

Lecture #2


Goal: from RAM to multi-disks in sorting obj \leftarrow atomic stripes

Quicksort (A, l, r)

~~while (r-l+1 > Mo)~~ while ($r-l+1 > Mo$)

$j = \text{pickPivotPos}(A, l, r);$

swap ($A[l], A[j]$);

$i = \text{Partition}(A, l, r);$ // pivot is in $A[l]$ initially, ~~the~~ 
 // $A[l..i-1]$ less or equal than pivot
 // $A[i..r]$ more or equal than pivot

if $i < \frac{l+r}{2}$ then { Quicksort ($A, l, i-1$); $l = i$; }

else { Quicksort (A, i, r); $r = i-1$; }

insertion sort (A, l, r);

Comments:

1) Do exist various Partition procedures



2) Pivot selection makes a difference

1 random
 always first
 median out of 3
 exact median
 skewed pivot $n/10$

• Last is faster because the reduced cost of "branch misprediction" (hence, less cost of flushing the pipeline) outweighs the cost induced by more recursive calls: $\log_{\frac{1}{1-\alpha}} n \approx \frac{1}{\alpha} \cdot \log n$

3) Code is structured to eliminate "tail recursion", which ensures only $\log_2 n/10$ recursive calls at most (this was not true on classic quicksort.)

•) Use of Insertion Sort is justified by cache and small no.

SELECTION: The quick-scheme can be adapted to implement a linear time (avg) procedure to select the k -th smallest element.

select(s, k)

// assert: $|s| \geq k$

pick $p \in s$ uniformly at random // First lecture

$a = \langle e \in s \mid e < p \rangle$;

if $|a| \geq k$ then return select(a, k);

$b = \langle e \in s \mid e = p \rangle$;

if $|a| + |b| \geq k$ then return p ; // this is the element

$c = \langle e \in s \mid e > p \rangle$;

return select($c, k - |a| - |b|$);

Call a pivot "good" if neither $|a|$ nor $|c|$ is larger than $\frac{2n}{3}$

$\delta = \mathbb{P}(\text{pivot is "good"}) \geq \frac{1}{3}$ because ~~the pivot is chosen among the ones having rank~~

it must be chosen among the ones having rank $[\frac{n}{3}, \frac{2n}{3}]$.

Let $T(n)$ be the average execution time over n elements:

$$T(n) \leq cn + \delta T\left(\frac{2n}{3}\right) + (1-\delta)T(n)$$

→ Solving by $T(n)$, you get $\leq 9cn$

■ Question: How do you select the smallest k -th element via one pass and $O(k)$ additional space? Heap

Sample Sort: Quicksort for hierarchical memories and multi-processor CPUs

It has the same performance guarantees of multi-way mergesort but easier to code and good for multi-disks, multi-CPU's, and for strings.

Mostly IN-PLACE w/ mergesort

Key idea: Instead of a single pivot, we have $(K-1)$ splitters s_1, \dots, s_{K-1} which allow to break the original sequence in K pieces (buckets)

$$B_i = \langle e \in S \mid s_{i-1} \leq e < s_i \rangle$$

where $s_0 = -\infty, s_K = +\infty$
also assume that all keys are different

Goal select splitters so that groups are "balanced"

\Rightarrow Random sample of $(a+1)K-1$ elements
sort it internally (using your favourite order, few elems)
take $s_i = \text{Sample}[(a+1)i]$

Proof: Samples are equally spaced, ~~and~~ and "a" is an oversampling factor which increases the probability that the splitting process was successful.

$a=0 \rightarrow$ just the median elems \rightarrow poor partition w.h.p.

$a \approx \frac{n}{K} \rightarrow$ exact splitters \rightarrow too costly

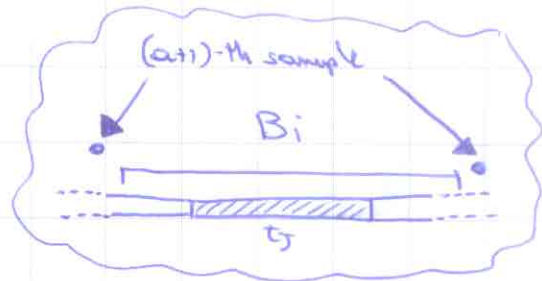
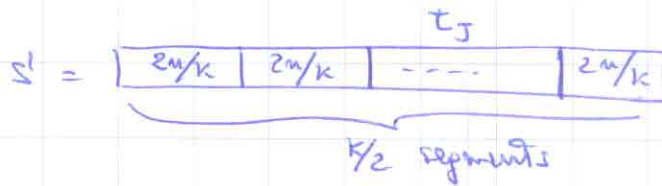
Our following analysis show that $a = \Theta(\log K)$ is a good choice.

Taking $\approx K \log K$ samples, buckets sizes are $\approx \frac{n}{K}$ (avg size).
Actually $\forall i \mid B_i \mid < \frac{4n}{K}$ with probability $\geq 1/2$ [few exceptions are enough to perturb it]

Let $s' = \langle e'_1, e'_2, \dots, e'_n \rangle$ sorted sequence

By contradiction, assume that $\exists B_i$ s.t. $|B_i| \geq \frac{4n}{k}$.

Recall that we are oversampling, so ~~the~~ the condition above may occur even if some sampled element lies in B_i - so we take another "street".



If $|B_i| \geq \frac{4n}{k} \rightarrow$ one segment ~~is~~ is fully included in B_i (so that $\mathbb{P}(|B_i| \geq \frac{4n}{k}) \leq \dots$)

that segment ~~does not~~ does not include a $(a+1)$ th ranked sample [by definition of B_i]

less than $(a+1)$ samples are taken from that segment

$$\mathbb{P}(\text{sample in } t_j) = \frac{2n/k}{n} = \frac{2}{k}$$

$$\mathbb{E}[\text{samples in } t_j] = \left[(a+1)k - 1 \right] \frac{2}{k} \geq \frac{3(a+1)}{2}$$

$$\mathbb{P}(\# \text{samples} < a+1) \leq \mathbb{P}(\# \text{samples} < (1 - \frac{1}{3}) \mathbb{E}[\# \text{samples}])$$

\uparrow Chernoff bound $\quad \quad \quad \mathbb{P}(X < (1-\delta)\mu) \leq e^{-\delta^2 \mu / 2}$

$$\leq e^{-\frac{1}{9} \frac{\mathbb{E}[X]}{2}} = 1/k$$

Applying the union bound over the $k/2$ -segments, we find the probability that "at least one fails, and is larger than $\frac{4n}{k}$ "

$$\frac{1}{k} \cdot \frac{k}{2} = 1/2$$

SMALL, few thousands

Picking $k = \Theta\left(\min\left(\frac{n}{M}, \frac{M}{B}\right)\right)$ we obtain the optimal I/O-complexity over 1-disk. $\frac{M}{B} \approx \frac{10^9}{10^3} = 10^6 \Rightarrow \frac{n}{M} \Rightarrow \frac{n}{10^6}$

Notice that one merge-pass is enough if $n < \frac{M^2}{B} \approx \frac{10^{18}}{10^3} = 10^{15}$
 10^3 TBs

PROBLEM: How to distribute elements in buckets very just
 another array, plus small additional memory.
 // We do not know the bucket sizes, indeed.

1st pass \Rightarrow compute bucket sizes
 2nd pass \rightarrow distribute

No in-place
 it may become short

1st pass is interesting:



load down as an heap

```
J = 2J ; if (>=) J++;
```

Predicated instruction has an additional predicate register as input and it is exec iff the boolean value in this reg is 1

~~predicated~~
~~on some~~
~~CPUs~~
~~(Intel Pentium)~~

No branch misprediction, because they do not affect the instr. flow.

Multiway-mergesort

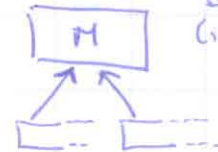
The most famous (word-case)

Binary mergesort \rightarrow 2 ~~problems~~

① starts merging short runs

independent

because memory is large
 after few pages random
 seek \equiv page here
 (if bulked)



② When runs are very long, no help in fetching them
 at once in internal memory

① \Rightarrow form run of M items in one shot

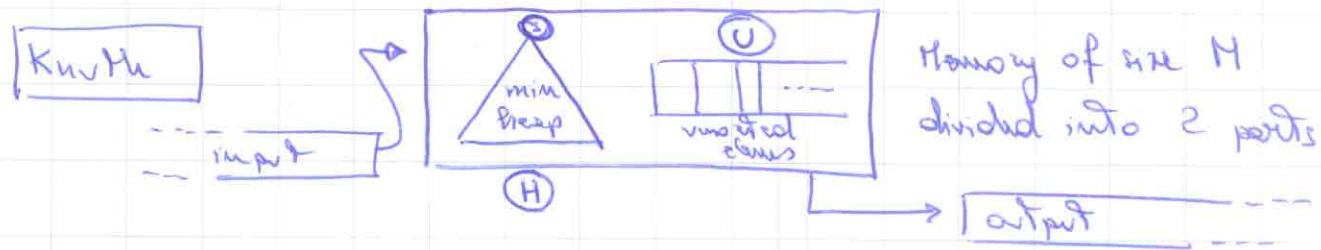
② \rightarrow multi-way merge: $k = \Theta(M/B)$

$$\#I/Os = O\left(\frac{M}{B} \cdot \log_{\frac{M}{B}} \frac{M}{M}\right)$$

• clearly if first-round runs are longer (say $2M$) we may further reduce the tree depth.

Snow Plow by Knuth

• We can speed up scanning via comprehension.



• \forall step: $\text{min-heap} \Rightarrow \text{output}$ (say s)

$t = \text{next input element}$

if $(t > s) \rightarrow \text{add } t \text{ to min-heap}$ (run is still increasing)

else add t to unsorted elements

Prop. ① $H + U = M$ obvious

② time ϕ : $H_0 = M, U_0 = 0$

③ time t : $H_t = 0, U_t = M$ restart

\hookrightarrow revolution time.

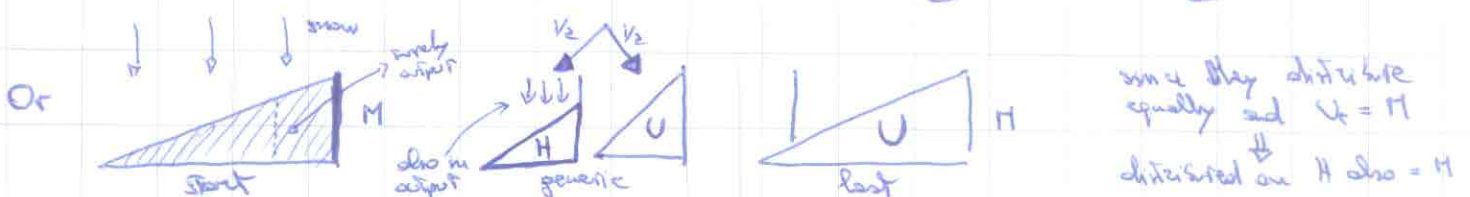
At time t we have output

① all H_0 , hence M elements

② $\frac{t}{2}$ elements, i.e. the ones that landed in H
(assume uniform distn.)

since $\frac{t}{2}$ landed in U too $\Rightarrow \frac{t}{2} \leq M \Rightarrow t \leq 2M$

Hence all output elements are $H_0 + \frac{t}{2} \approx M + \frac{2M}{2} = 2M$



PERMUTING

What about moving around obj?

sorting = sorted order comp + permuting

$$\text{Perm} < \text{Sort in RAM} \iff N \leq N \log N$$

- this means that the bigger cost here is in determining the sorted order.

On disk, we could either permute as in RAM $\approx \Theta(N)$ I/Os
or forget the given permutation and sort.

If $N > \frac{N}{B} \log_{\frac{M}{B}} \frac{N}{M}$ we have that sorting from scratch is advantageous.

\implies On disk $\text{Perm} \approx_{\text{may be}} \text{Sort}$

- this means that on disk the bigger cost is on the movement of the obj (what it is called) I/O-bottleneck

LOWER BOUND

case, and then

We derive the lower bound for permuting in 1-disk case, and then easily extend it to D-disks by dividing by D.

Simple I/O = item does exist either on disk or in memory (no duplicates are created). I/Os can be made simple, cancelling the ~~item~~ ^{number} if it is not needed, or moving it simple if needed.

Goal: Bound the number $\sqrt{\tau}$ of I/Os needed to generate potentially $N!$ permutations of N obj.

Actually, we distinguish between $t_I + t_o = t$ I/Os (simple)

Let $P_i = \#$ permutations potentially generated by i I/Os

① $P_0 = 1;$

② P_{i+1} can be expressed from P_i depending on the I/O-op namely if it is INPUT or OUTPUT I/O.

NOTES

- items not necessarily contiguous (some NIL) \leftarrow int mem disk

OUTPUT

$\leq \frac{N}{B} + 0 - 1$ \leftarrow non-empty blocks before O-M generation

$P_{i+1} = P_i \times \left(\frac{N}{B} + 0 \right) \leq (N \log_2 N) P_i$
 # places the output block can go
 ↑
 brutally since $O(N \log_2 N)$

INPUT

Two cases:

- 1- items block with the output of a previous I/O together in memory before \Rightarrow fetched from $\frac{N}{B} + 0$ blocks
- 2- block never fetched in memory before
 \hookrightarrow items never together in memory $\Rightarrow \frac{N}{B}$ blocks

② has a factor $b!$ more

① $\binom{M+b}{b} \leq \binom{M}{b}$ way of ~~many~~ many fetched items with items in memory

$P_{i+1} = P_i \times N \log_2 N \times \binom{M}{B}$ (case 1)
 $= P_i \times \frac{N}{B} \times \binom{M}{B}$ (case 2)

Actually, the papers are not necessarily always full, so it may be the case that less than B items are fetched in an I/O $\Rightarrow \binom{M}{b_i}$

Calculations are more sophisticated but the main issue is that

$(N \log_2 N)^t \times \binom{M}{B}^t \times (B!)^{N/B} \geq N!$
 $\uparrow \quad \uparrow$
 $t \geq 1 \quad t \geq 0$

$\Rightarrow N \cdot \frac{\log N/B}{B \log M/B + \log N} < N$ if $\log N \geq B \log M/B$ **NEVER**

MULTI-DISKS

parameter to fetch D blocks at every I/O

- simplest approach is to look at all disks as one logical disk with page size $B' = DB$.

\Rightarrow Problem $\log_{M/B}$ \Rightarrow $\log_{M/DB} \frac{N}{M} > \log_{M/B} \frac{N}{M}$

and the better is D the worse is the approach

RANDOMISED DISTRIBUTION

- ① Scan
- ② Distribute (D blocks in output buffer)
- ③ write output buffer

• if scan is executed on a file whose blocks are equally spread among the D disks, then scan takes $\frac{S}{DB}$ I/Os to fetch those S blocks

Considered landing in not out

• Hence ② must ensure this balanced distribution

DIFFICULTY

is that the D blocks to write in ③ may belong to different runs (under distrib formation), and these are arbitrarily allocated on the D disks, and we wish to pay $O(1)$ for ③ so using "striping" on all runs might generate CONFLICTS in WRITE.

because if we have ≥ 2 blocks in the run.

SOLUTION

let N large, $\bullet = \bullet = \bullet$ so every run under formation contains at least $\Omega(D \log D)$ blocks

~~some can differ $\Omega(D)$~~
~~blocks per other disk~~
~~writing to out buffer~~

- Each I/O-write reads its D blocks according to a random perm of $\{1, \dots, D\}$
- Every run distributes randomly ~~its~~ blocks over D disks if $b = \Omega(D \log D) \rightarrow O(b/D)$ blocks per disk w.h.p. (aka: occupancy problem)

so we have an EVEN distribution, which is perfect to guarantee full throughput for the next ①

NOTE

$b = \Omega(D \log D)$ hence the assumptions ~~that~~ ~~the~~ ~~pages~~ ~~above~~ since $\# \text{ splitter} = \min \left\{ \frac{M}{DB}, \frac{N}{M} \right\}$ N

There blocks of a run which go to the same disk. (they are of different I/O-writes.) they are randomly landing at that disk, so they occupy problem

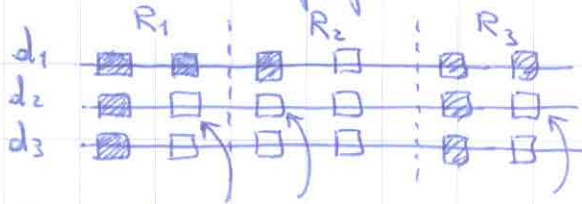
Greedy Sort

No randomization.

SKETCH

- Problem with merge-scheme: conflict in reading the next

D-blocks for merging.



FIG

~~merge-scheme~~

- Prediction sequence \Rightarrow which block to fetch next EASY
(Sort the min item of every block)

- Given the prediction sequence, we need to fetch D blocks at a time in $O(1)$, they could come from the same disk (notice round-robin helps on one run, but not when we have many of them to be merged).

- Again randomization could help

Randomized cycling is nice, used in ~~input~~

Greedy sort avoids every problem.

- Runs are written in a striped way
- Disks operate "independently", by fetching the best 2 blocks per disk of each I/O-read operation



smallest MIN $\rightarrow m_1$
smallest MAX $\rightarrow m_2$

- These 2 blocks are merged $\begin{cases} \text{smallest B items goes to output} \\ \text{largest B items goes to } m_1\text{-run} \\ \text{since } \text{max of block of } m_1 > m_2 \end{cases}$
- This is an "approximate" merge, provable that every item is $D\sqrt{MB}$ portions from its correct one.
- Typically $D\sqrt{MB} < M$, otherwise use ~~sort~~ Column sort to fix.