



SQL-ON-HADOOP FOR ANALYTICS + BI: WHAT ARE MY OPTIONS, WHAT'S THE FUTURE?

MARK RITTMAN, COO, RITTMAN MEAD

BUDAPEST DATA FORUM, JUNE 2016

Many Organisations are Running Big Data Initiatives

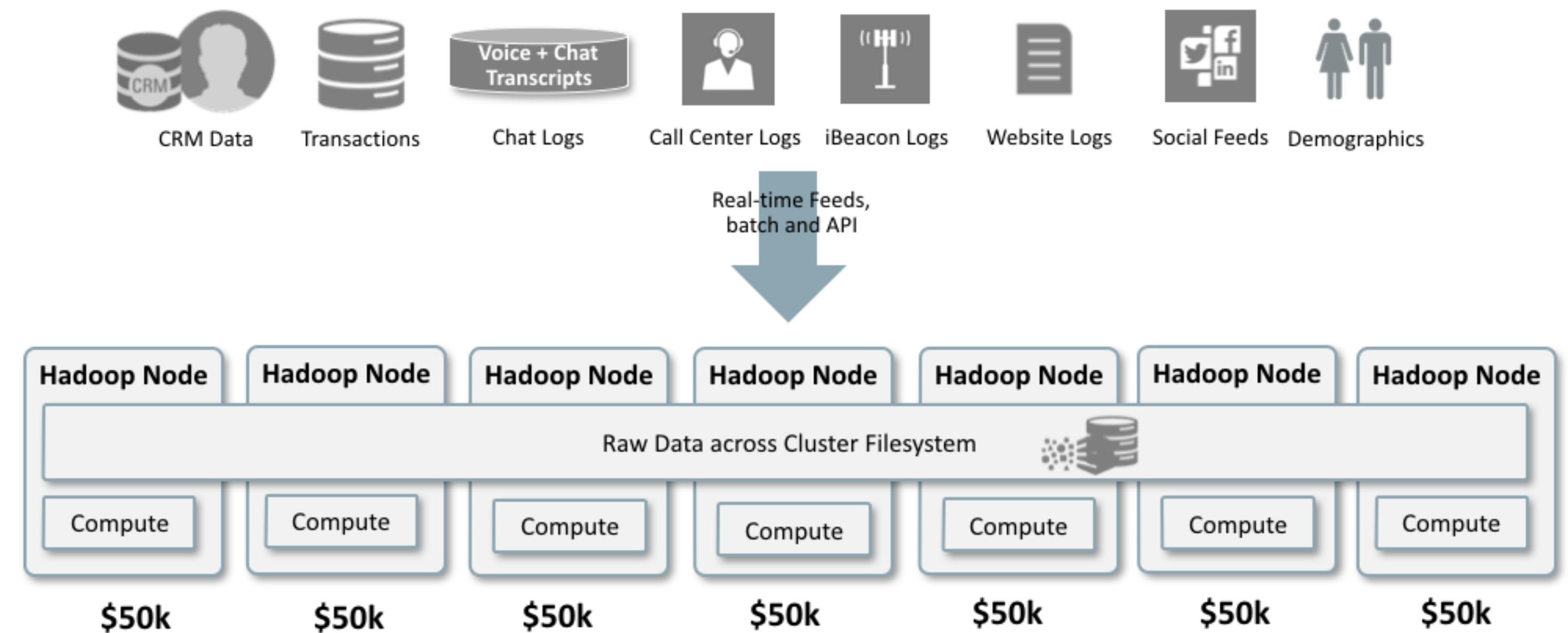
- Many customers and organisations are now running initiatives around “big data”
- Some are IT-led and are looking for cost-savings around data warehouse storage + ETL
- Others are “skunkworks” projects in the marketing department that are now scaling-up
- Projects now emerging from pilot exercises
- And design patterns starting to emerge

Digitization
× Datafication

Big Data
Disruption

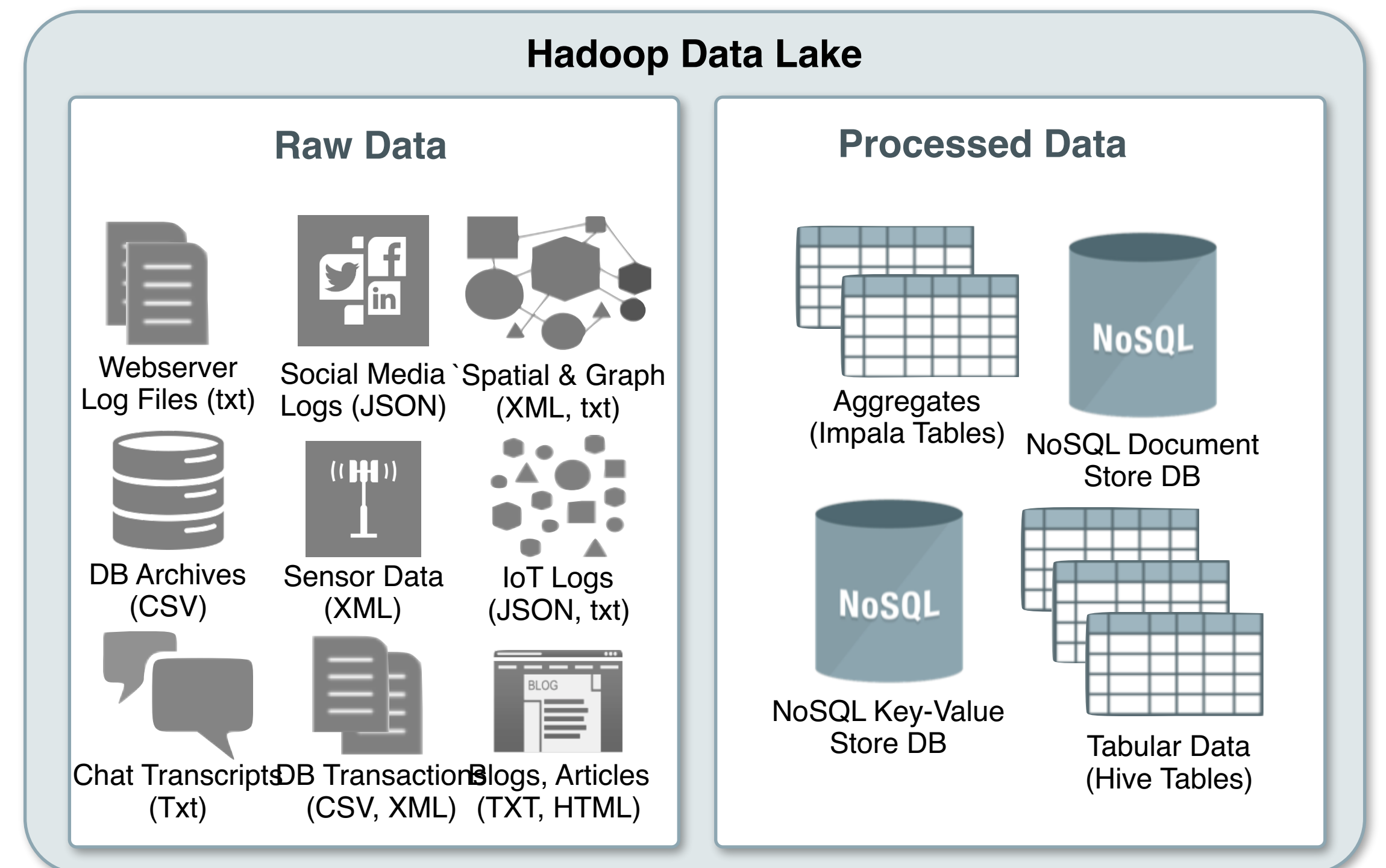
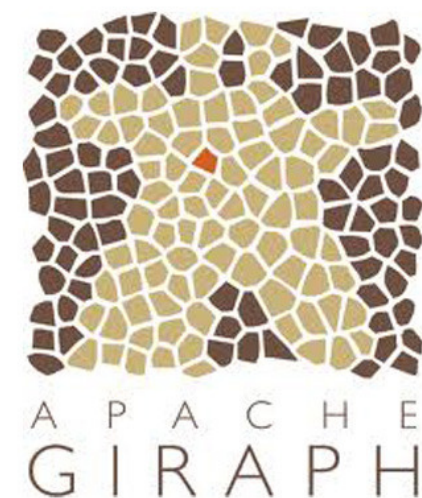
Highly Scalable (and Affordable) Cluster Computing

- Enterprise High-End RDBMSs such as Oracle can scale into the petabytes, using clustering
 - ▶ Sharded databases (e.g. Netezza) can scale further but with complexity / single workload trade-offs
- Hadoop was designed from outside for massive horizontal scalability - using cheap hardware
 - Anticipates hardware failure and makes multiple copies of data as protection
 - More nodes you add, more stable it becomes
 - And at a fraction of the cost of traditional RDBMS platforms



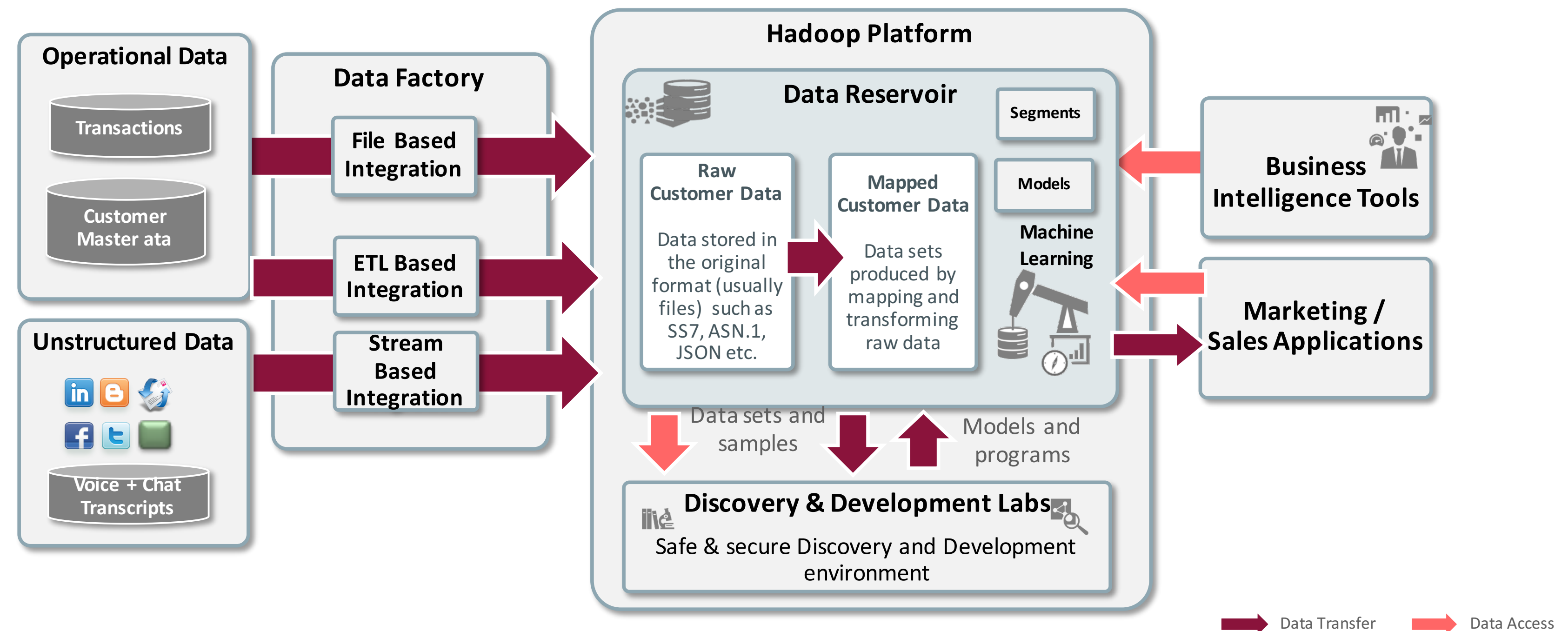
One Platform, Multiple Processing Frameworks

- We can now affordably create a single, massive archive of all our corporate data
- Separate to our OLTP and operational BI tools
- Leave it stored at the individual transaction level
- And then run multiple compute frameworks on it
 - ▶ SQL queries through Hive, Impala etc
 - ▶ Data processing using Spark, MR
 - ▶ Graph Analysis, etc



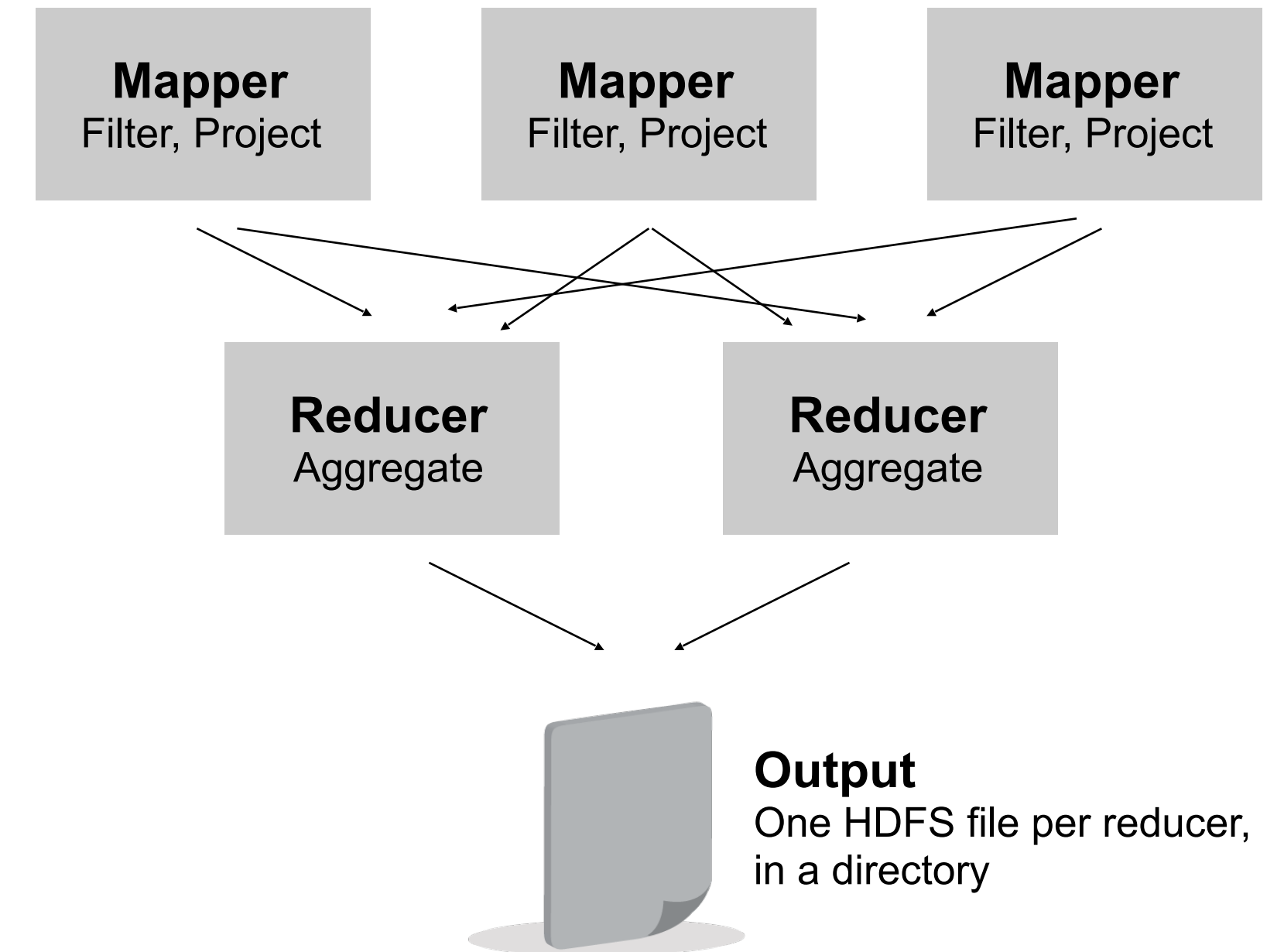
In the Context of BI & Analytics : The Data Reservoir

- Typical implementation of Hadoop and big data in an analytic context is the “data lake”
- Additional data storage platform with cheap storage, flexible schema support + compute
- Data lands in the data lake or reservoir in raw form, then minimally processed
- Data then accessed directly by “data scientists”, or processed further into DW



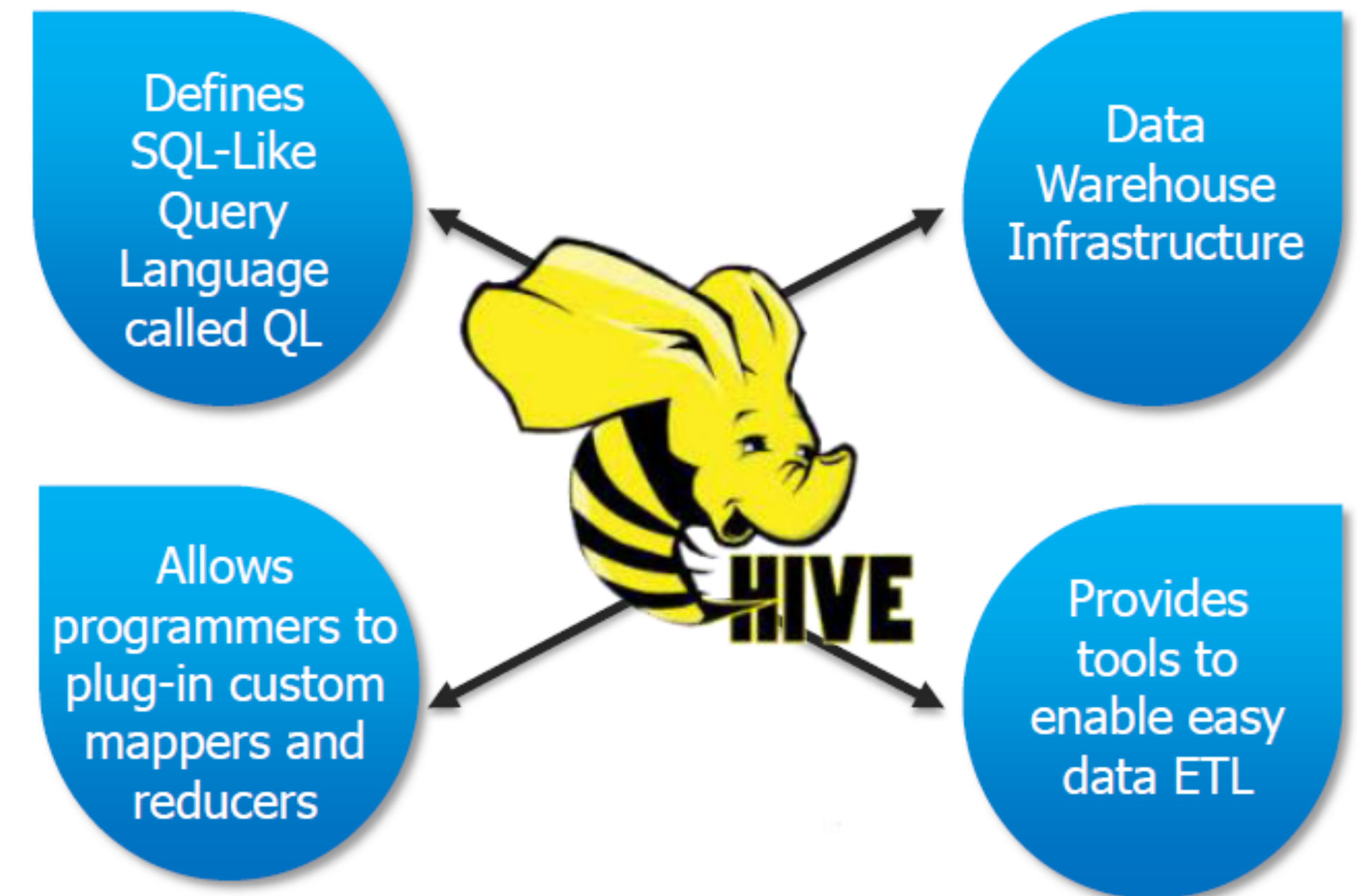
MapReduce - The Original Big Data Query Framework

- Programming model for processing large data sets in parallel on a cluster
- Not specific to a particular language, but usually written in Java
- Inspired by the **map** and **reduce** functions commonly used in functional programming
 - ▶ **Map ()** performs filtering and sorting
 - ▶ **Reduce ()** aggregates the output of mappers
 - ▶ and a **Shuffle ()** step to redistribute output by keys
- Resolved several complications of distributed computing:
 - ▶ Allows unlimited computations on unlimited data
 - ▶ Map and reduce functions can be easily distributed
 - ▶ Originated at Google; Hadoop was Yahoo's open-source implementation of MapReduce, + two are synonymous



Apache Hive : SQL Metadata + Engine over Hadoop

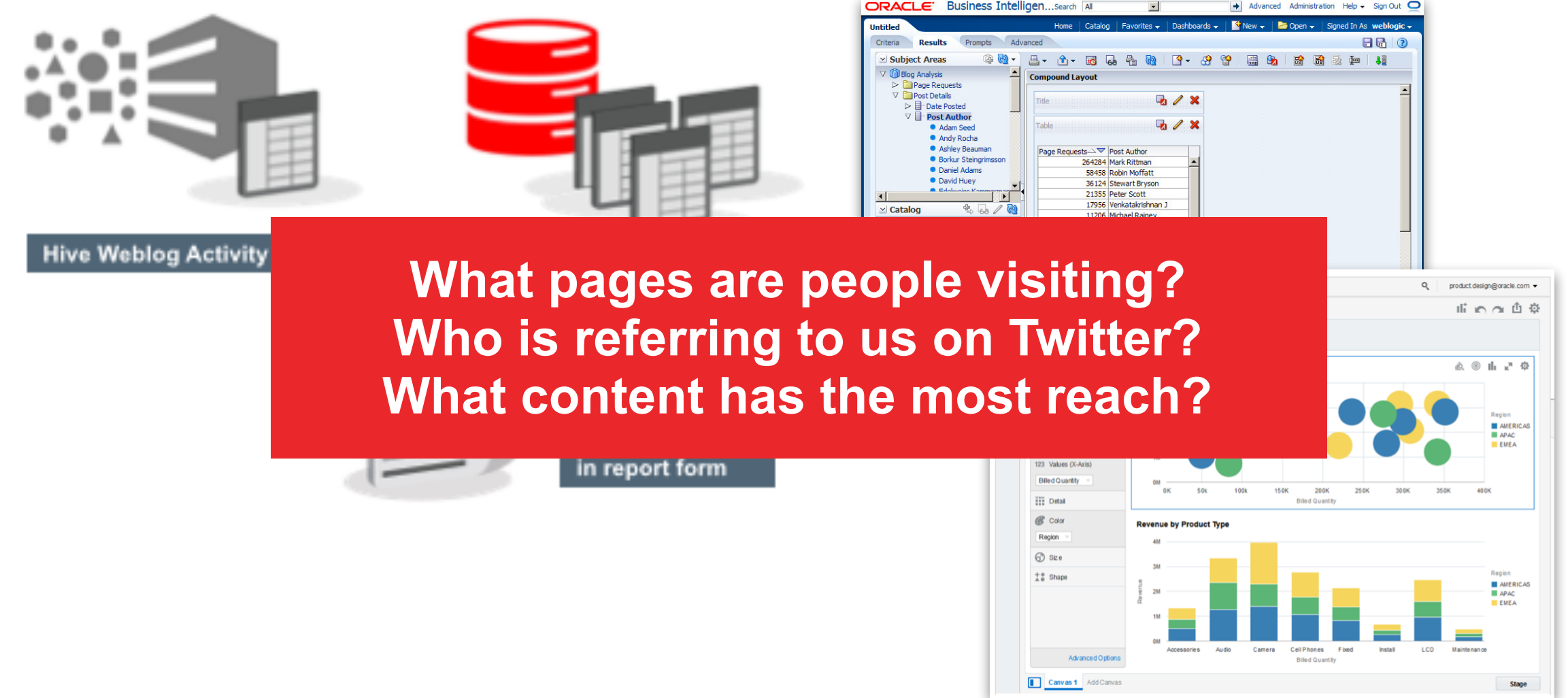
- Original developed at Facebook, now foundational within the Hadoop project
- Allows users to query Hadoop data using SQL-like language
- Tabular metadata layer that overlays files, can interpret semi-structured data (e.g. JSON)
- Generates MapReduce code to return required data
- Extensible through SerDes and Storage Handlers
- JDBC and ODBC drivers for most platforms/tools
- Perfect for set-based access + batch ETL work

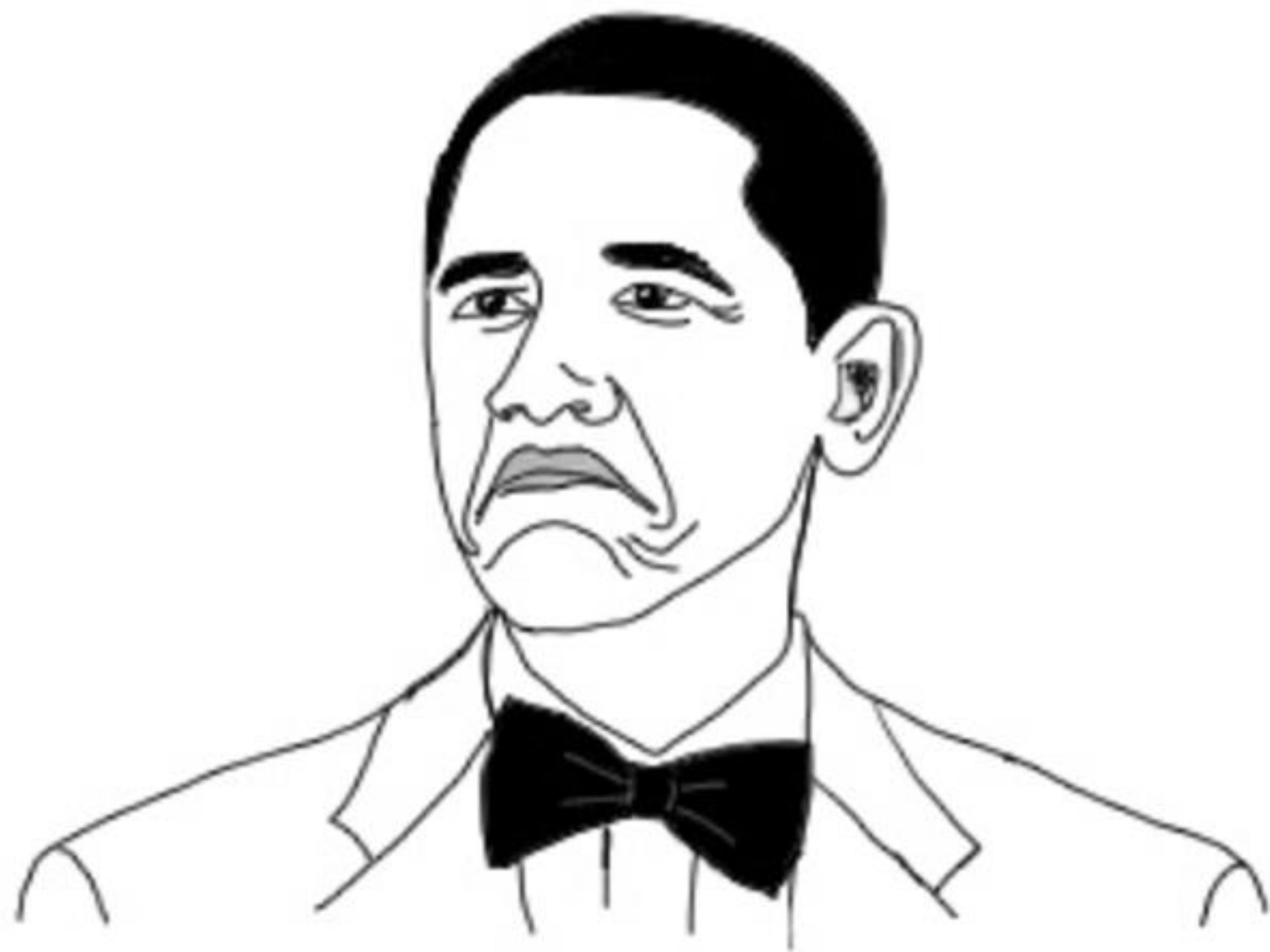


Hive Provides a SQL Interface for BI + ETL Tools

- Data integration tools could now load and process Hadoop data
- BI tools could treat Hadoop as just another data source
- Generally use MapReduce and Hive to access data
 - ▶ ODBC and JDBC access to Hive tabular data
 - ▶ Allows Hadoop unstructured/semi-structured data on HDFS to be accessed like RDBMS

Access direct Hive or extract using ODI12c for structured OBIEE dashboard analysis





NOT BAD



A close-up photograph of a snail moving across a mossy surface. The snail's shell is brown with dark, horizontal stripes. Its body is a pale, translucent yellowish color. The snail is positioned in the lower half of the frame, moving towards the left. The background is a soft-focus, green and brown mossy ground. The text "Hive is slow" is overlaid in white, bold, sans-serif font across the middle of the image, partially covering the snail's shell.

Hive is slow



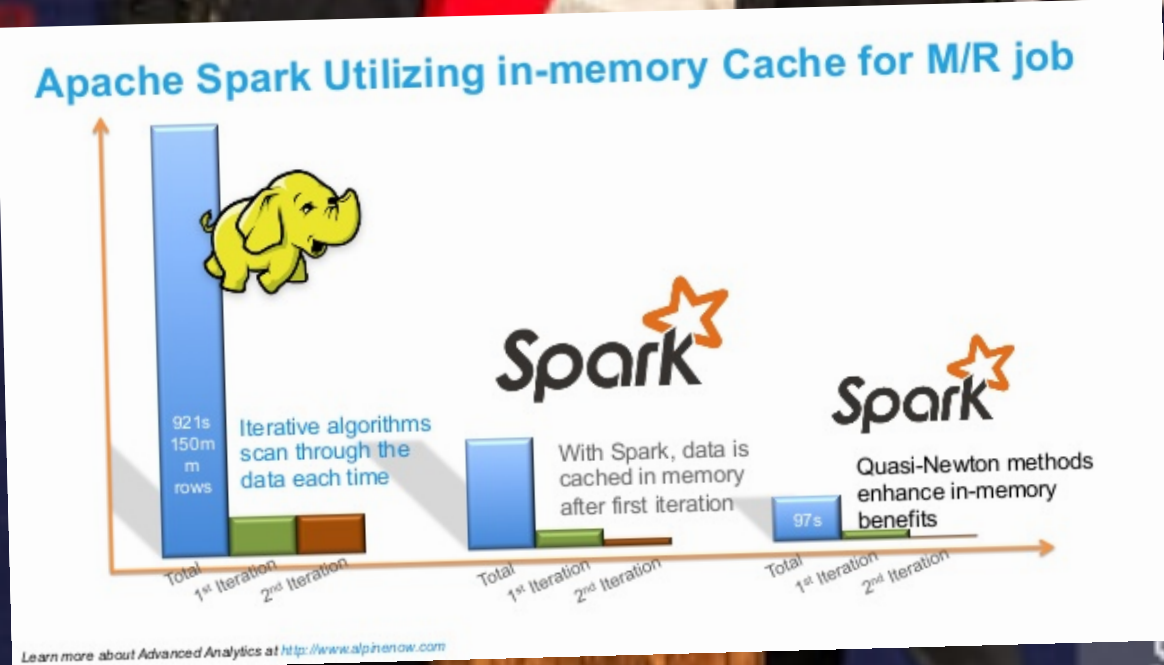
MapReduce is for losers



The background features a dark blue, isometric illustration of a digital workspace. On the left, there are several server racks. In the center, a desktop monitor displays a gear icon, with a keyboard in front of it. To the right, a laptop shows a line graph on its screen. Various other icons are scattered throughout, including folders, a magnifying glass, a cloud, a gear, and an envelope. A network of white lines with arrows connects these elements, suggesting data flow and connectivity.

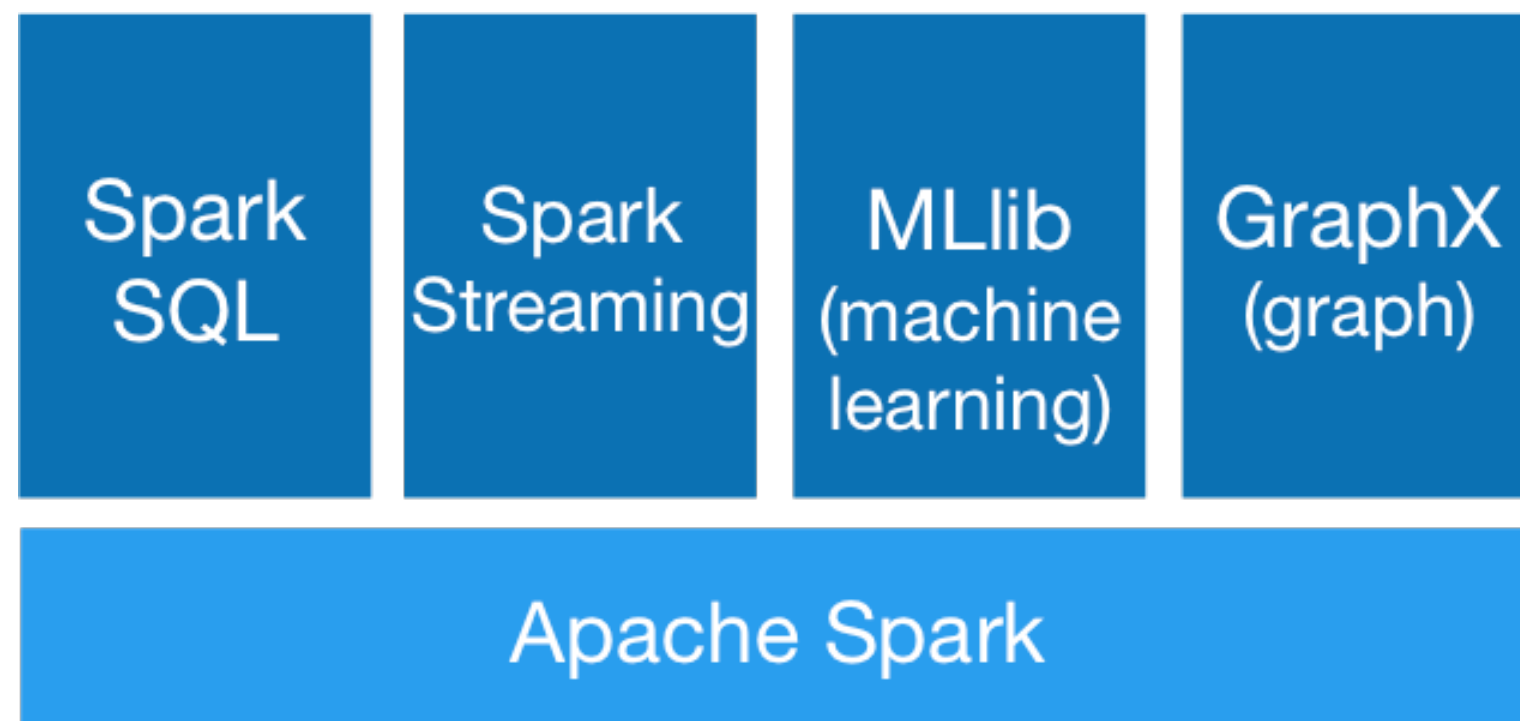
But the future is fast

MAKE AMERICA GREAT AGAIN!



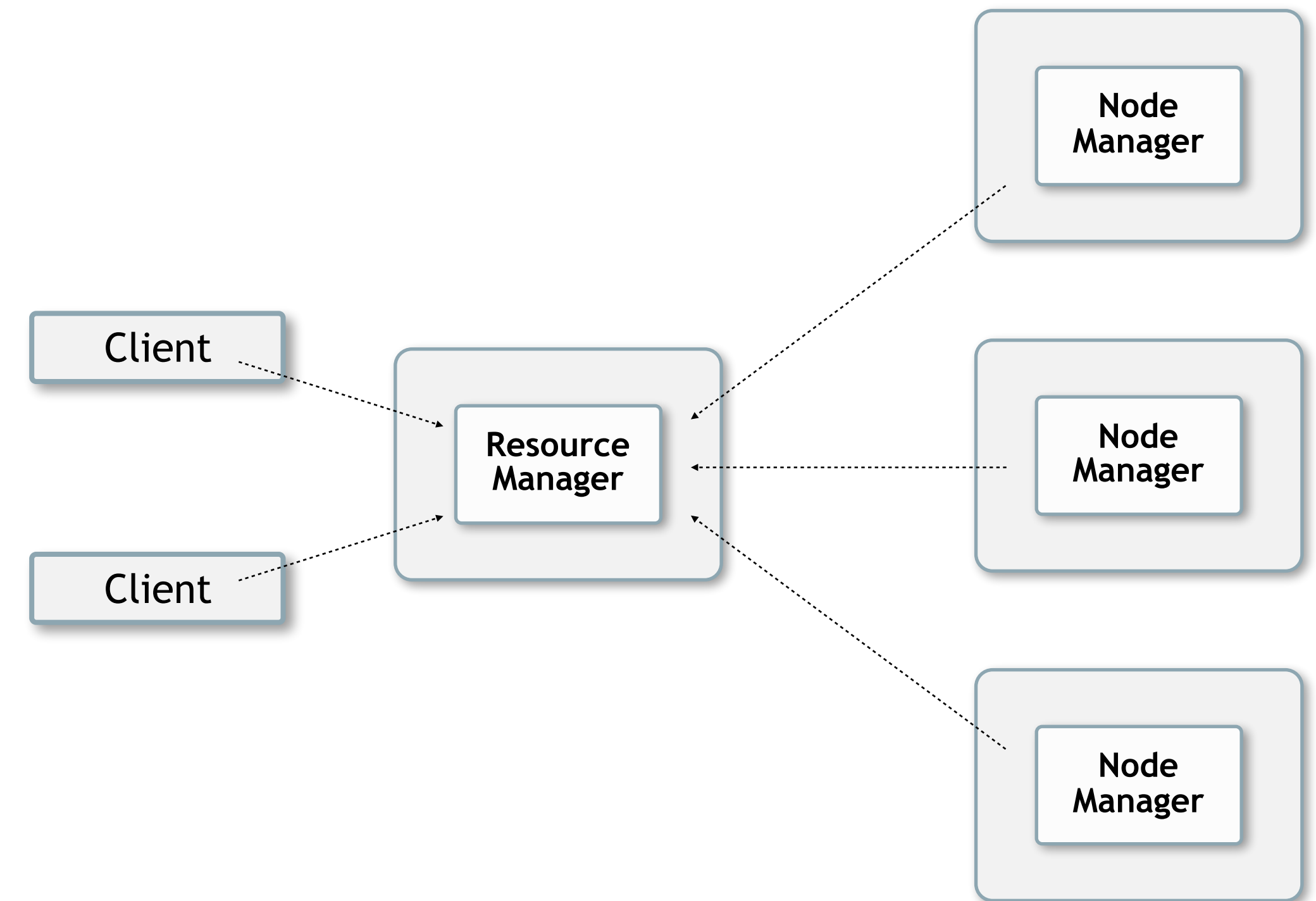
gettyimages®
Sean Rayford

Hadoop 2.0 Processing Frameworks + Tools



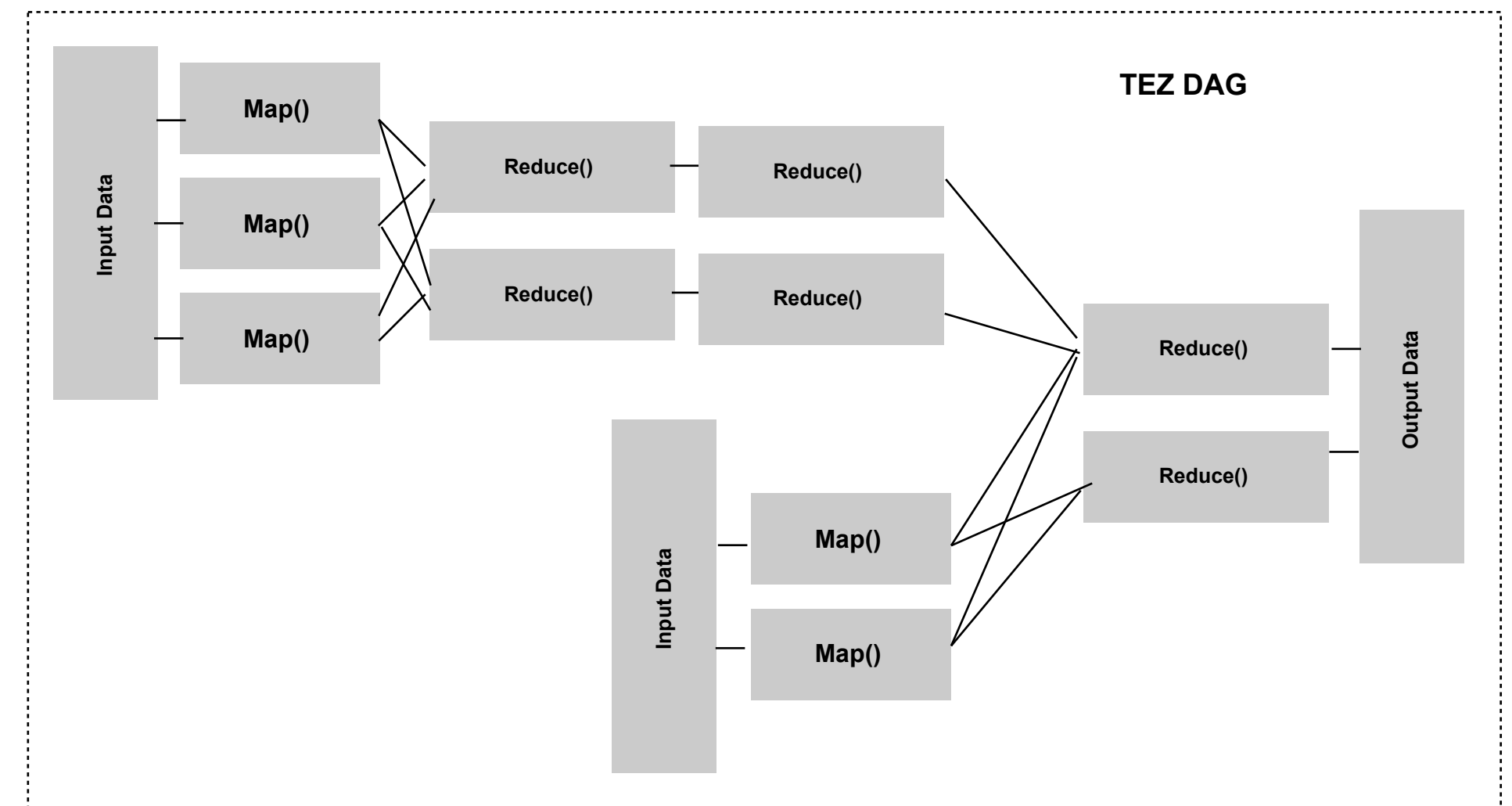
MapReduce 2 and YARN

- MapReduce 2 (MR2) splits the functionality of the JobTracker by separating resource management and job scheduling/monitoring
- Introduces YARN (Yet Another Resource Manager)
- Permits other processing frameworks to MR
 - ▶ For example, Apache Spark
- Maintains backwards compatibility with MR1
- Introduced with CDH5+

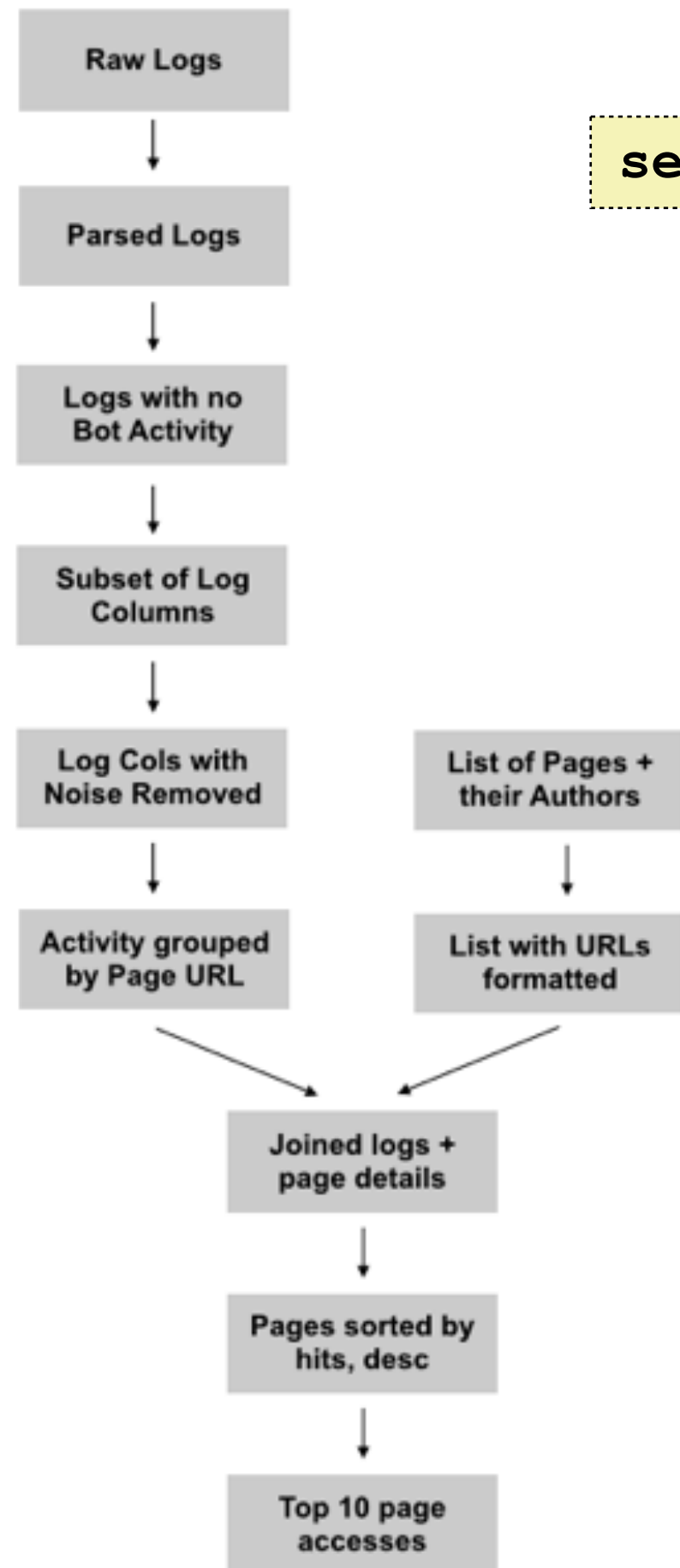


Apache Tez

- Runs on top of YARN, provides a faster execution engine than MapReduce for Hive, Pig etc
- Models processing as an entire data flow graph (DAG), rather than separate job steps
 - ▶ DAG (Directed Acyclic Graph) is a new programming style for distributed systems
 - ▶ Dataflow steps pass data between them as streams, rather than writing/reading from disk
- Supports in-memory computation, enables Hive on Tez (Stinger) and Pig on Tez
- Favoured In-memory / Hive v2 route by Hortonworks



Tez Advantage - Drop-In Replacement for MR with Hive, Pig



`set hive.execution.engine=mr`

JobId	Maps	Reduces	Alias
job_1417127396023_0145	12	2	logs_base, logs_base_nobots, logs_base_page, logs_
job_1417127396023_0146	2	1	pages_and_post_details, pages_and_posts_trim, pos
job_1417127396023_0147	1	1	pages_and_posts_sorted
job_1417127396023_0148	1	1	pages_and_posts_sorted
job_1417127396023_0149	1	1	pages_and_posts_sorted

4m 17s

```
markrittman -- root@bda4node1:~ -- ssh -- 97x11
Last login: Fri Dec 5 13:12:13 on ttys003
officeimac:~ markrittman$ ssh root@bda4node1
root@bda4node1's password:
Last login: Fri Dec 5 13:29:24 2014 from bda4node1.rittmandev.com
[root@bda4node1 ~]# pig -x tez analyzeblog.pig
14/12/06 17:16:37 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
14/12/06 17:16:37 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
14/12/06 17:16:37 INFO pig.ExecTypeProvider: Trying ExecType : TEZ_LOCAL
14/12/06 17:16:37 INFO pig.ExecTypeProvider: Trying ExecType : TEZ
14/12/06 17:16:37 INFO pig.ExecTypeProvider: Picked TEZ as the ExecType
2014-12-06 17:16:37,571 [main] INFO org.apache.pig.Main - Apache Pig version 0.14.0.2.2.0.0-2041
```

Dashboard metrics:

- CPU Usage: 100%
- Cluster Load: 12%
- NameNode Heap: 0.19 ms
- NameNode CPU: 0.0%
- NameNode Uptime: 1.1 d
- ResourceManager Heap: 5%
- ResourceManager Uptime: 1.1 d
- NodeManager Uptime: 5%

`set hive.execution.engine=tez`

2m 25s

Cloudera Impala - Fast, MPP-style Access to Hadoop Data

- Cloudera's answer to Hive query response time issues
- MPP SQL query engine running on Hadoop, bypasses MapReduce for direct data access
 - Mostly in-memory, but spills to disk if required
- Uses Hive metastore to access Hive table metadata
- Similar SQL dialect to Hive - not as rich though and no support for Hive SerDes, storage handlers etc



Enabling Hive Tables for Impala

- Log into Impala Shell, run INVALIDATE METADATA command to refresh Impala table list
- Run SHOW TABLES Impala SQL command to view tables available
- Run COUNT(*) on main ACCESS_PER_POST table to see typical response time

```
[oracle@bigdatalite ~]$ impala-shell
Starting Impala Shell without Kerberos authentication

[bigdatalite.localdomain:21000] > invalidate metadata;
Query: invalidate metadata

Fetched 0 row(s) in 2.18s
[bigdatalite.localdomain:21000] > show tables;
Query: show tables
+-----+
| name          |
+-----+
| access_per_post
| access_per_post_cat_author
| ...
| posts
+-----+
Fetched 45 row(s) in 0.15s
```

```
[bigdatalite.localdomain:21000] > select count(*)
                                from access_per_post;
Query: select count(*) from access_per_post
+-----+
| count(*) |
+-----+
| 343      |
+-----+
Fetched 1 row(s) in 2.76s
```

Apache Parquet - Column-Orientated Storage for Analytics

- Beginners usually store data in HDFS using text file formats (CSV) but these have limitations
- Apache AVRO often used for general-purpose processing
 - ▶ Splitability, schema evolution, in-built metadata, support for block compression
- Parquet now commonly used with Impala due to column-orientated storage
 - ▶ Mirrors work in RDBMS world around column-store
 - ▶ Only return (project) the columns you require across a wide table

APACHE PARQUET

Columnar storage for the people



Collocating entire columns allows for efficient column projection

Read off disk only the columns you need

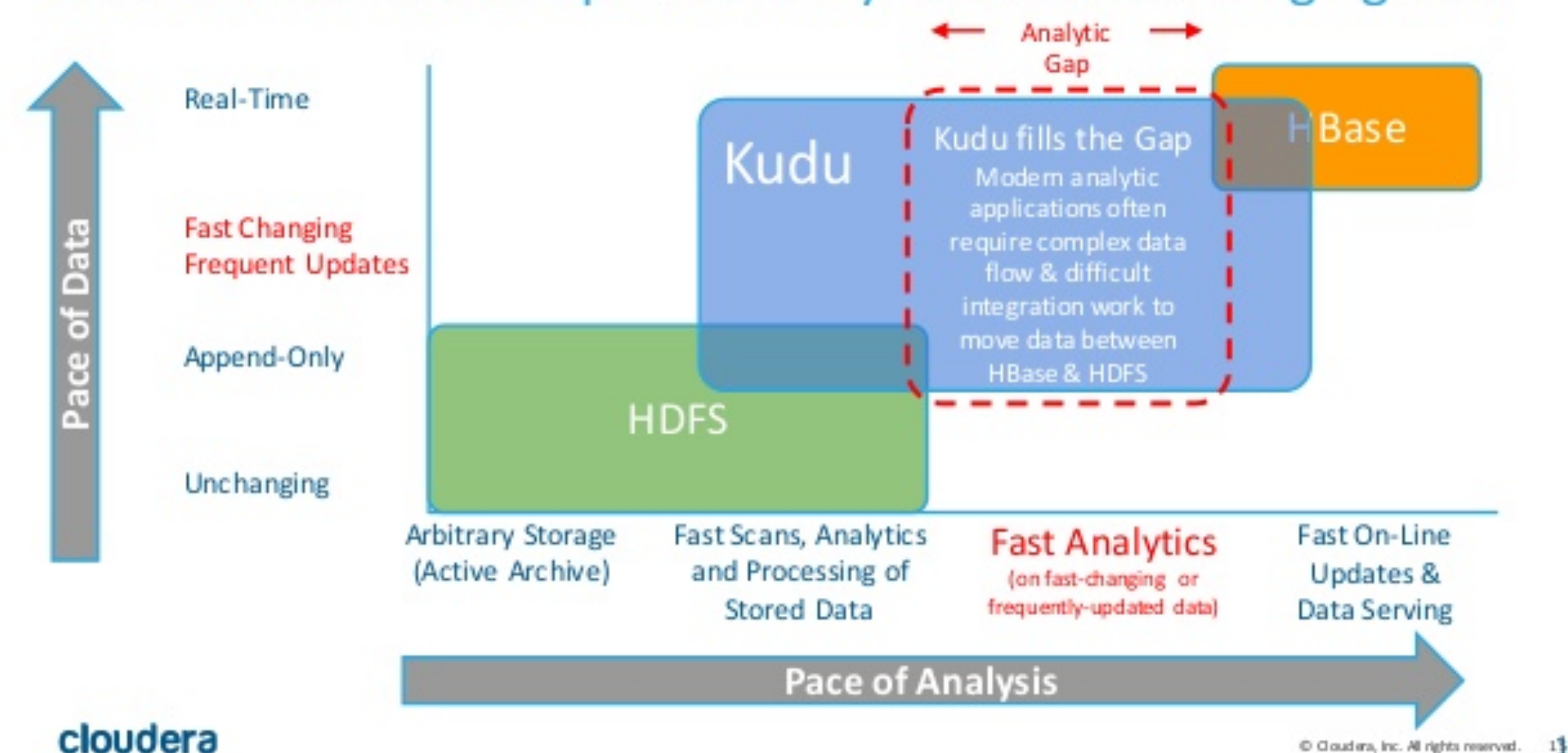
Possibly more importantly: deserialize only the columns you need

a	b	c
a1	b1	c1
a2	b2	c2
a3	b3	c3
a4	b4	c4
a5	b5	c5

Cloudera Kudu - Combining Best of HBase and Column-Store

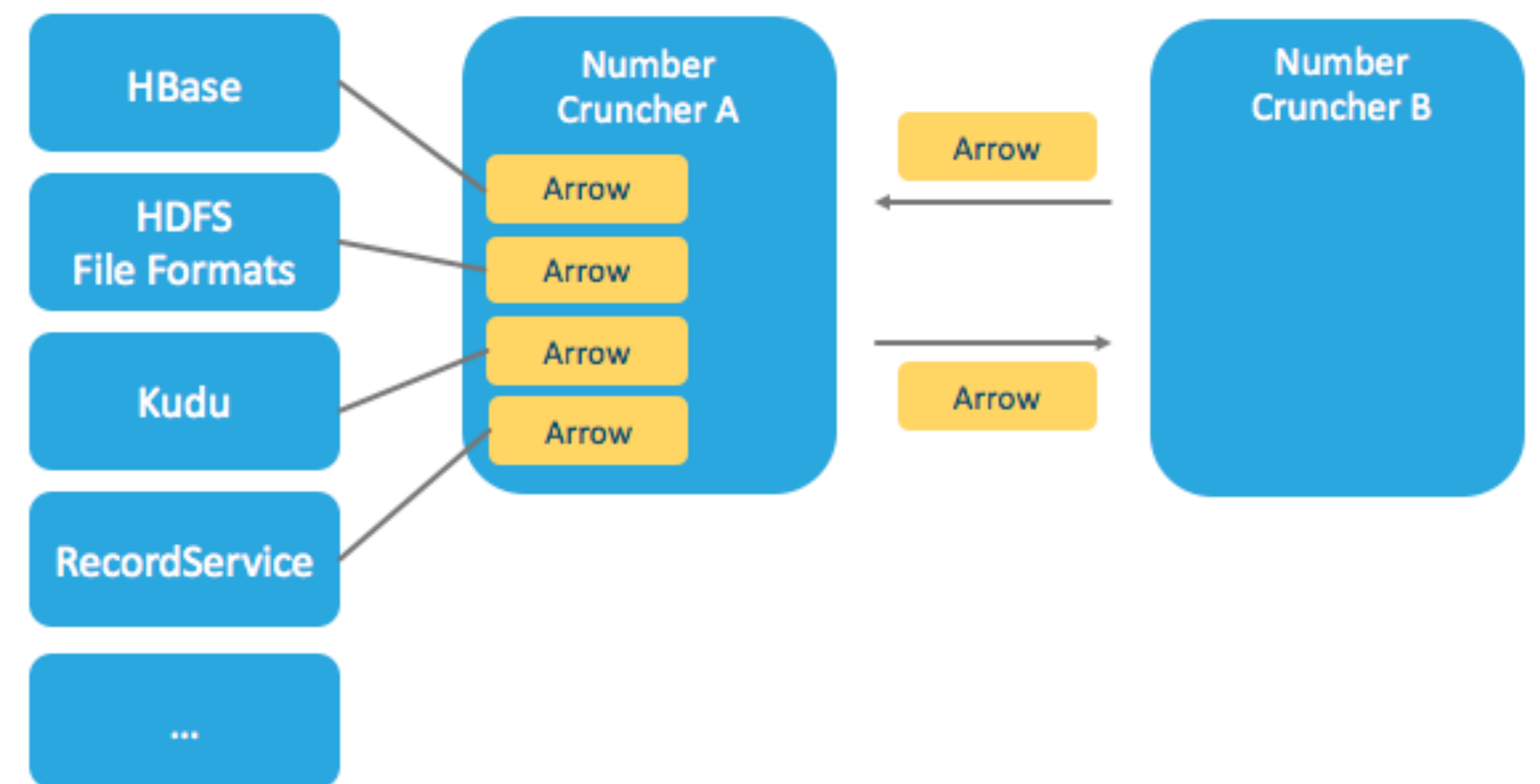
- But Parquet (and HDFS) have significant limitation for real-time analytics applications
 - ▶ Append-only orientation, focus on column-store makes streaming ingestion harder
- Kudu aims to combine best of HDFS + HBase
 - ▶ Real-time analytics-optimised
 - ▶ Supports updates to data
 - ▶ Fast ingestion of data
 - ▶ Accessed using SQL-style tables and get/put/update/delete API

Kudu Fills a Critical Gap: Fast Analytics on Fast Changing Data



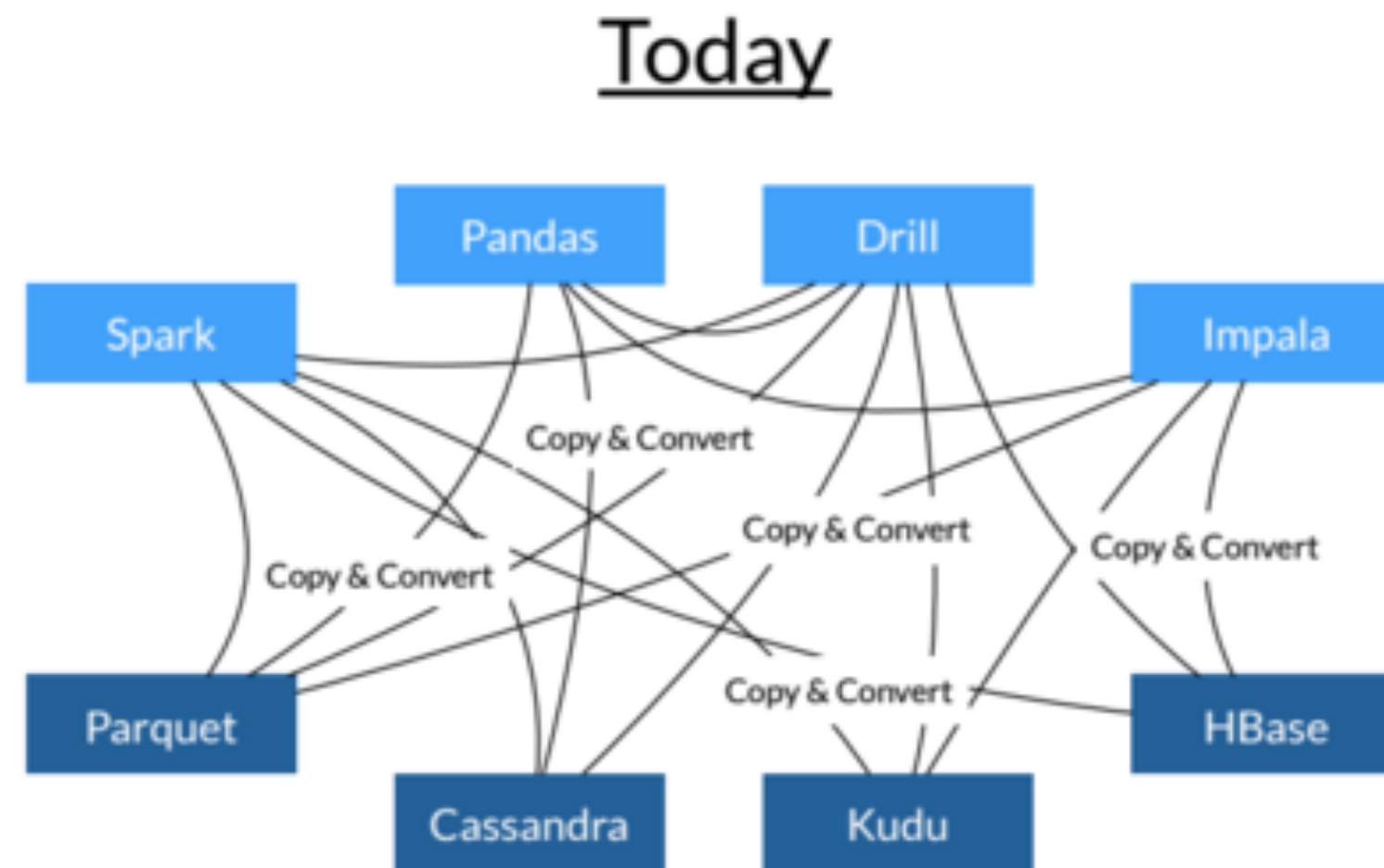
Apache Arrow - Standardising In-Memory Structures

- Many Hadoop tools now use in-memory processing (R, Python, Spark etc)
- But they all work to their own standards
- Considerable overhead in serialising / deserialising data between tools
- Apache Arrow standardised how in-memory data is held
- Considerably reduced overhead and latency between tools

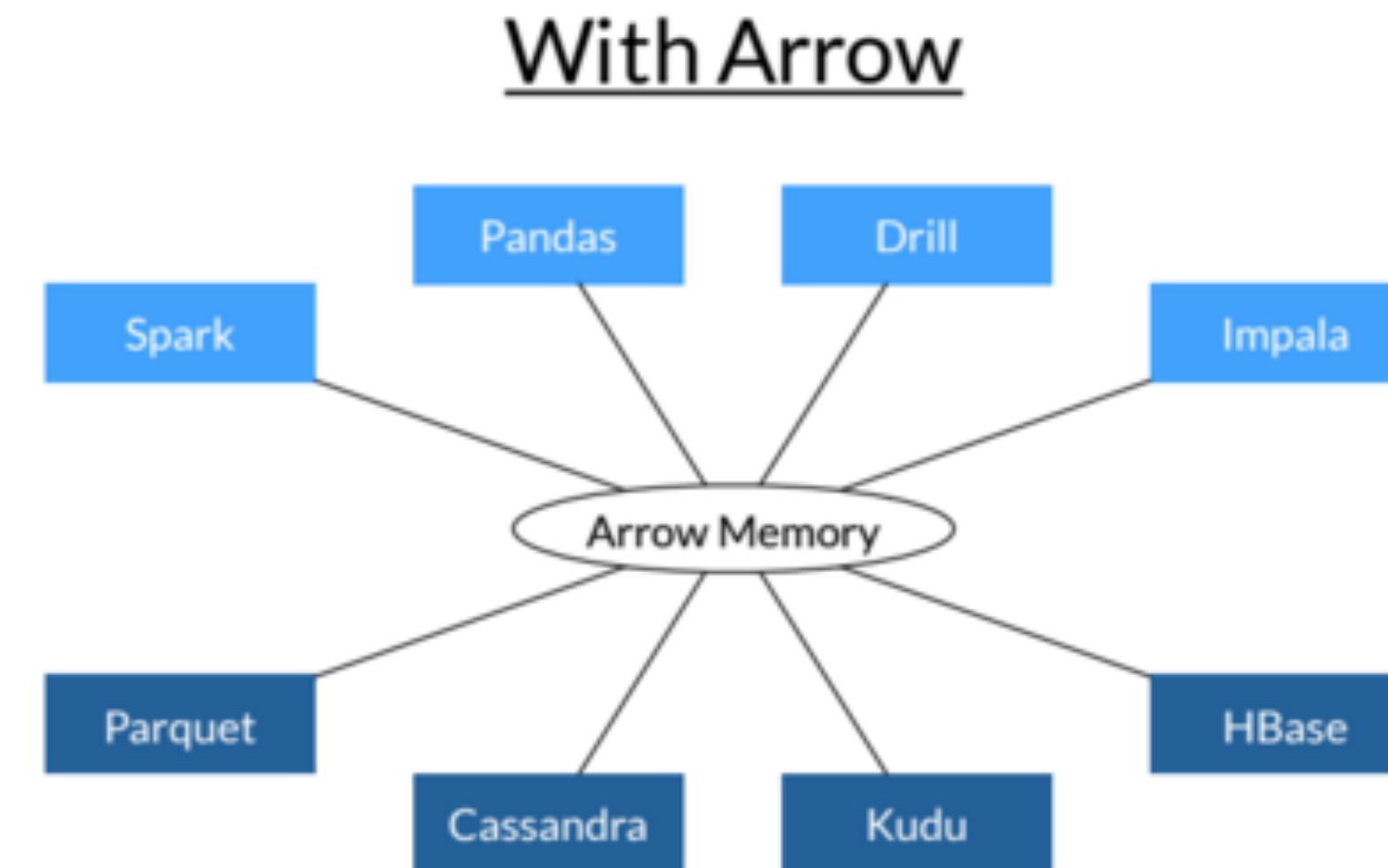


Apache Arrow : Common In-Memory Layer for Hadoop

Advantages of a Common Data Layer



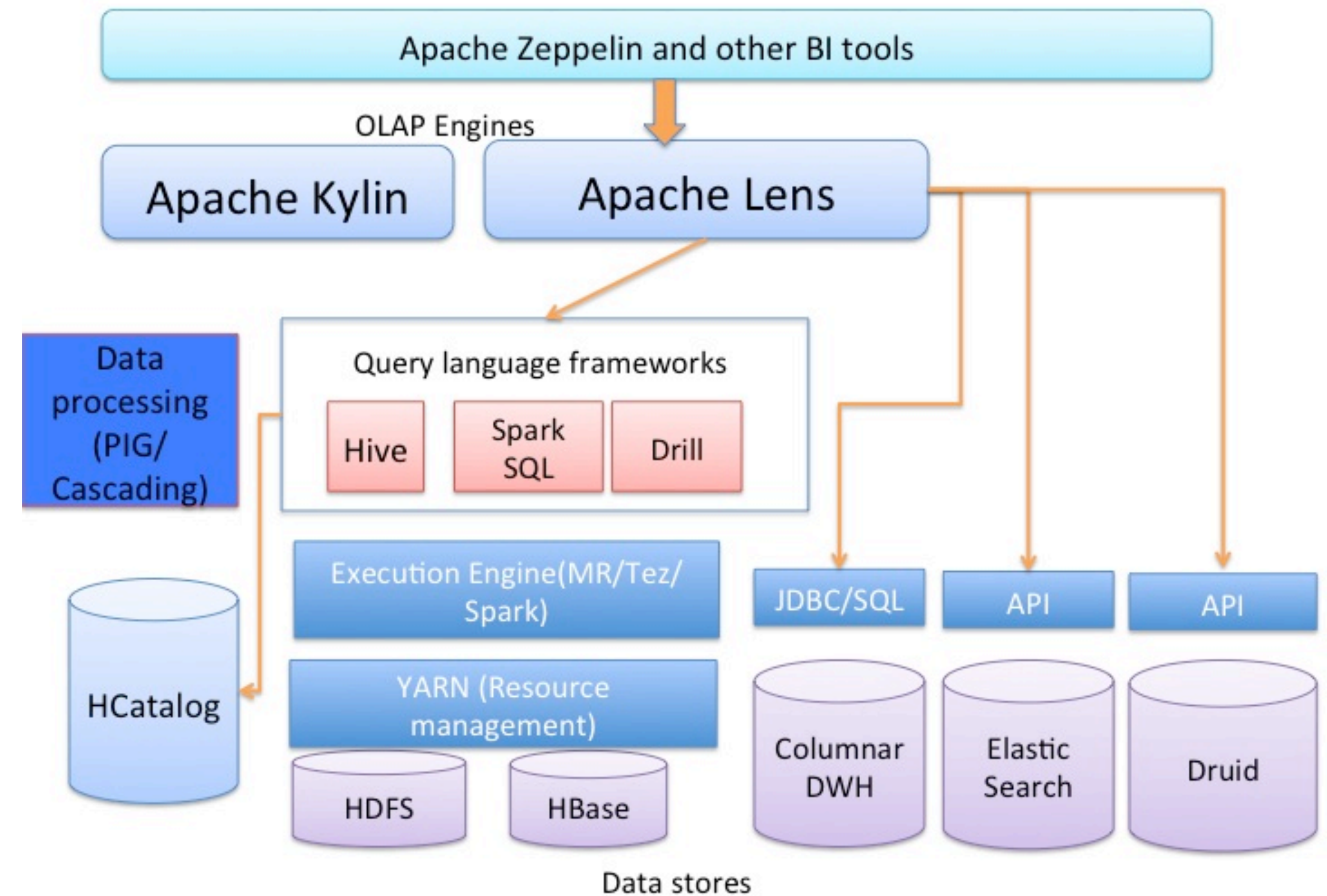
- Each system has its own internal memory format
- 70-80% CPU wasted on serialization and deserialization
- Similar functionality implemented in multiple projects



- All systems utilize the same memory format
- No overhead for cross-system communication
- Projects can share functionality (eg, Parquet-to-Arrow reader)

Apache Lens : Logical Dimensional Model for Hadoop

- Similar concept to Oracle BI Server
 - ▶ Logical star scheme for business model
 - ▶ Maps to federated data sources
 - ▶ Integrates and models data for query
 - ▶ Can run alongside Kylin MOLAP server
- Part of wider next-gen Hadoop BI stack
 - ▶ Apache Zeppelin web-based notebook
 - ▶ Spark/MR/Tez execution engines
 - ▶ Hive and Drill for ETL/ad-hoc SQL
- Slowly the BI stack gets built out...





In Oracle we have PL/SQL

The background features a dark blue, isometric illustration of a server room and network infrastructure. It includes several server racks, a central computer monitor with a gear icon, a keyboard, a laptop displaying a line graph, and various icons representing data storage (cylinders), folders, and network connections (arrows).

In SQL Server we have TSQL

An isometric illustration of a data center and computing environment. It features several server racks on the left, a central computer monitor with a gear icon, a keyboard, a laptop on the right displaying a line graph, and various icons representing data storage (cylinders), folders, and network connections. The background is a dark blue gradient with white lines and arrows indicating data flow.

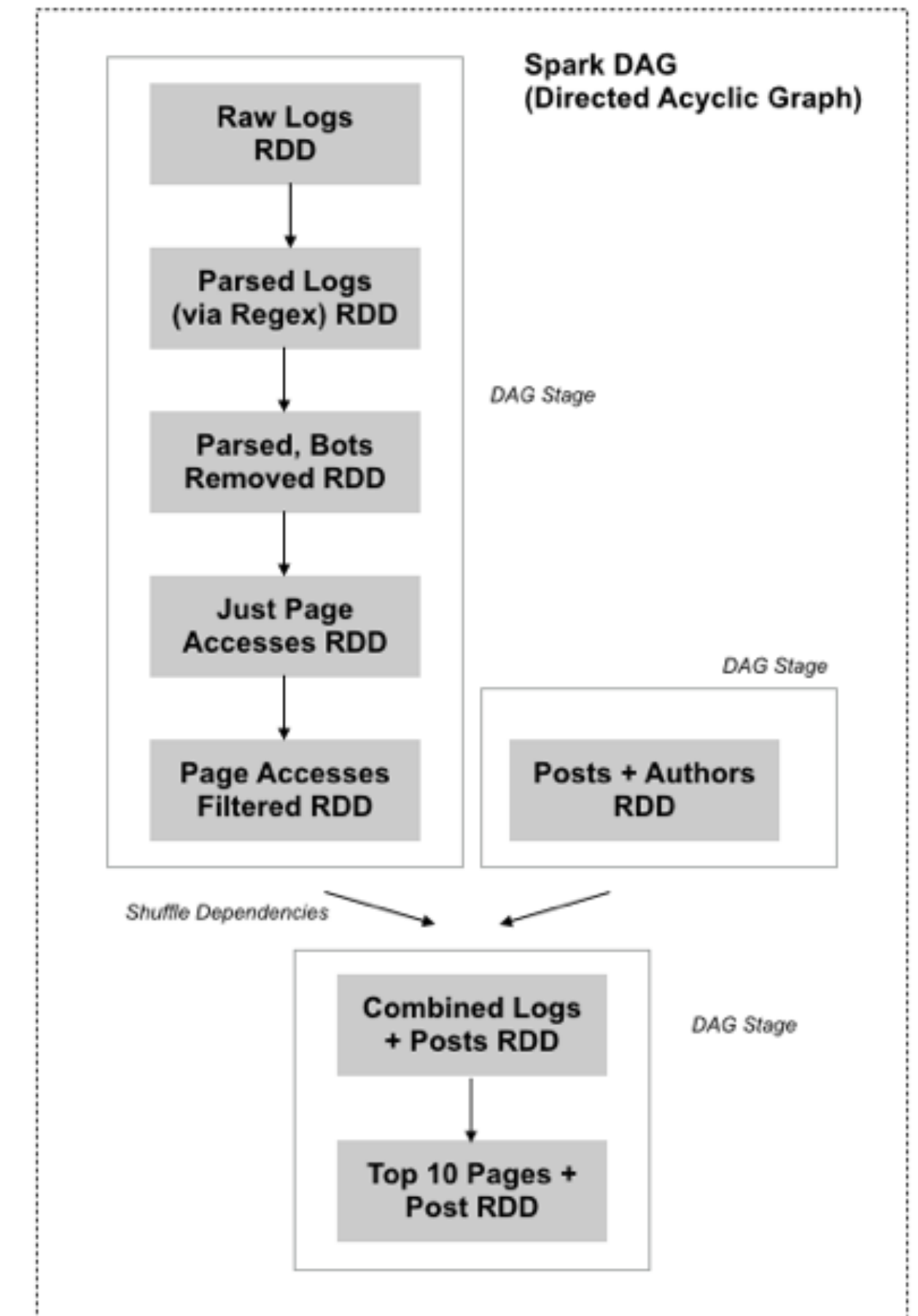
In Hadoop we have ... Spark



Lightning-Fast Cluster Computing

Apache Spark

- Another DAG execution engine running on YARN
- More mature than TEZ, with richer API and more vendor support
- Uses concept of an RDD (Resilient Distributed Dataset)
 - ▶ RDDs like tables or Pig relations, but can be cached in-memory
 - ▶ Great for in-memory transformations, or iterative/cyclic processes
- Spark jobs comprise of a DAG of tasks operating on RDDs
- Access through Scala, Python or Java APIs
- Related projects include
 - ▶ Spark SQL
 - ▶ Spark Streaming



Rich Developer Support + Wide Developer Ecosystem

- Native support for multiple languages with identical APIs
 - ▶ Python - prototyping, data wrangling
 - ▶ Scala - functional programming features
 - ▶ Java - lower-level, application integration
- Use of closures, iterations, and other common language constructs to minimize code
- Integrated support for distributed + functional programming
- Unified API for batch and streaming

```
scala> val logfile = sc.textFile("logs/access_log")
14/05/12 21:18:59 INFO MemoryStore: ensureFreeSpace(77353)
called with curMem=234759, maxMem=309225062
14/05/12 21:18:59 INFO MemoryStore: Block broadcast_2
stored as values to memory (estimated size 75.5 KB, free 294.6 MB)
logfile: org.apache.spark.rdd.RDD[String] =
MappedRDD[31] at textFile at <console>:15

scala> logfile.count()
14/05/12 21:19:06 INFO FileInputFormat: Total input paths to process : 1
14/05/12 21:19:06 INFO SparkContext: Starting job: count at <console>:1
...
14/05/12 21:19:06 INFO SparkContext: Job finished:
count at <console>:18, took 0.192536694 s
res7: Long = 154563
```

```
scala> val logfile = sc.textFile("logs/access_log").cache
scala> val biapps11g = logfile.filter(line => line.contains("/biapps11g/"))
biapps11g: org.apache.spark.rdd.RDD[String] = FilteredRDD[34] at filter at <console>:17
scala> biapps11g.count()
...
14/05/12 21:28:28 INFO SparkContext: Job finished: count at <console>:20, took 0.387960876 s
res9: Long = 403
```


Spark SQL - SQL within Apache Spark

- Spark SQL, and Data Frames, allow RDDs in Spark to be processed using SQL queries
- Bring in and federate additional data from JDBC sources
- Load, read and save data in Hive, Parquet and other structured tabular formats

```
val accessLogsFilteredDF = accessLogs
    .filter( r => ! r.agent.matches(".*(spider|robot|bot|slurp).*"))
    .filter( r => ! r.endpoint.matches(".*(wp-content|wp-admin).*")) .toDF()
    .registerTempTable("accessLogsFiltered")

val topTenPostsLast24Hour = sqlContext.sql("SELECT p.POST_TITLE, p.POST_AUTHOR, COUNT(*)
    as total
    FROM accessLogsFiltered a
    JOIN posts p ON a.endpoint = p.POST_SLUG
    GROUP BY p.POST_TITLE, p.POST_AUTHOR
    ORDER BY total DESC LIMIT 10 ")

// Persist top ten table for this window to HDFS as parquet file

topTenPostsLast24Hour.save("/user/oracle/rm_logs_batch_output/topTenPostsLast24Hour.parquet"
    , "parquet", SaveMode.Overwrite)
```

CREATING TABLES AND DDL COMMANDS



**ISN'T THIS JUST THE SAME AS
WE DO NOW?**

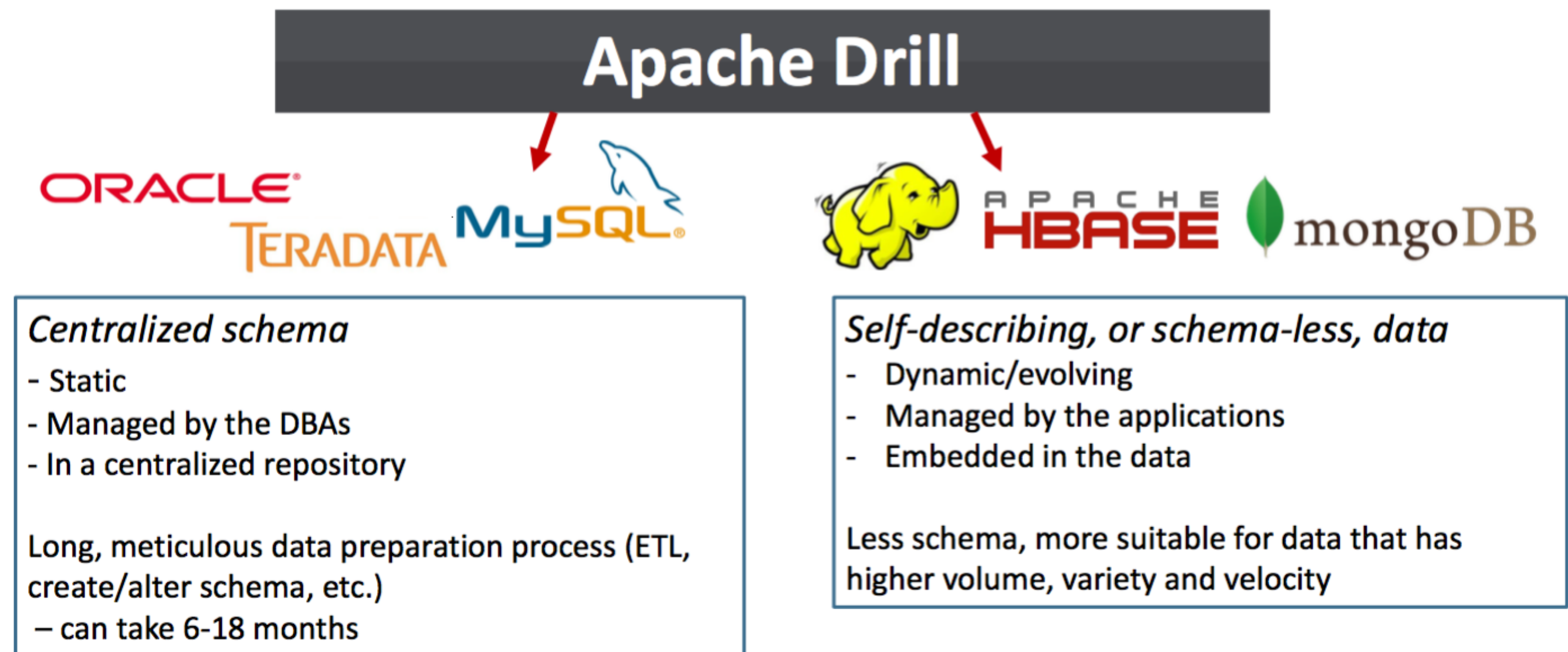
Apache Drill - Leverage Embedded Metadata in Files

- Apache Drill is another SQL-on-Hadoop project that focus on schema-free data discovery
- Inspired by Google Dremel, innovation is querying raw data with schema optional
 - Automatically infers and detects schema from semi-structured datasets and NoSQL DBs
 - Join across different silos of data e.g. JSON records, Hive tables and HBase database
 - Aimed at different use-cases than Hive - low-latency queries, discovery (think Endeca vs OBIEE)



Self-Describing Data - Parquet, AVRO, JSON etc

- Most modern datasource formats embed their schema in the data (“schema-on-read”)
- Apache Drill makes these as easy to join to traditional datasets as “point me at the data”
 - Cuts out unnecessary work in defining Hive schemas for data that’s self-describing
 - Supports joining across files, databases, NoSQL etc



Apache Drill and Text Files

- Files can exist either on the local filesystem, or on HDFS
- Connection to directory or file defined in storage configuration
 - Can work with CSV, TXT, TSV etc
 - First row of file can provide schema (column names)

```
SELECT * FROM dfs.`/tmp/csv_with_header.csv2`;
```

name	num1	num2	num3
hello	1		
hello	1		
hello	1		
hello	1		
hello	1		
hello	1		
hello	1		

```
SELECT * FROM dfs.`/tmp/csv_no_header.csv`;
```

columns
["hello", "1", "2", "3"]
["hello", "1", "2", "3"]
["hello", "1", "2", "3"]
["hello", "1", "2", "3"]
["hello", "1", "2", "3"]
["hello", "1", "2", "3"]
["hello", "1", "2", "3"]

7 rows selected (0.112 seconds)

7 rows selected (0.112 seconds)

The screenshot shows the Apache Drill web interface. At the top, there are navigation tabs: Apache Drill, Query, Profiles, Storage, Metrics, and Threads. Below these, the 'Enabled Storage Plugins' section is visible, showing a table with one entry: 'cp'. To the right of 'cp' are buttons for 'Update' and 'Disable'. A 'Configuration' modal window is open, displaying the configuration for the 'cp' plugin. The configuration is a JSON object with the following structure:

```
{
  "type": "file",
  "enabled": true,
  "connection": "file:///",
  "config": null,
  "workspaces": {
    "iot": {
      "location": "/home/iot/iot_demo/comms_travel",
      "writable": true,
      "defaultInputFormat": null
    }
  },
  "formats": {
    "csv": {
      "type": "text",
      "extensions": [
        "csv2"
      ],
      "extractHeader": true,
      "delimiter": ","
    }
  }
}
```

At the bottom of the configuration modal, there are buttons for 'Back', 'Update', 'Disable', and 'Delete'.

Apache Drill and JSON Documents

- JSON (Javascript Object Notation) documents are often used for data interchange
 - Exports from Twitter and other consumer services
 - Web service responses and other B2B interfaces
- A more lightweight form of XML that is “self-describing”
- Handles evolving schemas, and optional attributes
- Drill treats each document as a row, and has features to
 - Flatten nested data (extract elements from arrays)
 - Generate key/value pairs for loosely structured data

Configuration

```
{
  "type": "file",
  "enabled": true,
  "connection": "file:///",
  "config": null,
  "workspaces": {
    "iot": {
      "location": "/home/iot/iot_demo/comms_t",
      "writable": true,
      "defaultInputFormat": null
    }
  },
  "formats": {
    "json": {
      "type": "json",
      "extension": "json"
    }
  }
}
```

Back Update

Apache Drill Query Profiles Storage Metrics Threads

Enabled Storage Plugins

cp	Update	Disable
dfs	Update	Disable

Disabled Storage Plugins

hbase	Update	Enable
hive	Update	Enable
kudu	Update	Enable

```
use dfs.iot;
show files;
select in_reply_to_user_id, text from `all_tweets.json`
limit 5;
```

```
+-----+
| in_reply_to_user_id | text |
+-----+
| null | BI Forum 2013 in Brighton has now sold-out |
| null | "Football has become a numbers game |
```

```
select name, flatten(fillings) as f
from dfs.users.`/donuts.json`
where f.cal < 300;
```

```
Book |
e Past" |
miming |
```

Apache Drill and Hive, HBase, Parquet Sources etc

- Drill can connect to Hive to make use of metastore (incl. multiple Hive metastores)
- NoSQL databases (HBase etc)
- Parquet files (native storage format - columnar + self describing)

```
USE hbase;  
SELECT * FROM students;
```

row_key	account	
[B@e6d9eb7	{"name": "QWxpY2U="}	{"state"
[B@2823a2b4	{"name": "Qm9i"}	{"state"
[B@3b8eec02	{"name": "RnJhbms="}	{"state"
[B@242895da	{"name": "TWfyeQ=="}	{"state"

```
SELECT * FROM dfs.`iot_demo/geodata/region.parquet`;
```

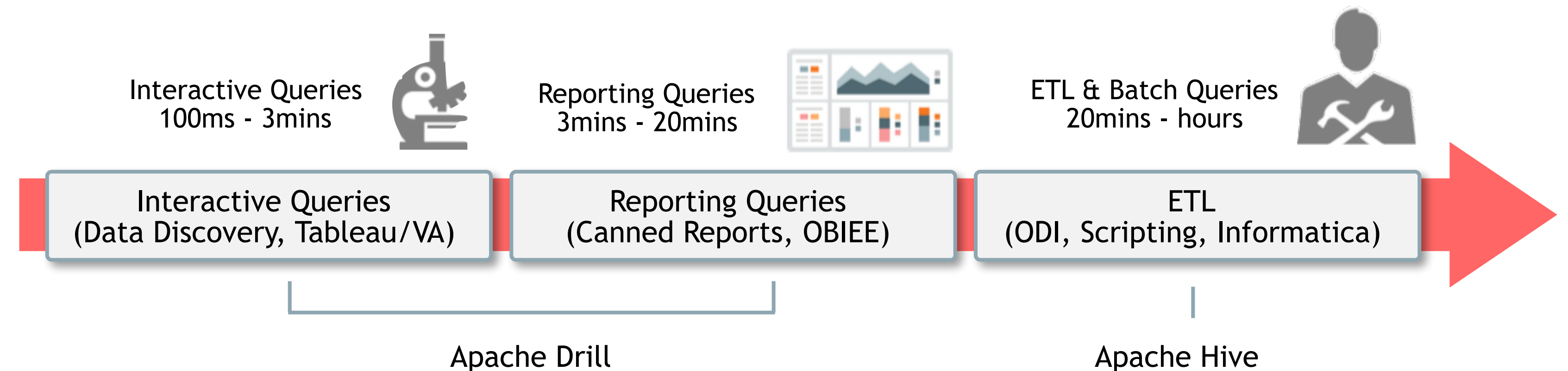
R_REGIONKEY	R_NAME	R_COMMENT
0	AFRICA	lar deposits. blithe
1	AMERICA	hs use ironic, even
2	ASIA	ges. thinly even pin
3	EUROPE	ly final courts cajo
4	MIDDLE EAST	uickly special accou

```
SELECT firstname,lastname FROM  
hiveremote.`customers` limit 10;`
```

firstname	lastname
Essie	Vaill
Cruz	Roudabush
Billie	Tinnes
Zackary	Mockus
Rosemarie	Fifield

Apache Drill vs. Apache Hive

- Drill developed for **real-time, ad-hoc data exploration** with **schema discovery on-the-fly**
 - Individual analysts exploring new datasets, leveraging corporate metadata/data to help
 - Hive is more about large-scale, centrally curated set-based big data access
- **Drill models conceptually as JSON**, vs. Hive's tabular approach
- **Drill introspects schema** from whatever it connects to, vs. formal modeling in Hive







www.rittmanmead.com



SQL-ON-HADOOP FOR ANALYTICS + BI: WHAT ARE MY OPTIONS, WHAT'S THE FUTURE?

MARK RITTMAN, COO, RITTMAN MEAD

BUDAPEST DATA FORUM, JUNE 2016