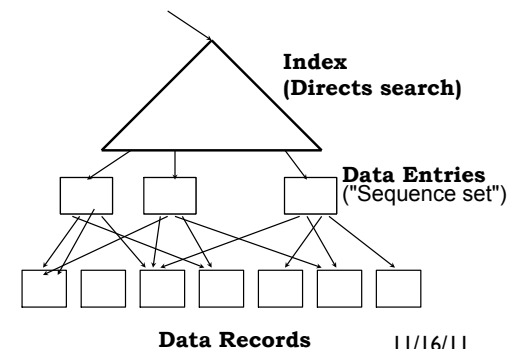
11/16/11



Unclustered B+ Tree Used for Sorting

10

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



11/16/11



External Sorting vs. Unclustered

11

| 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.

11/16/11



Summary

12

- 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.

11/16/11



Summary, cont.

13

- 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.

11/16/11

