# SQL on everything, in memory

**Apache Calcite**

Julian Hyde

Strata, NYC

October 16th, 2014

Hortonworks

# About me

**Julian Hyde**

**Architect at Hortonworks**

**Open source:**

- Founder & lead, Apache Calcite (query optimization framework)
- Founder & lead, Pentaho Mondrian (analysis engine)
- Committer, Apache Drill
- Contributor, Apache Hive
- Contributor, Cascading Lingual (SQL interface to Cascading)

**Past:**

- SQLstream (streaming SQL)
- Broadbase (data warehouse)
- Oracle (SQL kernel development)

Hortonworks

# SQL: in and out of fashion

**1969 — CODASYL (network database)**

**1979 — First commercial SQL RDBMSs**

**1990 — Acceptance — transaction processing on SQL**

**1993 — Multi-dimensional databases**

**1996 — SQL EDWs**

**2006 — Hadoop and other "big data" technologies**

**2008 — NoSQL**

**2011 — SQL on Hadoop**

**2014 — Interactive analytics on {Hadoop, NoSQL, DBMS}, using SQL**

**SQL remains popular.**

**But why?**

# "SQL inside"

**Implementing SQL well is hard**

- System cannot just "run the query" as written
- Require relational algebra, query planner (optimizer) & metadata

**…but it's worth the effort**

**Algebra-based systems are more flexible**

- Add new algorithms (e.g. a better join)
- Re-organize data
- Choose access path based on statistics
- Dumb queries (e.g. machine-generated)
- Relational, schema-less, late-schema, non-relational (e.g. key-value, document)

# Apache Calcite

# Apache Calcite

**Apache incubator project since May, 2014**

- Originally named Optiq

**Query planning framework**

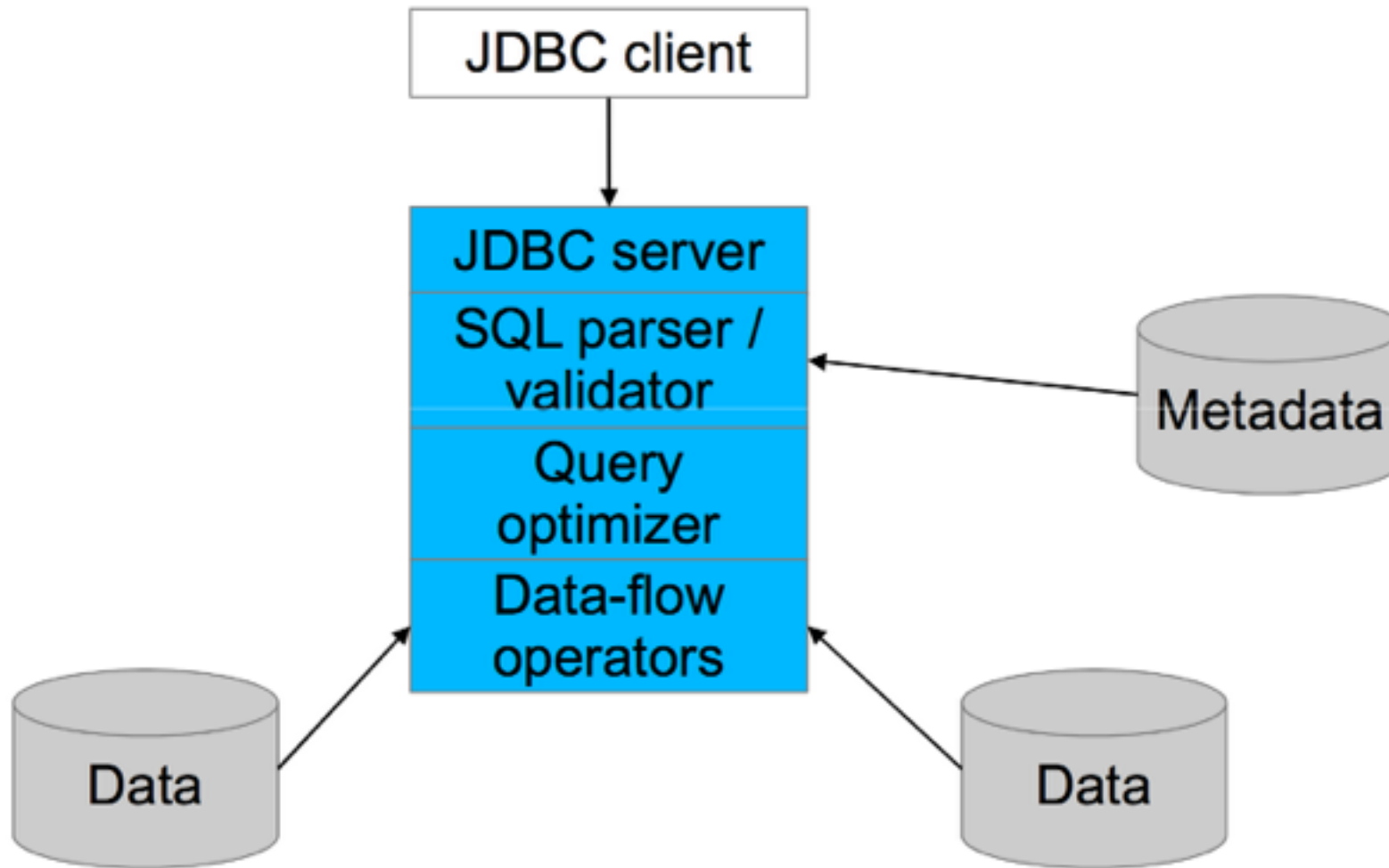- Relational algebra, rewrite rules, cost model

- Extensible

**Packaging**

- Library (JDBC server optional)
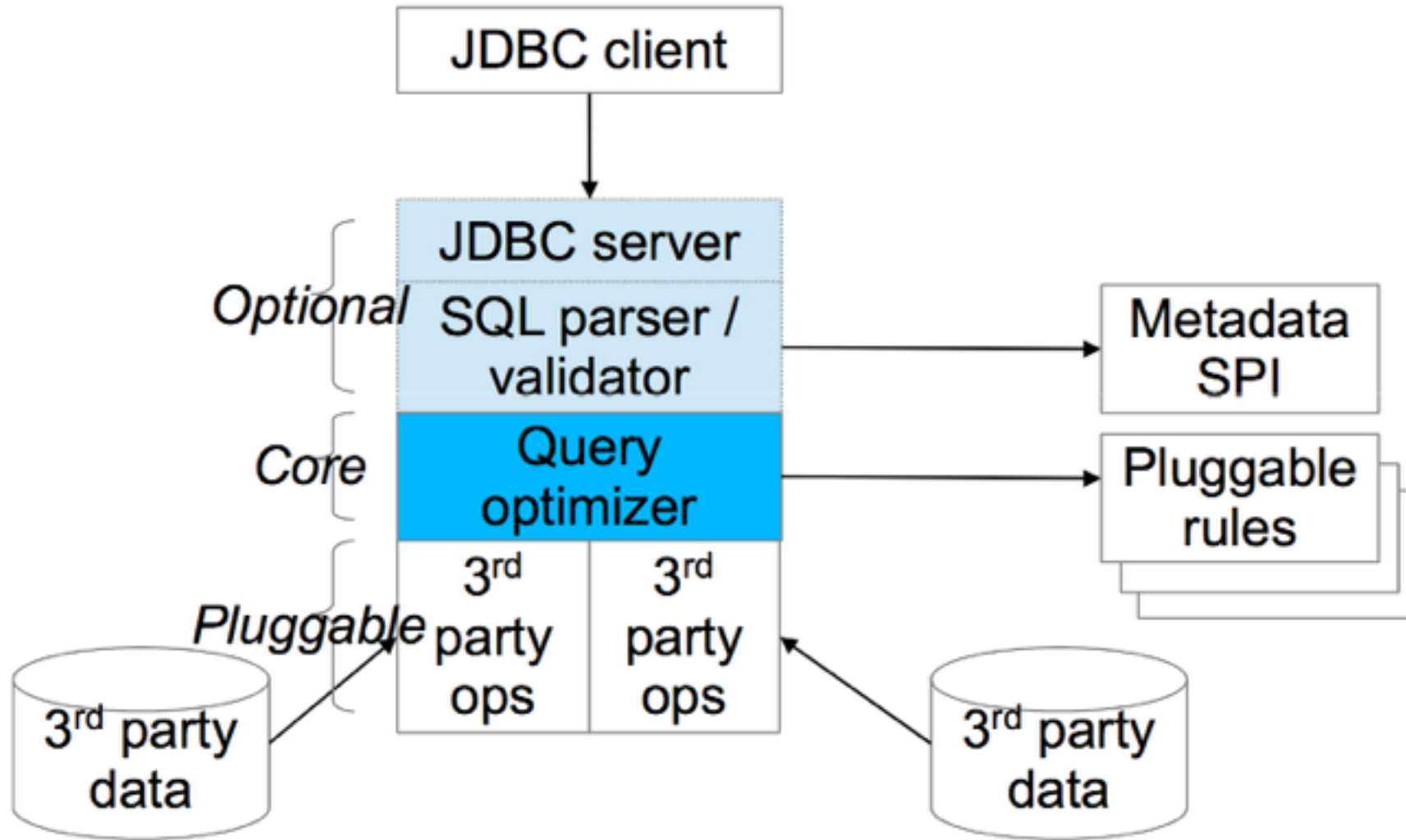
- Open source

- Community-authored rules, adapters

**Adoption**

- **Embedded**: Lingual (SQL interface to Cascading), Apache Drill, Apache Hive, Kylin OLAP

- **Adapters**: Splunk, Spark, MongoDB, JDBC, CSV, JSON, Web tables, In-memory data

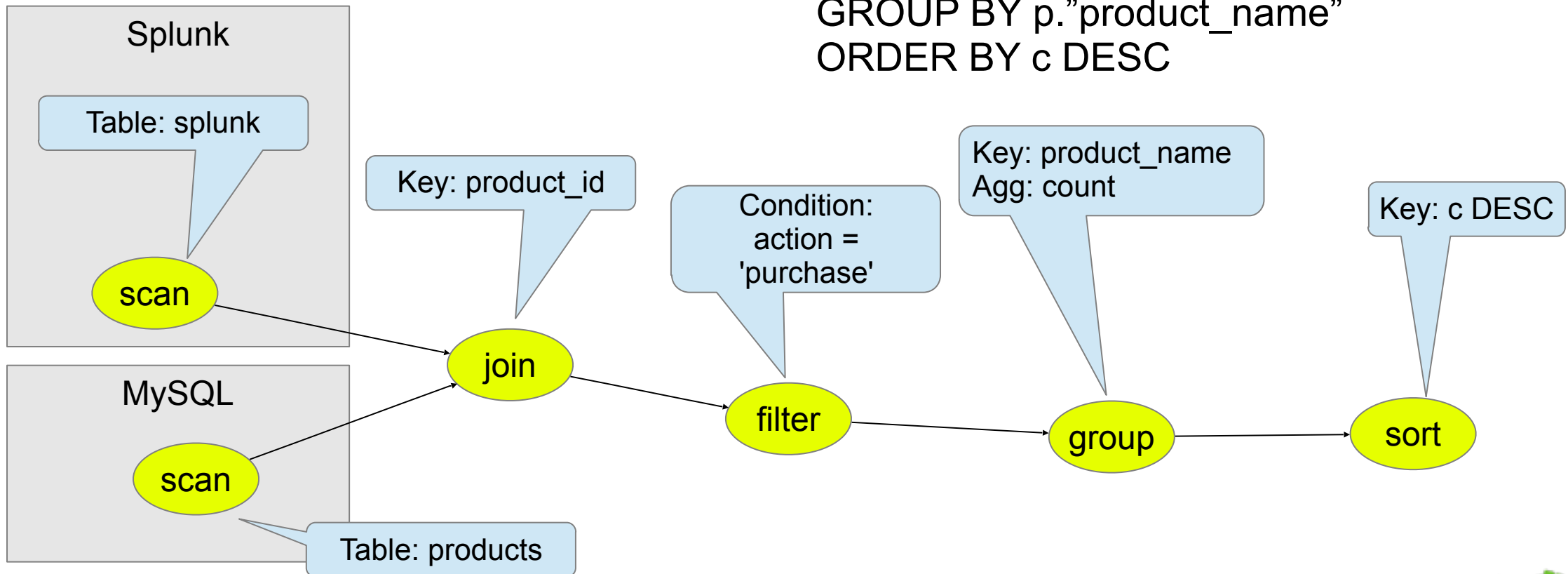# Conventional DB architecture

# Calcite architecture
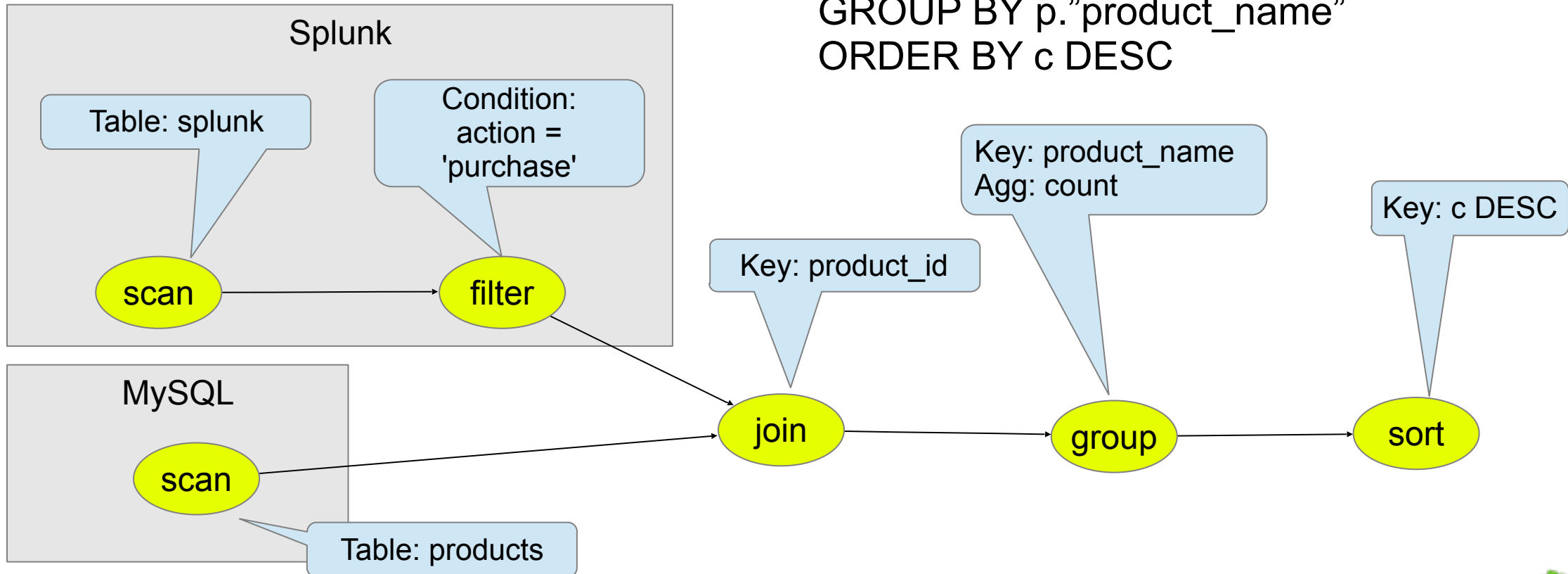
# Demo

{sqlline, apache-calcite-0.9.1, .csv}

# Expression tree

SELECT p."product_name", COUNT(*) AS c
FROM "splunk"."splunk" AS s
    JOIN "mysql"."products" AS p
    ON s."product_id" = p."product_id"
WHERE s."action" = 'purchase'
GROUP BY p."product_name"
ORDER BY c DESC

# Expression tree (optimized)

SELECT p."product_name", COUNT(*) AS c
FROM "splunk"."splunk" AS s
    JOIN "mysql"."products" AS p
    ON s."product_id" = p."product_id"
WHERE s."action" = 'purchase'
GROUP BY p."product_name"
ORDER BY c DESC

# Calcite – APIs and SPIs

**Relational algebra**

RelNode (operator)
- TableScan
- Filter
- Project
- Union
- Aggregate
- …

RelDataType (type)
RexNode (expression)
RelTrait (physical property)
- RelConvention (calling-convention)
- RelCollation (sortedness)
- TBD (bucketedness/distribution)

**SQL parser**

SqlNode
SqlParser
SqlValidator

**Metadata**

Schema
Table
Function
- TableFunction
- TableMacro
Lattice

**JDBC driver**

**Transformation rules**

RelOptRule
- MergeFilterRule
- PushAggregateThroughUnionRule
- 100+ more
Global transformations
- Unification (materialized view)
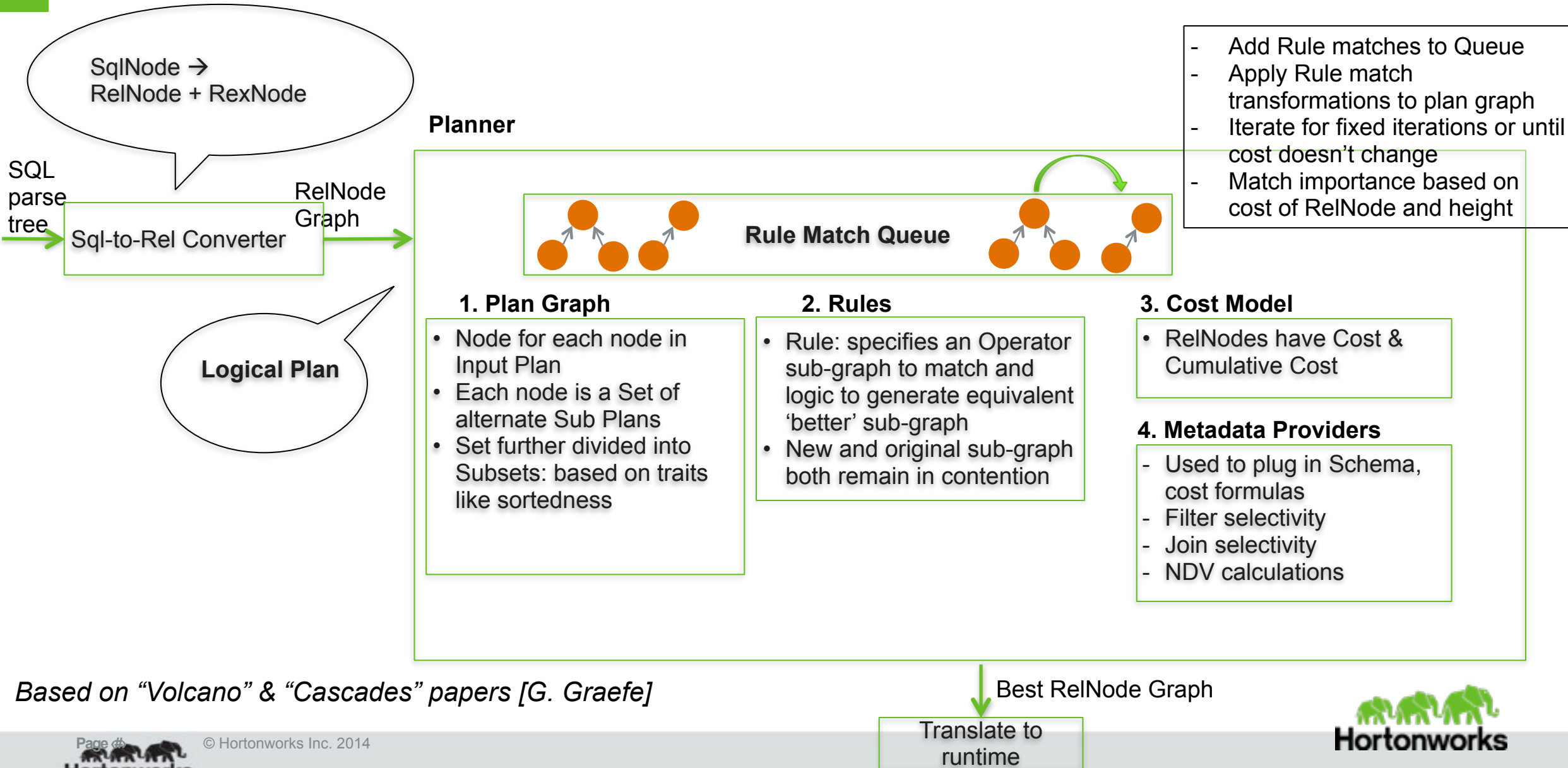- Column trimming
- De-correlation

**Cost, statistics**

RelOptCost
RelOptCostFactory
RelMetadataProvider
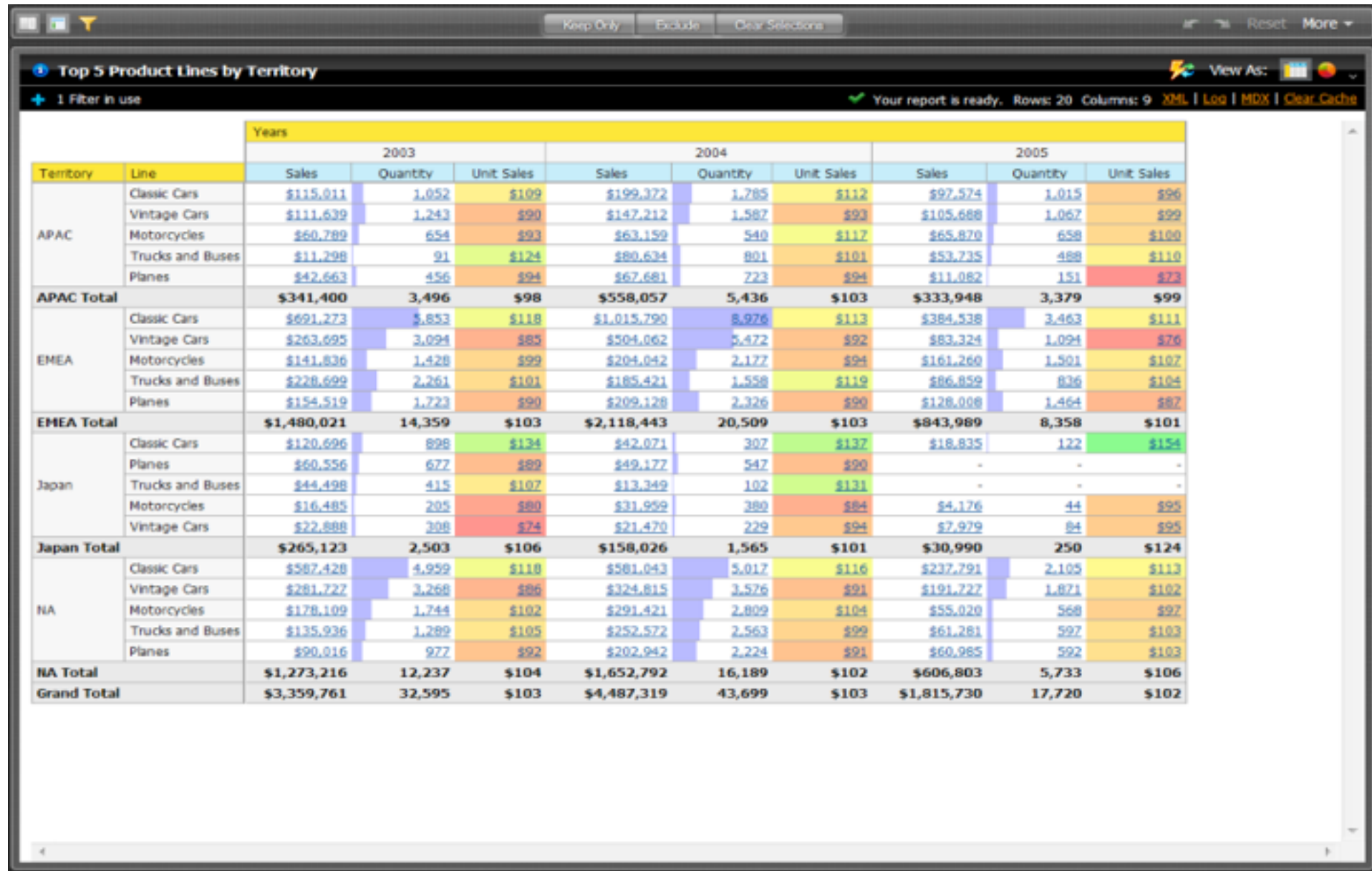- RelMdColumnUniquensss
- RelMdDistinctRowCount
- RelMdSelectivity

HORTONWORKS
Hortonworks

# Calcite Planning Process

SqlNode →
RelNode + RexNode

SQL
parse
tree

**Planner**

Sql-to-Rel Converter

RelNode
Graph

Logical Plan

- Add Rule matches to Queue
- Apply Rule match transformations to plan graph
- Iterate for fixed iterations or until cost doesn't change
- Match importance based on cost of RelNode and height



**Rule Match Queue**

### 1. Plan Graph

- Node for each node in Input Plan
- Each node is a Set of alternate Sub Plans
- Set further divided into Subsets: based on traits like sortedness

### 2. Rules

- Rule: specifies an Operator sub-graph to match and logic to generate equivalent 'better' sub-graph
- New and original sub-graph both remain in contention

### 3. Cost Model

- RelNodes have Cost & Cumulative Cost

### 4. Metadata Providers

- Used to plug in Schema, cost formulas
- Filter selectivity
- Join selectivity
- NDV calculations

*Based on "Volcano" & "Cascades" papers [G. Graefe]*

Best RelNode Graph

Translate to runtime

**Hortonworks**

# Demo

{sqlline, apache-calcite-0.9.1, .csv, CsvPushProjectOntoTableRule}

# Analytics

# Mondrian OLAP (Saiku user interface)

# Interactive queries on NoSQL

## Typical requirements

NoSQL operational database (e.g. HBase, MongoDB, Cassandra)

Analytic queries aggregate over full scan

Speed-of-thought response (< 5 seconds first query, < 1 second second query)

Data freshness (< 10 minutes)

## Other requirements

Hybrid system (e.g. Hive + HBase)

Star schema

# Star schema



**Key**
Fact table
Dimension table
→ Many-to-one relationship

# Simple analytics problem?

**System**

100M US census records

1KB each record, 100GB total

4 SATA 3 disks, total read throughput 1.2GB/s

**Requirement**

Count all records in < 5s

**Solution #1**

It's not possible! It takes 80s just to read the data

**Solution #2**

Cheat!

**Hortonworks**

# How to cheat

## Multiple tricks

Compress data

Column-oriented storage

Store data in sorted order

Put data in memory

Cache previous results

Pre-compute (materialize) aggregates

## Common factors

Make a copy of the data

Organize it in a different way

Optimizer chooses the most suitable data organization

SQL query is unchanged

# Filter-join-aggregate query

SELECT product.id,  sum(sales.units), sum(sales.price), count(*)
FROM sales …
JOIN customer ON …
JOIN time ON …
JOIN product ON …
JOIN product_class ON …
WHERE time.year = 2014
AND time.quarter = 'Q1'
AND product.color = 'Red'
GROUP BY …

# Materialized view, lattice, tile

**Materialized view**

A table whose contents are guaranteed to be the same as
executing a given query.

**Lattice**

Recommends, builds, and recognizes summary
materialized views (tiles) based on a star schema.

A query defines the tables and many:1 relationships in the
star schema.

## Tile

A summary materialized view that belongs to a lattice.

A tile may or may not be materialized.

Materialization methods:

• Declare in lattice

• Generate via recommender algorithm

• Created in response to query

(FAKE SYNTAX)

```
CREATE MATERIALIZED VIEW t AS
SELECT * FROM emps
WHERE deptno = 10;


CREATE LATTICE star AS
SELECT *
FROM sales_fact_1997 AS s
JOIN product AS p ON …
JOIN product_class AS pc ON …
JOIN customer AS c ON …
JOIN time_by_day AS t ON …;


CREATE MATERIALIZED VIEW zg IN star
SELECT gender, zipcode,
  COUNT(*), SUM(unit_sales)
FROM star
GROUP BY gender, zipcode;
```

HORTONWORKS

# Lattice

```
()1
```

```
> select count(*) as c, sum(unit_sales) as s
> from star;
+------------+------------+
|     C      |         S  |
+------------+------------+
| 1,000,000  | 266,773.0  |
+------------+------------+
1 row selected

> select * from star;
1,000,000 rows selected
```

```
raw 1m
```

Hortonworks

# Lattice - top tiles



Key

z zipcode (43k)
s state (50)
g gender (2)
y year (5)
m month (12)

() 1

(z) 43k    (s) 50    (g) 2    (y) 5    (m) 12

```
> select zipcode, count(*) as c,
>   sum(unit_sales) as s
> from star
> group by zipcode;
+---------+---------+---------+
| ZIPCODE |       C |       S |
+---------+---------+---------+
| 10000   |      23 |    31.5 |
  …
+---------+---------+---------+
43,000 rows selected
```

```
> select state, count(*) as c,
>   sum(unit_sales) as s
> from star
> group by state;
+-------+---------+---------+
| STATE |       C |       S |
+-------+---------+---------+
| AL    | 201,693 | 5,520.0 |
  …
+-------+---------+---------+
50 rows selected
```

raw 1m

Hortonworks

# Lattice - more tiles



Key

z zipcode (43k)
s state (50)
g gender (2)
y year (5)
m month (12)

# Lattice - complete



() 1

(z) 43k    (s) 50    (g) 2    (y) 5    (m) 12

(z, s) 43.4k    (g, y) 10    (g, m) 24    (y, m) 60

(z, s, g) 87k    (g, y, m) 120

(z, s, g, y) 391k    (z, s, g, m) 643k    (z, s, y, m) 830k    (z, g, y, m) 909k    (s, g, y, m) 6k

(z, s, g, y, m) 910k

raw 1m

Key

z zipcode (43k)
s state (50)
g gender (2)
y year (5)
m month (12)

Hortonworks

Hortonworks

# Lattice - optimized



() 1

(z) 43k    (s) 50    (g) 2    (y) 5    (m) 12

(z, s) 43.4k    (g, y) 10    (g, m) 24    (y, m) 60

(z, s, g) 87k    (g, y, m) 120

(z, s, g, y) 392k    (z, s, g, m) 644k    (z, s, y, m) 831k    (z, g, y, m) 909k    (s, g, y, m) 6k

(z, s, g, y, m) 912k

raw 1m

Key

z zipcode (43k)
s state (50)
g gender (2)
y year (5)
m month (12)

Hortonworks

# Lattice - optimized



Diagram nodes:
- () 1
- (z) 43k
- (s) 50
- (g) 2
- (y) 5
- (m) 12
- (z, s) 43.4k
- (g, y) 10
- (g, m) 24
- (y, m) 60
- (z, s, g) 87k
- (g, y, m) 120
- (z, s, g, y) 392k
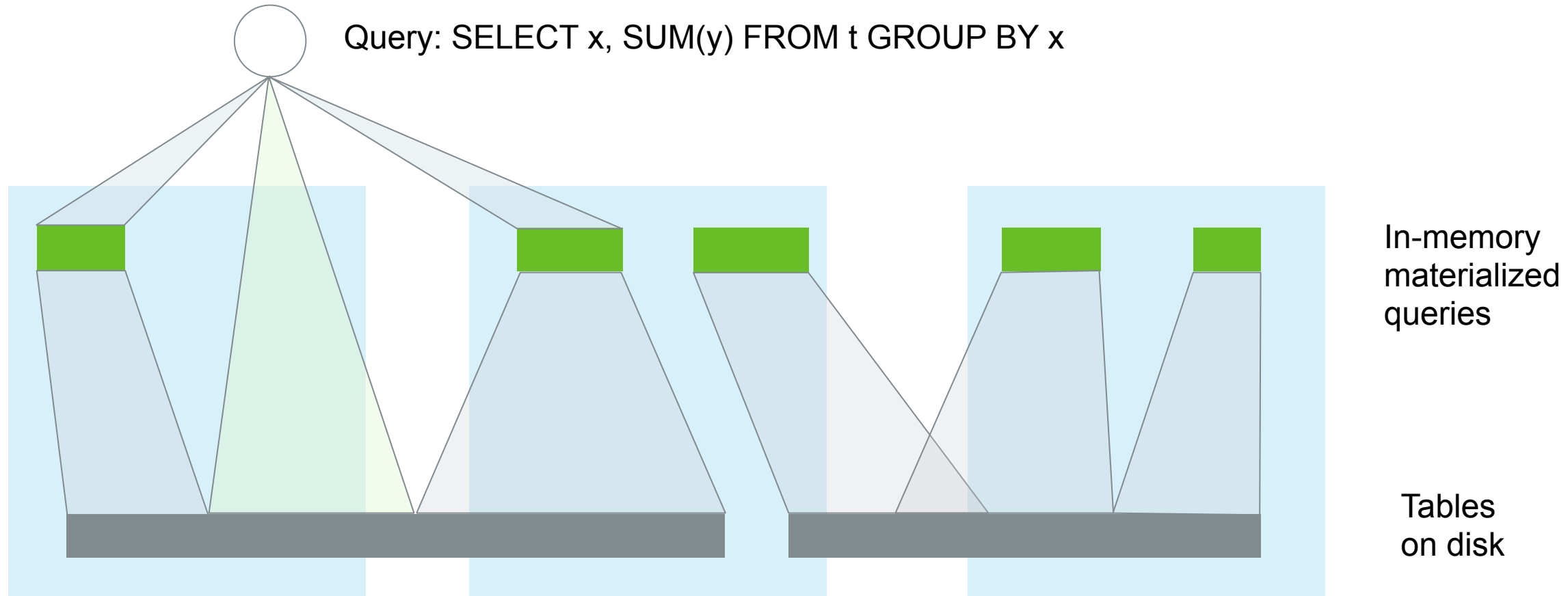- (s, g, y, m) 6k

Key

z zipcode (43k)
s state (50)
g gender (2)
y year (5)
m month (12)

| Aggregate | Cost (rows) | Benefit (query rows saved) | % queries |
|-----------|-------------|----------------------------|-----------|
| s, g, y, m | 6k | 497k | 50% |
| z, s, g | 87k | 304k | 33% |
| g, y | 10 | 1.5k | 25% |
| g, m | 24 | 1.5k | 25% |
| s, g | 100 | 1.5k | 25% |
| y, m | 60 | 1.5k | 25% |

Hortonworks

# Demo

{mysql-foodmart-lattice-model.json}

# Tiled, in-memory materializations



Query: SELECT x, SUM(y) FROM t GROUP BY x

In-memory materialized queries

Tables on disk

**Where we're going… smart, distributed memory cache & compute framework**
**http://hortonworks.com/blog/dmmq/**

Hortonworks

# Kylin OLAP engine

# Thank you!

**Apache Calcite**

**@julianhyde**

**http://calcite.incubator.apache.org**

**http://www.kylin.io**

**Hortonworks**