

CMU SCS

Carnegie Mellon Univ.  
Dept. of Computer Science  
15-415 - Database Applications

Lecture 12: external sorting  
(R&G ch. 13)

Faloutsos 15-415 1

---

---

---

---

---

---

---



CMU SCS

Why Sort?

Faloutsos 15-415 2

---

---


---

---

---

---

---



CMU SCS

Why Sort?

- select ... order by
  - e.g., find students in increasing *gpa* order
- bulk loading B+ tree index.
- duplicate elimination (select distinct)
- select ... group by
- Sort-merge join algorithm involves sorting.

Faloutsos 15-415 3

---

---

---

---

---

---

---

CMU SCS

## Outline

- ➔ two-way merge sort
  - external merge sort
  - fine-tunings
  - B+ trees for sorting

Faloutsos 15-415 4

---

---

---

---

---

---

---

---

CMU SCS

## 2-Way Sort: Requires 3 Buffers

- Pass 0: Read a page, sort it, write it.
  - only one buffer page is used
- Pass 1, 2, 3, ..., etc.: requires 3 buffer pages
  - merge pairs of **runs** into runs twice as long
  - three buffer pages used.

Faloutsos 15-415 5

---

---

---

---

---

---

---

---

CMU SCS

## Two-Way External Merge Sort

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

Faloutsos 15-415 6

---

---

---

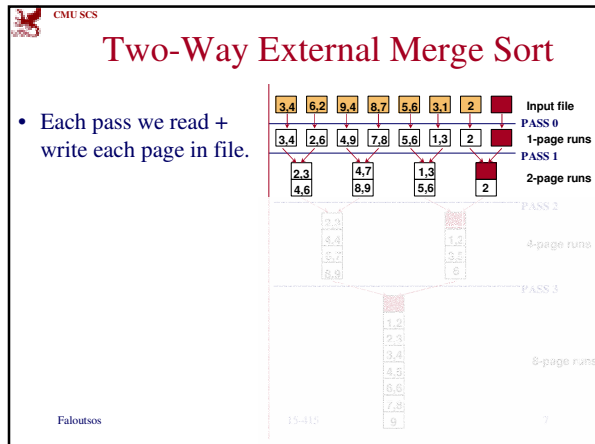
---

---

---

---

---




---

---

---

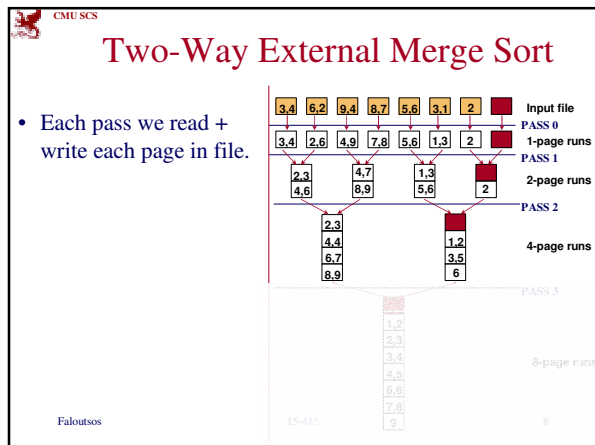
---

---

---

---

---




---

---

---

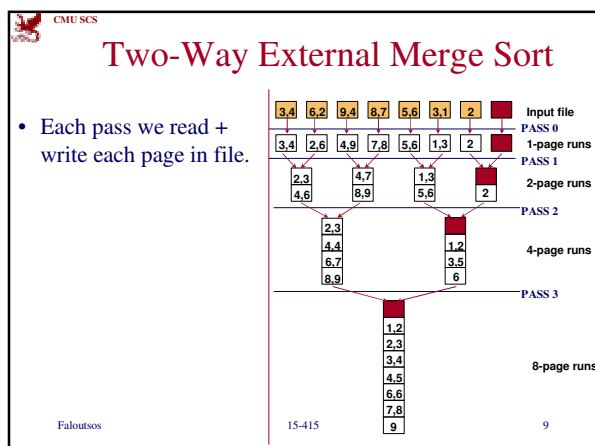
---

---

---

---

---




---

---

---

---

---

---

---

---

CMU SCS

15-415

## Two-Way External Merge Sort

- Each pass we read + write each page in file.
- N pages in the file =>  $\lceil \log_2 N \rceil + 1$
- So total cost is:  
 $2N(\lceil \log_2 N \rceil + 1)$
- Idea: Divide and conquer:* sort subfiles and merge

Input file  
PASS 0  
1-page runs  
PASS 1  
2-page runs  
PASS 2  
4-page runs  
PASS 3  
8-page runs

15-415 10

---

---

---

---

---

---

---

---

CMU SCS

15-415

## Outline

- two-way merge sort
- ➡ external merge sort
- fine-tunings
- B+ trees for sorting

Faloutsos 15-415 11

---

---

---

---

---

---

---

---

CMU SCS

15-415

## External merge sort

B > 3 buffers

- Q1: how to sort?
- Q2: cost?

Faloutsos 15-415 12

---

---

---

---

---

---

---

---

CMU SCS

## General External Merge Sort

$B > 3$  buffer pages. How to sort a file with  $N$  pages?

Faloutsos 15-415 13

---

---

---

---

---

---

---

---

CMU SCS

## General External Merge Sort

- Pass 0: use  $B$  buffer pages. Produce  $\lceil N / B \rceil$  sorted runs of  $B$  pages each.
- Pass 1, 2, ..., etc.: merge  $B-1$  runs.

Faloutsos 15-415 14

---

---

---

---

---

---

---

---

CMU SCS

## Sorting

- create sorted runs of size  $B$  (how many?)
- merge them (how?)

Faloutsos 15-415 15

---

---

---

---

---

---

---

---

CMU SCS

## Sorting

- create sorted runs of size B
- merge first B-1 runs into a sorted run of  $(B-1) * B, \dots$

Faloutsos 15-415 16

---

---

---

---

---

---

---

---

CMU SCS

## Sorting

- How many steps we need to do?  
‘i’, where  $B * (B-1)^i > N$
- How many reads/writes per step?  $N+N$

Faloutsos 15-415 17

---

---

---

---

---

---

---

---

CMU SCS

## Cost of External Merge Sort

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

Faloutsos 15-415 18

---

---

---


---

---

---

---

---



CMU SCS

### Cost of External Merge Sort

- 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

Formula check:  $\lceil \log_4 22 \rceil = 3 \dots + 1 \rightarrow$  4 passes ✓

Faloutsos

15-415

19

---

---


---

---

---

---

---



CMU SCS

### Number of Passes of External Sort

(I/O cost is 2N times number of passes)

| 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     |

Faloutsos

15-415

20

---

---


---

---

---

---

---



CMU SCS

### Outline

- two-way merge sort
- external merge sort
- ➡ fine-tunings
- B+ trees for sorting

Faloutsos

15-415

21

---

---


---

---

---

---

---



CMU SCS

Outline

- two-way merge sort
- external merge sort
- fine-tunings
  - ➡ – which internal sort for Phase 0?
  - blocked I/O
- B+ trees for sorting

Faloutsos

15-415

22

---

---


---

---

---

---

---



CMU SCS

Internal Sort Algorithm

- Quicksort is a fast way to sort in memory.
- But: we get B buffers, and produce 1 run of length B.
- Can we produce longer runs than that?

Faloutsos

15-415

23

---

---


---

---

---

---

---



CMU SCS

Internal Sort Algorithm

- Quicksort is a fast way to sort in memory.
- But: we get B buffers, and produce 1 run of length B.
- Can we produce longer runs than that?

B=3

B=3

**Heapsort:**

- Pick smallest
- Output
- Read from **next** buffer

Faloutsos

15-415

24

---

---

---

---

---

---

---



CMU SCS

## Internal Sort Algorithm

- Quicksort is a fast way to sort in memory.
- But: we get B buffers, and produce 1 run of length B.
- Can we produce longer runs than that?
- Alternative: “tournament sort” (a.k.a. “heapsort”, “replacement selection”)
- Produces runs of length  $\sim 2*B$
- Clever, but not implemented, for subtle reasons: tricky memory management on variable length records

Faloutsos 15-415 25

---

---

---

---

---

---

---

---

CMU SCS

## Reminder: Heapsort

pick smallest, write to output buffer:

```

graph TD
    10[10] --> 14[14]
    10 --> 11[11]
    14 --> 15[15]
    14 --> 17[17]
    11 --> 18[18]
    11 --> 16[16]
  
```

Faloutsos 15-415 26

---

---

---

---

---

---

---

---

CMU SCS

## Heapsort:

pick smallest, write to output buffer:

```

graph TD
    10[10] --> Root[...]
    Root --> 14[14]
    Root --> 11[11]
    14 --> 15[15]
    14 --> 17[17]
    11 --> 18[18]
    11 --> 16[16]
  
```

Faloutsos 15-415 27

---

---

---

---

---

---

---

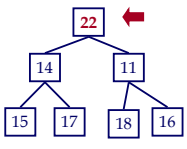
---

CMU SCS

details

## Heapsort:

get next key; put at top and 'sink' it



```
graph TD; 22[22] --> 14[14]; 22 --> 11[11]; 14 --> 15[15]; 14 --> 17[17]; 11 --> 18[18]; 11 --> 16[16];
```

Faloutsos 15-415 28

---

---

---

---

---

---

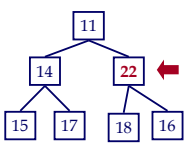
---

CMU SCS

details

## Heapsort:

get next key; put at top and 'sink' it



```
graph TD; 11[11] --> 14[14]; 11 --> 22[22]; 14 --> 15[15]; 14 --> 17[17]; 22 --> 18[18]; 22 --> 16[16];
```

Faloutsos 15-415 29

---

---

---

---

---

---

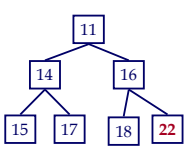
---

CMU SCS

details

## Heapsort:

get next key; put at top and 'sink' it



```
graph TD; 11[11] --> 14[14]; 11 --> 16[16]; 14 --> 15[15]; 14 --> 17[17]; 16 --> 18[18]; 16 --> 22[22];
```

Faloutsos 15-415 30

---

---

---

---

---

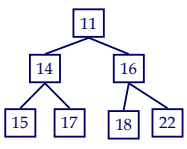
---

---

CMU SCS

details

## Heapsort:



```

graph TD
    11[11] --> 14[14]
    11 --> 16[16]
    14 --> 15[15]
    14 --> 17[17]
    16 --> 18[18]
    16 --> 22[22]
  
```

When done, pick top (= smallest) and output it, if 'legal' (ie.,  $\geq 10$  in our example)

This way, we can keep on reading new key values (beyond the B ones of quicksort)

Faloutsos 15-415 31

---

---

---

---

---

---

---

---

CMU SCS

## Outline

- two-way merge sort
- external merge sort
- fine-tunings
  - which internal sort for Phase 0?
- ➡ – blocked I/O
- B+ trees for sorting

Faloutsos 15-415 32

---

---

---

---

---

---

---

---

CMU SCS

## Blocked I/O & double-buffering

- So far, we assumed random disk access
- Cost changes, if we consider that runs are written (and read) sequentially
- What could we do to exploit it?

Faloutsos 15-415 33

---

---

---

---

---

---

---

---

CMU SCS

## Blocked I/O & double-buffering

- So far, we assumed random disk access
- Cost changes, if we consider that runs are written (and read) sequentially
- What could we do to exploit it?
- A1: Blocked I/O (exchange a few r.d.a for several sequential ones)
- A2: double-buffering

Faloutsos 15-415 34

---

---

---

---

---

---

---

---

CMU SCS

## 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*.

Faloutsos 35

---

---

---

---

---

---

---

---

CMU SCS

## Outline

- two-way merge sort
- external merge sort
- fine-tunings
- ➔ • B+ trees for sorting

Faloutsos 15-415 36

---

---

---


---

---

---

---

---



CMU SCS

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
  - B+ tree is not clustered

Faloutsos

15-415

37

---

---


---

---

---

---

---



CMU SCS

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!

Faloutsos

15-415

38

---

---


---

---

---

---

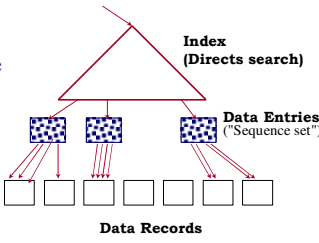
---



CMU SCS

Clustered B+ Tree Used for Sorting

- Cost: root to the left-most leaf, then retrieve all leaf pages (Alternative 1)



Faloutsos

15-415

39

Always better than external sorting!

---

---


---

---

---

---

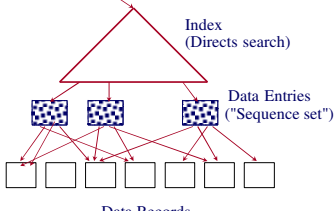
---



CMU SCS

## Unclustered B+ Tree Used for Sorting

- Alternative (2) for data entries; each data entry contains *rid* of a data record. In general, *one I/O per data record!*



Faloutsos

40

---

---

---


---

---

---

---

---



CMU SCS

## 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. <sup>41</sup>

Faloutsos

---

---

---


---

---

---

---

---



CMU SCS

## Summary

- External sorting is important
- External merge sort minimizes disk I/O cost:
  - Pass 0: Produces sorted *runs* of size *B* (# buffer pages).
  - Later passes: *merge* runs.
- Clustered B+ tree is good for sorting; unclustered tree is usually very bad.

Faloutsos

15-415

42

---

---

---

---

---

---

---

---

14