

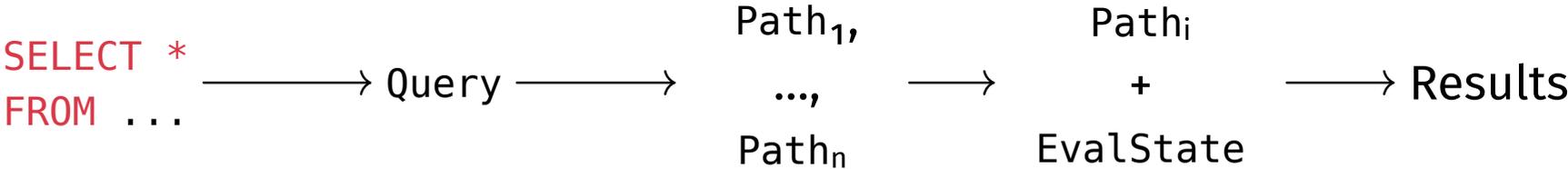
COST-BASED OPTIMIZATION

CSE 4/562: Database Systems | Lecture 11

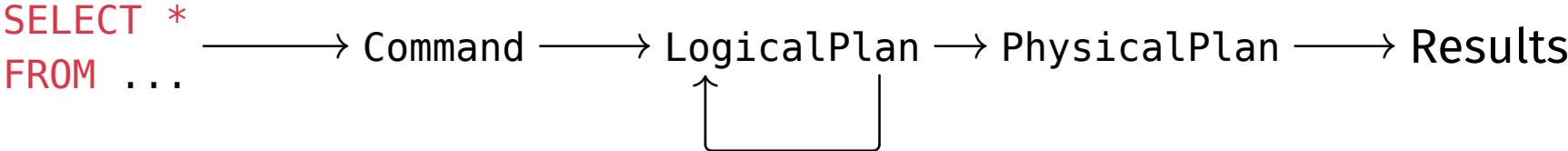
DB. Sys.: T.C.B.: Ch. 16

Quiz

Postgresql



Apache Spark



Input

- $\text{Path}_1, \dots, \text{Path}_n$
- Previously collected statistics

Output

- $\text{argmin}_i \text{Cost}(\text{Path}_i)$

We need a way to estimate $\text{Cost}(\text{Path}_i)$

Idea 1: $\text{Eval}(\text{Plan}_i)$

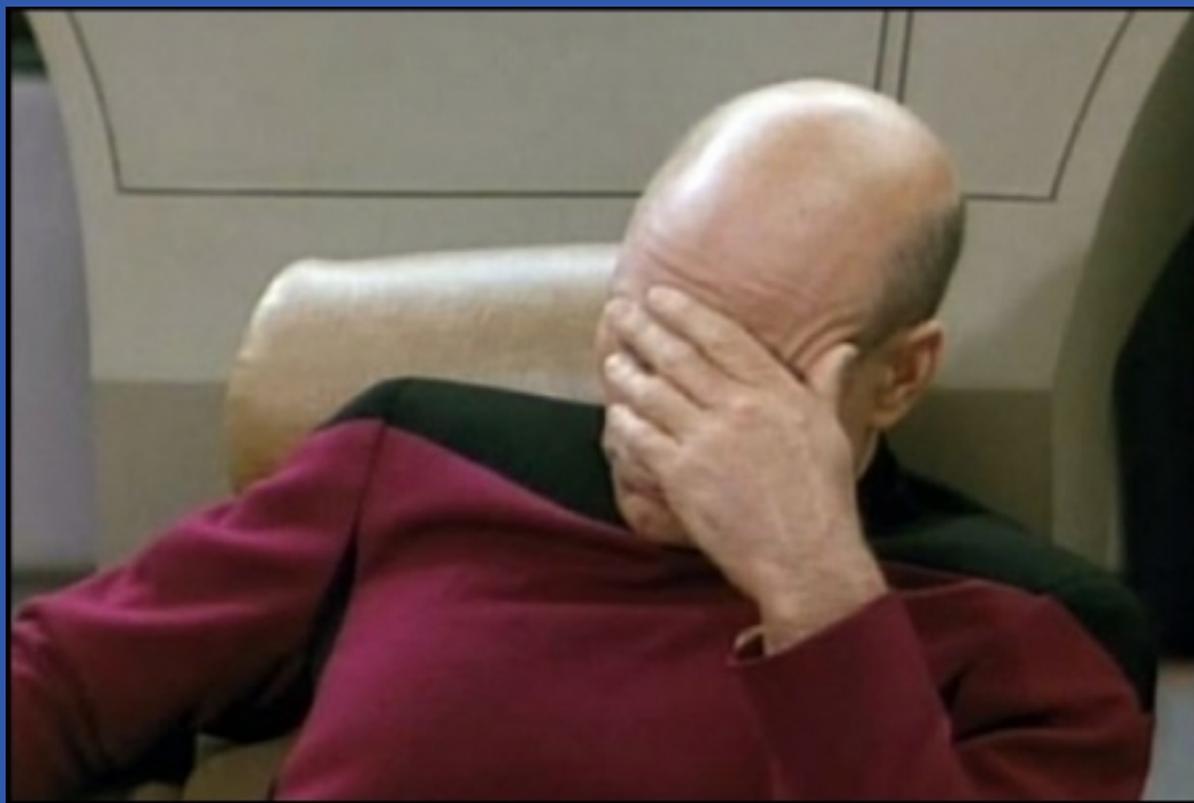


Image credit: © Paramount Pictures

If we can't get the exact cost of a plan, what can we do?

Idea 2: Run each plan on a small sample of the data.

Idea 3: Analytically estimate the cost of a plan.

Idea 2: Run each plan on a small sample of the data.

Idea 3: Analytically estimate the cost of a plan.

CPU Time

- How much time is spent processing data?

IO Cost

- How many pages are read/written?
- How much spatial locality is the workload?

Memory Cost

- What is the smallest amount of memory required (working set size)

CPU Time

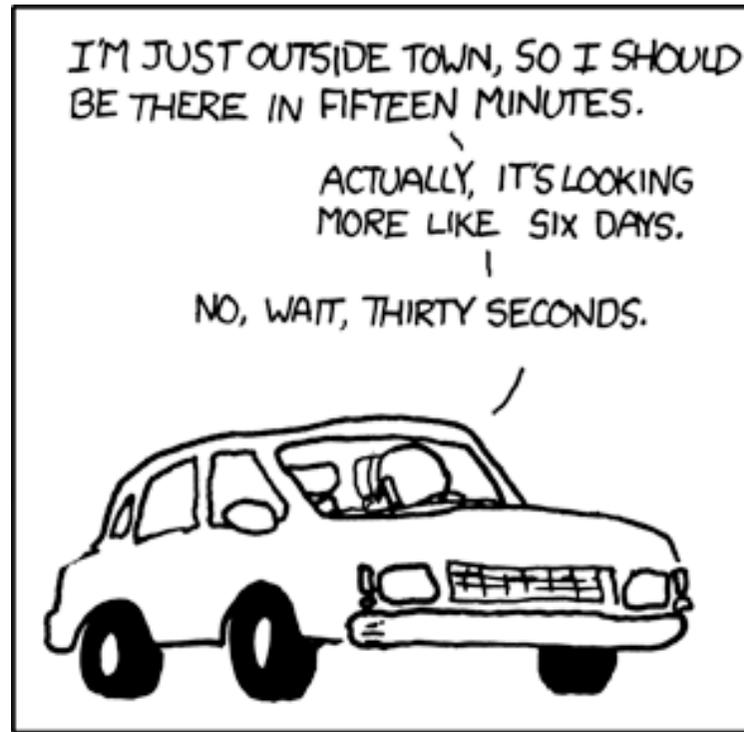
- How much time is spent processing data?

IO Cost

- How many pages are read/written?
- How much spatial locality is the workload?

Memory Cost

- What is the smallest amount of memory required (working set size)



THE AUTHOR OF THE WINDOWS FILE COPY DIALOG VISITS SOME FRIENDS.

Remember the real goals

1. Accurately **rank** the plans.
2. Don't spend more time optimizing than you get back.
3. Don't pick a plan that uses more memory than you have.

- $|R|$: The number of data pages in R
- B : Pages of buffer
- K : Keys per index page

Table Scan

Table Scan

$$IO(R) = |R|$$

Table Scan

$$\text{IO}(R) = |R|$$

$$\text{Mem}(R) = O(1)$$

Projection

Projection

$$\text{IO}(\pi R) = 0 \text{ (+ IO}(R))$$

Projection

$$\text{IO}(\pi R) = 0 \text{ (+ IO}(R))$$

$$\text{Mem}(\pi R) = O(1)$$

Selection

Selection

$$\text{IO}(\sigma R) = 0 \text{ (+ IO}(R))$$

Selection

$$\text{IO}(\sigma R) = 0 \text{ (+ IO}(R))$$

$$\text{Mem}(\sigma R) = O(1)$$

Union

Union

$$IO(R \cup S) = O(+ IO(R) + IO(S))$$

Union

$$\text{IO}(R \cup S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \cup S) = O(1)$$

Sort

Sort

In Memory

Sort

In Memory

$$\text{IO}(\tau R) = 0 \text{ (+ IO}(R))$$

Sort

In Memory

$$\text{IO}(\tau R) = 0 \text{ (+ IO}(R))$$

$$\text{Mem}(\tau R) = O(|R|)$$

Sort

In Memory

$$\text{IO}(\tau R) = 0 \text{ (+ IO}(R)\text{)}$$

$$\text{Mem}(\tau R) = O(|R|)$$

On Disk

Sort

In Memory

$$\text{IO}(\tau R) = 0 \text{ (+ IO}(R))$$

$$\text{Mem}(\tau R) = O(|R|)$$

On Disk

$$\text{IO}(\tau R) = 2 \cdot \lfloor \log_B(|R|) \rfloor \text{ (+ IO}(R))$$

Sort

In Memory

$$\text{IO}(\tau R) = 0 \text{ (+ IO}(R))$$

$$\text{Mem}(\tau R) = O(|R|)$$

On Disk

$$\text{IO}(\tau R) = 2 \cdot \lfloor \log_B(|R|) \rfloor \text{ (+ IO}(R))$$

$$\text{Mem}(\tau R) = O(B)$$

Index Scan;

B+ Tree

Index Scan;

B+ Tree

$$\text{IO}(\text{IndexScan}(R, \theta)) = \log_K(|R|) + |\sigma_c(R)|$$

Index Scan;

B+ Tree

$$\text{IO}(\text{IndexScan}(R, \theta)) = \log_K(|R|) + |\sigma_c(R)|$$

$$\text{Mem}(\text{IndexScan}(R, \theta)) = O(1)$$

Index Scan;

B+ Tree

$$\text{IO}(\text{IndexScan}(R, \theta)) = \log_K(|R|) + |\sigma_c(R)|$$

$$\text{Mem}(\text{IndexScan}(R, \theta)) = O(1)$$

Hash

Index Scan;

B+ Tree

$$\text{IO}(\text{IndexScan}(R, \theta)) = \log_K(|R|) + |\sigma_c(R)|$$

$$\text{Mem}(\text{IndexScan}(R, \theta)) = O(1)$$

Hash

$$\text{IO}(\text{IndexScan}(R, \theta)) = 1$$

Index Scan;

B+ Tree

$$\text{IO}(\text{IndexScan}(R, \theta)) = \log_K(|R|) + |\sigma_c(R)|$$

$$\text{Mem}(\text{IndexScan}(R, \theta)) = O(1)$$

Hash

$$\text{IO}(\text{IndexScan}(R, \theta)) = 1$$

$$\text{Mem}(\text{IndexScan}(R, \theta)) = O(1)$$

Block Nested Loop Join

Buffer S in Memory

Block Nested Loop Join

Buffer S in Memory

$$IO(R \times S) = 0 \text{ (+ } IO(R) \text{ + } IO(S))$$

Block Nested Loop Join

Buffer S in Memory

$$\text{IO}(R \times S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(|S|)$$

Block Nested Loop Join

Buffer S in Memory

$$\text{IO}(R \times S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(|S|)$$

Buffer S on Disk

Block Nested Loop Join

Buffer S in Memory

$$\text{IO}(R \times S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(|S|)$$

Buffer S on Disk

$$\text{IO}(R \times S) = \left(1 + \frac{|R|}{B}\right) \cdot |S| \text{ (+ IO}(R) \text{ + IO}(S))$$

Block Nested Loop Join

Buffer S in Memory

$$\text{IO}(R \times S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(|S|)$$

Buffer S on Disk

$$\text{IO}(R \times S) = \left(1 + \frac{|R|}{B}\right) \cdot |S| \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(B)$$

Block Nested Loop Join

Buffer S in Memory

$$IO(R \times S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(|S|)$$

Buffer S on Disk

$$IO(R \times S) = \left(1 + \frac{|R|}{B}\right) \cdot |S| \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(B)$$

Rewind S

Block Nested Loop Join

Buffer S in Memory

$$\text{IO}(R \times S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(|S|)$$

Buffer S on Disk

$$\text{IO}(R \times S) = \left(1 + \frac{|R|}{B}\right) \cdot |S| \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(B)$$

Rewind S

$$\text{IO}(R \times S) = \frac{|R|}{B} \cdot \text{IO}(S) \text{ (+ IO}(R))$$

Block Nested Loop Join

Buffer S in Memory

$$\text{IO}(R \times S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(|S|)$$

Buffer S on Disk

$$\text{IO}(R \times S) = \left(1 + \frac{|R|}{B}\right) \cdot |S| \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \times S) = O(B)$$

Rewind S

$$\text{IO}(R \times S) = \frac{|R|}{B} \cdot \text{IO}(S) \text{ (+ IO}(R))$$

$$\text{Mem}(R \times S) = O(B)$$

Hash Join

1-Pass (In-Memory)

Hash Join

1-Pass (In-Memory)

$$\text{IO}(R \bowtie S) = 0 \text{ (+ IO}(R) + \text{IO}(S))$$

Hash Join

1-Pass (In-Memory)

$$\text{IO}(R \bowtie S) = 0 \text{ (+ IO}(R) \text{ + IO}(S))$$

$$\text{Mem}(R \bowtie S) = O(|S|)$$

Hash Join

1-Pass (In-Memory)

$$\text{IO}(R \bowtie S) = 0 \text{ (+ IO}(R) + \text{IO}(S))$$

$$\text{Mem}(R \bowtie S) = O(|S|)$$

2-Pass (On Disk)

Hash Join

1-Pass (In-Memory)

$$\text{IO}(R \bowtie S) = 0 \text{ (+ IO}(R) + \text{IO}(S))$$

$$\text{Mem}(R \bowtie S) = O(|S|)$$

2-Pass (On Disk)

$$\text{IO}(R \bowtie S) = 2|R| + 2|S| \text{ (+ IO}(R) + \text{IO}(S))$$

Hash Join

1-Pass (In-Memory)

$$\text{IO}(R \bowtie S) = 0 \text{ (+ IO}(R) + \text{IO}(S))$$

$$\text{Mem}(R \bowtie S) = O(|S|)$$

2-Pass (On Disk)

$$\text{IO}(R \bowtie S) = 2|R| + 2|S| \text{ (+ IO}(R) + \text{IO}(S))$$

$$\text{Mem}(R \bowtie S) = O\left(\frac{|R| + |S|}{\text{\#Buckets}}\right)$$

Sort Merge Join

Sort Merge Join

$\text{IO}(R \bowtie S) = [\text{As Sort}]$

$\text{Mem}(R \bowtie S) = [\text{As Sort}]$

Aggregate

In Memory

Aggregate

In Memory

$$\text{IO}\left(\sum_A R\right) = 0 \text{ (+ IO}(R))$$

Aggregate

In Memory

$$\text{IO}\left(\sum_A R\right) = 0 \text{ (+ IO}(R))$$

$$\text{Mem}\left(\sum_A R\right) = O(\text{domain}(A))$$

Aggregate

In Memory

$$\text{IO}\left(\sum_A R\right) = 0 \text{ (+ IO}(R))$$

$$\text{Mem}\left(\sum_A R\right) = O(\text{domain}(A))$$

Sort / Merge

Aggregate

In Memory

$$\text{IO}\left(\sum_A R\right) = 0 \text{ (+ IO}(R))$$

$$\text{Mem}\left(\sum_A R\right) = O(\text{domain}(A))$$

Sort / Merge

$$\text{IO}\left(\sum_A R\right) = [\text{As Sort}]$$

$$\text{Mem}\left(\sum_A R\right) = [\text{As Sort}]$$

$|R|$: **Size of R**

B : **Pages of Buffer**

K : **Keys per Index Page**

$\text{domain}(A)$: **Number of distinct values of A**

$|R|$: **Size of R**

- Precomputed if R is a table
- ??? otherwise

B : **Pages of Buffer**

K : **Keys per Index Page**

$\text{domain}(A)$: **Number of distinct values of A**

$|R|$: **Size of R**

- Precomputed if R is a table
- ??? otherwise

B : **Pages of Buffer**

- Configurable Parameter

K : **Keys per Index Page**

$\text{domain}(A)$: **Number of distinct values of A**

$|R|$: **Size of R**

- Precomputed if R is a table
- ??? otherwise

B : **Pages of Buffer**

- Configurable Parameter

K : **Keys per Index Page**

- Given Parameter

$\text{domain}(A)$: **Number of distinct values of A**

$|R|$: Size of R

- Precomputed if R is a table
- ??? otherwise

B : Pages of Buffer

- Configurable Parameter

K : Keys per Index Page

- Given Parameter

$\text{domain}(A)$: Number of distinct values of A

- Can be precomputed ($|\delta_{A(R)}|$)

Estimating IO requires being able to estimate $|Q|$ and $|\delta_A(Q)|$

Unlike estimating IOs, cardinality estimation doesn't care about the algorithm.
We can work with simple RA.

Operator	RA	Estimated Size
Table	R	
Projection	$\pi(Q)$	
Union	$Q_1 \cup Q_2$	
Cross Product	$Q_1 \times Q_2$	
Sort	$\tau(Q)$	
Limit	$L_N(Q)$	
Selection	$\sigma_\theta(Q)$	
Join	$Q_1 \bowtie_\theta Q_2$	
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	
Union	$Q_1 \cup Q_2$	
Cross Product	$Q_1 \times Q_2$	
Sort	$\tau(Q)$	
Limit	$L_N(Q)$	
Selection	$\sigma_\theta(Q)$	
Join	$Q_1 \bowtie_\theta Q_2$	
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	
Cross Product	$Q_1 \times Q_2$	
Sort	$\tau(Q)$	
Limit	$L_N(Q)$	
Selection	$\sigma_\theta(Q)$	
Join	$Q_1 \bowtie_\theta Q_2$	
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	$ Q_1 + Q_2 $
Cross Product	$Q_1 \times Q_2$	
Sort	$\tau(Q)$	
Limit	$L_N(Q)$	
Selection	$\sigma_\theta(Q)$	
Join	$Q_1 \bowtie_\theta Q_2$	
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	$ Q_1 + Q_2 $
Cross Product	$Q_1 \times Q_2$	$ Q_1 \times Q_2 $
Sort	$\tau(Q)$	
Limit	$L_N(Q)$	
Selection	$\sigma_\theta(Q)$	
Join	$Q_1 \bowtie_\theta Q_2$	
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	$ Q_1 + Q_2 $
Cross Product	$Q_1 \times Q_2$	$ Q_1 \times Q_2 $
Sort	$\tau(Q)$	$ Q $
Limit	$L_N(Q)$	
Selection	$\sigma_\theta(Q)$	
Join	$Q_1 \bowtie_\theta Q_2$	
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	$ Q_1 + Q_2 $
Cross Product	$Q_1 \times Q_2$	$ Q_1 \times Q_2 $
Sort	$\tau(Q)$	$ Q $
Limit	$L_N(Q)$	$\min(N, Q)$
Selection	$\sigma_\theta(Q)$	
Join	$Q_1 \bowtie_\theta Q_2$	
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	$ Q_1 + Q_2 $
Cross Product	$Q_1 \times Q_2$	$ Q_1 \times Q_2 $
Sort	$\tau(Q)$	$ Q $
Limit	$L_N(Q)$	$\min(N, Q)$
Selection	$\sigma_\theta(Q)$	$ Q \times \text{SEL}(\theta, Q)$
Join	$Q_1 \bowtie_\theta Q_2$	
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	$ Q_1 + Q_2 $
Cross Product	$Q_1 \times Q_2$	$ Q_1 \times Q_2 $
Sort	$\tau(Q)$	$ Q $
Limit	$L_N(Q)$	$\min(N, Q)$
Selection	$\sigma_\theta(Q)$	$ Q \times \text{SEL}(\theta, Q)$
Join	$Q_1 \bowtie_\theta Q_2$	$(Q_1 \times Q_2) \times \text{SEL}(\theta, Q_1 \times Q_2)$
Distinct	$\delta_A(Q)$	
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	$ Q_1 + Q_2 $
Cross Product	$Q_1 \times Q_2$	$ Q_1 \times Q_2 $
Sort	$\tau(Q)$	$ Q $
Limit	$L_N(Q)$	$\min(N, Q)$
Selection	$\sigma_\theta(Q)$	$ Q \times \text{SEL}(\theta, Q)$
Join	$Q_1 \bowtie_\theta Q_2$	$(Q_1 \times Q_2) \times \text{SEL}(\theta, Q_1 \times Q_2)$
Distinct	$\delta_A(Q)$	$\text{UNIQ}(A, Q)$
Aggregate	$\sum_{A, \text{agg}}(Q)$	

Operator	RA	Estimated Size
Table	R	$ R $
Projection	$\pi(Q)$	$ Q $
Union	$Q_1 \cup Q_2$	$ Q_1 + Q_2 $
Cross Product	$Q_1 \times Q_2$	$ Q_1 \times Q_2 $
Sort	$\tau(Q)$	$ Q $
Limit	$L_N(Q)$	$\min(N, Q)$
Selection	$\sigma_\theta(Q)$	$ Q \times \text{SEL}(\theta, Q)$
Join	$Q_1 \bowtie_\theta Q_2$	$(Q_1 \times Q_2) \times \text{SEL}(\theta, Q_1 \times Q_2)$
Distinct	$\delta_A(Q)$	$\text{UNIQ}(A, Q)$
Aggregate	$\sum_{A, \text{agg}}(Q)$	$\text{UNIQ}(A, Q)$

1. $\text{SEL}(\theta, Q)$: The selectivity of c on $Q = \frac{|\sigma_{\theta}(Q)|}{|Q|}$
2. $\text{UNIQ}(A, Q)$: The number of distinct values of $Q.A$

Idea 1: Assume each selection filters down to 10% of the data.

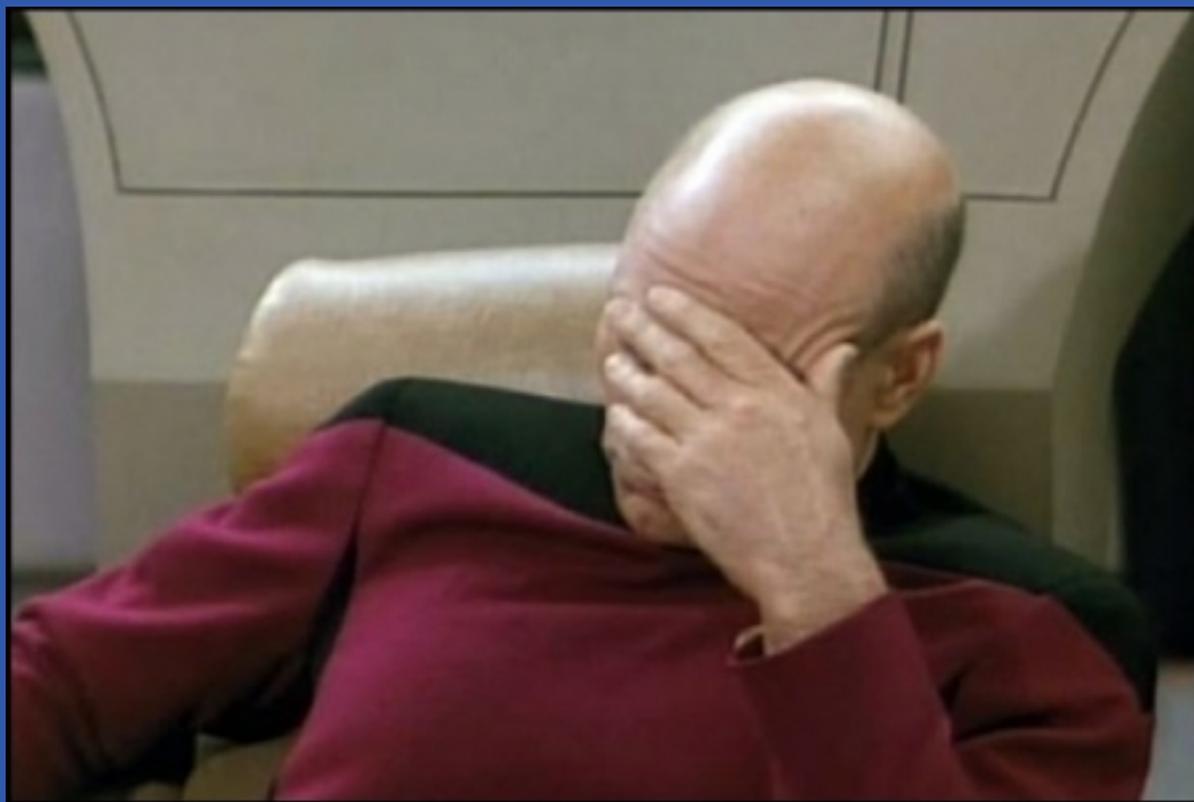


Image credit: © Paramount Pictures

No... really.

(Postgres and many other databases do this)

There are Problems

There are Problems

- Inconsistent estimation: $|\sigma_{\theta_1}(\sigma_{\theta_2}(R))| = |\sigma_{\theta_1 \wedge \theta_2}(R)|$

There are Problems

- Inconsistent estimation: $|\sigma_{\theta_1}(\sigma_{\theta_2}(R))| = |\sigma_{\theta_1 \wedge \theta_2}(R)|$
- Too consistent estimation: $|\sigma_{id=1}(\text{STUDENTS})| = |\sigma_{\text{residence} = \text{NY}}(\text{STUDENTS})|$

There are Problems

- Inconsistent estimation: $|\sigma_{\theta_1}(\sigma_{\theta_2}(R))| = |\sigma_{\theta_1 \wedge \theta_2}(R)|$
- Too consistent estimation: $|\sigma_{id=1}(\text{STUDENTS})| = |\sigma_{\text{residence} = \text{NY}}(\text{STUDENTS})|$

... remember that all we need is to **rank** plans.

Many major databases (Oracle, Postgres, Teradata, etc..) use something like the 10% rule if they have nothing better.

e.g. Teradata uses 10% for the first AND clause, 75% for every subsequent clause.

The 10% rule

- Rules of thumb if you have no other options

Uniform Prior

- Use basic statistics to make a rough guess

Sampling / History

- Small, quick sampling runs (or prior executions of the query)

Histograms

- More detailed statistics for improved guesses

Constraints

- Rules about the data for improved guesses

Assume that for $\sigma_\theta(Q)$ or $\delta_A(Q)$...

1. Basic statistics are known about Q :
 - `COUNT(*)`
 - `COUNT(DISTINCT A)` (for each A)
 - `MIN(A)`, `MAX(A)` (for each numeric A)
2. Attribute values are uniformly distributed.
3. No inter-attribute correlations.

- If necessary statistics aren't available, fall back to the 10% rule.
- If we make assumptions that are not perfectly correct, it will still likely be more accurate than the 10% rule.

$$\text{UNIQ}(A, \pi_{A, \dots}(Q)) = \text{UNIQ}(A, Q)$$

$$\text{UNIQ}(A, \pi_{A, \dots}(Q)) = \text{UNIQ}(A, Q)$$

$$\text{UNIQ}(A, \sigma(Q)) = \text{UNIQ}(A, Q)$$

$$\text{UNIQ}(A, \pi_{A, \dots}(Q)) = \text{UNIQ}(A, Q)$$

$$\text{UNIQ}(A, \sigma(Q)) = \text{UNIQ}(A, Q)$$

$$\text{UNIQ}(A, Q_1 \times Q_2) = \text{UNIQ}(A, Q_1) \text{ or } \text{UNIQ}(A, Q_2)$$

$$\text{UNIQ}(A, \pi_{A, \dots}(Q)) = \text{UNIQ}(A, Q)$$

$$\text{UNIQ}(A, \sigma(Q)) = \text{UNIQ}(A, Q)$$

$$\text{UNIQ}(A, Q_1 \times Q_2) = \text{UNIQ}(A, Q_1) \text{ or } \text{UNIQ}(A, Q_2)$$

$$\max(\text{UNIQ}(A, Q_1), \text{UNIQ}(A, Q_2)) \leq \text{UNIQ}(A, Q_1 \cup Q_2) \leq \text{UNIQ}(A, Q_1) + \text{UNIQ}(A, Q_2)$$

$$\min_A (\pi_{A,\dots}(Q)) = \min_A(Q)$$

$$\min_A(\pi_{A,\dots}(Q)) = \min_A(Q)$$

$$\min_A(\sigma(Q)) \approx \min_A(Q)$$

$$\min_A(\pi_{A,\dots}(Q)) = \min_A(Q)$$

$$\min_A(\sigma(Q)) \approx \min_A(Q)$$

$$\min_A(Q_1 \times Q_2) = \min_A(Q_1) \text{ or } \min_A(Q_2)$$

$$\min_A(\pi_{A,\dots}(Q)) = \min_A(Q)$$

$$\min_A(\sigma(Q)) \approx \min_A(Q)$$

$$\min_A(Q_1 \times Q_2) = \min_A(Q_1) \text{ or } \min_A(Q_2)$$

$$\min_A(Q_1 \cup Q_2) = \min(\min_A(Q_1) + \min_A(Q_2))$$

Estimating $\delta_A(Q)$ requires only `COUNT(DISTINCT A)`

Selectivity is a probability: $\text{SEL}(\theta, Q) = P[\theta(r) \mid r \in Q]$

Selectivity is a probability: $\text{SEL}(\theta, Q) = P[\theta(r) \mid r \in Q]$

$$P(A = c_1) = \frac{1}{\text{COUNT}(\text{DISTINCT } A)}$$

Selectivity is a probability: $\text{SEL}(\theta, Q) = P[\theta(r) \mid r \in Q]$

$$P(A = c_1) = \frac{1}{\text{COUNT}(\text{DISTINCT } A)}$$

$$P(A \in (c_1, \dots, c_n)) = \frac{n}{\text{COUNT}(\text{DISTINCT } A)}$$

Selectivity is a probability: $\text{SEL}(\theta, Q) = P[\theta(r) \mid r \in Q]$

$$P(A = c_1) = \frac{1}{\text{COUNT}(\text{DISTINCT } A)}$$

$$P(A \in (c_1, \dots, c_n)) = \frac{n}{\text{COUNT}(\text{DISTINCT } A)}$$

$$P(A \leq c_1) = \frac{c_1 - \text{MIN}(A)}{\text{MAX}(A) - \text{MIN}(A)}$$

Selectivity is a probability: $\text{SEL}(\theta, Q) = P[\theta(r) \mid r \in Q]$

$$P(A = c_1) = \frac{1}{\text{COUNT}(\text{DISTINCT } A)}$$

$$P(A \in (c_1, \dots, c_n)) = \frac{n}{\text{COUNT}(\text{DISTINCT } A)}$$

$$P(A \leq c_1) = \frac{c_1 - \text{MIN}(A)}{\text{MAX}(A) - \text{MIN}(A)}$$

$$P(c_1 \leq A \leq c_2) = \frac{c_2 - c_1}{\text{MAX}(A) - \text{MIN}(A)}$$

$$P(A = B) = \min\left(\frac{1}{\text{COUNT}(\text{DISTINCT } A)}, \frac{1}{\text{COUNT}(\text{DISTINCT } B)}\right)$$

$$P(A = B) = \min\left(\frac{1}{\text{COUNT}(\text{DISTINCT } A)}, \frac{1}{\text{COUNT}(\text{DISTINCT } B)}\right)$$

$$P(\theta_1 \wedge \theta_2) = P(\theta_1) \times P(\theta_2)$$

$$P(A = B) = \min\left(\frac{1}{\text{COUNT}(\text{DISTINCT } A)}, \frac{1}{\text{COUNT}(\text{DISTINCT } B)}\right)$$

$$P(\theta_1 \wedge \theta_2) = P(\theta_1) \times P(\theta_2)$$

$$P(\theta_1 \vee \theta_2) = 1 - (1 - P(\theta_1)) \times (1 - P(\theta_2))$$

We don't always have statistics for Q

$$\pi_{A=(B \times C)}(R)$$

We don't always have clear rules for θ

$$\sigma_{\text{FitsModel}(A, B, C)}(R)$$

Attribute values are not always uniformly distributed

$$\sigma_{\text{major} = \text{CSE}}(\text{STUDENTS})$$

Attribute values are sometimes correlated

$$\sigma_{\text{major} = \text{CSE} \wedge \text{school} = \text{SEAS}}(\text{STUDENTS})$$

- Quiz 3 results posted
- 1 week to Checkpoint 2 due date
 - Don't forget to sign up for a code review