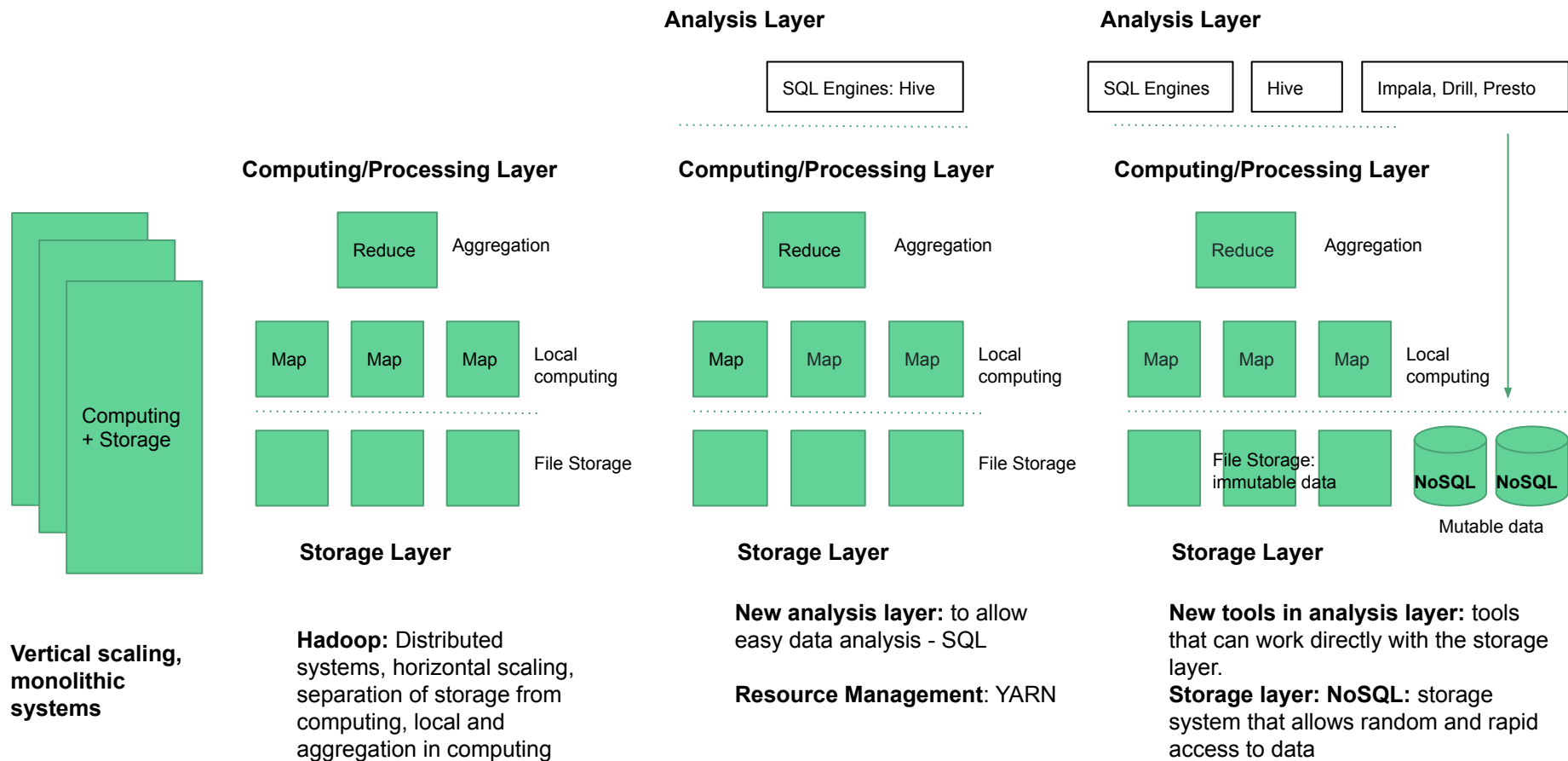


Intro Big Data Hadoop & Architecture Evolution

Evolution of the systems towards big data architecture



A new computing fw: Apache Spark

Analysis Layer

Hive

Spark SQL

Impala, Drill, ..

Computing/Processing Layer

Spark/
TEZ

Master

Spark/
TEZ

Spark/
TEZ

Spark/
TEZ

Local
Executors

File Storage

NoSQL

NoSQL

Storage Layer

Spark Core came as a replacement for old Map Reduce and to improve computing times. **But, Spark evolved into a framework that consists of Spark Core, SQL, Streaming, Graph, Machine Learning..**

Big Data Architecture

Visualization Layer (Zoomdata, Tableau, Zeppelin,..)

Analysis Layer

Hive

Spark SQL

Impala, Drill, ..

Computing/Processing Layer

Spark/TEZ

Master

Spark/
TEZ

Spark/
TEZ

Spark/
TEZ

Local
Executors

File Storage

NoSQL

NoSQL

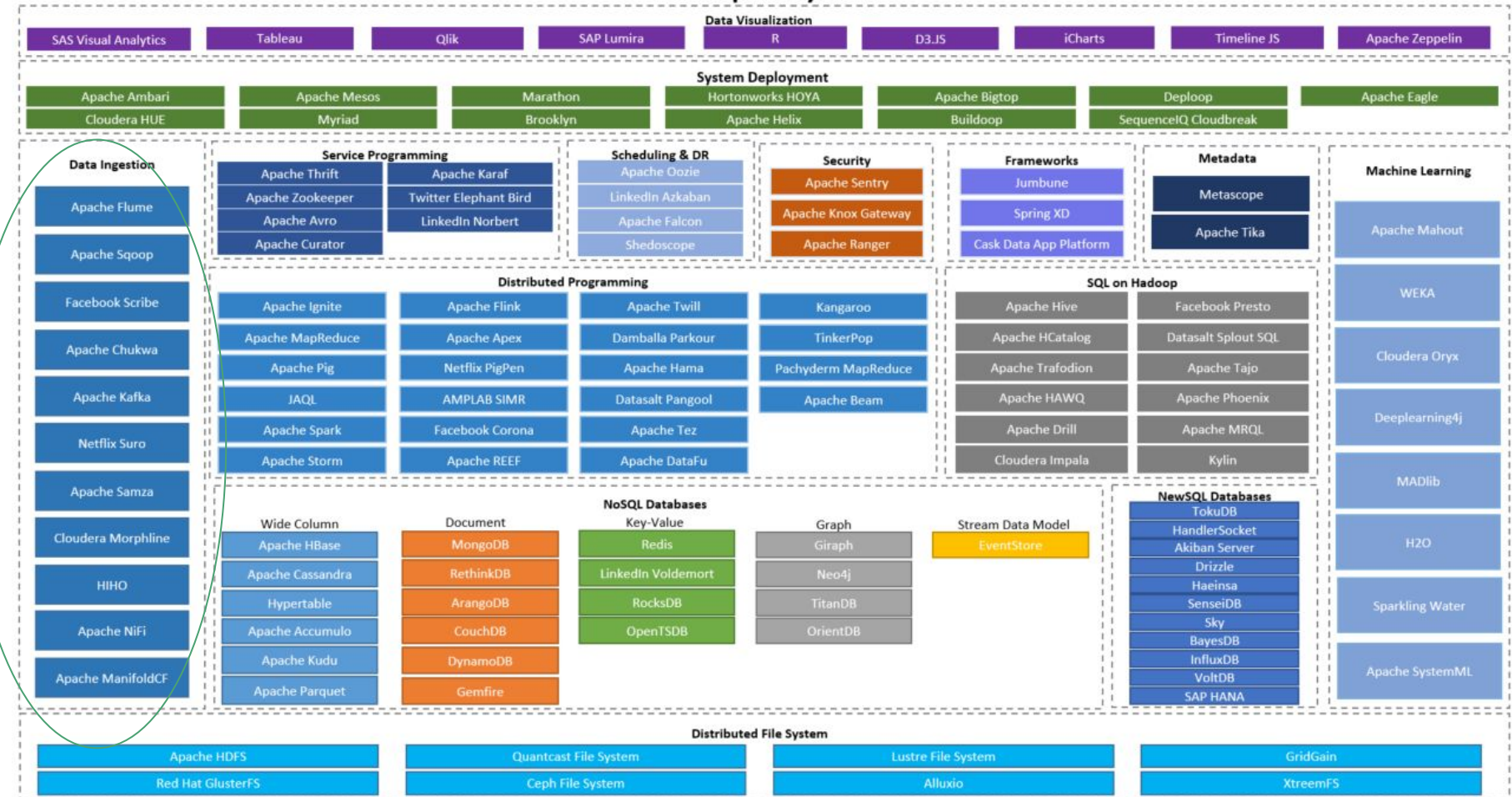
Storage Layer

Other layers: Ingestion, Visualization, Security, ...

**Ingestion
Layer
(Kafka,
Nifi, ...)**

**Resources
Mgmt
(Mesos)**

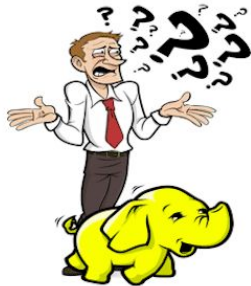
Hadoop Ecosystem



Ingestion layer



Hadoop Best Practices for Data Ingestion



Here are the 100 databases
where our data is scattered
in. Bring it all in the data
lake ASAP.



Data ingestion

- **Collection of data from different sources - Ingestion** (“some data is stored in Hadoop, some in Oracle, and a little bit is in Excel”)
 - IOT like case (loads of small sources of data): tool that has a mini client version
 - Specialized tools: like Sqoop for moving data between RDBMS and Hadoop
- **Data Flow management**
 - Optimal movement of data between different systems
- **Messaging bus/broker**
 - Reliable storage of ingestion data
 - Reprocessing of data

Solutions in data ingestion space

- **Sqoop -- HDFS focused tool**

- uses JDBC drivers to extract data from relational databases and then directly into the Hadoop stack, typically to Hive (or another metastore governed tool) or HDFS. Has options to schedule incremental data loads. Extracts are run in Mappers using Hadoop Map-Reduce and highly parallelized.

- **Flume -- HDFS focused tool**

- The original log collector for Hadoop. It is made up of sources (where is the data coming from), channels (what are you doing with the data once you've grabbed it), and sink (where is it going). It is effectively replaced nowadays by a combination of Kafka/Storm. One cool thing Flume does do is that it can bind to a source (such as a web stream) and push data all the way to HDFS without having to write any code, all in one tool. Also supports contextual routing.

- **Streamsets & Streamsets Data Collector**

- SDC was started by a California-based startup in 2014 as an open source ETL project available on GitHub. The first release was published in June 2015. Cloudera supported solution for dataflow management: builds, tests, runs and maintains dataflow pipelines connecting a variety of batch and streaming data sources and compute platforms.

Solutions in data ingestion space

- **Nifi**
 - Hortonworks driven dataflow management solution. NiFi can be used in mission-critical data flows with rigorous security & compliance requirements, where we can visualize the entire process and make changes immediately, in real-time.
- **Apache Spark:** Non specific to ingestion layer, its computing framework can be used for ingestion data aggregation
- **Kafka**
 - Distributed messaging bus
- **Others: Kinesis**

Flume, NiFi, Kafka

Ingestion Tool	Out-of-the-box	Limits	Use Cases
Flume	<ul style="list-style-type: none">- Configuration-based- Sources, channels & sinks- Interceptors	<ul style="list-style-type: none">- Data loss scenarios when not using Kafka Channel- Data size (KB)- No data replication	<ul style="list-style-type: none">- collecting, aggregating, and moving high-volume streaming events into Hadoop.
Kafka	<ul style="list-style-type: none">- Back-pressure- Reliable stream data storage- Kafka Streams- Sources/sinks with Kafka-Connect	<ul style="list-style-type: none">- Custom coding often needed- Data size (KB)- Fixed protocol/format/schema	<ul style="list-style-type: none">- Streaming data- Messaging- Systems integration- Commit log
NiFi	<ul style="list-style-type: none">- Configuration-based UI- Many drag & drop processors- Back-pressure- Prioritized queuing- Data provenance- Flow templates <p>can handle messages with arbitrary sizes</p>	<ul style="list-style-type: none">- Not for CEP or windowed computations- No data replication	<ul style="list-style-type: none">- Dataflow management with visual control- Data routing between disparate systems- Arbitrary data size

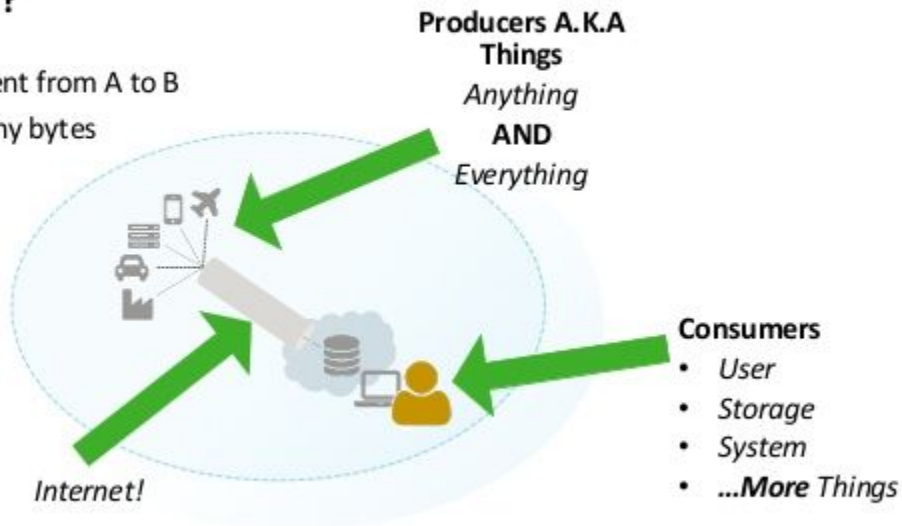
Streamsets



Data Flow Management

What is dataflow?

- Moving some content from A to B
- Content could be any bytes
 - Logs
 - HTTP
 - XML
 - CSV
 - Images
 - Video
 - Telemetry



Apache NiFi

NiFi was donated by the NSA to the Apache Foundation in 2014 and current development and support is provided mostly by Hortonworks.

Used mainly for

- Reliable and secure transfer of data between systems
- Delivery of data from sources to analytics platforms
- Enrichment and preparation of data:
 - Conversion between formats
 - extraction/parsing
 - Routing decision

Note: Nifi doesnt have a built-in fault-tolerant mechanism.
Meaning if one node goes down then there has to be a manual intervention to prevent data loss/redirect flow etc

What is not used for

- Distributed computing
- Complex event processing, Joins/Complex Rolling Windows operations
- Data reprocessing - e.g. Kafka like

User Interface

Less of this... ... more of this

