CSE462/562: Database Systems (Spring 22) Lecture 13: External sorting 4/5/2022



What is external sorting/hashing?

- Problem: sort or hashing 1TB of data over 1GB of RAM
 - Why not virtual memory?
 - Swaps involve expensive random I/Os
 - Why not using B-Tree/extendible hashing/linear hashing?
 - Dynamic structures carry additional overhead for maintenance (not needed in QP)
 - Missing optimization opportunities with hybrid approach (see later)
- General wisdom:
 - I/O cost dominates the total cost
 - Design algorithms to reduce the number of I/Os

Two-way merge-sort: a starting point

- Recall the two-way merge-sort
 - given a list of items in A[0..n-1]
 - recursively divide and conquer the problem
 - divide the list into two halves $A_1\left[0, \left\lfloor \frac{n}{2} \right\rfloor\right]$, $A_2\left[\left\lfloor \frac{n}{2} \right\rfloor + 1, n-1\right]$
 - merge-sort A_1 and A_2 individually
 - merge the two sorted list A_1, A_2



External two-way merge sort

- Needs 3 buffers
- Instead of recursion
 - works bottom up from the input



External two-way merge sort

- Needs 3 buffers
- Instead of recursion
 - works bottom-up from the input





External two-way merge sort

• Input: N pages

Disk file

- Cost for a pass: reading & writing N pages once •
- # of passes: height of the tree = $[\log_2 N] + 1$ •
- Total cost: $2N(\lfloor \log_2 N \rfloor + 1) \rfloor / Os$
 - Transfer cost: $2t_T N([\log_2 N] + 1)$
 - Seek cost: $2t_S N([log_2 N] + 1)$
 - $total = 2(t_T + t_S)N([\log_2 N] + 1)$



External multi-way merge sort

- How do we fully utilize all the *M* buffers?
 - Solution: (M-1)-way merge-sort
- Pass 0: internal sort to produce initial runs
 - read every *M* pages into memory
 - use some internal sorting algorithm (e.g., quick sort)
 - can produce even larger runs (later)
 - write all the *M* pages as a run

INPUT 1

INPUT 2

• • •

INPUT M

Disk

N pages in input $\left[\frac{N}{M}\right]$ runs after pass 0 Cost: 2N pages read/written + $2\left[\frac{N}{M}\right]$ seeks i.e. $2Nt_T + 2\left[\frac{N}{M}\right]t_S$ Input file 6 6,2 9,4 8,7 5,6 3. PASS 0 2,6 7 1,3 5,6 7,8 6,9



Disk

- Pass 1, 2, ...: merge as many runs as possible from previous pass into a sorted run
 - maintain *a min-heap/max-heap* (aka priority queue)
 - supports O(log M) time insertion of any item and deletion of the smallest/largest item
 - a complete binary tree where parent is smaller/larger than both children
 - how to implement
 - numbering nodes level by level sequentially from 1, store in an array A[1..n]
 - (how to translate 1-based index to 0-based in C/C++?)
 - parent of A[i] is A[i/2], left child of A[i] is A[i * 2], right child of A[i] is A[i * 2 + 1]
 - push-down or push-up to maintain the variant



- Pass 1, 2, ...: merge as many runs as possible from previous pass into a sorted run
 - maintain *a min-heap*
 - load one page from each of the M 1 runs
 - and maintain pointers of next page to read
 - for each loaded page
 - insert the first key into the min-heap
 - maintain next slot ids for each page
 - Repeatedly remove the smallest item from the min heap
 - and replace it with the next key in its run
 - write out the output page once it's full



PASS 1

- Pass 1, 2, ...: merge as many runs as possible from previous pass into a sorted run
 - maintain *a min-heap*
 - load one page from each of the M 1 runs
 - and maintain pointers of next page to read
 - for each loaded page
 - insert the first key into the min-heap
 - maintain next slot ids for each page
 - Repeatedly remove the smallest item from the min heap
 - and replace it with the next key in its run
 - write out the output page once it's full







2,3

4,4

6,9

next



- Pass 1, 2, ...: merge as many runs as possible from previous pass into a sorted run
 - maintain *a min-heap*
 - load one page from each of the M 1 runs
 - and maintain pointers of next page to read
 - for each loaded page
 - insert the first key into the min-heap
 - maintain next slot ids for each page
 - Repeatedly remove the smallest item from the min heap





- Pass 1, 2, ...: merge as many runs as possible from previous pass into a sorted run
 - maintain *a min-heap*
 - load one page from each of the M 1 runs
 - and maintain pointers of next page to read
 - for each loaded page
 - insert the first key into the min-heap
 - maintain next slot ids for each page
 - Repeatedly remove the smallest item from the min heap
 - and replace it with the next key in its run
 - write out the output page once it's full





2,3

- Pass 1, 2, ...: merge as many runs as possible from previous pass into a sorted run
 - maintain *a min-heap*
 - load one page from each of the M-1 runs
 - and maintain pointers of next page to read
 - for each loaded page
 - insert the first key into the min-heap
 - maintain next slot ids for each page
 - Repeatedly remove the smallest item from the min heap
 - and replace it with the next key in its run
 - write out the output page once it's full



N pages to read/write per pass $\begin{bmatrix} log_{M-1} \\ \hline{M} \end{bmatrix}$ merge passes Cost per merge pass: 2N pages read/written + 2N seeks Total cost for merge passes: $2(t_T + t_S)N[log_{M-1} [\frac{N}{M}]]$



Cost analysis

- Cost analysis:
 - Pass 0: $2Nt_T + 2\left[\frac{N}{M}\right]t_S$

- gain of utilizing all available buffers
- importance of a high fan-in during merging
- Pass 1, 2, ... combined: $2(t_T + t_S)N[\log_{M-1}[\frac{N}{M}]]$
- Total = $2t_T N\left(\left[log_{M-1}\left[\frac{N}{M}\right]\right] + 1\right) + 2t_S\left(\left[\frac{N}{M}\right] + N\left[log_{M-1}\left[\frac{N}{M}\right]\right]\right)$

Ν	M=3	=5	=9	=17	=129	=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

• Can we do it better?

Batching I/Os for merge sort

- Refinement 1
 - reducing random I/Os by reading/writing *B* pages per run during merge
 - using (M 1)-way merge sort
 - memory usage increases to *MB* pages
 - number of pages transferred do not change
 - but the number of random seeks per merge pass reduced to approximately $2\left[\frac{N}{P}\right]$
 - total cost reduced to $2t_T N\left(\left[log_{M-1}\left[\frac{N}{MB}\right]\right] + 1\right) + 2t_S\left(\left[\frac{N}{MB}\right] + \left[\frac{N}{B}\right]\left[log_{M-1}\left[\frac{N}{MB}\right]\right]\right)^{\frac{1}{2}}$



Exercise: what if we only have M pages instead of MB pages and still read/write pages in B-page batches?

$$2t_T N\left(\left\lceil \log_{\lfloor\frac{M}{B}\rfloor-1} \left\lceil \frac{N}{M} \right\rceil \right\rceil + 1\right) + 2t_S\left(\left\lceil \frac{N}{M} \right\rceil + \lceil \frac{N}{B} \rceil \lceil \log_{\lfloor\frac{M}{B}\rfloor-1} \lceil \frac{N}{M} \rceil \rceil\right)$$

Pipelining output

- Refinement 2
 - in most cases, do not need to write the final file
 - pipelining to the next operator
 - or output to user
 - Hence, no need to count the write of the final pass
 - total cost reduced to $t_T N\left(2\left[log_{\left\lfloor\frac{M}{B}\right\rfloor-1}\left\lceil\frac{N}{M}\right\rceil\right]+1\right)+t_S\left(2\left\lceil\frac{N}{M}\right\rceil+\left\lceil\frac{N}{B}\right\rceil(2\left\lceil log_{\left\lfloor\frac{M}{B}\right\rfloor-1}\left\lceil\frac{N}{M}\right\rceil\right\rceil-1)\right)$



Tournament sort

- Refinement 3
 - producing initial runs as large as possible in pass 0
 - Alternative to quick-sort: "tournament sort" (a.k.a. "heapsort", "replacement selection")
- Keep two heaps in memory, H1 and H2, reserve an input buffer page and an output buffer page read M-2 pages of records, inserting into H1; while (records left) { m = H1.removemin(); put m in output buffer; if (H1 is empty) swap H1 and H2 (pointer swap only!); start new output run; else read in a new record r (use 1 buffer for input pages); if (r < m) H2.insert(r); else H1.insert(r); H1.output(); start new run; H2.output();

Tournament sort

• Tournament sort explained:



- <u>1 input, 1 output, M 2 for current and next set (min heaps)</u>
- Main idea: ensure the *smallest* key in the <u>current set (H1)</u> is *greater* than any key that has been written to this output run.
 - If it can't be satisfied, write to the next set (H2), which goes into the next run.
- Memory usage of the min-heaps combined never exceeds the M-2 pages

Tournament sort

- Fact: average length of a run is 2(M-2)
- Total cost reduced to on average

$$t_T N\left(2\left[\log_{\left\lfloor\frac{M}{B}\right\rfloor-1}\left\lceil\frac{N}{2M-4}\right\rceil\right]+1\right)+t_S\left(2\left\lceil\frac{N}{2M-4}\right\rceil+\lceil\frac{N}{B}\right](2\lceil\log_{\left\lfloor\frac{M}{B}\right\rfloor-1}\lceil\frac{N}{2M-4}\rceil)-1)\right)$$

- Worst-Case:
 - What is min length of a run?
 - How does this arise?
- Best-Case:
 - What is max length of a run?
 - How does this arise?
- Quicksort is faster, but ... longer runs often means fewer passes!

- 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

Good idea since it's already available! Could be a very bad idea! (Random I/O) unless all columns are included in the key

Certain basic operator implementation w/ sorting

- Some basic operators can be implemented on top of sorting
 - Can use pipelining over the sort results
- Examples
 - deduplication (projection in standard RA)
 - maintain the last key
 - for each output from the sort
 - emit it if it is different from the last key
 - otherwise, discard it
 - aggregation
 - maintain the aggregation state
 - for each output from the sort
 - emit the finalized aggregation value if it is different from the last key (unless this is the first)
 - otherwise, accumulate it to the state
 - exercise: work out the details of $\cup, \cap, -$
- No additional I/O due to pipelining
 - can support rewinding (why?)

This lecture

- Summary:
 - External sorting (multi-way merge-sort)
 - Certain operator implementation using sorting
- Next lecture
 - join algorithms
 - nested loop
 - index nested loop
 - hash join and hybrid hashing