

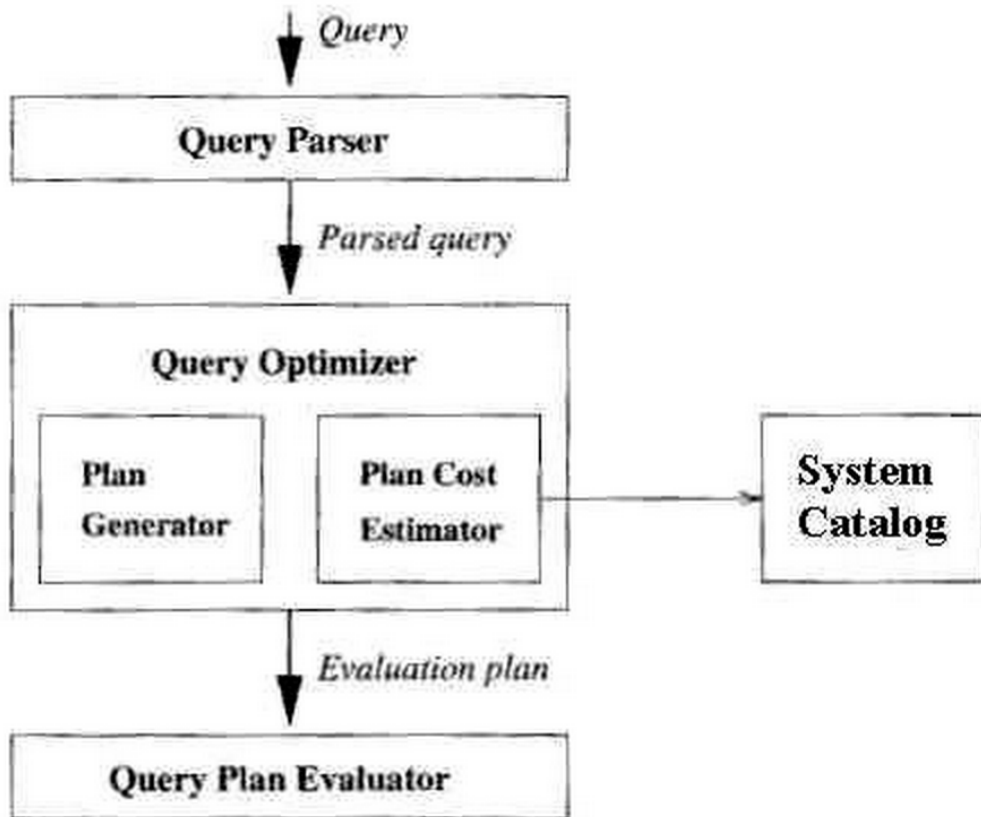
HELLO!

Rate-Based Query Optimization for Streaming Information Sources

Stratis D. Viglas

Jeffrey F. Naughton

Query Optimization



Cardinality Based vs. Rate Based Cost Estimation

Let us consider two select operations A and B. Assume that the selectivity for A is 0.1 and B is 0.2 and that the input size is 500.

$$\text{Cost (A} \rightarrow \text{B)} = 500 * c_{\{A\}} + 500 * 0.1 * c_{\{B\}}$$

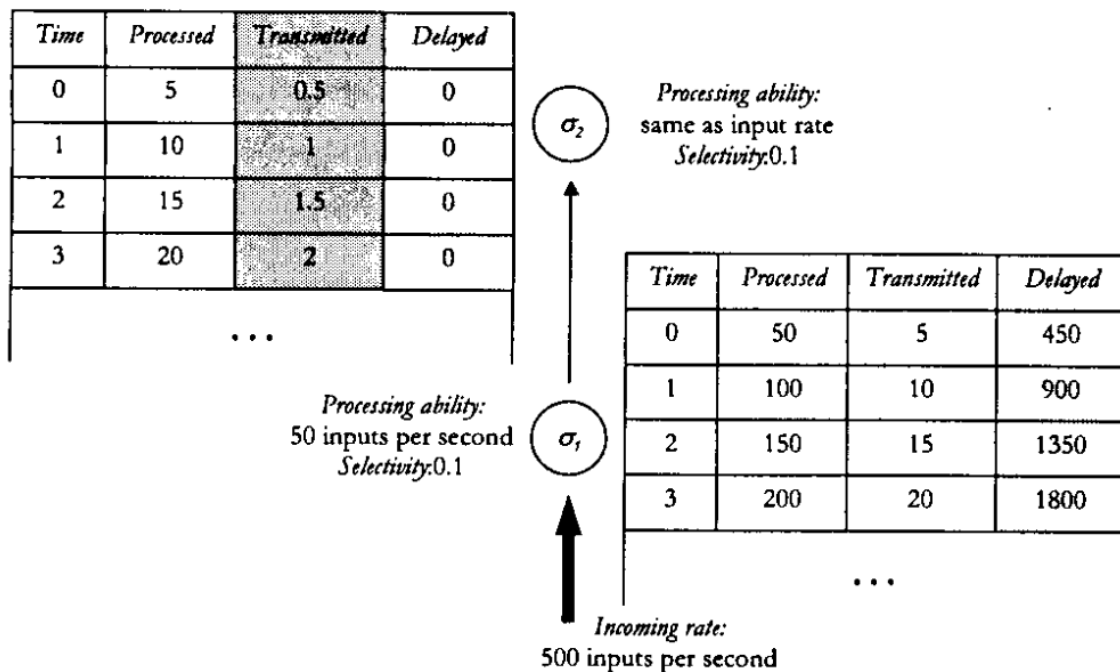
$$\text{Cost (B} \rightarrow \text{A)} = 500 * c_{\{B\}} + 500 * 0.2 * c_{\{A\}}$$

Assume that the selectivity of each of A and B is 0.1; input arrives at 500 tuples per second; A can process 50 inputs per second and B can process data as fast as it receives it.

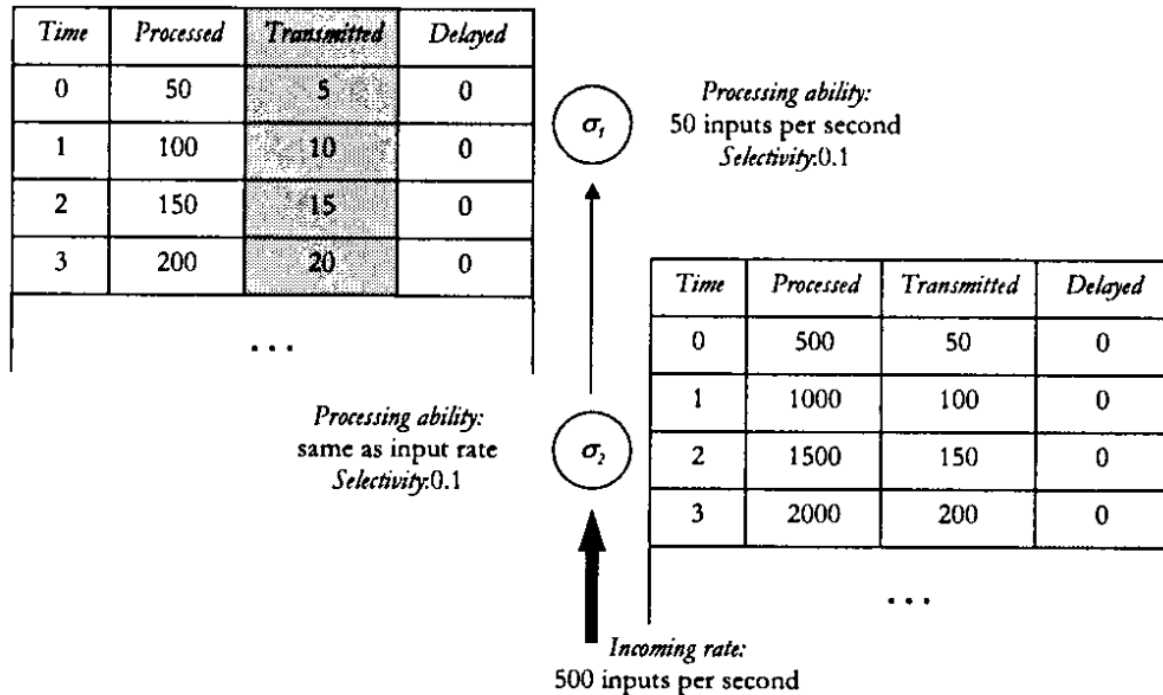


Size of input is infinite

⇒ Cost of each plan is infinite



(a) Output rate = 0.5 outputs per second



(b) Output rate = 5 outputs per second

“

~~*“What is the cost of this query plan?”*~~

“What is the expected output rate of this query plan?”

Estimating Output Rates

$$\text{Output rate} = \frac{\text{Number of outputs transmitted}}{\text{Time needed to make the transmission}}$$

Table 1: Cost variables used in the estimation of output rates

<i>Cost Variable</i>	<i>Meaning</i>	
C_{π}	Cost of projecting parts of an input object	r_o Output Rate
C_{σ}	Cost of performing a selection on an input object	r_i Input Rate
C_l	Cost of handling an input coming from the left-hand side of a join	r_r Right Input Rate
C_r	Cost of handling an input coming from the right-hand side of a join	r_l Left Input Rate
T	Cost of making a single transmission	

Projections

$$r_o = r_i$$

Selections

$$r_o = f \cdot r_i$$

Joins

$$r_o = \frac{f \cdot r_l \cdot r_r \cdot t}{r_l \cdot C_l + r_r \cdot C_r}$$

**Optimize for a specific time point in the execution process
using local rate maximization**

**Optimize for output production size using local time
minimization**

Experimental Validation

Rate Based Cost Model

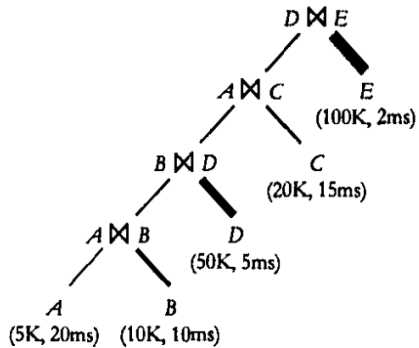
Does the cost model correctly estimate individual plan performance?

Is the framework capable of providing correct decisions regarding the best choice among a set of plans?

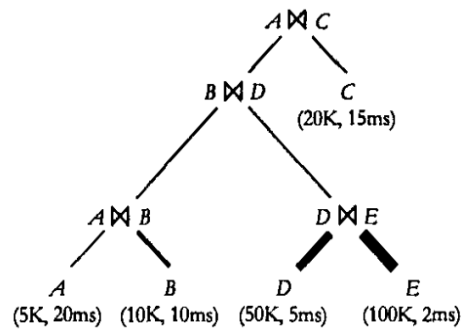
5 XML data sources
Wide range of selectivities

<i>Source</i>	<i>Number of tuples</i>	<i>Size</i>
<i>A</i>	5,000	0.7 MB
<i>B</i>	10,000	1.5 MB
<i>C</i>	20,000	1.8 MB
<i>D</i>	50,000	5.9 MB
<i>E</i>	100,000	9.3 MB

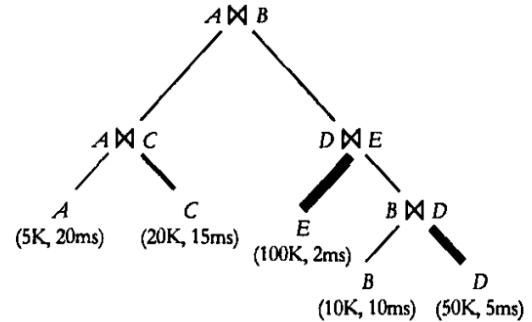
5 Way Equi Join



(a) Left Deep



(b) Fast Leaves



(c) Evenly Spread

Comparison to Traditional Cost
Model

<i>Plan</i>	<i>Traditional Estimation</i>	<i>Rate-Based estimation</i>
<i>Left Deep</i>	10^4	$1.3 \cdot 10^3$
<i>Fast Leaves</i>	$2 \cdot 10^3$	$9.7 \cdot 10^2$
<i>Evenly Spread</i>	$5 \cdot 10^3$	$8.8 \cdot 10^2$



**Rate Based Estimation
is the way to go!**

THANKS!

Any questions?