



# Self-Service Data Exploration with Apache Drill

Tomer Shiran, VP Product Management



Empowering “as it happens”  
businesses by speeding up the  
data-to-action cycle



Top-Ranked **Hadoop**  
Distribution

---

Top-Ranked **NoSQL**

---

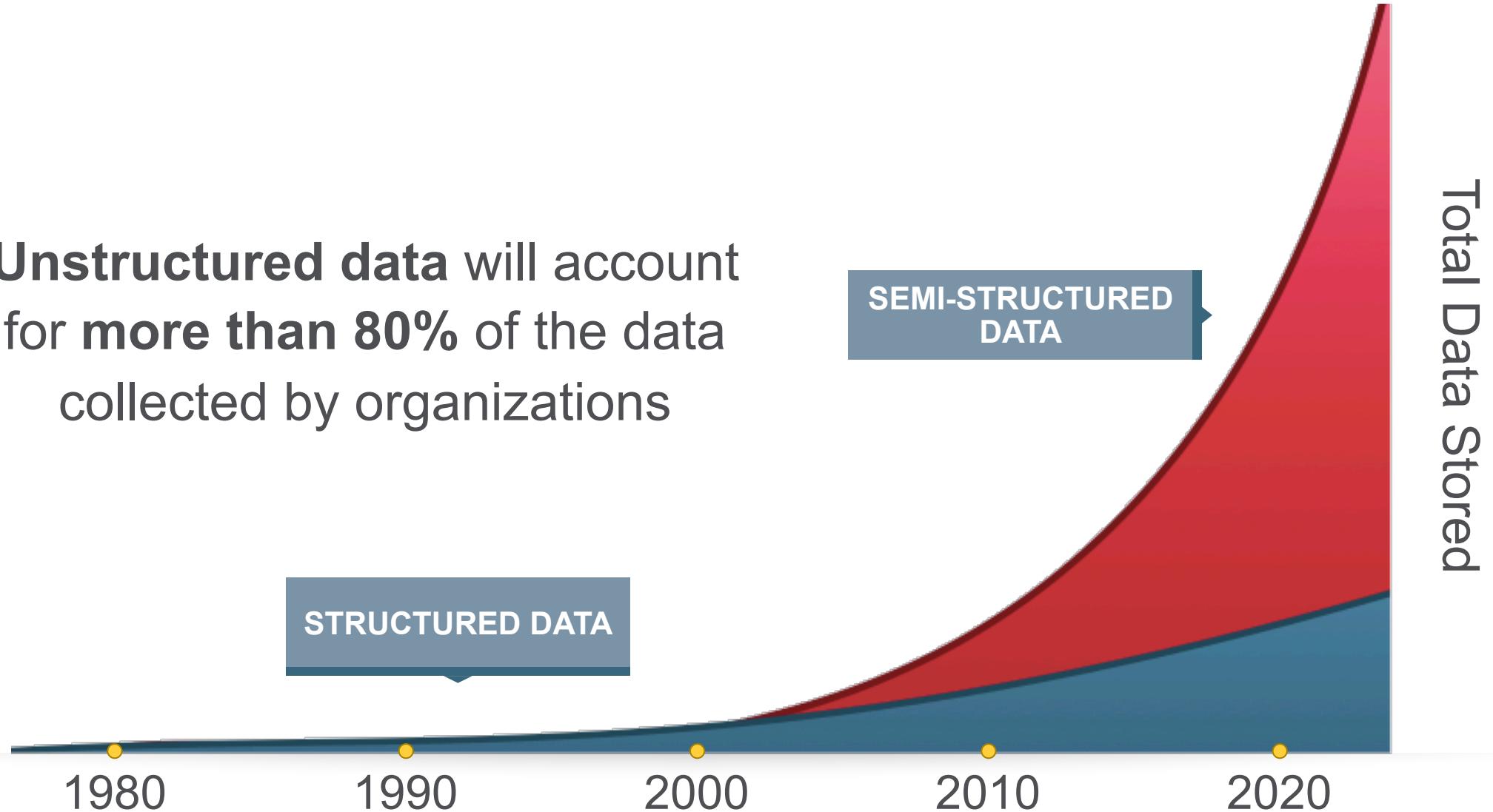
Top-Ranked **SQL-on-Hadoop**  
Solution

---

**Data is doubling in  
size every two years**



**Unstructured data will account for more than 80% of the data collected by organizations**



# Data Increasingly Stored in Non-Relational Datastores

Volume

MBs-TBs

Structure

Structured

Development

Planned (release cycle = months-years)



## RELATIONAL DATABASES

**ORACLE**  Microsoft SQL Server **MySQL** 

**Fixed schema**  
DBA controls structure

Database

1980

1990

2000

2010

2020



TBs-PBs

Structured, semi-structured and unstructured

Iterative (release cycle = days-weeks)



## NON-RELATIONAL DATASTORES

 **hadoop**  **HBASE**

**Dynamic schema (schema-free)**  
Application controls structure

# SQL in a Non-Relational World

## DON'T WANT

- Create and maintain schemas on:
  - HDFS (Parquet, JSON, etc.)
  - HBase
  - ...
- Transform or copy data

## WANT

- SQL
- BI (Tableau, MicroStrategy, etc.)
- Low latency
- Scalability



# APACHE DRILL

- Data exploration for Hadoop and NoSQL
- Low latency at scale
- Point-and-query vs. schema-first
- Extreme ease of use
- Industry-standard APIs: ANSI SQL, ODBC/JDBC, RESTful APIs

## Agility

- Explore big data in its native format without IT intervention

## Flexibility

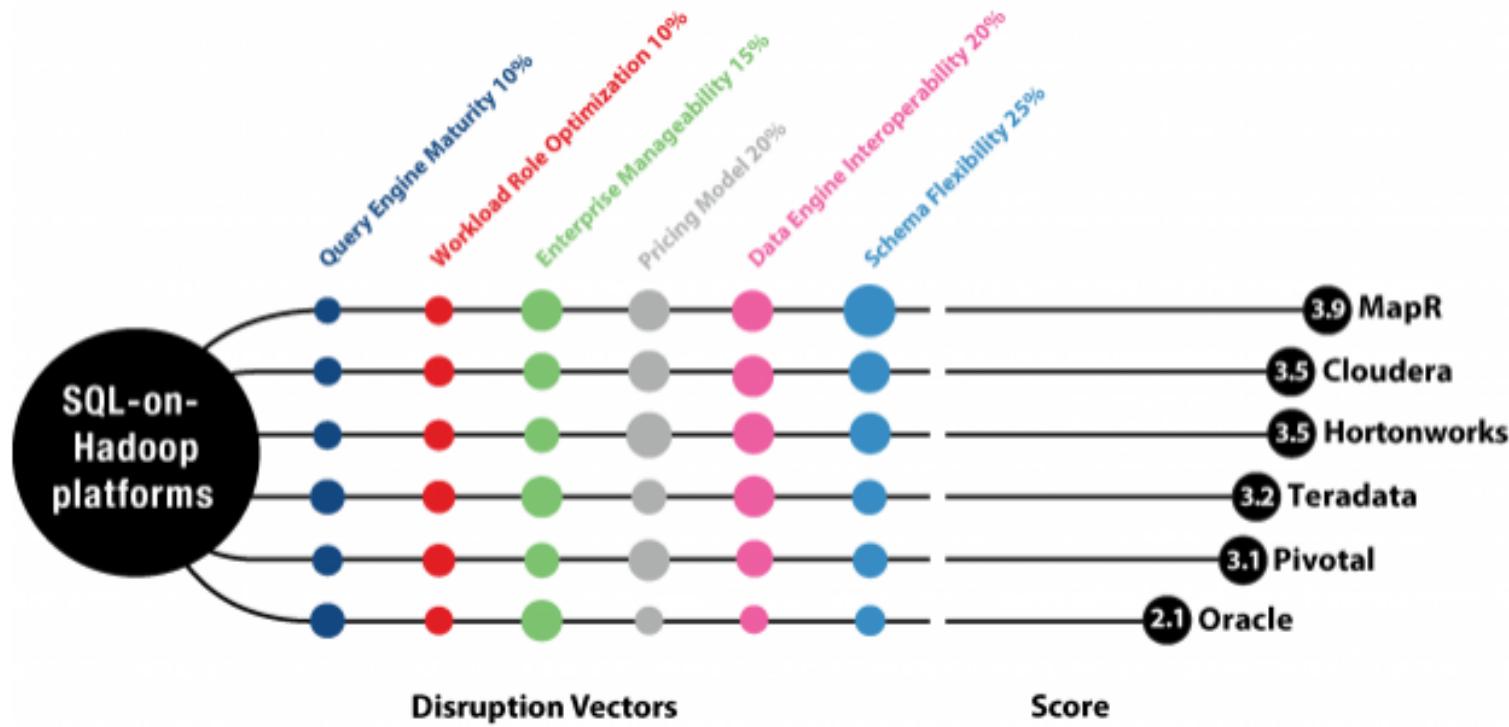
- Analyze semi-structured/nested data coming from NoSQL applications and JSON sources

## Plug-and-Play

- Leverage existing SQL skillsets, BI tools and Apache Hive deployments



# Drill is the Top-Ranked SQL-on-Hadoop



## Key:

- Number indicates companies relative strength across all vectors
- Size of ball indicates company's relative strength along individual vector

**GIGAOM**

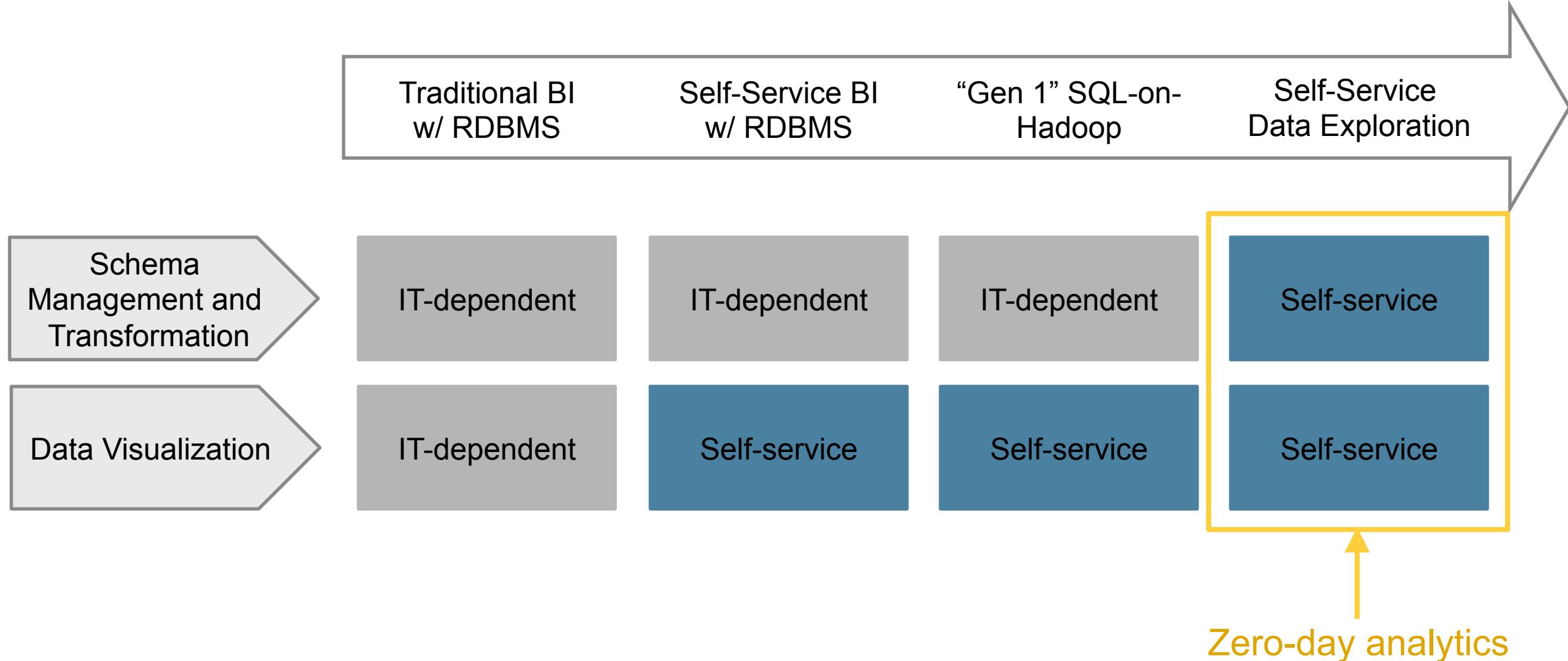


Source: Gigaom Research

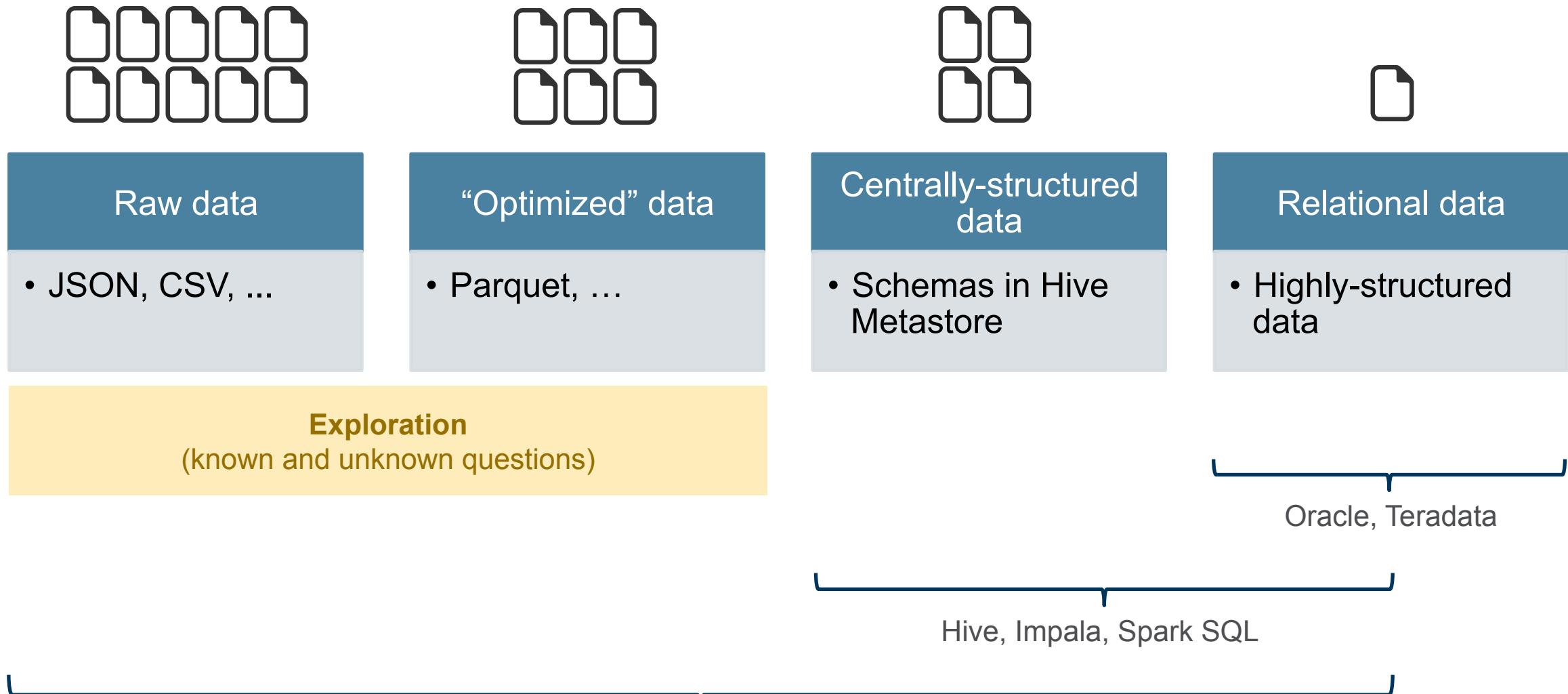
© 2014 MapR Technologies

**MAPR**

# Evolution Towards Self-Service Data Exploration



# Drill's Role in the Enterprise Data Architecture



# Leverage Existing SQL Tools and Skills



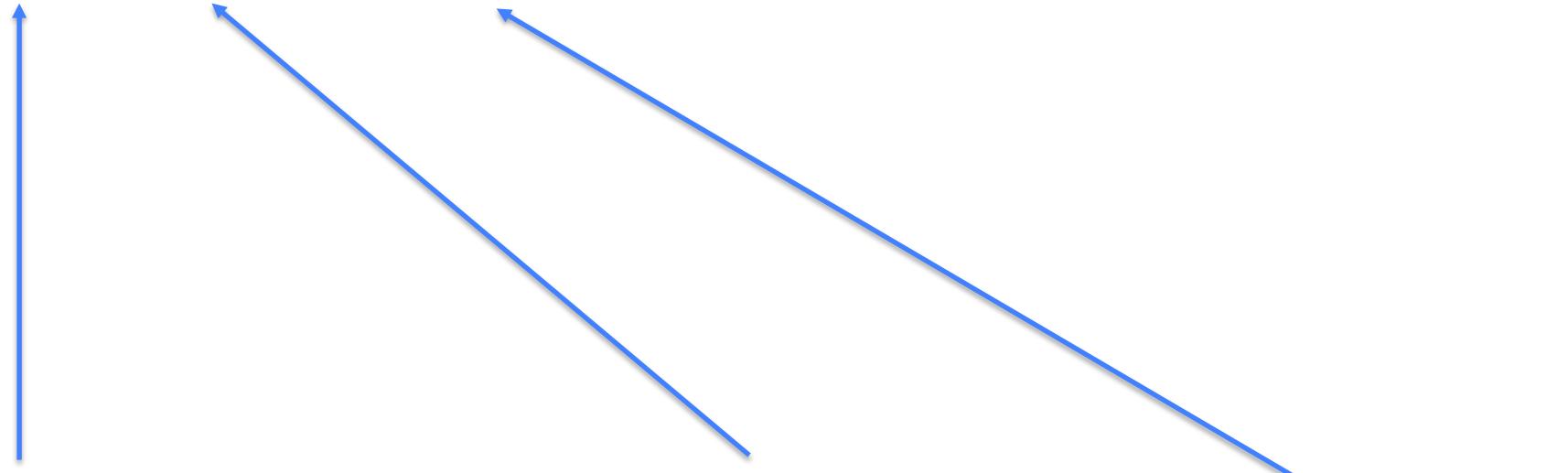
Leverage SQL-compatible tools (BI, query builders, etc.) via Drill's standard ODBC, JDBC and ANSI SQL support

Enable business analysts, technical analysts and data scientists to explore and analyze large volumes of real-time data



# Combine Data from Multiple Sources on the Fly

```
SELECT * FROM dfs.demo.`yelp/business.json`
```



A storage plugin instance

- DFS
- Hbase/MapRDB
- Hive Metastore/HCatalog

A workspace

- Sub-directory
- Hive database
- HBase namespace

A table

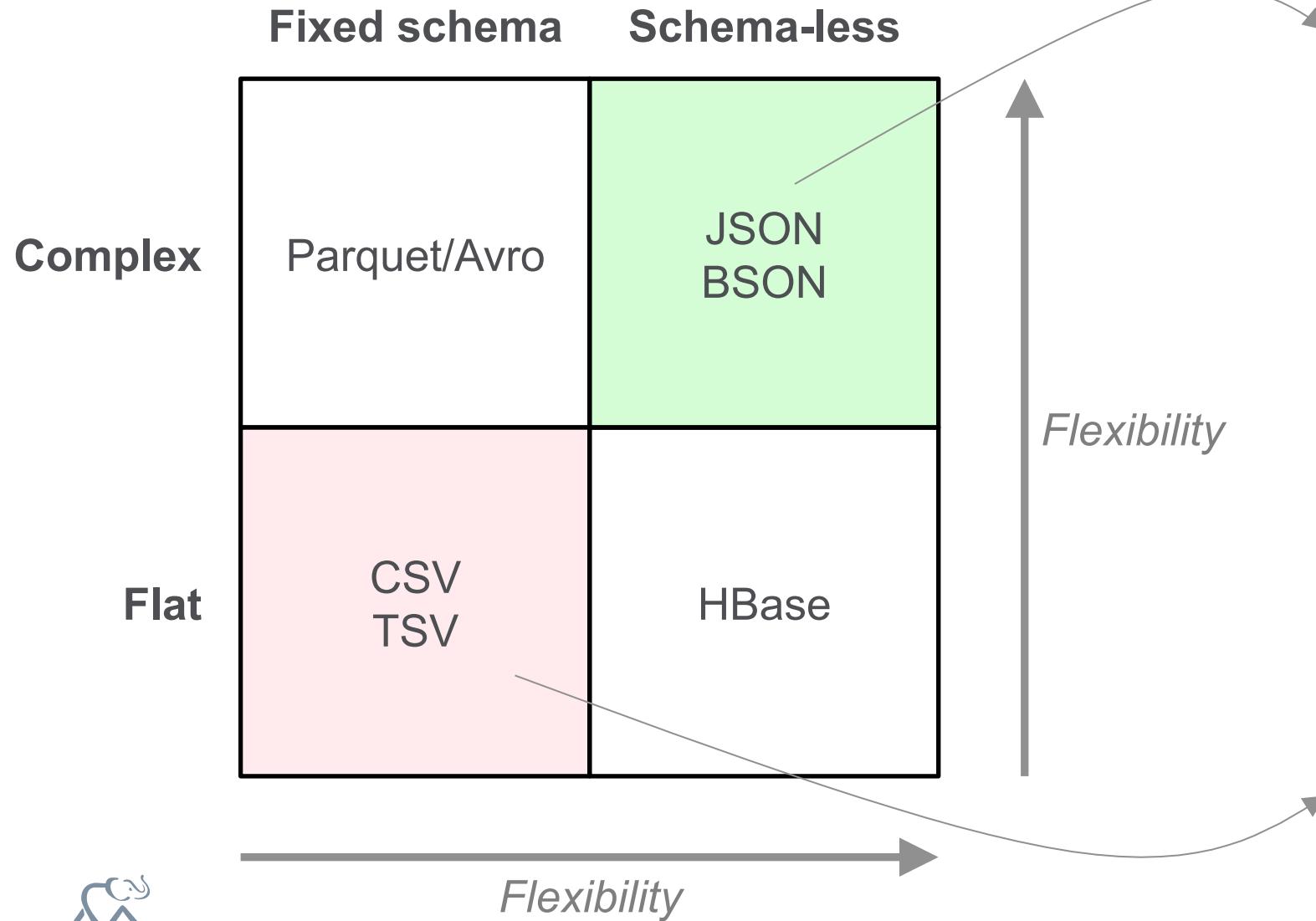
- pathnames
- HBase table
- Hive table



# How Drill Achieves Data Agility



# Drill's Data Model is Flexible



Apache Drill table

```
{  
  name: {  
    first: Michael,  
    last: Smith  
  },  
  hobbies: [ski, soccer],  
  district: Los Altos  
}  
{  
  name: {  
    first: Jennifer,  
    last: Gates  
  },  
  hobbies: [sing],  
  preschool: CCLC  
}
```

RDBMS/SQL-on-Hadoop table

| Name     | Gender | Age |
|----------|--------|-----|
| Michael  | M      | 6   |
| Jennifer | F      | 3   |
|          |        |     |



# Drill Supports *Schema Discovery On-The-Fly*

## Schema Declared In Advance

- Fixed schema
- Leverage schema in centralized repository (Hive Metastore)

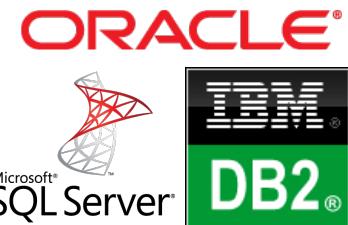
SCHEMA ON  
WRITE

SCHEMA  
BEFORE READ

## Schema Discovered On-The-Fly

- Fixed schema, evolving schema or schema-less
- Leverage schema in centralized repository or self-describing data

SCHEMA ON THE  
FLY



**MAPR**<sup>®</sup>

\$50M  
in Free Training



Free on-demand Hadoop training  
leading to certification

Start becoming an expert now  
[mapr.com/training](http://mapr.com/training)



# Q&A

Engage with us!

@mapr



maprtech

mapr-technologies



MapR

tshiran@mapr.com



maprtech



# Under the Hood



# High Level Architecture

## Cluster of commodity servers

- Daemon (drillbit) on each node

## ZooKeeper maintains ephemeral cluster membership information

- Drillbit uses ZooKeeper to find other drillbits in the cluster
- Client uses ZooKeeper to find drillbits

## Built-in, optimistic query execution engine. Doesn't require a particular storage or execution system (MapReduce, Spark, Tez)

- Better performance and manageability

## Data processing unit is *columnar record batches*

- Enables schema flexibility with negligible performance impact

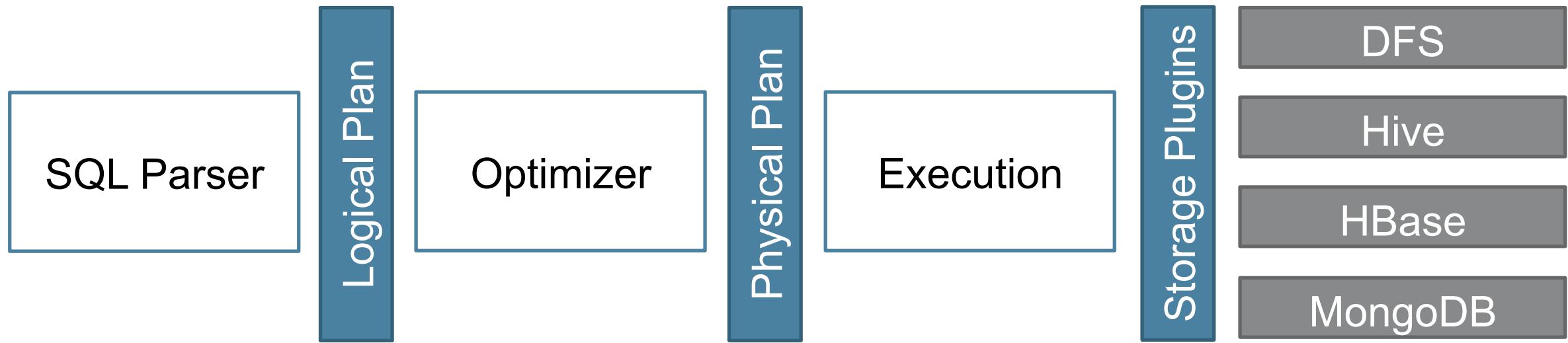


# Drill Maximizes Data Locality

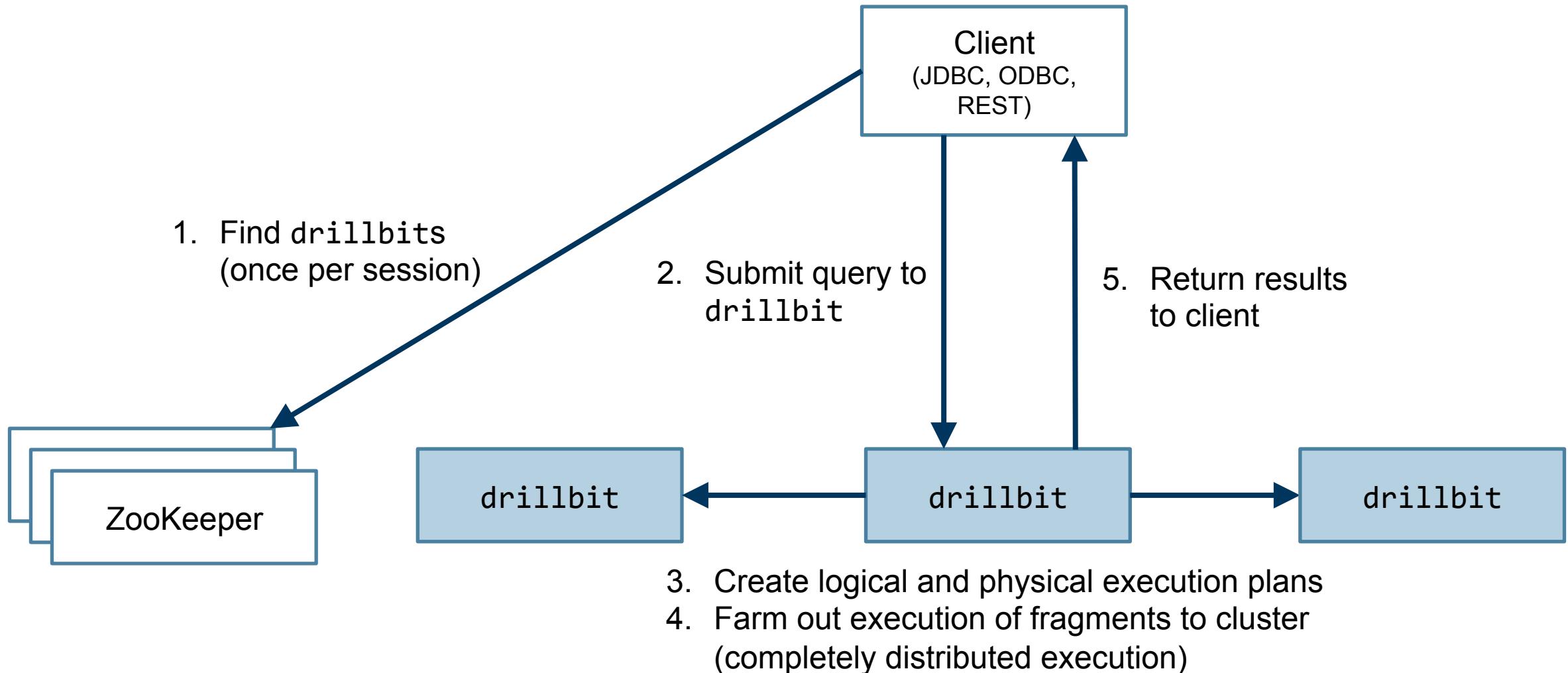


| Data Source      | Best Practice  |
|------------------|--|
| HDFS or MapR-FS  | drillbit on each DataNode  |
| HBase or MapR-DB | drillbit on each RegionServer  |
| MongoDB          | drillbit on each mongod node (when using replicas, run it on the replica node) |

# Core Modules within drillbit



# SELECT \* Query Execution



\* CTAS (CREATE TABLE AS SELECT) queries include steps 1-4