

Data Engineering

April 2019 - STL Big Data User Group



Discussion Topics

Introductions

Background - Chris

Distributed systems - Chris

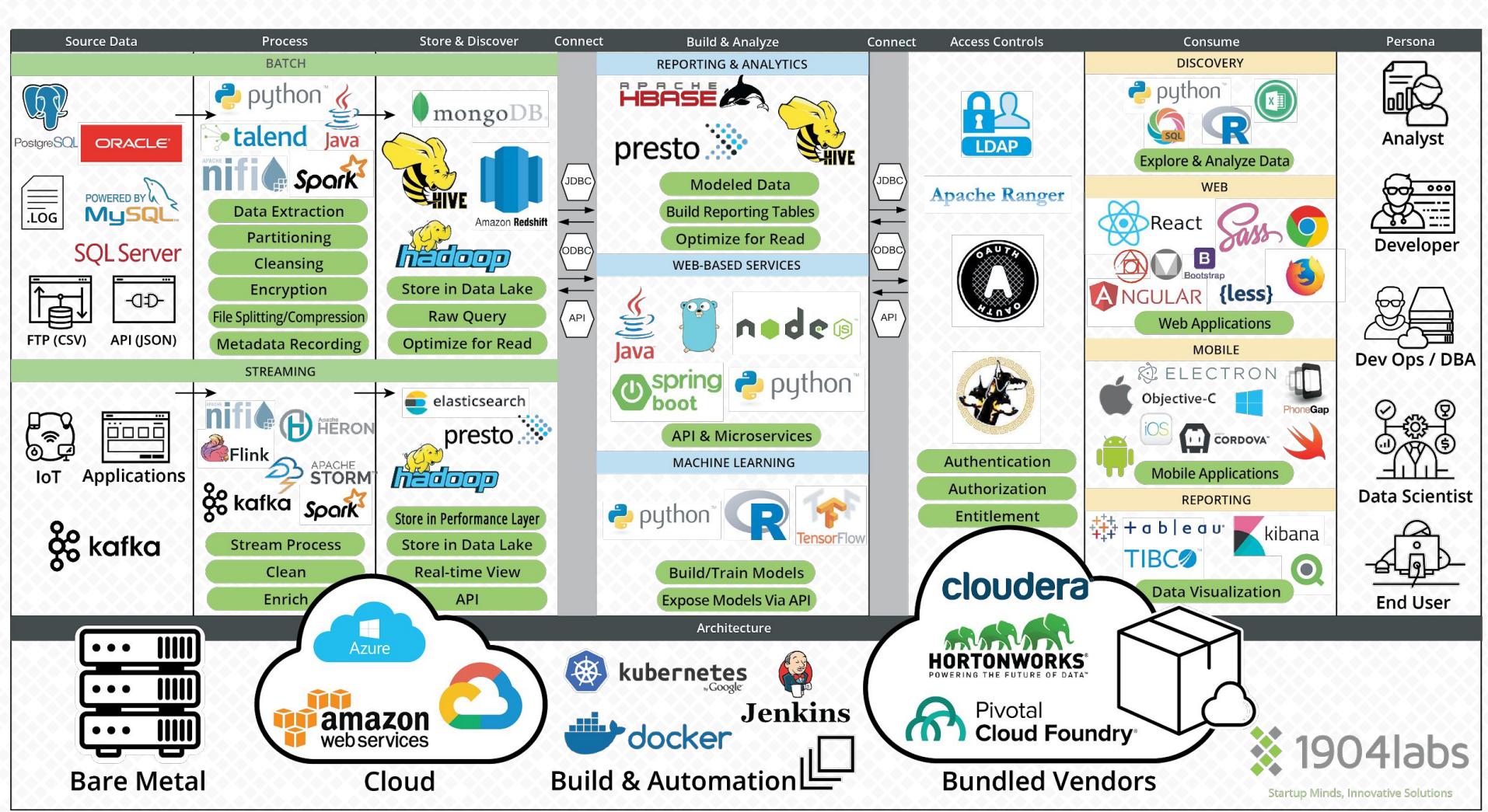
Day in the life - Tim

Demo - Kit



Why Data Engineering?

- It's not just about analytics.
- Treating data as an asset.
- Data science is built upon the work of good data engineering.
- The data we collect is changing. The way we store and access it must as well.
- Automate, automate, automate.

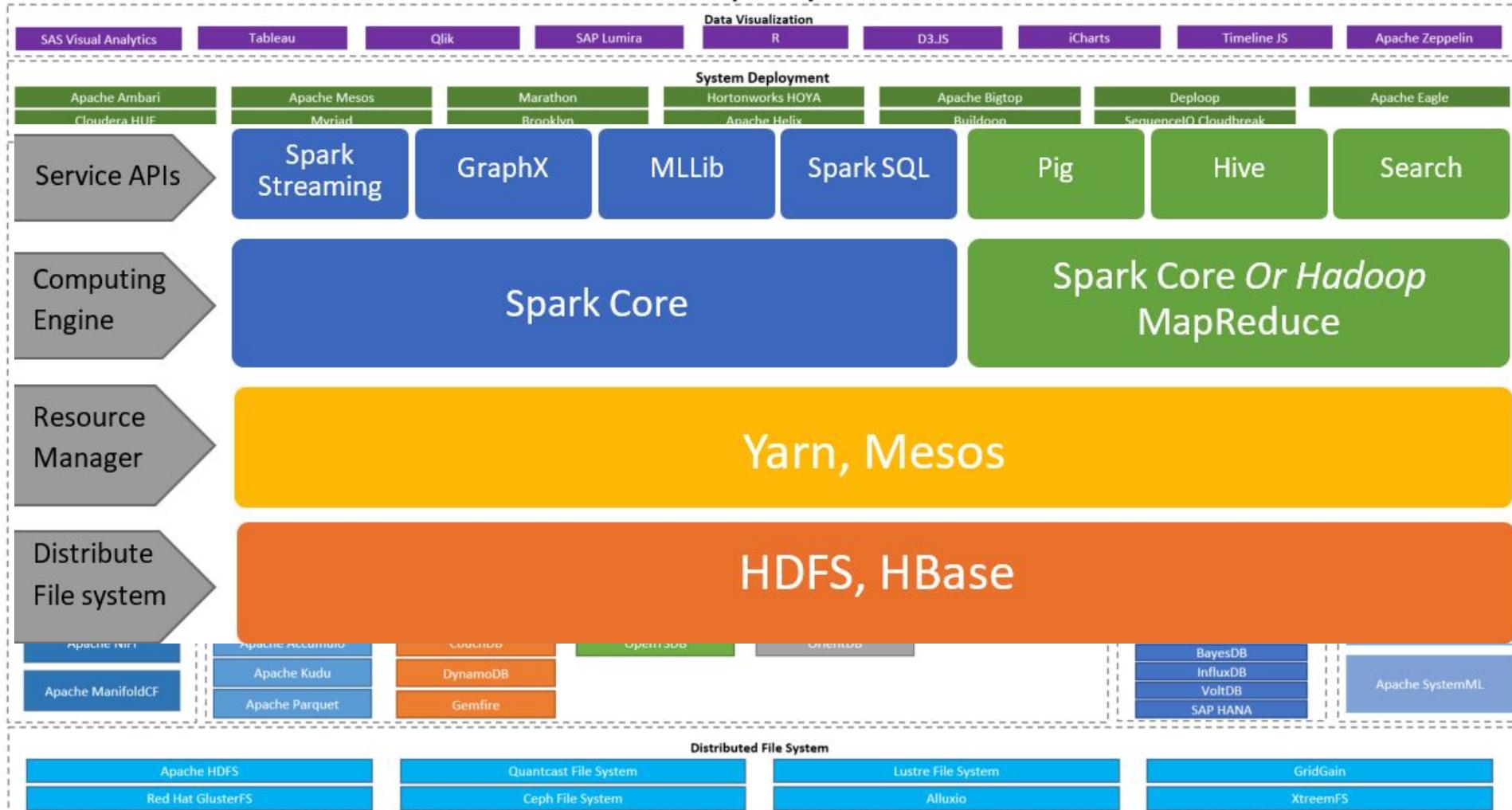


Distributed Systems

What does it mean?

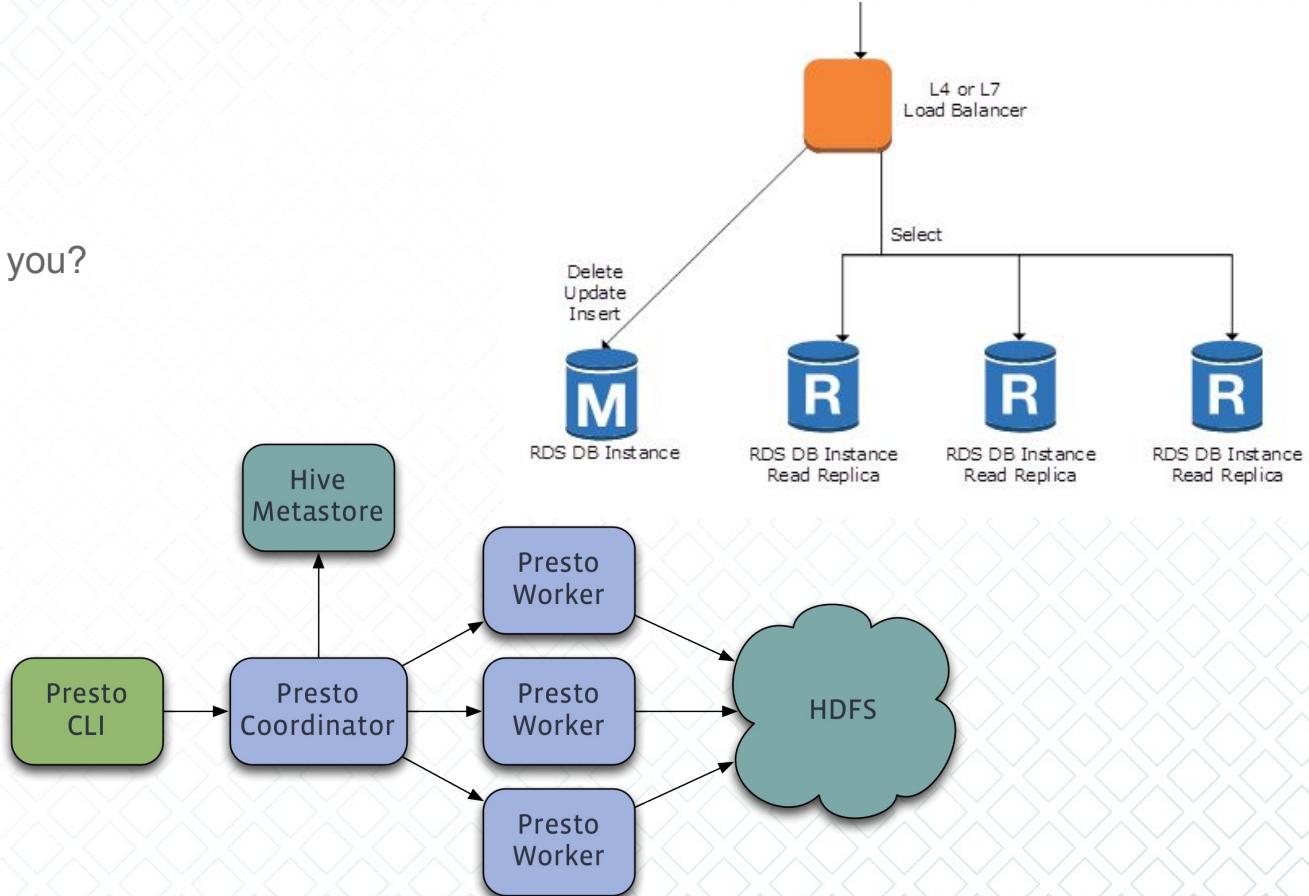
A distributed system is a system whose components are located on different networked computers, which communicate and coordinate their actions by passing messages to one another.

Hadoop Ecosystem



Distributed Compute

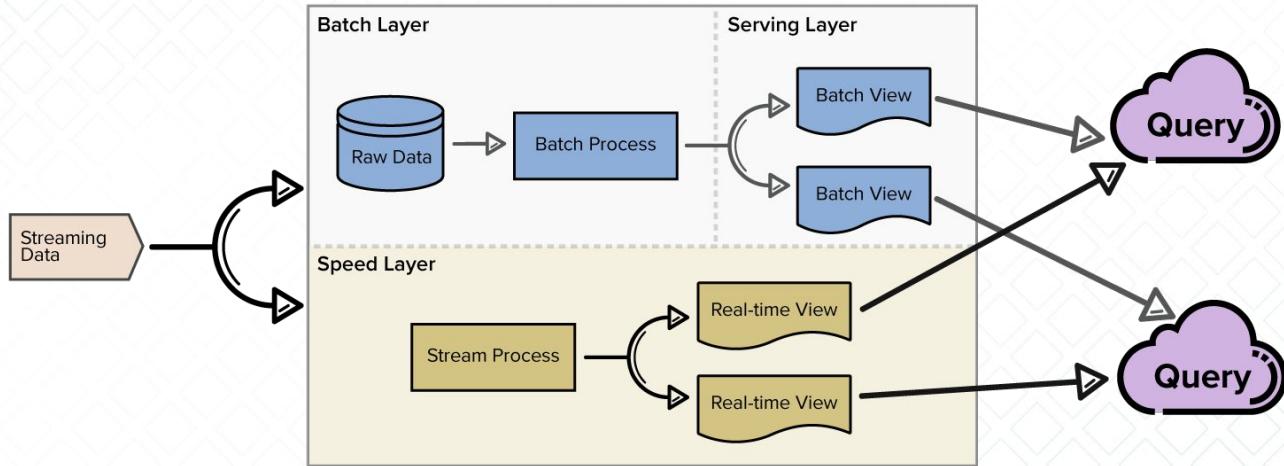
1. Overview
2. Scaling
3. What does it mean for you?



Data Architecture - Lambda & Kappa

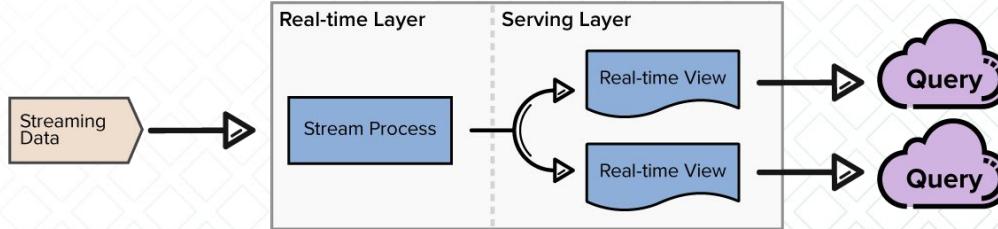
Lambda: Batch and Speed Layer

- Stores raw data
- Computes tables in batch
- Has streaming component
- Performs data modeling
- Exposes end tables to users



Kappa: Speed Layer

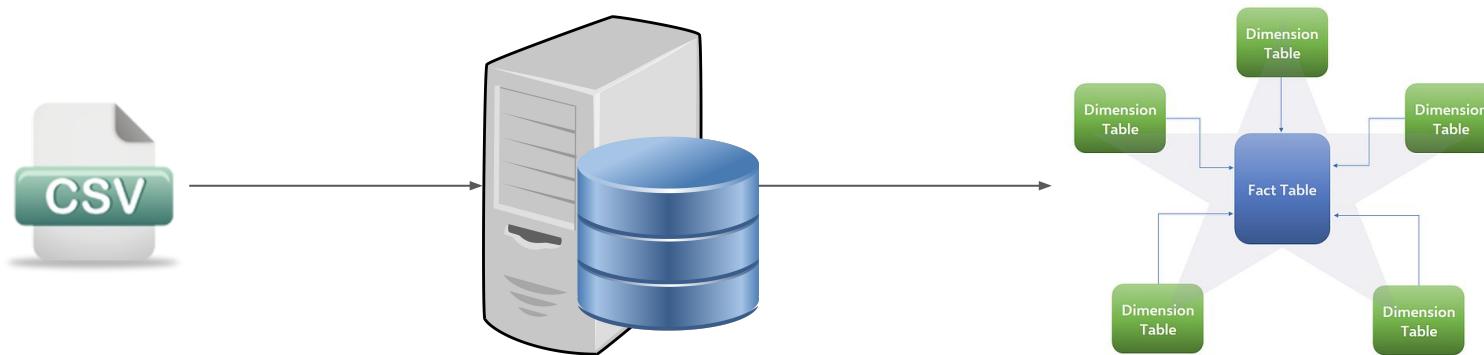
- All data as a stream
- In stream processing
- Real-time views



A Day In The Life

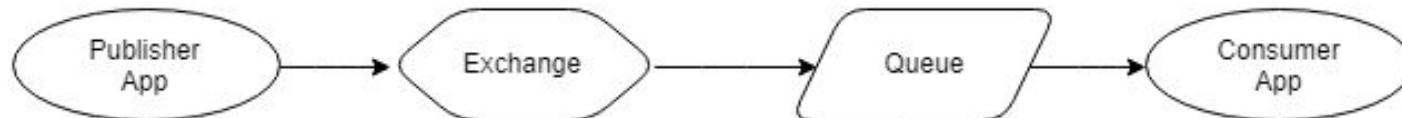
What's not modern about this architecture?

- Cron job runs once a day to check for presence of a file on an FTP server
- CSV file is downloaded and loaded into an RDBMS star schema



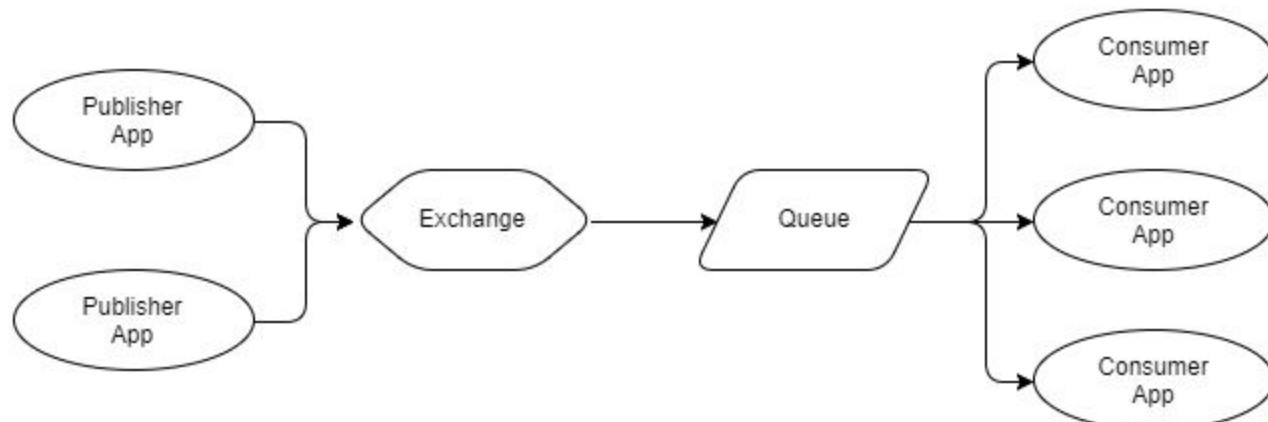
Message Queues

- Publishers
 - Exchanges
 - Queues
 - Consumers
- Messages are removed from queues once a consumer reads them



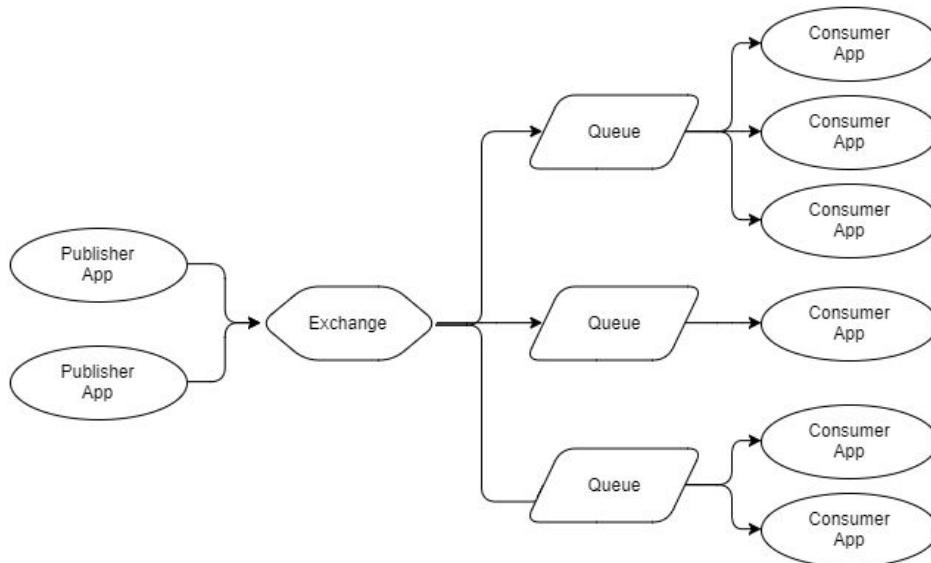
Scaling Up Message Queues

What if a single consumer can't keep up with the volume?



Scaling Them Up Even Further...

What happens if each message needs to go to more than one set of consumers?



Distributed Logs

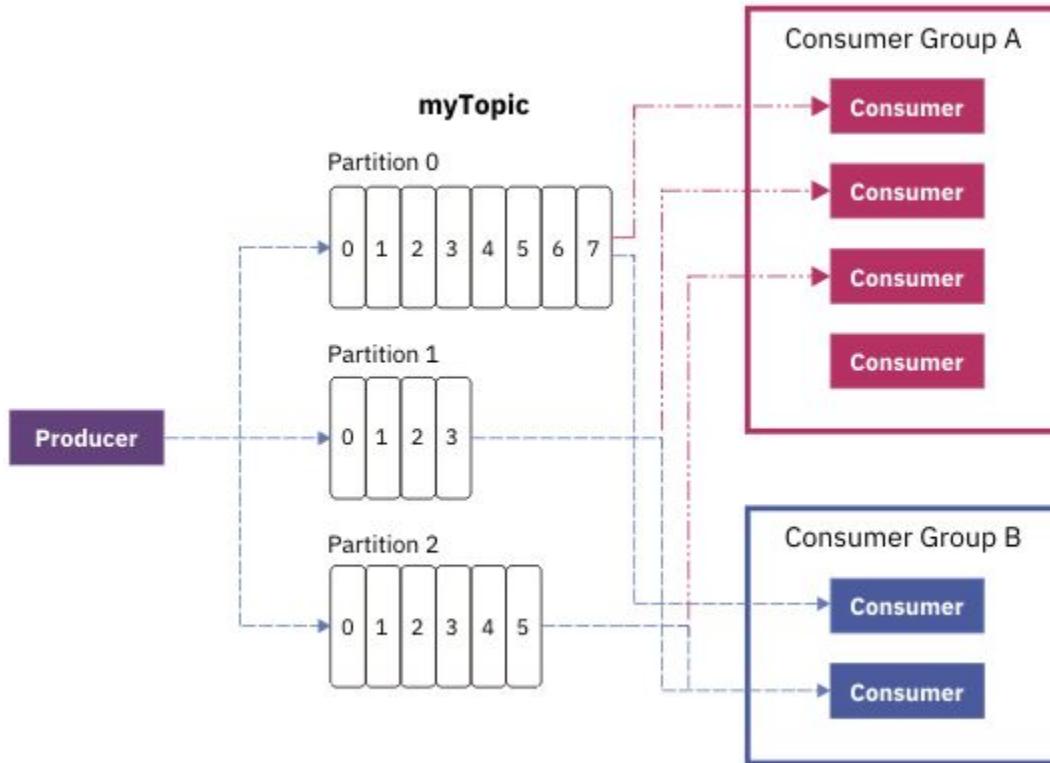
- Topic - a named stream of messages
- Partition - an ordered list of message. A topic is made up of a set of partitions
- Producer - a process that publishes messages to a topic
- Consumer - a process that consumes messages from a topic
- Consumer group - a named group of consumers who together consume messages from a set of topics

Topics can be consumed by multiple consumer groups independently.

Each consumer group is responsible for tracking its own position in message processing

Messages stay in the log for the defined retention period (are not removed once all consumer groups have processed them)

Distributed Logs



Data Formats



Apache Avro

- Row oriented
- Fast at writes
- Rich support for schema evolution

```
{  
    "type" : "record",  
    "namespace" : "STL Big Data",  
    "name" : "Inventory",  
    "fields" : [  
        { "name" : "date" , "type" : "date" },  
        { "name" : "price" , "type" : "int" },  
        { "name" : "size" , "type" : "int" }  
    ]  
}
```

<u>date</u>	<u>price</u>	<u>size</u>
2011-01-20	10.1	10
2011-01-21	10.3	20
2011-01-22	10.5	40
2011-01-23	10.4	5
2011-01-24	11.2	55
2011-01-25	11.4	66
...
2013-03-31	17.3	100

Apache Avro

- Row oriented
- Fast at writes
- Rich support for schema evolution

```
{  
    "type" : "record",  
    "namespace" : "STL Big Data",  
    "name" : "Inventory",  
    "fields" : [  
        { "name" : "date" , "type" : "date" },  
        { "name" : "price" , "type" : "int" },  
        { "name" : "size" , "type" : "int" },  
        { "name" : "qty", "type" : [ "null", "int" ], "default" : null }  
    ]  
}
```

<u>date</u>	<u>price</u>	<u>size</u>
2011-01-20	10.1	10
2011-01-21	10.3	20
2011-01-22	10.5	40
2011-01-23	10.4	5
2011-01-24	11.2	55
2011-01-25	11.4	66
...
2013-03-31	17.3	100

Column Oriented Storage

- Instead of storing all columns for a row together on disk, store each column together on disk
- Enables higher compression ratios
- What kind of queries will you be executing against the data?

<u>date</u>	<u>price</u>	<u>size</u>
2011-01-20	10.1	10
2011-01-21	10.3	20
2011-01-22	10.5	40
2011-01-23	10.4	5
2011-01-24	11.2	55
2011-01-25	11.4	66
...
2013-03-31	17.3	100

<u>date</u>	<u>price</u>	<u>size</u>
2011-01-20	10.1	10
2011-01-21	10.3	20
2011-01-22	10.5	40
2011-01-23	10.4	5
2011-01-24	11.2	55
2011-01-25	11.4	66
...
2013-03-31	17.3	100

Apache Parquet

- Column oriented
- Enables fast reading
- Compresses well
- Commonly used with Impala

Apache ORC

- Column oriented
- Enables fast reading
- Compresses extremely well
- Commonly used with Hive and Presto

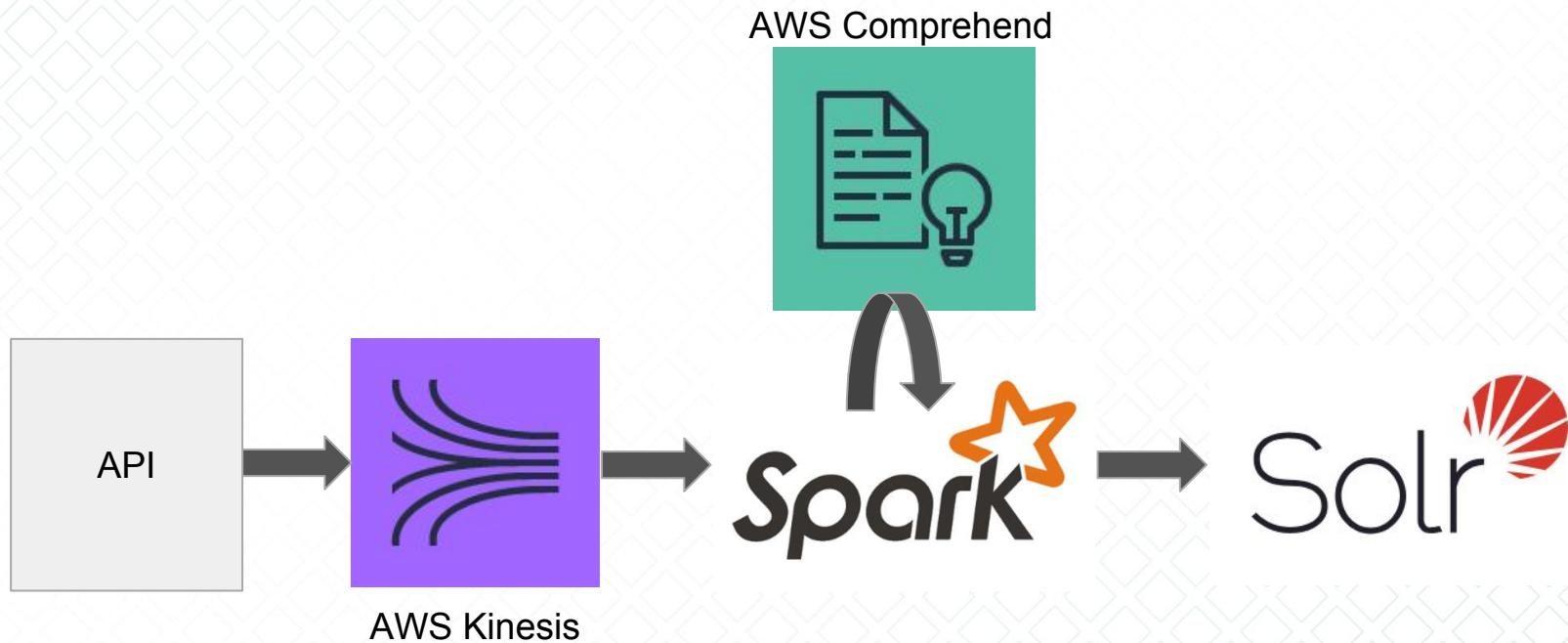
Compression Comparison

Format	Size
CSV	2.5 GB
Avro	1.2 GB
Parquet	817 MB
ORC	796 MB

Source: Netflix Prize Data <https://www.kaggle.com/netflix-inc/netflix-prize-data>

Demo

<https://github.com/kitmenke/spark-hello-world>



Appendix



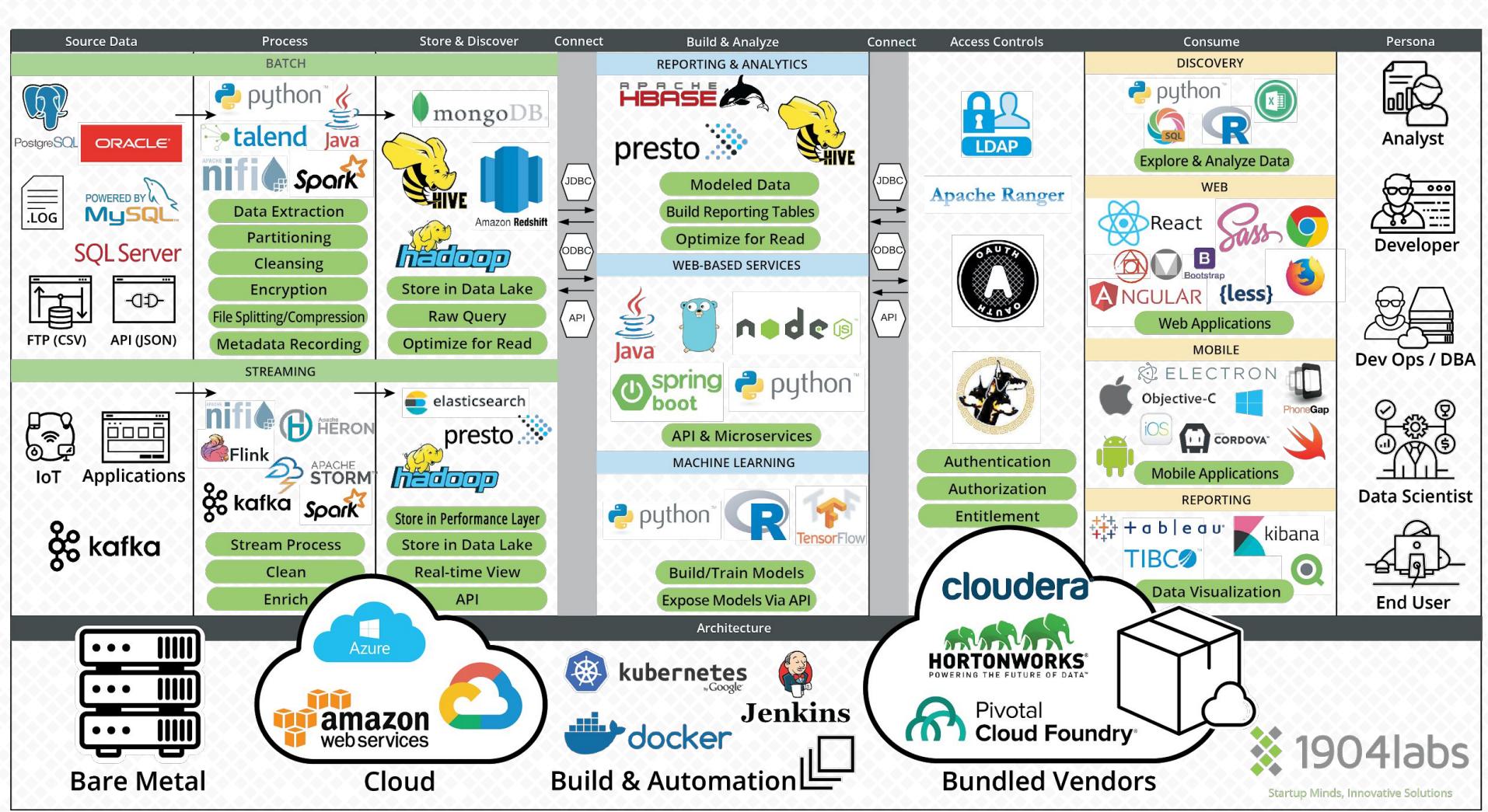
Building the Right Thing the Right Way

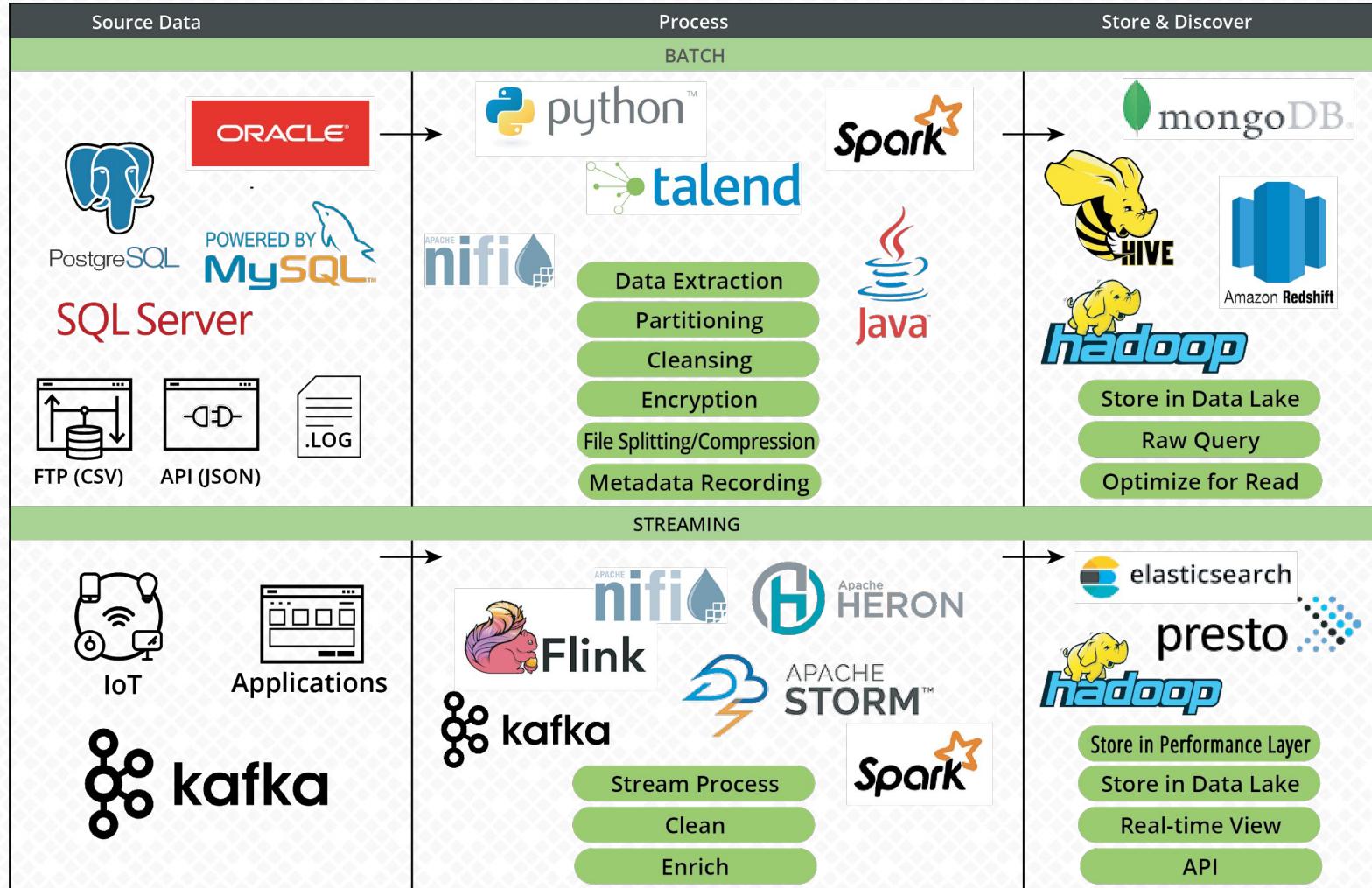
We listen to and understand our clients' most complex business challenges and bring an innovative approach to building modern technology solutions.



Best in class modern architecture

... At 1904labs, we have a passion for open source technologies, strive to build cloud first applications, and are motivated by our desire to transform businesses into data-driven enterprises.

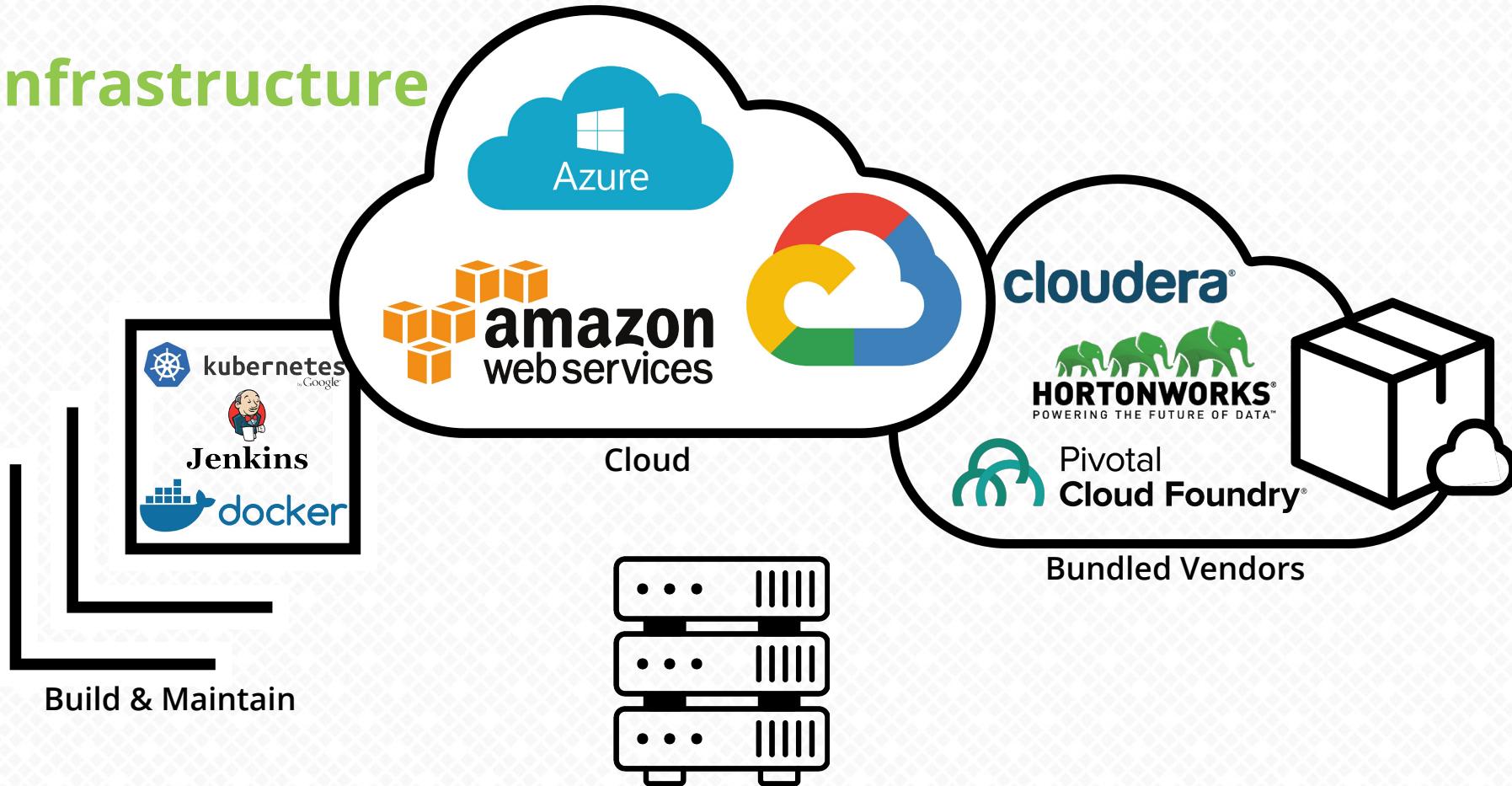








Infrastructure



Source Data - Data as it exists in your current company infrastructure

Process - Pulling source data while enriching, cleansing, and encrypting the data in transit.

Store & Discover - Inserting data into modern architecture. Raw data can be queried at this level.

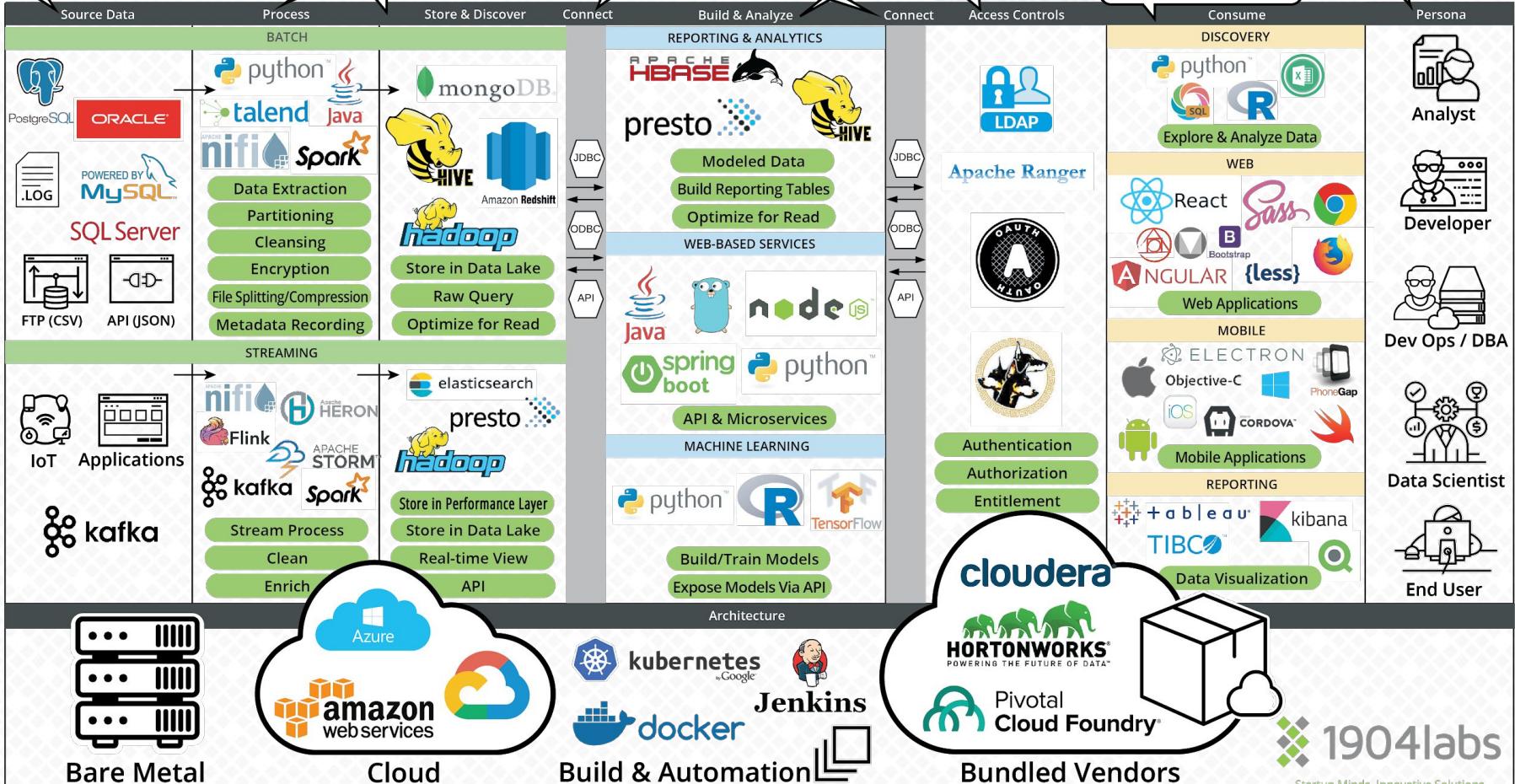
Build & Analyze - The storage layer, data models, and web services are built and consumed by reporting services and end users.

Connect - Both external users & internal components will typically connect through JDBC, ODBC, and API.

Access Controls - Standard security layer used for authentication & authorization of all services.

Consume - The applications which connect to the underlying modern architecture.

Persona - All people involved with the modern architecture. Some will be international developers, while others are external.



Open-Source Adoption in the Industry

These open source tools are already used broadly across industries today

1904labs' primary focus is to help our clients design, build and operate world class and modern data, development operations and decision science capabilities...

...Utilizing best of breed open source and/or commercial open source tools and platforms.

