



# Hive, Impala, and Spark, Oh My: SQL-on-Hadoop in Cloudera 5.5

Justin Erickson | Director of Product Management |  
Cloudera



# Agenda

- History of SQL-on-Hadoop technologies
- Picking the right tool for the job
- What's new with Cloudera 5.5
- Real-world use cases
- Future of SQL-on-Hadoop

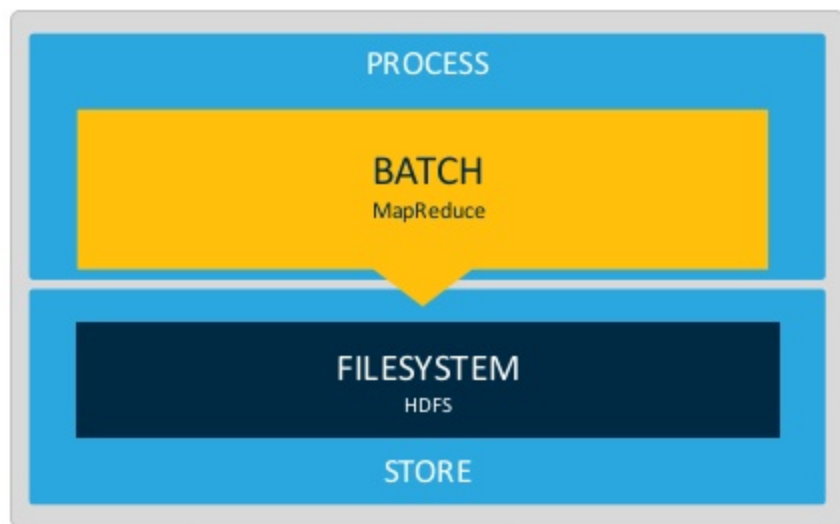
# MapReduce: The Early Years

The original processing engine for Hadoop

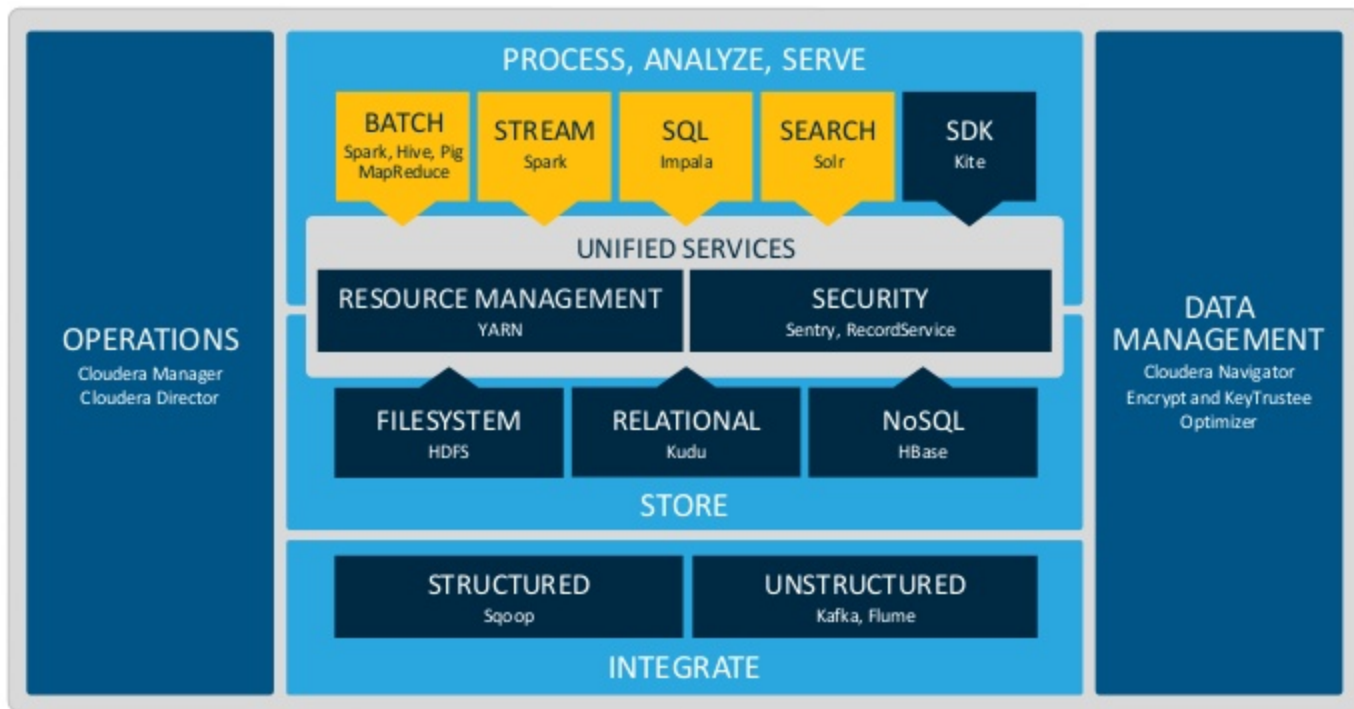
- Process any type of data in any format
- Scale infinitely for multiple, large jobs
- Pioneer of bringing compute to data

*But...*

- Difficult to program
- Slow processing
- Limited expressivity



# One Platform, Many Workloads



## Batch, Interactive, and Real-Time.

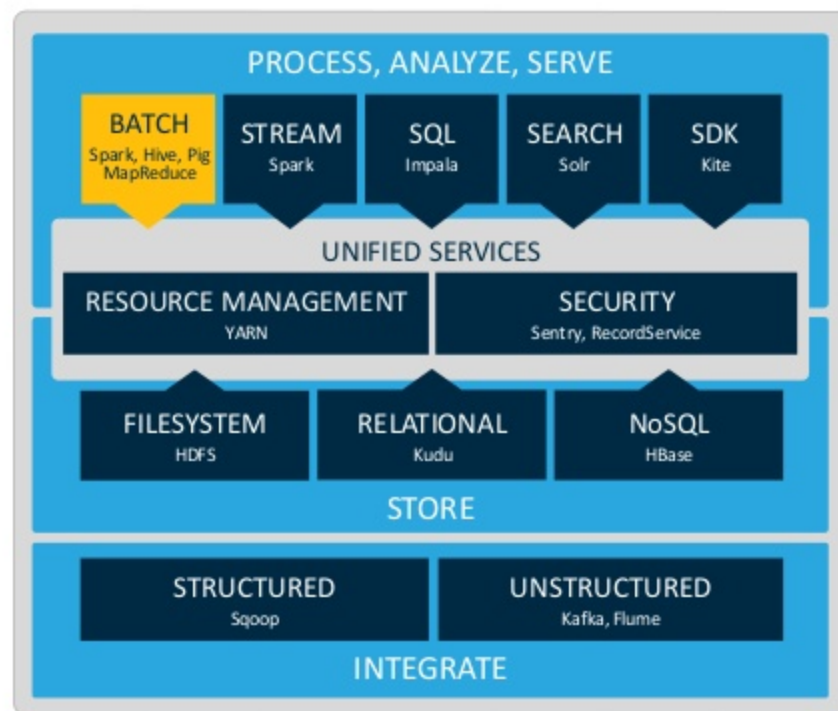
Leading performance and usability in one platform.

- End-to-end analytic workflows
- Access more data
- Work with data in new ways
- Enable new users

# The Need for SQL for Batch Processing

## Apache Hive

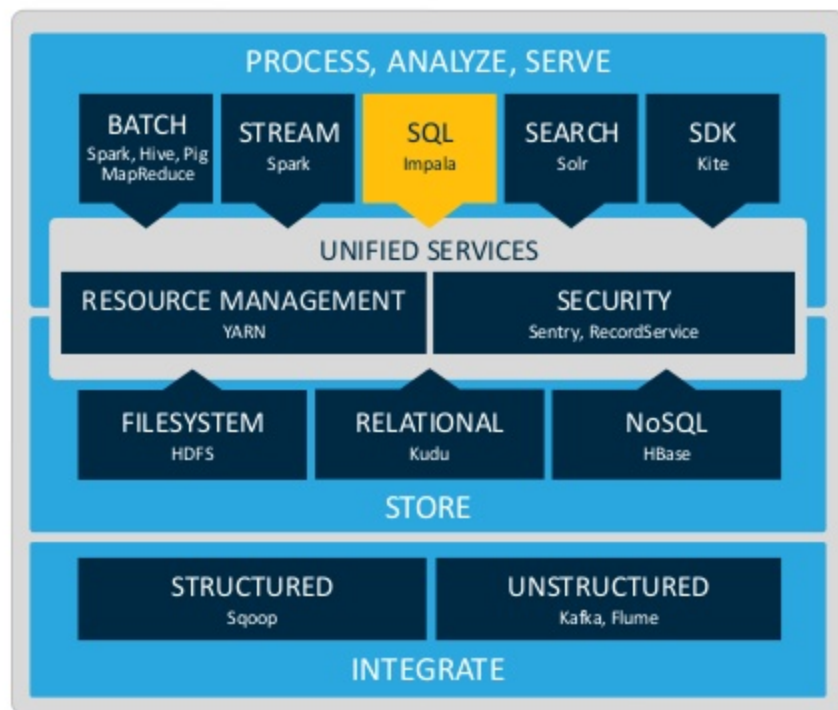
- Eases development on MapReduce with familiar SQL
- Built for long-running ETL, data preparation, and batch processing
- Shared data structures across Hadoop tools



# The Need for Interactive SQL for BI

## Apache Impala (incubating)

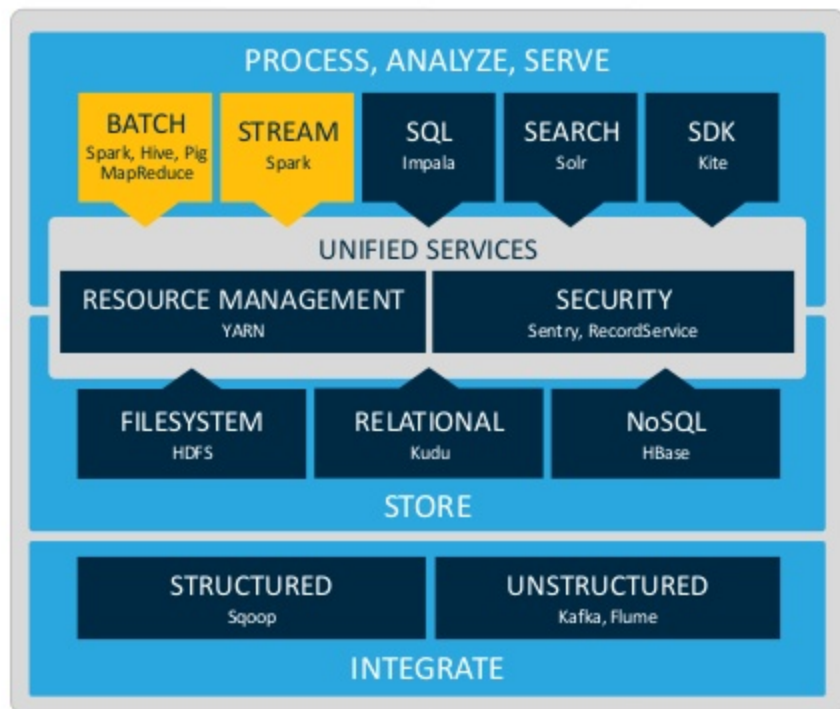
- Low latency for interactive performance
- Built for multi-user workloads
- Compatible with SQL and leading BI partner tools



# The Need for Flexible Data Processing

## Apache Spark (and Spark SQL)

- Easy development
- Flexible, extensible API across multiple workload types
- In-memory batch and stream processing performance boost





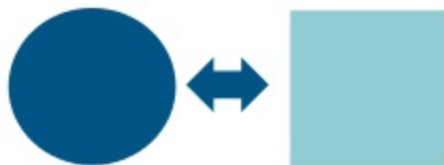
# Focus on Open Source Standards

Open source does not guarantee a future-proof investment



## Long-Term Architecture

Only open standards get continuing, long-term investment from across the ecosystem.



## Avoidance of Lock-in

Open standards have multi-vendor support, giving customers choices and preventing lock-in.



## Ecosystem Compatibility

Open standards attract more third-party connectors/certifications due to broad adoption.



# Choosing the Right SQL Engine

Know Your Audience, Know Your Use Case



Batch  
Processing



BI and  
SQL Analytics



Procedural  
Development

# SQL-on-Hadoop in Cloudera 5.5

	Apache Hive	Apache Impala (incubating)	Apache Spark SQL
Audience	ETL Developers	Business Analysts	Data Engineers & Data Scientists
Strengths	<ul style="list-style-type: none"><li>• Built for very long-running ETL, data preparation, or batch processing</li><li>• Supports custom file formats</li><li>• Handles massive ETL sorts with joins</li></ul>	<ul style="list-style-type: none"><li>• Scales to high-concurrency</li><li>• Supports high-performance interactive SQL</li><li>• Compatible with BI tools &amp; skills</li><li>• Hadoop integration &amp; usability</li></ul>	<ul style="list-style-type: none"><li>• Easily embed SQL into Java, Scala, or Python applications</li><li>• Simple language for common operations</li><li>• Seamlessly mix SQL and Spark code within a single application</li></ul>
New Features	<ul style="list-style-type: none"><li>• Hive in the cloud (S3)</li><li>• Hive-on-Spark beta</li><li>• Governance &amp; Lineage</li></ul>	<ul style="list-style-type: none"><li>• Nested data types</li><li>• Column-level security</li><li>• Integration with Kudu (beta)</li></ul>	<ul style="list-style-type: none"><li>• Support for Spark SQL &amp; DataFrames</li><li>• Hive integration</li><li>• Automatic performance optimizations</li></ul>

# SQL-on-Hadoop Benchmark

Impala, Spark SQL, Hive-on-Tez

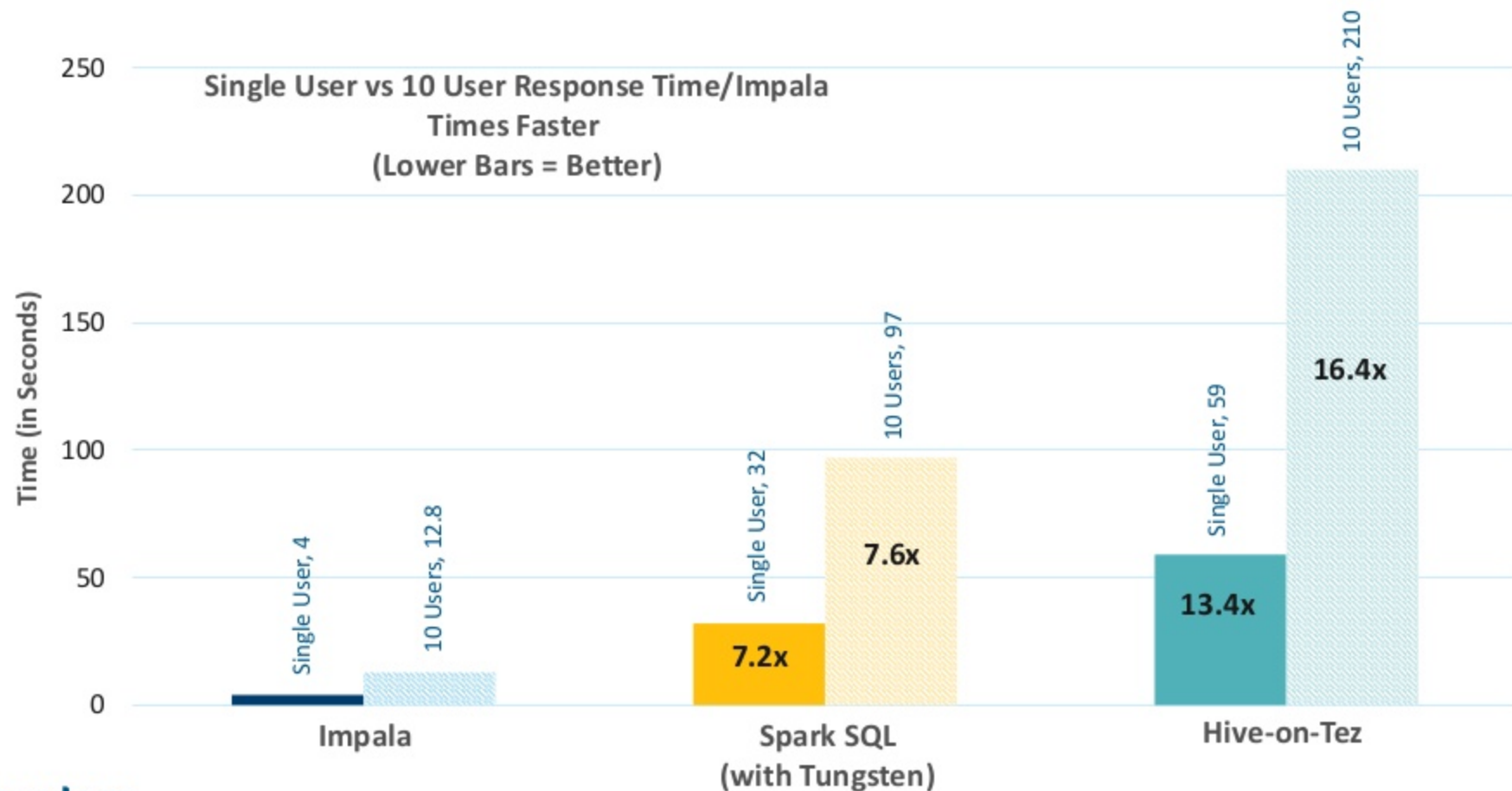
## Versions:

- Impala 2.3
- Hive 2.0 on Tez 0.5.2 (aka “Stinger”)
- Spark SQL 1.5 with Tungsten

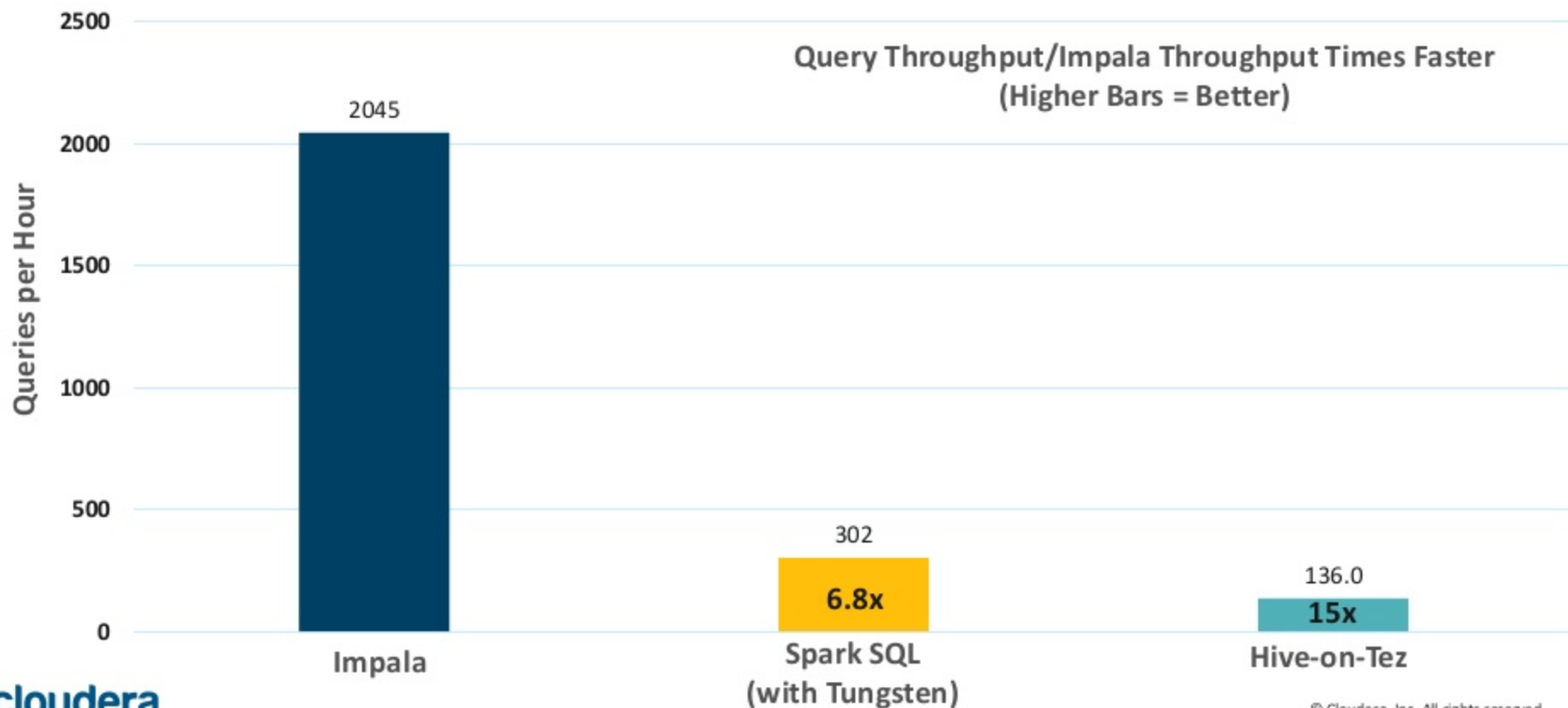
## Benchmark Details

- Based on industry standards (TPC)
- Repeatable
- Methodical testing with multiple runs on same hardware
- Help competing software do well
  - Run on optimal file formats for each
  - Tune query engines appropriately

# Impala Multi-User Performance Over 7x Faster



# Impala Enables Nearly 7x Throughput



# Performance Benchmark Takeaways

- **Impala unlocks BI usage directly on Hadoop**
  - Meets BI **low-latency** and **multi-user** requirements
  - Advantage expands for single-user vs just 10 users
- **Hive is designed (and still great) for batch processing**
  - Most Impala customers use Hive for data preparation
  - Hive is the most commonly used ETL framework
- **Spark SQL enables easier Spark application development**
  - Enables mixed procedural Spark (Java/Scala) and SQL job development
- **Mid-term trends will further favor Impala's design approach for latency and concurrency**
  - More data sets **move to memory** (HDFS caching, in-memory joins, Intel joint roadmap)
  - **CPU efficiency** will increase in importance
  - Native code enables **easy optimizations for CPU** instruction sets
  - **Intel joint roadmap** support these opportunities

# Use Cases



## PROBLEM

Needed to efficiently collect, process, and analyze data from growing hospital network

- EDW couldn't meet scale and unstructured data demands
- Processing too slow for actionable decisions
- Limited, time consuming supply chain matching

## SOLUTION

Integrated 1000s of hospital systems through unified enterprise data hub

- Ingest and process **45% more spend data**
- **Faster analytics on \$41B** through end-user healthcare spend dashboard
- Unprecedented **matching of 98%** of supply chain data
- **Better TCO** through unification and licensing costs for **new opportunities**



## PROBLEM

Clients had limited insights to thousands of marketing campaigns across channels

- Clients want real-time campaign updates with 3-sec SLA
- Existing system couldn't meet scaling or data type demands
- Limited self-service BI

## SOLUTION

Built next-generation digital marketing platform for 360-degree customer view

- **Improved query performance** from minutes to seconds to meet SLAs
- **Enhanced modeling** with combined online and offline data
- Real-time optimizations through **interactive, self-service access**



## PROBLEM

Couldn't support data integration across 20+ brands

- Existing systems couldn't scale for data consolidation
- Siloed access based on workload
- No real-time data ingestion or access

## SOLUTION

Brought all data directly to the business to lower costs and open up new use cases fast

- **Reduced TCO by 50%** by consolidating over 1PB of data, adding 200M rows daily
- Enabled **real-time vs hourly** updates on ad performance
- **Optimized inventory management** through data matching and consolidation





# Impala Roadmap

## 2H 2015

- SQL Support & Usability
  - Nested structures
  - Kudu updates (beta)
- Management & Security
  - Record reader service (beta)
  - Finer-grained security (Sentry)
- Integration
  - Isilon support
  - Python interface (Ibis)
- Performance & Scale
  - Improved predictability under concurrency

## 1H 2016

- Performance & Scale
  - Continued scalability and concurrency
  - Initial perf/scale improvements
- Management & Security
  - Improved admission control
  - Resource utilization and showback
- SQL Support & Usability
  - Dynamic partitioning
  - Improved timestamp compatibility

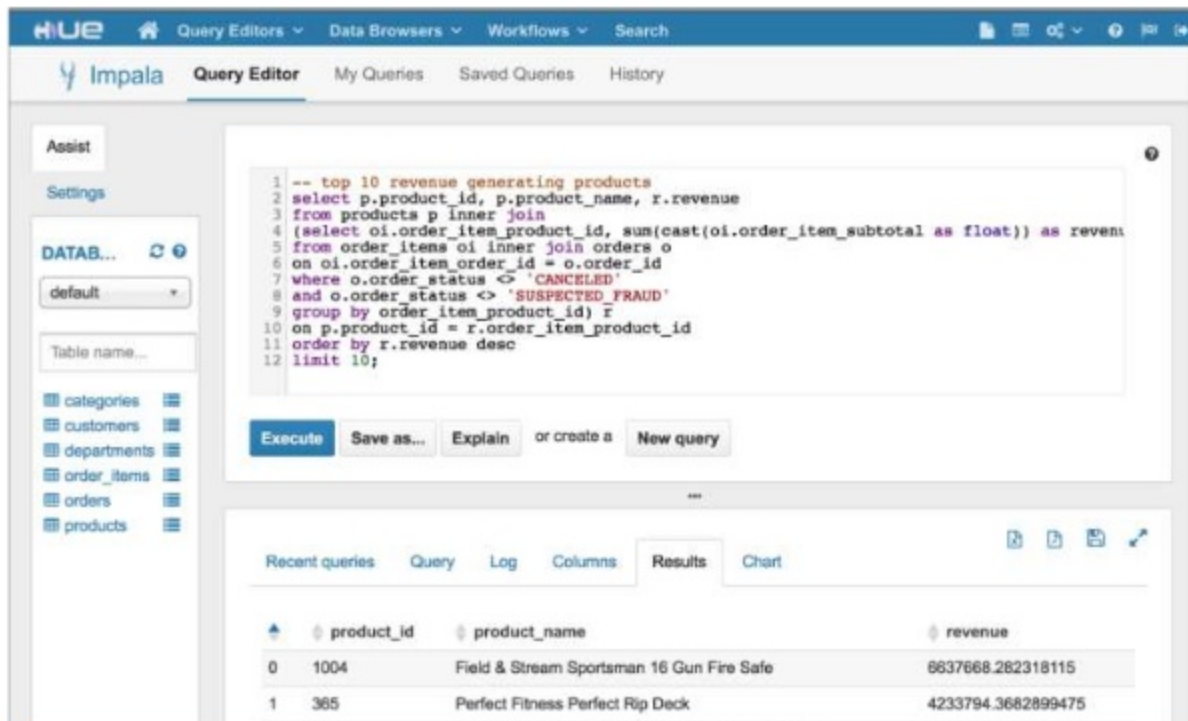
## 2016

- Performance & Scale
  - >20x performance
  - Multi-threaded joins/aggregations
  - Continued scale work
- Management & Security
  - Improved YARN integration
  - Automated metadata
- Integration
  - S3 support
- SQL Support & Usability
  - Nested types with Avro
  - Date type
  - Added SQL extensions

# Download Cloudera 5.5

[cloudera.com/downloads](https://cloudera.com/downloads)

# Try It With Cloudera Live



The screenshot shows the Cloudera Hue web interface. The top navigation bar includes 'HUE', 'Query Editors', 'Data Browsers', 'Workflows', and 'Search'. Below this, the 'Impala Query Editor' is active, showing a SQL query to find the top 10 revenue-generating products. The left sidebar has 'Assist' and 'Settings' tabs, and a 'DATAB...' section with a dropdown menu and a table name input field. The main area displays the SQL query and buttons for 'Execute', 'Save as...', 'Explain', 'or create a', and 'New query'. Below the query editor, the 'Results' tab is selected, showing a table with columns 'product\_id', 'product\_name', and 'revenue'.

```
1 -- top 10 revenue generating products
2 select p.product_id, p.product_name, r.revenue
3 from products p inner join
4 (select oi.order_item_product_id, sum(cast(oi.order_item_subtotal as float)) as revenue
5 from order_items oi inner join orders o
6 on oi.order_item_order_id = o.order_id
7 where o.order_status <> 'CANCELLED'
8 and o.order_status <> 'SUSPECTED_FRAUD'
9 group by order_item_product_id) r
10 on p.product_id = r.order_item_product_id
11 order by r.revenue desc
12 limit 10;
```

	product_id	product_name	revenue
0	1004	Field & Stream Sportsman 16 Gun Fire Safe	6637668.282318115
1	365	Perfect Fitness Perfect Rip Deck	4233794.3682899475

Featuring tutorials on:

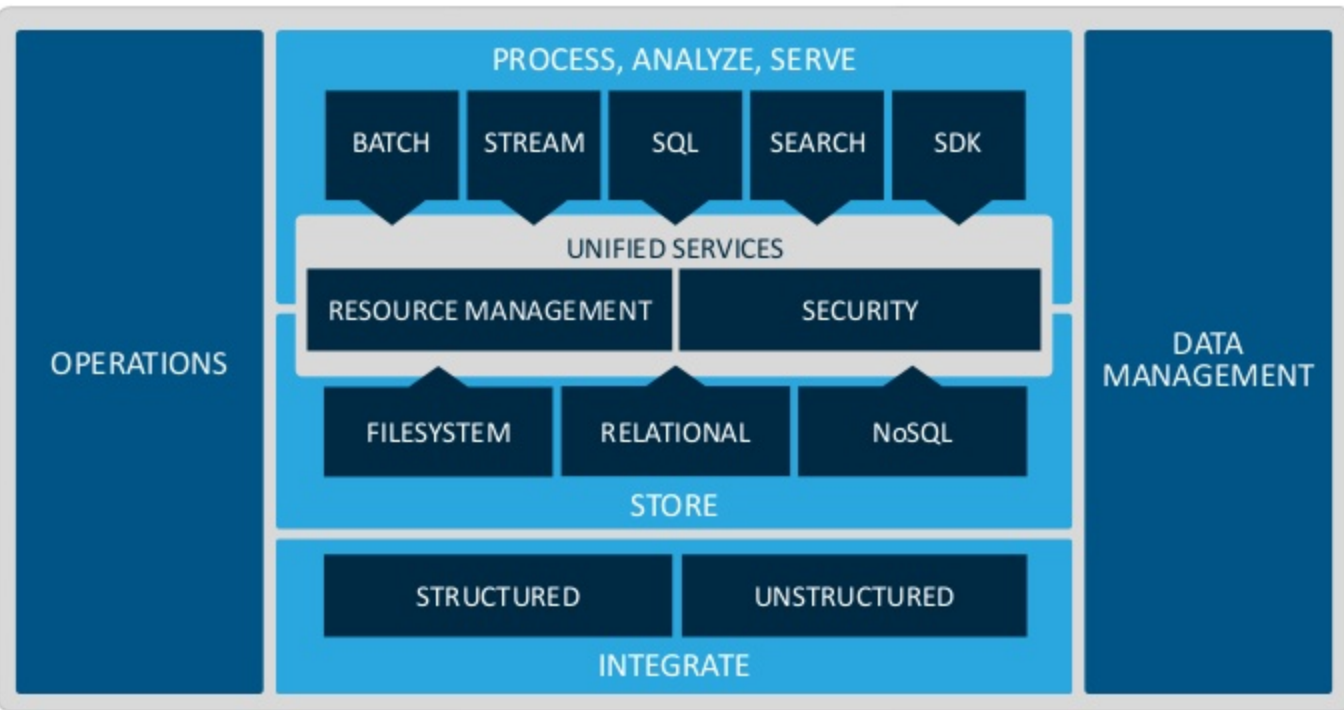
**cloudera®**



[cloudera.com/live](https://cloudera.com/live)

# Cloudera Enterprise

Making Hadoop Fast, Easy, and Secure



## A new kind of data platform:

- One place for unlimited data
- Unified, multi-framework data access

## Cloudera makes it:

- **Fast** for business
- **Easy** to manage
- **Secure** without compromise





**cloudera**  
Thank You!