



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



CSC3170

Tutorial 8

School of Data Science

The Chinese University of Hong Kong, Shenzhen

Outline

- 2-Way External Merge Sort
- General External Merge Sort
- Example
- Practice Questions

2-Way External Merge Sort

- We will start with a simple example of a 2-way external merge sort.
 - “2” is the number of runs that we are going to merge into a new run for each pass.
- Data is broken up into N pages.
- The DBMS has a finite number of B buffer pool pages to hold input and output data.

Simplified 2-Way External Merge Sort

- **Pass #0**

- Read one page of the table into memory
- Sort page into a “run” and write it back to disk
- Repeat until the whole table has been sorted into runs

- **Pass #1,2,3,...**

- Recursively merge pairs of runs into runs twice as long
- Need at least 3 buffer pages (2 for input, 1 for output)

General External Merge Sort

- **Pass #0**
 - Use B buffer pages
 - Produce $\lceil N / B \rceil$ sorted runs of size B
- **Pass #1,2,3,...**
 - Merge $B-1$ runs (i.e., K-way merge)
- Number of passes = $1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil$
- Total I/O Cost = $2N \cdot (\# \text{ of passes})$

Example

- Determine how many passes it takes to sort 108 pages with 5 buffer pool pages: $N=108$, $B=5$
 - **Pass #0:** $[N / B] = [108 / 5] = 22$ sorted runs of 5 pages each (last run is only 3 pages).
 - **Pass #1:** $[N' / B-1] = [22 / 4] = 6$ sorted runs of 20 pages each (last run is only 8 pages).
 - **Pass #2:** $[N'' / B-1] = [6 / 4] = 2$ sorted runs, first one has 80 pages and second one has 28 pages.
 - **Pass #3:** Sorted file of 108 pages.
- $1 + \lceil \log_{B-1} [N / B] \rceil = 1 + \lceil \log_4 22 \rceil = 1 + \lceil 2.229... \rceil = 4 \text{ passes}$

Practice Questions

1. You are trying to sort the Students table which has 1960 pages with 8 available buffer pages.
 - a. How many sorted runs will be produced after each pass?
 - b. How many pages will be in each sorted run for each pass?
 - c. How many I/Os does the entire sorting operation take?
 - Pass #0: $[N/B] = [1960/8] = 245$ sorted runs of 8 pages each
 - Pass #1: $[N'/B - 1] = [245/7] = 35$ sorted runs of $8*7 = 56$ pages each
 - Pass #2: $[N''/B - 1] = [35/7] = 5$ sorted runs of $56*7 = 392$ pages each
 - Pass #3: The next merging pass can merge all remaining sorted runs (because there are <7 sorted runs) so it will produce 1 sorted run of all 1960 pages.

 - Therefore the answer to a is: **245, 35, 5, 1** and the answer to b is: **8, 56, 392, 1960**.
 - Each pass takes $2 * N$ I/Os where N is the total number of data pages because each page gets read and written in a pass. This means that the answer to c is: $4 * 2 * 1960 = 15,680$ I/Os.

Practice Questions

2. You are trying to sort the Sailors table which has 200 pages. Suppose that during Pass 0, you have 10 buffer pages available to you, but for Pass 1 and onwards, you only have 5 buffer pages available.
- How many sorted runs will be produced after each pass?
 - How many pages will be in each sorted run for each pass?
 - How many I/Os does the entire sorting operation take?
- Pass #0: $[N/B_1] = [200/10] = 20$ sorted runs of 10 pages each
 - Pass #1: $[N'/B_2 - 1] = [20/4] = 5$ sorted runs of $10 * 4 = 40$ pages each
 - Pass #2: $[N''/B_2 - 1] = [5/4] = 2$ sorted runs (1 with $40 * 4 = 160$ pages, 1 with 40 pages)
 - Pass #3: The next merging pass can merge all remaining sorted runs (because there are <4 sorted runs) so it will produce 1 sorted run of all 200 pages.
-
- Therefore the answer to a is: **20, 5, 2, 1** and the answer to b is: **10, 40, (160 and 40), 200**.
 - Each pass takes $2 * N$ I/Os where N is the total number of data pages because each page gets read and written in a pass. This means that the answer to c is: $4 * 2 * 200 = 1600$ I/Os.

Practice Questions

1. What is the minimum number of buffer pages that we need to sort 1000 data pages in two passes?

$$\frac{1000}{B(B - 1)} \leq 1$$

B = 32.1 which means we need **33 buffer pages**.

2. True or False: Increasing the number of buffer pages does not affect the number of I/Os performed in Pass 0 of an external sort.

True. Regardless of the number of buffer pages, every pass of an external sort (including Pass 0) performs the same number of IOs.

3. True or False: Sorting an already sorted table takes fewer IOs than sorting a randomly-arranged table.

False. Regardless of how a table (of fixed length) is sorted, the sorting algorithm always takes the same number of IOs.

Q&A