# SQL in Hadoop: To Boldly Go Where No Data Warehouse has Gone Before

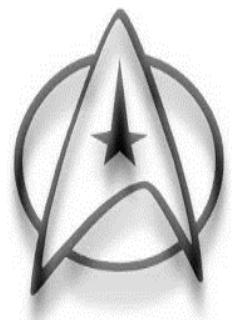Emma McGrattan

SVP Engineering, Actian Corp

# Who is Actian?

**$140M Revenues + Profitable**

**10,000+ Customers**

**Global Presence: 8 world-wide offices, 7x 24 multinational support model**

*"Actian is now very **powerfully positioned** in the big data and analytics markets."* Bloor

*"Fast becoming a **big data powerhouse** to challenge the market."* Forrester

# Actian Vortex

## Application Development and Tools

**Vortex - Analytic Workbench**

Read → Blend & Enrich → Data Science & Analytics → Load Actian

**MicroStrategy**

**COGNOS**

**+ + + tableau** SOFTWARE

**Actian Management Console**

## SOURCE DATA

Databases / Marts Warehouses

Structured & Unstructured Data

Enterprise Applications

Cloud / SaaS Applications

DATA PLATFORM

**Vector in Hadoop**

SQL Analytics

**DataFlow**

**DataFlow**

Elastic Data Prep

Predictive Analytics

**SPARQLverse**

Graph Analytics

Library of Analytic Blueprints

hadoop

**Hortonworks**

cloudera

MAPR TECHNOLOGIES

SQL | Java, C/++, Pyhtn

## ANALYTIC APPS

Financial Services

Health Care

Other Verticals

**Deployment Options**

**rackspace.** the open cloud company

**amazon** web services

# Actian Vector in Hadoop Architecture



**Client application**

SQL query     result

**SQL Processing**

- SQL parser
  - parsed tree
- Optimizer
  - query plan
- Cross compiler

X100 algebra

**Master node**

**X100**

- Distributed rewriter
  - annotated query tree
- Builder
  - operator tree
- Execution engine
  - data request    data
- Buffer manager

I/O

HDFS

HDFS namenode

MPI
annotated tree

MPI
partial result set

MPI inter-node communication

X100
X100
X100
X100

**Worker node [1..n] (datanodes)**

**X100**

- Rewriter
  - annotated query tree
- Builder
  - partial operator tree
- Execution engine
  - data request    data
- Buffer manager

I/O

HDFS

HDFS datanode

HDFS   X100
HDFS   X100
HDFS   X100
HDFS   X100
HDFS   X100

# Vector: Built for Warp Speed

**1** **Vectorized End-to-End**



Single
Instruction
Multiple
Data

**2** **Exploiting Chip Cache**



Process data on chip – not in RAM

**3** **Update Capability**


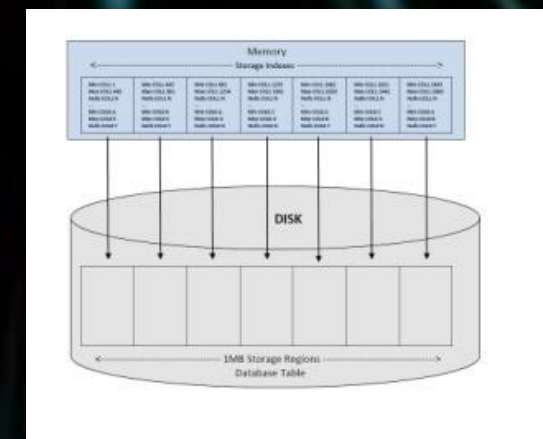
Limit I/O
Efficient real time updates
Update & Delete individual records

**4** **Smart Compression**



Maximize throughput
Vectorized decompression

**5** **Storage Indexes**
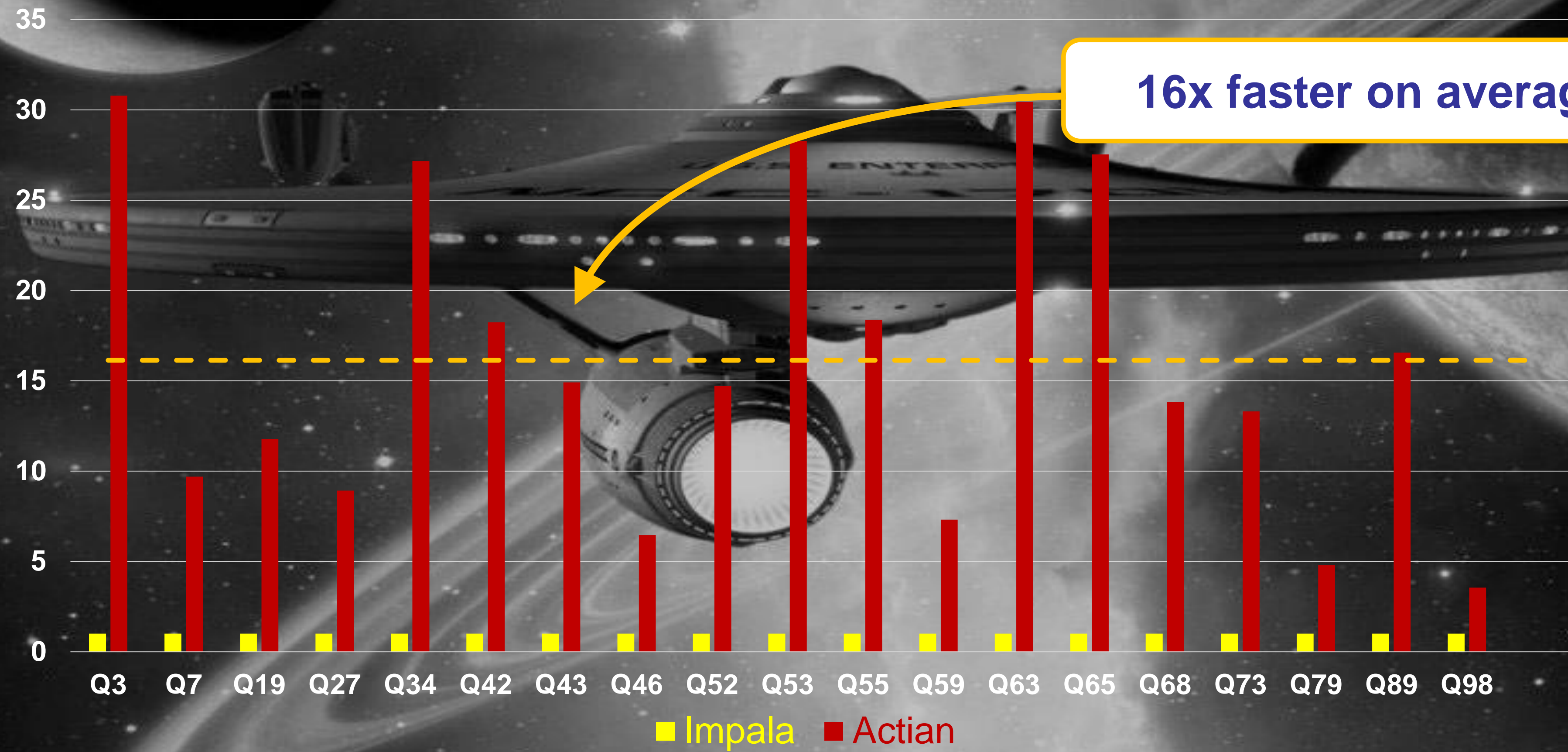


Quickly identify candidate
data blocks
Minimize I/O

**6** **YARN Integration**

Intelligent Block Placement
Dynamic Resource Management

# Vector:  To Boldly Go…

**"Impala Subset" of TPC-DS Queries at Scale Factor 3000 (3TB)**
**Speedup vs Impala**



**16x faster on average**

Number of times faster than Impala

35
30
25
20
15
10
5
0

Q3  Q7  Q19  Q27  Q34  Q42  Q43  Q46  Q52  Q53  Q55  Q59  Q63  Q65  Q68  Q73  Q79  Q89  Q98

■ Impala   ■ Actian
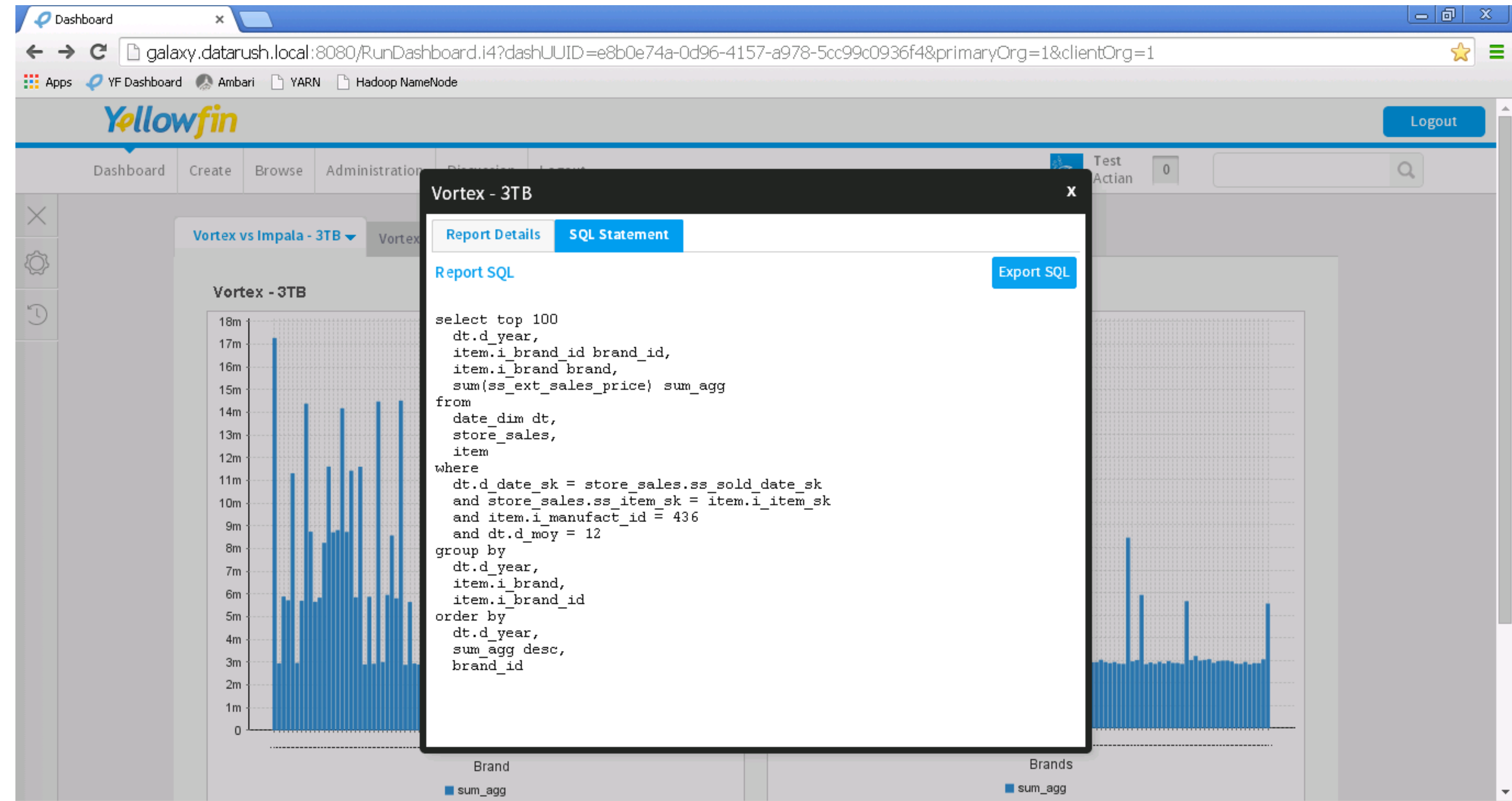
Both Executed on the same hardware and software environment:
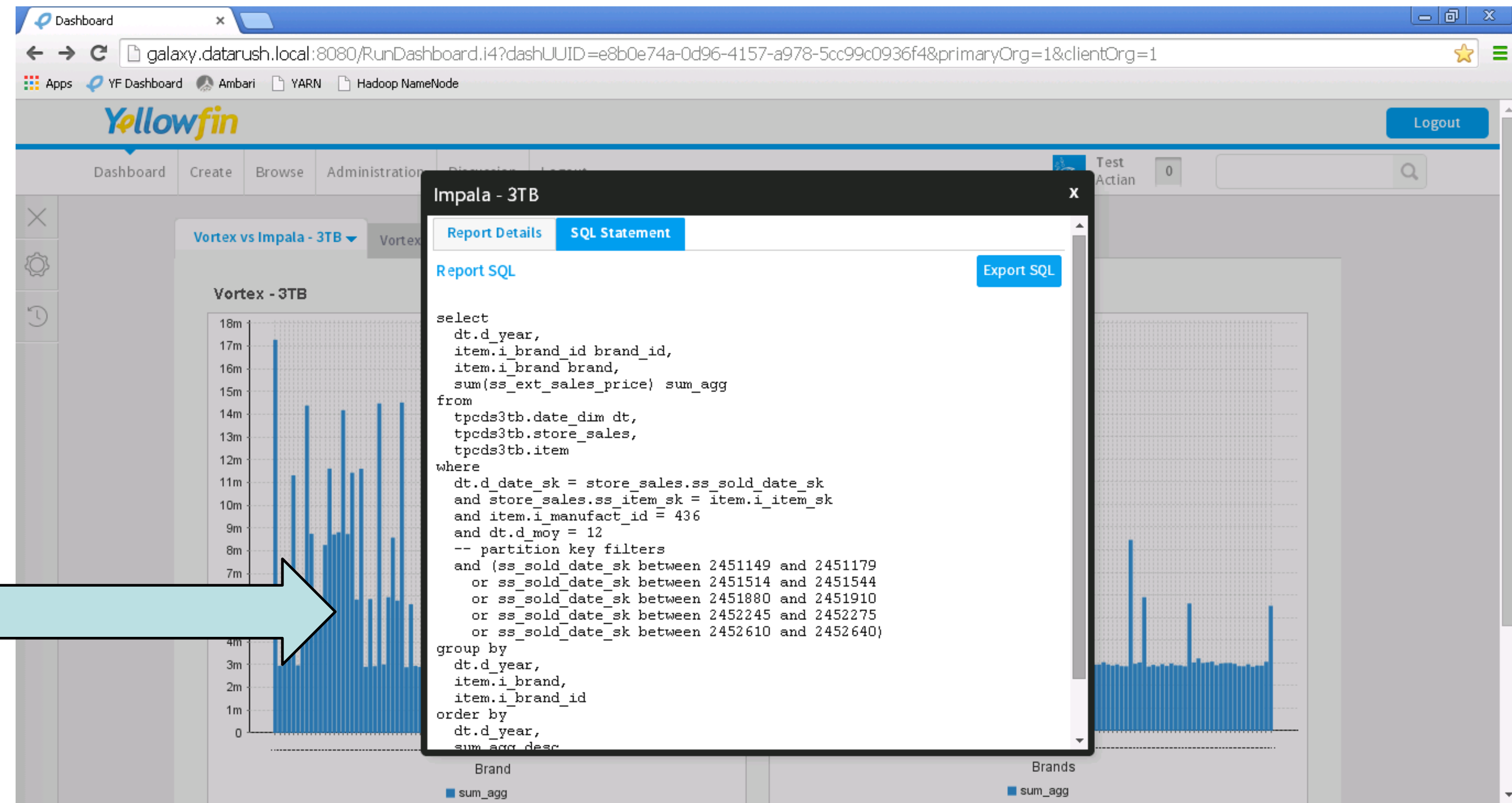5 Node Cluster with 64GB of RAM per node and 24 x 1TB Hard Disks.

# The SQL Behind the Actian Numbers

# The Impala Equivalent Uses "Hints"

Note the use of partition keys



```
select
    dt.d_year,
    item.i_brand_id brand_id,
    item.i_brand brand,
    sum(ss_ext_sales_price) sum_agg
from
    tpcds3tb.date_dim dt,
    tpcds3tb.store_sales,
    tpcds3tb.item
where
    dt.d_date_sk = store_sales.ss_sold_date_sk
    and store_sales.ss_item_sk = item.i_item_sk
    and item.i_manufact_id = 436
    and dt.d_moy = 12
    -- partition key filters
    and (ss_sold_date_sk between 2451149 and 2451179
        or ss_sold_date_sk between 2451514 and 2451544
        or ss_sold_date_sk between 2451880 and 2451910
        or ss_sold_date_sk between 2452245 and 2452275
        or ss_sold_date_sk between 2452610 and 2452640)
group by
    dt.d_year,
    item.i_brand,
    item.i_brand_id
order by
    dt.d_year,
    sum_agg desc
```

CAPTAIN JAMES T. KIRK

I'M SORRY, I CAN'T HEAR YOU OVER THE
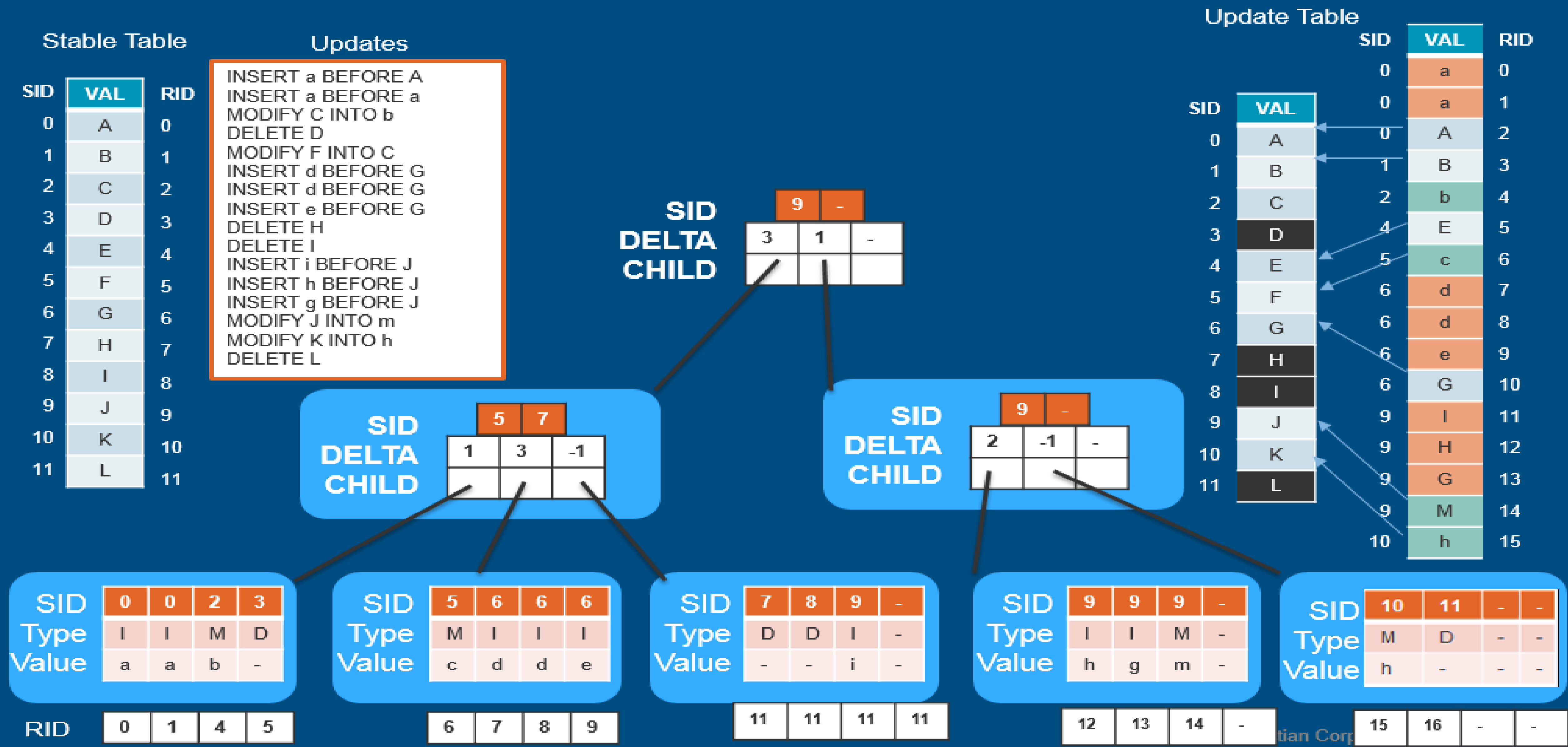SOUND OF HOW AWESOME I AM.

# Trickle Update Support

- The Kobayashi Maru of Hadoop

  - The design paradigm for HDFS is for data to be written once and read ever after.

  - Appending updated records to the end of a column/table or rewriting the entire table significantly impacts system performance

- The Solution – Positional Delta Trees

  - Enable on-line updates, without impacting read performance

  - Keep track of the tuple position of Inserts/Modifies/Deletes

  - Designed to make merging in of these updates fast by providing the tuple positions where differences have to be applied at update time.

# Positional Delta Trees

# Data Security

- Access Control

- Role Separation
  - System Administrator & Database Administrator should not have access to all data

- Security Auditing
  - Ability to audit who accessed, or attempted to access, what and when

- Encryption
  - Data at rest – minimize performance impact by enabling at column level
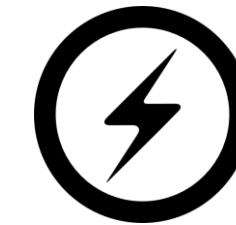  - Data in motion
  - File system

# It's SQL in Hadoop, Jim, but not as we know it

## Highest Performing and Fully Industrialized SQL in Hadoop

*Full ANSI SQL 92 support* – enables use of ALL standard BI tools and apps

*Fully ACID compliant* – brings transactional integrity to Hadoop to prevent inaccurate results

*Hadoop distribution agnostic -* avoids vendor lock-in and provides customer flexibility

*Update Capability* – provides ability to update without impacting read performance

*Highest Concurrency* – allows your customers to have simultaneous users and tasks run without long wait times

*Highly Performant –* up to 30x faster than our closest competitor, Impala

*Native DBMS Security* - authentication, user and role-based security, data protection, and encryption

*Mature, proven planner and fastest optimizer* ensures customers can maximize number of nodes, CPU, memory and cache

*Native in-Hadoop YARN* – manage Hadoop resources automatically to prevent inefficiencies

*Collaborative architecture* - query native Hadoop file formats (like Parquet) without ingestion

*Open APIs -* allow read access to our block format

# Beam it Down, Scotty!