

# Relational Query Optimization I: The Plan Space

Alvin Cheung

Fall 2022

Reading: R&G 15



# Architecture of a DBMS



- Completed →
- We are here →
- Completed →
- Completed →
- Completed →
- Completed →



# Query Optimization is Magic



- The bridge between a *declarative* domain-specific language...
  - “What” you want as an answer
- ... and custom *imperative* computer programs
  - “How” to compute the answer
- A lot of smart people and a lot of time has been spent on this problem!
- Reminiscent of many cutting-edge “AI” problems
  - Similar tricks: optimization + heuristic pruning
  - Robot path planning, natural language understanding, etc

# Invented in 1979 by Pat Selinger et al.



- We'll focus on “System R” (“Selinger”) optimizers
  - From IBM Research in Almaden
- “Cascades” optimizer is the other common one
  - Later, with notable differences, but similar big picture

The image shows the front cover of a technical paper. At the top, the title "Access Path Selection in a Relational Database Management System" is printed in a serif font. Below the title, the authors' names are listed: P. Griffiths Selinger, M. M. Astrahan, D. D. Chamberlin, R. A. Lorie, and T. G. Price. Further down, the location "IBM Research Division, San Jose, California 95193" is mentioned. To the right of the text, there is a black and white portrait photograph of a woman with long, light-colored hair, smiling. The entire document is presented within a rectangular frame.

Access Path Selection  
in a Relational Database Management System

P. Griffiths Selinger  
M. M. Astrahan  
D. D. Chamberlin  
R. A. Lorie  
T. G. Price

IBM Research Division, San Jose, California 95193

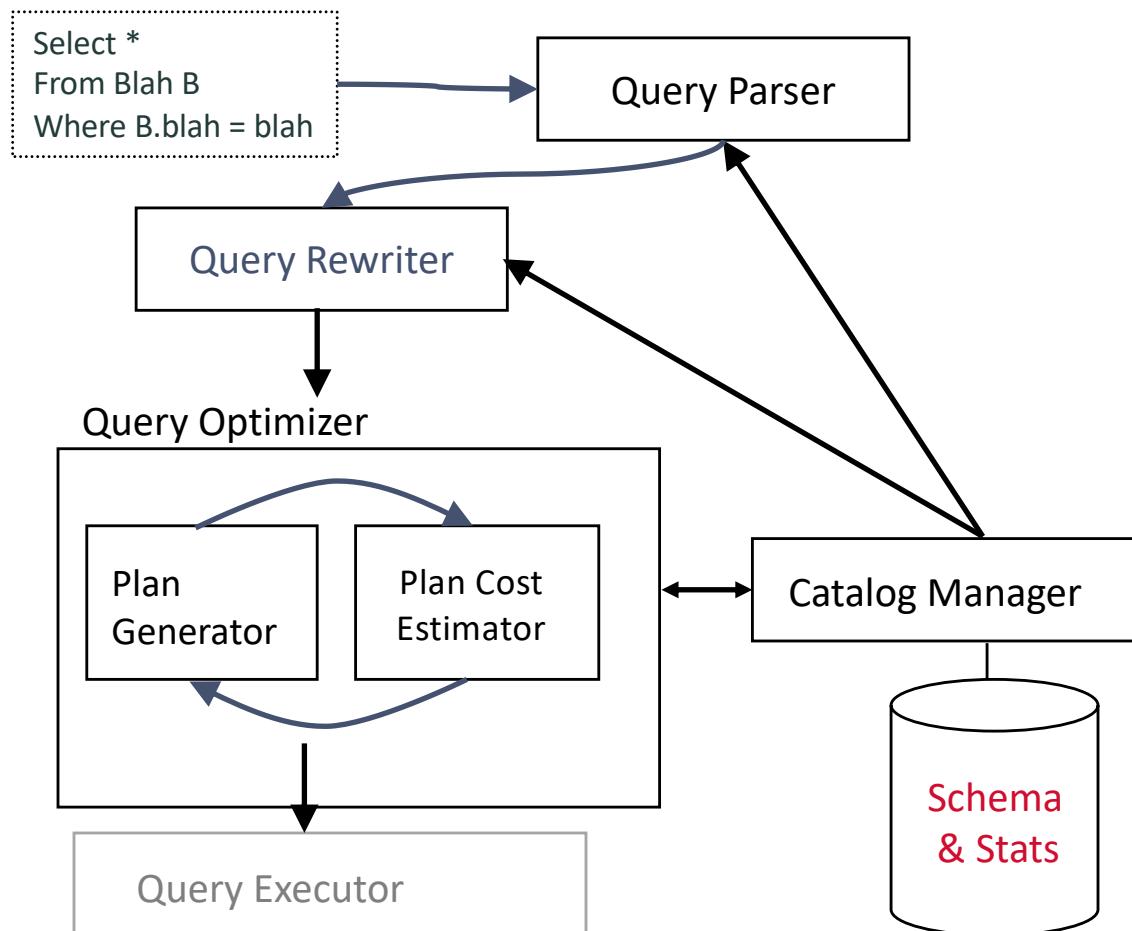
ABSTRACT: In a high level query and data manipulation language such as SQL, requests

retrieval. Nor does a user specify in what order joins are to be performed. The

# Query Parsing & Optimization: Query Lifecycle



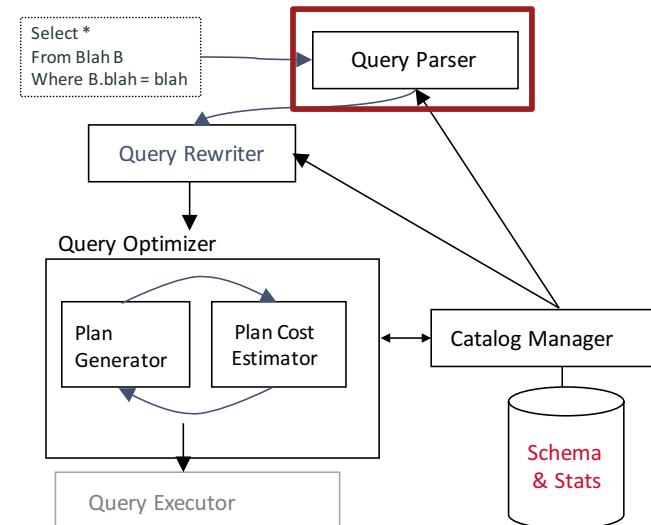
IMPORTANT



# Query Parsing & Optimization Part 2



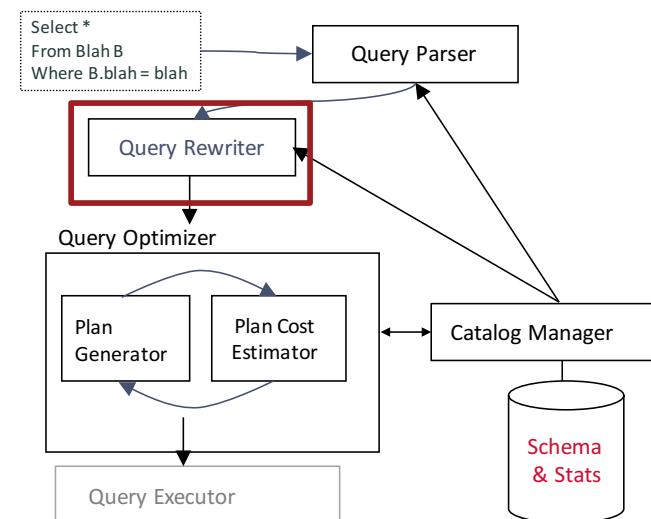
- **Query parser** *Component 1*
  - Checks correctness, authorization
  - Generates a parse tree
  - Straightforward
  - Not our focus



# Query Parsing & Optimization Part 3



- **Query rewriter** | *Component 2*
  - Converts queries to canonical form
    - flatten views
    - subqueries into fewer query blocks
      - e.g., by replacing w/ joins
  - Not our focus



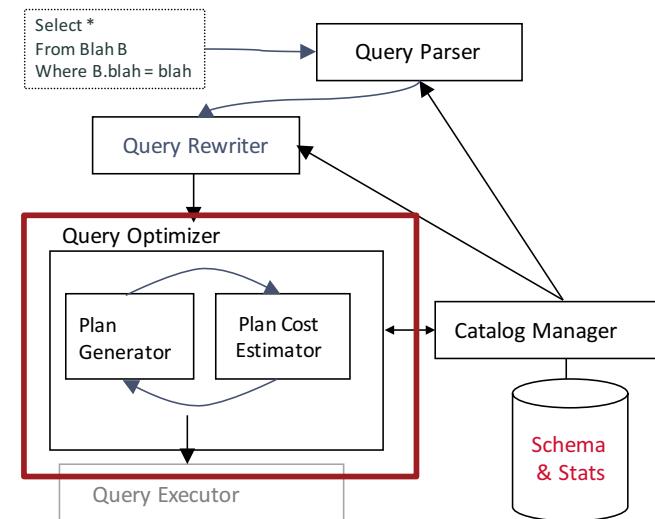
# Query Parsing & Optimization Part 4



*Component 3*

- “Cost-based” Query Optimizer
  - Our focus!
  - Optimizes 1 query block at a time
    - Select, Project, Join
    - GroupBy/Agg
    - Order By (if top-most block)
  - Uses catalog stats to find least-“cost” plan per query block
  - Often not truly “optimal”, lots of heuristic rules and magic

*of the relations*



# Query Optimization: The Components



- Three (mostly) orthogonal concerns:
    - Plan space:
      - for a given query, what plans are considered?
      - larger the plan space, more likely to find a cheaper plan, but harder to search
    - Cost estimation:
      - how is the cost of a plan estimated?
      - want to find the cheapest plan
    - Search strategy:
      - how do we “search” in the “plan space”?

# Query Optimization: The Goal



- Optimization goal:
  - Ideally: Find the plan with least actual cost = one that runs fastest
  - Reality: Find the plan with least estimated cost.
    - At the very least, try to avoid really bad actual plans!

Kinda like Dijkstra-greedy

# Today



- We will get a feel for the plan space
- Explore one simple example query

# Plan Space



- To generate a space of candidate plans, we need to think about how to rewrite relational algebra expressions into other ones
- Therefore, need a set of equivalence rules

# Relational Algebra Equivalences: Selections



$c_1 \text{ and } c_2 \text{ and } \dots$

- Selections:

- $\sigma_{c_1 \wedge \dots \wedge c_n}(R) \equiv \sigma_{c_1}(\dots(\sigma_{c_n}(R))\dots)$  (cascade)  $\sigma$  property 1
  - Intuitively, RHS says check  $c_n$  first on all tuples, then  $c_{n-1}$  etc.
- $\sigma_{c_1}(\sigma_{c_2}(R)) \equiv \sigma_{c_2}(\sigma_{c_1}(R))$  (commute)  $\sigma$  property 2

# Relational Algebra Equivalences: Projections



- Selections:
  - $\sigma_{c_1 \wedge \dots \wedge c_n}(R) \equiv \sigma_{c_1}(\dots(\sigma_{c_n}(R))\dots)$  (cascade)
  - $\sigma_{c_1}(\sigma_{c_2}(R)) \equiv \sigma_{c_2}(\sigma_{c_1}(R))$  (commute)
- Projections: *a<sub>1</sub> stays in all expressions*
  - $\pi_{\underline{a_1}}(\dots(R)\dots) \equiv \pi_{\underline{a_1}}(\dots(\pi_{\underline{a_1}, \dots, a_{n-1}}(R))\dots)$  (cascade) *π property*
    - Essentially, allows partial projection earlier in the expression
    - As long as we're keeping a<sub>1</sub> (and everything else we need outside) we're OK
  - Q: Are there any commute rules for projections?

## Relational Algebra Equivalences: Cartesian Product



- Selections:
  - $\sigma_{c_1 \wedge \dots \wedge c_n}(R) \equiv \sigma_{c_1}(\dots(\sigma_{c_n}(R))\dots)$  (cascade)
  - $\sigma_{c_1}(\sigma_{c_2}(R)) \equiv \sigma_{c_2}(\sigma_{c_1}(R))$  (commute)
- Projections:
  - $\pi_{a_1}(\dots(R)\dots) \equiv \pi_{a_1}(\dots(\pi_{a_1, \dots, a_{n-1}}(R))\dots)$  (cascade)
- Cartesian Product
  - $R \times (S \times T) \equiv (R \times S) \times T$  (associative)
  - $R \times S \equiv S \times R$  (commutative)
    - Recall that the ordering of attributes doesn't matter

# Are Joins Associative and Commutative?

Berkeley  
cs186

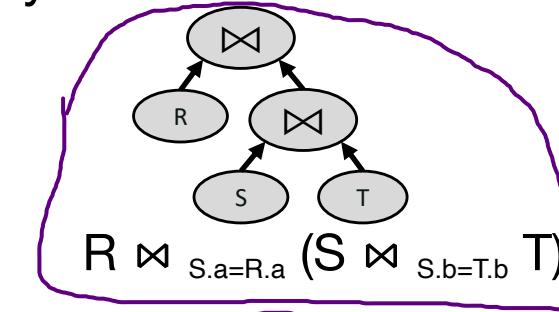
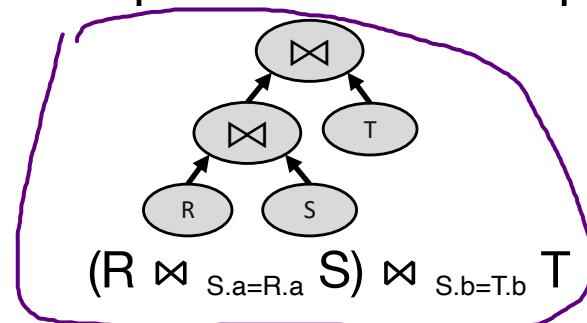
- After all, just Cartesian Products with Selections (*conditions*)
- You can think of them as associative and commutative... "X" vs. " $\bowtie_{S.a=R.b}$ "  
*SOMETIMES*
- But beware of join turning into cross-product!
  - Consider  $R(a,z)$ ,  $S(a,b)$ ,  $T(b,y)$
  - Attempt 1: Does this work? Why?
    - $(S \bowtie_{S.b=T.b} T) \bowtie_{S.a=R.a} R \neq S \bowtie_{S.b=T.b} (T \bowtie_{S.a=R.a} R)$ 
      - not legal!!  $T$  does not have attribute "S.a"
  - Attempt 2: Does this work? Why?
    - $(S \bowtie_{S.b=T.b} T) \bowtie_{S.a=R.a} R \neq S \bowtie_{S.b=T.b} (T \times R)$ 
      - not the same!! Join on  $S.a=R.a$  now gone
  - Attempt 3: Does this work?
    - $(S \bowtie_{S.b=T.b} T) \bowtie_{S.a=R.a} R \equiv S \bowtie_{S.b=T.b \wedge S.a=R.a} (T \times R)$

*Yes, but left side cheaper*

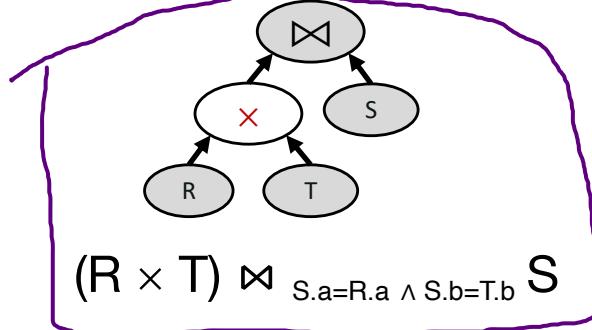
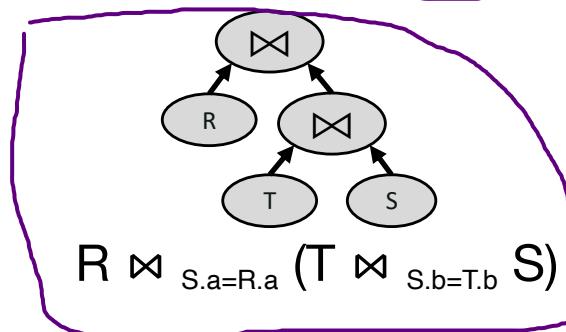
# Join ordering, again



- Similarly, note that some join orders have cross products, some don't
- Equivalent for the query above:



```
SELECT *
  FROM R, S, T
 WHERE R.a = S.a
   AND S.b = T.b;
```



# Plan Space



- To generate a space of candidate plans, we need to think about how to rewrite relational algebra expressions into other ones
- Therefore, need a set of equivalence rules – **done**
- Next, will discuss a **set of heuristics that are used to restrict attention to plans that are mostly better:**
  - we've already seen one of these in the relational alg lectures.

$\sigma \equiv \text{where}$

## Some Common Heuristics: Selections



- Selection cascade and pushdown
  - Apply selections as soon as you have the relevant columns
  - Ex:
    - $\pi_{\text{sname}} (\sigma_{(\text{bid}=100 \wedge \text{rating} > 5)} (\text{Reserves} \bowtie_{\text{Reserves.sid}=\text{Sailors.sid}} \text{Sailors}))$
    - $\pi_{\text{sname}} (\sigma_{\text{bid}=100} (\text{Reserves}) \bowtie_{\text{Reserves.sid}=\text{Sailors.sid}} \sigma_{\text{rating} > 5} (\text{Sailors}))$  ← better
  - Why is this an improvement?
    - Selection is essentially free, joins are expensive
    - Take care of selections early -- side effect is that the intermediate inputs to joins are smaller

$\pi^{= \text{selected}}$



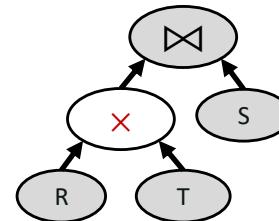
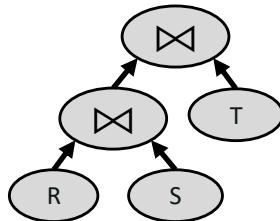
## Some Common Heuristics: Projections

- Projection cascade and pushdown
  - Keep only the columns you need to evaluate downstream operators
  - Reserves(sid, bid, day), Sailors (sid, rating, sname)
  - Ex:
    - $\pi_{\text{sname}} \sigma_{(\text{bid}=100 \wedge \text{rating} > 5)}$  (Reserves  $\bowtie_{\text{Reserves.sid}=\text{Sailors.sid}}$  Sailors)
    - Q: How might we cascade and push projections and selections down?
    - $\pi_{\text{sname}} (\underline{\pi_{\text{sid}}} (\sigma_{\text{bid}=100} (\text{Reserves})) \bowtie_{\text{Reserves.sid}=\text{Sailors.sid}} \underline{\pi_{\text{sname}, \text{sid}}} (\sigma_{\text{rating} > 5} (\text{Sailors})))$
    - Other rewritings exist! (reorder selection and projection) Just make all columns are retained

# Some Common Heuristics



- **Avoid Cartesian products**
  - Given a choice, do theta-joins rather than cross-products
  - Consider  $R(a,b)$ ,  $S(b,c)$ ,  $T(c,d)$
  - Favor  $(R \bowtie S) \bowtie T$  over  $(R \times T) \bowtie S$  (*in general - heuristic*)
  - Case where this doesn't quite improve things:
    - if  $R \times T$  is small (e.g.,  $R$  &  $T$  are very small and  $S$  is relatively large)
    - Still, it's a good enough heuristic that we will use it



# Plan Space



- To generate a space of candidate plans, we need to think about how to rewrite relational algebra expressions into other ones
- Therefore, need a set of equivalence rules – **done**
- Next, will discuss a set of heuristics that are used to restrict attention to plans that are mostly better – **done**
- Both of these were logical equivalences, will also quickly discuss physical equivalences, next.

# Physical Equivalences



- Base table access
  - Heap scan
  - Index scan (if available on referenced columns)
- Equijoins
  - Block (Chunk) Nested Loop: simple, exploits extra memory
  - Index Nested Loop: often good if 1 rel small and the other indexed properly
  - Sort-Merge Join: good with small memory, equal-size tables
  - Grace/Hybrid Hash Join: even better than sort with 1 small table
- Non-Equijoins (*no join conditions*)
  - Block (Chunk) Nested Loop

# Schema for Examples

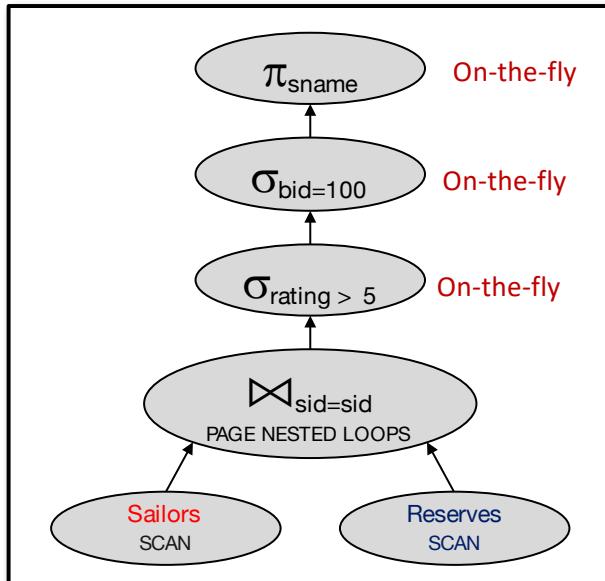


Sailors (*sid*: integer, *sname*: text, *rating*: integer, *age*: real)  
Reserves (*sid*: integer, *bid*: integer, *day*: date, *rname*: text)

- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have  $B = 5$  pages to use for joins
- Remember: just counting IOs

# Motivating Example: Plan 1

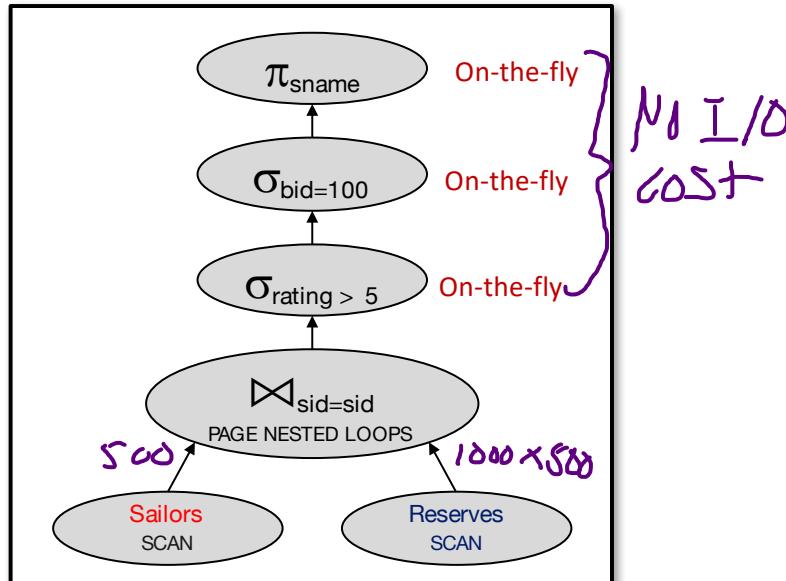
- Here's a reasonable query plan:



```
SELECT S.sname  
FROM Reserves R, Sailors S  
WHERE R.sid=S.sid  
AND R.bid=100  
AND S.rating>5
```

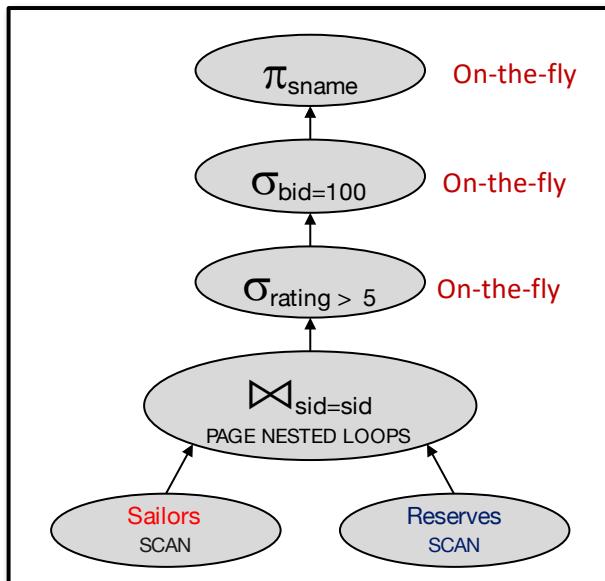
- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have  $B = 5$  pages to use for joins

# Motivating Example: Plan 1 Cost



- Let's estimate the cost:
  - Scan Sailors (500 IOs)
  - For each page of Sailors, Scan Reserves (1000 IOs)
  - Total:  $500 + 500 * 1000$ 
    - 500,500 IOs
- 3 nested join*
- Reserves:
    - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
    - Assume there are 100 boats (each equally likely)
  - Sailors:
    - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
    - Assume there are 10 different ratings (each equally likely)
  - Assume we have  $B = 5$  pages to use for joins

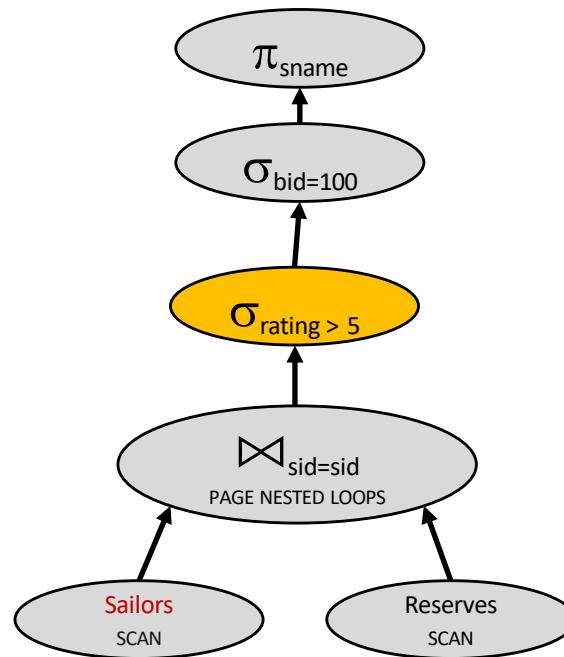
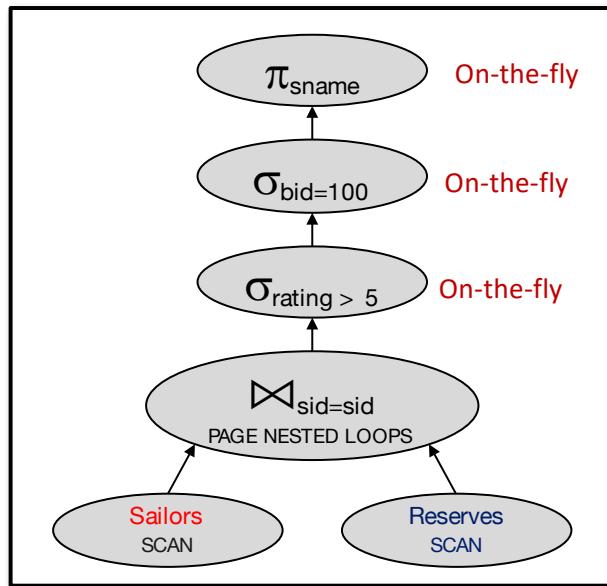
# Motivating Example: Plan 1 Cost Analysis



- Cost:  $500 + 500 * 1000$  I/Os
- By no means the worst plan!
- Misses several opportunities:
  - selections could be ‘pushed’ down
  - no use of indexes
- Goal of optimization:
  - Find faster plans that compute the same answer.

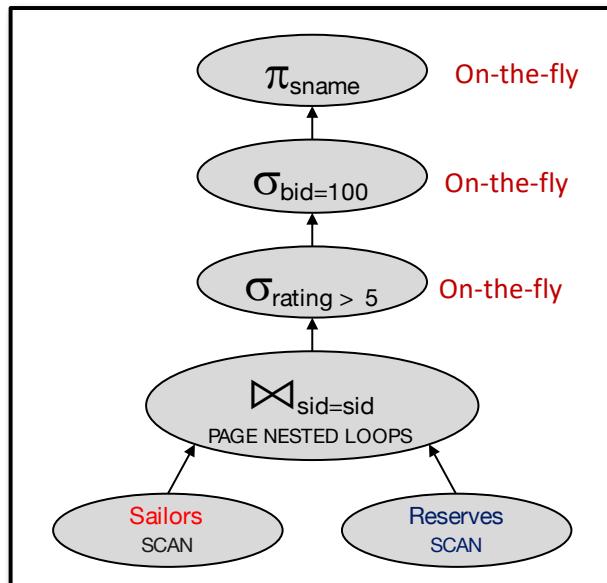
- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have  $B = 5$  pages to use for joins

# Selection Pushdown

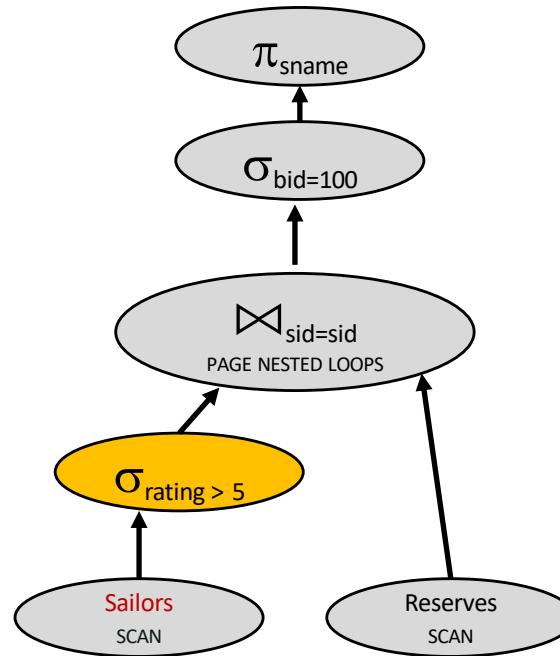


500,500 IOs

# Selection Pushdown, cont



500,500 IOs

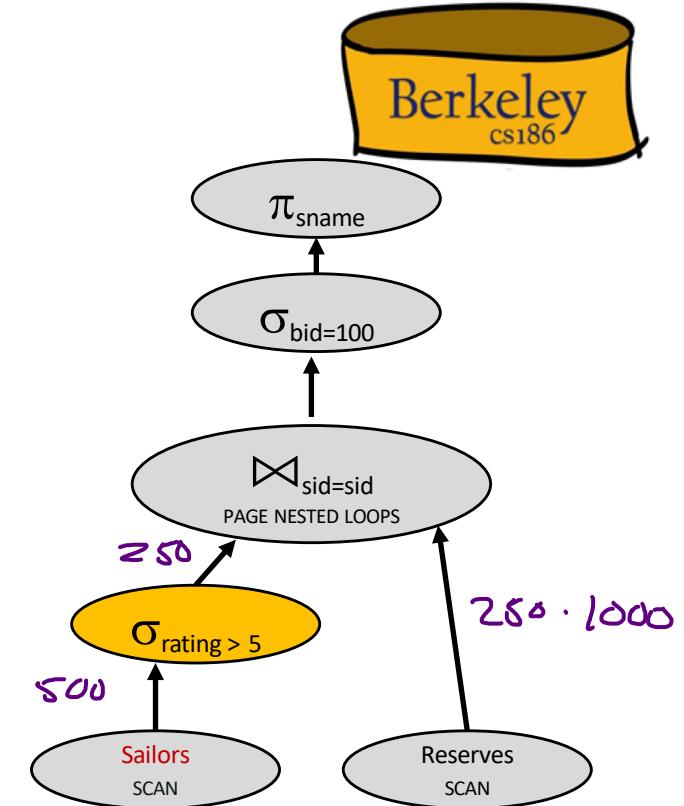


Cost?

# Query Plan 2 Cost

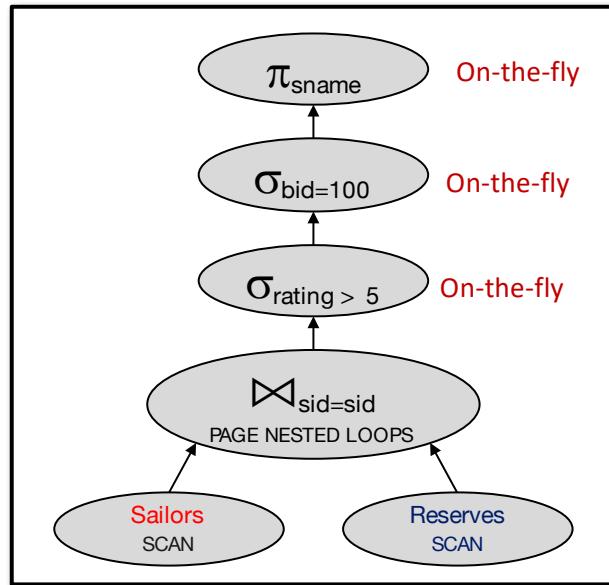
- Let's estimate the cost:
- Scan Sailors (500 IOs)
- For each pageful of high-rated Sailors,  
Scan Reserves (1000 IOs)
- Total:  $500 + ??? * 1000$
- Remember: 10 ratings, all equally likely
- Total:  $500 + \underline{250} * 1000$**

$\text{rating} > 5$   
 is about half the  
 $\text{rating} 5$

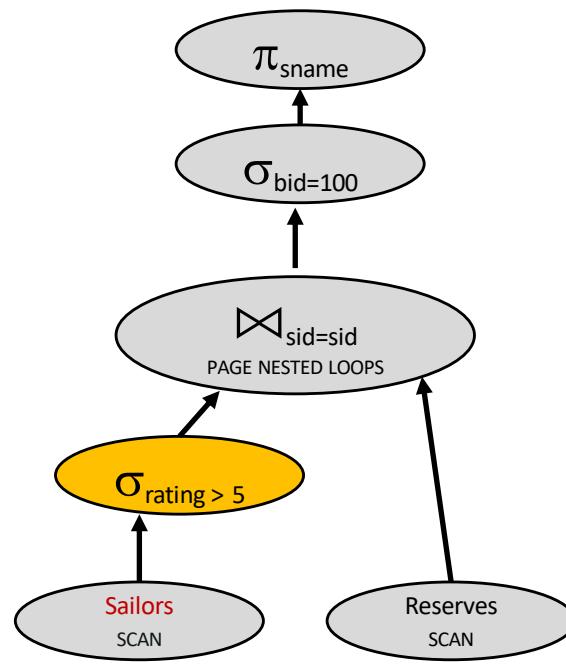


- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have  $B = 5$  pages to use for joins

# Decision?

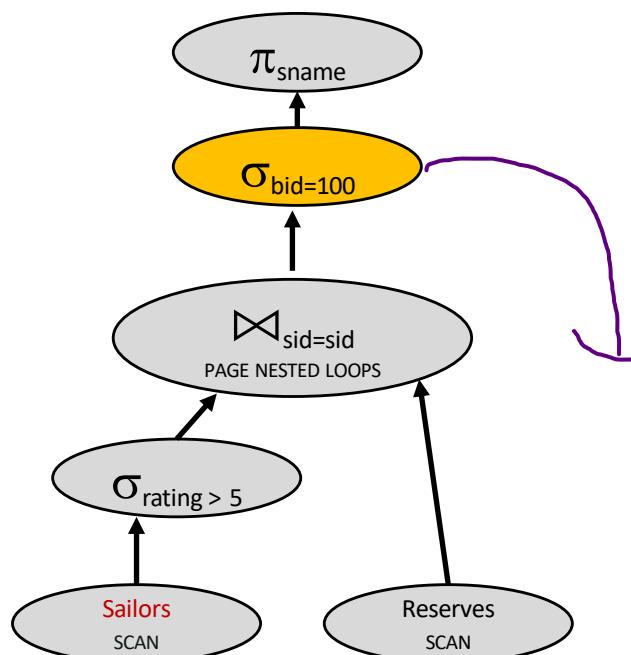
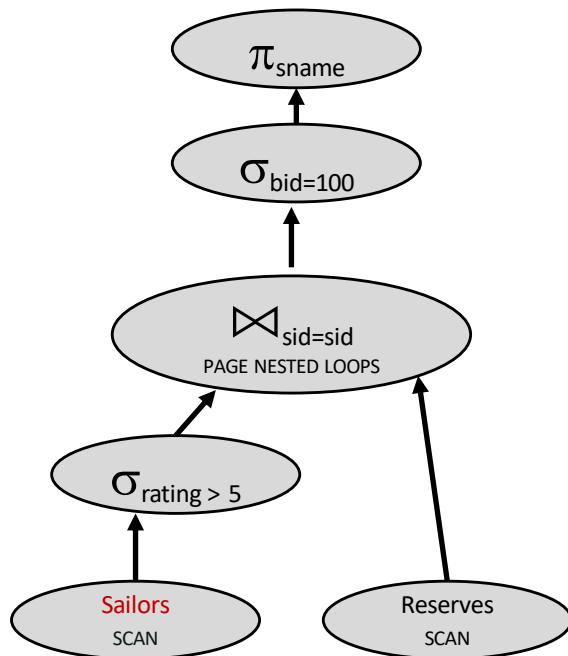


500,500 IOs



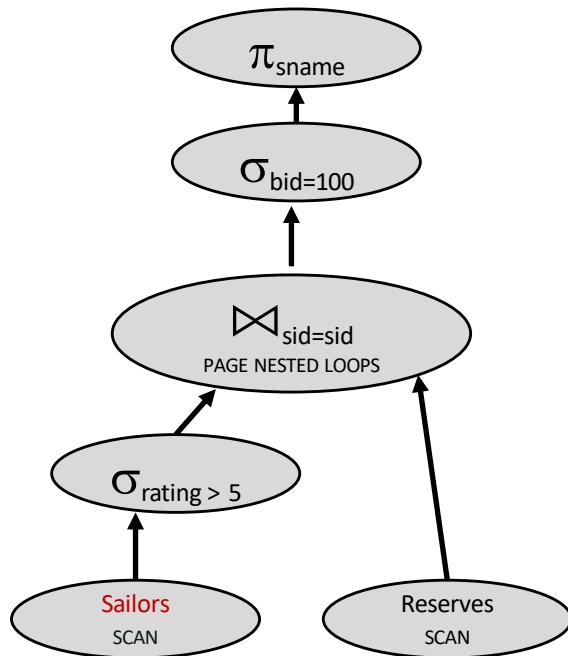
250,500 IOs

# More Selection Pushdown

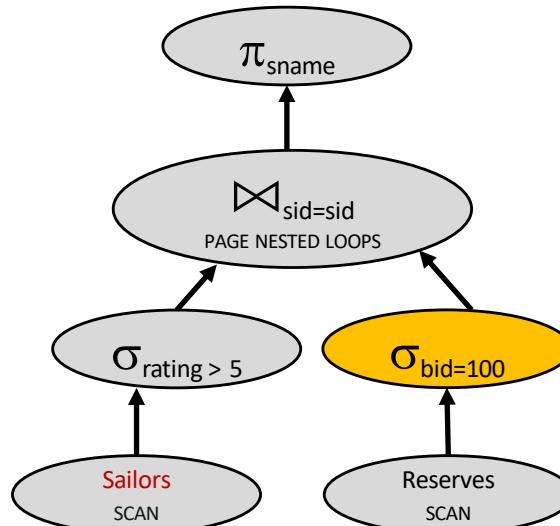


250,500 IOs

# More Selection Pushdown, cont



250,500 IOs

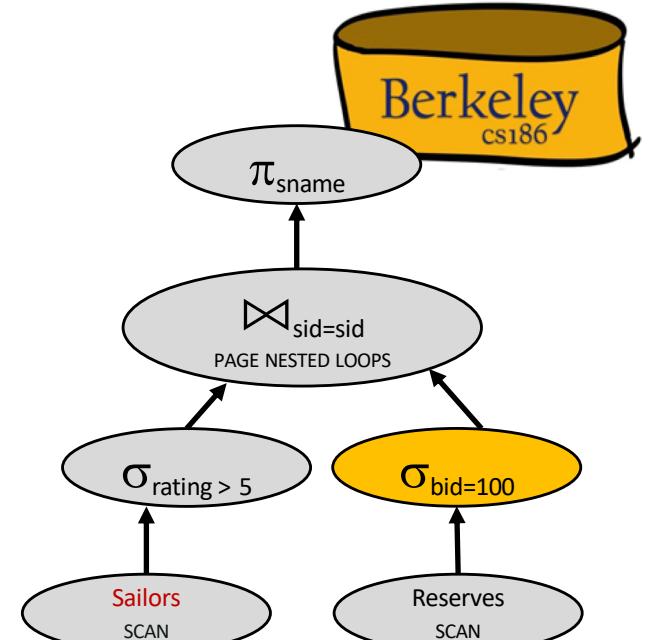


Cost???

# Query Plan 3 Cost Analysis

Let's estimate the cost:

- Scan Sailors (500 IOs)
- For each pageful of high-rated Sailors,  
Read through Reserves tuples that match
- Total:  $500 + 250 * (???)$
- For each scan of Reserves, we apply filter on tuples  
on the fly
- Problem: this doesn't actually save any IOs – to  
determine the Reserves tuples that match, we end  
up scanning Reserves the same # of times.
- Total:  $500 + 250 * 1000!$



Reserves:

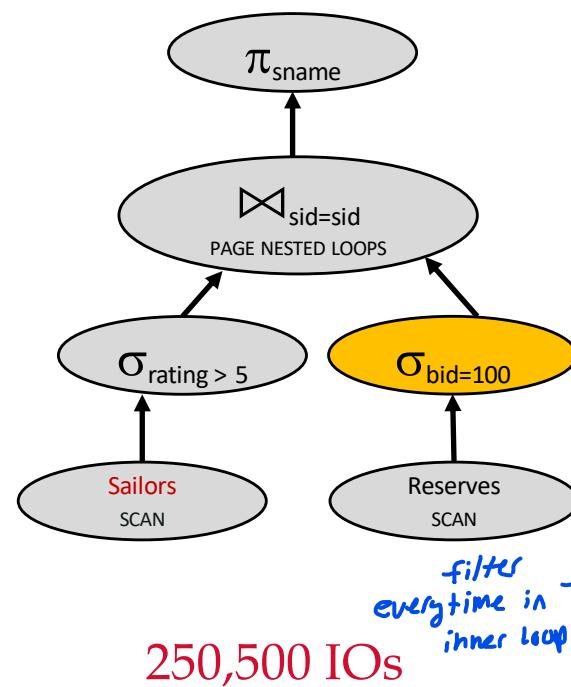
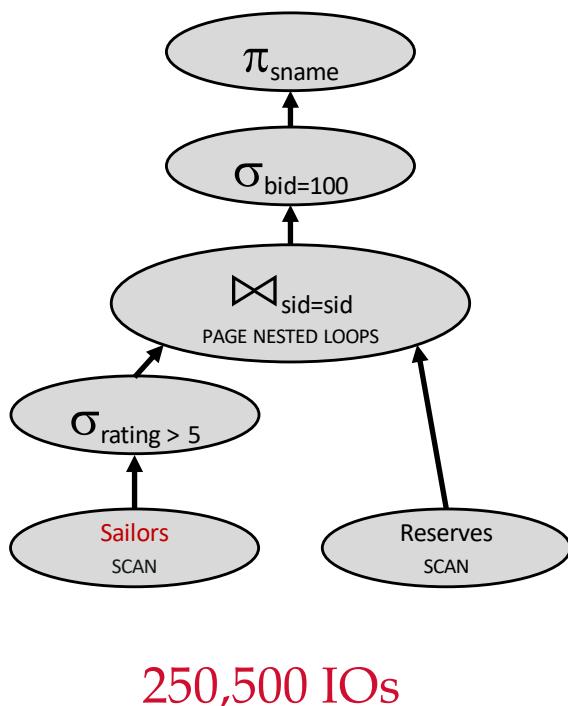
- Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
- Assume there are 100 boats (each equally likely)

Sailors:

- Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
- Assume there are 10 different ratings (each equally likely)

• Assume we have  $B = 5$  pages to use for joins

# More Selection Pushdown Analysis



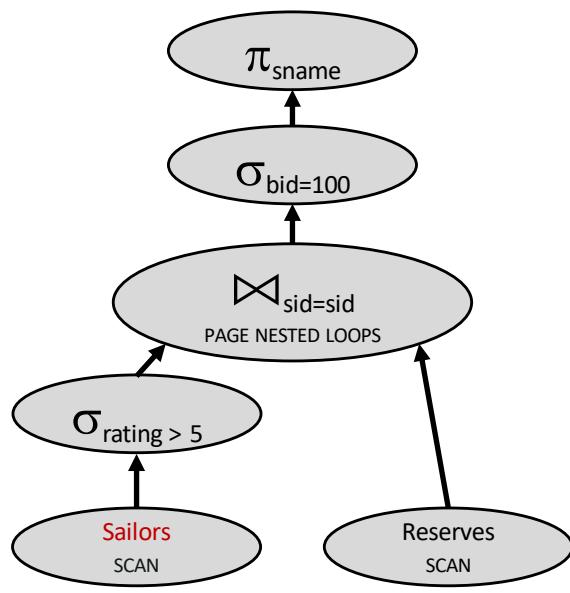
Pushing a selection into the inner loop of a nested loop join doesn't save I/Os! Essentially equivalent to having the selection above.

pseudo code:

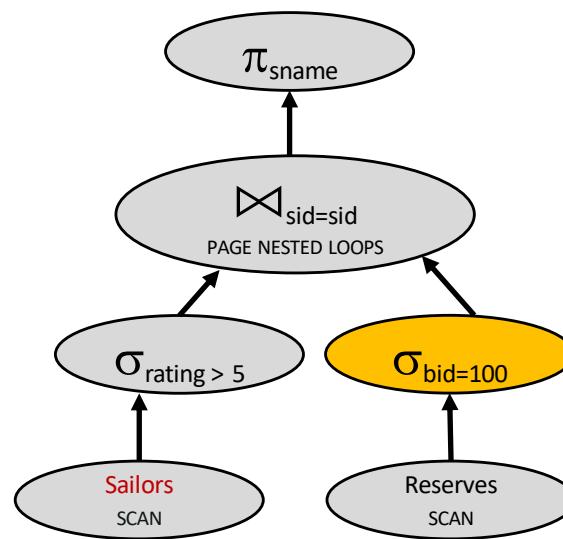
sailors  $\rightarrow$  high rating sailors  
for page in high rating sailors:  
 reserves  $\rightarrow$  bid=100 reserves  
 join (reserve, high rating sailor)

filter  
everytime in  
inner loop

# Decision 2



250,500 IOs



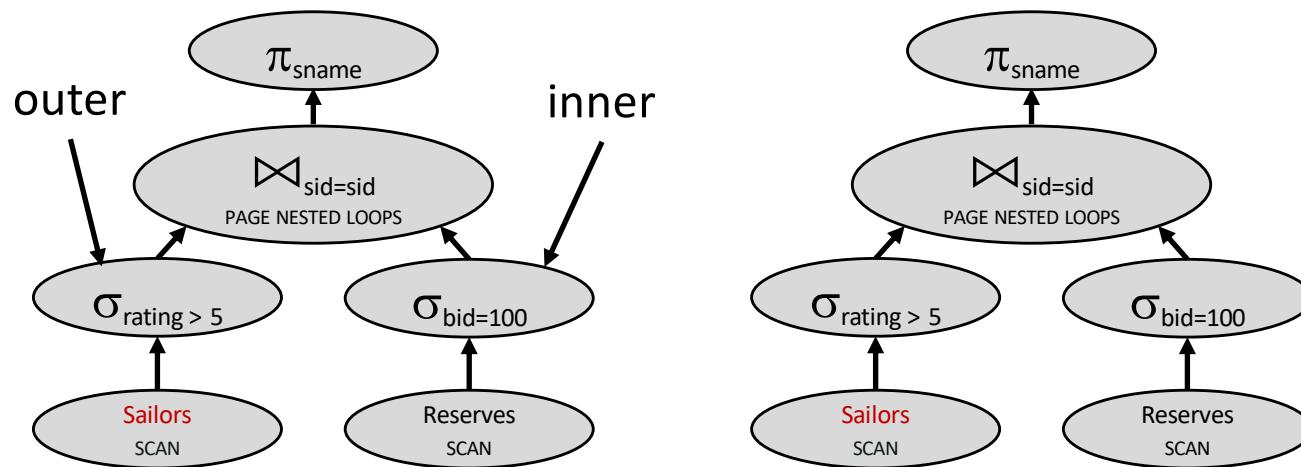
250,500 IOs

# So far, we've tried



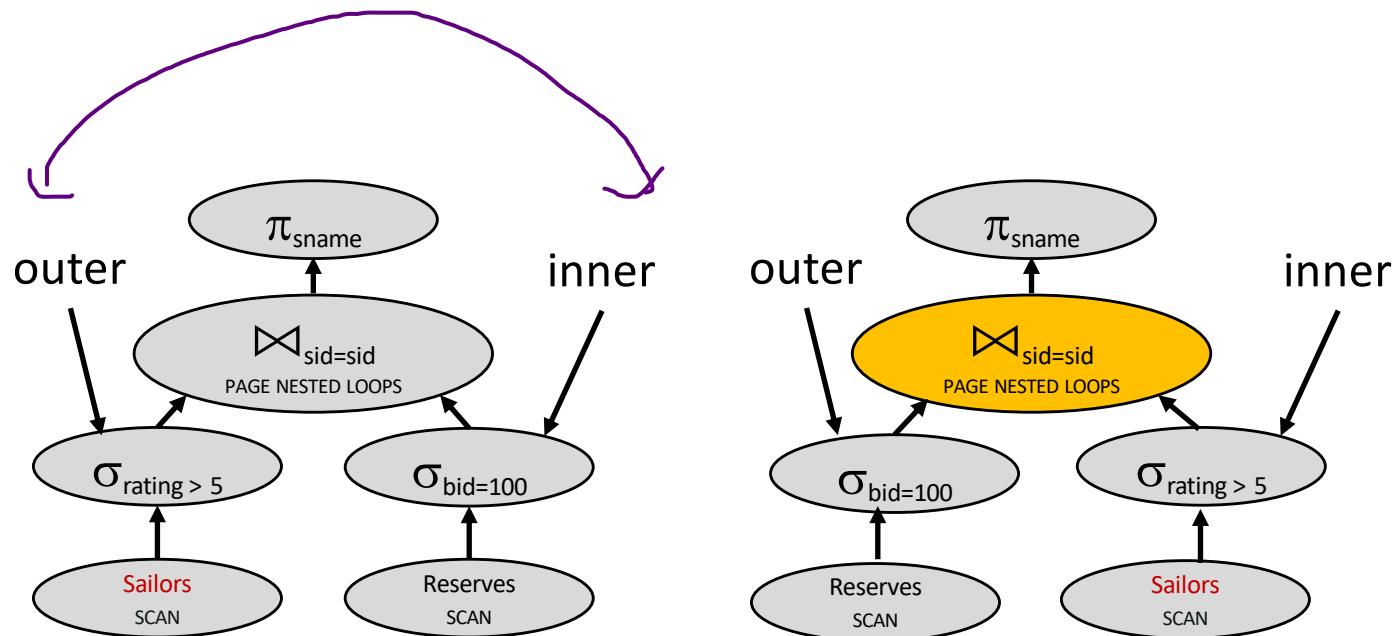
- Basic page nested loops (500,500)
- Selection pushdown on left (250,500)
- More selection pushdown on right (250,500)
- Next up, join ordering

## Next up: Join Ordering



250,500 IOs

# Join Ordering, cont



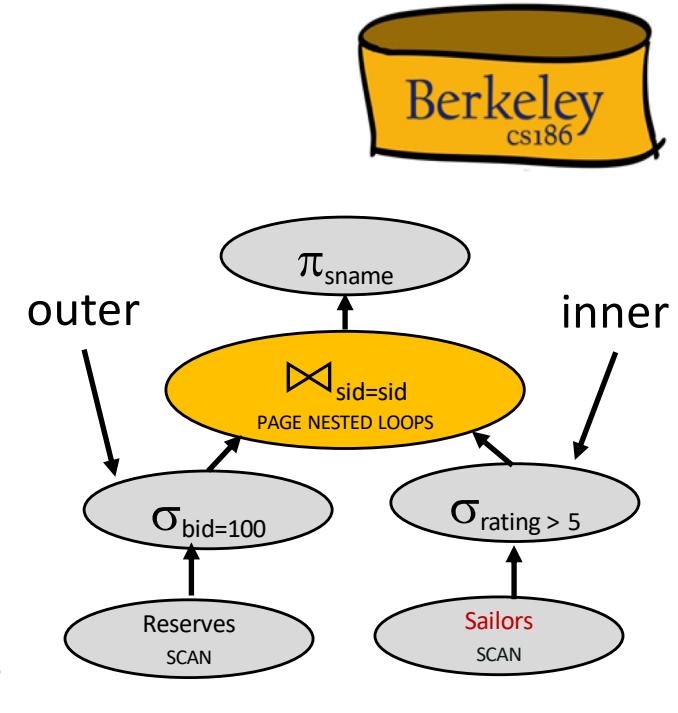
250,500 IOs

# Query Plan 4 Cost

- Let's estimate the cost:
- Scan Reserves (1000 IOs)
- For each pageful of Reserves for bid 100, Scan Sailors (500 IOs)
- Total:  $1000 + ??? * 500$
- Uniformly distributed across 100 boat values
- Total:  $1000 + 10 * 500$

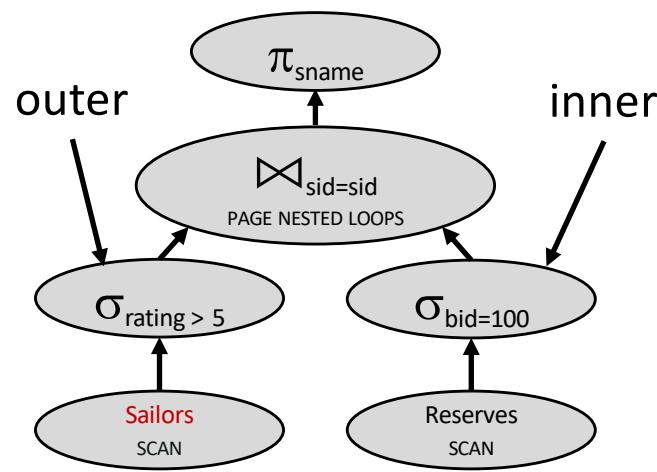
$$\frac{\text{num}(\text{bid} = 100)}{100 \text{ diff boats}} = \frac{1}{100}$$

$$\frac{1}{100} \cdot 100 \quad \text{boats} = 10$$

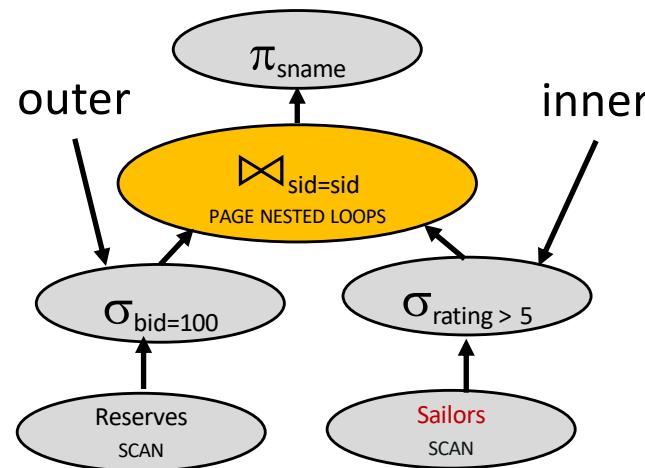


- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have B = 5 pages to use for joins

# Decision 3



250,500 IOs



6000 IOs

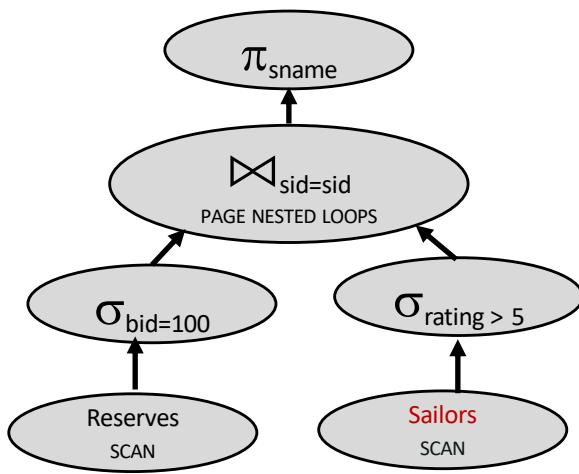
# So far, we've tried



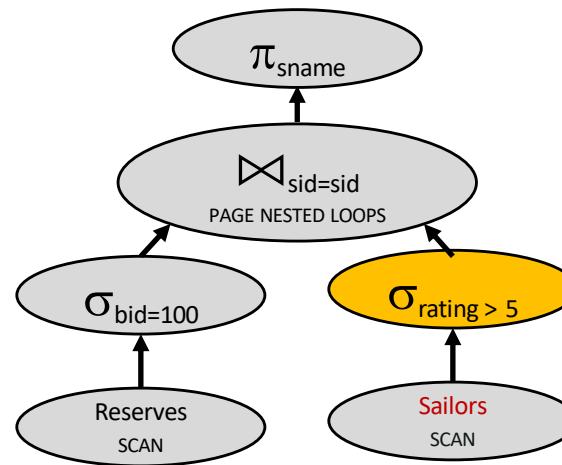
- Basic page nested loops (500,500)
- Selection pushdown on left (250,500)
- More selection pushdown on right (250,500)
- Join ordering (6000)
- Next up, materialization ...

# Query Plan 5

## Materializing Inner Loops



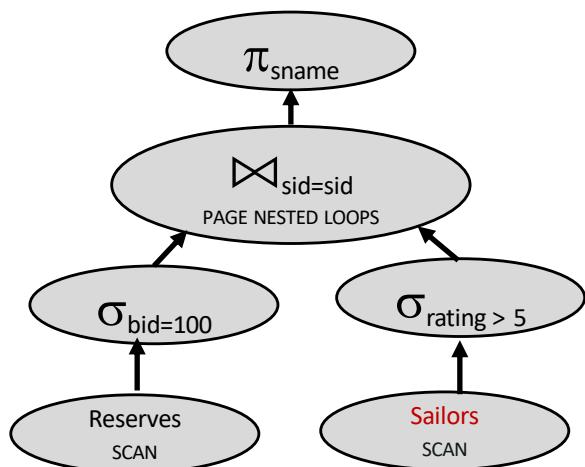
6000 IOs



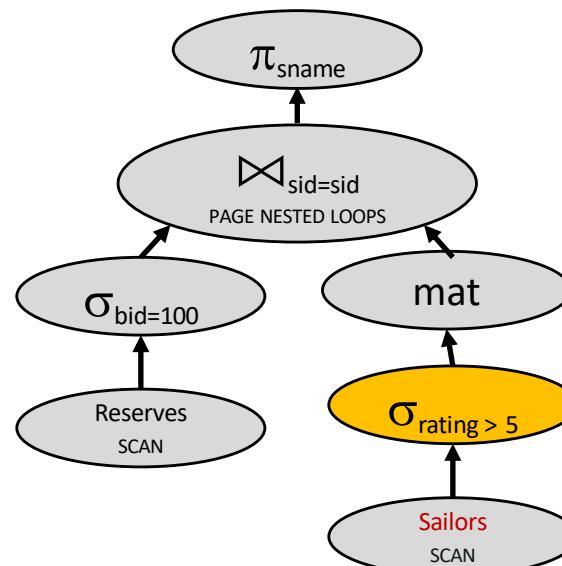
If you recall,  
selection pushdown  
on the right doesn't  
help because it is  
done on the fly.

What if we  
materialize the  
result after the  
selection?

# Materializing Inner Loops, cont



6000 IOs

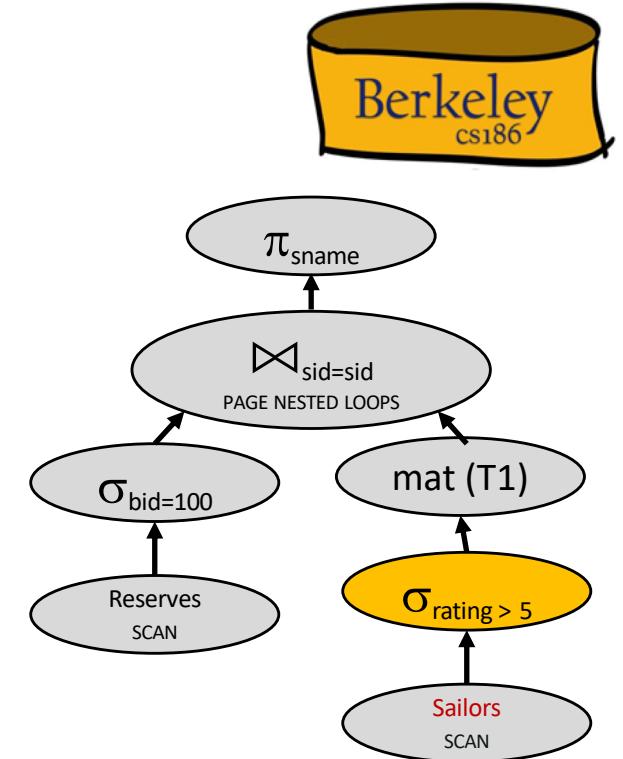


Cost???

Now don't need  
to rescan everytime  
(see query plan 3)

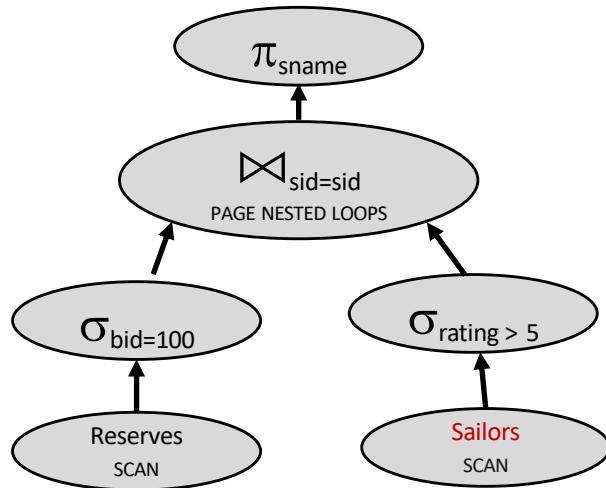
# Plan 5 Cost Analysis

- Let's estimate the cost:
- Scan Reserves (1000 IOs)
- Scan Sailors (500 IOs)
- Materialize Temp table T1 (??? IOs)
- For each pageful of Reserves for bid 100,  
Scan T1 (??? IOs)
- Total:  $1000 + 500 + ??? + 10 * ???$
- **1000 + 500+ 250 + (10 \* 250)**

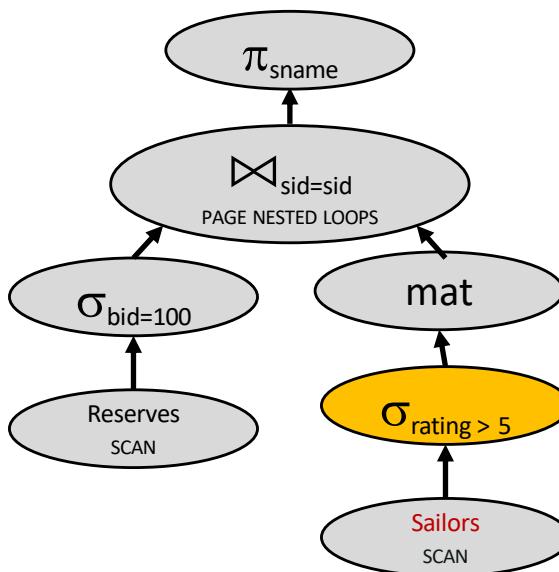


- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have  $B = 5$  pages to use for joins

# Materializing Inner Loops, cont.

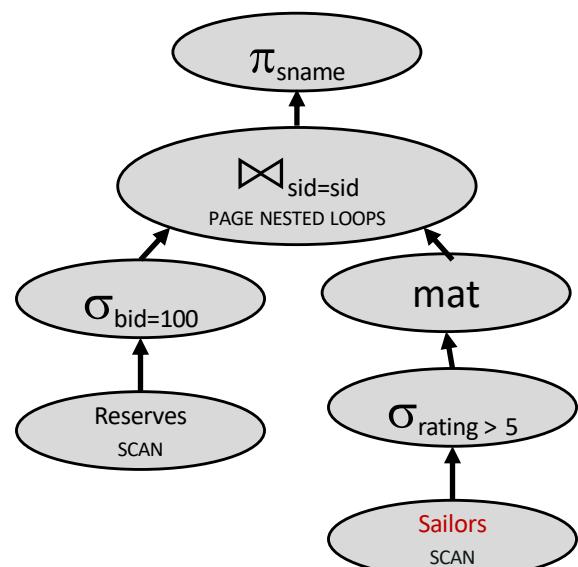


6000 IOs

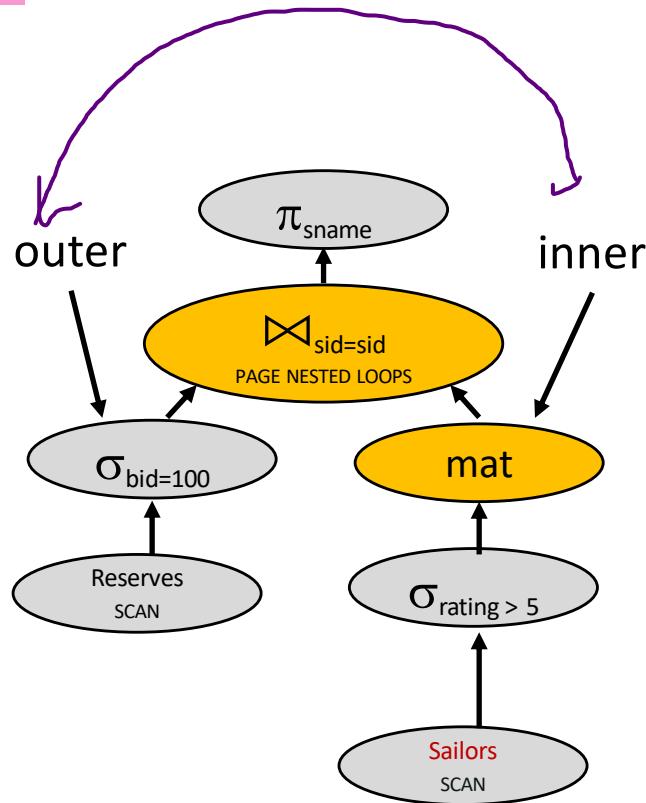


4250 IOs

## Join Ordering Again



4250 IOs



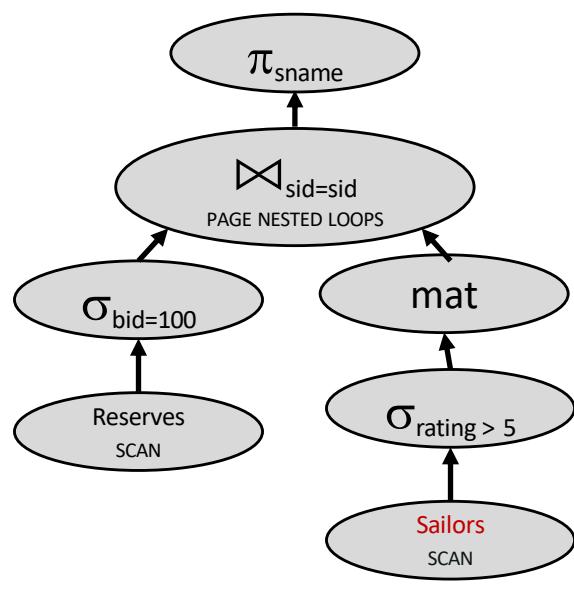
Let's try flipping  
the join order  
again with  
materialization  
trick



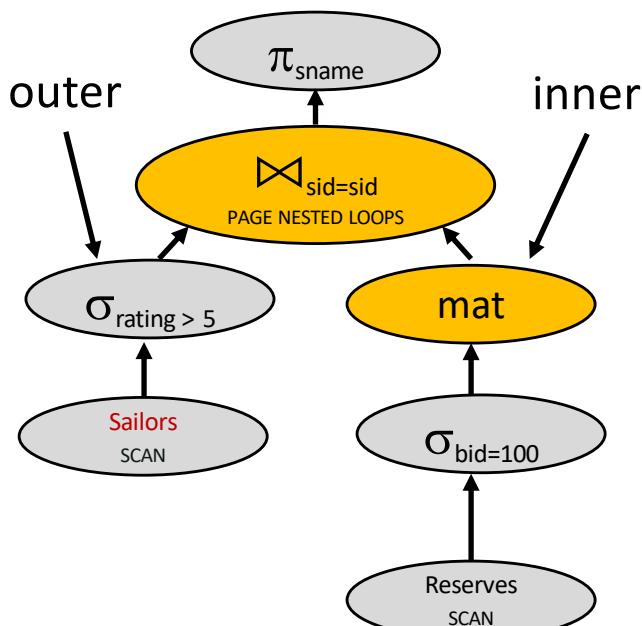
# Join Ordering Again, Cont



Let's try flipping  
the join order  
again with  
materialization  
trick



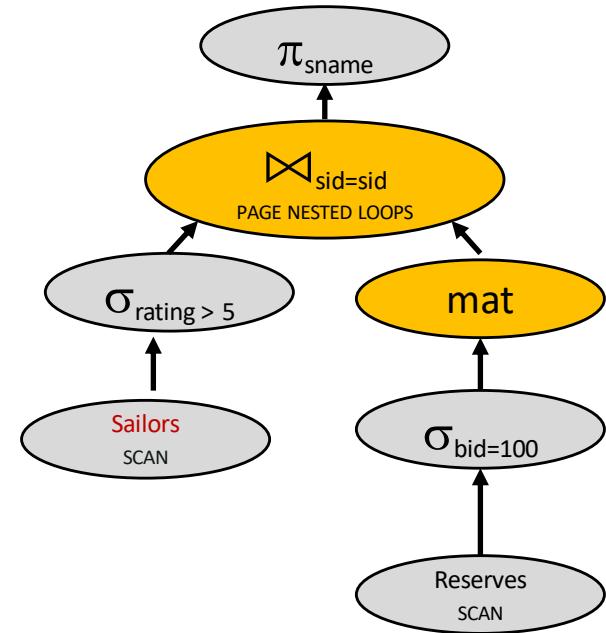
4250 IOs



Cost???

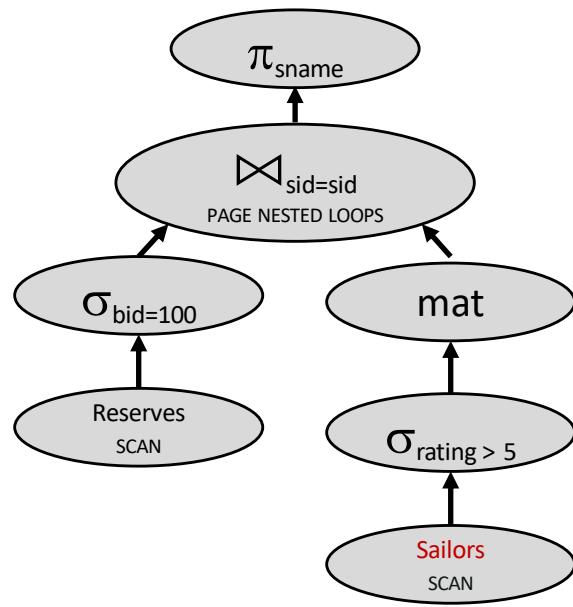
# Plan 6 Cost Analysis

- Let's estimate the cost:
- Scan Sailors (500 IOs)
- Scan Reserves (1000 IOs)
- Materialize Temp table T1 (??? IOs)
- For each pageful of high-rated Sailors,  
Scan T1 (??? IOs)
- Total:  $500 + 1000 + ??? + 250 * ???$
- **500 + 1000 +10 +(250 \*10)**

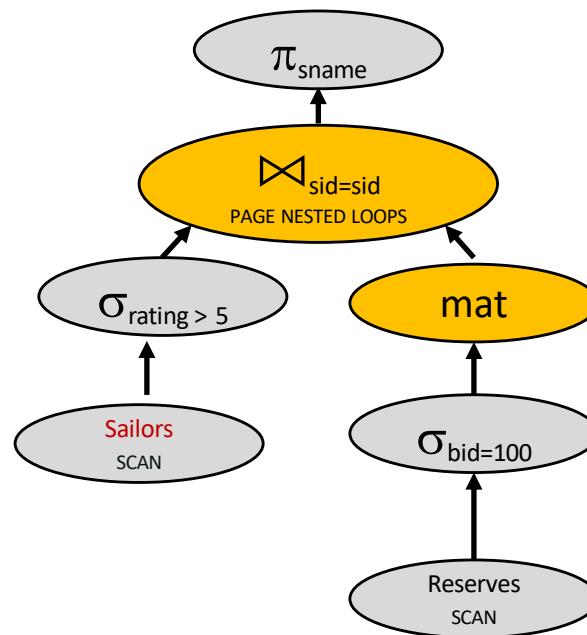


- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have  $B = 5$  pages to use for joins

# Decision 4



4250 IOs



4010 IOs

# So far, we've tried

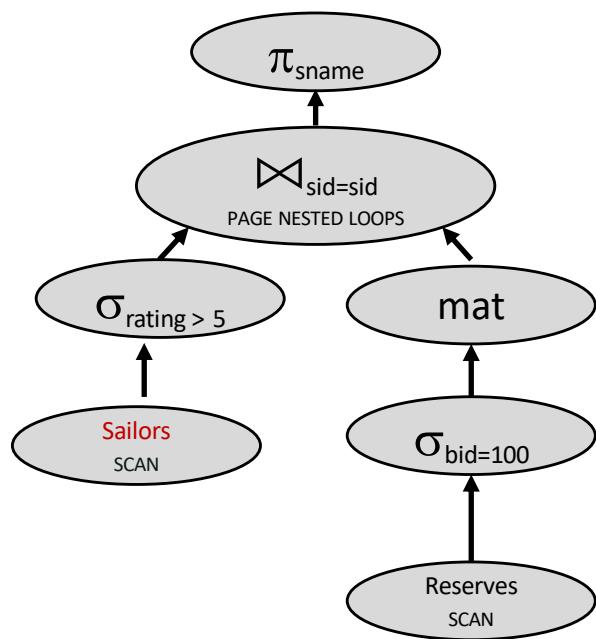


- Basic page nested loops (500,500)
- Selection pushdown on left (250,500)
- More selection pushdown on right (250,500)
- Join ordering (6000)
- Materializing inner loop (4250)
- Join ordering again with materialization (4010)
- Next up, sort merge ...

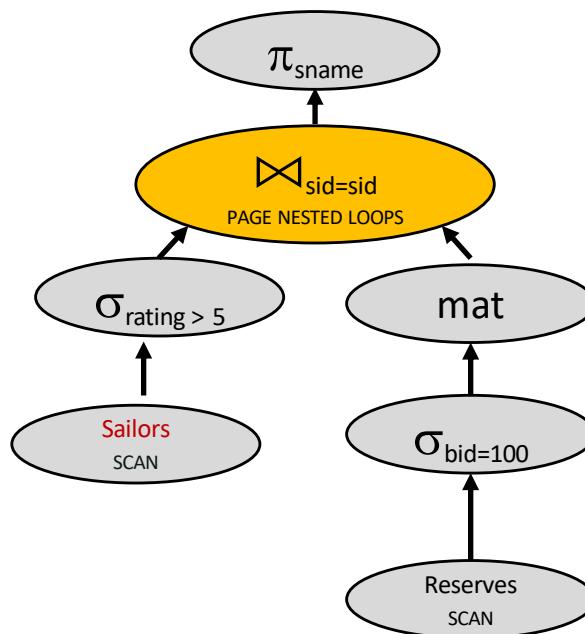
# Join Algorithm



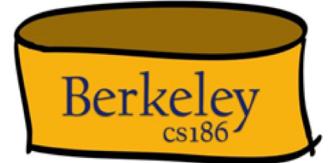
What if we  
change the join  
algorithm?



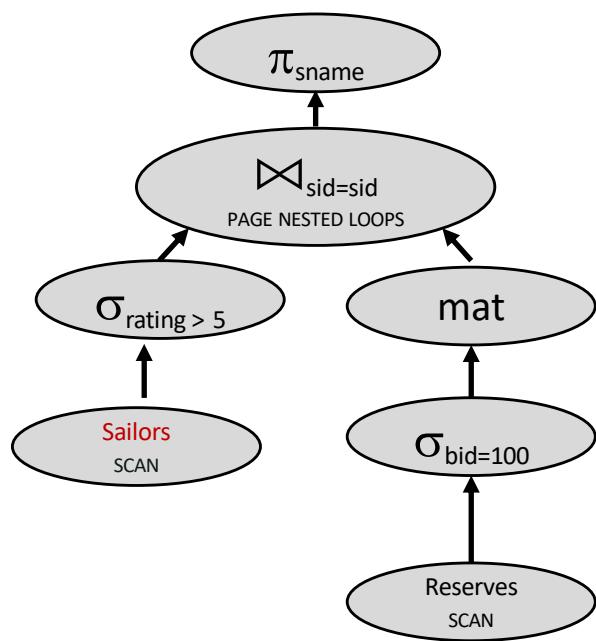
4010 IOs



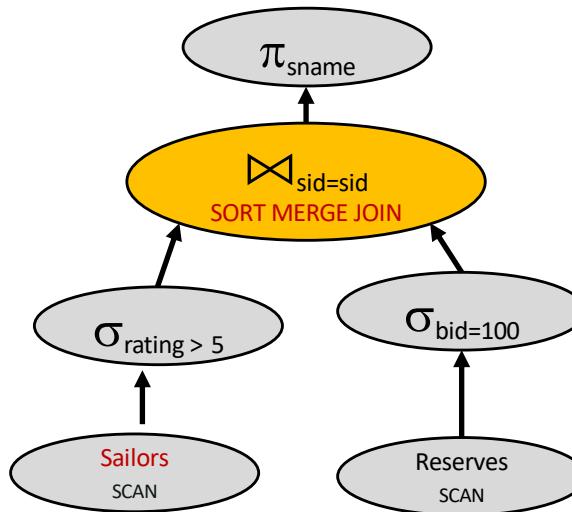
# Join Algorithm, cont.



What if we  
change the join  
algorithm?



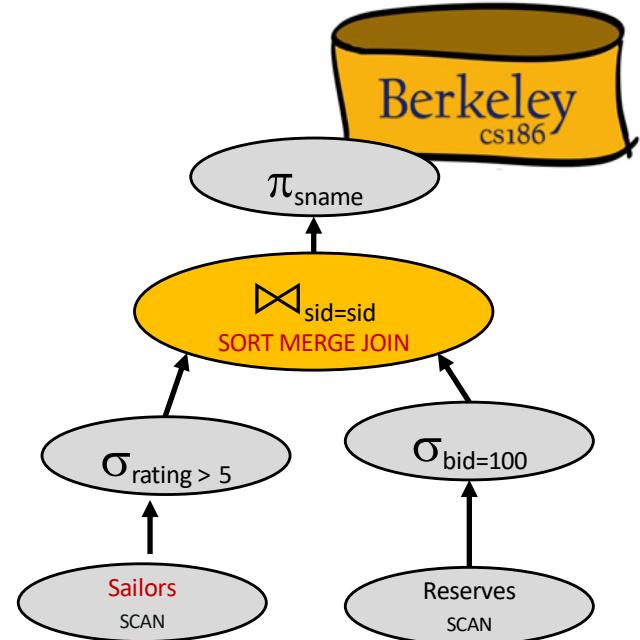
4010 IOs



Cost???

# Query Plan 7 Cost Analysis

- With 5 buffers, cost of plan:
- Scan Reserves (1000)
- Scan Sailors (500)
- Sort high-rated sailors (???)  
Note: pass 0 doesn't do read I/O, just gets input from select.
- Sort reservations for boat 100 (???)  
Note: pass 0 doesn't do read I/O, just gets input from select.
- How many passes for each sort?
- Merge (10+250) = 260
- Total:



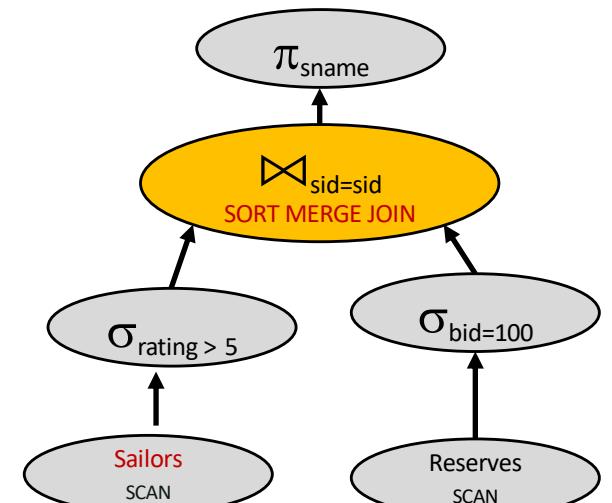
- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have B = 5 pages to use for joins

# Query Plan 7 Cost Analysis Part 2



- With 5 buffers, cost of plan:
- Scan Reserves (1000)
- Scan Sailors (500)
- Sort
  - 2 passes for reserves  
pass 0 = 10 to write, pass 1 =  $2 \times 10$  to read/write
  - 4 passes for sailors  
pass 0 = 250 to write, pass 1,2,3 =  $2 \times 250$  to read/write
- Merge ( $10 + 250$ ) = 260

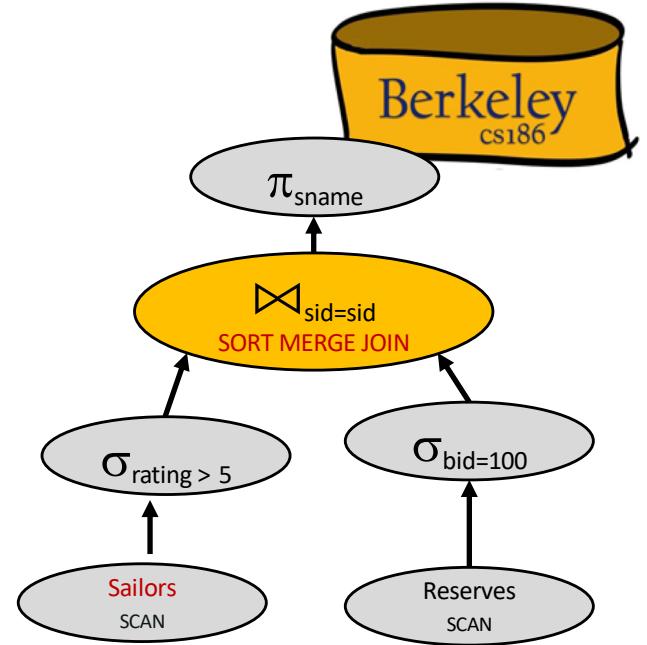
Scan both  $(1000 + 500)$  + sort reserves( $10 + 2 \times 10$ ) +  
sort sailors ( $250 + 3 \times 2 \times 250$ ) + merge ( $10 + 250$ ) = **3540**



- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have  $B = 5$  pages to use for joins

# Query Plan 7 Cost Analysis Part 2

- With 5 buffers, cost of plan:
  - Scan Reserves (1000)
  - Scan Sailors (500)
  - Sort
    - 2 passes for reserves  
pass 0 = 10 to write (2 runs of 5 each); pass 1 =  $2 \times 10$  to read/write (one representative from 2 runs)
    - 4 passes for sailors  
pass 0 = 250 to write (50 runs of 5 each); pass 1 (merging to give  $50/4=13$  runs of  $4 \times 5$  size each); pass 2 (merging to give  $13/4=4$  runs of  $4 \times 4 \times 5$  size each);  
pass 3 (merging to give one run of 250 in total) pass 1,2,3 =  $2 \times 250$  to read/write
  - Merge ( $10+250$ ) = 260
- Scan both ( $1000 + 500$ ) + sort reserves( $10 + 2 \times 10$ ) + sort sailors ( $250 + 3 \times 2 \times 250$ ) + merge ( $10+250$ ) = **3540**



## Reserves:

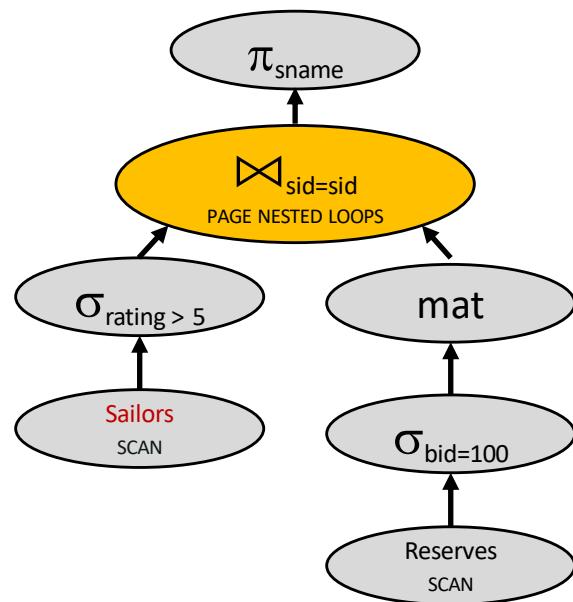
- Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
- Assume there are 100 boats (each equally likely)

## Sailors:

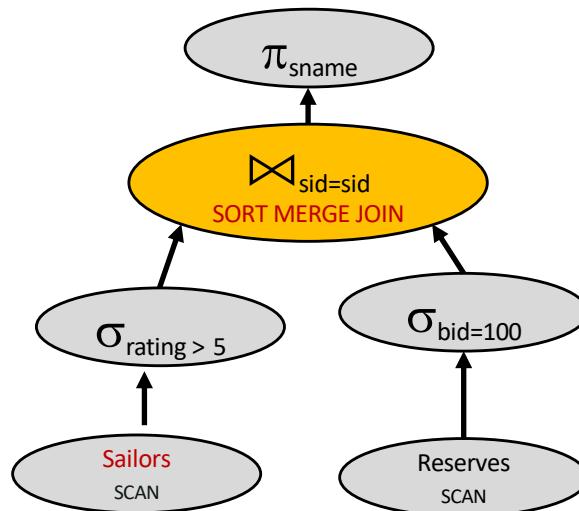
- Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
- Assume there are 10 different ratings (each equally likely)

- Assume we have B = 5 pages to use for joins

# Decision 5



4010 IOs



3540 IOs



We didn't go through the rest of the slides in class, but they are included here to show that we *really have a lot of query plans to choose from!!!*

# So far, we've tried

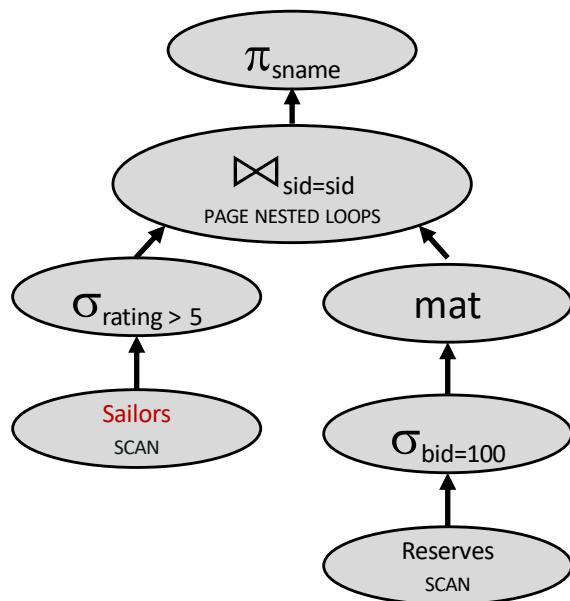


- Basic page nested loops (500,500)
- Selection pushdown on left (250,500)
- More selection pushdown on right (250,500)
- Join ordering (6000)
- Materializing inner loop (4250)
- Join ordering again with materialization (4010)
- Sort-merge join (3540)
- Next up, block nested ...

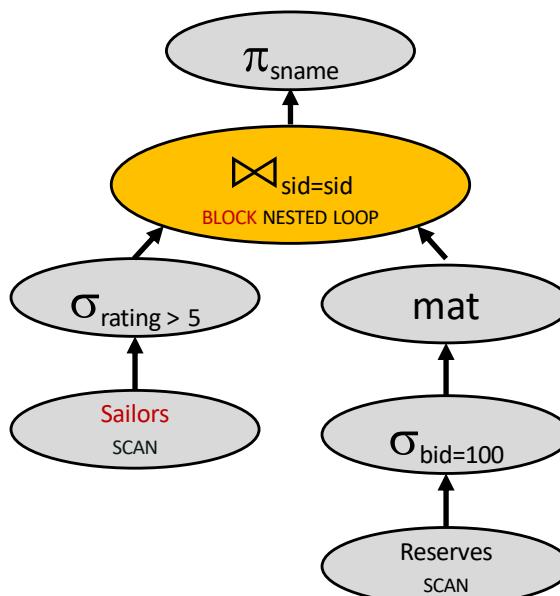
# Join Algorithm Again, Again



Returning to our best (so far) page nested loops plan again...



4010 IOs  
(And Sort-Merge at 3540 IOs)

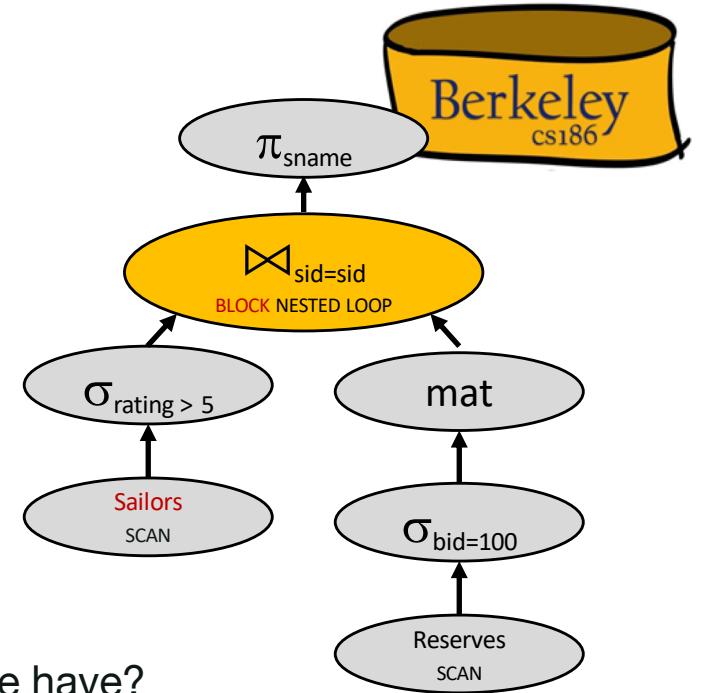


Cost???

# Query 8 Cost Analysis

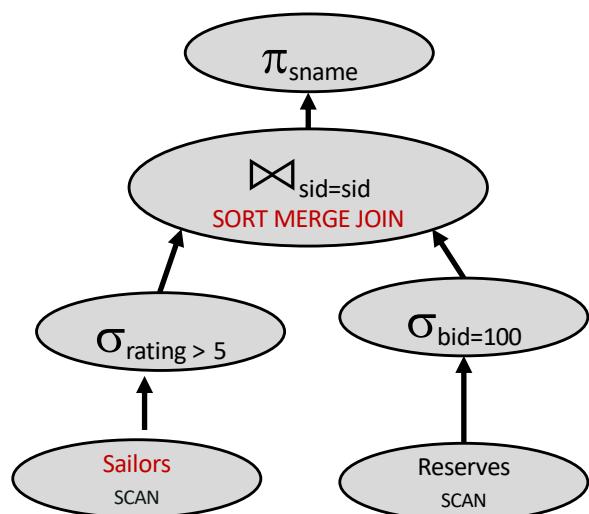
- With 5 buffers, cost of plan:
- Scan Sailors (500)
- Scan Reserves (1000)
- Write Temp T1 (10)
- For each blockful of high-rated sailors
  - Loop on T1 (**???** \* 10)
    - What is the block size? How many blocks (**???**) will we have?
    - 3 pages;  $\text{ceil}(250/3)$
- Total:
- Scan both(500 + 1000) + write out T1(10) + BNLJ ( $\text{ceil}(250/3) * 10$ )**

$$= 500 + 1000 + 10 + (84 * 10) = 2350$$
  - Reserves:
    - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
    - Assume there are 100 boats (each equally likely)
  - Sailors:
    - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
    - Assume there are 10 different ratings (each equally likely)
  - Assume we have B = 5 pages to use for joins

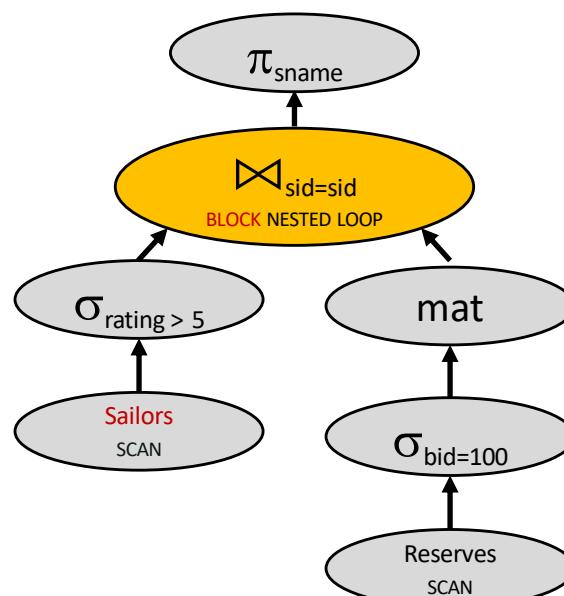


- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have B = 5 pages to use for joins

# Decision 6



3540 IOs



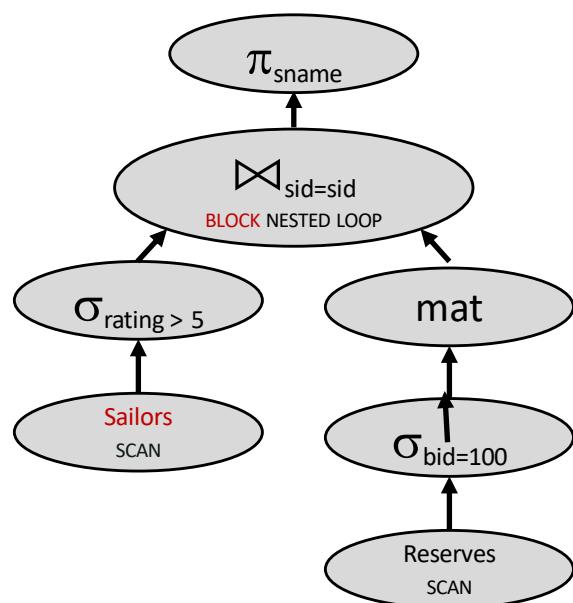
2350 IOs

# So far, we've tried

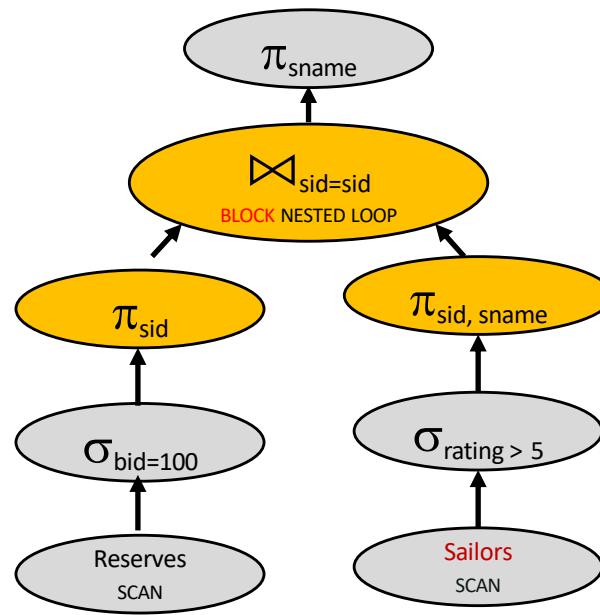


- Basic page nested loops (500,500)
- Selection pushdown on left (250,500)
- More selection pushdown on right (250,500)
- Join ordering (6000)
- Materializing inner loop (4250)
- Join ordering again with materialization (4010)
- Sort-merge join (3540)
- Block nested loops (2350)
- Next up, projection cascade

# Projection Cascade



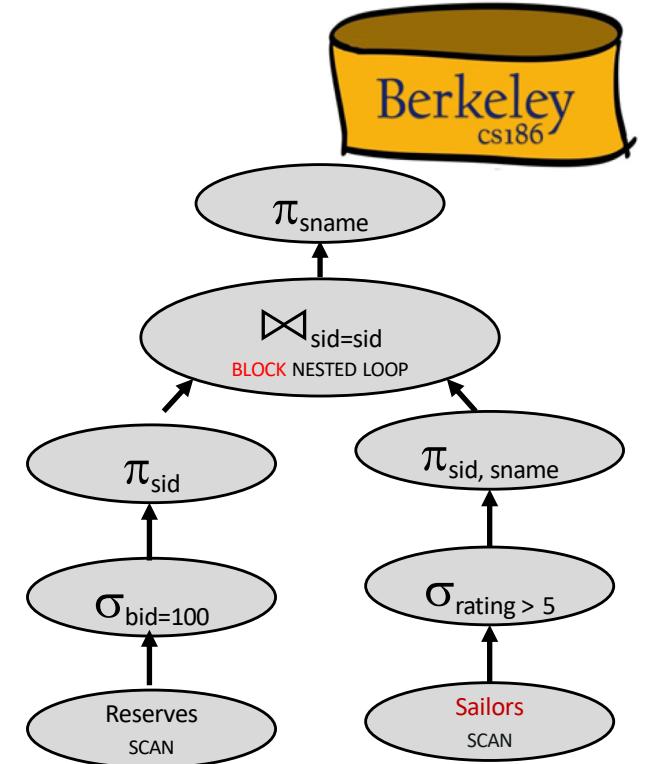
2350 IOs



Let's look at this plan with project cascade

# Query Plan 9 Cost Analysis

- With 5 buffers, cost of plan:
- Scan Reserves (1000)
- For each blockful of sids that rented boat 100
  - (recall Reserve tuple is 40 bytes, assume sid is 4 bytes)
    - 10 pages down to 1 page
    - So just one chunk for the BNLJ!
- Loop on Sailors ( $??? * 500 = 1 * 500$ )
- Total: 1500
- Can't do much better w/o indexes! Why?



- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have B = 5 pages to use for joins

# So far, we've tried

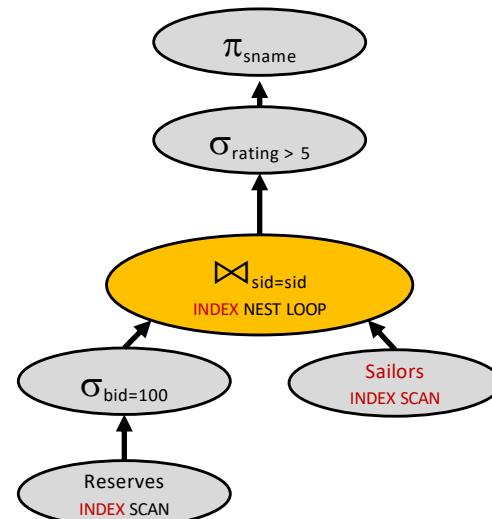


- Basic page nested loops (500,500)
- Selection pushdown on left (250,500)
- More selection pushdown on right (250,500)
- Join ordering (6000)
- Materializing inner loop (4250)
- Join ordering again with materialization (4010)
- Sort-merge join (3540)
- Block nested loops (2350)
- Projection cascade, plus reordering again (1500)
- Next up, indexes

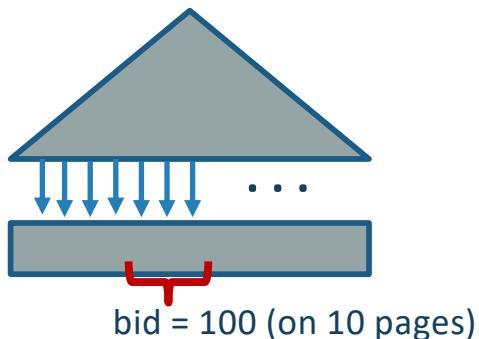
# How About Indexes?



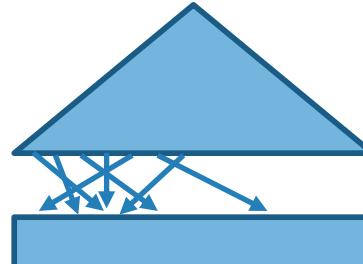
- Indexes:
  - Reserves.bid clustered
  - Sailors.sid unclustered
- Assume alt 2 indexes fit in memory



Reserves: bid



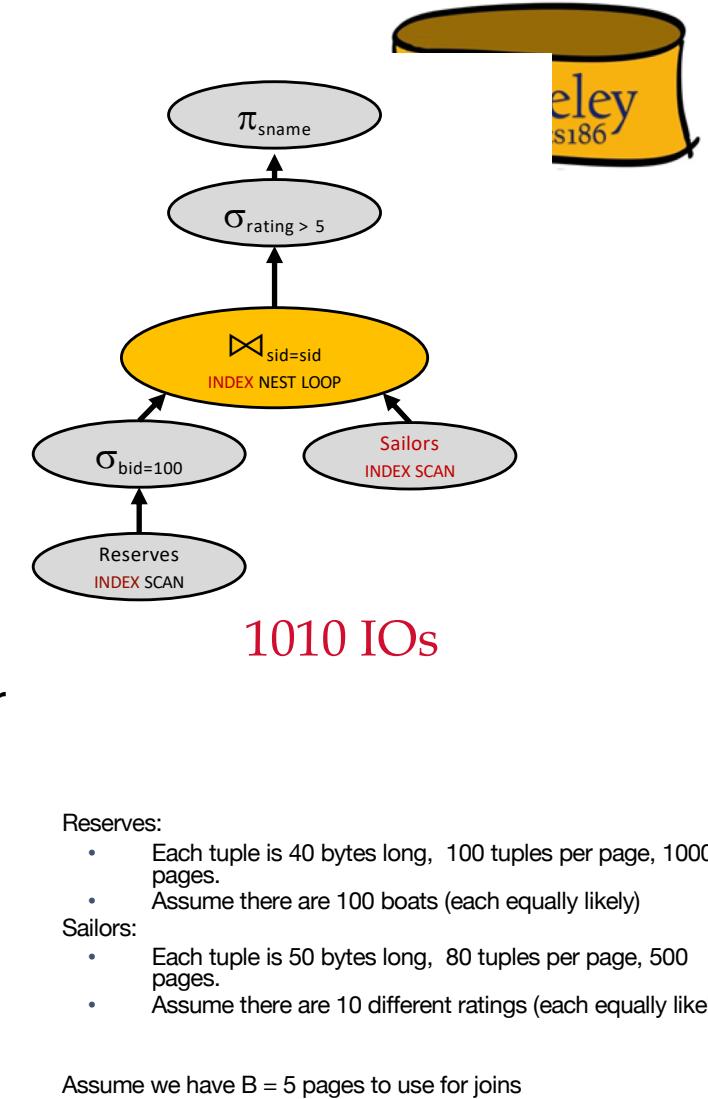
Sailors



- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats (each equally likely)
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings (each equally likely)
- Assume we have B = 5 pages to use for joins

# Index Cost Analysis Part 2

- With clustered index on bid of Reserves, we access how many pages of Reserves?:
  - $100,000/100 = 1000$  tuples on  $1000/100 = 10$  pages.
- for each Reserves tuple (1000 such tuples) get matching Sailors tuple (1 IO)
- $10 + 1000*1$
- Cost: Selection of Reserves tuples (10 I/Os); then, for each, must get matching Sailors tuple (1000); total 1010 I/Os.



# The Entire Story



- Basic page nested loops (500,500)
- Selection pushdown on left (250,500)
- More selection pushdown on right (250,500)
- Join ordering (6000)
- Materializing inner loop (4250)
- Join ordering again with materialization (4010)
- Sort-merge join (3540)
- Block nested loops (2350)
- Projection cascade, plus reordering again (1500)
- Index nested loops (1010)
- Still only a subset of the possible plans for this query!!!

# Summing up



- There are *lots* of plans
  - Even for a relatively simple query
- Engineers often think they can pick good ones
- Not so clear that's true!
  - Manual query planning can be tedious, technical
  - Machines are better at enumerating options than people
  - We will see soon how optimizers make simplifying assumptions to examine a reasonable set of plans