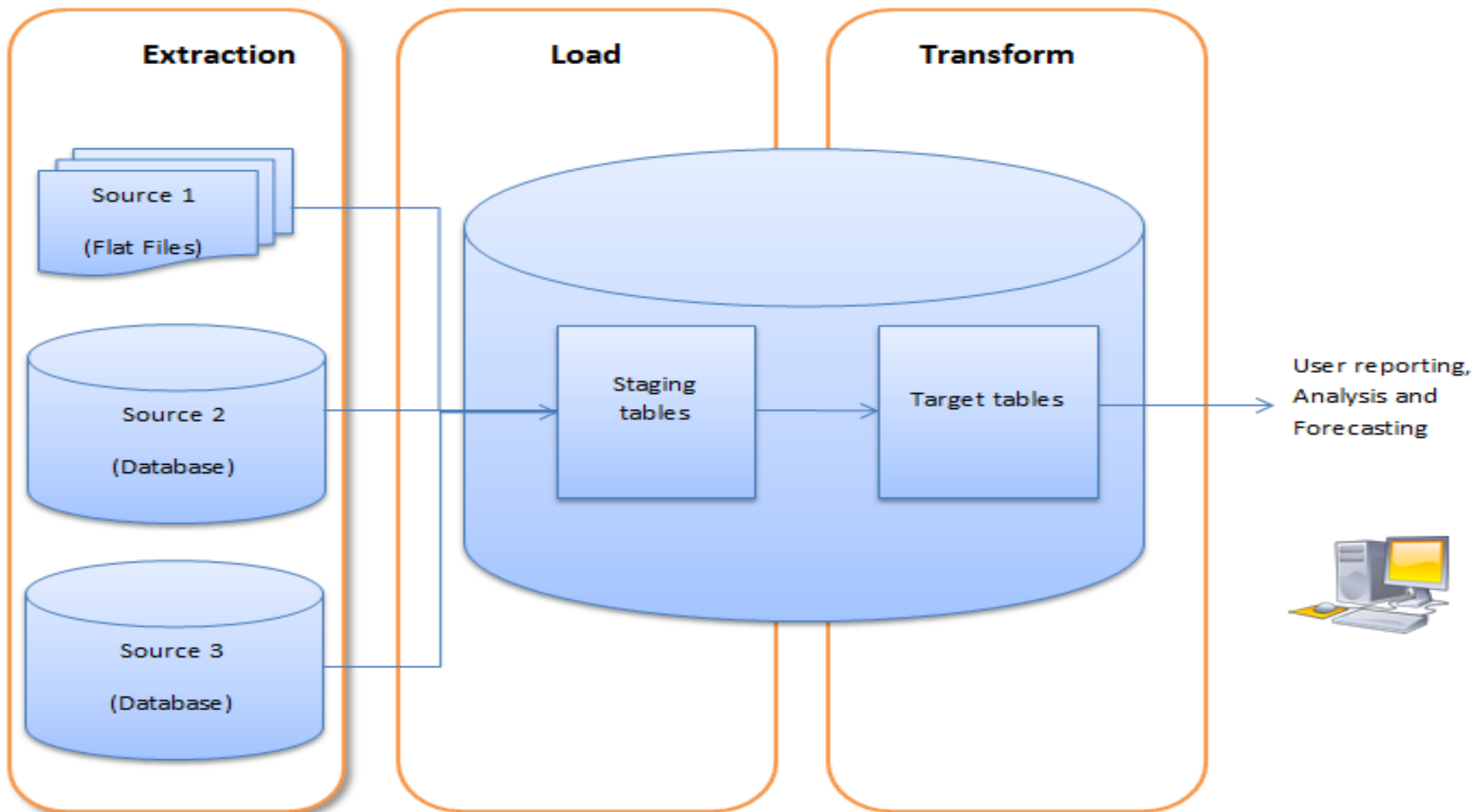


The ETL Process

Extract, Transform, and Load operations
for data warehouses

ETL (extract, transform, load)

ELT based Data Warehouse Architecture diagram



THE ETL PROCESS

- Capture/Extract
- Scrub or data cleansing
- Transform
- Load and Index

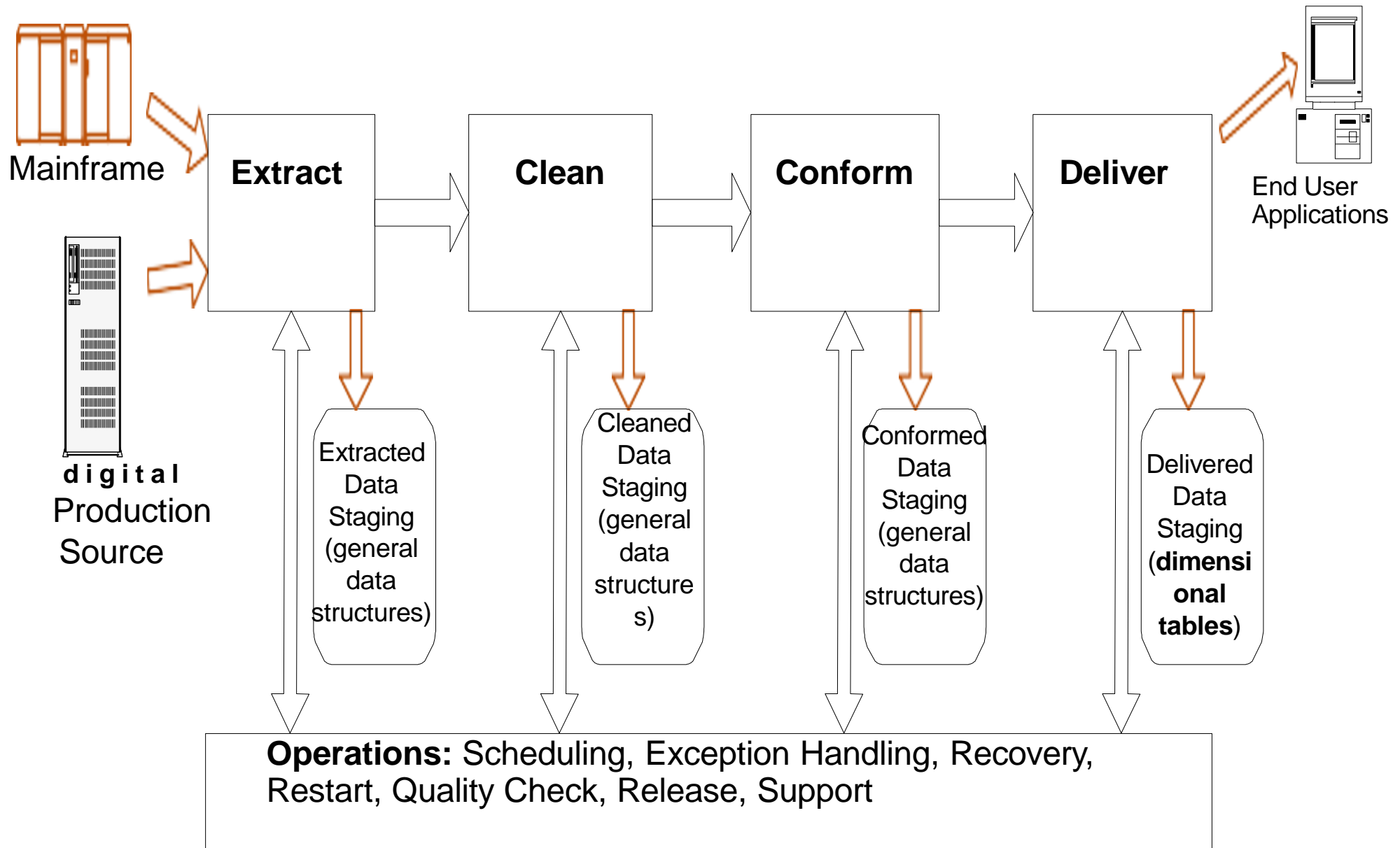
Pre ETL Tasks

- Profiling the data:
 - Know your data
 - Will give direct insight in the data quality of the source systems.
 - Find out how many rows have missing or invalid values, or what is the distribution of values in a specific column.
- This knowledge will help to specify rules (data standards / quality checks) in order to cleanse the data and to keep really bad data out of the repository.
- Doing data profiling before designing the ETL process helps you to better design a system that is correct, robust and has a clear structure.

ETL

- Data moved from source to target databases
- Focus: Preparing the data for reporting / analysis
- ETL = Extract -> Transform -> Load
 - Extract: Get the data from source(s) as efficiently as possible
 - Transform: Perform calculation / map data / clean data
 - Load: Load data to the target storage

The Four Staging Steps of a Data Warehouse



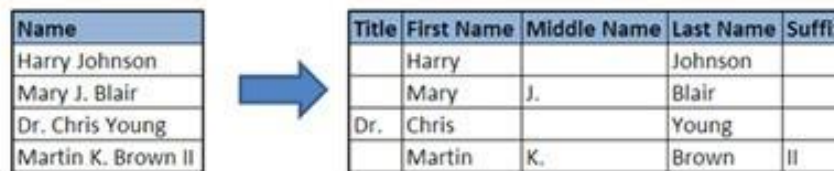
Dirty Data

- Absence of Data / Missing Data
- Multipurpose Fields
- Cryptic Data
- Contradicting Data
- Violation of Data Rules
- Reused Primary Keys
- Non-Unique Identifiers
- Data Integration Problems

Data Cleaning: Parsing / Combining

- Parsing locates and identifies individual data elements in the source files and then *isolates* these data elements in the target files.

- Examples include full name stored in one column



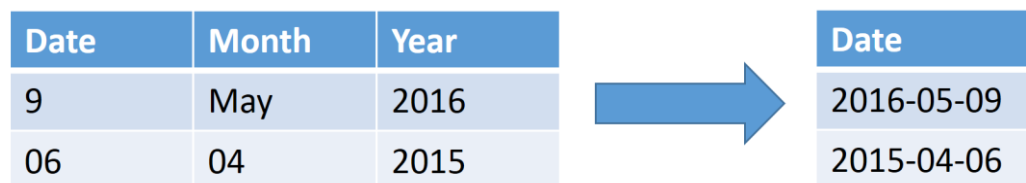
The diagram illustrates the parsing process. On the left, a table with a single column 'Name' contains four entries: 'Harry Johnson', 'Mary J. Blair', 'Dr. Chris Young', and 'Martin K. Brown II'. A blue arrow points to the right, where a new table structure is shown. This new table has five columns: 'Title', 'First Name', 'Middle Name', 'Last Name', and 'Suffix'. The data from the original 'Name' column is distributed across these columns: 'Harry Johnson' becomes 'Harry' (First Name) and 'Johnson' (Last Name); 'Mary J. Blair' becomes 'Mary' (First Name), 'J.' (Middle Name), and 'Blair' (Last Name); 'Dr. Chris Young' becomes 'Dr.' (Title), 'Chris' (First Name), and 'Young' (Last Name); and 'Martin K. Brown II' becomes 'Martin' (First Name), 'K.' (Middle Name), 'Brown' (Last Name), and 'II' (Suffix).

Name
Harry Johnson
Mary J. Blair
Dr. Chris Young
Martin K. Brown II

Title	First Name	Middle Name	Last Name	Suffix
	Harry		Johnson	
	Mary	J.	Blair	
Dr.	Chris		Young	
	Martin	K.	Brown	II

- Combining locates and identifies individual data elements in the source files and then *combines* these data elements in the target files

- Examples include event date stored in different columns as date, month and year




The diagram illustrates the combining process. On the left, a table has three columns: 'Date', 'Month', and 'Year'. The first row contains '9', 'May', and '2016'. The second row contains '06', '04', and '2015'. A blue arrow points to the right, where a new table structure is shown. This new table has a single column 'Date'. The data from the original columns is combined into a single date string: '2016-05-09' for the first row and '2015-04-06' for the second row.

Date	Month	Year
9	May	2016
06	04	2015

Date
2016-05-09
2015-04-06

Data Cleaning: Standardizing

- Standardizing applies conversion routines to transform data into its preferred (and consistent) format using both standard and custom data rules.



Sex	Sex
Male	1
F	2
M	1
Female	2
NULL	0
f	2
m	1

Data Cleaning: Matching

- Searching and matching records within and across the parsed, combined, corrected and standardized data based on predefined data rules to eliminate duplications, sequences

Pregnancy Number	Outcome Date	Outcome
1	2011-06-07	Twin

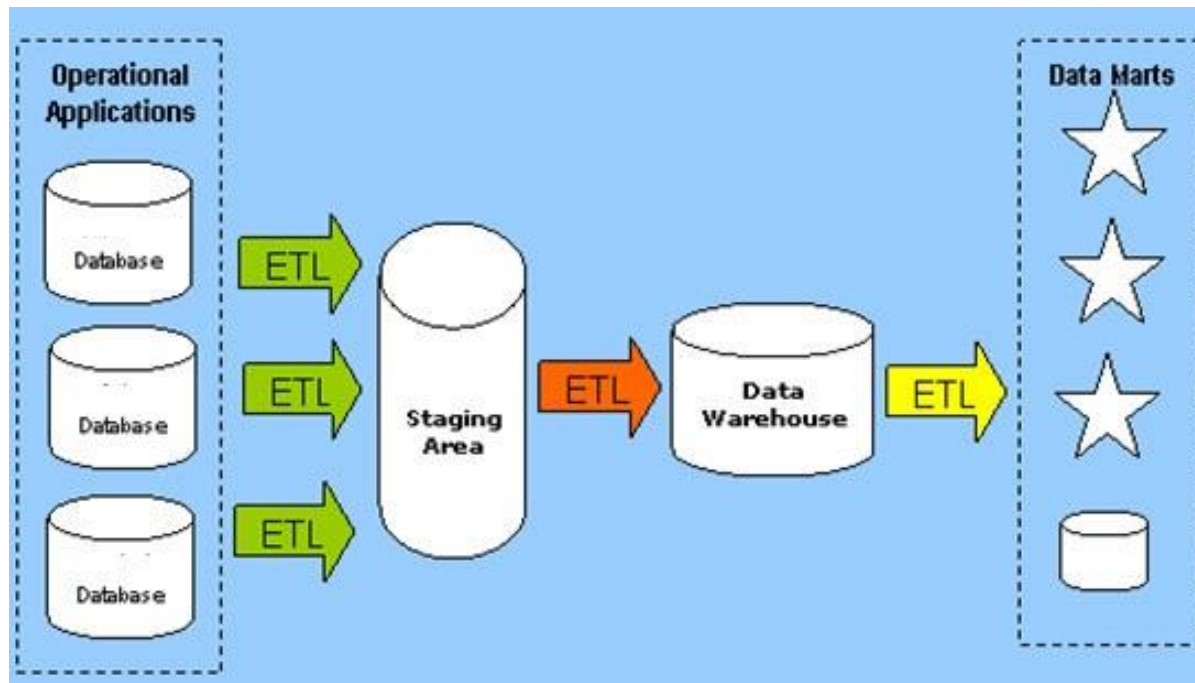
DoB	Name
2011-06-07	Child1
2011-06-07	Child2



Pregnancy Number	Birth Date	Name
1	2011-06-07	Child1
1	2011-06-07	Child2

Data Staging

- Used as an interim step between data extraction and later steps
- Accumulates data from asynchronous sources using native interfaces, flat files, FTP sessions, or other processes
- Data in the staging file is transformed and loaded to the warehouse
- There is no end user access to the staging file
- An operational data store may be used for data staging

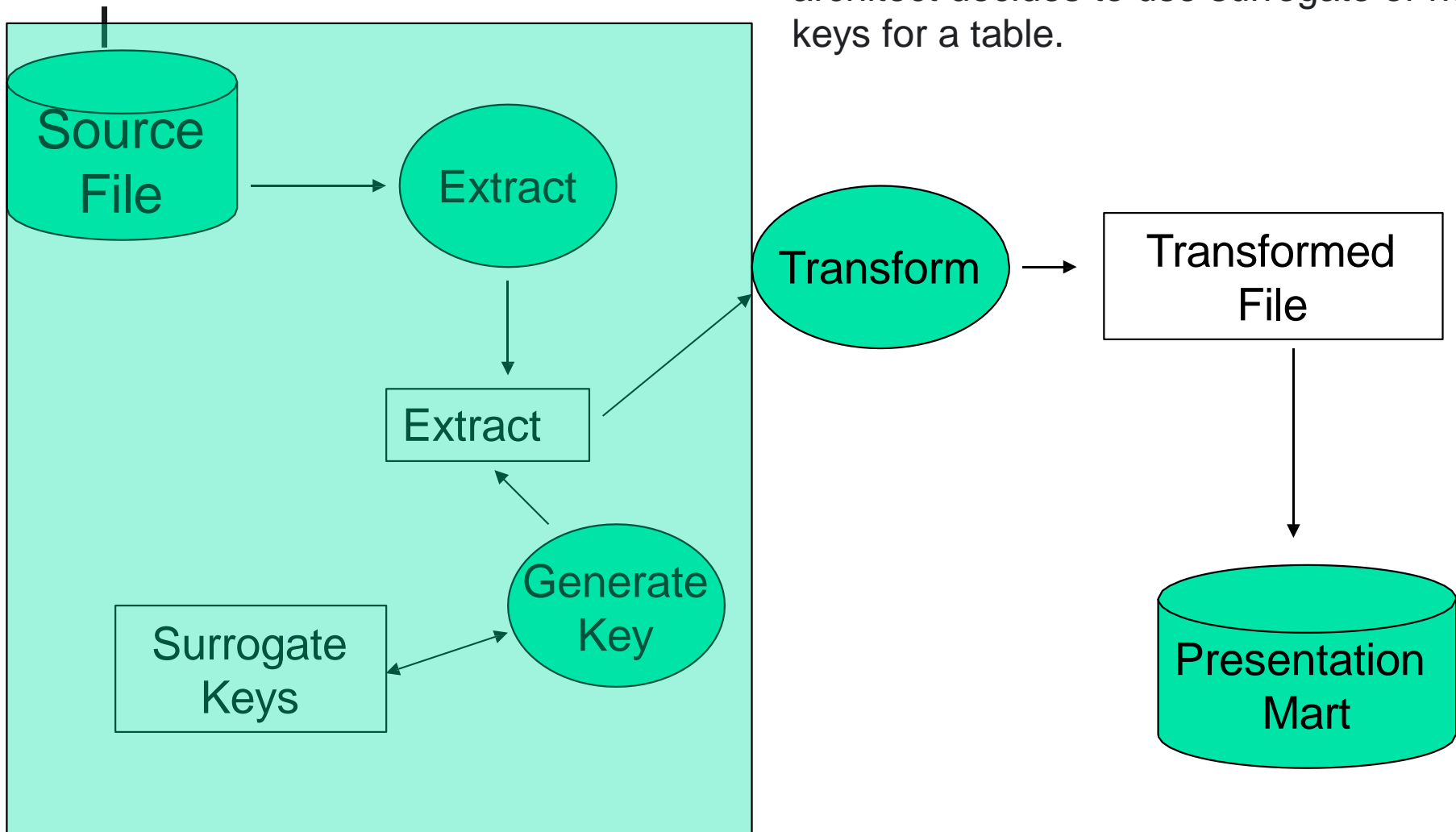


ETL Tools

- Few popular commercial and freeware(open-sources) ETL Tools
- Commercial ETL Tools:
 - IBM InfoSphere DataStage
 - Informatica PowerCenter
 - Oracle Warehouse Builder (OWB)
 - Oracle Data Integrator (ODI)
 - SAS ETL Studio
 - Business Objects Data Integrator(BODI)
 - Microsoft SQL Server Integration Services(SSIS)
 - Ab Initio
- Freeware, open source ETL tools:
 - Pentaho Data Integration (PDI)
 - Talend Integrator Suite
 - CloverETL
 - Jasper ETL

The Flow

A surrogate key is a **unique key for an entity in the client's business or for an object in the database**. Sometimes natural keys cannot be used to create a unique primary key of the table. This is when the data modeler or architect decides to use surrogate or helping keys for a table.



Add Surrogate Keys

in case we do not have a natural primary key in a table, then we need to artificially create one in order to uniquely identify a row in the table , this key is called the surrogate key or synthetic primary key of the table.

Example of Surrogate Key

registration_no	name	percentage
-----------------	------	------------

210101	Harry	90
--------	-------	----

210102	Maxwell	65
--------	---------	----

210103	Lee	87
--------	-----	----

210104	Chris	76
--------	-------	----



MERGE

surr_no	registration_no	name	percentage
---------	-----------------	------	------------

1	210101	Harry	90
---	--------	-------	----

2	210102	Maxwell	65
---	--------	---------	----

3	210103	Lee	87
---	--------	-----	----

4	210104	Chris	76
---	--------	-------	----

5	CS107	Taylor	49
---	-------	--------	----

6	CS108	Simon	86
---	-------	-------	----

7	CS109	Sam	96
---	-------	-----	----

8	CS110	Andy	58
---	-------	------	----

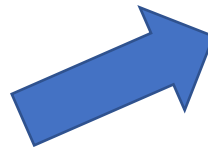
registration_no	name	percentage
-----------------	------	------------

CS107	Taylor	49
-------	--------	----

CS108	Simon	86
-------	-------	----

CS109	Sam	96
-------	-----	----

CS110	Andy	58
-------	------	----

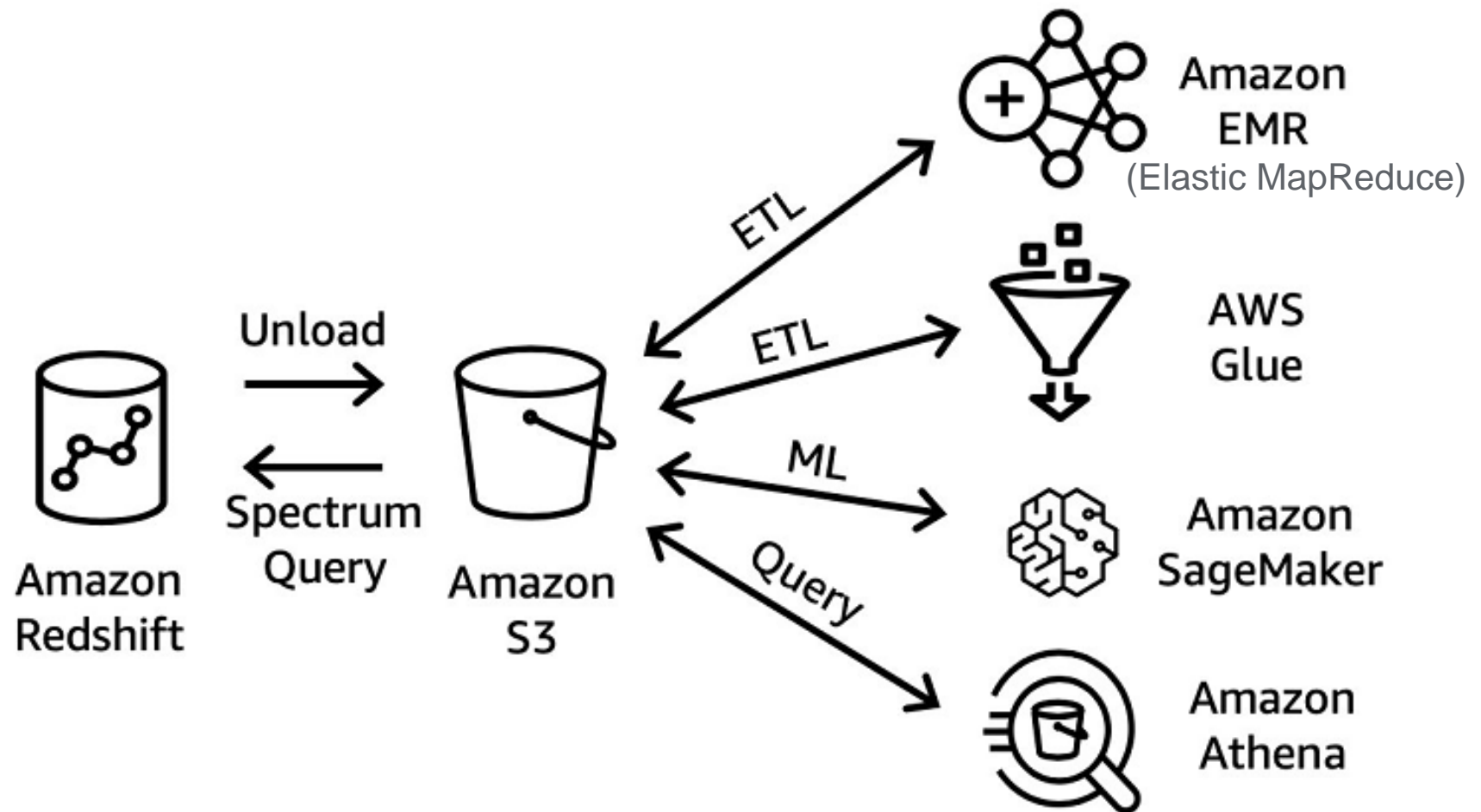


IDENTITY keyword to perform an auto-increment feature.

Create a Surrogate Key Table

```
create table surrogate (  
    AddressID int identity(1,1),  
    NOT NULL,  
    ...,  
    PRIMARY KEY (AddressID)  
);
```


Amazon ETL



ETL on Redshift



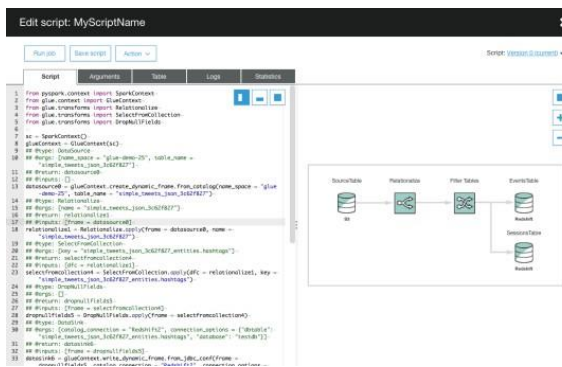
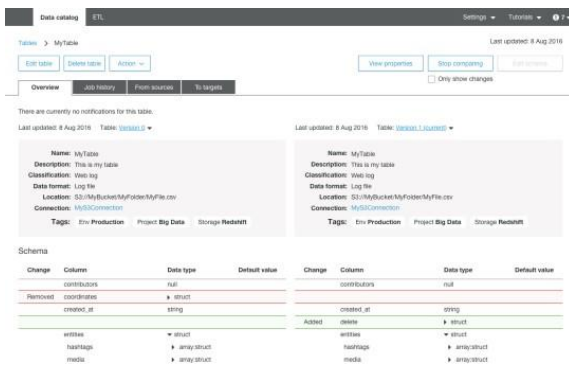
AWS Glue

Fully-managed data catalog and ETL service

Integrated with:

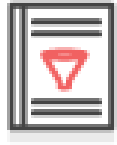


S3, RDS, Redshift & any
JDBC-compliant data store



Amazon Web Services, Inc.

AWS Glue: Components



Data Catalog

- Hive Metastore compatible with enhanced functionality
- Crawlers automatically extracts metadata and creates tables
- Integrated with Amazon Athena, Amazon Redshift Spectrum



Job Authoring

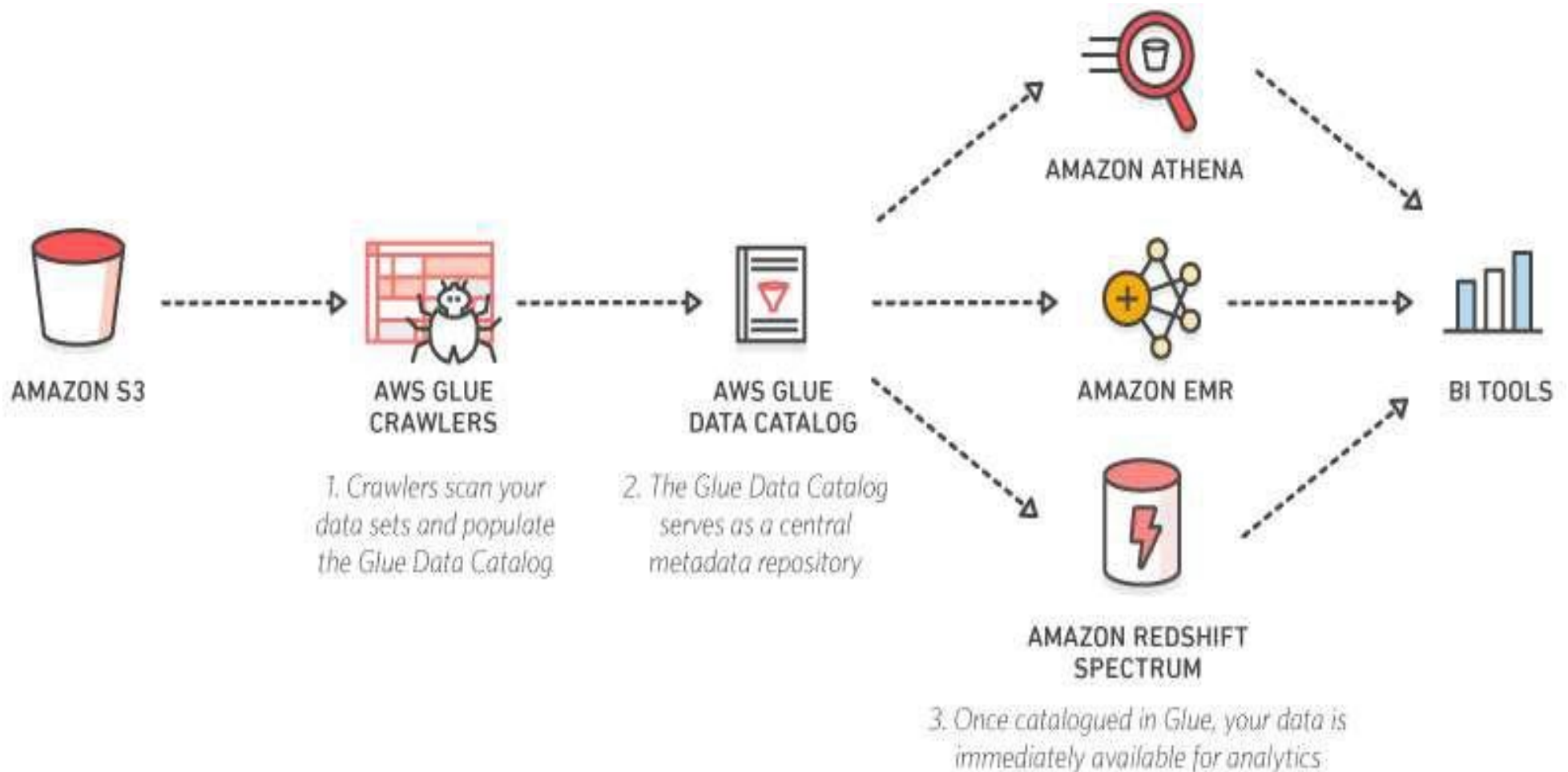
- Auto-generates ETL code
- Build on open frameworks – Python and Spark
- Developer-centric – editing, debugging, sharing



Job Execution

- Run jobs on a serverless Spark platform
- Provides flexible scheduling
- Handles dependency resolution, monitoring and alerting

Instantly query your data lake on Amazon S3



Data Catalog: Crawlers

Crawlers automatically build your Data Catalog and keep it in sync



Automatically discover new data, extracts schema definitions

- Detect schema changes and version tables
- Detect Hive style partitions on Amazon S3



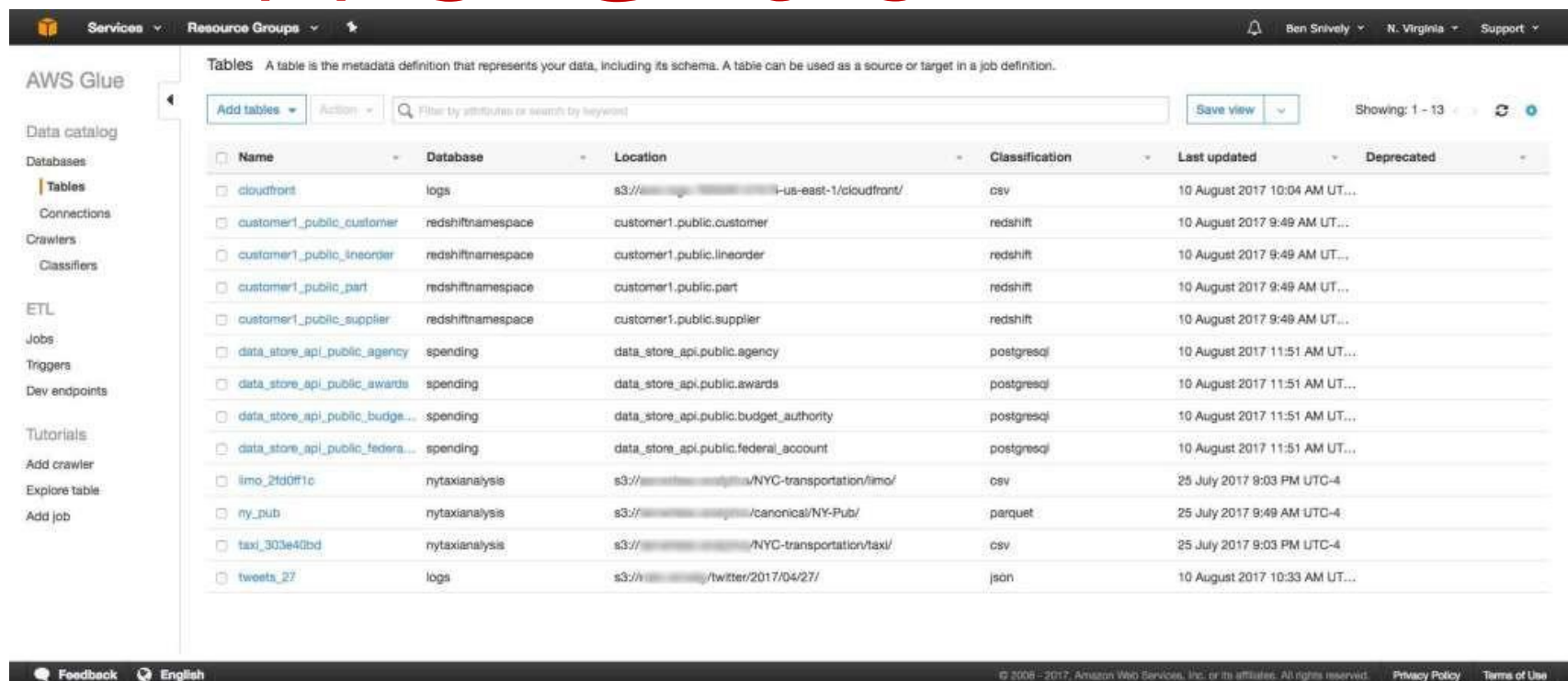
Built-in classifiers for popular types; custom classifiers using Grok expressions



Run ad hoc or on a schedule; serverless – only pay when crawler runs

Bring in metadata from a variety of data sources (Amazon S3, Amazon Redshift, etc.) into a single categorized list that is searchable

AWS Glue

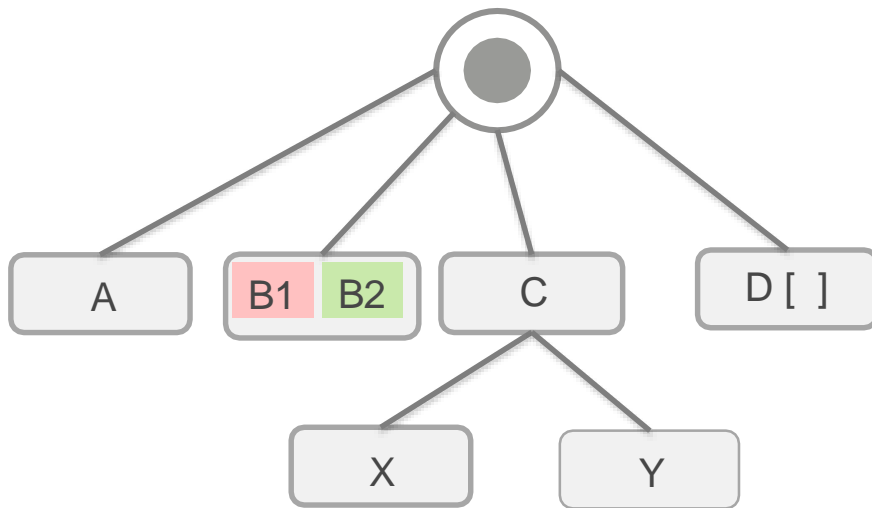


The screenshot displays the AWS Glue console interface. On the left, a navigation sidebar lists various components: Data catalog, Databases, Tables (selected), Connections, Crawlers, Classifiers, ETL, Jobs, Triggers, Dev endpoints, Tutorials, Add crawler, Explore table, and Add job. The main content area is titled 'Tables' and includes a descriptive text: 'A table is the metadata definition that represents your data, including its schema. A table can be used as a source or target in a job definition.' Below this, there are controls for 'Add tables', 'Action', a search filter, and a 'Save view' button. A table lists 13 tables with columns for Name, Database, Location, Classification, Last updated, and Deprecated. The tables include 'cloudfront', 'customer1_public_customer', 'customer1_public_lineorder', 'customer1_public_part', 'customer1_public_supplier', 'data_store_api_public_agency', 'data_store_api_public_awards', 'data_store_api_public_budge...', 'data_store_api_public_federa...', 'lmo_2td0ff1c', 'ny_pub', 'taxi_303e40bd', and 'tweets_27'.

Name	Database	Location	Classification	Last updated	Deprecated
cloudfront	logs	s3://[redacted]-us-east-1/cloudfront/	csv	10 August 2017 10:04 AM UT...	
customer1_public_customer	redshiftnamespace	customer1.public.customer	redshift	10 August 2017 9:49 AM UT...	
customer1_public_lineorder	redshiftnamespace	customer1.public.lineorder	redshift	10 August 2017 9:49 AM UT...	
customer1_public_part	redshiftnamespace	customer1.public.part	redshift	10 August 2017 9:49 AM UT...	
customer1_public_supplier	redshiftnamespace	customer1.public.supplier	redshift	10 August 2017 9:49 AM UT...	
data_store_api_public_agency	spending	data_store_api.public.agency	postgresql	10 August 2017 11:51 AM UT...	
data_store_api_public_awards	spending	data_store_api.public.awards	postgresql	10 August 2017 11:51 AM UT...	
data_store_api_public_budge...	spending	data_store_api.public.budget_authority	postgresql	10 August 2017 11:51 AM UT...	
data_store_api_public_federa...	spending	data_store_api.public.federal_account	postgresql	10 August 2017 11:51 AM UT...	
lmo_2td0ff1c	nytaxianalysis	s3://[redacted]/NYC-transportation/lmo/	csv	25 July 2017 9:03 PM UTC-4	
ny_pub	nytaxianalysis	s3://[redacted]/canonical/NY-Pub/	parquet	25 July 2017 9:49 AM UTC-4	
taxi_303e40bd	nytaxianalysis	s3://[redacted]/NYC-transportation/taxi/	csv	25 July 2017 9:03 PM UTC-4	
tweets_27	logs	s3://[redacted]/twitter/2017/04/27/	json	10 August 2017 10:33 AM UT...	

Job Authoring: Glue Dynamic Frames

Dynamic frame schema



Like Spark's Data Frames, but better for:

- Cleaning and (re)-structuring **semi-structured** data sets, e.g. JSON, Avro, Apache logs ...

No upfront schema needed:

- Infers schema on-the-fly, enabling transformations in a **single pass**

Easy to handle the unexpected:

- Tracks new fields, and inconsistent changing data types with **choices**, e.g. integer or string
- Automatically mark and separate error records

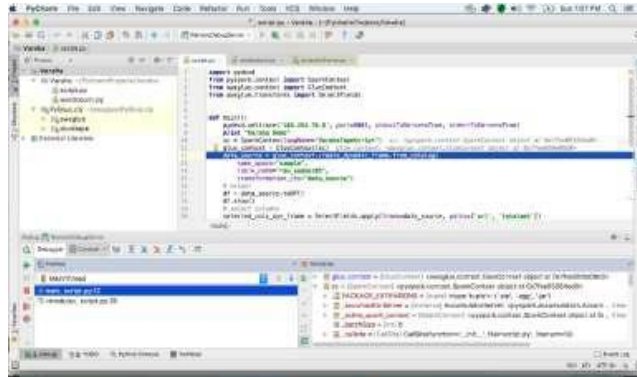
Job authoring: Glue transformations

Add transform

Name	Description
<input type="radio"/> DropFields	Drop fields from a DynamicFrame
<input type="radio"/> DropNullFields	DynamicFrame without null fields
<input type="radio"/> Join	Join two DynamicFrames
<input type="radio"/> MapToCollection	Apply a transform to each DynamicFrame in this DynamicFrameCollection
<input type="radio"/> Relationalize	Flatten nested schema and pivot out array columns from the flattened frame
<input type="radio"/> RenameField	Rename a field within a DynamicFrame
<input type="radio"/> SelectFields	Select fields from a DynamicFrame
<input type="radio"/> SelectFromCollection	Select one DynamicFrame from a DynamicFrameCollection
<input type="radio"/> SplitFields	Split fields within a DynamicFrame
<input type="radio"/> SplitRows	Split rows within a DynamicFrame based on comparators
<input type="radio"/> Unbox	Unbox a string field

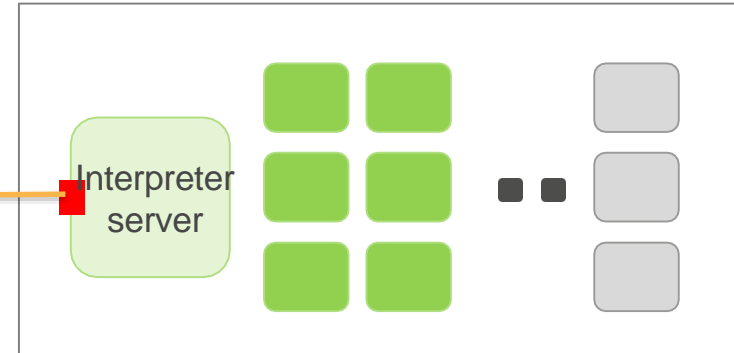
- **Prebuilt transformation:** Click and add to your job with simple configuration
- **Spigot** writes sample data from DynamicFrame to S3 in JSON format
- **Expanding...** more transformations to come

Job authoring: Developer endpoints

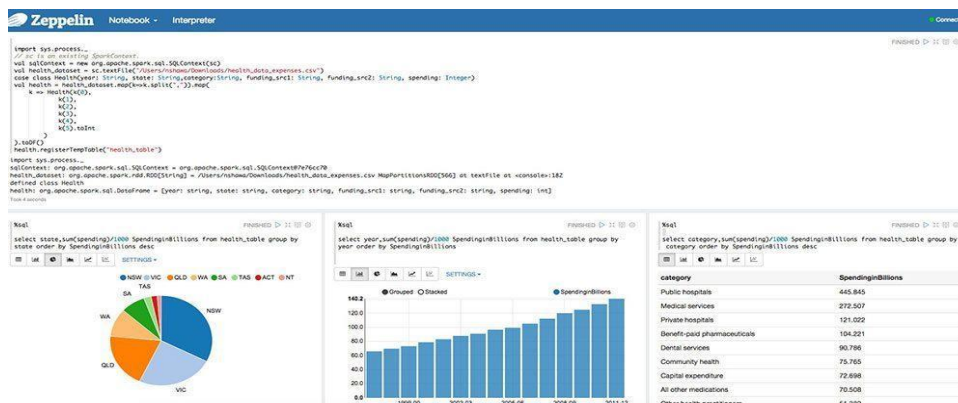


Remote
interpreter

Glue Spark environment



- Environment to **iteratively develop** and test ETL code.
- Connect your IDE or notebook (e.g. **Zeppelin**) to a Glue development endpoint.
- When you are satisfied with the results you can create an ETL job that runs your code.



<https://www.cloudera.com/products/open-source/apache-hadoop/apache-zeppelin.html>

Amazon Web Services, Inc.

Job execution: Scheduling and monitoring

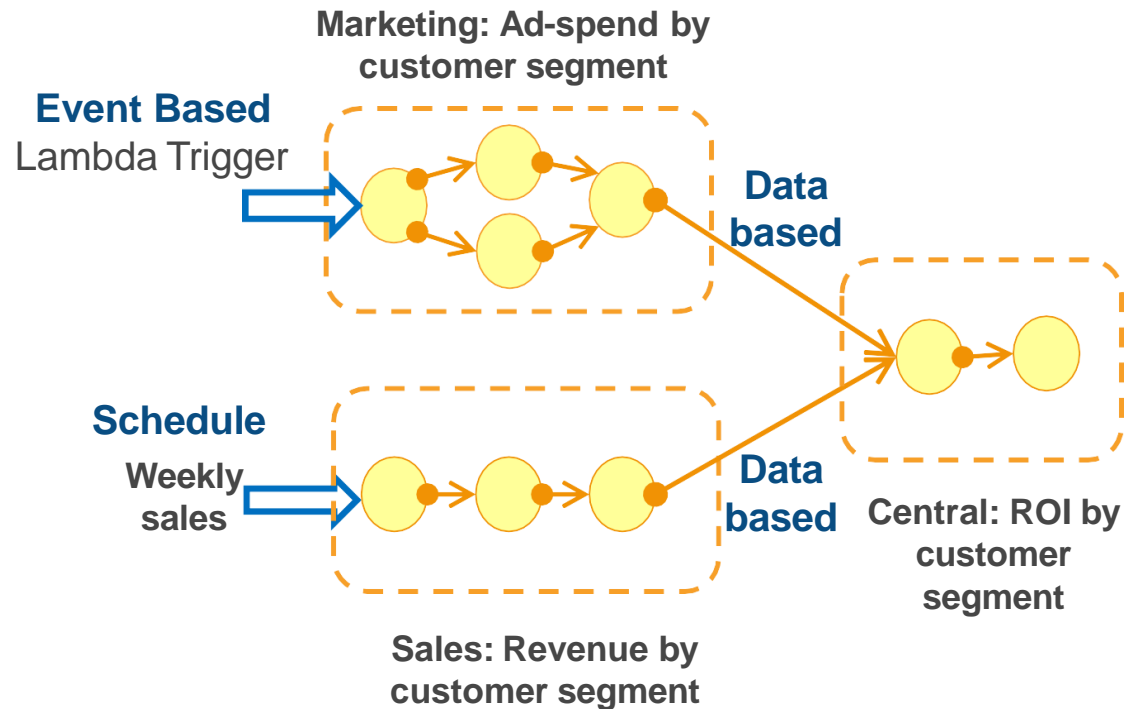
Compose jobs globally with event-based dependencies

- Easy to reuse and leverage work across organization boundaries

Multiple triggering mechanisms

- **Schedule-based:** e.g., time of day
- **Event-based:** e.g., job completion
- **On-demand:** e.g., AWS Lambda
- **More coming soon:** Data Catalog based events, S3 notifications and Amazon CloudWatch events

Logs and alerts are available in Amazon CloudWatch

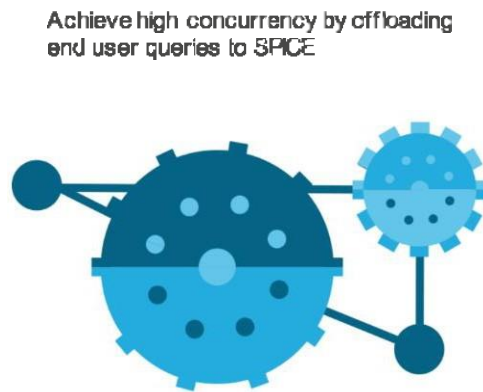


QuickSight for BI on Redshift

Amazon QuickSight allows everyone in your organization to understand your data by asking questions in natural language, exploring through interactive dashboards, or automatically looking for patterns and outliers powered by machine learning.



Amazon Redshift



Achieve high concurrency by offloading end user queries to SPICE

Calculations can be done in SPICE reducing the load on the underlying database



<https://aws.amazon.com/quicksight/>



Apache Flume

What is Flume

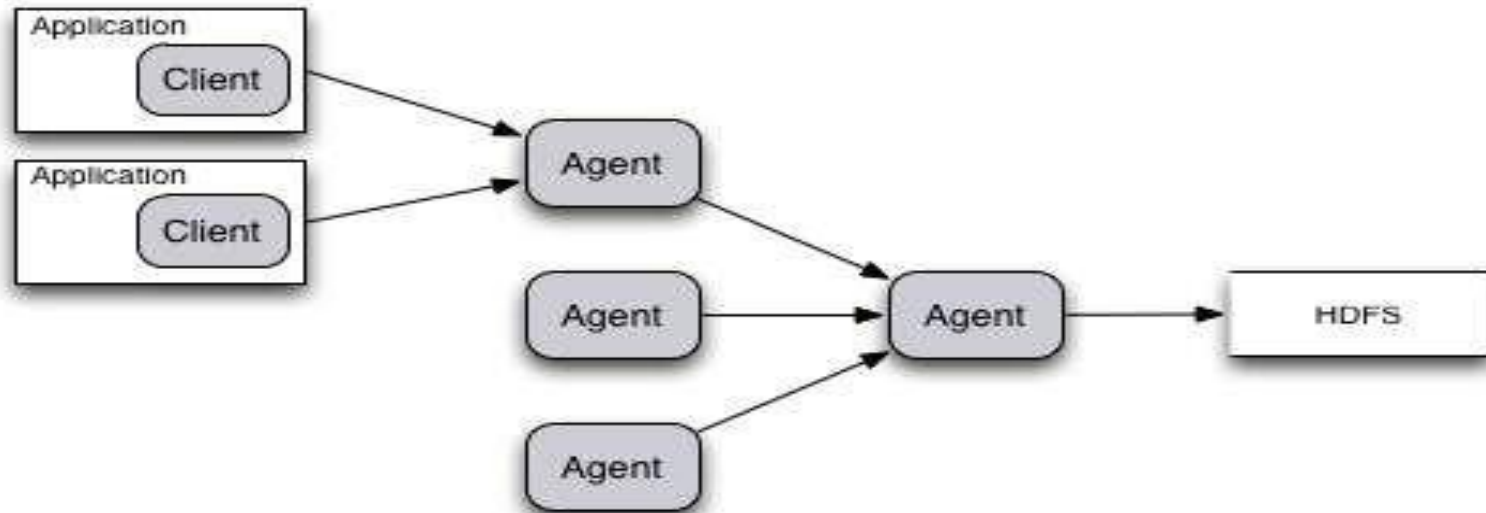
- Collection, Aggregation of streaming Event Data
 - Typically used for log data
- Significant advantages over ad-hoc solutions
 - Reliable, Scalable, Manageable, Customizable and High Performance
 - Declarative, Dynamic Configuration
 - Contextual Routing
 - Feature rich
 - Fully extensible

Core Concepts: Event

An Event is the fundamental unit of data transported by Flume from its point of origination to its final destination. Event is a byte array payload accompanied by optional headers.

- Payload is opaque to Flume
- Headers are specified as an unordered collection of string key-value pairs, with keys being unique across the collection
- Headers can be used for contextual routing

Typical Aggregation Flow



$[Client]^+ \rightarrow Agent \rightarrow Agent^* \rightarrow Destination$

Core Concepts: Channel

A passive component that buffers the incoming events until they are drained by Sinks.

- Different Channels offer different levels of persistence:
 - Memory Channel: volatile
 - Data lost if JVM or machine restarts
 - File Channel: backed by WAL implementation
 - Data not lost unless the disk dies.
 - Eventually, when the agent comes back data can be accessed.

Core Concepts: Sink

An active component that removes events from a Channel and transmits them to their next hop destination.

- Different types of Sinks:
 - Terminal sinks that deposit events to their final destination. For example: HDFS, HBase, Morphline-Solr, Elastic Search
 - Sinks support serialization to user's preferred formats.
 - HDFS sink supports time-based and arbitrary bucketing of data while writing to HDFS.
 - IPC sink for Agent-to-Agent communication: Avro, Thrift
- Require exactly one channel to function

Flume Agent



- Installed on each node
- Collects events

Flume Agent



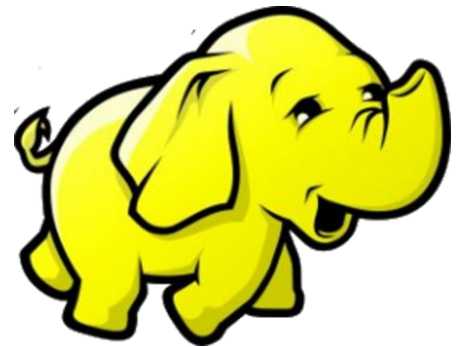
Flume Agent



Flume Agent



Flume Agent



Flume Agent

Interceptor



- Filters useless events
- Decorates events by adding metadata
 - e.g. timestamp, hostname, UUID, static markers

Flume Agent

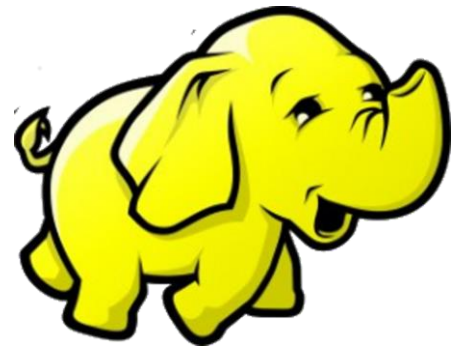
Flume Agent

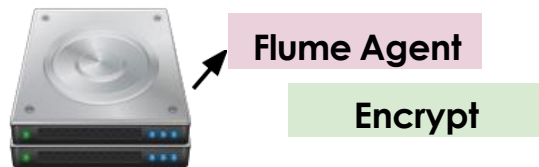
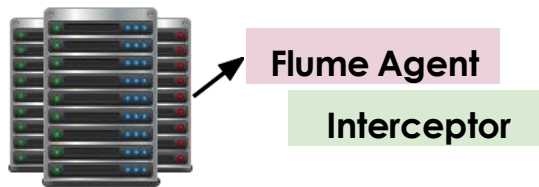


Flume Agent

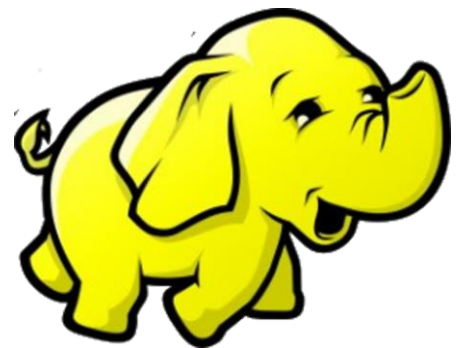
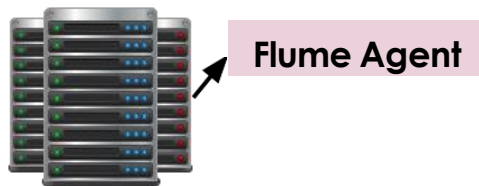
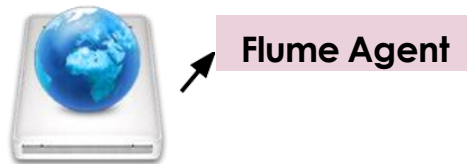


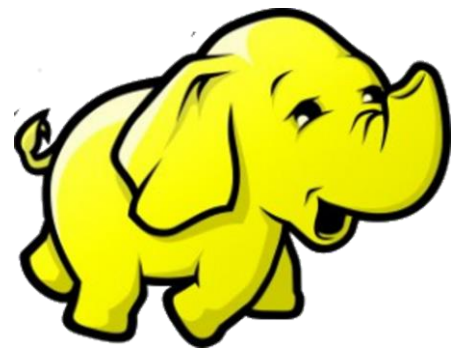
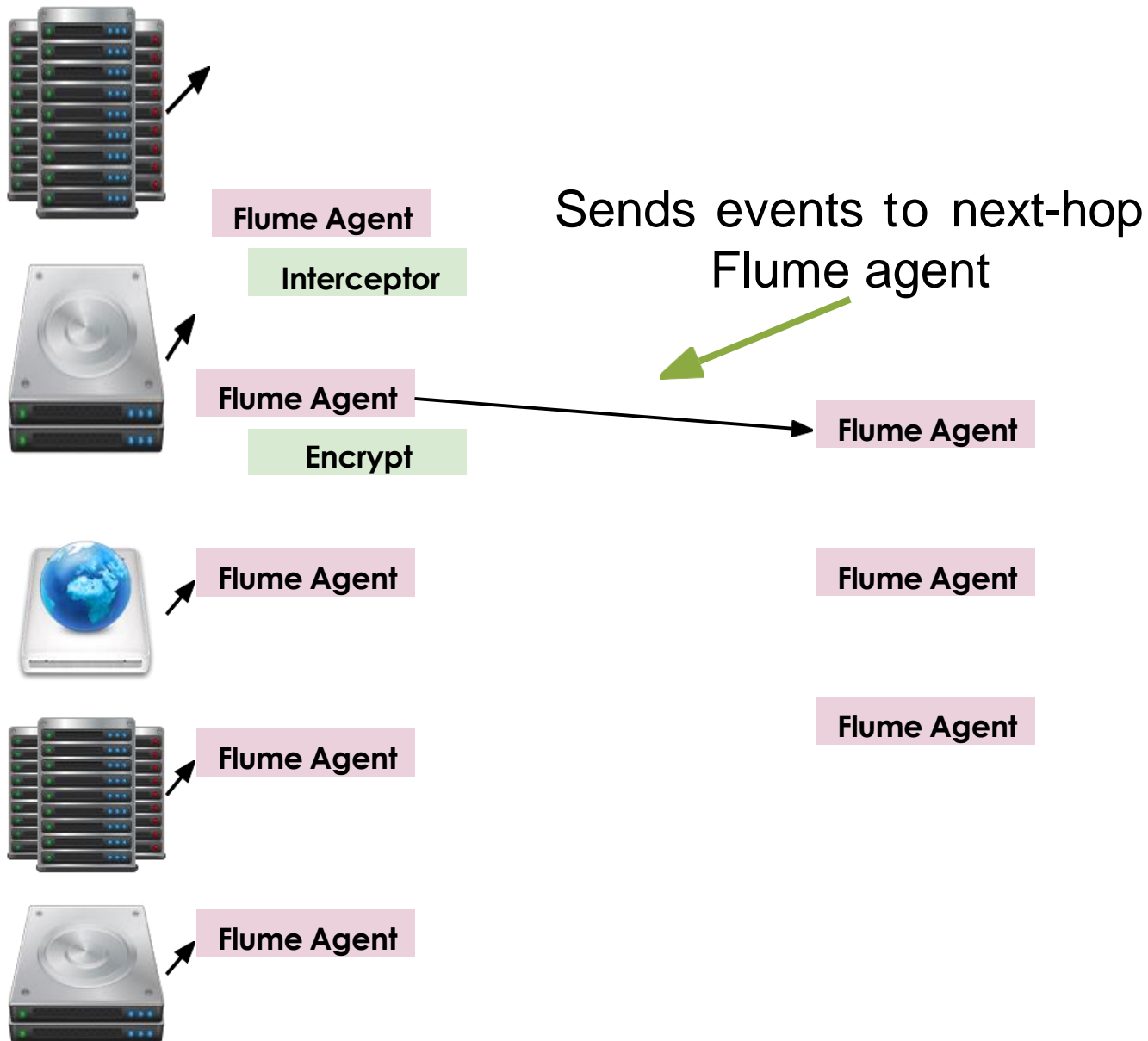
Flume Agent

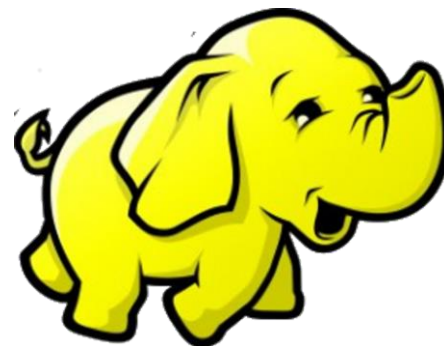
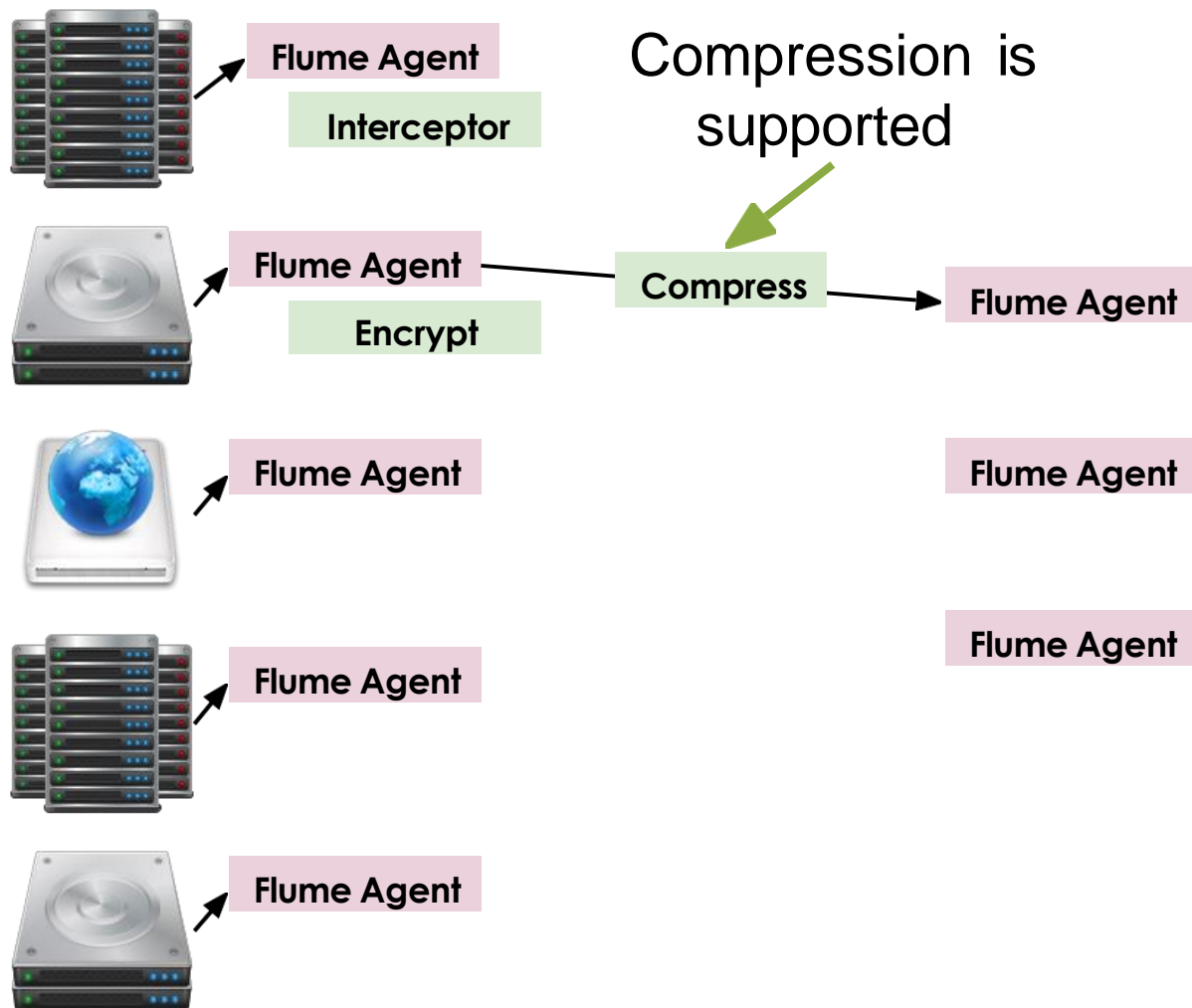


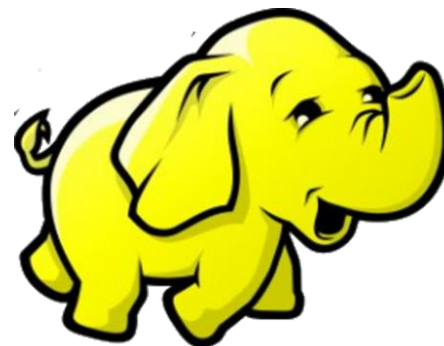
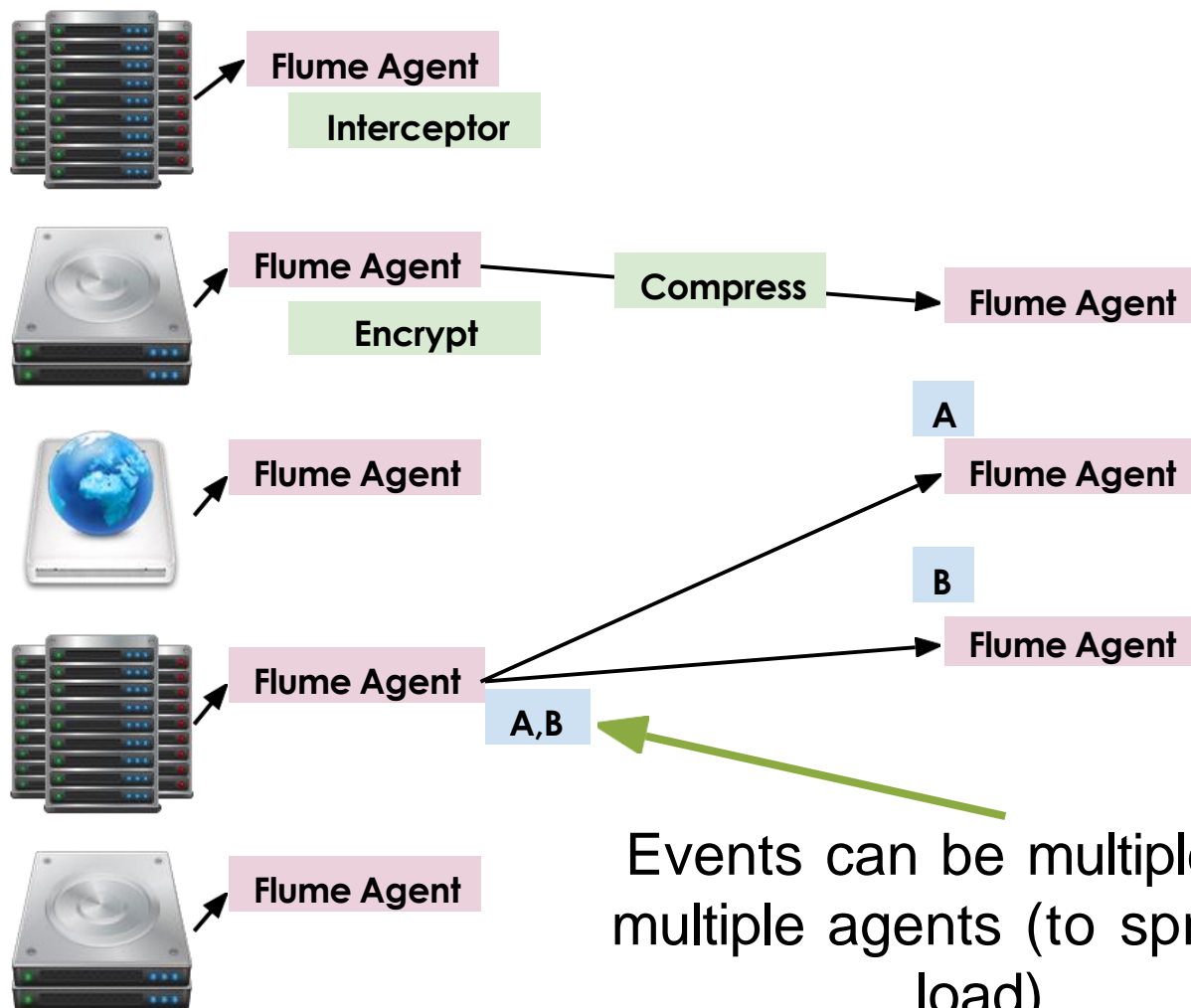


Encrypts events in
a file on disk

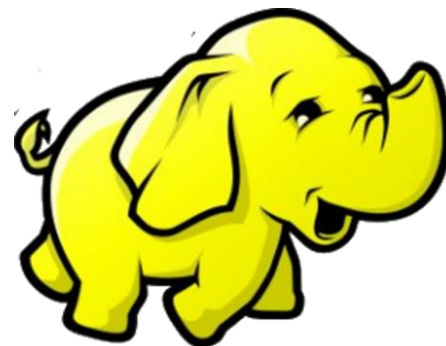
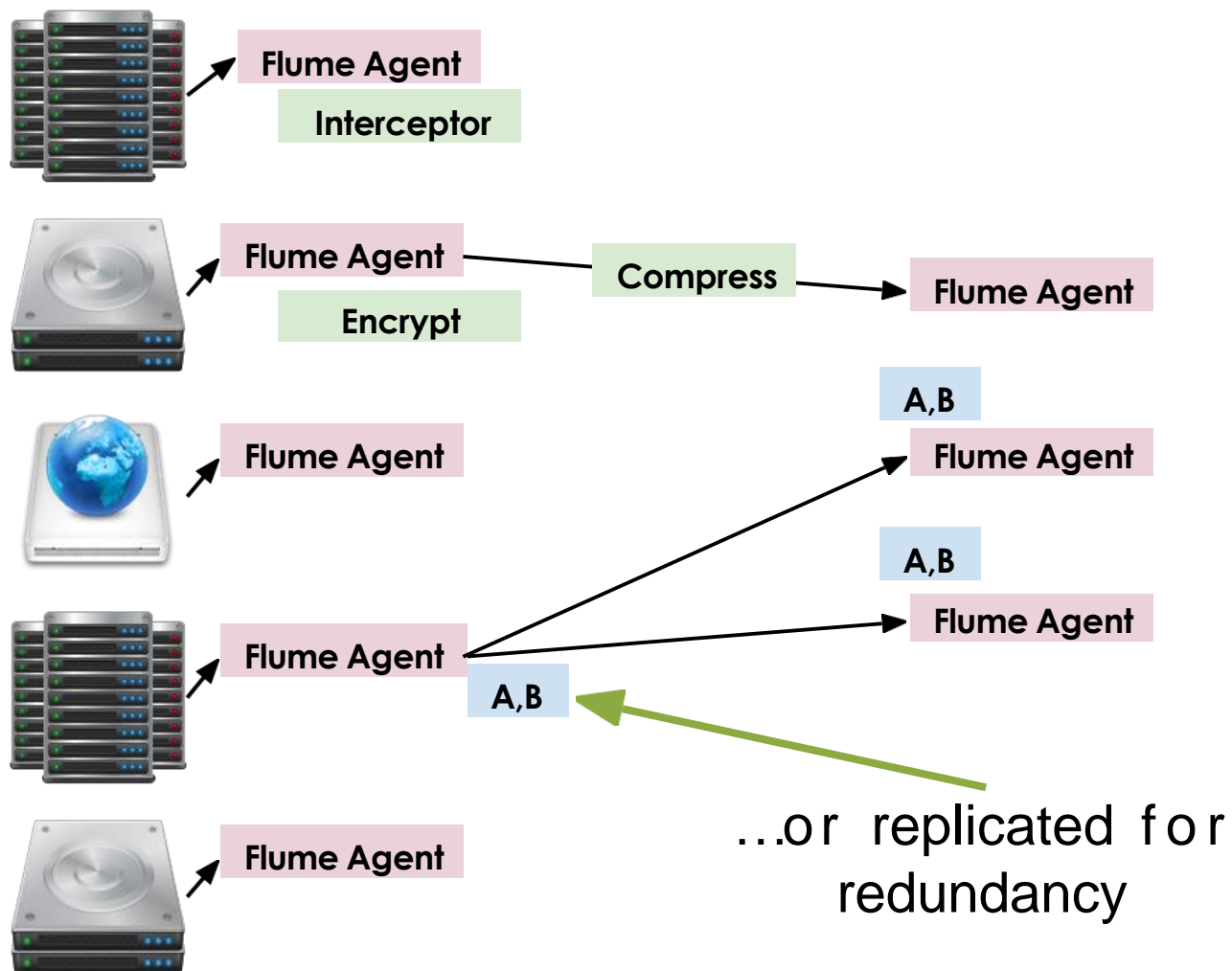


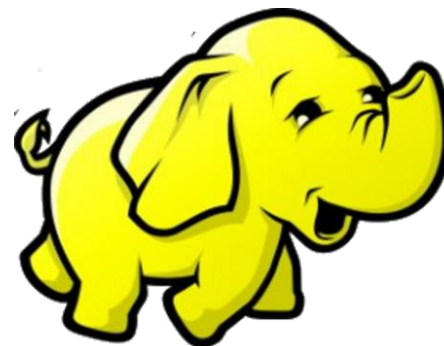
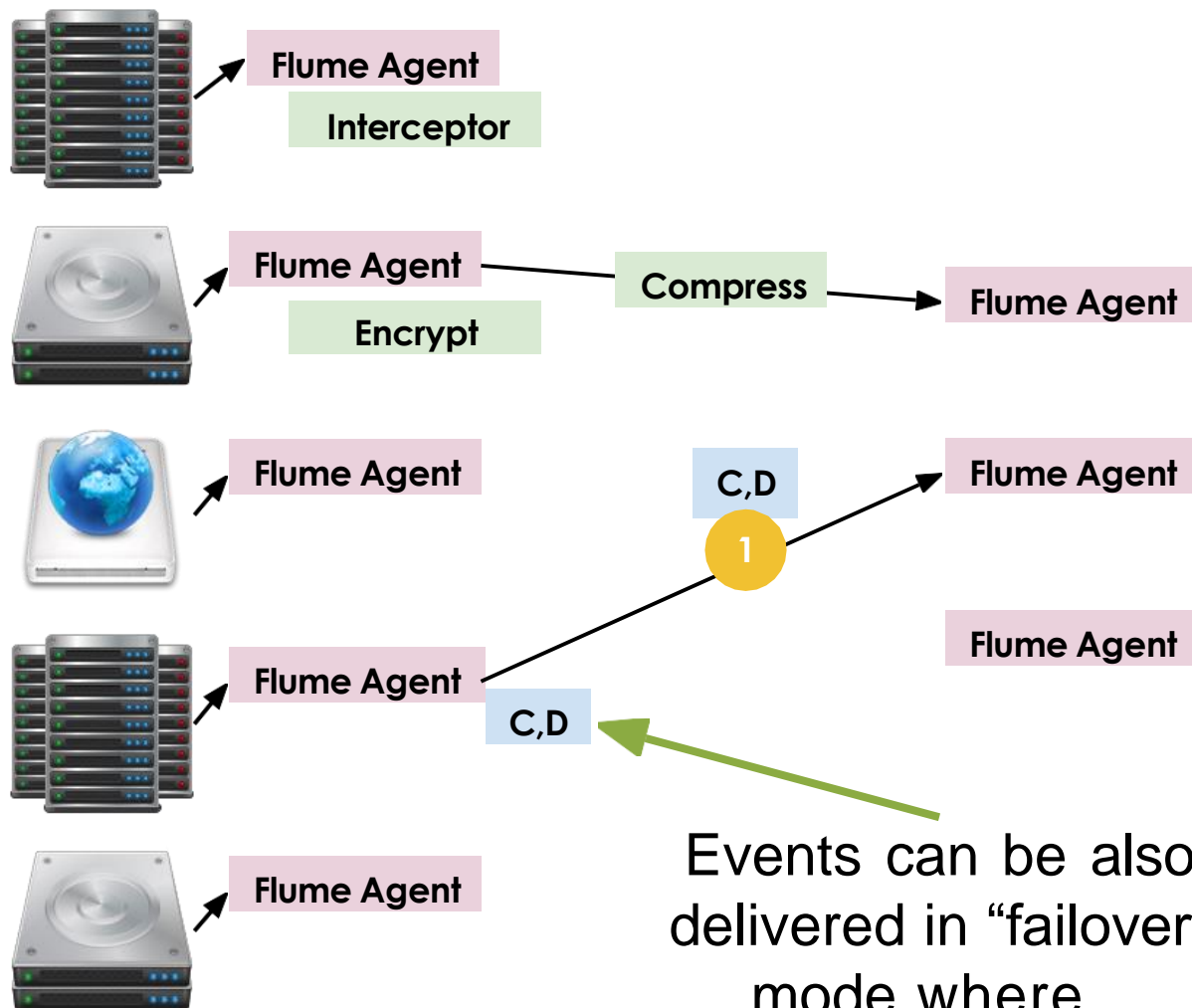


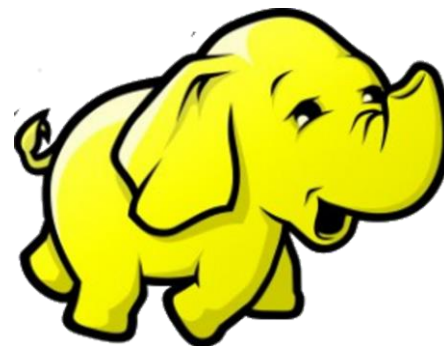
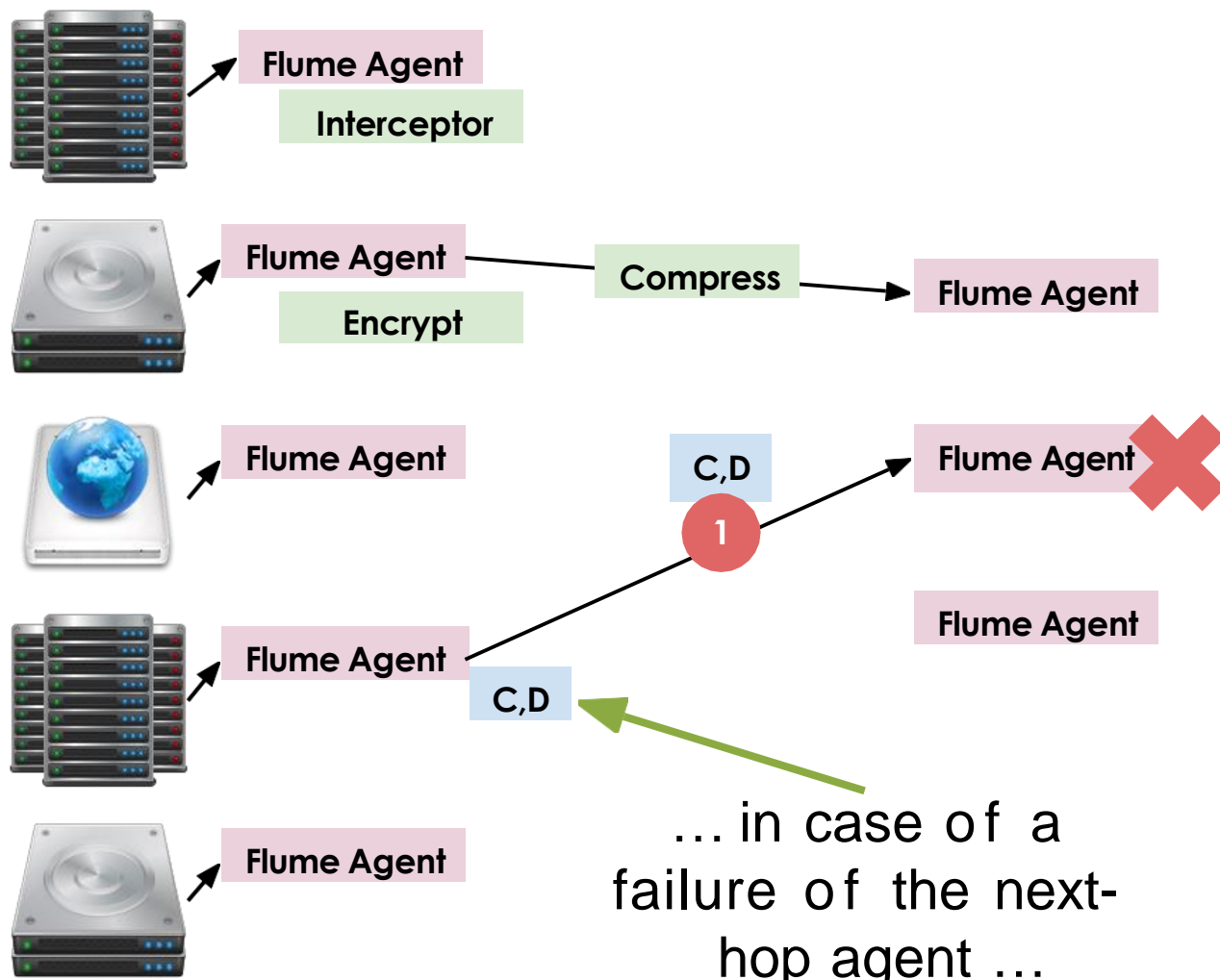


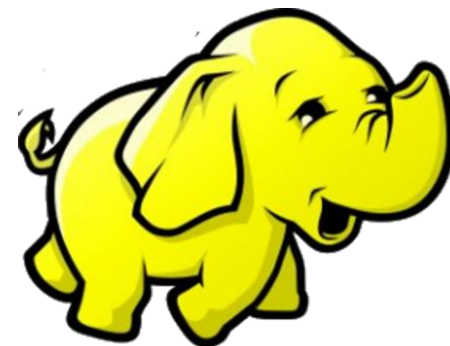
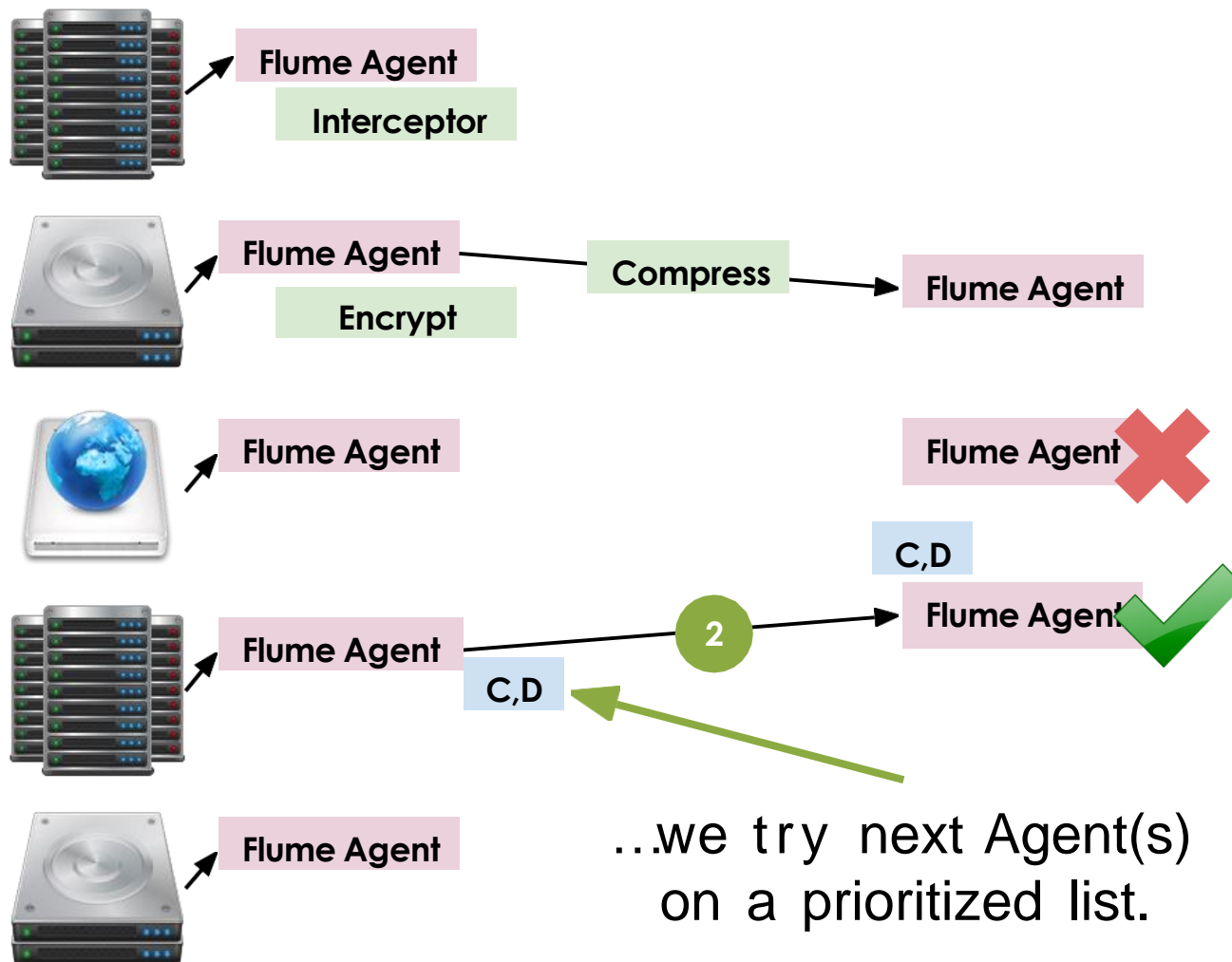


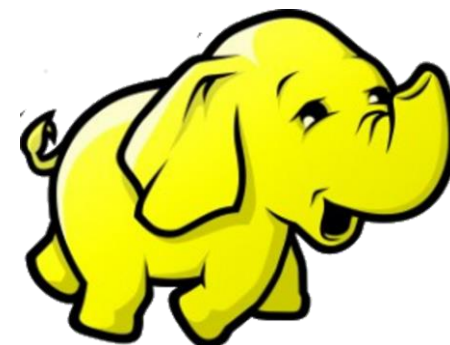
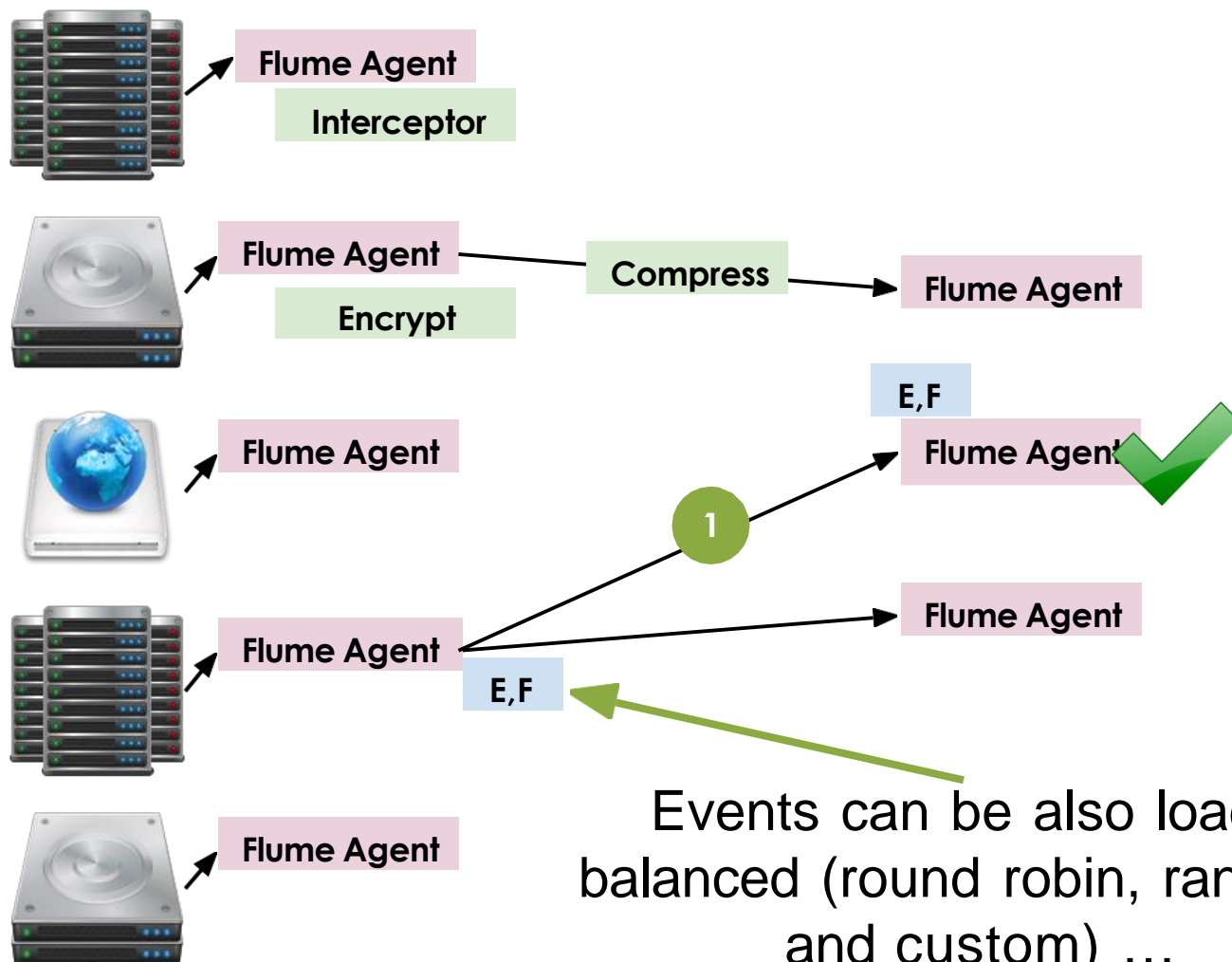
Events can be multiplexed to multiple agents (to spread the load) ...

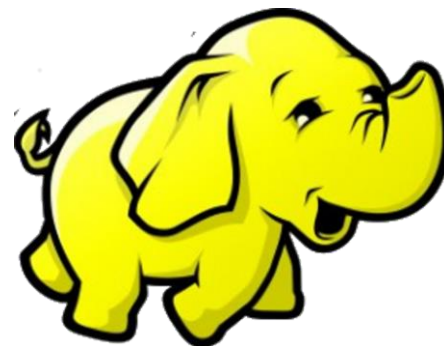
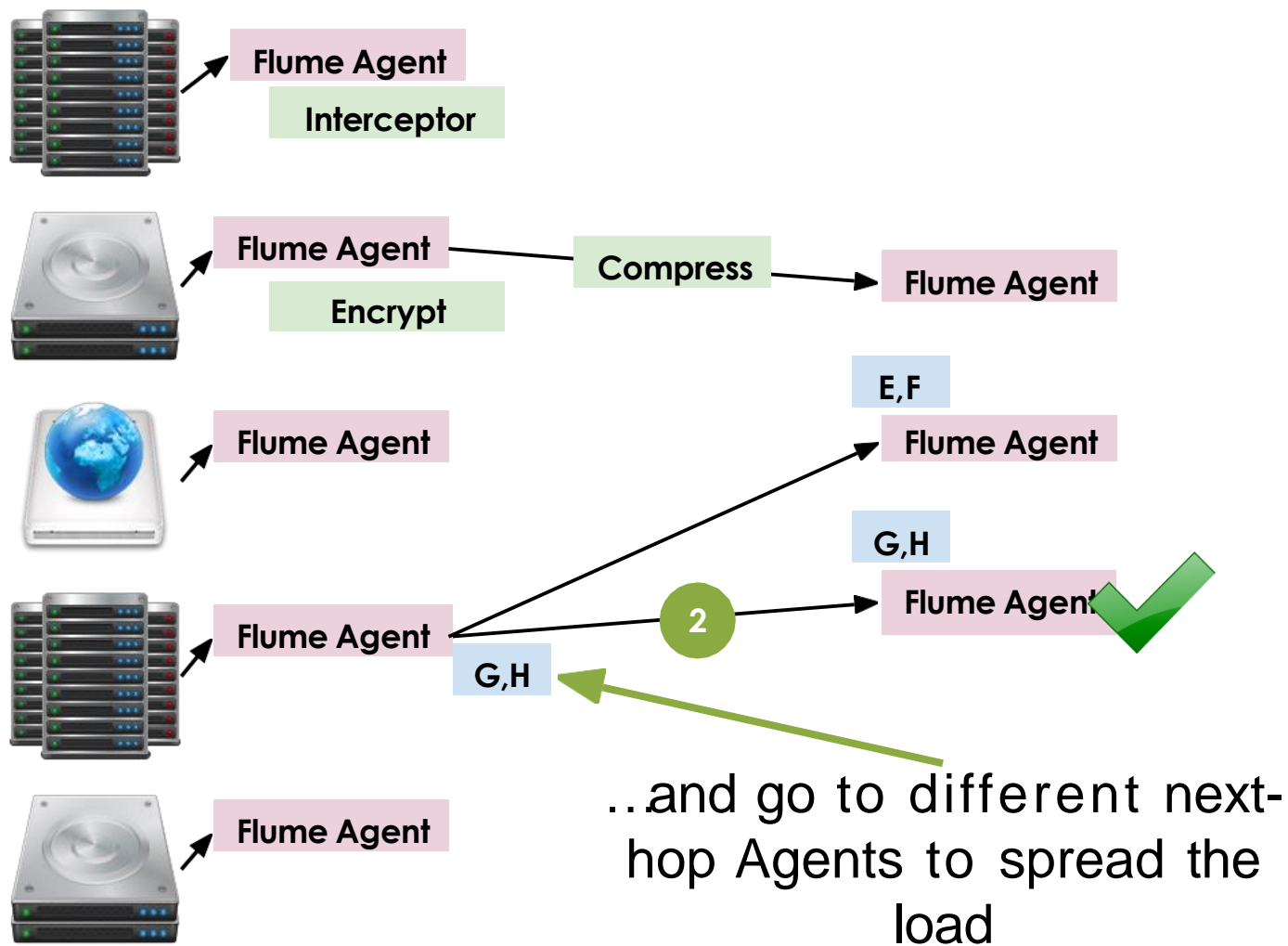


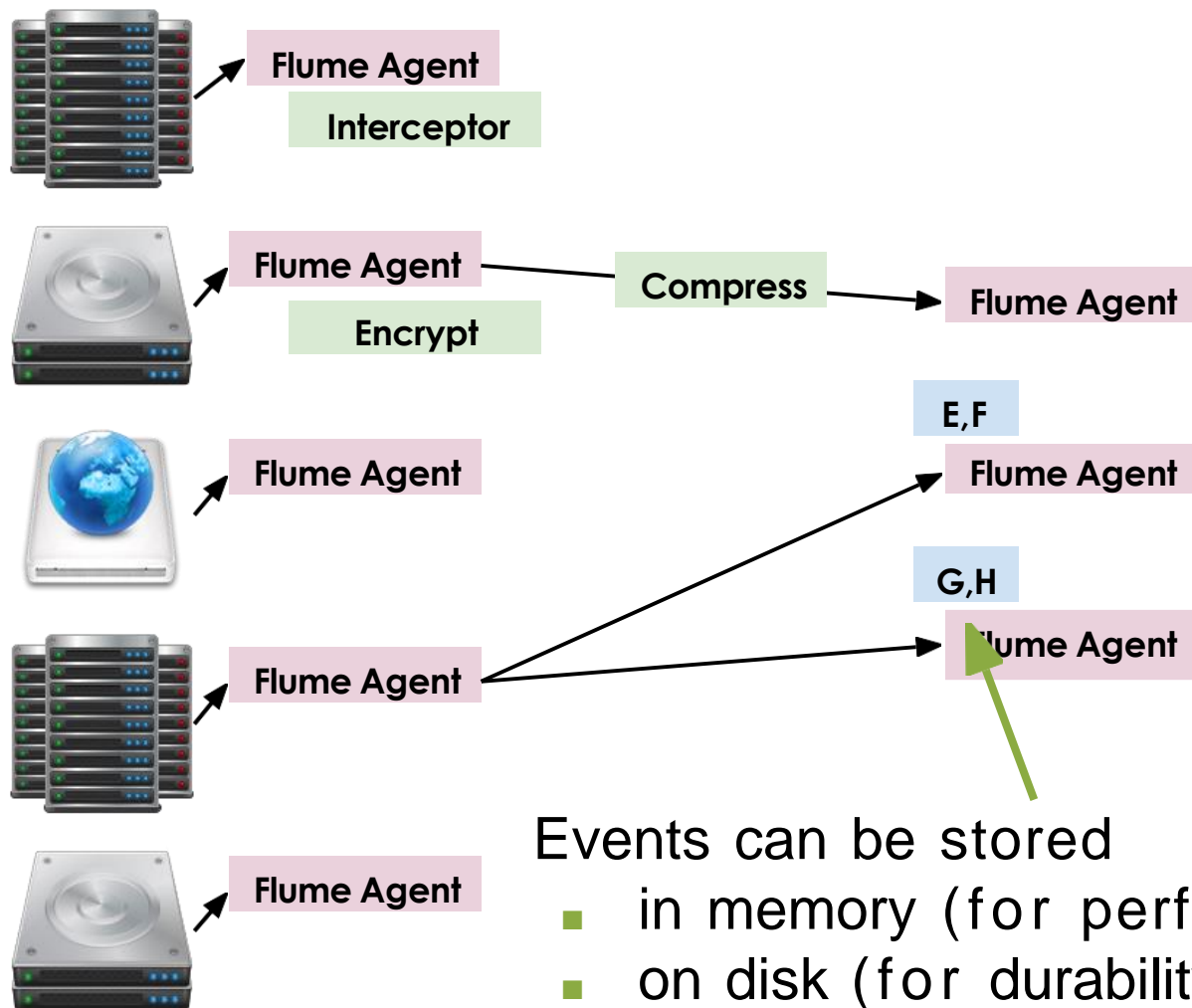






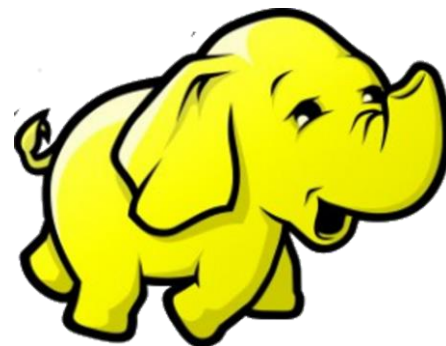


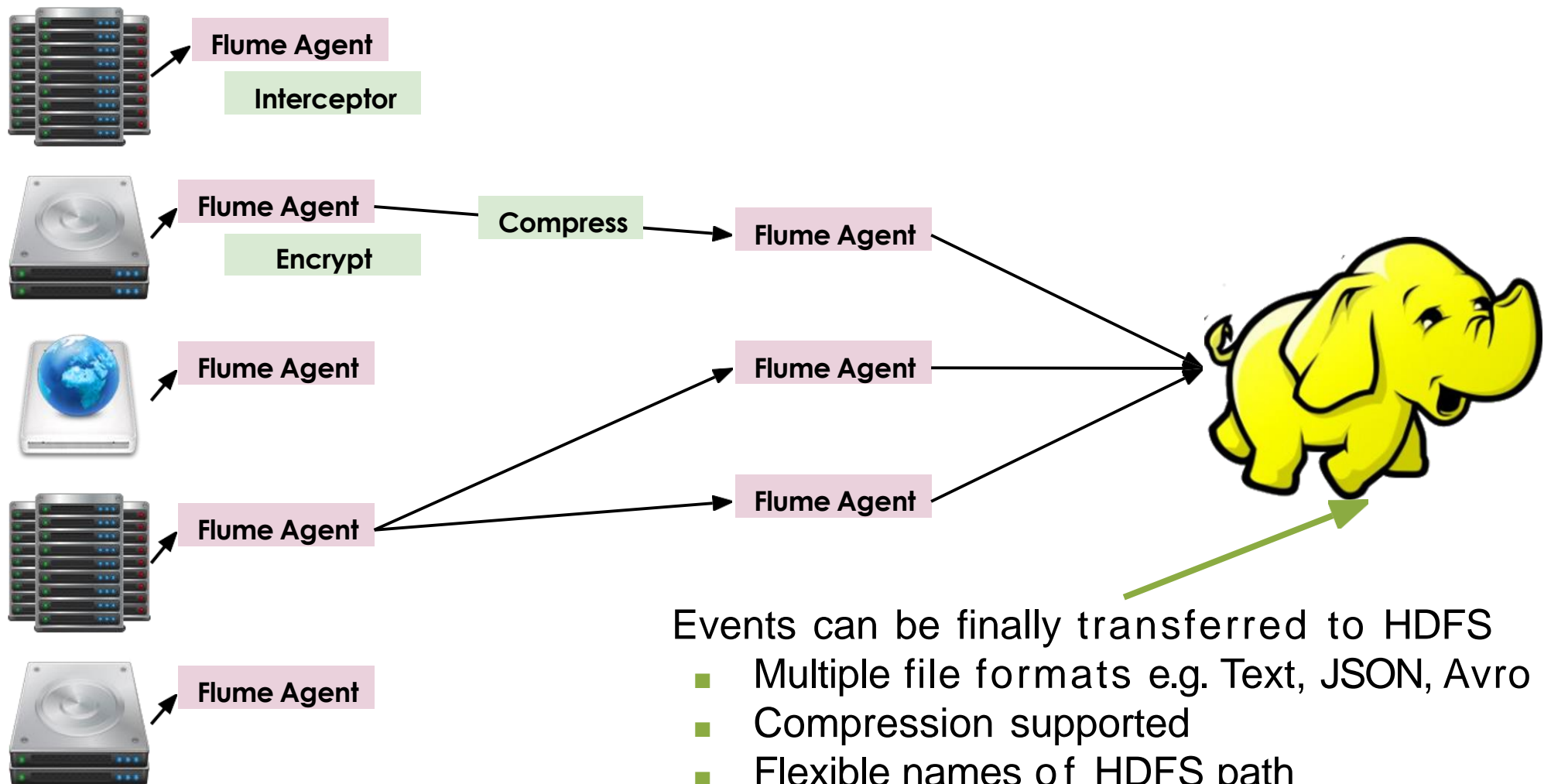




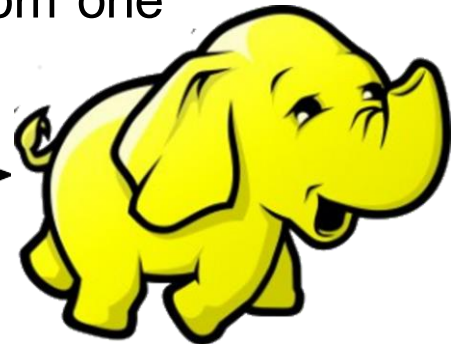
Events can be stored

- in memory (for performance)
- on disk (for durability)





- Many destinations are supported
- One can implement a custom one



APACHE
HBASE



elasticsearch.



Flume

- **Distributed**
 - Agents installed on many machines
- **Scalable**
 - Add more machines to transfer more events
- **Reliable**
 - Durable storage, failover and/or replication
- **Manageable**
 - Easy to install, configure, reconfigure and run



Flume

- **Nicely integrated with the Hadoop Ecosystem**
 - Various destinations e.g. HDFS, HBase
 - Various file formats e.g. Avro, SequenceFile
- **Extensible**
 - Possibility to add new functionality e.g. source and destination for events

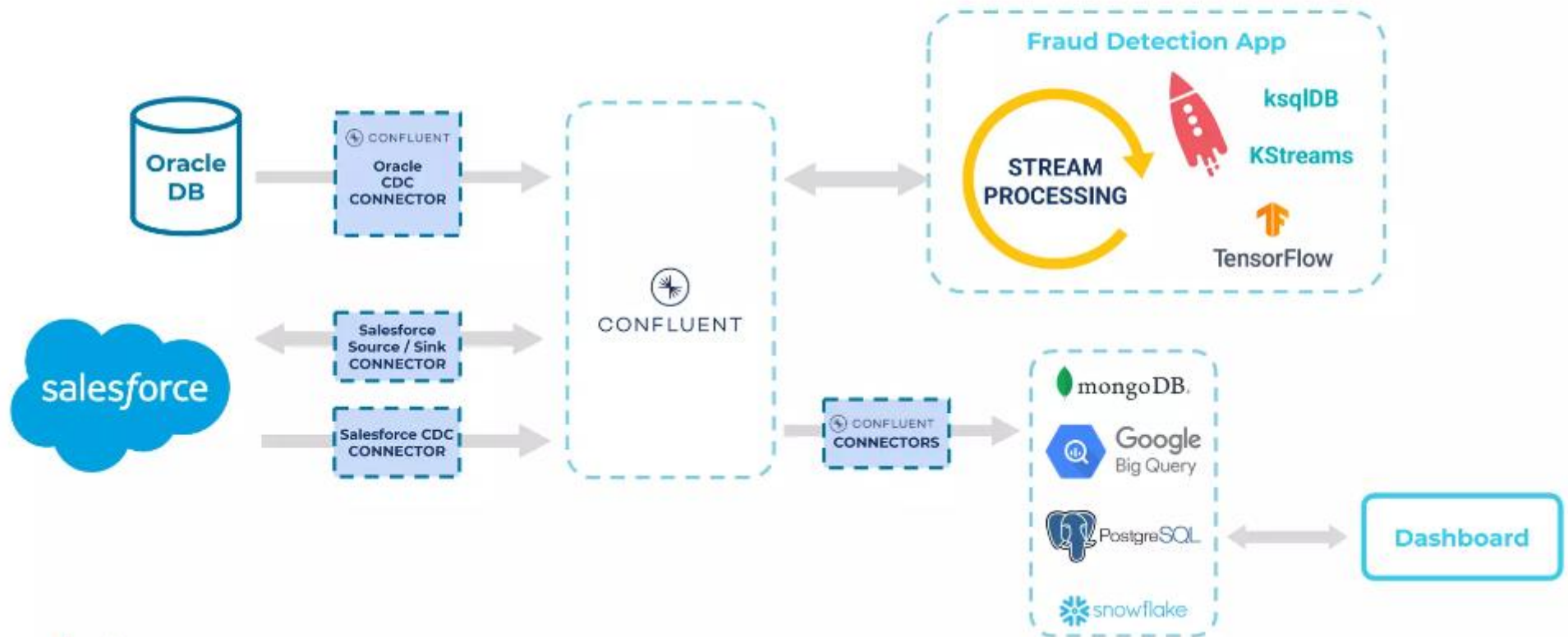
An Event Streaming Platform

The Underpinning of Data in Motion



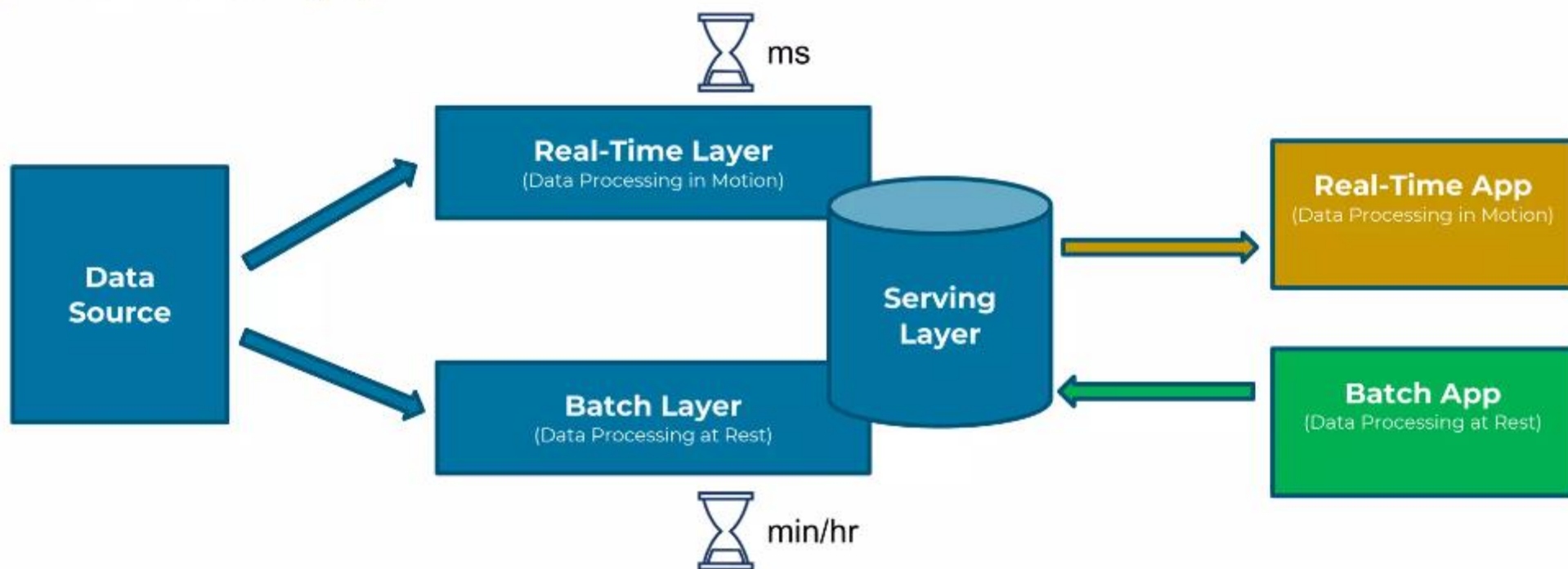
Example Architecture for Data in Motion

Real-time decision making for claim processing and fraud detection



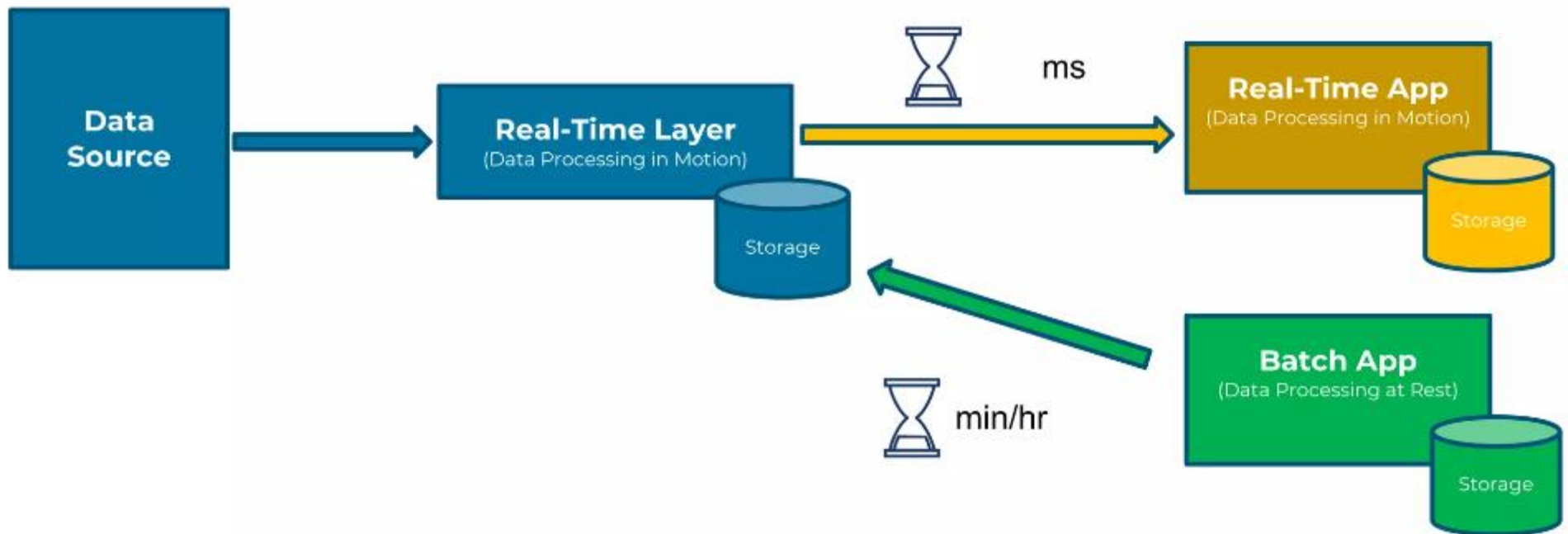
Lambda Architecture

Option 1: Unified serving layer

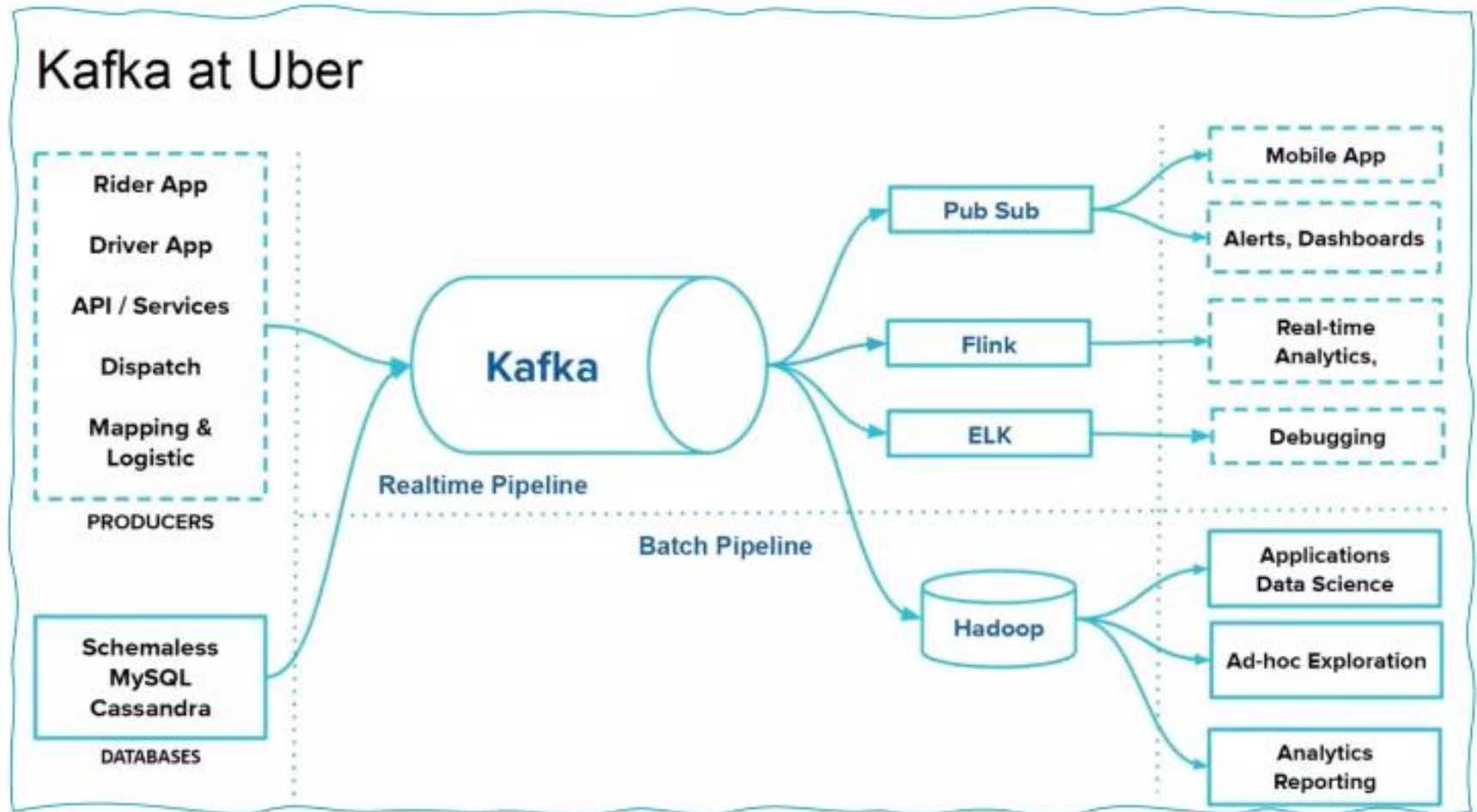


Kappa Architecture

One pipeline for real-time and batch consumers

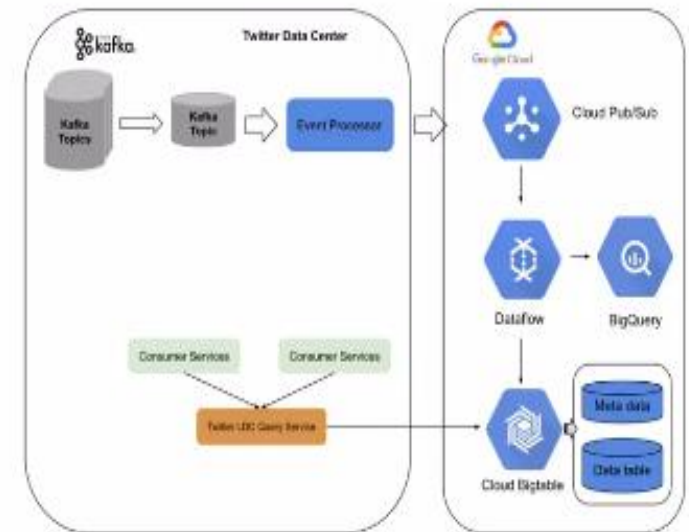
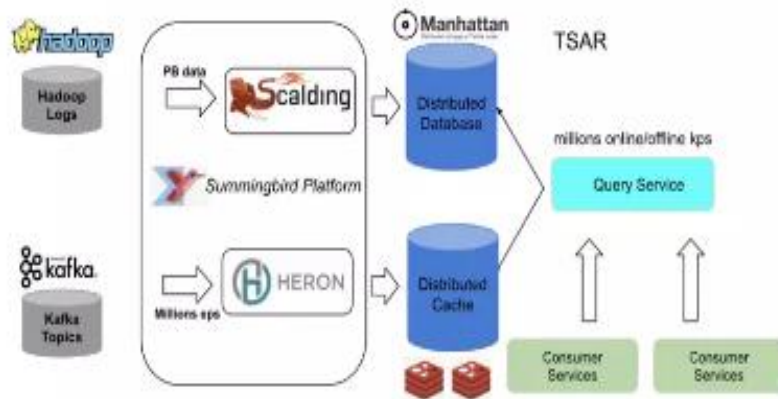


Kappa @ Uber



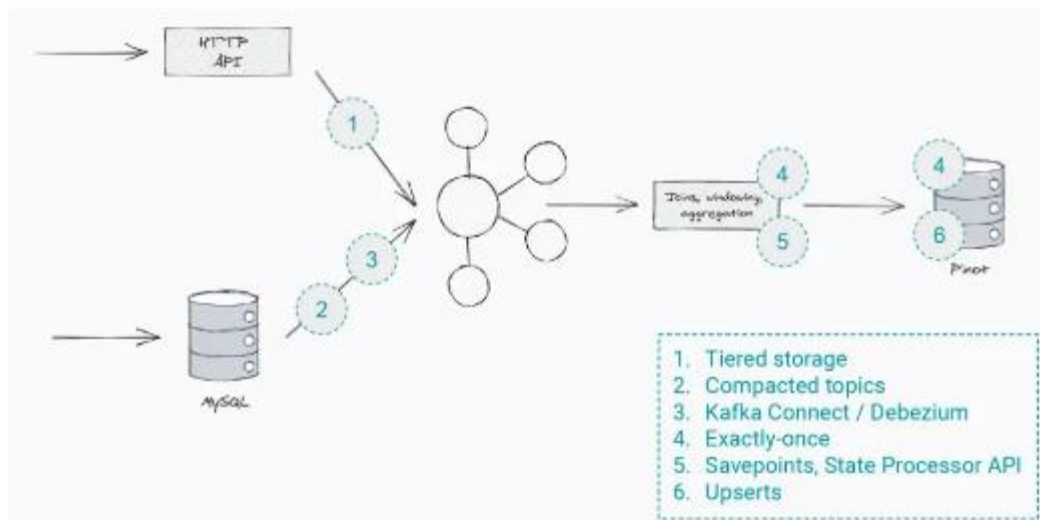
Flink ML is a library which provides machine learning (ML) APIs and infrastructures that simplify the building of ML pipelines

Kappa @ Twitter



Migration from Hadoop and Kafka to a hybrid architecture on both Twitter data center and Google Cloud Platform with Kafka and GCP, Twitter is able to process billions of events in real-time and achieve low latency, high accuracy, stability, architecture simplicity, and reduced operation cost

Kappa @ Shopify



Kappa Building Blocks

The Log (Kafka)

- Durability with Topic Compaction and Tiered Storage
- Consistency via Exactly-Once Semantics (EOS)
- Data Integration via Kafka Connect
- Elasticity via dynamic Kafka clusters

Streaming Framework (Kafka Streams / Flink)

- Reliability and scalability
- Fault tolerance
- State management

Sinks

- Update/Upsert for simplified design:
 - RDBMS, NoSQL, Compacted Kafka Topics
- Append-only: Regular Kafka Topics, Time Series