# Hash Index, Equality Selection

cost will ne cost to find correct bucket.
+ # of I/Os to return correct tuples from that
bucket (might be in several pages)

So $\log 2$ I/Os to find bucket

# General Selects

## Selection w/o Disjunction

Could
① do file scan to retrieve tuples
or use an index that matches some
of conjuncts to scan records then
cut from output according to non-primary
conjuncts

② try to use several indexes

Idea   sort rids satisfying conditions
represent using bit map
and maps

## Selection w/ Disjunction

Consider
$$\sigma_{(R)}$$
$$\underbrace{day < 8/9/02}_{\text{no index}} \vee \underbrace{rname = 'Joe'}_{\text{index}} = \varphi$$

Because of this need to do file scan so no pt use index on
rname.

However if have $\sigma_{\varphi \wedge sid=3}(R)$
$\underbrace{\quad}_{\text{index}}$

Would use sid index to first winnow results
down. Then scan to calculate $\varphi$.

If had index on both day & name then could
use both indexes + bitmap strategy.

# Projections Based on Sorting

Consider:

```
SELECT DISTINCT R.sid, R.bid
FROM Reserves R
```

## Using Sorting:

① Scan R and produce a set of tuples that contain only the desired attributes

② Sort this set

③ Scan the sorted result, comparing adjacent tuples and eliminate duplicates

Cost — Scan Reserves 1000 I/Os. If result tuples 10 nytes then tmp file write is 250 I/Os. Sort step 1000 I/Os. So 2500 I/Os

Scan sort set ↑ 250

# Projections Based on Hashing

Suppose # of buffers pages $B$ is large compared to # of pages in R.

Have one input buffer page $B-1$ out buffer

Read R into input project out unwanted attributes. Apply a hash $f^{\underline{1}}$ to remainder to pick one of the $B-1$ output buffer

~~When an out buffer fills write to it~~

When an out buffer fills write to an out file

How does this help? Each out file partition contains different tuples. Within a partition tuples may or may not be different. So we read in each partition now and apply a second different hash $f^2$. Provide $B > \sqrt{T}$ can hash whole partition into existing pages in buffer. Then sort these pages & output.

What is cost Read R + 2 Read tmp = 1500 I/Os

$\left(\frac{T}{B-1}\right)$ of ↑ pages in partition

# Join Operation

Consider   SELECT *
           FROM Reserves R, Sailors S
           WHERE  R.sid = S.sid

simple join →

## Nested Loop Join

foreach page $P_R$ of R
  foreach page $P_S$ of S
    for all matching r, s in
      $P_R$ and $P_S$ add <r,s> to output

Cost if R is M pages
and S is N pages
then get $M + M \cdot N$

In running
example   $1000 + 1000 \cdot 500 = 501,000$ I/s

At 10 mS/I/O would take 1.4 hours.

● Note   $R \bowtie B = B \bowtie R$
So what if used S in outer loop?
then get  $N + M \cdot N = 500 + 1000 \cdot 500$
                         $= 500,500$ I/Os

So smaller relation should be in outer
loop.