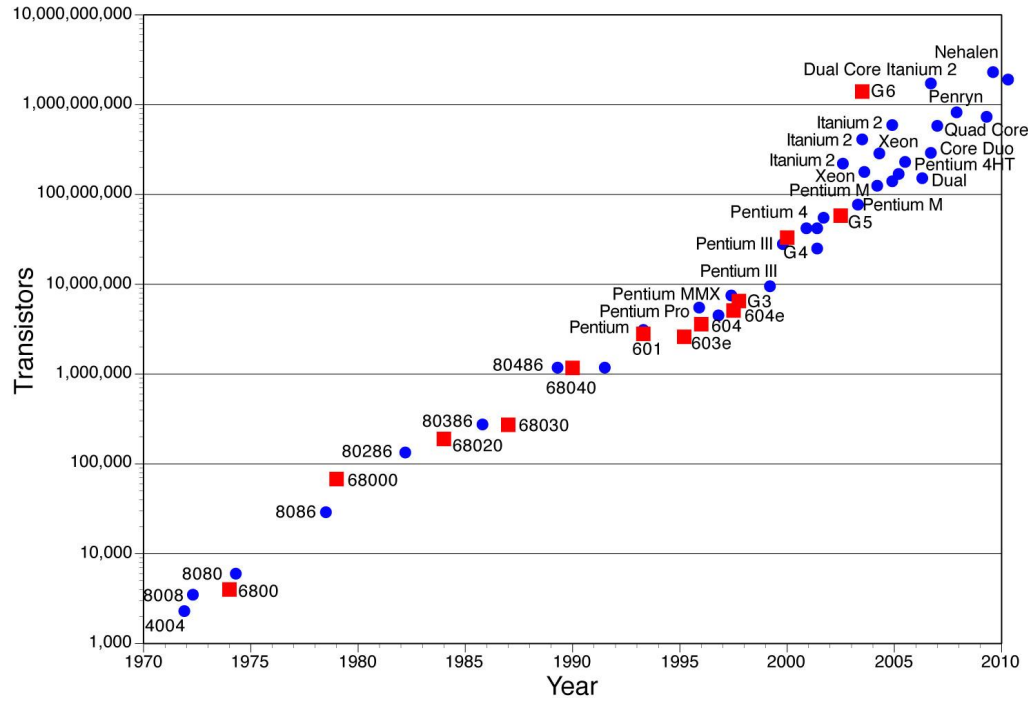


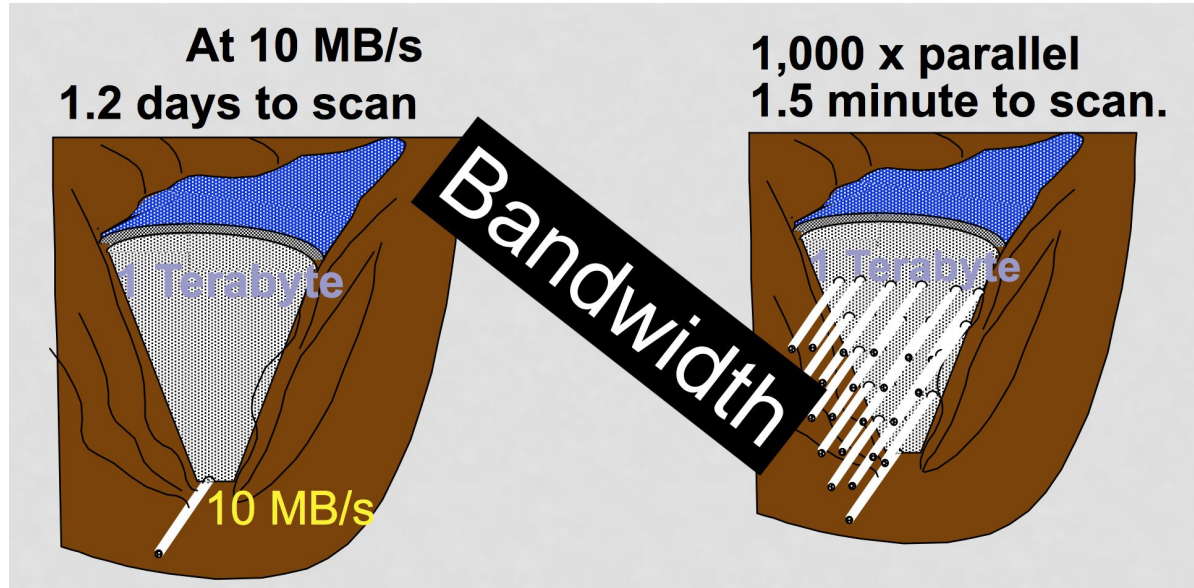
SQL-on-Hadoop

Aron Szanto and Jack Dent

Why do we need to parallelize data analysis?



Why do we need to parallelize data analysis?



Why do we need to parallelize data analysis?

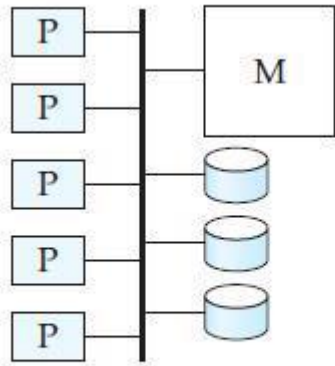
d = data size (GB)

b = bandwidth of single machine (GB/s)

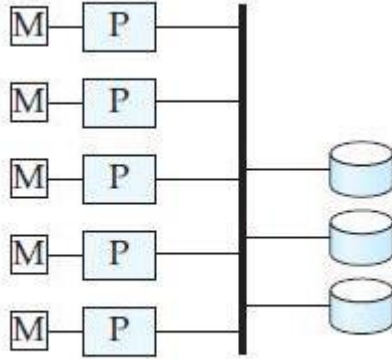
Time on single machine architecture = d/b

Time on n -machine architecture = d/nb
(assumes perfect horizontal scalability)

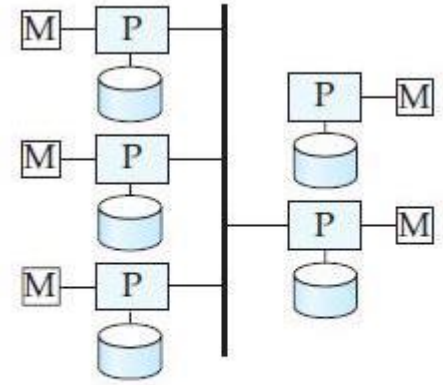
Parallel database architectures



(a) shared memory

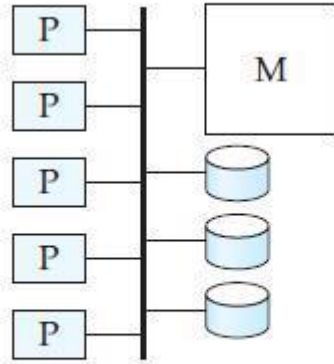


(b) shared disk



(c) shared nothing

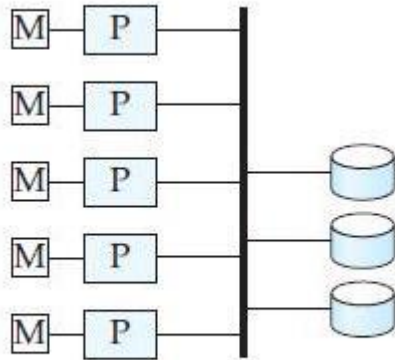
Shared-memory architectures



(a) shared memory

Definition: there is a single memory address-space for all processors, but each processor can have its own disk, local memory, and cache

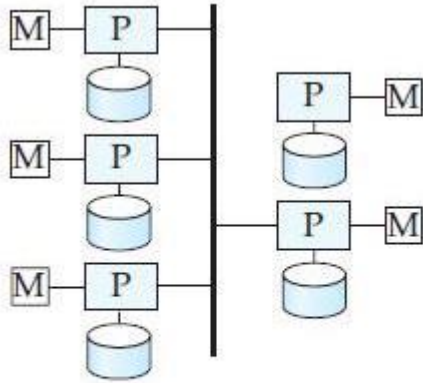
Shared-disk architectures



(b) shared disk

Definition: “every processor has its own memory (not accessible by others), and all machines can access all disks in the system”

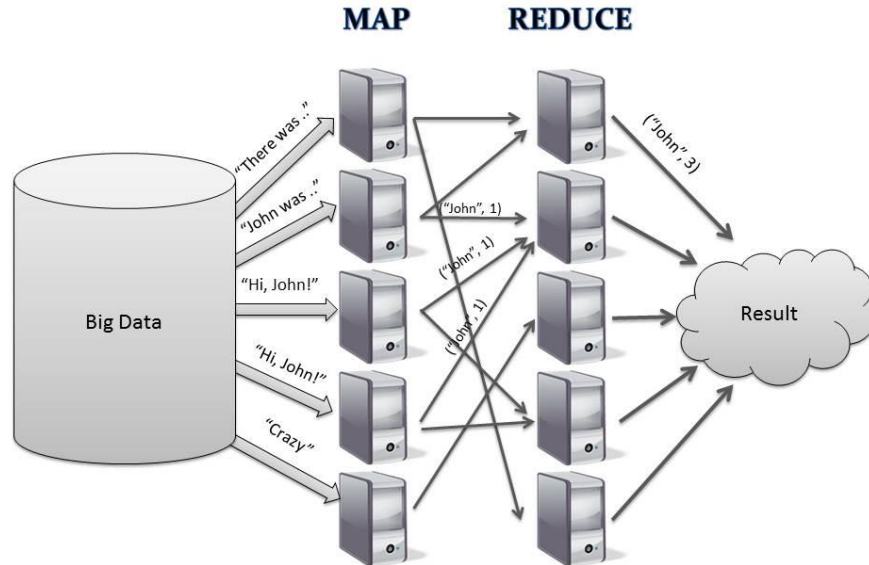
Shared-nothing architectures



(c) shared nothing

Definition: “a collection of independent, possibly virtual, machines, each with local disk and local main memory, connected together on a high-speed network”

MapReduce: shared-nothing data analysis

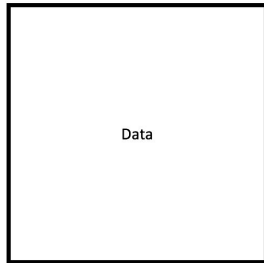


Key paper: "MapReduce: Simplified Data Processing on Large Clusters", Dean and Ghemawat, Google, 2004

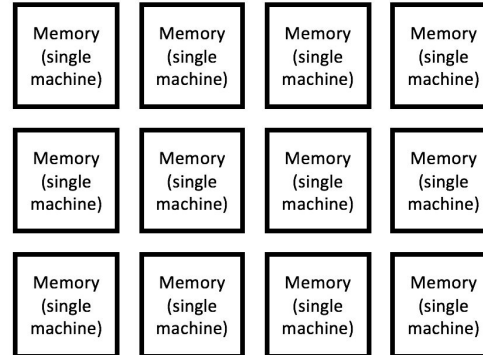
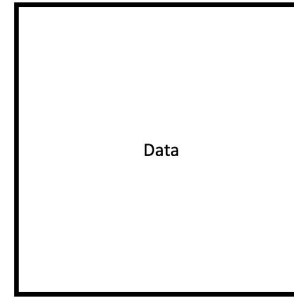
Open source implementation in Apache Hadoop suite

Scaling main memory

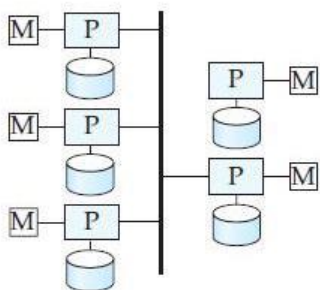
Single machine



Parallel machines



Challenge: SQL queries on shared-nothing architectures?



(c) shared nothing

+



SQL

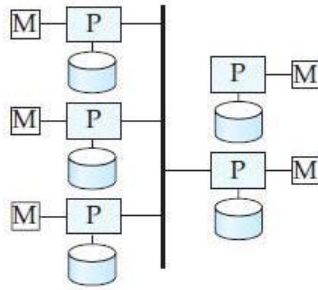
Scale out to 1000s of machines
Fault tolerant
Support heterogeneous environments

... but difficult to program, and not
performant for structured data

Scale up (fast queries over structured
data)
Flexible query language

... but do not scale out well

Challenge: SQL queries on shared-nothing architectures?



(c) shared nothing

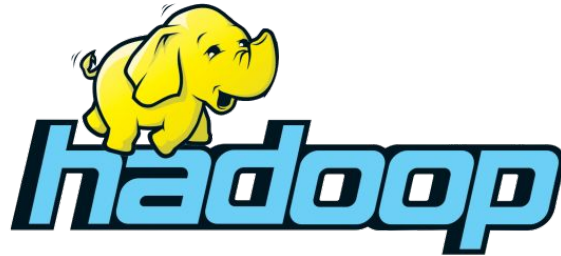
+



SQL

Can we combine the positive features (*performance, flexible query interface*) of shared-architecture parallel databases with the positive features (*fault tolerance, horizontal scalability*) of shared-nothing architectures?

HadoopDB (background)



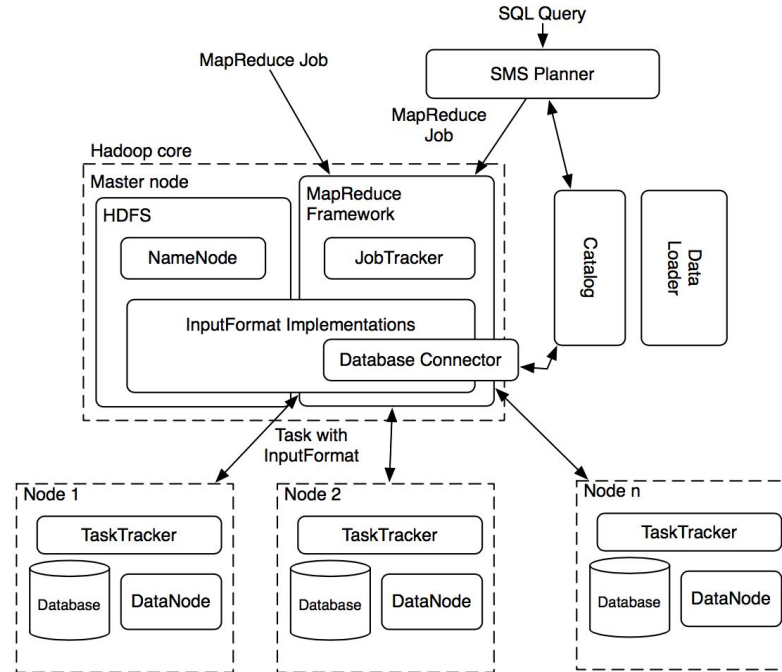
HDFS + MapReduce
inter-node

+



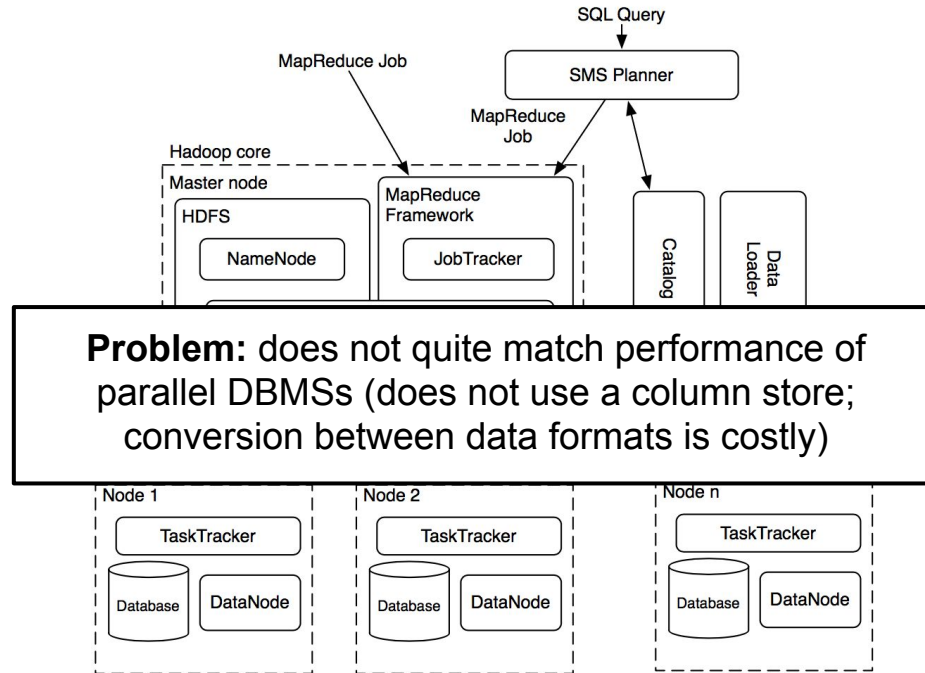
SQL query execution
intra-node

HadoopDB (background)



Source(s): “HadoopDB: An Architectural Hybrid of MapReduce and DBMS Technologies for Analytical Workloads”

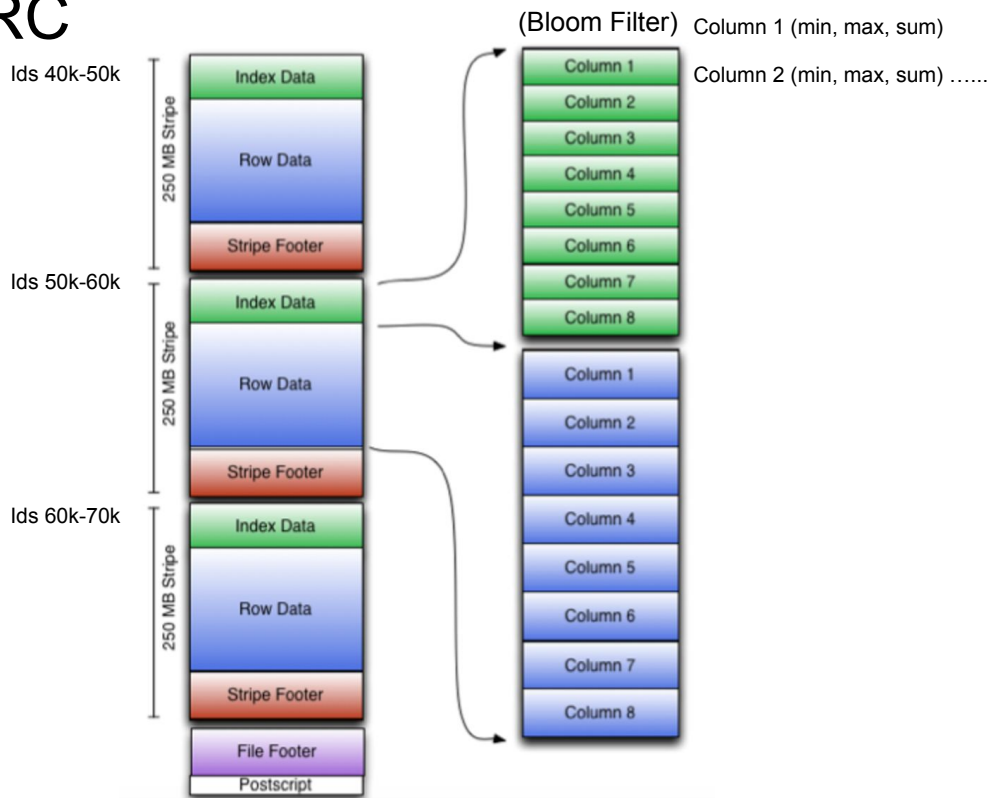
HadoopDB (background)



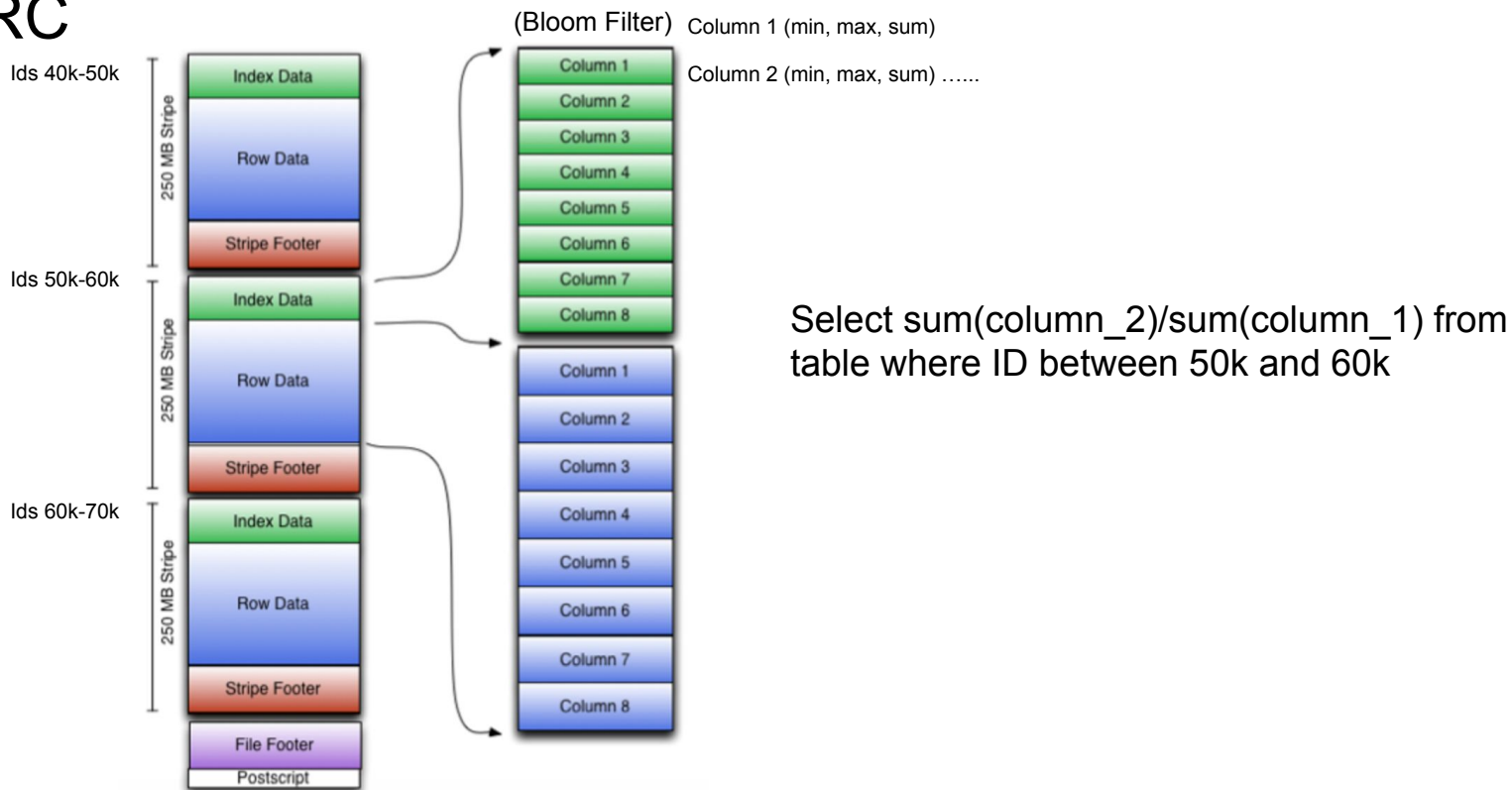
SQL with shared-nothing architectures

	File system	File format	Query language	Distributed runtime
Apache Hive	Apache HDFS	Optimized Row Columnar (ORC)	HiveQL	MapReduce or Tez
Cloudera Impala	Apache HDFS	Parquet	Impala SQL	<code>impalad</code>

Hive file format: ORC

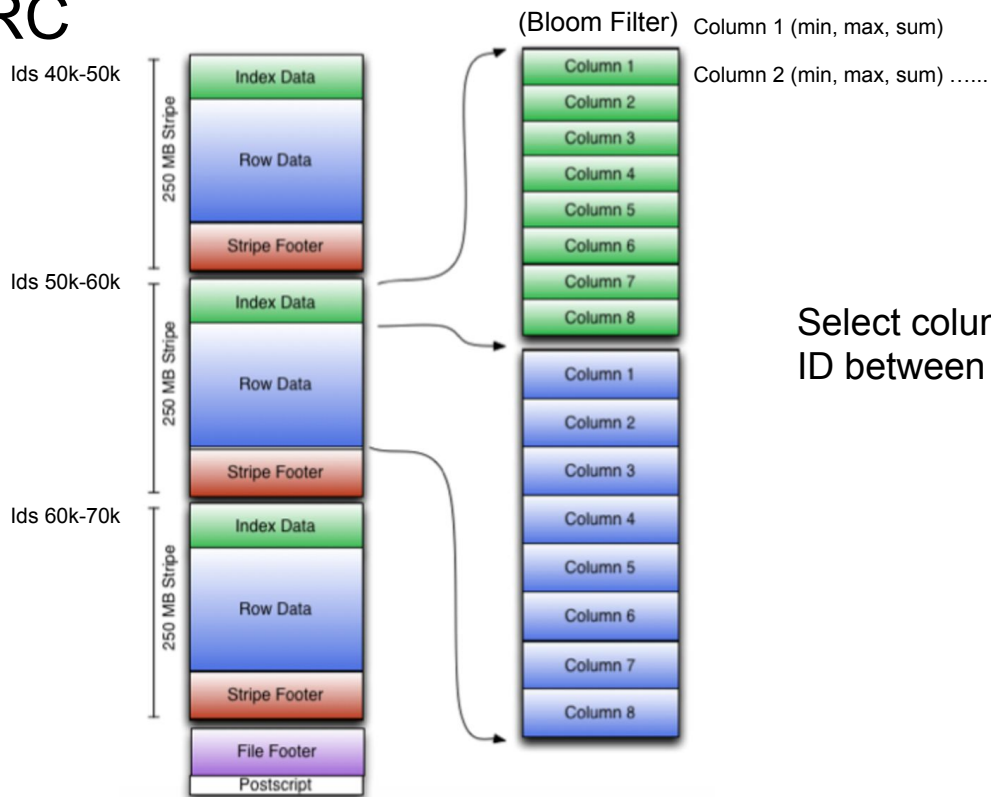


Hive file format: ORC



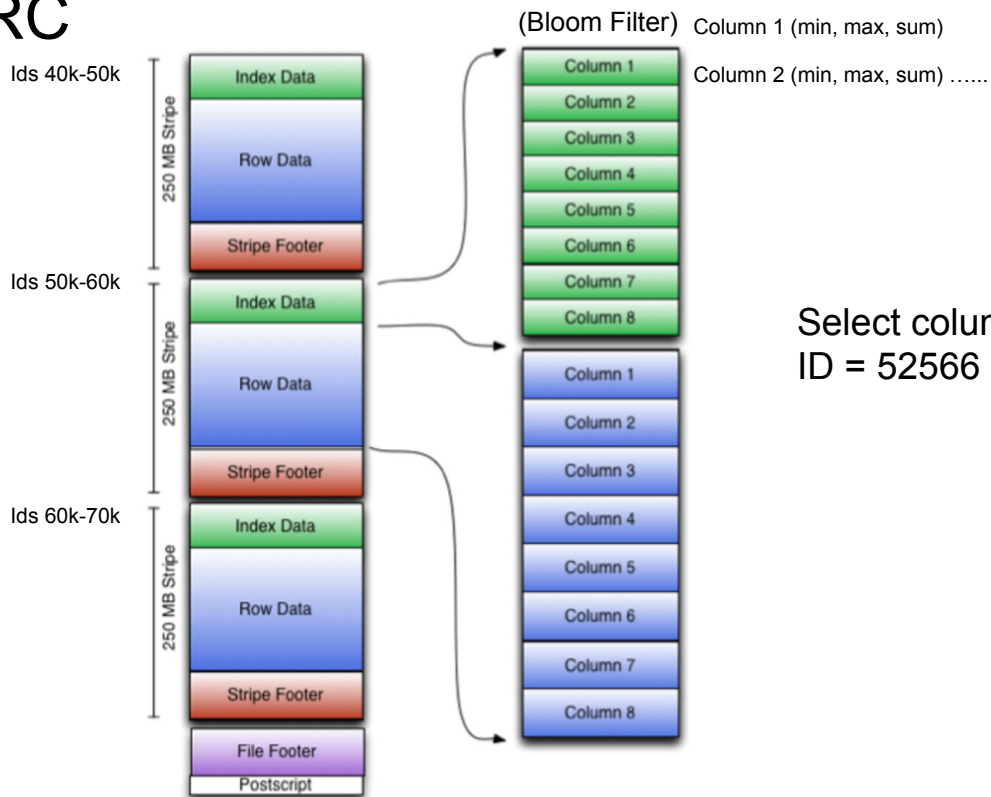
Select $\text{sum}(\text{column_2})/\text{sum}(\text{column_1})$ from table where ID between 50k and 60k

Hive file format: ORC



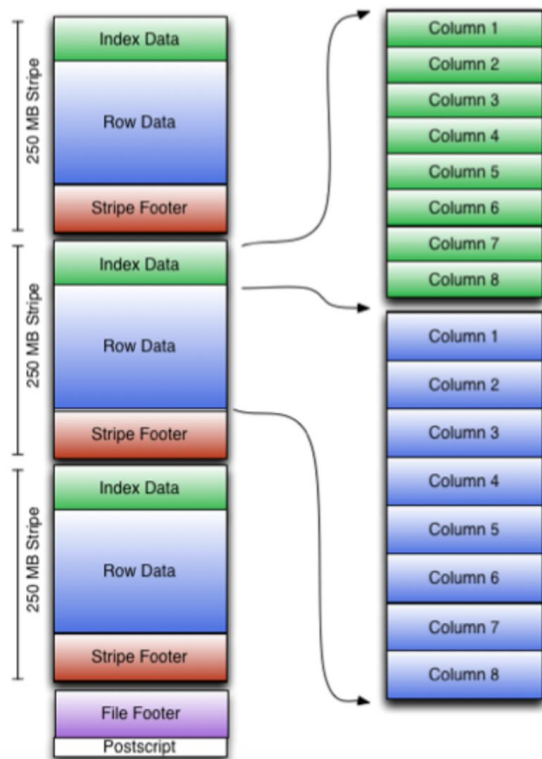
Select column_2, column_4 from table where
ID between 52k and 57k

Hive file format: ORC



Select column_2, column_4 from table where
ID = 52566 (which doesn't exist!)

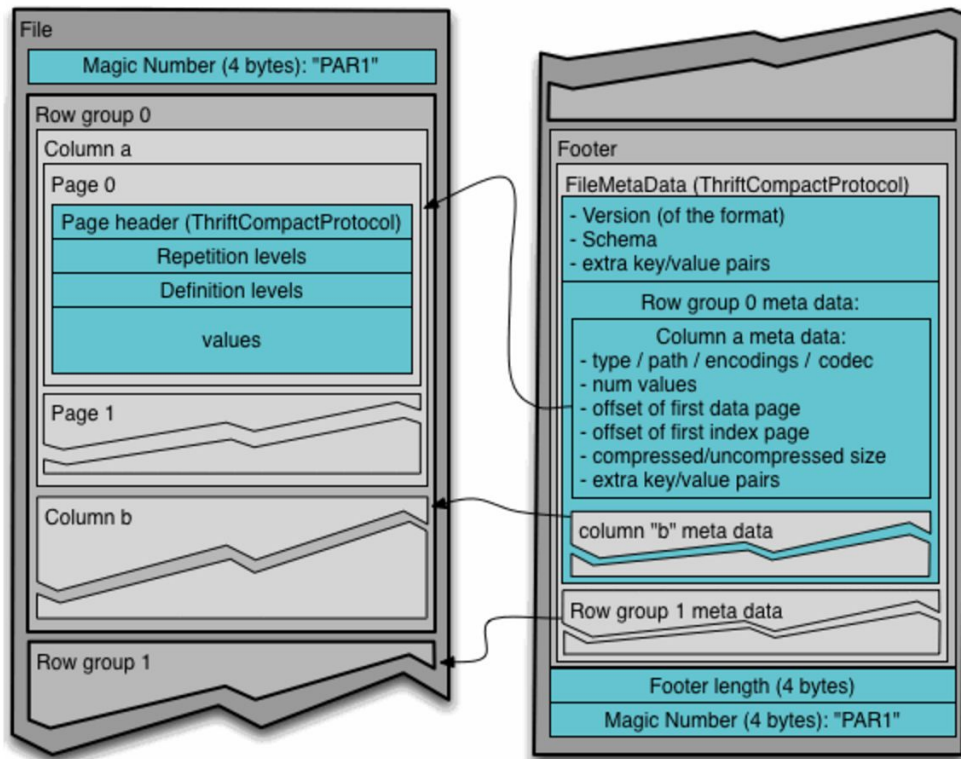
Hive file format: ORC



Is this a “good” architecture



Impala file format: Parquet

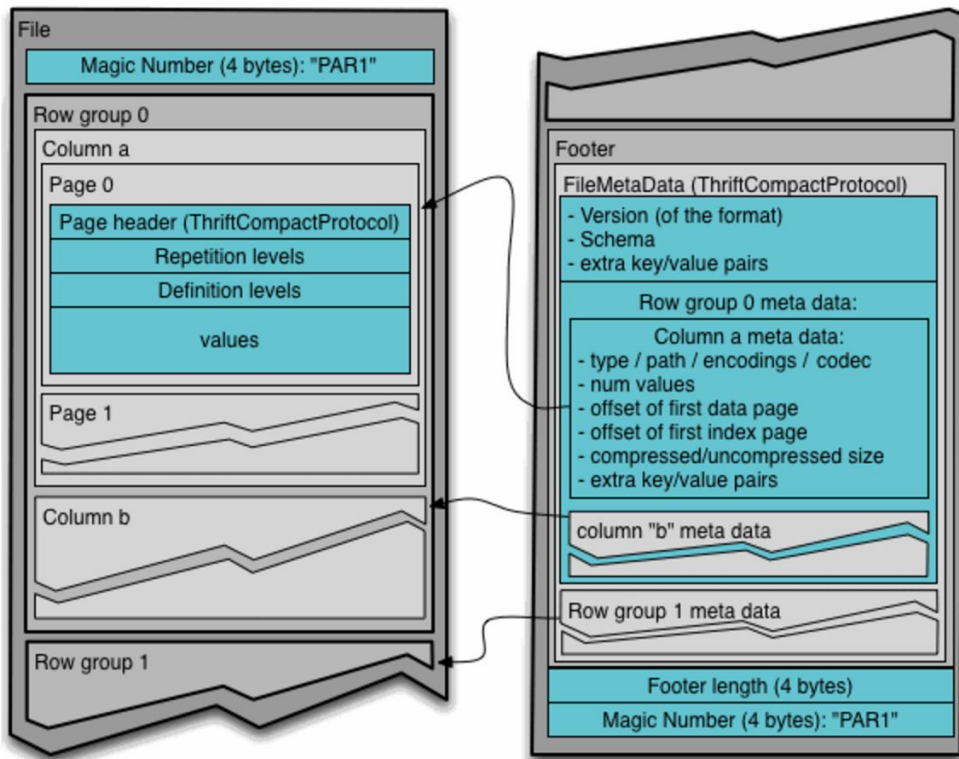


What's the big difference

Why does it matter



Impala file format: Parquet

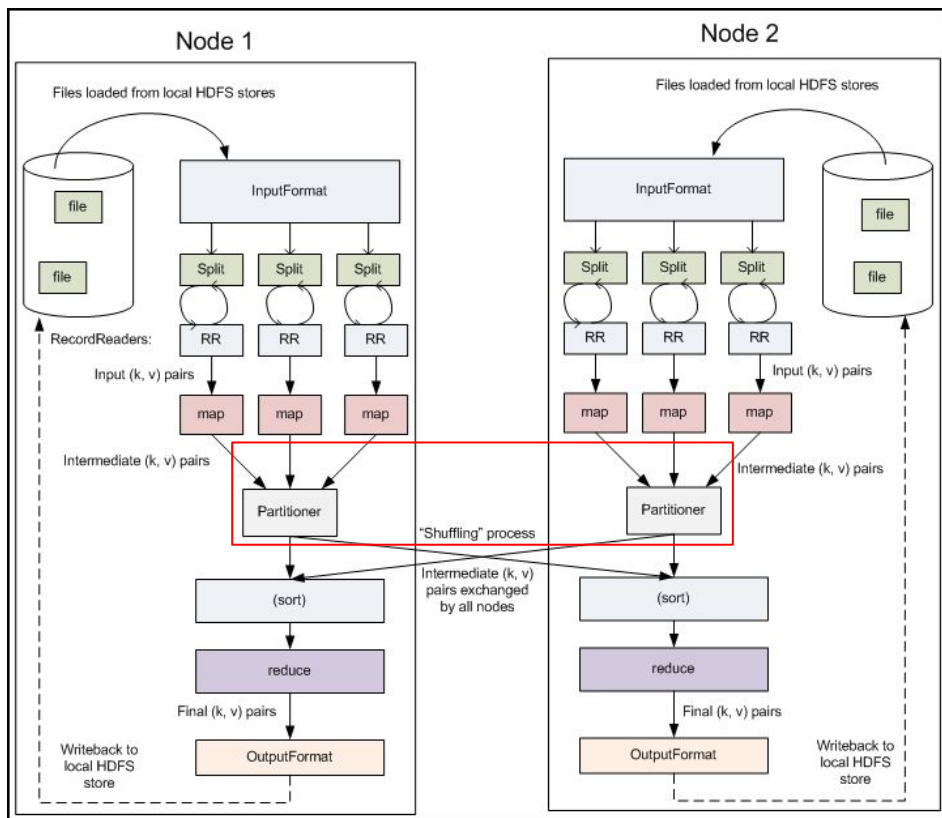


What's the big difference

Why does it matter



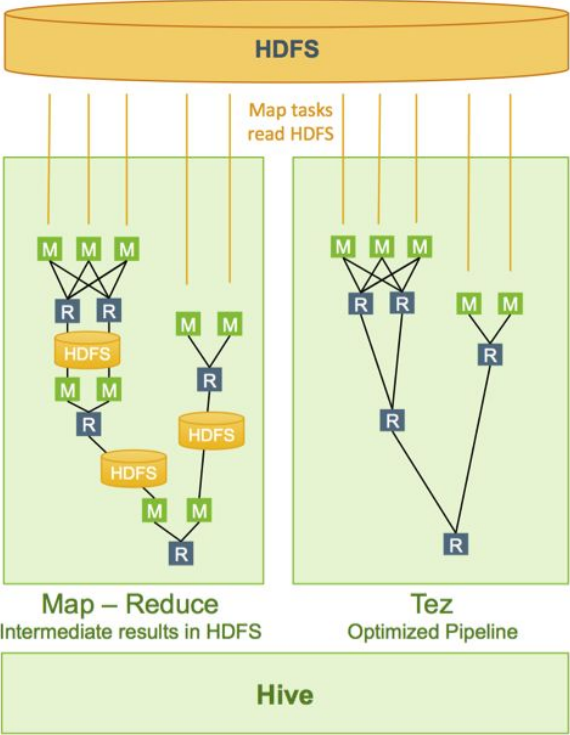
Hive runtime: MapReduce



Hive-MapReduce **materializes** intermediate results and writes to disk

Why is this bad? Why is this good?

Hive runtime: From MR to Tez

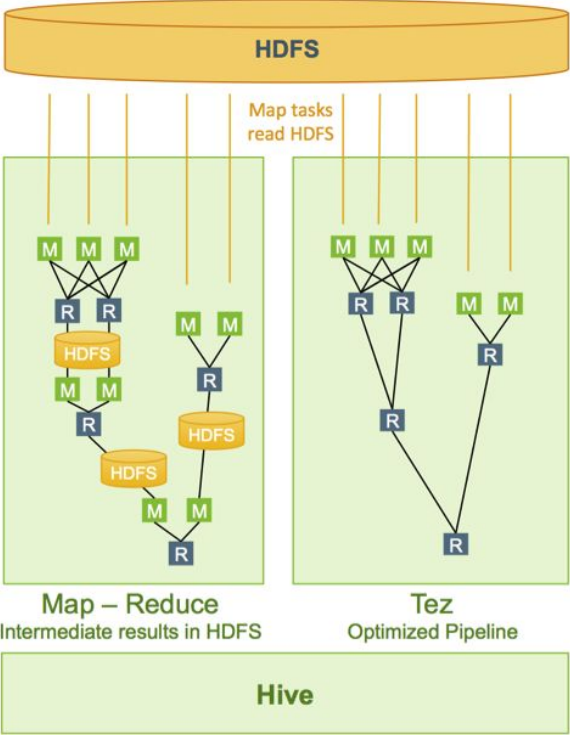


What's the big difference?

Why does it matter?



Hive runtime: From MR to Tez

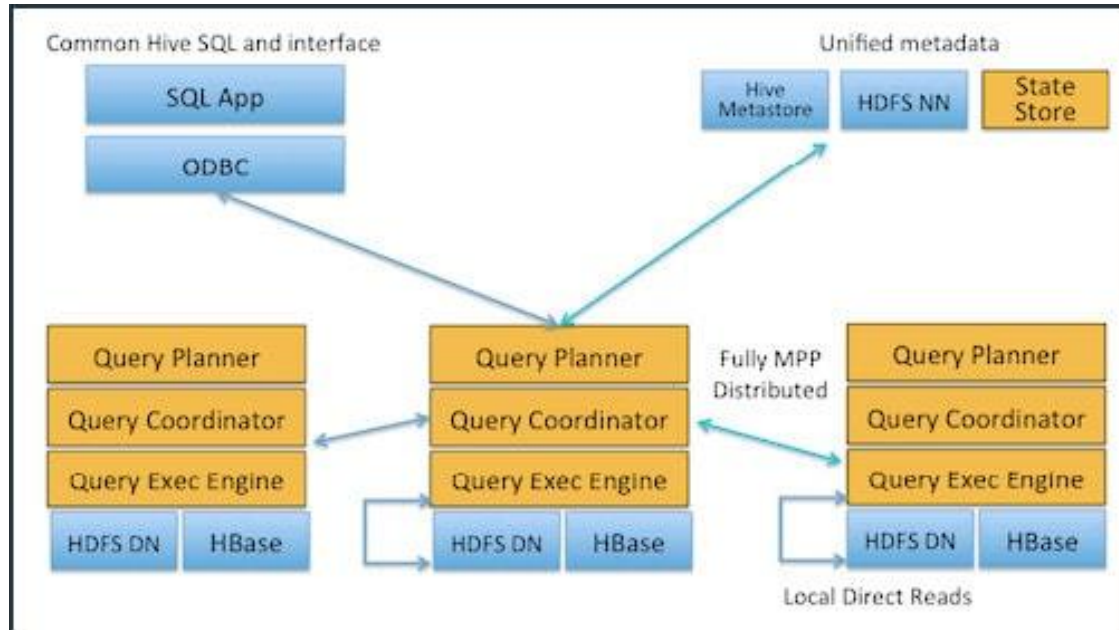


What's the big difference?

Why does it matter?



Impala runtime



Fully shared-nothing architecture with no intermediate materialization




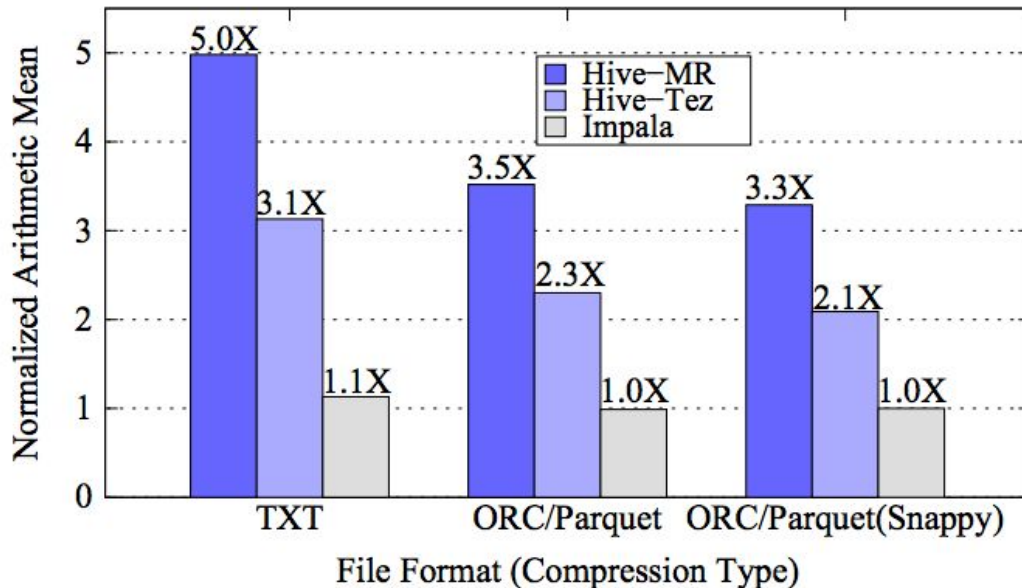
How Fast is Really Fast?

Benchmarks: Loading Time

Task: Load 1TB data

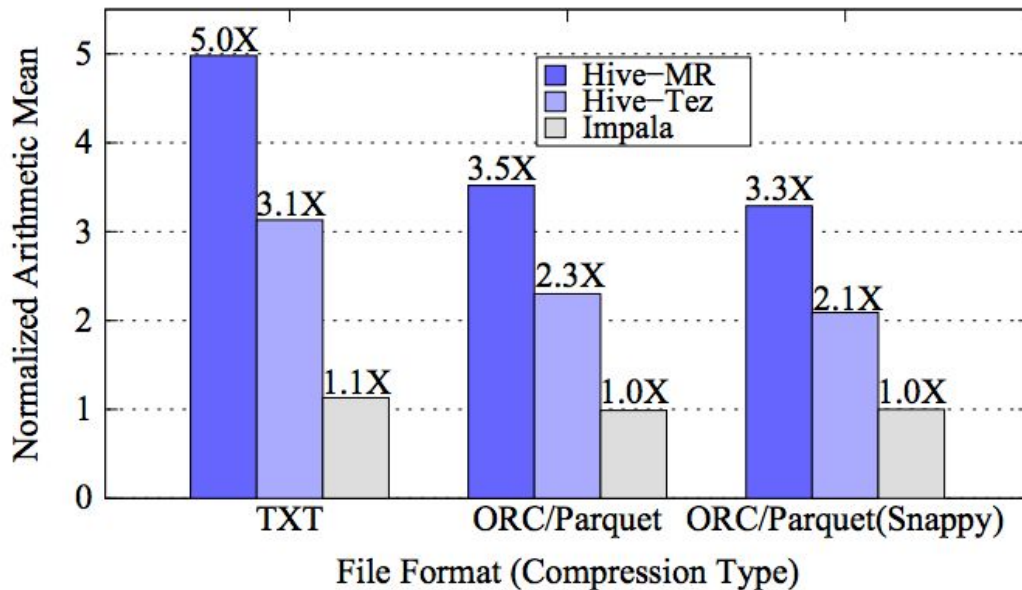
Vary: Compression
and data system

Result: 

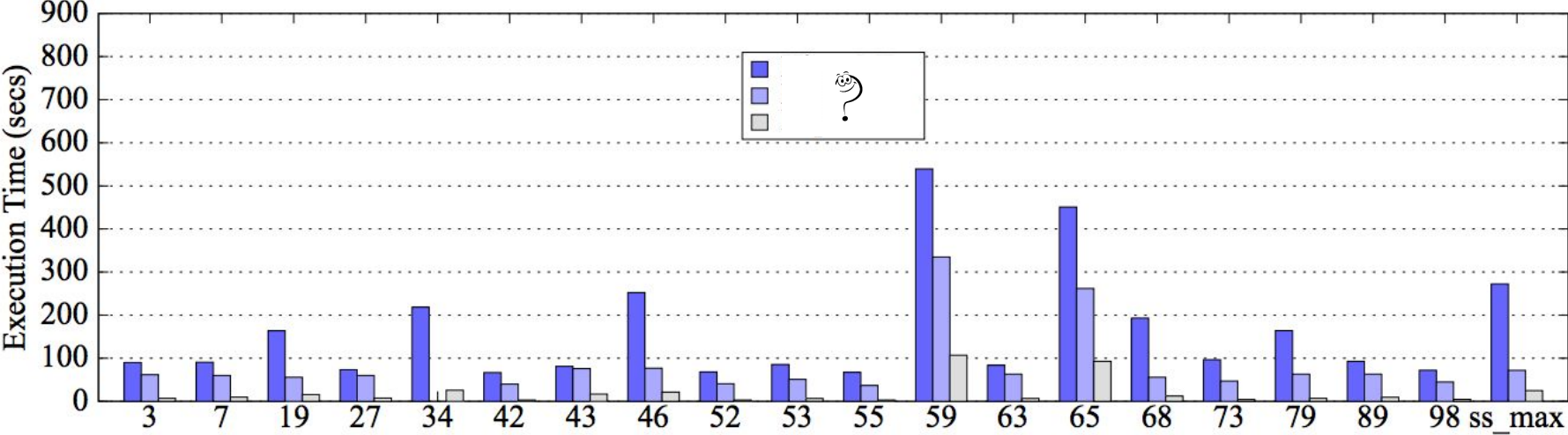


Benchmarks: Loading Time

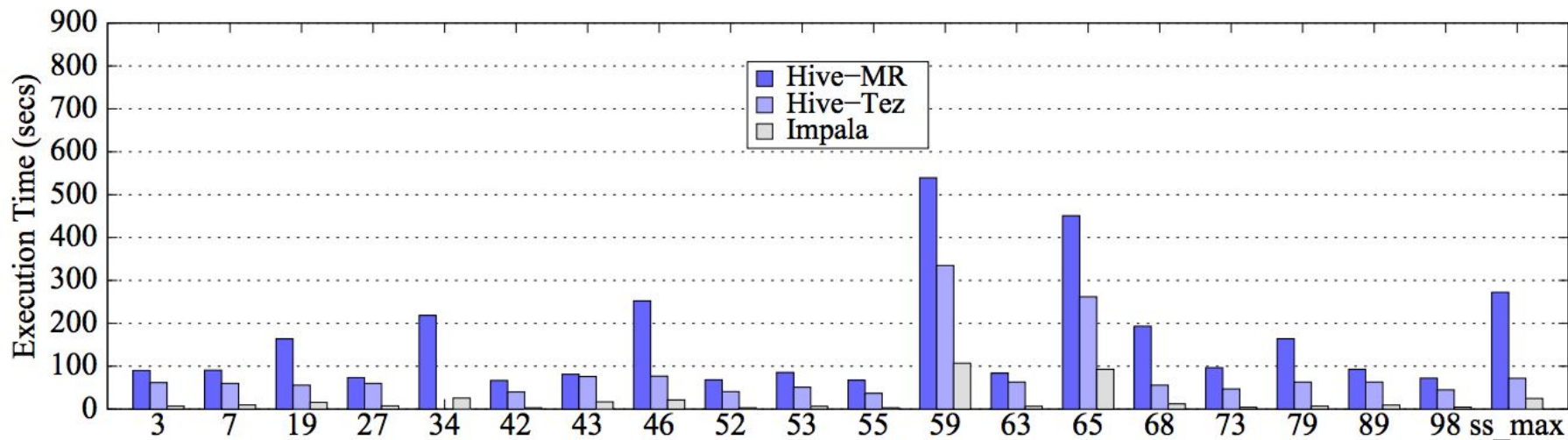
Why the difference?



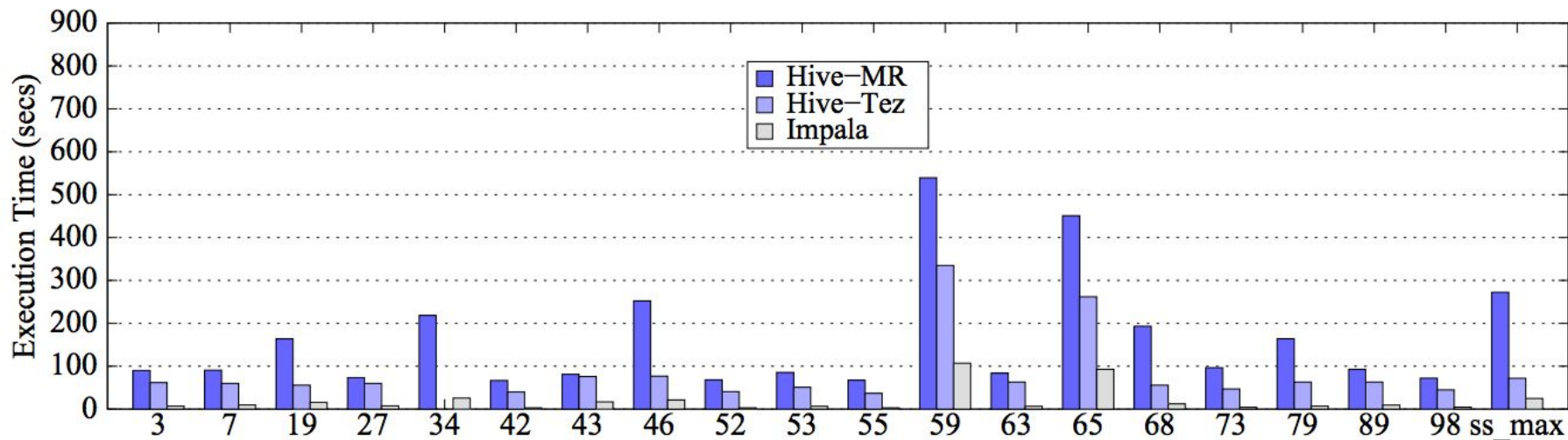
Benchmarks: Query Execution Time



Benchmarks: Query Execution Time



Benchmarks: Query Execution Time



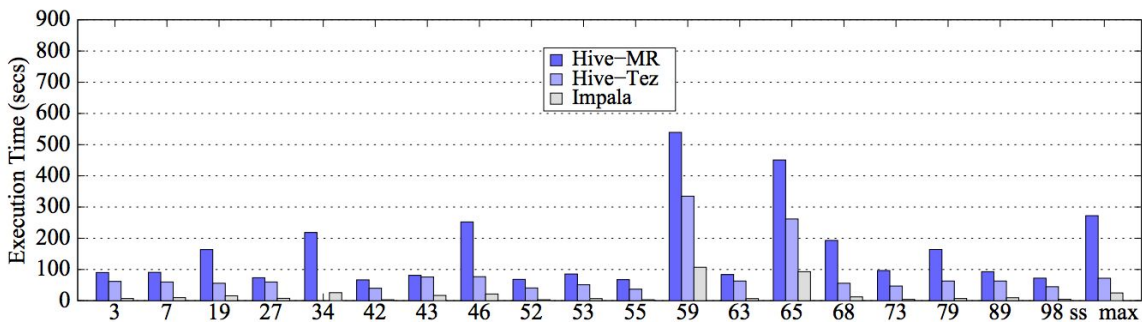
Why is Impala so much faster



Benchmarks: Query Execution Time

Quiz: which of these is responsible?

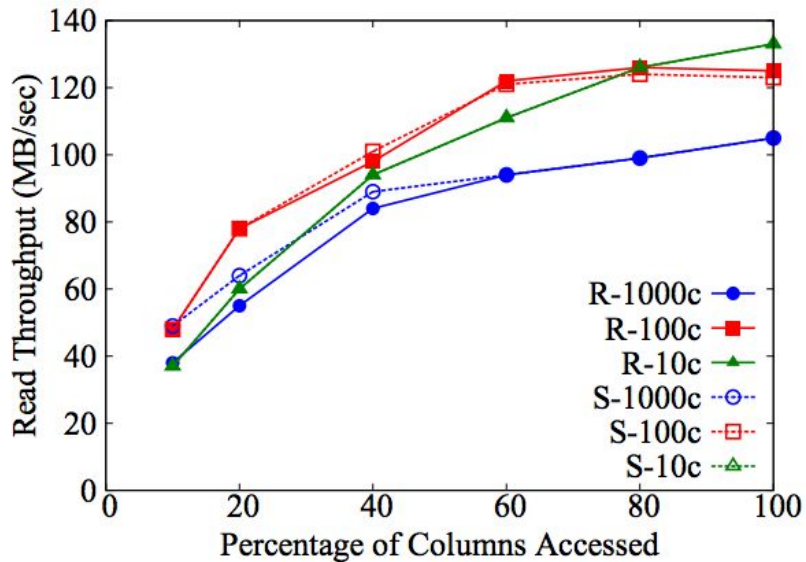
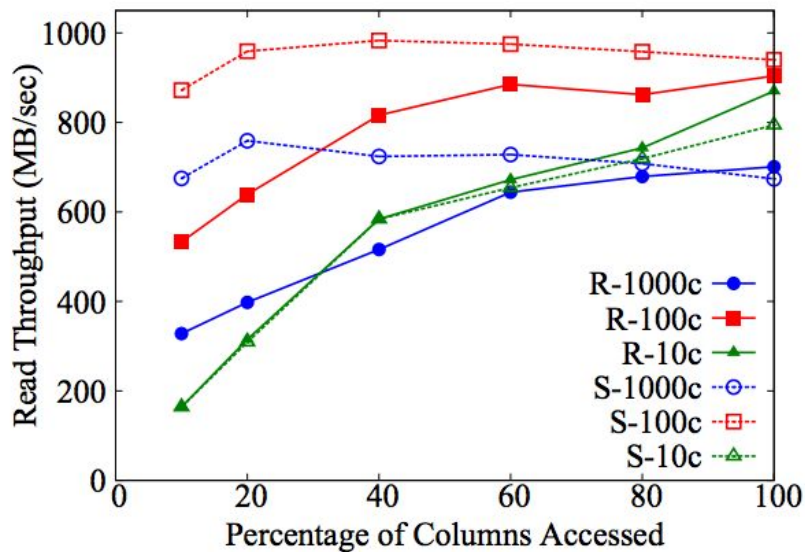
- (a) efficient I/O
- (b) no initialization overhead
- (c) pipelined rather than materialized intermediaries
- (d) magic??



Why is Impala so much faster



Benchmarks: Data Access



How similar are these graphs?



Future work

Failure recovery for Impala

Workloads that exceed the size of main memory (e.g. backpressure, or buffer intermediate results to disk)

Caching common sub-DAG query results