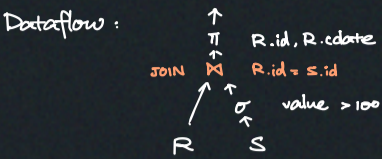Put two tables together — often the most expensive operation
Most typical : join primary & foreign key

# Inner equijoin      ⋈

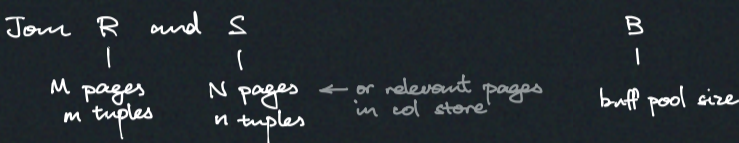Optimisation :   try to make smaller table the outer table

Dataflow :

$$
\begin{array}{c}
\uparrow \\
\pi \quad R.id, R.cdate \\
\text{JOIN} \quad \bowtie \quad R.id = s.id \\
\sigma \quad value > 100 \\
R \quad\quad S
\end{array}
$$

Semantic :    for row r in R :                         } Equivalently, do R×S,
                    for row s in S :                     } then filter
                        if R.id = s.id :
                            concat r and s, add to output

Efficiency :   depends on col vs row store and OLAP vs OLTP
                              └─ usually more IO
                                  └─ usually more work to track which records
                                      e.g. track ids, use bitmaps, ...

                    tradeoff btwn early vs late materialisation
                       └ can hybrid       └ hard to refer back to ids/positions
                                        easier to code

Cost analysis

     Join R and S                                    B
           |          |                              |
      M pages    N pages  ← or relevant pages    buff pool size
      m tuples   n tuples   in col store

# Algorithms

▷ Naïve  for loop                    M + (mN)
                                     or N + (nM)
                "outer table"        by swapping
     for  r ∈ R :                    loop order
        for s ∈ S :
           if r,s match , emit

▷ Block nested loop                  M + (MN)
         └ block of multiple (or one) pages    if 3 buff poll pages

     for B_R in R :                  sequential flooding! if S
        for B_S in S :               has more blocks than buff pool size
           for r ∈ B_R :             we can't keep them all pinned for
              for s ∈ B_S :          the whole scan
                 if r,s match , emit    if more pages available
                                     pin more pages in R in a time
     Multi-page block                reduce S scan count
            ↙ page count
     for B_R^{B-2} in R :  ← try use smaller table   → Use B-2 buff for outer table,
        for B_S' in S :          as outer               1 buf to scan inner, and 1 to
           for r ∈ B_R^{B-2} :                          write output
              for s ∈ B_S' :                         ↳ M + (⌈M/(B-2)⌉ N)
                 if r,s match , emit
                                     if inner table fits in B-2 bufs,
                                     just keep it in mem works too

▷ Index nested loop                  M + (m C)
                                             ↑
     for r ∈ R :                       index lookup cost
        lookup r in index of S

▷ Sort - Merge Join

     Idea :  sort both independently and traverse both with two cursors
             keep track of where mini scans for jumping back if
             hitting duplicate

     Optimisation :  combine the merge & traverse steps

     Worse case :  many duplicates — need to go back all the time

     Good case :  already sorted tables

▷ Hash Join

     1. Build hash table for small table use hash func h.
              └ value depends on early vs late materialisation
     2. Probe the hash table when traversing other table

     Worse case :  all duplicate

     Optimisation :  bloom filter on R. Probe bloom filter before hash table
                     lookup

     Try fit hash table in mem. Otherwise partition ...

▷ GRACE hash join

     1. Hash R into k buckets with h₁   ] write to disk if necessary
     2.  "    S  "  k  "   "    h₁   ]
     3. Join each partition pair
     4. Merge the results

     Corner cases :  everything go into same partition ...   some partition
                     too big ...                             too big ...
                        → recursively partition by another hash func

     Cost .   3 ( M + N )
                └ read, write partitions , read partitions

     Optimisation :  hybrid hash join
                     if only need two partitions, keep 0ᵗʰ part. in mem
                     etc.

     Observations
     - The probe-side table can be any size