

Scaleable Online Statistical Processing

Chris Jermaine

Computer and Information Sciences and Engineering
Department

University of Florida, Gainesville

Performance: Aren't DBs Fast Enough?

You decide:

- Quick check of latest TPC-H results
- Spend \$0.75 million to store 300GB (17TB of 36GB disks)
- Q21 *still* takes 18 minutes to answer in throughput test
- Got 10TB?
- You can spend \$6.7 million to wait almost four hours for Q18

My Conclusion:

- Current AP solutions are not there
- Result: difficult to treat the DBMS as a sandbox

How To Address This Problem?

Not easy:

- AP studied intensively for 15 years
- We've only managed two serious solutions
 1. Pre-computed cubes - "OLAP" - often too restrictive
 2. "Glue-on" solution; add bitmaps to your DBMS
- And one interesting "new" idea
 - "Column-oriented" DBs
 - Though this really just fixes the problems with 2. above
- Is there another way to go?

Use Randomization!

Three key observations:

1. AP almost always statistical
2. Not always clear that \$1.3745m differs from \$1.3757m
3. Most exploratory queries are “wrong”

So, imagine the following DBMS:

- You ask **any** AP-style query, it gives you an immediate guess
- Guess bracketed by error guarantees
- Confidence region shrinks throughout computation
- Zero-width at query completion
- Total execution time same as classic RDBMS

Hasn't This Been Done?

UC Berkeley Control project:

- Developed online aggregation, ripple joins
- Got second two bullets to work (shrinking conf. region)
- But not scalability, generality

Aside: are synopsis structures relevant?

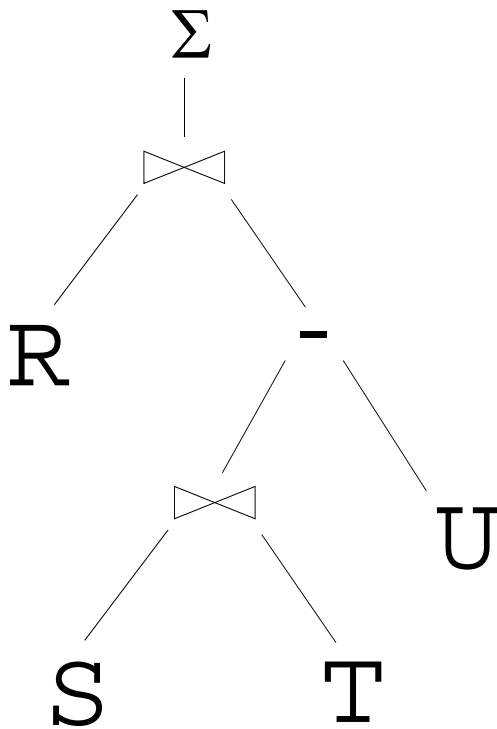
- Wavelets, histograms, sketches
- Great to study, but can they replace existing DBMS tech.?
 - Probably not; generally don't handle arbitrary queries
 - Bigger issue is that they are fixed precision

Our Goal

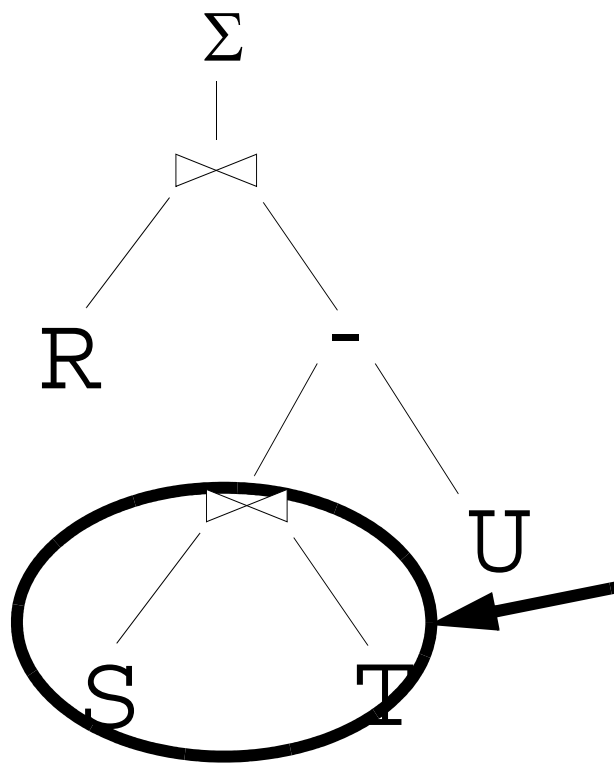
- Re-design database from ground up based on randomization
- Try to meet each of the five goals given previously:
 1. Any AP-style query → immediate guess
 2. Guess bracketed by error guarantees
 3. Shrinking confidence region
 4. Zero-width at completion
 5. Fast as classic RDBMS
- Resulting system is called *DB-Online*, or DBO for short

Seems Pretty Hard To Do...

Given an arbitrary plan like this one, how are you ever going to guess the answer from start to finish?



Seems Pretty Hard To Do...



Given an arbitrary plan like this one, how are you ever going to guess the answer from start to finish?

Our first (imperfect) idea was to tackle a 2-table join, like this

Uses ripple join as basic building block...

The Ripple Join

- Scan randomly-permuted input relations in parallel

		(Tom, 12)	(Jan, 3)	(Sam, 5)	(Jon, 10)	(Joe, 8)	(Tim, 4)
(Jan)							
(Tim)							
(Sue)							
(Joe)							
(Jon)							
(Nat)							

SALES

$$\frac{6}{1} \times \frac{6}{1} \times (0) = 0$$

```
SELECT SUM SALES.b
FROM EMP, SALES
WHERE EMP.a = SALES.a
```

The Ripple Join

- Scan randomly-permuted input relations in parallel

		(Tom, 12)	(Jan, 3)	(Sam, 5)	(Jon, 10)	(Joe, 8)	(Tim, 4)
(Jan)		■	■				
(Tim)		■	■				
(Sue)							
(Joe)							
(Jon)							
(Nat)							

SALES

$$\frac{6}{2} \times \frac{6}{2} \times (3) = 27$$

```
SELECT SUM SALES.b
FROM EMP, SALES
WHERE EMP.a = SALES.a
```

The Ripple Join

- Scan randomly-permuted input relations in parallel

	(Tom, 12)	(Jan, 3)	(Sam, 5)	(Jon, 10)	(Joe, 8)	(Tim, 4)
(Jan)						
(Tim)						
(Tom)						
(Joe)						
(Jon)						
(Sue)						

SALES

$$\frac{6}{4} \times \frac{6}{4} \times (3 + 12) = 33.75$$

```
SELECT SUM SALES.b
FROM EMP, SALES
WHERE EMP.a = SALES.a
```

The Ripple Join

- Scan randomly-permuted input relations in parallel

	(Tom, 12)	(Jan, 3)	(Sam, 5)	(Jon, 10)	(Joe, 8)	(Tim, 4)
(Jan)						
(Tim)						
(Tom)						
(Joe)						
(Jon)						
(Sue)						

SALES

$$\frac{6}{5} \times \frac{6}{5} \times (3 + 12 + 10) = 36$$

```
SELECT SUM SALES.b
FROM EMP, SALES
WHERE EMP.a = SALES.a
```

The Ripple Join

- Scan randomly-permuted input relations in parallel

EMP	(Tom, 12)	(Jan, 3)	(Sam, 5)	(Jon, 10)	(Joe, 8)	(Tim, 4)
(Jan)		■				
(Tim)						
(Tom)	■					
(Joe)						
(Jon)				■		
(Sue)						

SALES

Final answer

$$\frac{6}{6} \times \frac{6}{6} \times (3 + 12 + 10) = 25$$

```
SELECT SUM SALES.b
FROM EMP, SALES
WHERE EMP.a = SALES.a
```

The Ripple Join

- Issue: only hashed RJ is practical
- But requires $e \in \text{EMP}'$, $s \in \text{SALES}'$ in memory

What do you do when you run out of memory?

Sort-Merge-Shrink Join

- 3-stage join algorithm

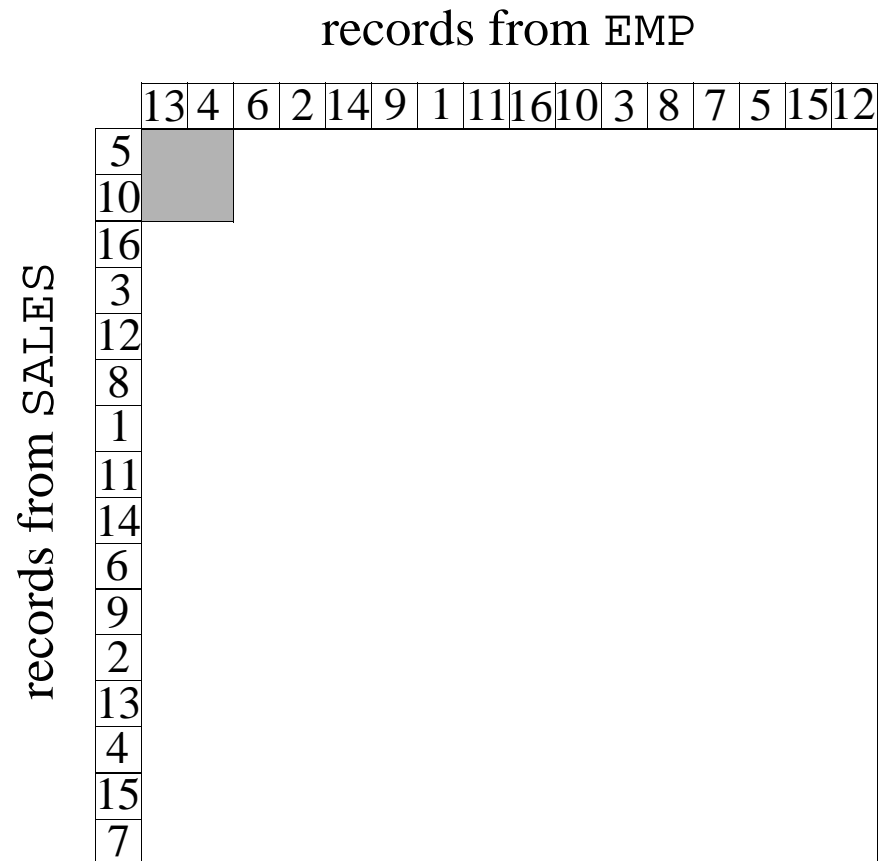
records from EMP

	13	4	6	2	14	9	1	11	16	10	3	8	7	5	15	12
5																
10																
16																
3																
12																
8																
1																
11																
14																
6																
9																
2																
13																
4																
15																
7																

records from SALES

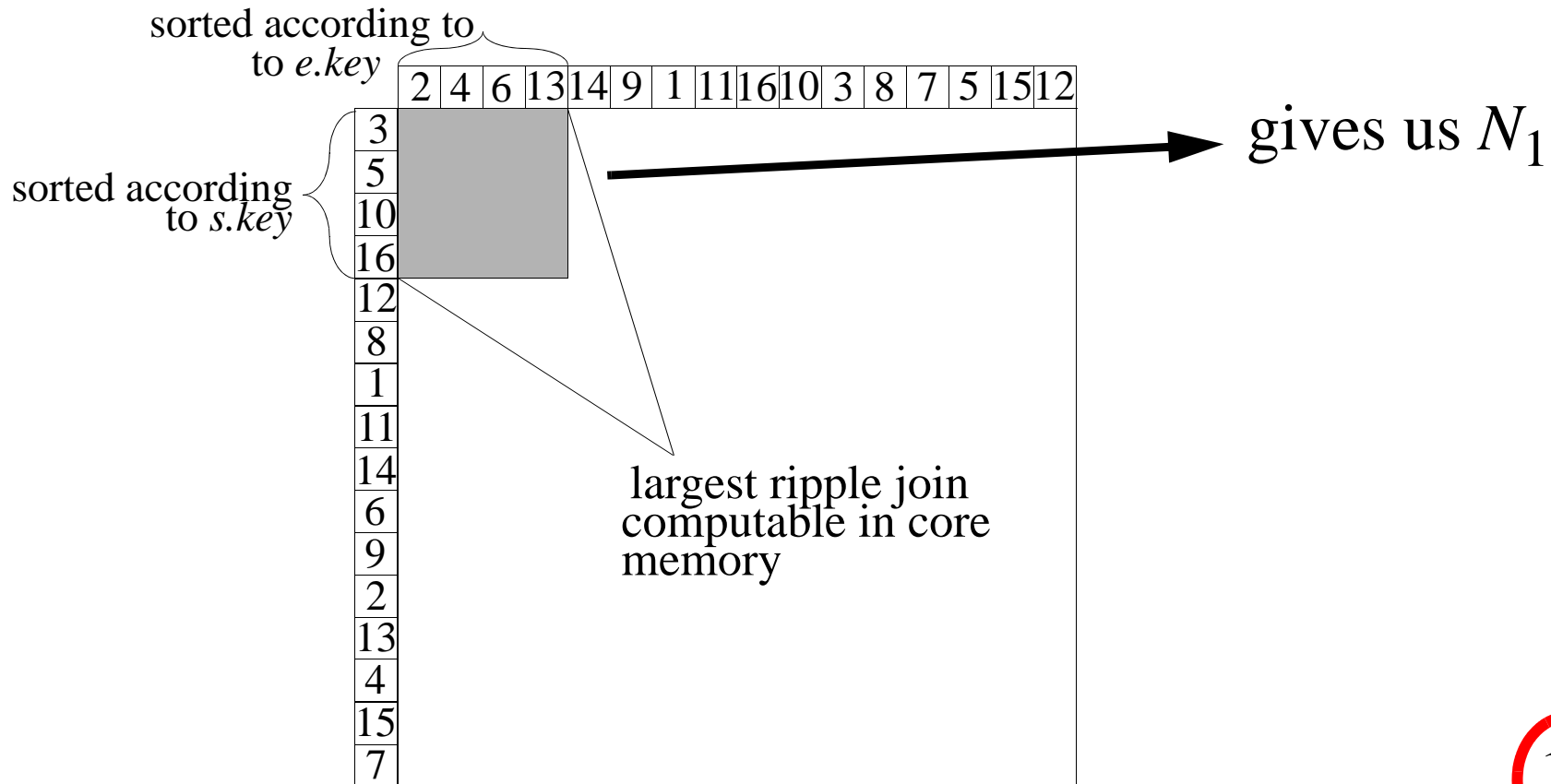
Sort-Merge-Shrink Join

- Sort phase: starts just like a ripple join



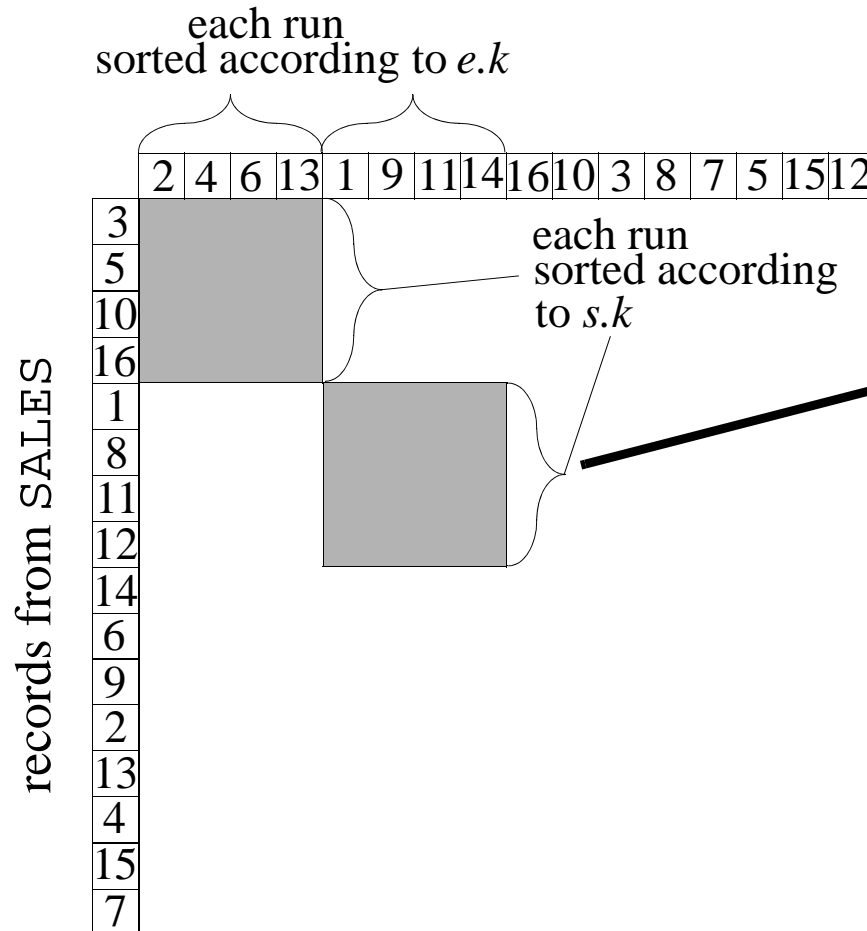
Sort-Merge-Shrink Join

- When memory fills, records from first RJ sorted and written back to disk



Sort-Merge-Shrink Join

- Process is then repeated



gives us N_2

$$N = \sum_i^n w_i N_i$$

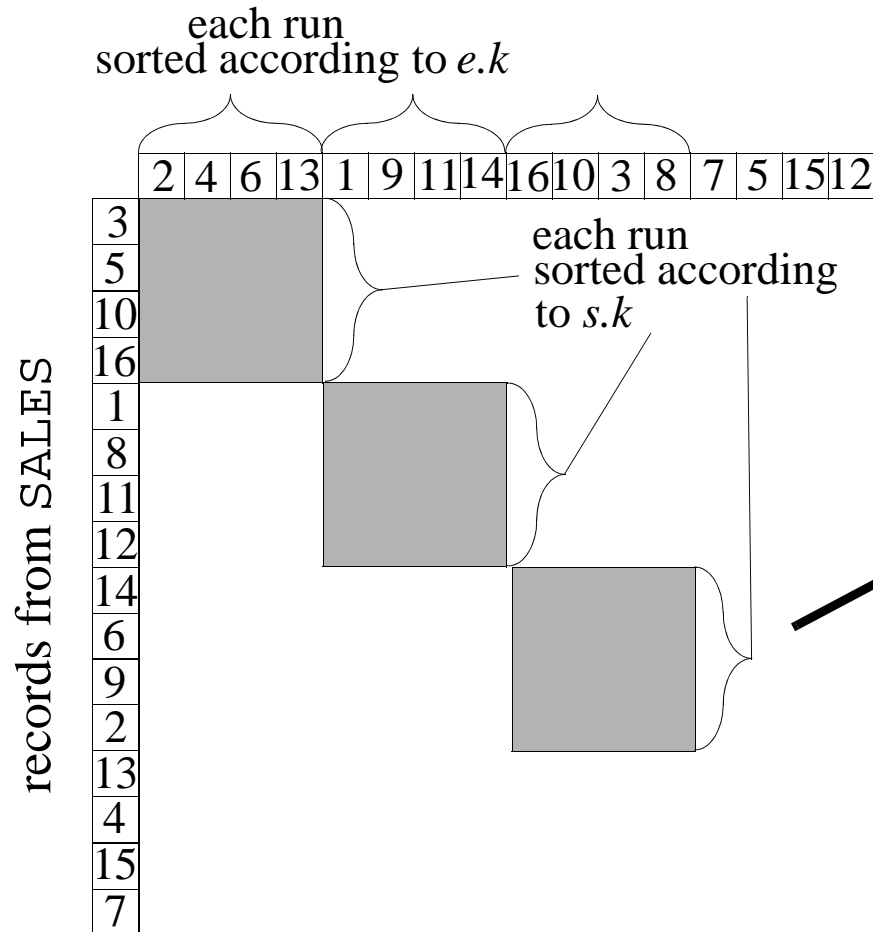
with

$$\sum_i^n w_i = 1$$

is unbiased

Sort-Merge-Shrink Join

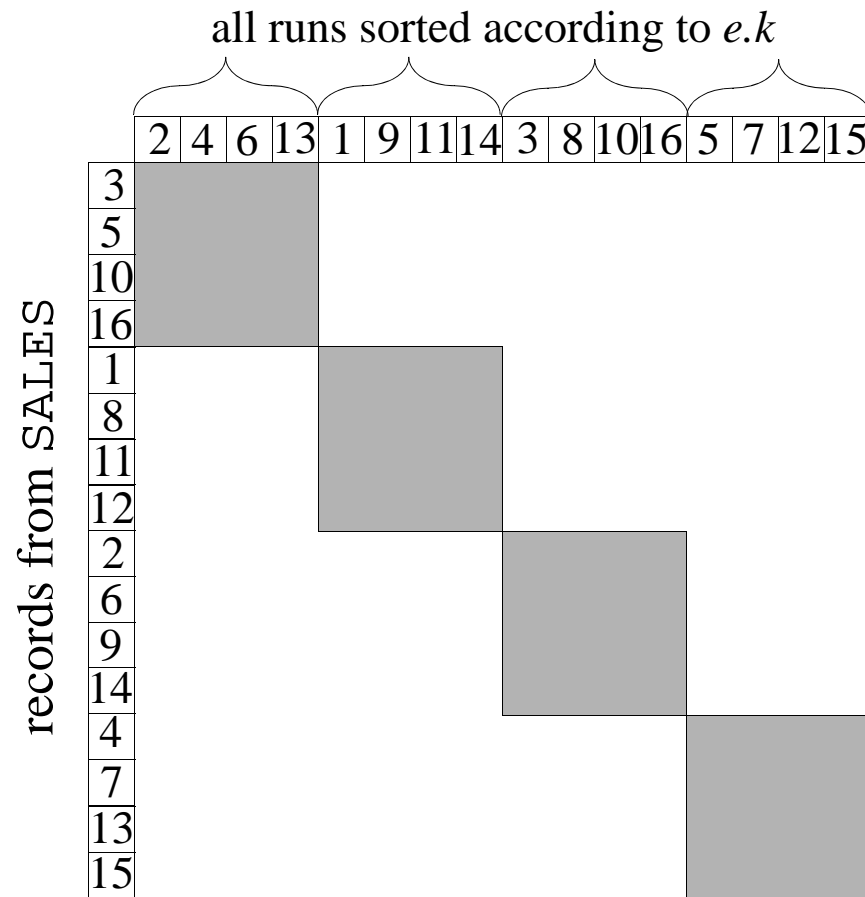
- Process is then repeated



gives us N_3

Sort-Merge-Shrink Join

- Until sort phase completes and all database records have participated in one RJ



Now we have N_1 through N_4

gives us N_4

Sort-Phase Statistical Considerations

- During sort phase, let N_i be estimate from i th ripple join - N_i unbiased so N is

- Know $Var(N) =$

$$\sum_i^n w_i^2 Var(N_i) + \sum_j^n \sum_{i \neq j}^n w_i w_j Cov(N_i, N_j)$$

- Use this quantity to give bounds via CLT or other appropriate result

Computing the RJ Variance

$$\sum_i^n w_i^2 \text{Var}(N_i) + \sum_j^n \sum_{i \neq j}^n w_i w_j \text{Cov}(N_i, N_j)$$

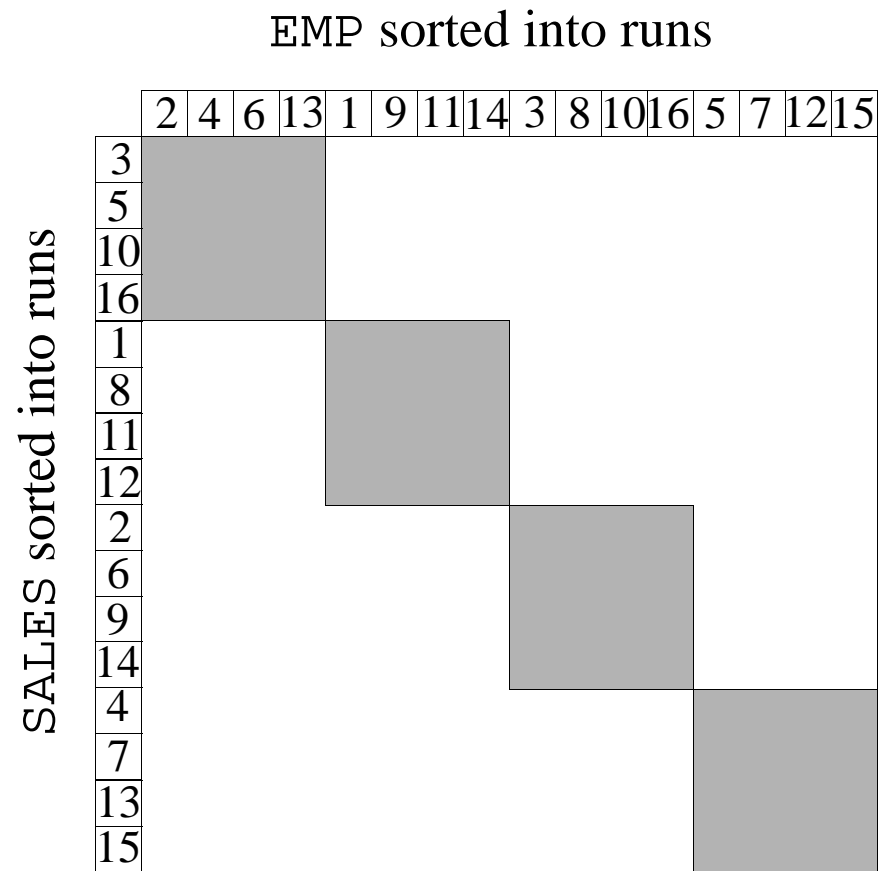
- $\text{Var}(N_i)$ can be estimated using formulas from HH99 (large sample w. replacement)
- Or can use the nasty (yet exact) permutational formula we derived
- w_i computed via a straightforward quadratic optimization problem

Computing the RJ Variance

$$\sum_i^n w_i^2 \text{Var}(N_i) + \sum_j^n \sum_{i \neq j}^n w_i w_j \text{Cov}(N_i, N_j)$$

- How about the covariance?
- Can be estimated using the nasty formulas we derived
- Or can possibly ignore it... usually negative anyway

How To Finish the Join?

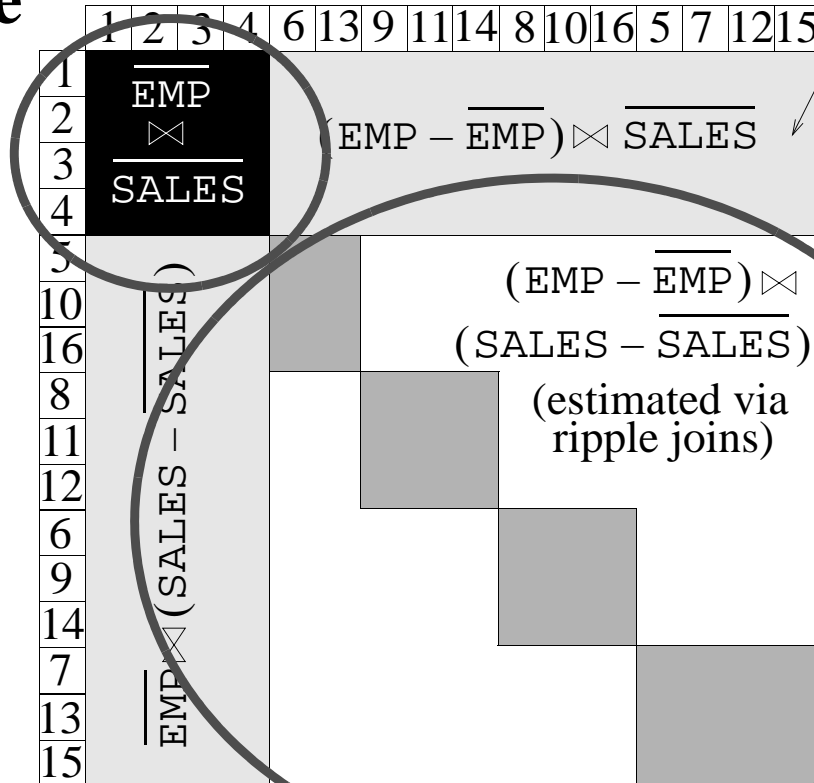


Sort-Merge-Shrink Join

- The merge/shrinking phases begin...

region contains no record pairs where $pred(e, s)$ evaluates to true

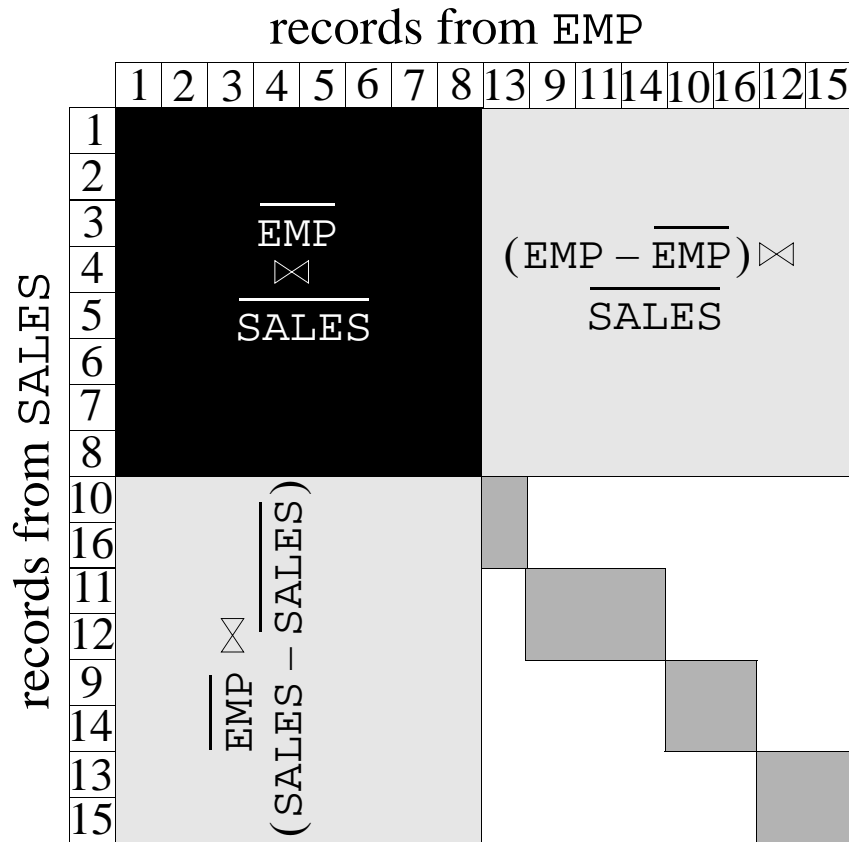
Merge Phase



Shrinking Phase

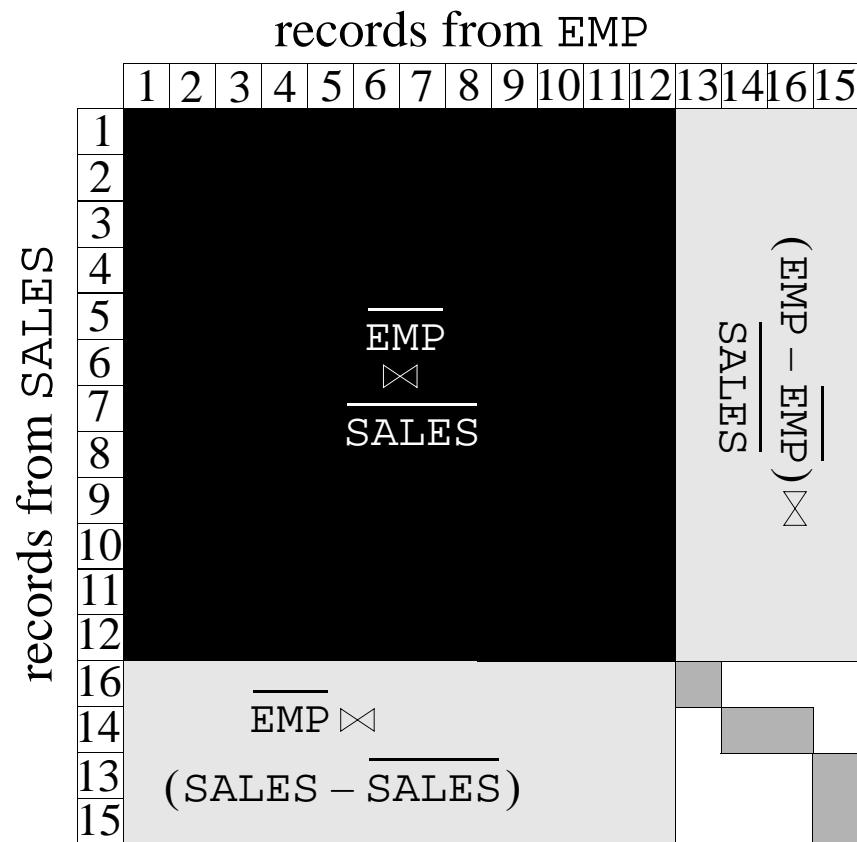
Sort-Merge-Shrink Join

- Merged portion of data space increases



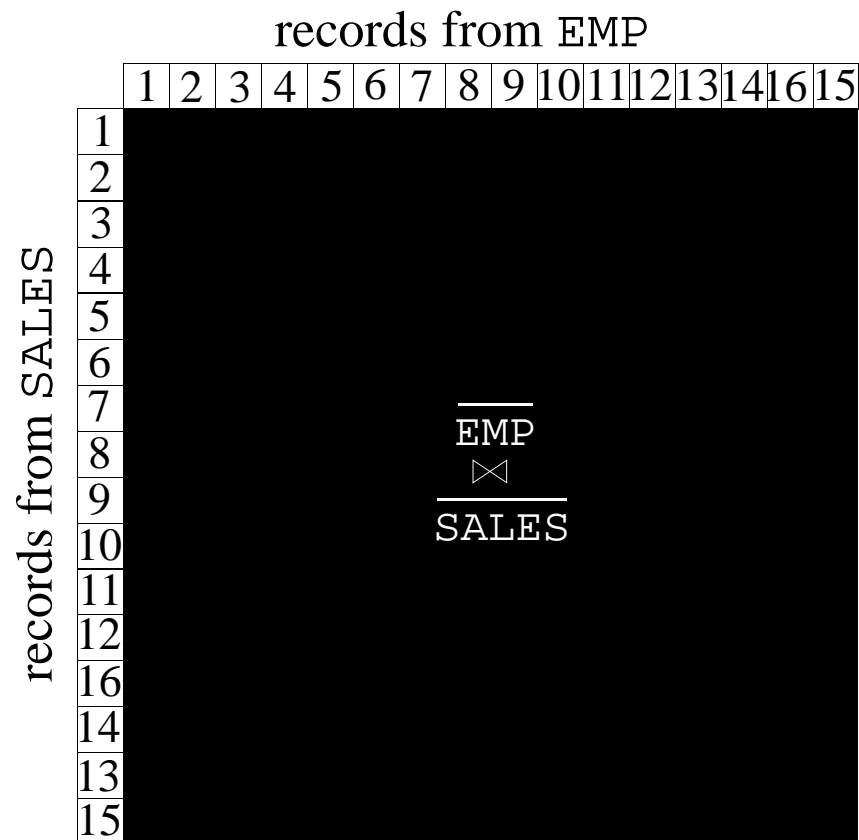
Sort-Merge-Shrink Join

- Until it dominates the data space



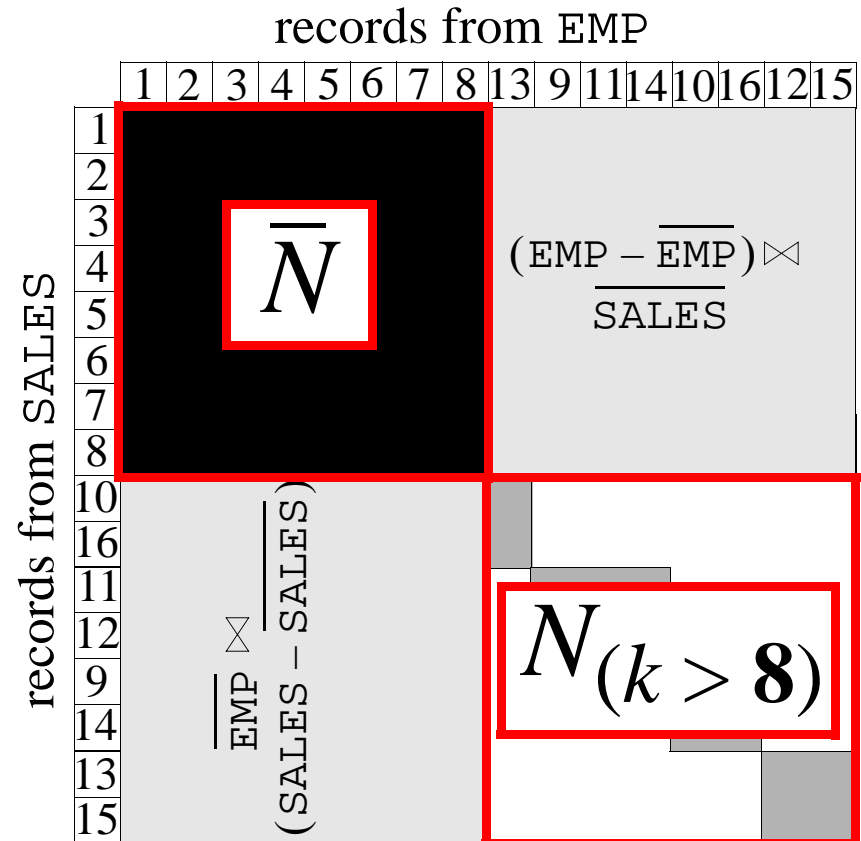
Sort-Merge-Shrink Join

- Then the join completes



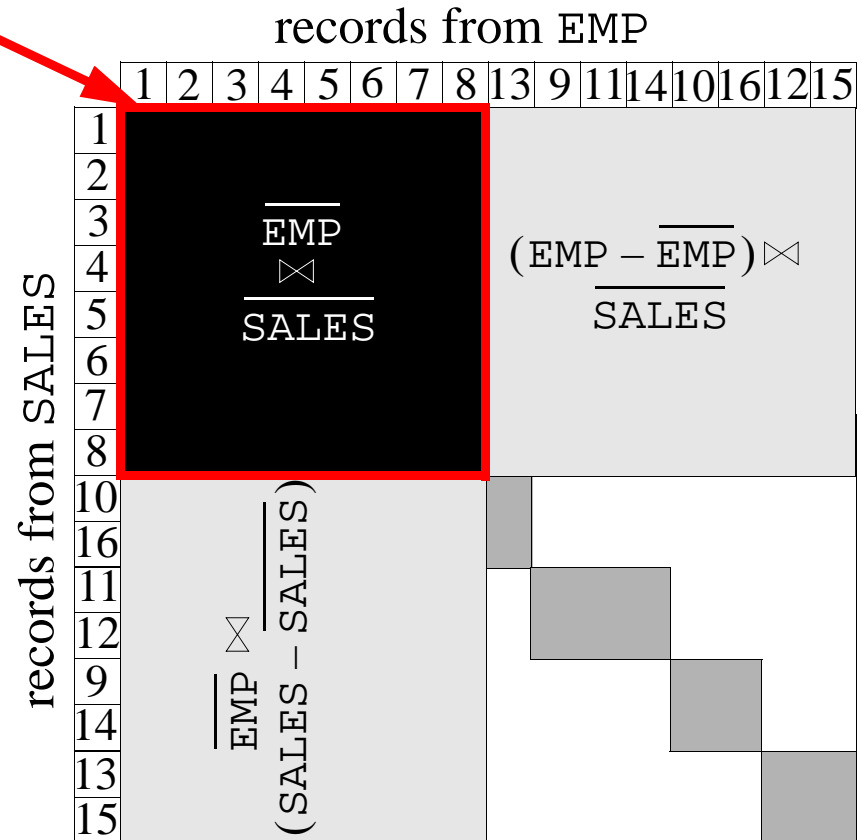
Merge/Shrink Statistics

- $\bar{N} + N_{(k > \mathbf{k})}$ is unbiased for query result
- $Var(\bar{N} + N_{(k > \mathbf{k})})$ is $Var(N_{(k > \mathbf{k})})$



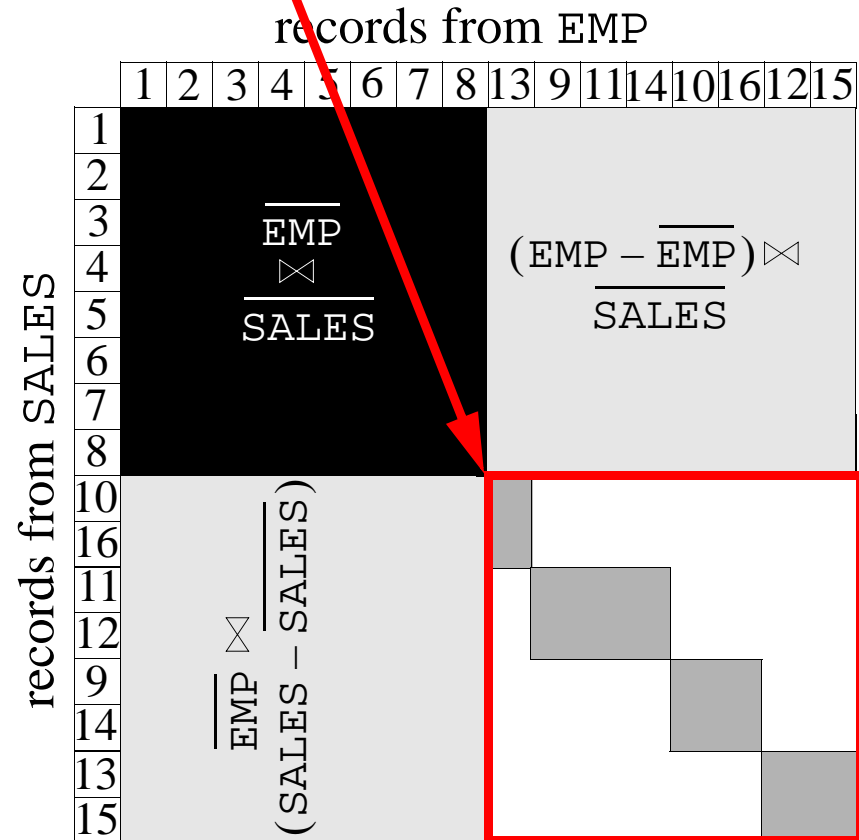
Merge/Shrink Statistics

- Getting \bar{N} is easy



Merge/Shrink Statistics

- But $N_{(k > \mathbf{k})}$, $Var(N_{(k > \mathbf{k})})$ is harder. Why?
- Each RJ has been written to disk, so no way to compute $N_{(k > \mathbf{k})}$ w/o re-reading data



Merge/Shrink Statistics

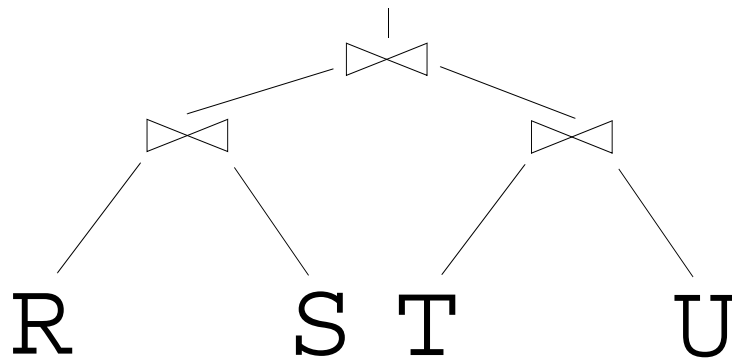
- Solution: precompute and store $N_{(k > \mathbf{k})}$ and $Var(N_{(k > \mathbf{k})})$ for “enough” \mathbf{k} values
- Output new estimate every time you merge past one of those values
- Done via reverse ripple join during sort phase, before each run written back to disk
- Small storage: $\sim 3600 \times 100 \times 12$ bytes

What Do We Find When We Run This?

- Example, joining two, 20GB tables:
 - Memory gone ~20 seconds, w. one disk at 50MB/sec
 - Confidence bounds halve after 100 secs
 - Confidence bounds halve again 1000 secs
 - 1/8 to 1/10 as wide after 3000 seconds

SMS Join Is Nice, Won't Work in DBO

In retrospect, a neat algorithm of questionable utility... say you have:

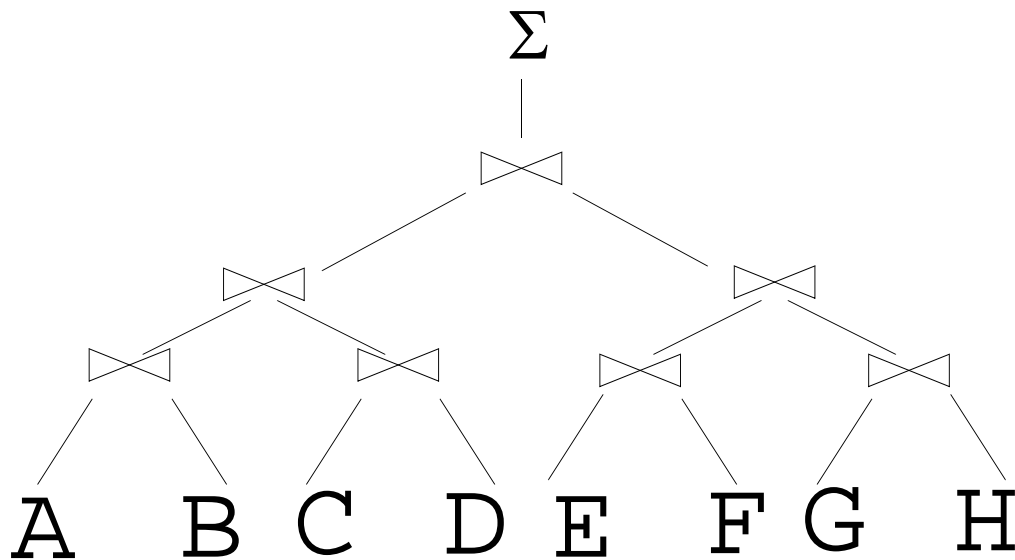


Can SMS join R and S, but rather useless for guessing the final query answer...

- Output tuples not produced randomly, so no subsequent use
- Can do 3+ tables in one SMS join, but only if can use the same sort order for each

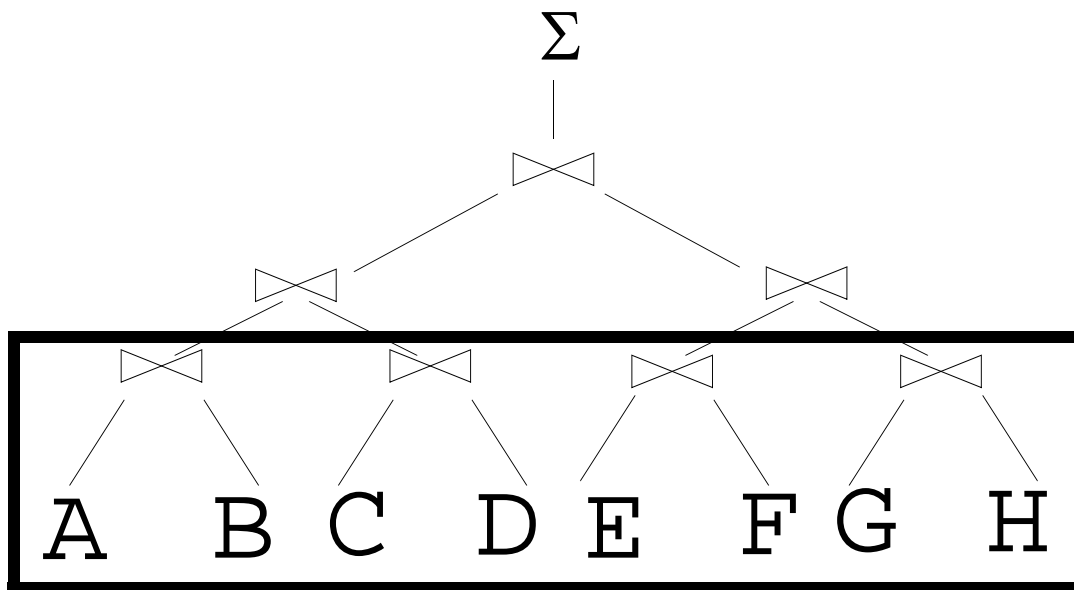
Basic Query Processing in DBO

Imagine a complex query plan:



Basic Query Processing in DBO

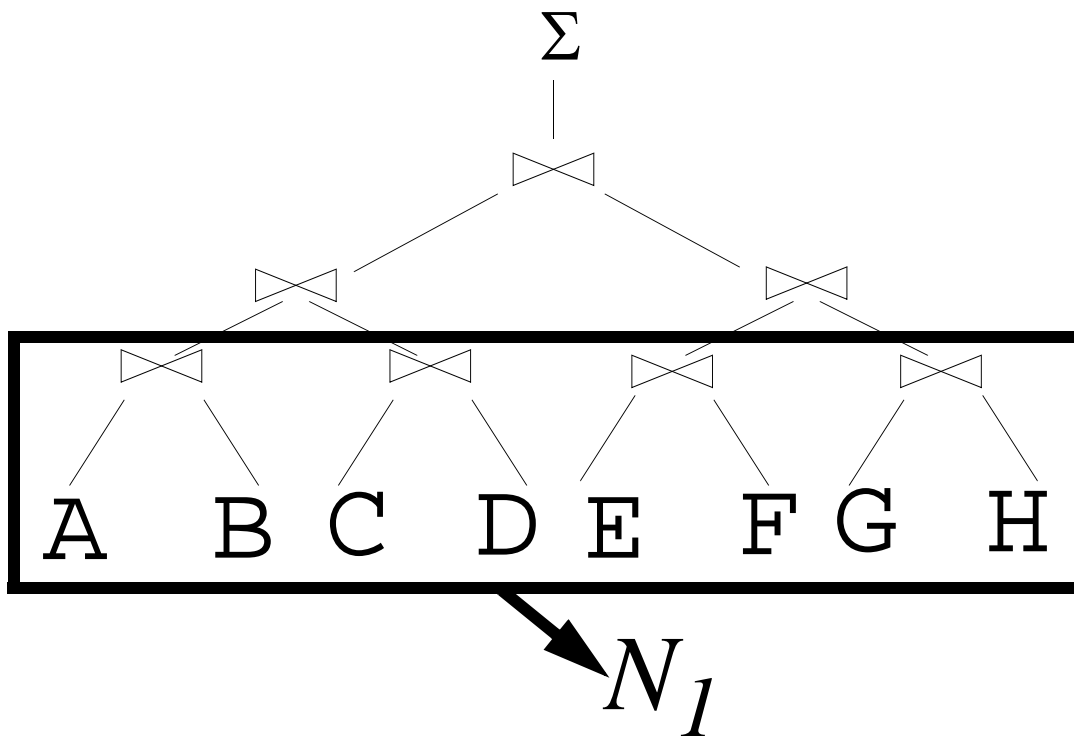
To process the plan:



First, all bottom-level joins are processed concurrently, kind of like one huge SMS sort phase; this is called a *levelwise step*

Basic Query Processing in DBO

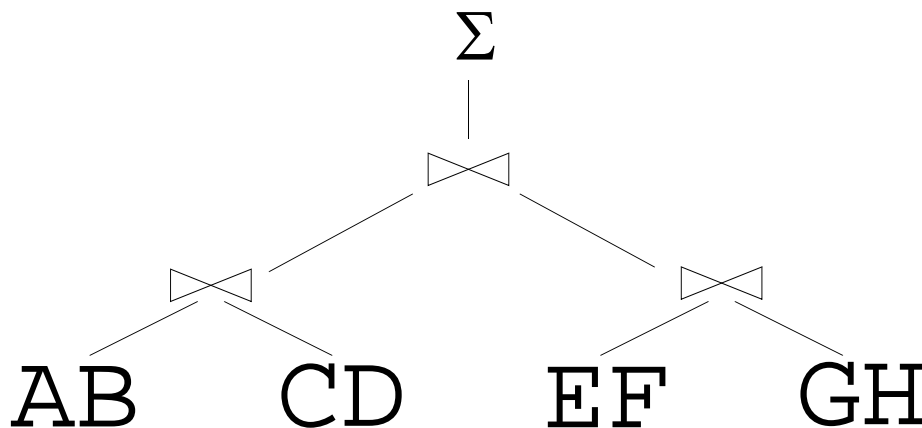
To process the plan:



The first levelwise step (at all times) checks for “lucky” answer tuples in order to maintain an online guess at the query result

Basic Query Processing in DBO

To process the plan:

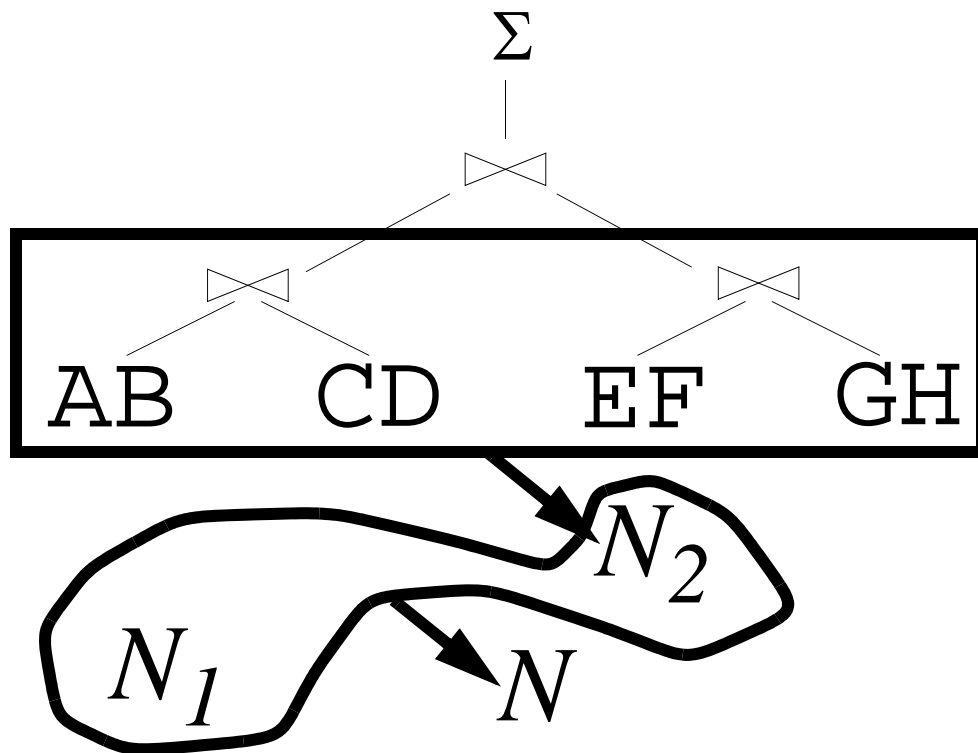


N_1

Eventually the first levelwise step completes all of the “lowest” joins, and N_1 is frozen

Basic Query Processing in DBO

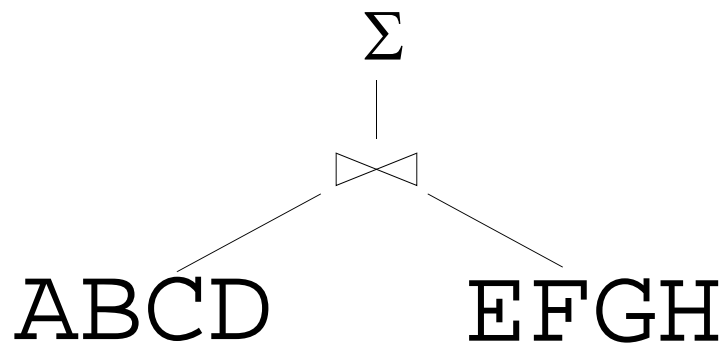
To process the plan:



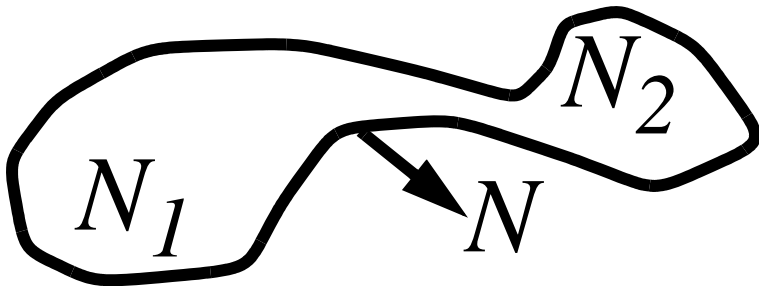
Then levelwise step two begins... N_2 will soon be far more accurate than N_1 ; they are combined to give the user an estimate

Basic Query Processing in DBO

To process the plan:

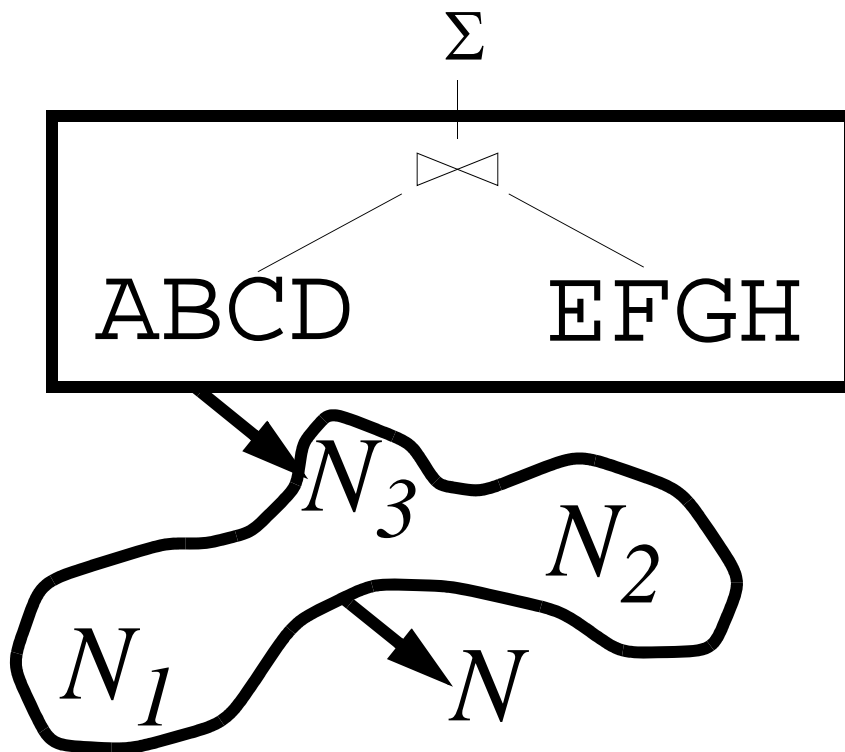


Eventually levelwise
step two finishes and N_2
is frozen as well.



Basic Query Processing in DBO

To process the plan:



Then levelwise step three begins... as it proceeds, the variance of N_3 goes to zero

Basic Query Processing in DBO

To process the plan:

Σ
|
ABCDEFGH

When levelwise step three finishes, the query has completed execution.

Details, Details, Details...

Each levelwise step has two phases: the *scan phase* and the *merge phase*

The scan phase is analogous to the sort phase of a SMJ, except:

1. There is one scan phase for all joins at a levelwise step
2. Any “answer tuples” discovered are used to update N_i
3. Round-robin processing of runs
4. Makes use of a randomized sort order (provided by H)

Scan Phase in Detail

To concurrent join (R_1 and R_2), (R_3 and R_4)

WHERE $R_1 \cdot B = R_2 \cdot C$

AND $R_2 \cdot E = R_3 \cdot F$

AND $R_3 \cdot G = R_4 \cdot H$

R_1			R_2				R_3			R_4		
A	B		C	D	E		F	G	H	I		
8	7	4	7	1	5	4	1	1	9	8	9	5
1	0	5	0	5	4	5	7	2	0	7	0	4
9	7	4	4	3	1	0	2	5	2	1	4	9
2	9	8	6	4	9	3	5	0	1	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Scan Phase in Detail

To concurrent join (R_1 and R_2), (R_3 and R_4)

Read one run
from *each relation* in the level-
wise step into
memory

R_1			R_2				R_3			R_4		
A	B		C	D	E		F	G	H	I		
8	7	4	7	1	5	4	1	1	9	8	9	5
1	0	5	0	5	4	5	7	2	0	7	0	4
9	7	4	4	3	1	0	2	5	2	1	4	9
2	9	8	6	4	9	3	5	0	1	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Scan Phase in Detail

To concurrent join $(R_1 \text{ and } R_2)$, $(R_3 \text{ and } R_4)$

Any output tuples are immediately discovered; used to produce N_i (unbiased!)

Result tuple discovered on-the-fly

R ₁			R ₂				R ₃			R ₄		
A	B		C	D	E		F	G		H	I	
8	7	4	7	1	5	4	1	1	9	8	9	5
1	0	5	0	5	4	5	7	2	0	7	0	4
9	7	4	4	3	1	0	2	5	2	1	4	9
2	9	8	6	4	9	3	5	0	1	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0

$H_1(B) \uparrow$ $H_1(C) \uparrow$ $H_2(G) \uparrow$ $H_2(H) \uparrow$

Scan Phase in Detail

To concurrent join $(R_1 \text{ and } R_2)$, $(R_3 \text{ and } R_4)$

Sort run from first relation on $H_1(B)$; write back to disk

Sorted on $H_1(B)$, written out to disk

R ₁			R ₂				R ₃		R ₄			
A	B		C	D	E	F	G	H	I			
8	7	4	7	1	5	4	1	1	9	8	9	5
9	7	4	0	5	4	5	7	2	0	7	0	4
1	0	5	4	3	1	0	2	5	2	1	4	9
2	9	8	6	4	9	3	5	0	1	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$		$H_2(H) \uparrow$			

Scan Phase in Detail

To concurrent join (R_1 and R_2), (R_3 and R_4)

Load second run
from first relation

R_1			R_2				R_3			R_4		
A	B		C	D	E		F	G		H	I	
8	7	4	7	1	5	4	1	1	9	8	9	5
9	7	4	0	5	4	5	7	2	0	7	0	4
1	0	5	4	3	1	0	2	5	2	1	4	9
2	9	8	6	4	9	3	5	0	1	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Scan Phase in Detail

To concurrent join $(R_1 \text{ and } R_2)$, $(R_3 \text{ and } R_4)$

Use any discovered result tuples to update estimate (still unbiased!)

R ₁			R ₂				R ₃			R ₄		
A	B		C	D	E	F	G		H	I		
8	7	4	7	1	5	4	1	1	9	8	9	5
9	7	4	0	5	4	5	7	2	0	7	0	4
1	0	5	4	3	1	0	2	5	2	1	4	9
2	9	8	6	4	9	3	5	0	1	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0

$H_1(B) \uparrow$ $H_1(C) \uparrow$ $H_2(G) \uparrow$ $H_2(H) \uparrow$

Scan Phase in Detail

To concurrent join (R_1 and R_2), (R_3 and R_4)

Sort run on $H_1(C)$
and write it back
to disk

R_1			R_2				R_3			R_4		
A	B		C	D	E		F	G		H	I	
8	7	4	4	3	1	0	1	1	9	8	9	5
9	7	4	6	4	9	3	7	2	0	7	0	4
1	0	5	7	1	5	4	2	5	2	1	4	9
2	9	8	0	5	4	5	5	0	1	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Scan Phase in Detail

To concurrent join (R_1 and R_2), (R_3 and R_4)

Load next run
from second rela-
tion into memory

R_1			R_2				R_3			R_4		
A	B		C	D	E		F	G		H	I	
8	7	4	4	3	1	0	1	1	9	8	9	5
9	7	4	6	4	9	3	7	2	0	7	0	4
1	0	5	7	1	5	4	2	5	2	1	4	9
2	9	8	0	5	4	5	5	0	1	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Scan Phase in Detail

To concurrent join $(R_1 \text{ and } R_2)$, $(R_3 \text{ and } R_4)$

No output tuples
discovered, so
update estimate
and write next run
back to disk

R_1			R_2				R_3			R_4		
A	B		C	D	E		F	G		H	I	
8	7	4	4	3	1	0	7	2	0	8	9	5
9	7	4	6	4	9	3	5	0	1	7	0	4
1	0	5	7	1	5	4	2	5	2	1	4	9
2	9	8	0	5	4	5	1	1	9	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Scan Phase in Detail

To concurrent join (R_1 and R_2), (R_3 and R_4)

Load next run
into memory

R_1			R_2				R_3			R_4			
A	B		C	D	E			F	G		H	I	
8	7	4	4	3	1	0				8	9	5	
9	7	4	6	4	9	3				7	0	4	
1	0	5	7	1	5	4				1	4	9	
2	9	8	0	5	4	5				0	5	2	
5	5	9	1	7	2	7				0	2	1	
2	4	0	8	2	7	1				3	7	3	
9	9	8	5	8	3	9				4	9	8	
4	8	1	9	1	8	8				2	5	0	
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$			

Scan Phase in Detail

To concurrent join $(R_1 \text{ and } R_2)$, $(R_3 \text{ and } R_4)$

Immediately discover any output tuples and update the estimate

R ₁			R ₂				R ₃			R ₄		
A	B		C	D	E		F	G		H	I	
8	7	4	4	3	1	0	7	2	0	8	9	5
9	7	4	6	4	9	3	5	0	1	7	0	4
1	0	5	7	1	5	4	2	5	2	1	4	9
2	9	8	0	5	4	5	1	1	9	0	5	2
5	5	9	1	7	2	7	8	0	1	0	2	1
2	4	0	8	2	7	1	6	9	7	3	7	3
9	9	8	5	8	3	9	4	3	3	4	9	8
4	8	1	9	1	8	8	3	4	8	2	5	0

$H_1(B) \uparrow$ $H_1(C) \uparrow$ $H_2(G) \uparrow$ $H_2(H) \uparrow$

Scan Phase in Detail

To concurrent join (R_1 and R_2), (R_3 and R_4)

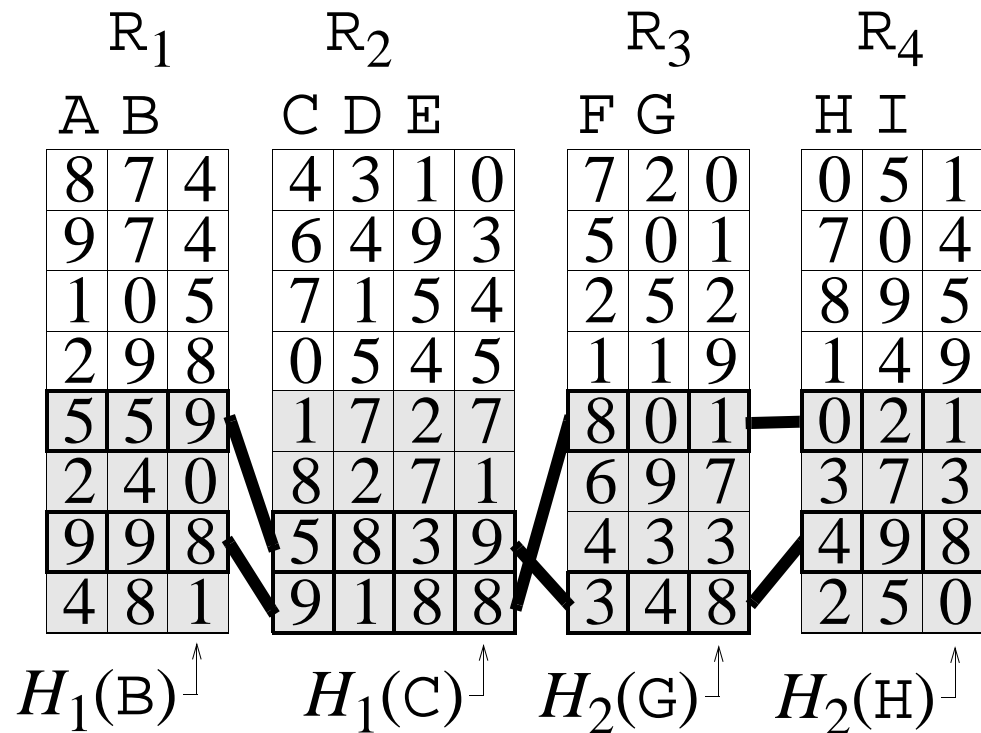
Sort and write
back to disk

R_1			R_2				R_3			R_4			
A	B		C	D	E			F	G		H	I	
8	7	4	4	3	1	0		7	2	0	0	5	1
9	7	4	6	4	9	3		5	0	1	7	0	4
1	0	5	7	1	5	4		2	5	2	8	9	5
2	9	8	0	5	4	5		1	1	9	1	4	9
5	5	9	1	7	2	7		8	0	1	0	2	1
2	4	0	8	2	7	1		6	9	7	3	7	3
9	9	8	5	8	3	9		4	3	3	4	9	8
4	8	1	9	1	8	8		3	4	8	2	5	0
$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$			

Scan Phase in Detail

To concurrent join $(R_1 \text{ and } R_2)$, $(R_3 \text{ and } R_4)$

Do the same for
the fourth relation



Scan Phase in Detail

To concurrent join (R_1 and R_2), (R_3 and R_4)

All tuples have
been read, so
write back all in-
memory tuples,
and we're done!

	R_1			R_2				R_3			R_4		
	A	B		C	D	E		F	G		H	I	
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8
	$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Key Points

- (In this example) Five different updates to N_i
- One update every time a run is processed
- N_i is unbiased assuming random input order
- Can characterize variance (very challenging!)
- Ready for merge phase...

Merge Phase

- Separate merge for each join in i th levelwise step
- Each merge a lot like merge phase of SMJ
- Recall that we use a random sort order
- So result tuples come out in (semi-) random order
- Output tuples pipelined directly into scan phase of *next* levelwise step
- Since tuples produced randomly, $(i + 1)$ th levelwise step valid

Merge Phase in Detail

	R ₁			R ₂				R ₃			R ₄		
	A	B		C	D	E		F	G		H	I	
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8
	$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Output so far:

R₁₂

R₃₄

Merge Phase in Detail

Head of each run of each relation is read into memory...

	R ₁			R ₂				R ₃			R ₄			Output so far:						
	A	B		C	D	E			F	G		H	I		R ₁₂			R ₃₄		
Run 1	8	7	4	4	3	1	0		7	2	0	0	5	1						
	9	7	4	6	4	9	3		5	0	1	7	0	4						
	1	0	5	7	1	5	4		2	5	2	8	9	5						
	2	9	8	0	5	4	5		1	1	9	1	4	9						
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0							
	4	8	1	1	7	2	7	4	3	3	0	2	1							
	9	9	8	9	1	8	8	6	9	7	3	7	3							
	5	5	9	5	8	3	9	3	4	8	4	9	8							
	$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$									

Merge Phase in Detail

Search for output tuples in first join...

	R ₁			R ₂				R ₃			R ₄		
	A	B		C	D	E		F	G		H	I	
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8

$H_1(B) \uparrow$ $H_1(C) \uparrow$ $H_2(G) \uparrow$ $H_2(H) \uparrow$
 Joined through key 1

Output so far:

R ₁₂					R ₃₄				
A	B	C	D	E	A	B	C	D	E
2	4	4	3	1					
4	8	8	2	7					

pipelined directly into next levelwise step

Merge Phase in Detail

Search for output tuples in second join...

	R ₁			R ₂				R ₃			R ₄		
	A	B		C	D	E		F	G		H	I	
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8
	$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		
	Joined through key 1						Joined through key 1						

Output so far:

R ₁₂					R ₃₄			
A	B	C	D	E	F	G	H	I
2	4	4	3	1	7	2	2	5
4	8	8	2	7	5	0	0	5
					5	0	0	2
					8	0	0	5
					8	0	0	2

Merge Phase in Detail

Exhausted run-heads are replaced in memory...

	R ₁			R ₂				R ₃			R ₄		
	A	B		C	D	E		F	G		H	I	
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8
	$H_1(B) \downarrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		

Joined through key 1 Joined through key 1

Output so far:

R ₁₂					R ₃₄			
A	B	C	D	E	F	G	H	I
2	4	4	3	1	7	2	2	5
4	8	8	2	7	5	0	0	5
					5	0	0	2
					8	0	0	5
					8	0	0	2

Merge Phase in Detail

New output tuples are written to output...

	R ₁			R ₂				R ₃			R ₄		
	A	B		C	D	E		F	G		H	I	
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8

$H_1(B) \uparrow$ $H_1(C) \uparrow$ $H_2(G) \uparrow$ $H_2(H) \uparrow$
 Joined through key 3 Joined through key 3

Output so far:

R ₁₂					R ₃₄			
A	B	C	D	E	F	G	H	I
2	4	4	3	1	7	2	2	5
4	8	8	2	7	5	0	0	5
					5	0	0	2
					8	0	0	5
					8	0	0	2
					4	3	3	7

Merge Phase in Detail

Exhausted runs are replenished, new output tuples discovered...

	R ₁			R ₂				R ₃		R ₄			
	A	B		C	D	E		F	G	H	I		
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8

$H_1(B) \uparrow$ $H_1(C) \uparrow$ $H_2(G) \uparrow$ $H_2(H) \uparrow$
 Joined through key 4 Joined through key 3

Output so far:

R ₁₂					R ₃₄			
A	B	C	D	E	F	G	H	I
2	4	4	3	1	7	2	2	5
4	8	8	2	7	5	0	0	5
8	7	7	1	5	5	0	0	2
9	7	7	1	5	8	0	0	5
					8	0	0	2
					4	3	3	7

Merge Phase in Detail

Again replace processed tuples in memory (perhaps skipping some steps!)

	R ₁			R ₂			R ₃		R ₄				
	A	B		C	D	E	F	G	H	I			
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8
	$H_1(B) \uparrow$			$H_1(C) \uparrow$			$H_2(G) \uparrow$		$H_2(H) \uparrow$				
	Joined through key 4						Joined through key 3						

Output so far:

R ₁₂					R ₃₄			
A	B	C	D	E	F	G	H	I
2	4	4	3	1	7	2	2	5
4	8	8	2	7	5	0	0	5
8	7	7	1	5	5	0	0	2
9	7	7	1	5	8	0	0	5
					8	0	0	2
					4	3	3	7

Merge Phase in Detail

New output tuples from first relation are written to output...

	R ₁			R ₂				R ₃			R ₄		
	A	B		C	D	E		F	G		H	I	
	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
Run 1	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8
	$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		
	Done!							Joined through key 3					

Output so far:

R ₁₂					R ₃₄			
A	B	C	D	E	F	G	H	I
2	4	4	3	1	7	2	2	5
4	8	8	2	7	5	0	0	5
8	7	7	1	5	5	0	0	2
9	7	7	1	5	8	0	0	5
1	0	0	5	4	8	0	0	2
2	9	9	1	8	4	3	3	7
9	9	9	1	8				
5	5	5	8	3				

Merge Phase in Detail

Same for second relation, and we are done!

	R ₁			R ₂				R ₃			R ₄		
	A	B		C	D	E		F	G		H	I	
Run 1	8	7	4	4	3	1	0	7	2	0	0	5	1
	9	7	4	6	4	9	3	5	0	1	7	0	4
	1	0	5	7	1	5	4	2	5	2	8	9	5
	2	9	8	0	5	4	5	1	1	9	1	4	9
Run 2	2	4	0	8	2	7	1	8	0	1	2	5	0
	4	8	1	1	7	2	7	4	3	3	0	2	1
	9	9	8	9	1	8	8	6	9	7	3	7	3
	5	5	9	5	8	3	9	3	4	8	4	9	8
	$H_1(B) \uparrow$			$H_1(C) \uparrow$				$H_2(G) \uparrow$			$H_2(H) \uparrow$		
	Done!							Done!					

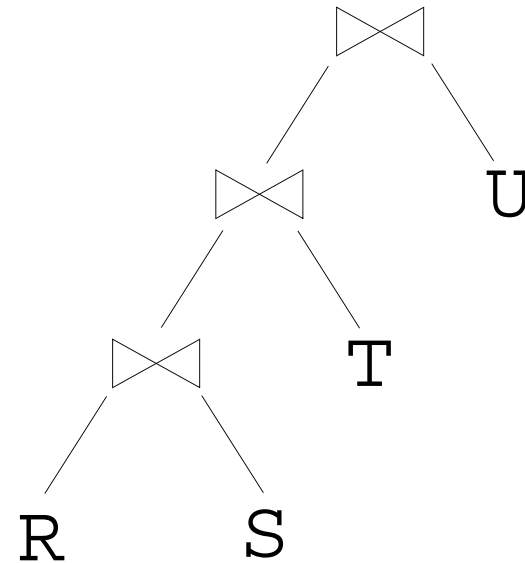
Output so far:

R ₁₂					R ₃₄			
A	B	C	D	E	F	G	H	I
2	4	4	3	1	7	2	2	5
4	8	8	2	7	5	0	0	5
8	7	7	1	5	5	0	0	2
9	7	7	1	5	8	0	0	5
1	0	0	5	4	8	0	0	2
2	9	9	1	8	4	3	3	7
9	9	9	1	8	3	4	4	9
5	5	5	8	3	1	1	1	4

A Few More Details...

Handling “inconvenient” queries

```
SELECT SUM (R.A)
FROM R, S, T, U
WHERE R.A = S.A AND
      R.B = T.B AND
      R.C = U.C
```

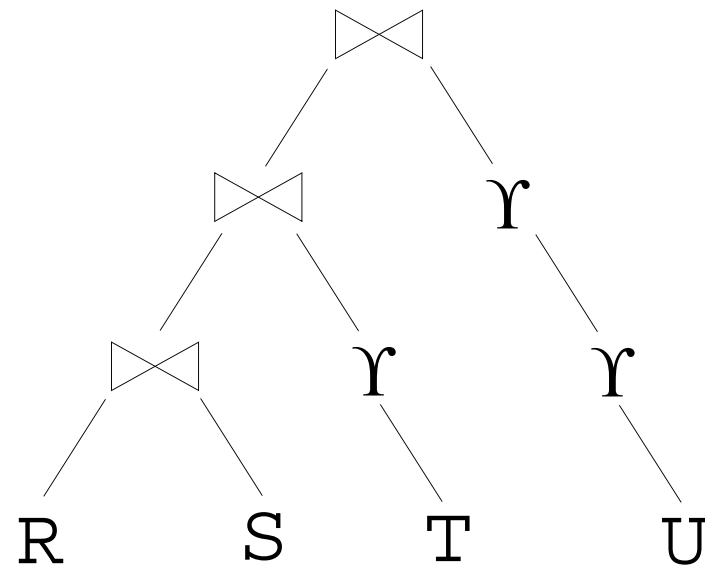


Issue: which relations are in the first levelwise step?

A Few More Details...

Handling “inconvenient” queries

```
SELECT SUM (R.A)
FROM R, S, T, U
WHERE R.A = S.A AND
      R.B = T.B AND
      R.C = U.C
```



Solution: introduce a “scan and reorder” operator

A Few More Details...

Computing the current estimate

- Use $N = \sum w_i N_i$; since independent, getting w_i is an easy optimization problem

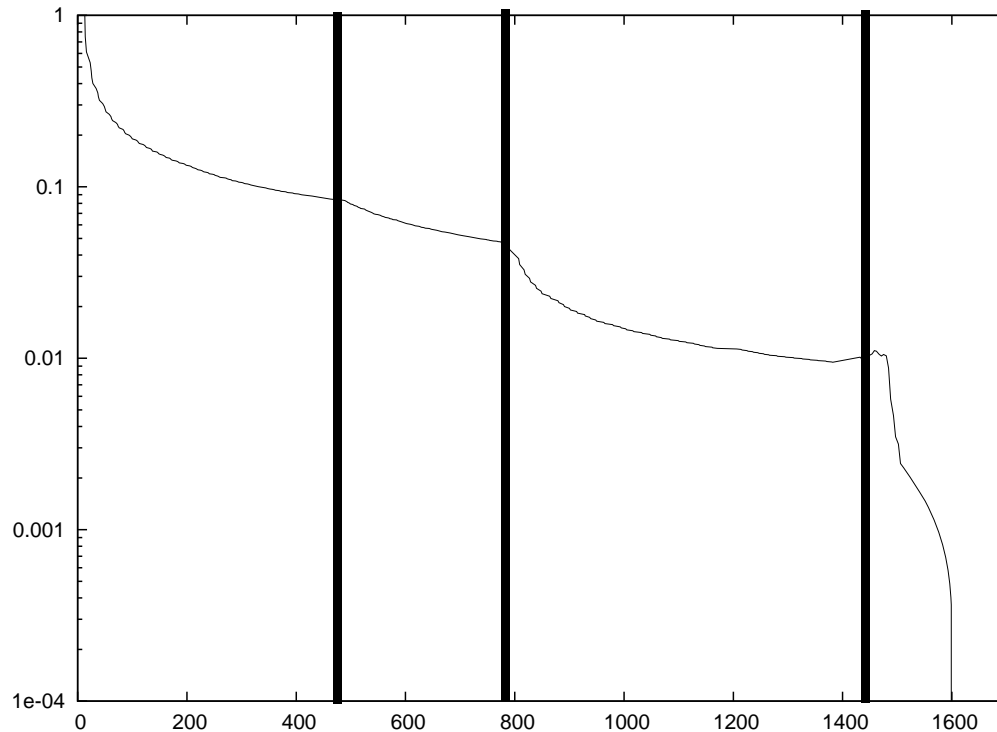
Computing accuracy guarantees

- Not easy!
 1. Two kinds of randomization
 2. “Chunks” of tuples from merge phase
 3. LOTS of algebra, recursive variance formulas
 4. Covariance among estimates making up each N_i

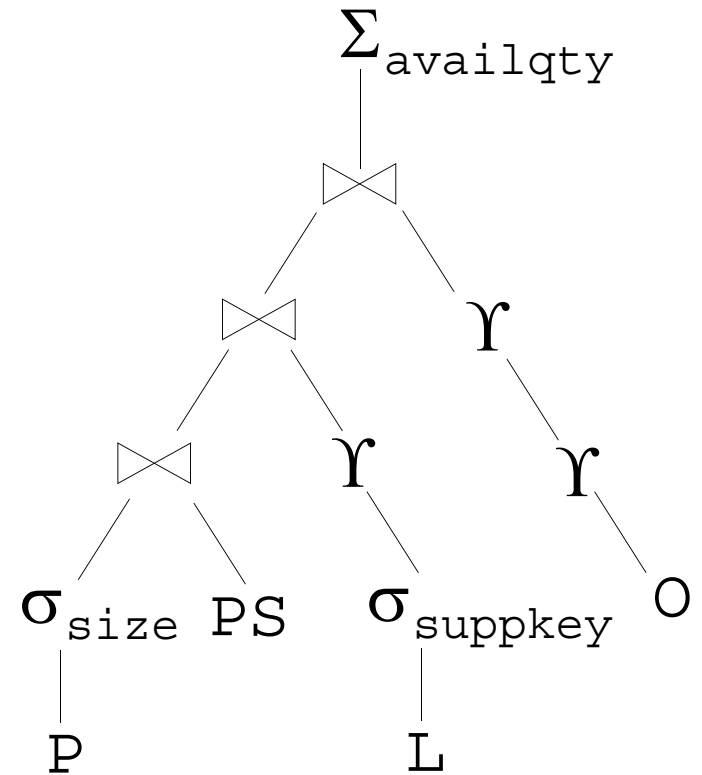
How Well Does This Work?

Ex: TPC-H query: DBO 26m42s
Postgres 43m47s

Bound Width



Query Duration



Are We Done Yet?

Of Course Not!

- Subtraction and related operators are a big problem
- Indexing (how to provide randomness?)
- Query optimization: what is the goal?
- Many more...

Thank You!

- Special thanks to my UF colleague Alin Dobra, and my students Abhijit Pol, Subramanian Arumugam, Shantanu Joshi
- Some papers on this stuff:
 1. “A Disk-Based Join with Probabilistic Guarantees,” Jermaine, Dobra, Arumugam, Joshi, Pol, SIGMOD 2005
 2. “The Sort-Merge-Shrink Join,” Jermaine, Dobra, Arumugam, Joshi, Pol, TODS 31(4), December 2006