# Lec 10 | Sorting and Aggregation
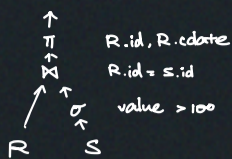
Preview          Query plan              Algorithm
                                         - Maximise sequential IO
                    ↑                    - Spill to disk if necessary
                    π    R.id, R.cdate   - Sort when ORDER BY
                    ↑                       ... DISTINCT, GROUP BY
                    ⋈    R.id = S.id
                  ↑   ↑
                  σ    σ    value > 100
                  ↑    ↑
                  R    S

## # In - memory   Sorting

→ Just   run sorting algorithm
  └ Quicksort , TimSort , Insertion Sort
        └ hybrid merge sort & insertion sort
  └ Perf depends on data distribution
    └ nearly sorted — simple ones may be faster

## # Top - n  heap  sort

If only want  top - n elems, iteratively put things in
sorted size-n heap ( discard things when appropriate when
heap is  full )
viz. keep  top - n so far while scanning input

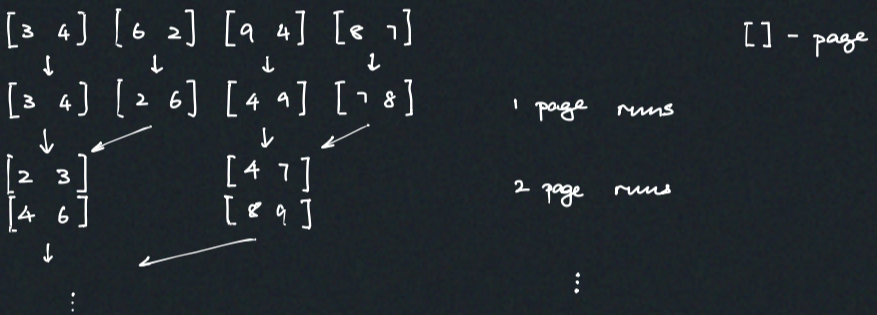## # External  merge  sort — sort things larger then mem

D&C into  multiple smaller runs , sort them, then
combine
→ Keep keys  & point to data , sort, then chase pointer
→ Keep keys  & data, then sort     (early materialisation)

▷ 2 way  external merge sort

  Sort each run, write to new tmp file, keep merging

  [3 4] [6 2] [9 4] [8 7]                    [] - page
   ↓     ↓     ↓     ↓
  [3 4] [2 6] [4 9] [7 8]        1 page  runs
   ↓      ↘    ↓    ↙
  [2 3]       [4 7]              2 page  runs
  [4 6]       [8 9]
   ↓       ↙
   ⋮                                    ⋮

  Whole algorithm    need 3 pages in buf pool
  Num  passes        1 + ⌈log₂N⌉
  IO cost            2N · num passes

  Generalise    - increase fanout , multi-way merge
  Optimise                   └ use heap to find min
                - go straight to higher page run
                                        ↗ step func! 🆒
                → Num passes  1 + ⌈log_{B-1} ⌈$\frac{N}{B}$⌉⌉    B = buffer size
                  IO cost     2N · num passes

                - Double buffering : prefetch data for next run in
                  background while merging. Requires double buff
                  pool size — should use if halving B doesn't bump
                  num passes

                - Code opt. : inline the comparison func

                - String sorting :  try use prefix

## # BTree   Sorting
  └ If clustered ( leaves in physical order ) , then good idea
    to  run  sorting on those pages . Else bad idea

## # Aggregation

▷ DISTINCT

  → filter , get column , sort , scan & remove duplicate
    └ still  sorted after  dedup
  → external  hashing
    1. Partition — two indep hash funcs h₁, h₂
       filter , get column, partition by  h₁ into buckets
    2. Rehashing  ( hopefully hash table fit in memory )
       rehash each bucket into  hash table via h₂ ( on RAM
       or with disk hash table )

▷ GROUP BY .   calc AVG

  → keep running aggregation in mem  as hash table .
    └ for AVG, keep running sum & count       └ hopefully in mem
        ∴ for COUNT, MIN, MAX, SUM, etc.