

Adaptive parallelization in multi-core systems

Mrunal Gawade, Martin Kersten

CWI, Amsterdam

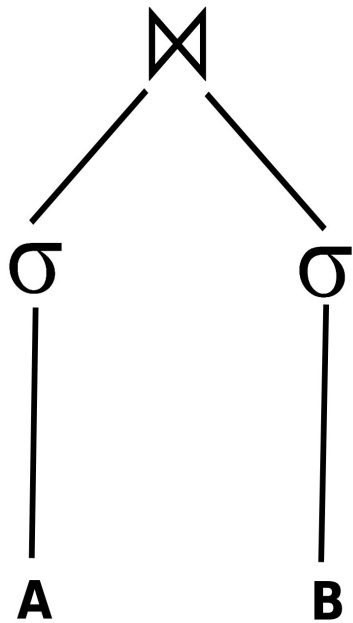
Erasmus University Rotterdam
29th Nov 2013

Why Parallelize Adaptively?

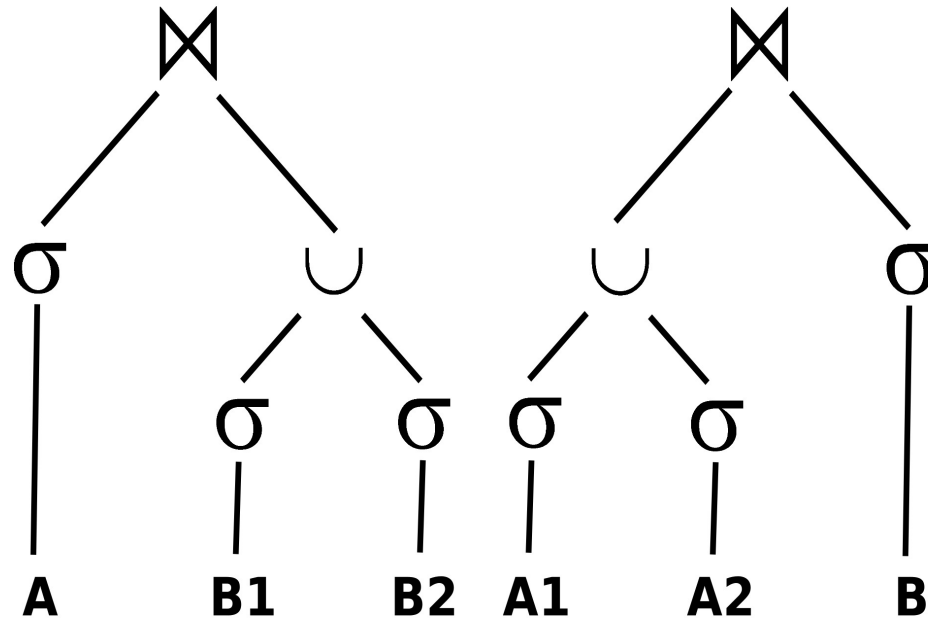
- Multi-core plan space is very large
(as compared to a serial plan space)
- Statistics maintenance overheads
(cost based optimizers)
- Optimal degree of parallelism is hard to find
(resource variations)

Multi-core Plan Space

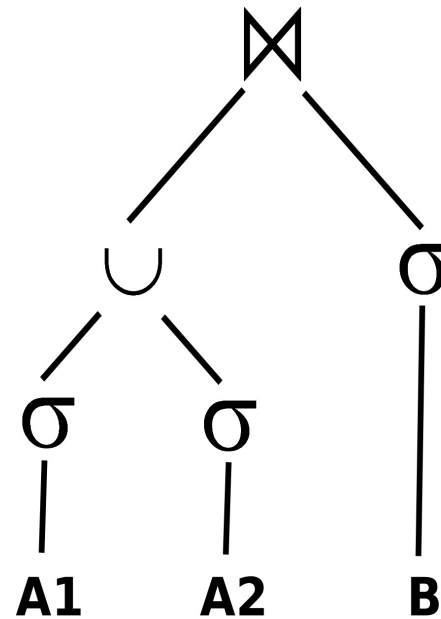
Join (Select(A), Select(B))



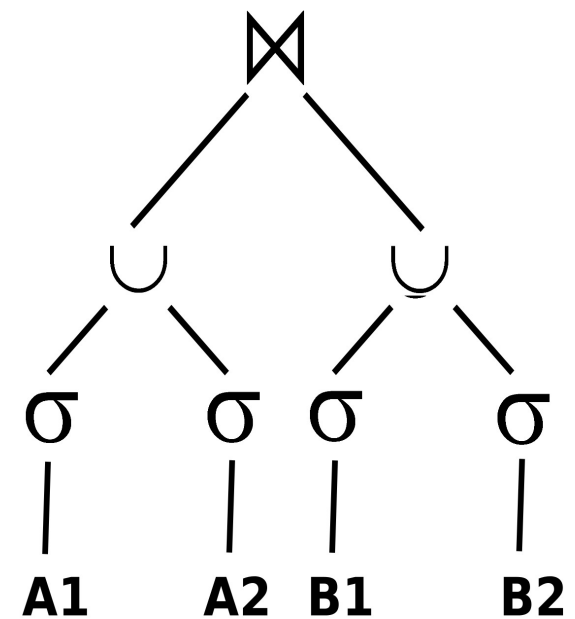
Plan 1



Plan 2

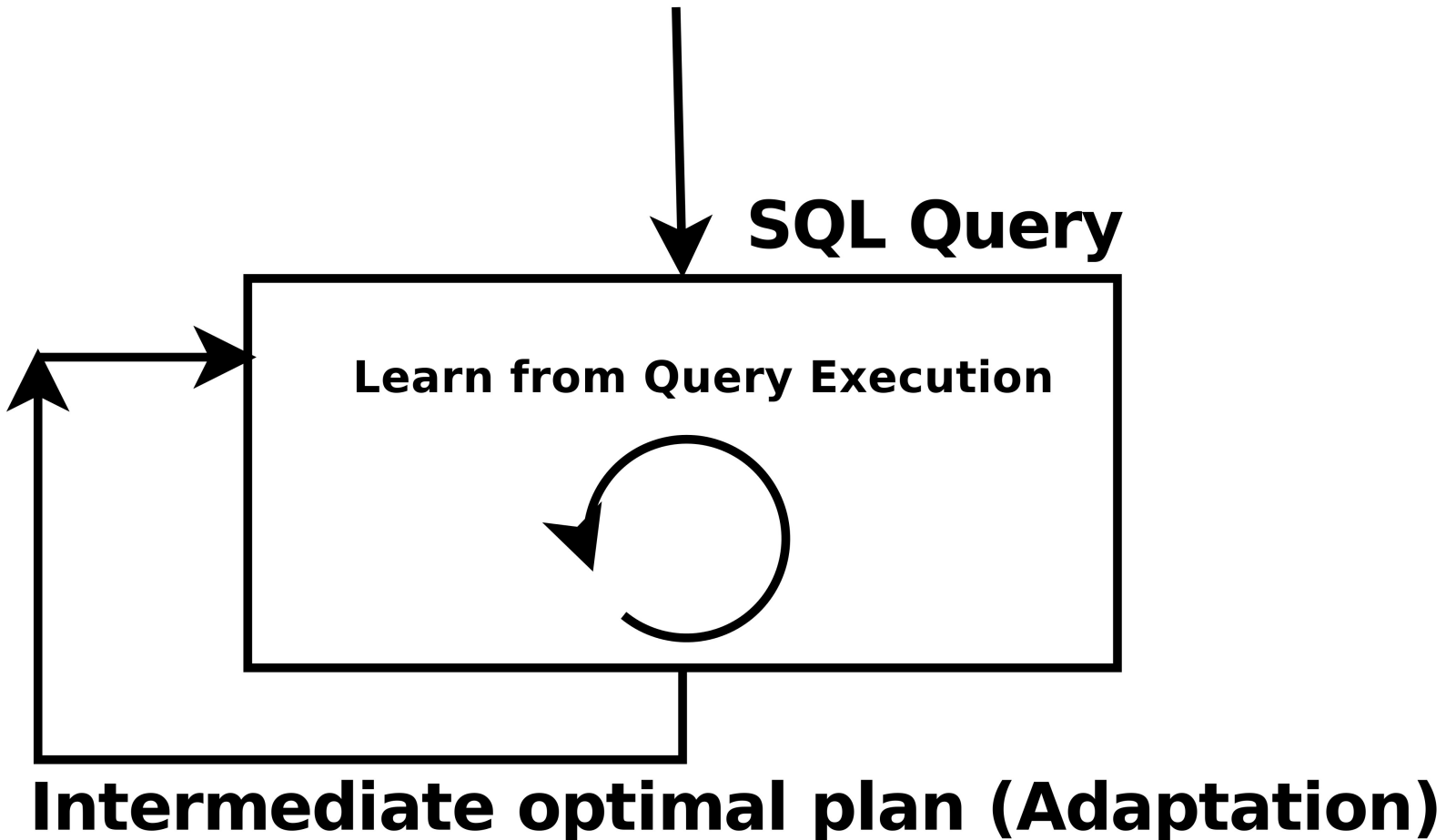


Plan 3

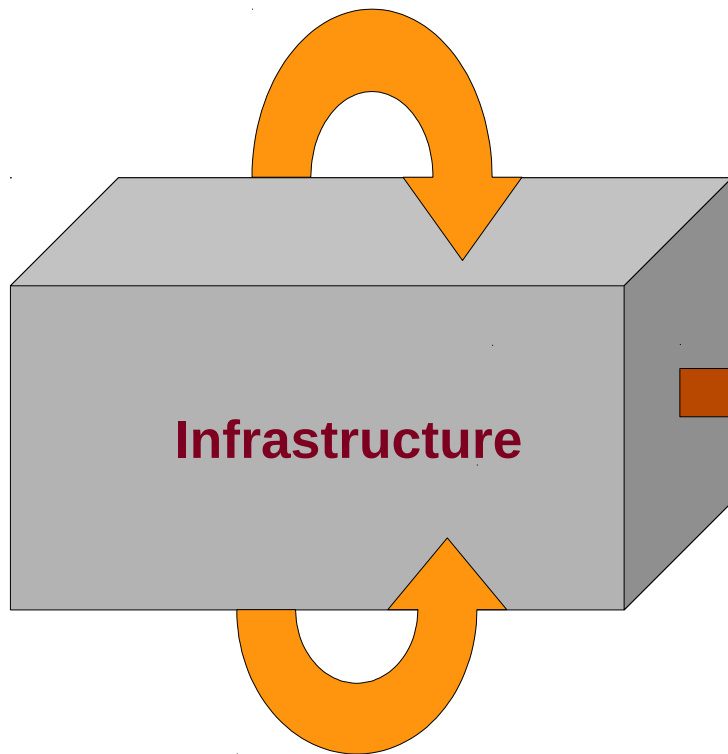


Plan 4

Bird's Eye View

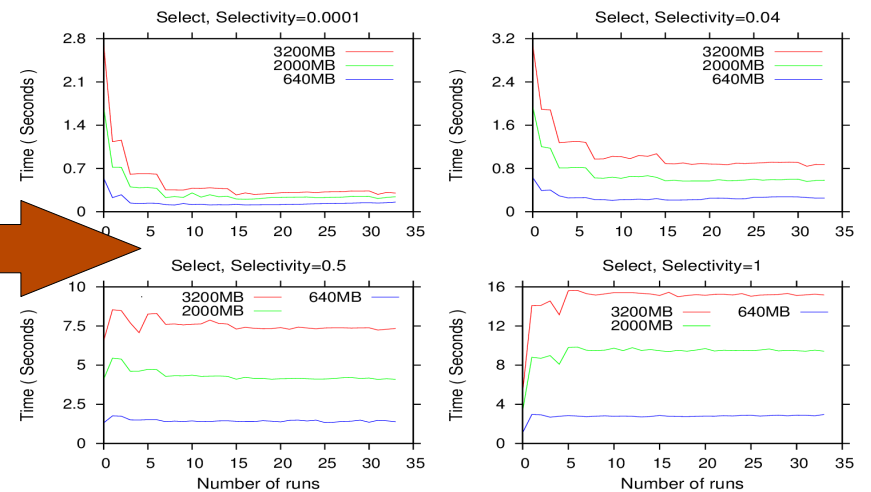


Contributions



Convergence Algorithm

Select & Join Operator Analysis



Infrastructure

Design

- Query execution time is the decision metric
- Store profiled history of plans
- Mutate an old plan to generate a new plan
- Parallelize the most expensive operator
(from previous query invocation)

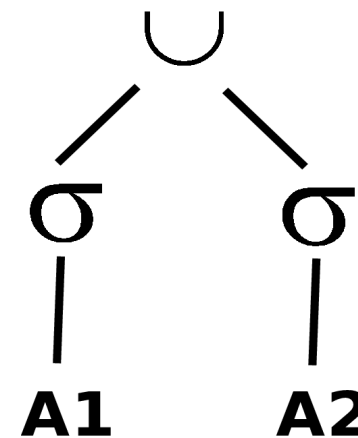
Select Operator Mutation

$X_1 = \text{select}(A, \text{low}, \text{high});$

- $A1 = \text{partition}(A, 2, 0);$
- $A2 = \text{partition}(A, 2, 1);$

- $Y_1 = \text{select}(A1, \text{low}, \text{high});$
- $Y_2 = \text{select}(A2, \text{low}, \text{high});$

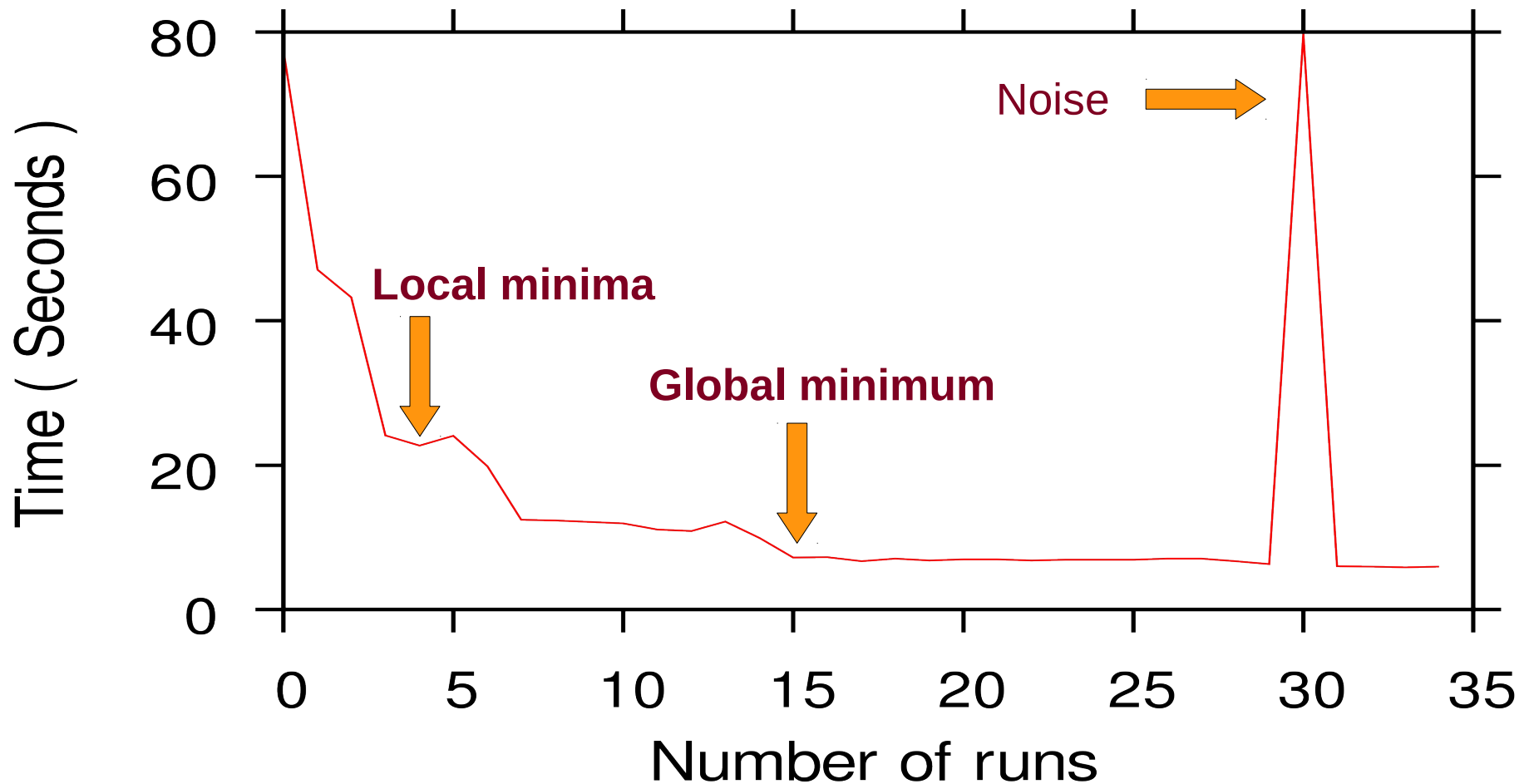
- $X_1 = \text{pack}(Y_1, Y_2);$



Convergence Algorithm

Sample Convergence Case

Convergence algorithm scenarios



Convergence Algorithm

ROI = Rate Of Improvement

CET = Current run Execution Time.

PET = Previous run Execution Time.

$$\text{ROI} = (\text{CET} - \text{PET}) / \text{MAX}(\text{CET}, \text{PET})$$

$$\text{Credit} = \text{Credit} + (\text{ROI} \times \text{Number Of Cores})$$

$$\text{Debit} = \text{Debit} + (\text{ABS}(\text{ROI}) \times \text{Number Of Cores})$$

$$\text{Convergence Runs} = \text{Credit} - \text{Debit} = f(\text{ROI})$$

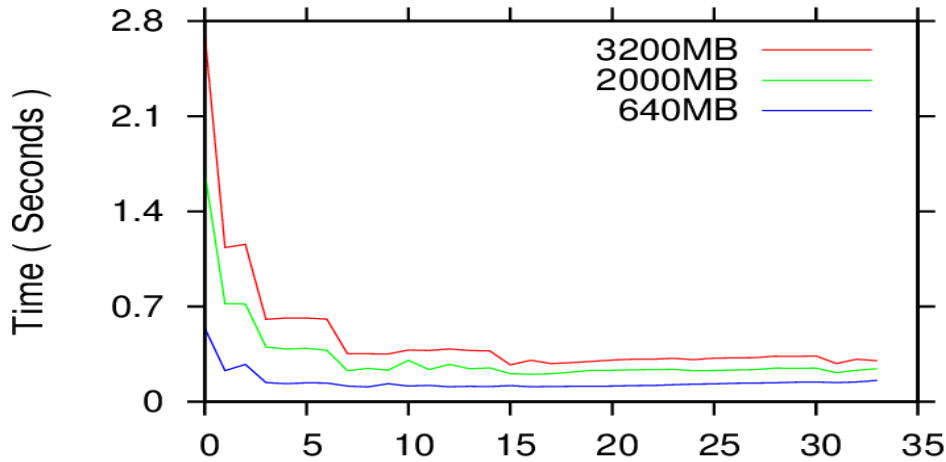
Select Operator Analysis

Experiment platform

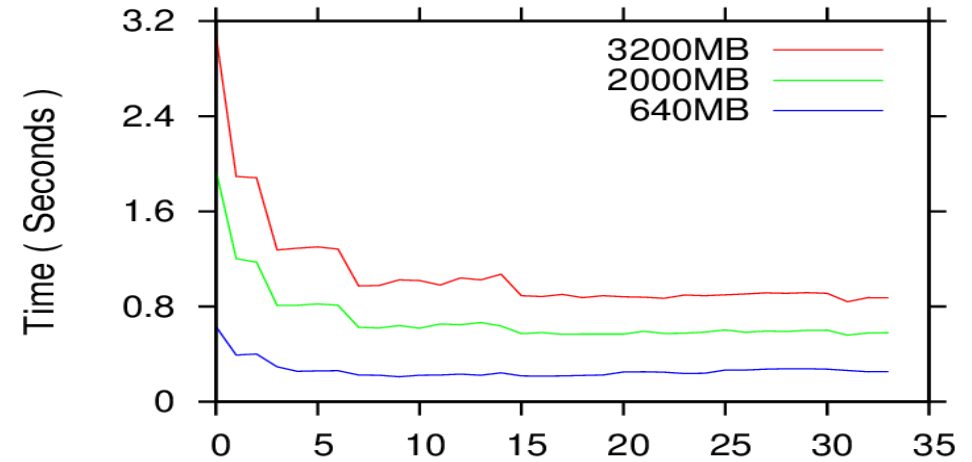
- Intel Xeon E5-2650@ 2.00GHz CPU (Dual socket)
16 physical cores, 32 threads with hyper-threading
- Cache shared L3 = 20MB.
- 256 GB of 4 channel DDR3 RAM.
- MonetDB on Fedora Core 16.

Adaptive Parallelization of Select

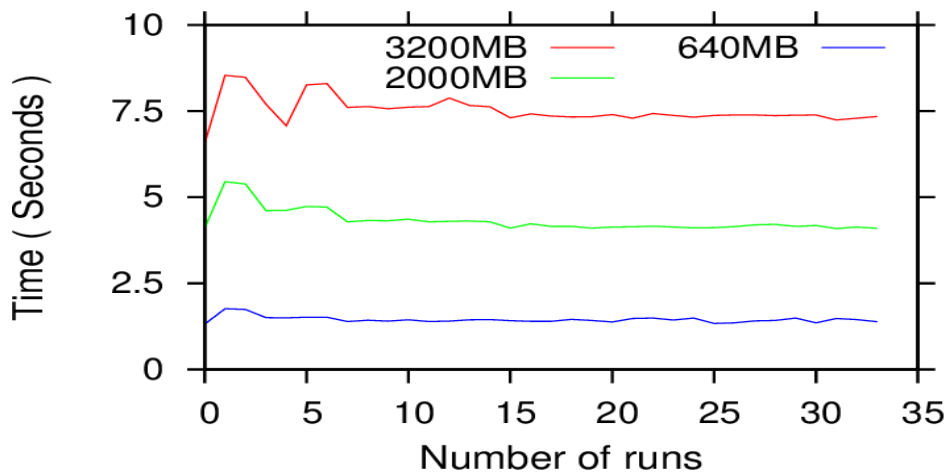
Select, Selectivity=0.0001



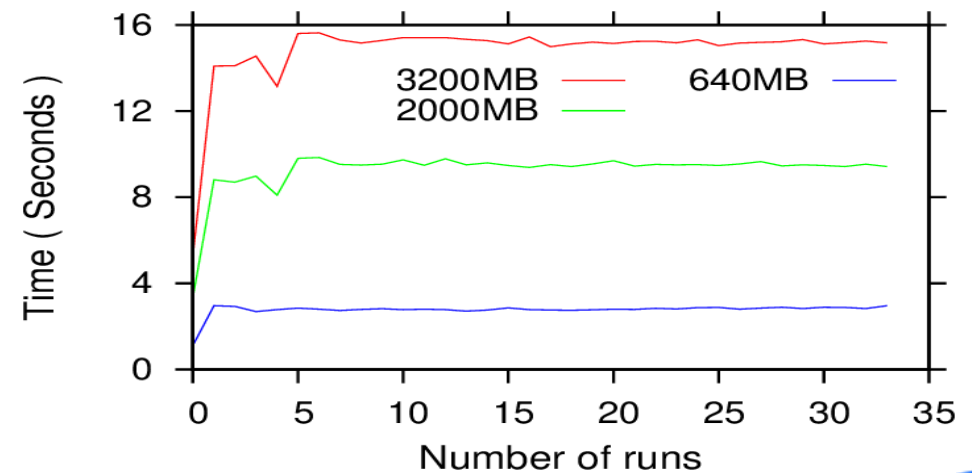
Select, Selectivity=0.04



Select, Selectivity=0.5



Select, Selectivity=1



Speed-up & Partitioning

Size(MB)	Selectivity		# Optimal Partitions
	0.0001	0.04	
	Speed-up		
3200	10	3.5	15
2000	7.5	3	7
640	4	2.5	3

TPC-H queries

Queries	Adaptive Parallelization	Heuristic Parallelization
Q 3	21.6 seconds	36 seconds
Q 5	25	28
Q 8	17.7	24
Q 14	4.1	4.1
Q 17	240	250
Q 20	42.5	45

- 11 queries show lower performance (mostly short queries)
- 5 queries not considered

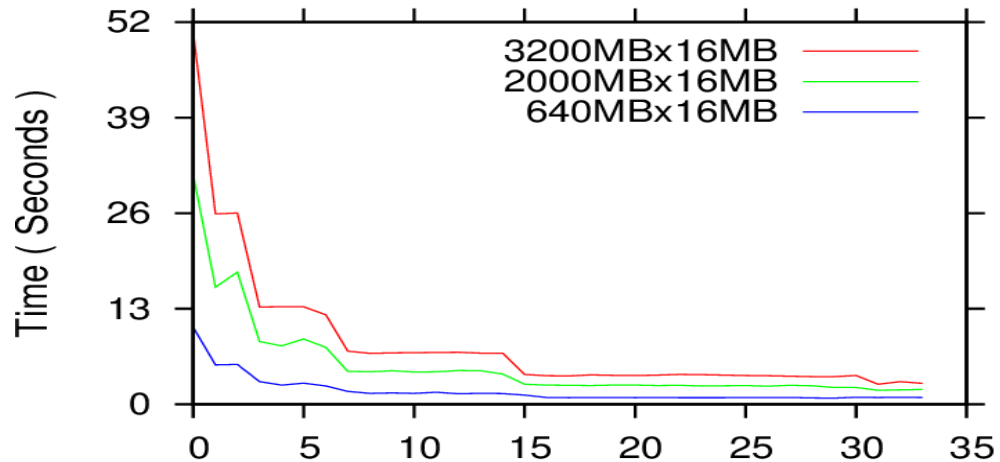
Summary

- Detailed analysis of Degree of Parallelism of plans
- Convergence algorithm adjusts dynamically
- High selectivity operator speed-up of an order of magnitude
- Long running TPC-H queries benefit

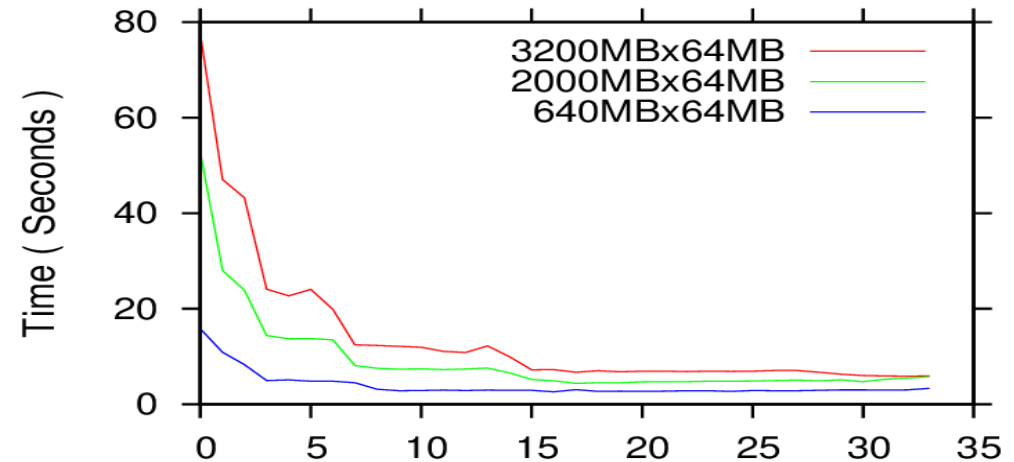
Thank you

Adaptive Parallelization of Join

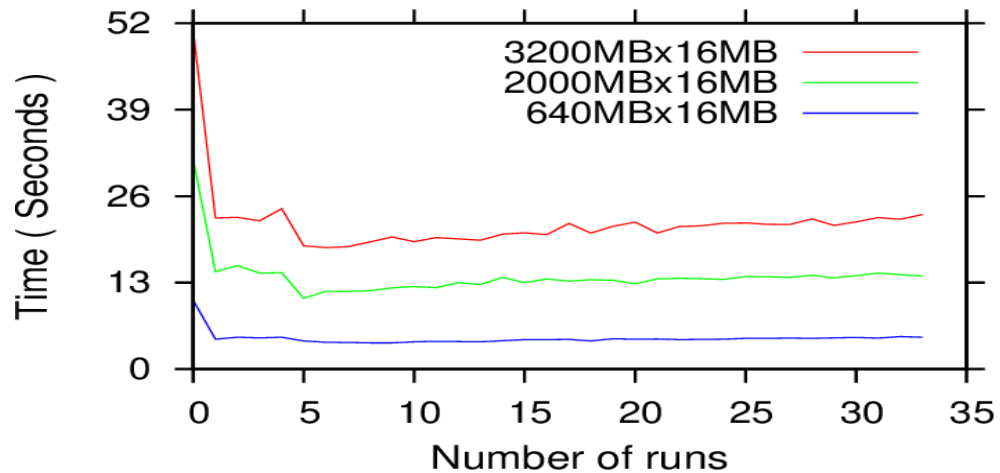
Case A-Join, larger input partitioned



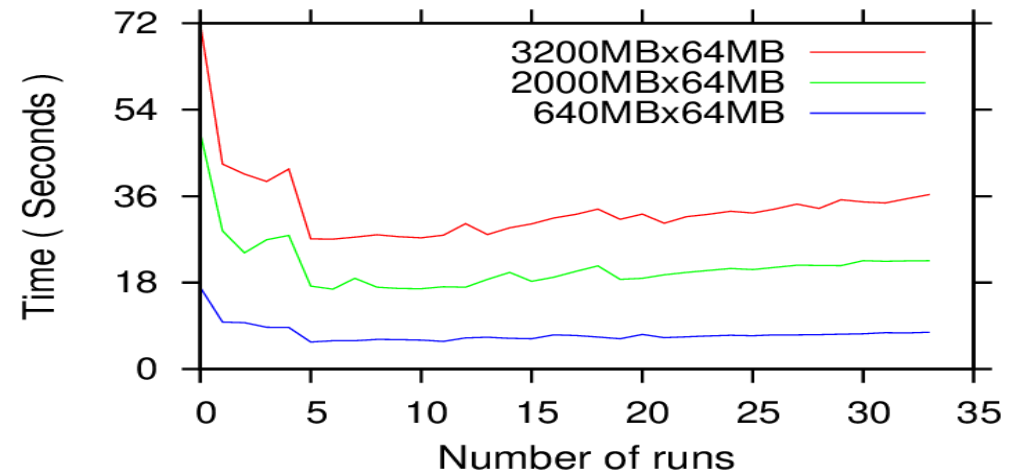
Case A-Join, larger input partitioned



Case B-Join, both input partitioned



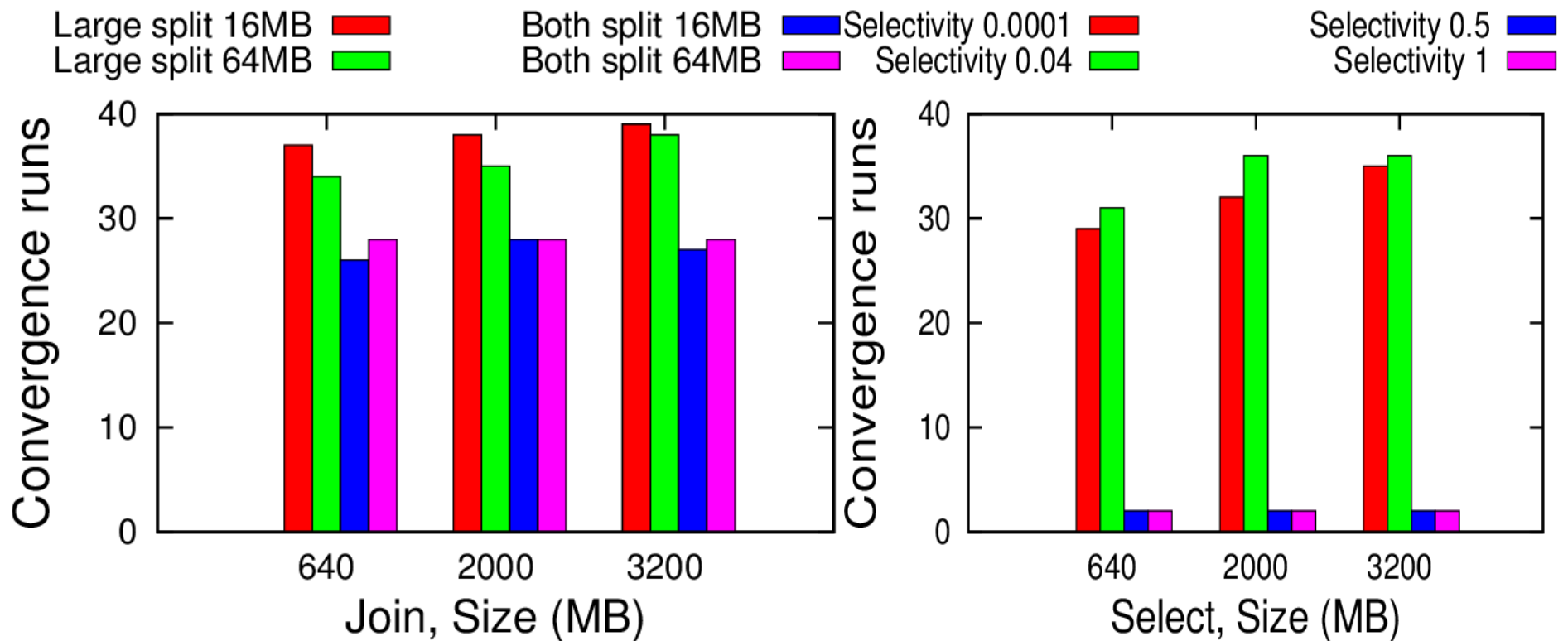
Case B-Join, both input partitioned



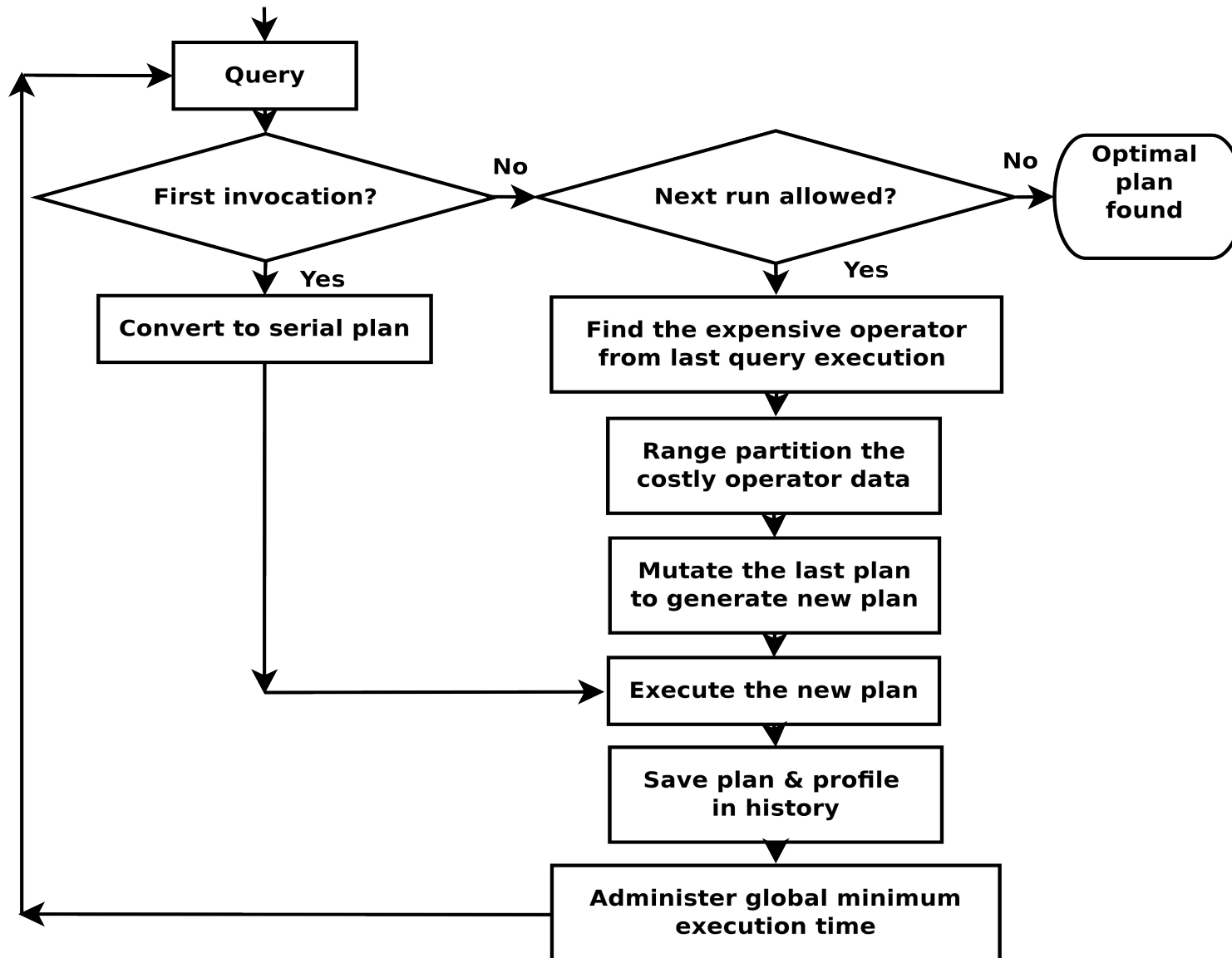
Speed-up & Partitioning

Size (MB)	16	64	16	64	16	64	16	64
	Case A		Case B		Case A		Case B	
	Speed-up				#Partitions			
3200	18.5	12	2.75	2.5	31	31	16	16
2000	16	12	3	3	31	15	16	16
640	11	6	2.5	3	15	15	28	16

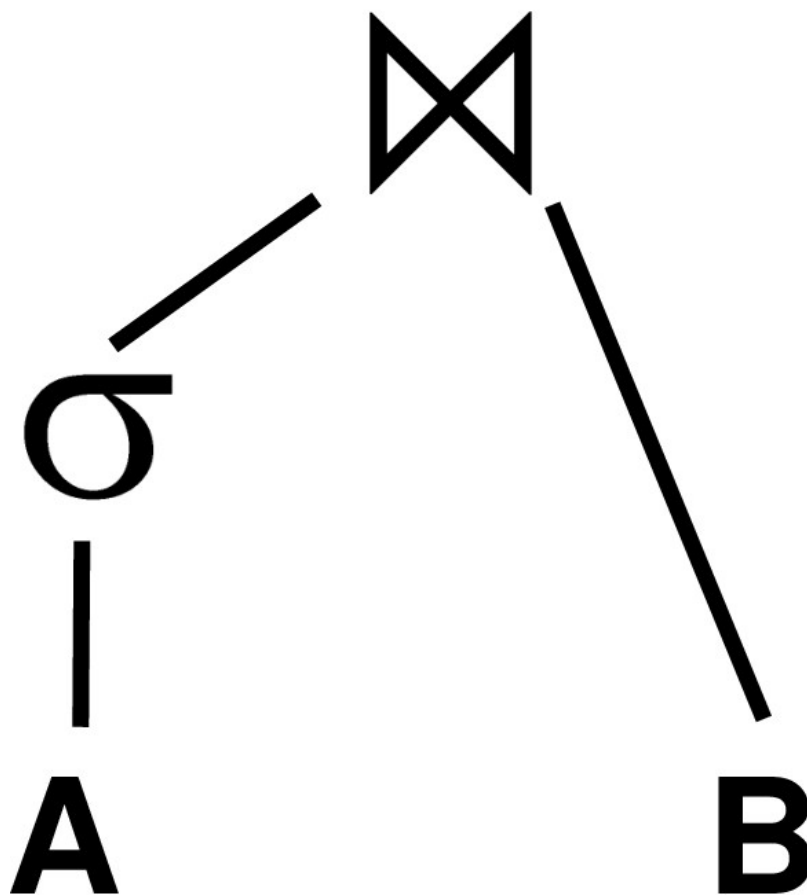
Convergence Runs Dynamic Adjustment



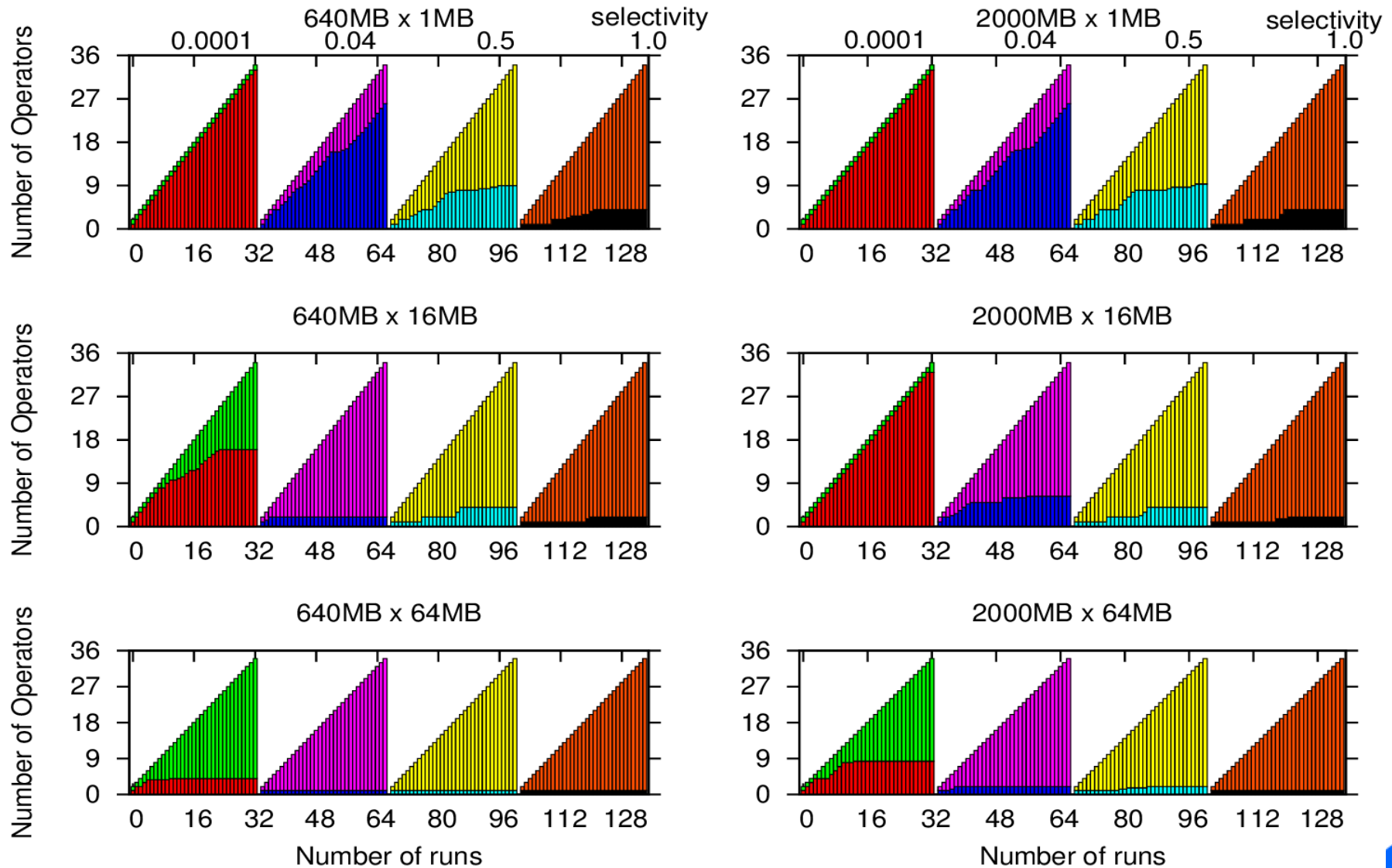
Work-flow



Select-Join Plan



Select-Join Degree of Parallelism



Select-Join execution

