# Apache Drill
# interactive, ad-hoc query at scale

Michael Hausenblas, Chief Data Engineer EMEA, MapR

Hadoop ecosystem - Open Source drives innovation and adoption in Big Data, 2013-01-05

Which workloads do **you** encounter in **your** environment?

APACHE DRILL

# Batch processing



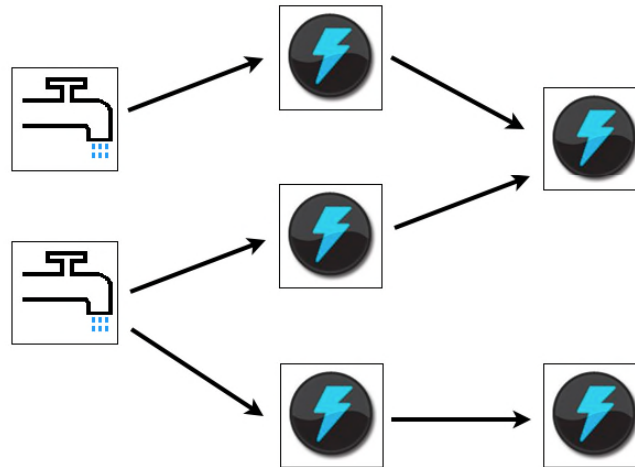… for recurring tasks such as large-scale data mining, aggregation, ETL offloading, etc.

# OLTP



… for example user-facing eCommerce transactions, real-time messaging at scale (FB) , etc.
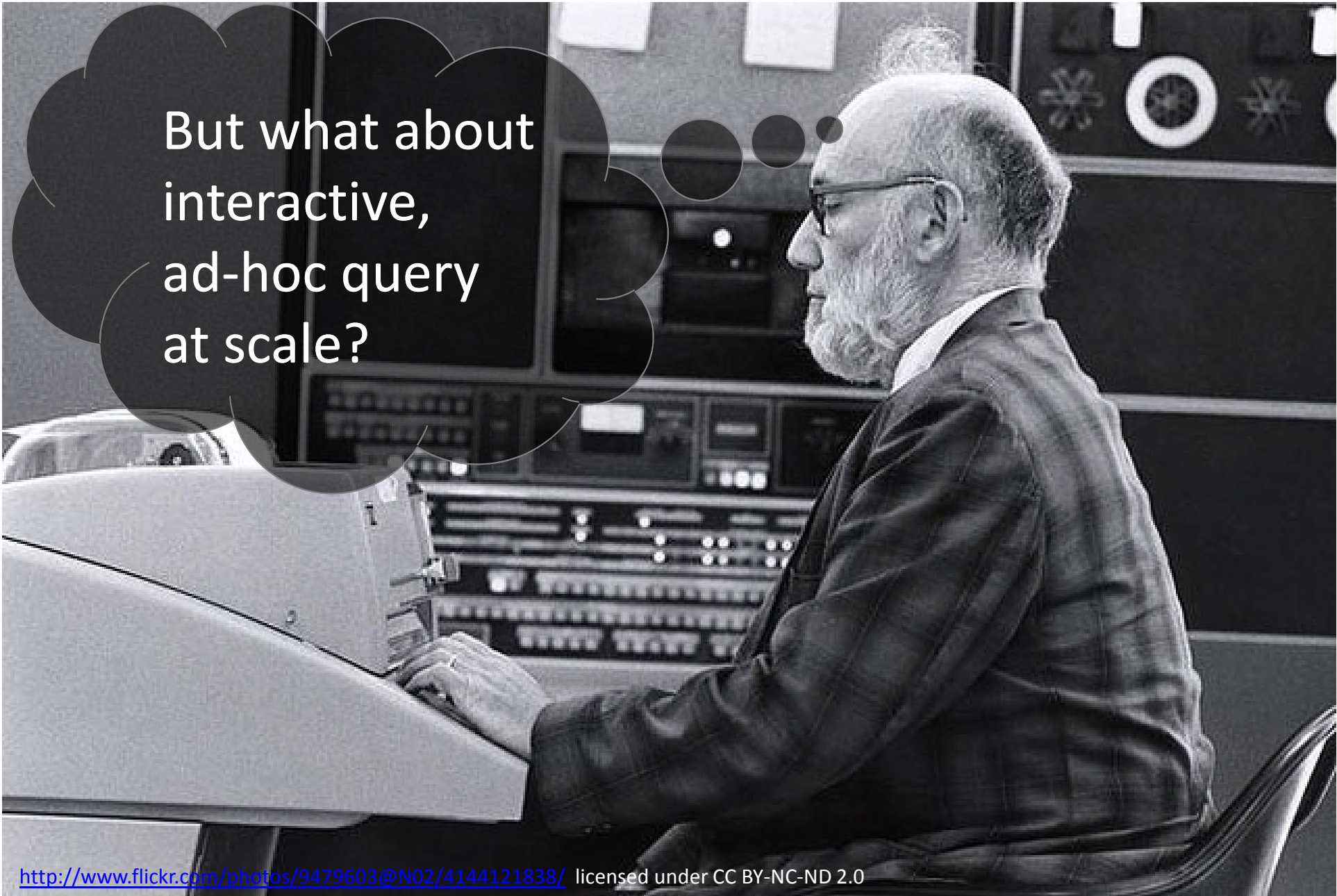
# Stream processing



... in order to handle stream sources such as social media feeds or sensor data (mobile phones, RFID, weather stations, etc.)

# Search



… retrieval of items from semi-structured data formats
(XML, JSON, etc.), documents (plain text, etc.) and
datastores (MongoDB, CouchDB, etc.)

But what about interactive, ad-hoc query at scale?

# Interactive Query (?)



low-latency

# Use Case I

- Jane, a marketing analyst

- Determine target segments

- Data from different sources


Transaction information — ORACLE


User profiles — mongoDB


Access logs — hadoop

APACHE DRILL

MAPR TECHNOLOGIES

# Use Case II

- Logistics – supplier status
- Queries
  - How many shipments from supplier X?
  - How many shipments in region Y?

| SUPPLIER_ID | NAME | REGION |
|---|---|---|
| ACM | ACME Corp | US |
| GAL | GotALot Inc | US |
| BAP | Bits and Pieces Ltd | Europe |
| ZUP | Zu Pli | Asia |

```
{
  "shipment": 100123,
  "supplier": "ACM",
  "timestamp": "2013-02-01",
  "description": "first delivery today"
},
{
  "shipment": 100124,
  "supplier": "BAP",
  "timestamp": "2013-02-02",
  "description": "hope you enjoy it"
}
...
```

APACHE DRILL

# Requirements

- Support for different data sources
- Support for different query interfaces
- Low-latency/real-time
- Ad-hoc queries
- Scalable, reliable

And now for something completely different …

# Google's Dremel

"

Dremel is a scalable, interactive ad-hoc query system for analysis of read-only nested data. By combining multi-level execution trees and columnar data layout, it is capable of running aggregation queries over trillion-row tables in seconds. The system scales to thousands of CPUs and petabytes of data, and has thousands of users at Google.
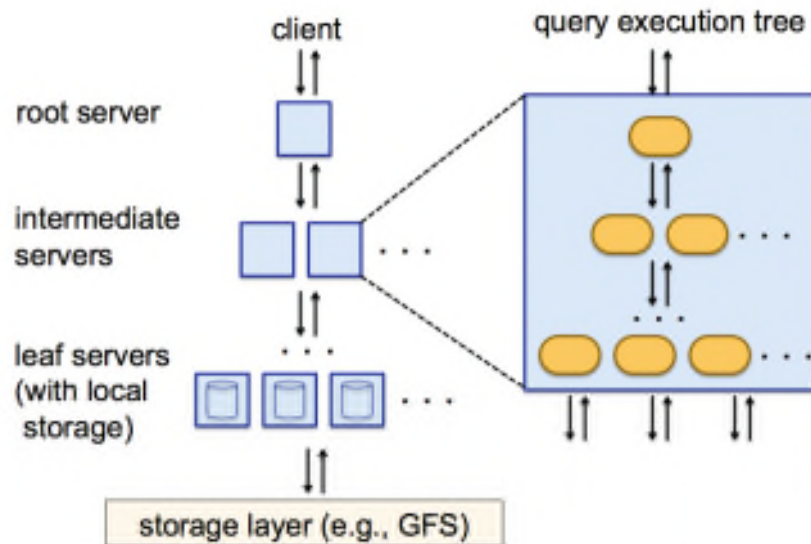
…

"

APACHE
DRILL

MAPR
TECHNOLOGIES

# Google's Dremel



multi-level execution trees

columnar data layout

# Google's Dremel



nested data + schema                    column-striped representation

**mapping nested data to tables**

# Google's Dremel

| Table name | Number of records | Size (unrepl., compressed) | Number of fields | Data center | Repl. factor |
|---|---|---|---|---|---|
| T1 | 85 billion | 87 TB | 270 | A | 3× |
| T2 | 24 billion | 13 TB | 530 | A | 3× |
| T3 | 4 billion | 70 TB | 1200 | A | 3× |
| T4 | 1+ trillion | 105 TB | 50 | B | 3× |
| T5 | 1+ trillion | 20 TB | 30 | B | 2× |

**experiments**:
datasets & query performance

Back to Apache Drill …

# Apache Drill–key facts

- Inspired by Google's **Dremel**

- Standard **SQL 2003** support

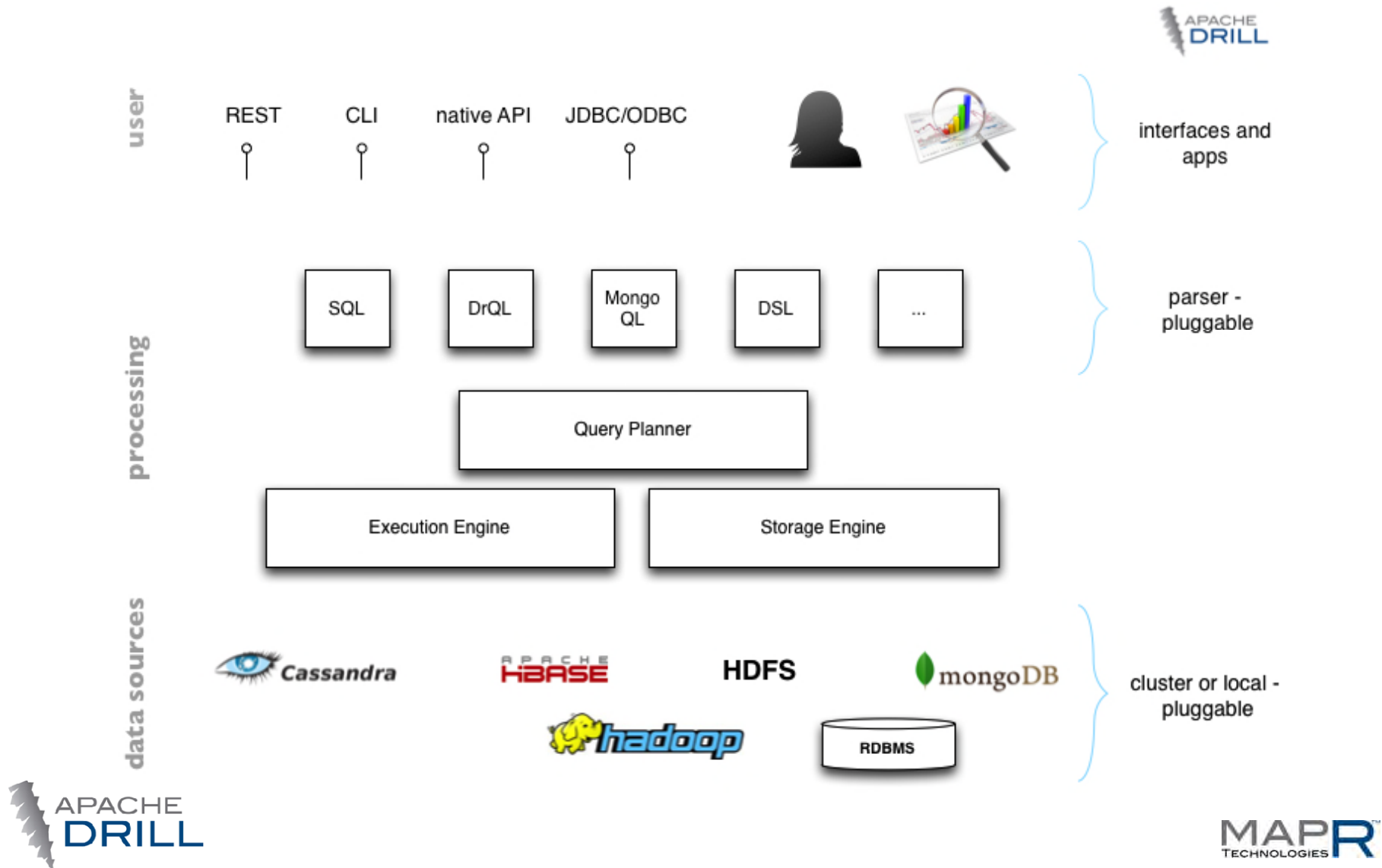- Plug-able **data sources**

- **Nested data** is a first-class citizen

- **Schema** is **optional**

- **Community** driven, **open**, 100's involved

APACHE
DRILL

MAPR
TECHNOLOGIES

# High-level Architecture

# Wire-level Architecture

- Each node: **Drillbit** - maximize data locality
- Co-ordination, query planning, execution, etc, are **distributed**
- By default Drillbits hold all roles
- Any node can act as endpoint for a query



APACHE DRILL

MAPR TECHNOLOGIES

# Wire-level Architecture

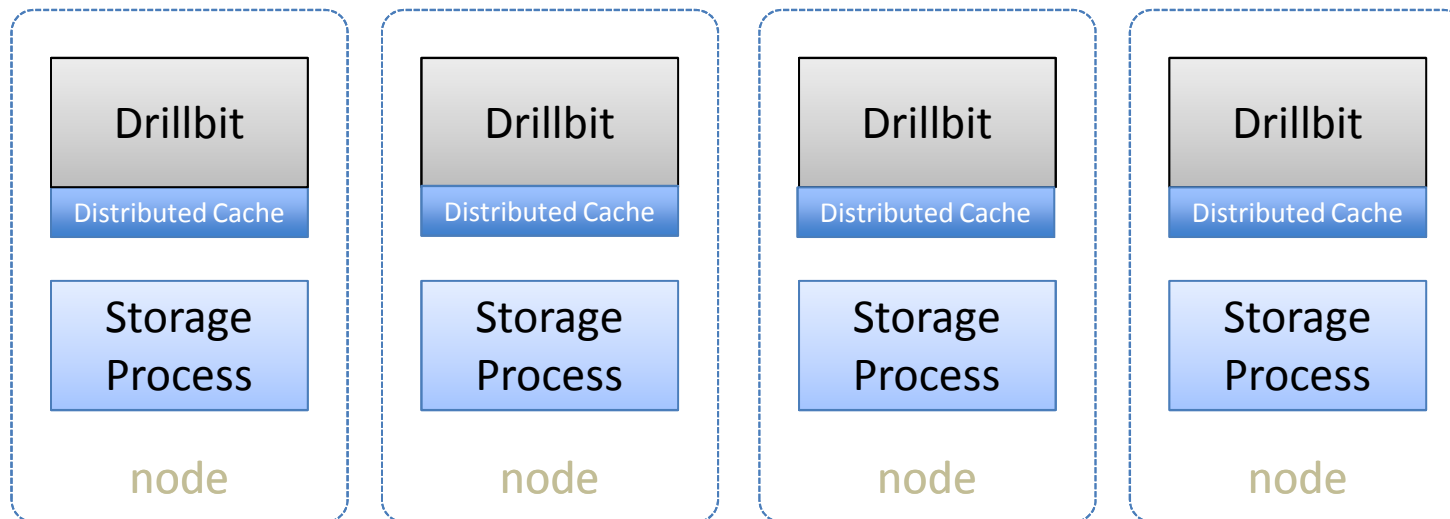- **Zookeeper** for ephemeral cluster membership info
- **Distributed cache** (Hazelcast) for metadata, locality information, etc.

Curator/Zk

| node | node | node | node |
|---|---|---|---|
| Drillbit | Drillbit | Drillbit | Drillbit |
| Distributed Cache | Distributed Cache | Distributed Cache | Distributed Cache |
| Storage Process | Storage Process | Storage Process | Storage Process |

APACHE DRILL

MAPR TECHNOLOGIES
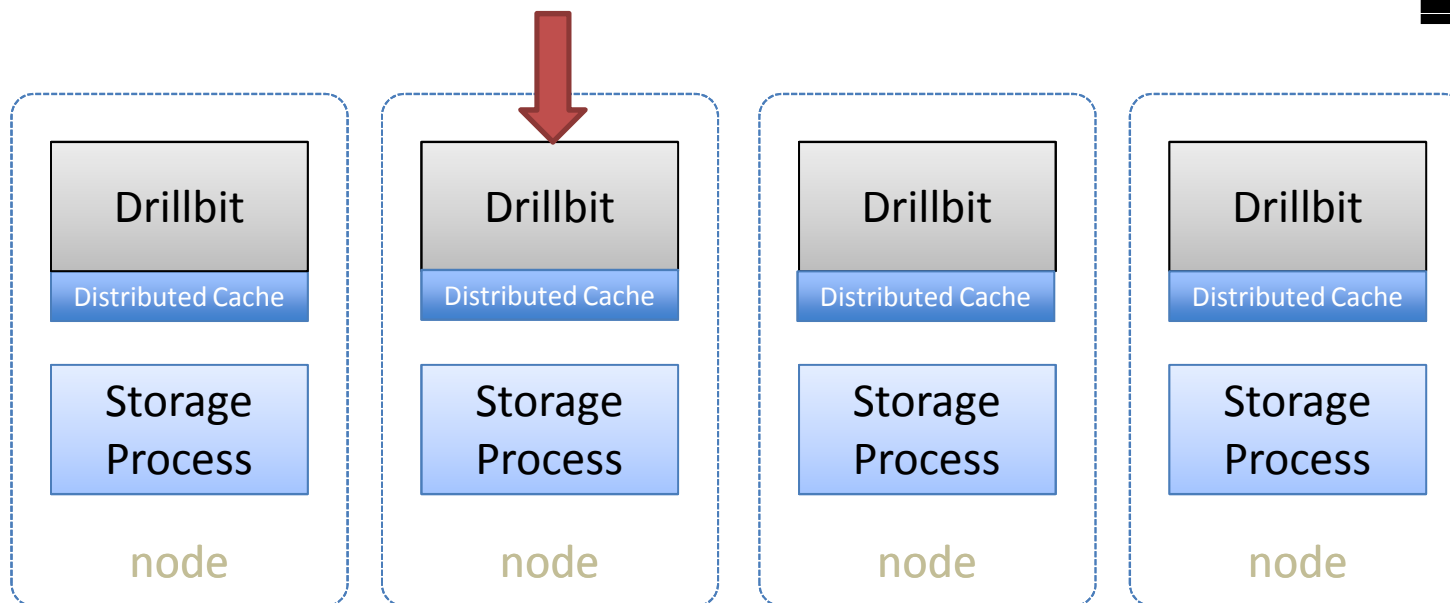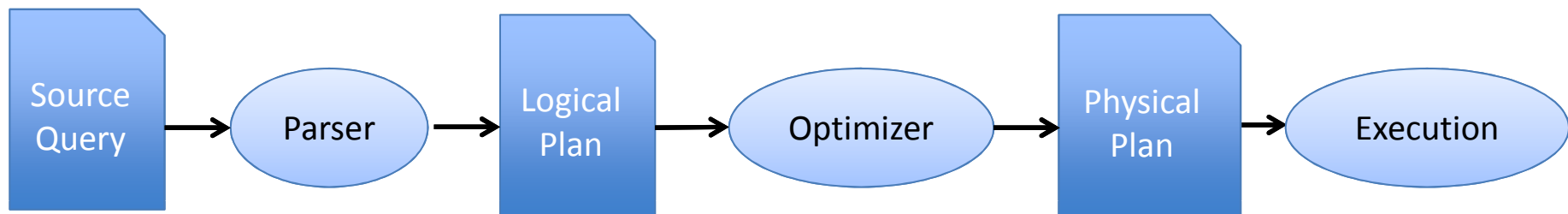
# Wire-level Architecture

- **Originating Drillbit** acts as foreman, manages query execution, scheduling, locality information, etc.

- Streaming data **communication** avoiding SerDe

# Principled Query Execution

Source Query → Parser → Logical Plan → Optimizer → Physical Plan → Execution
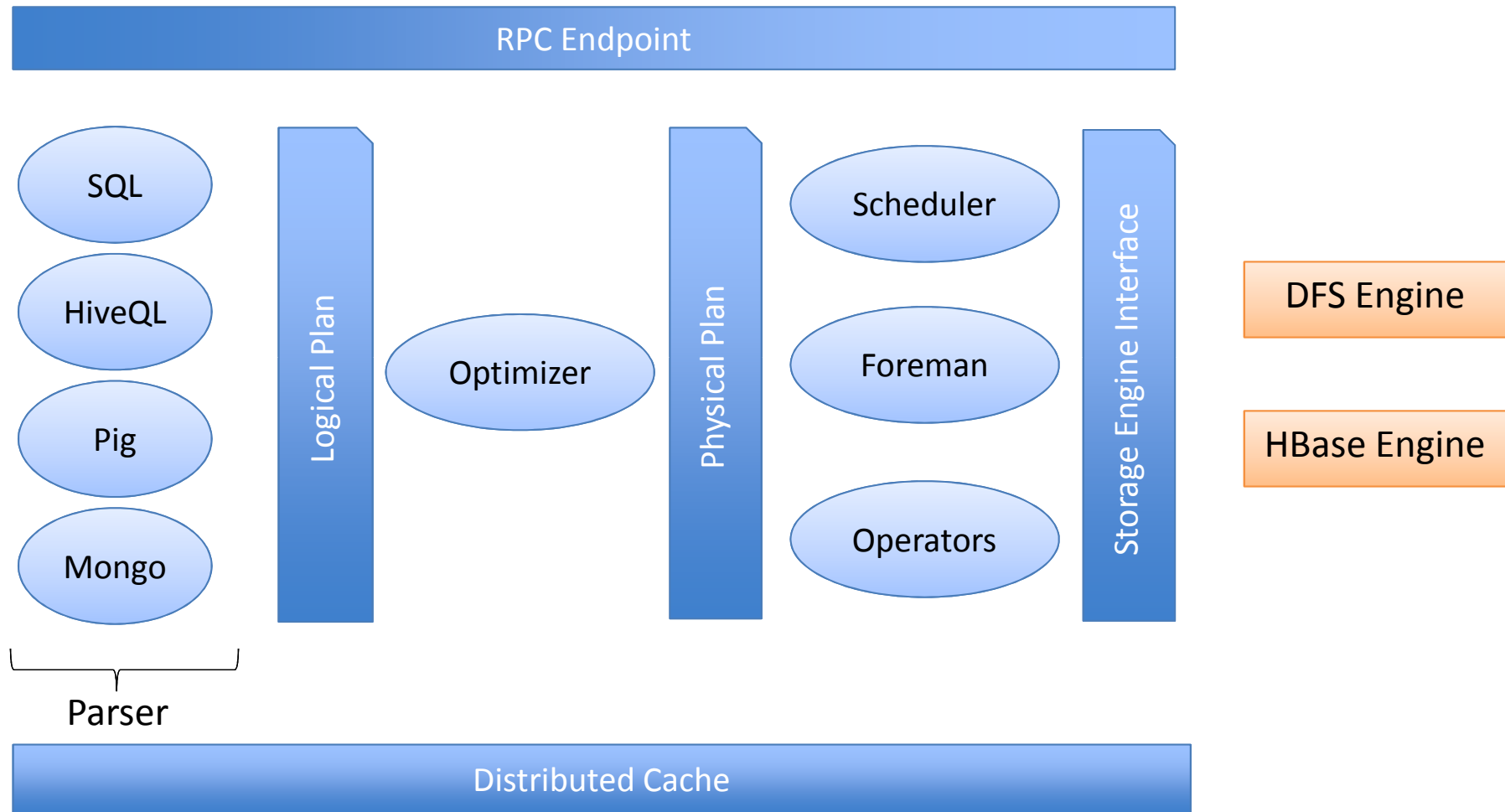
SQL 2003
DrQL
MongoQL
DSL

parser API

```
query: [
{
 @id: "log",
 op: "sequence",
 do: [
  {
   op: "scan",
   source: "logs"
  },
  {
   op:
    "filter",
   condition:
    "x > 3"
  },
```

topology

scanner API
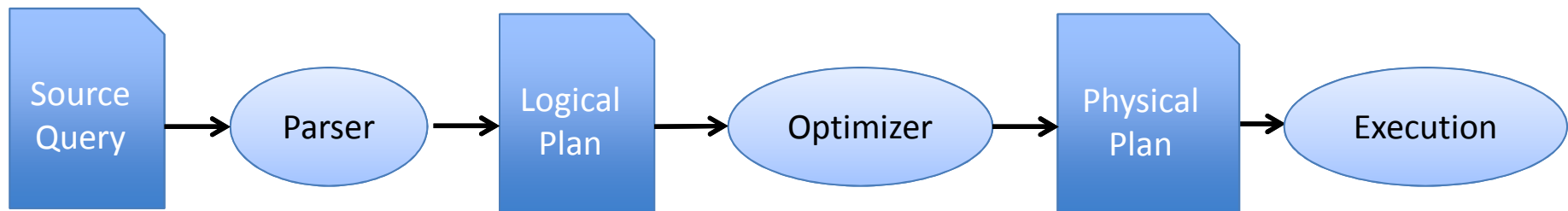
# Drillbit Modules

# Key features

- Full SQL – ANSI SQL 2003

- Nested Data as first class citizen

- Optional Schema

- Extensibility Points …

# Extensibility Points

- Source query → parser API
- Custom operators, UDF → logical plan
- Serving tree, CF, topology → physical plan/optimizer
- Data sources &formats → scanner API

# … and Hadoop?

- HDFS can be a data source

- Complementary use cases*

- … use Apache Drill
  - Find record with specified condition
  - Aggregation under dynamic conditions

- … use MapReduce
  - Data mining with multiple iterations
  - ETL

An Inside Look at Google BigQuery

Table of Contents

# Basic Demo

```
{
 "id": "0001",
 "type": "donut",
 "ppu": 0.55,
 "batters":
 {
  "batter":
  [
        { "id": "1001", "type": "Regular" },
        { "id": "1002", "type": "Chocolate" },
...
```

data source: **donuts.json**

```
query:[ {
    op:"sequence",
    do:[
        {
         op: "scan",
         ref: "donuts",
         source: "local-logs",
         selection: {data: "activity"}
        },
        {
         op: "filter",
         expr: "donuts.ppu < 2.00"
        },
...
```

logical plan: **simple_plan.json**

APACHE **DRILL**

APACHE **DRILL**

```
{
   "sales" : 700.0,
   "typeCount" : 1,
   "quantity" : 700,
   "ppu" : 1.0
}
 {
   "sales" : 109.71,
   "typeCount" : 2,
   "quantity" : 159,
   "ppu" : 0.69
}
 {
   "sales" : 184.25,
   "typeCount" : 2,
   "quantity" : 335,
   "ppu" : 0.55
}
```

result: **out.json**

https://cwiki.apache.org/confluence/display/DRILL/Demo+HowTo

MAP**R**
TECHNOLOGIES

# BE A PART OF IT!

# Status

- Heavy development by multiple organizations

- Available
  - Logical plan ([ADSP](#))
  - Reference interpreter
  - Basic SQL parser
  - Basic [demo](#)

APACHE
DRILL

MAPR
TECHNOLOGIES

# Status

May 2013

- Full SQL support (+JDBC)
- Physical plan
- In-memory compressed data interfaces
- Distributed execution
- HBase and MySQL storage engine
- WebUI client

# Contributing

Contributions appreciated (besides code drops)!

- Test data & test queries
- Use case scenarios (textual/SQL queries)
- Documentation
- Further schedule
  - Alpha Q2
  - Beta Q3

APACHE
DRILL

MAPR
TECHNOLOGIES

# Kudos to …

- Julian Hyde, Pentaho
- Lisen Mu, XingCloud
- Tim Chen, Microsoft
- Chris Merrick, RJMetrics
- David Alves, UT Austin
- Sree Vaadi, SSS/NGData
- Jacques Nadeau, MapR
- Ted Dunning, MapR

# Engage!

- Follow @ApacheDrill on Twitter

- Sign up at mailing lists (user | dev)
  http://incubator.apache.org/drill/mailing-lists.html

- Standing G+ hangouts every Tuesday at 5pm GMT
  http://j.mp/apache-drill-hangouts

- Keep an eye on http://drill-user.org/

Apache Drill User

KEEPING TRACK OF APACHE DRILL. FROM A GEEKS, FOR GEEKS.

APACHE DRILL

MAPR TECHNOLOGIES