

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>

# Data Modeling, Transformation and Serving

---



DeepLearning.AI

## Data Transformations



DeepLearning.AI

# Data Transformations

---

## Week 3 Overview



DeepLearning.AI

## Batch Transformations

---

**Batch Transformation Patterns  
and Use Cases**

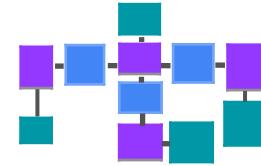
# Transformations for Data Modeling

Target Model

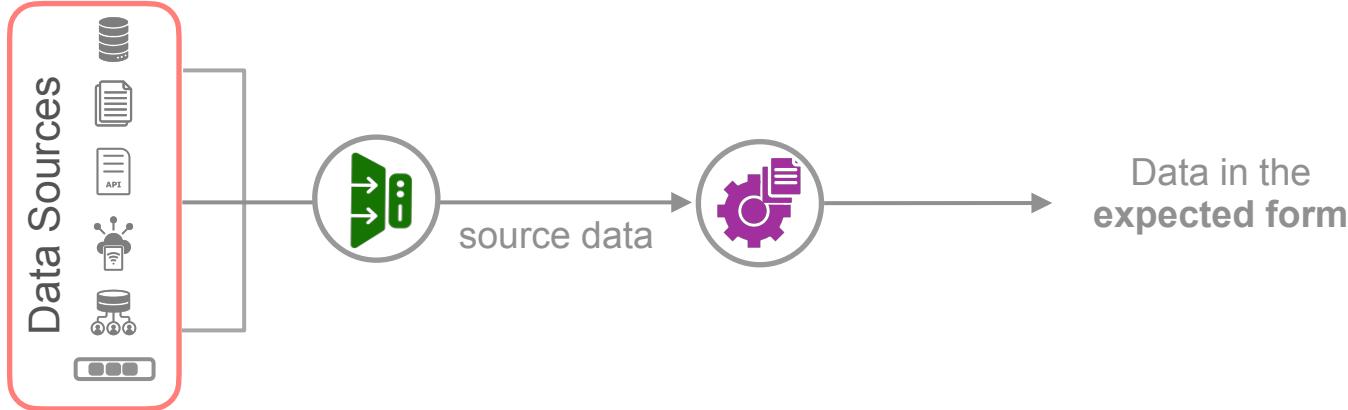
Star schemas

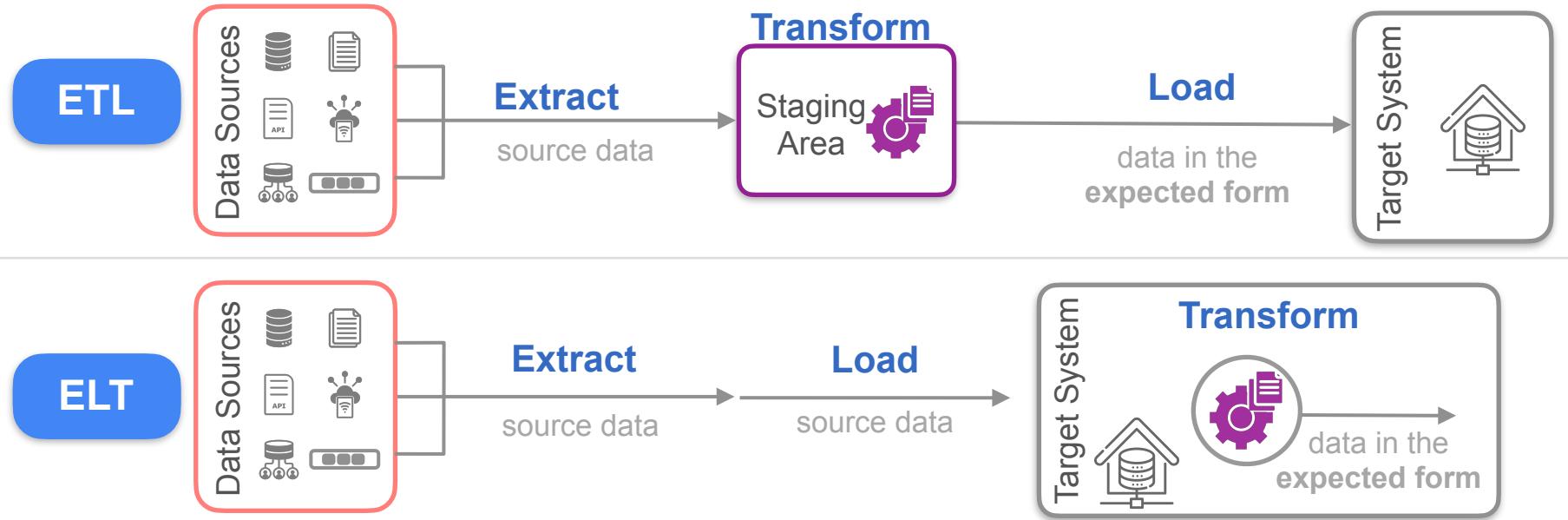


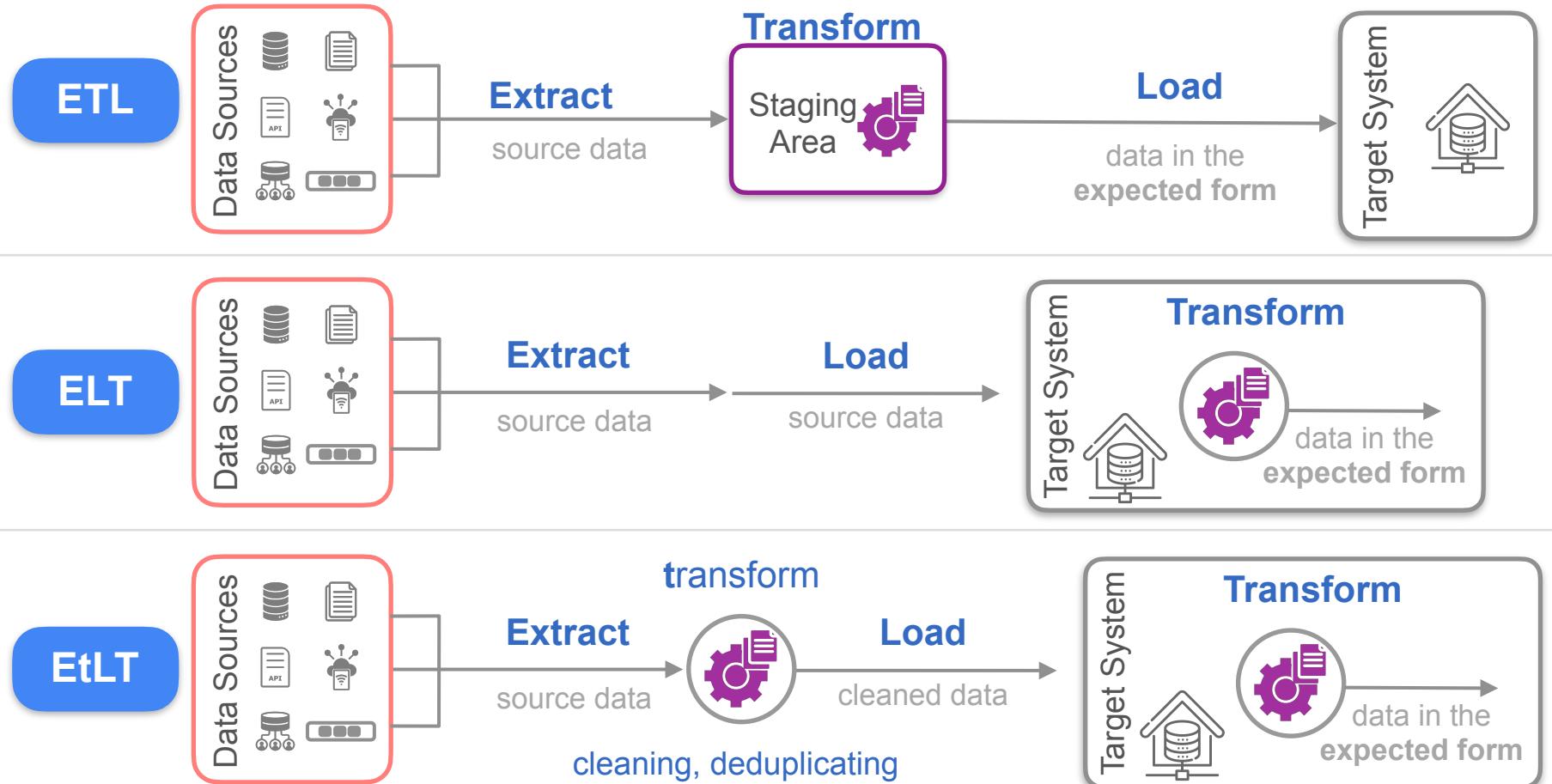
Data Vault



Restructure the source data into the **expected form**



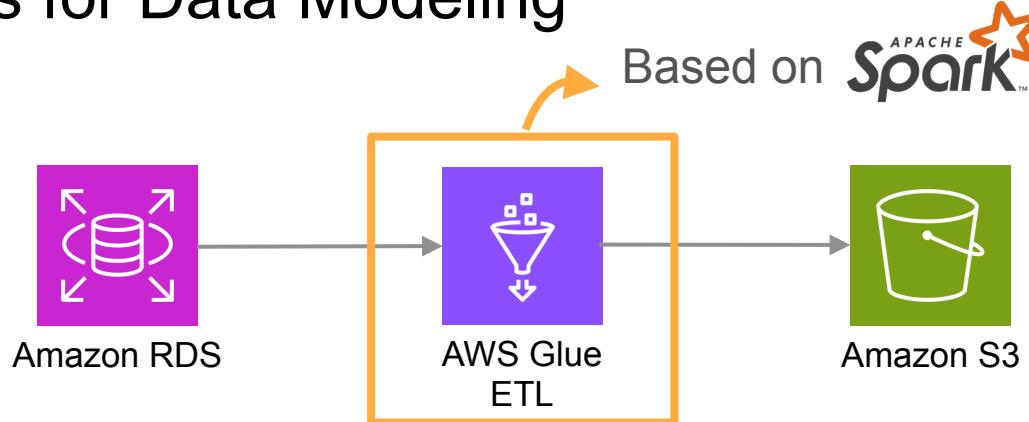




# Transformations for Data Modeling

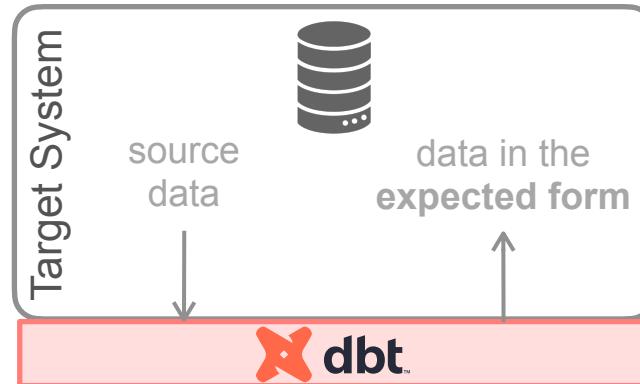
ETL

Course 1



ELT

Course 4

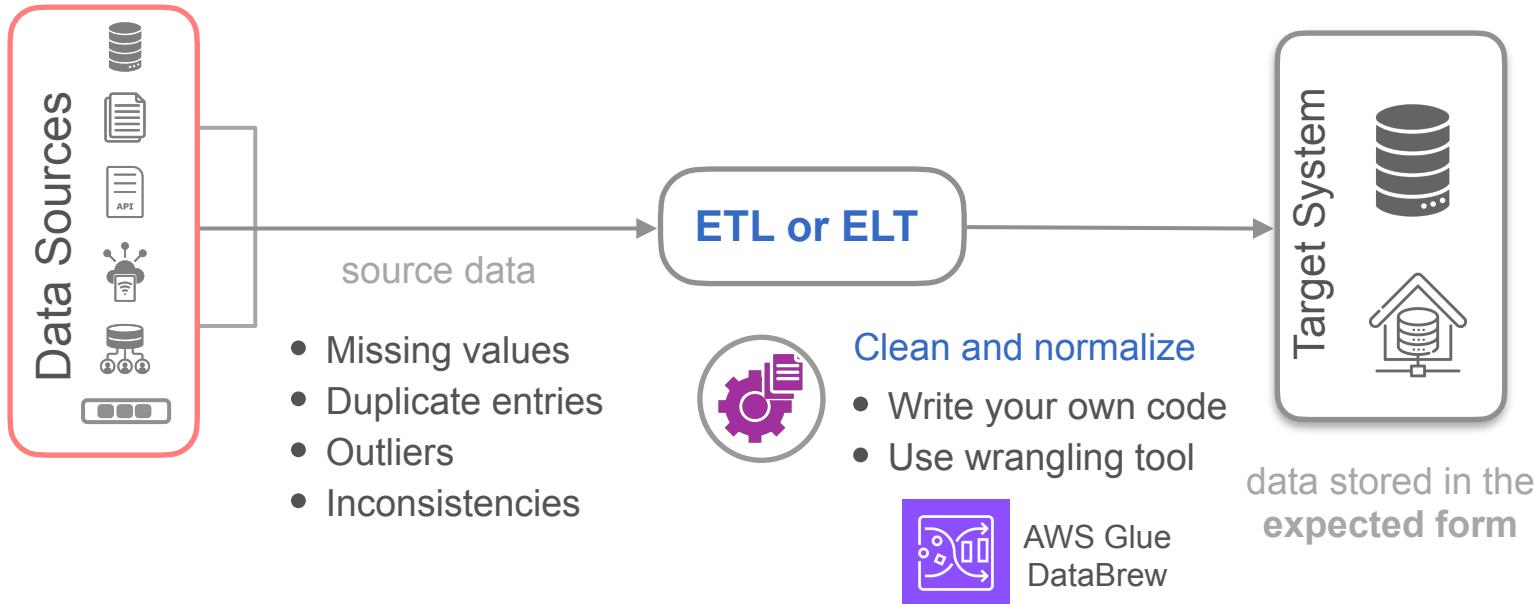


- dbt is not an execution tool like Spark
- SQL-tool that facilitates transformation within the database or data warehouse

# Transformations for Data Cleaning

## Data Wrangling

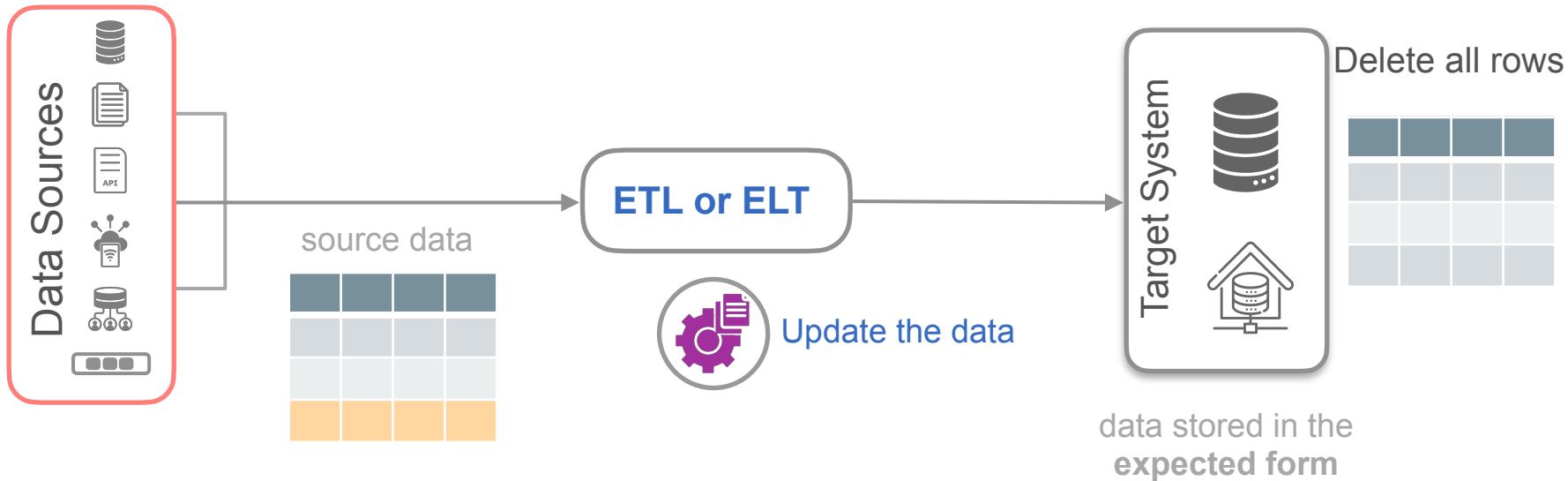
Take messy, malformed data and turn it into clean data



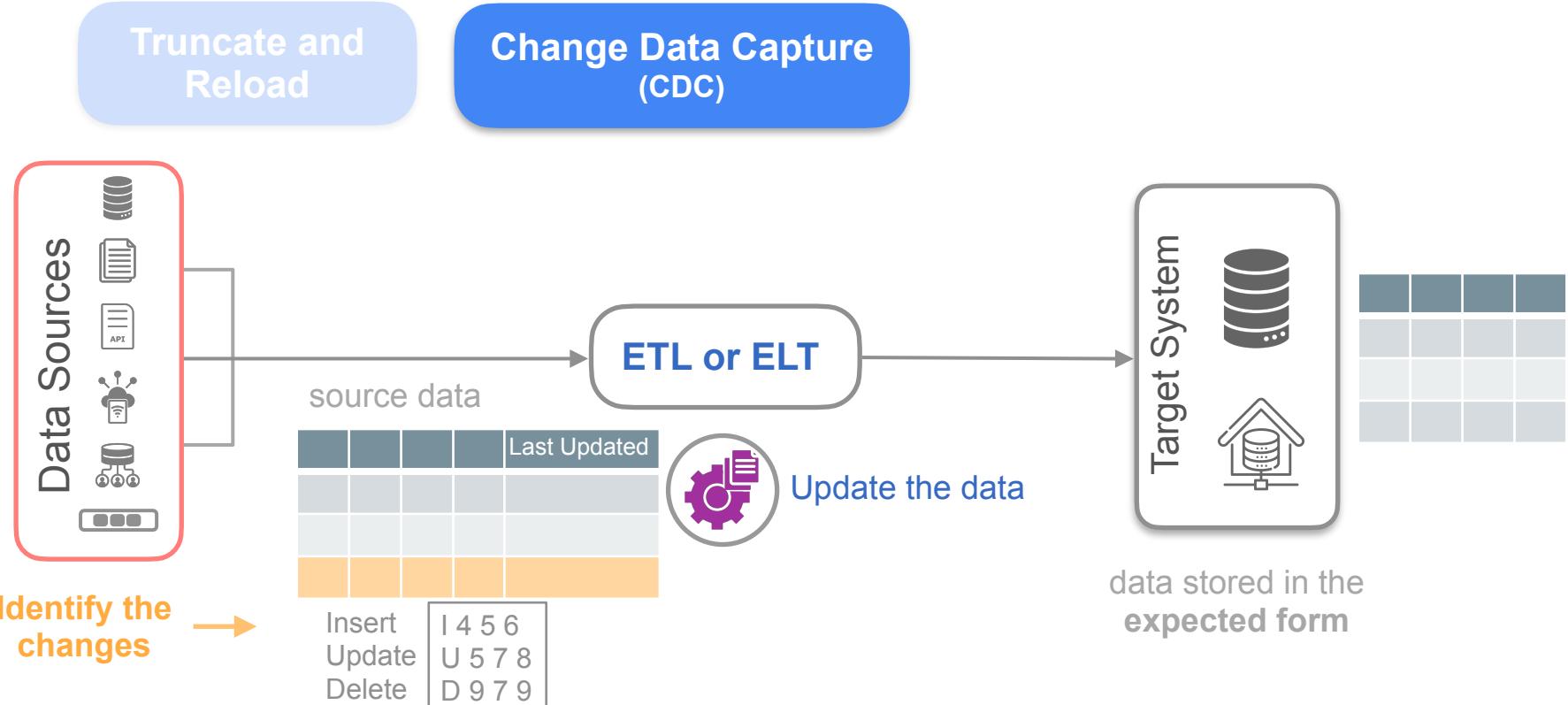
# Transformations for Data Updating

## Truncate and Reload

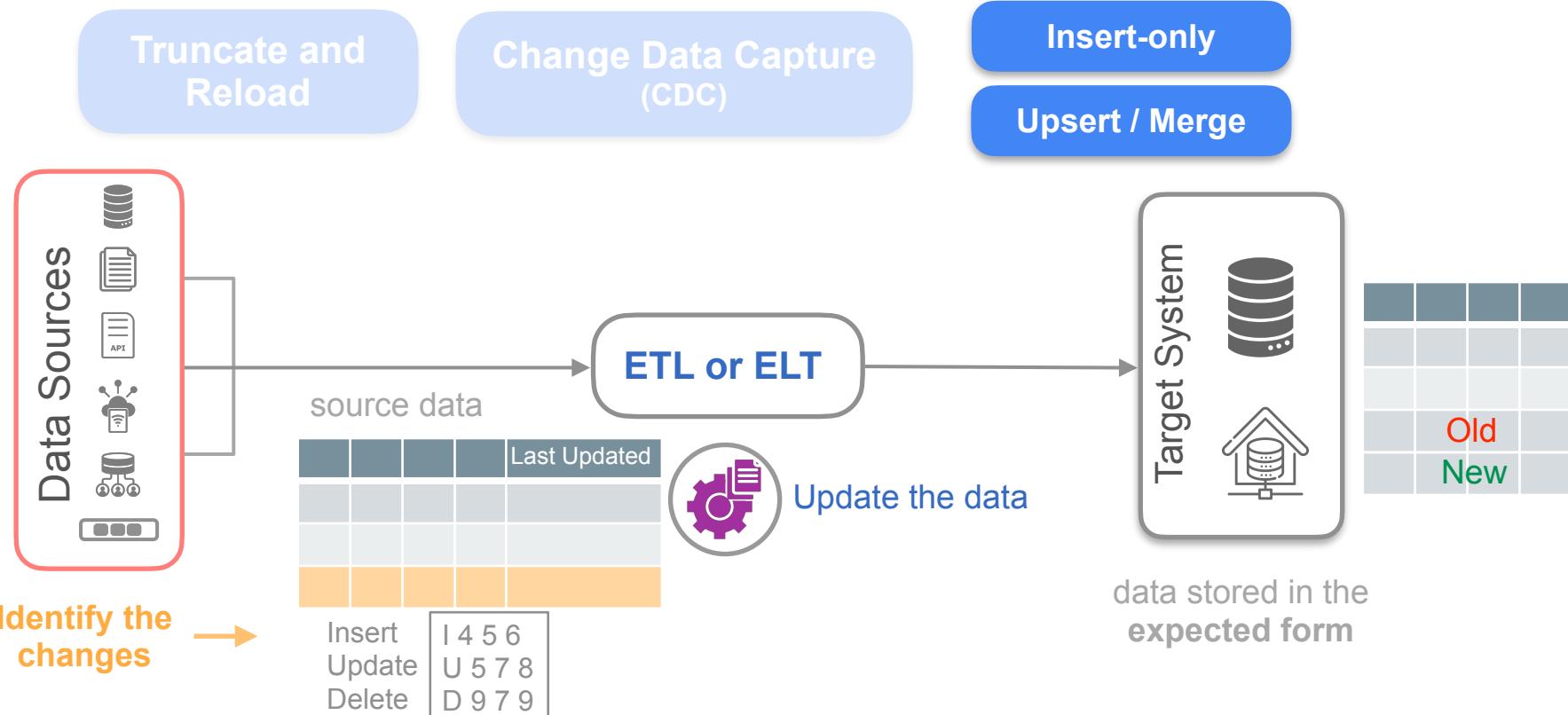
- Delete all records in target system and reload the data from the source
- Have a small dataset
- Update the data once in a while



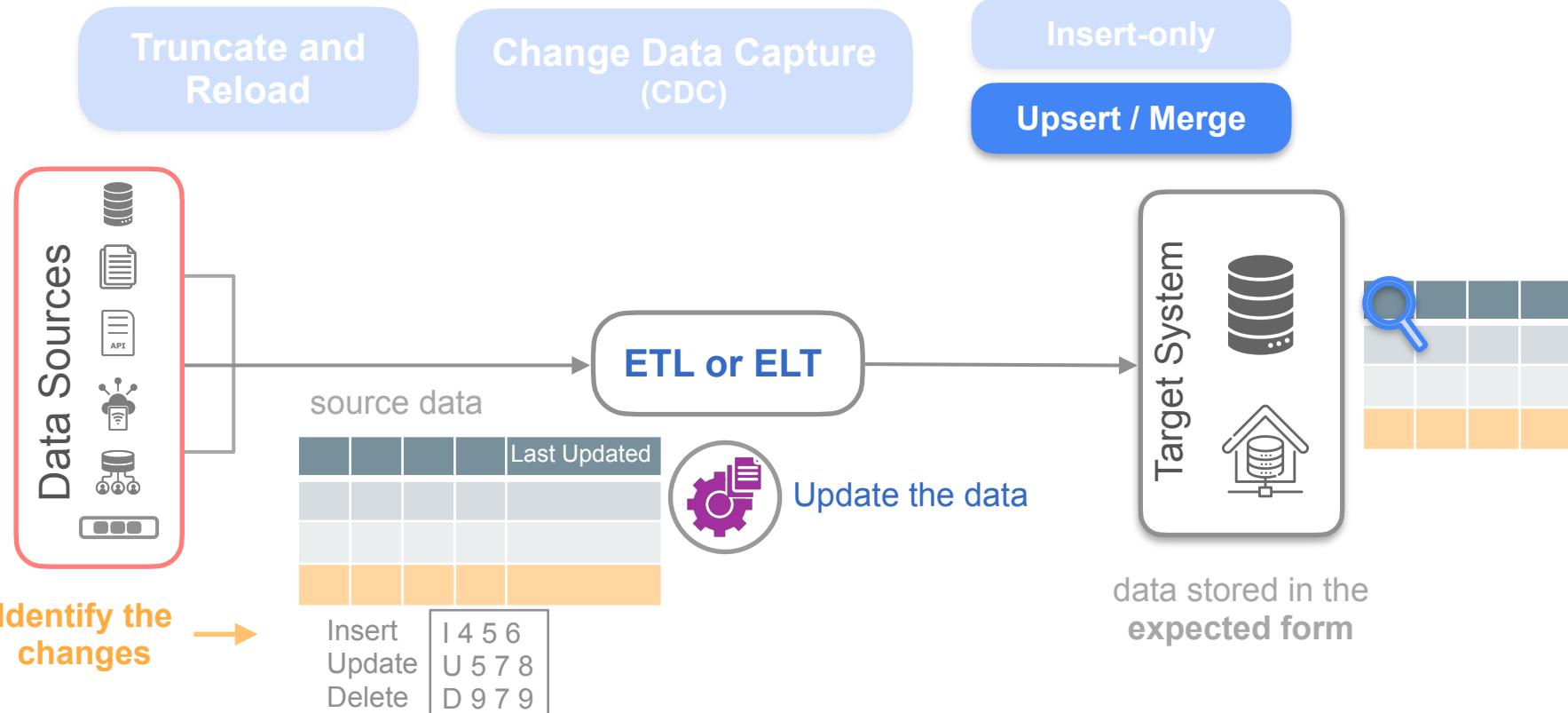
# Transformations for Data Updating



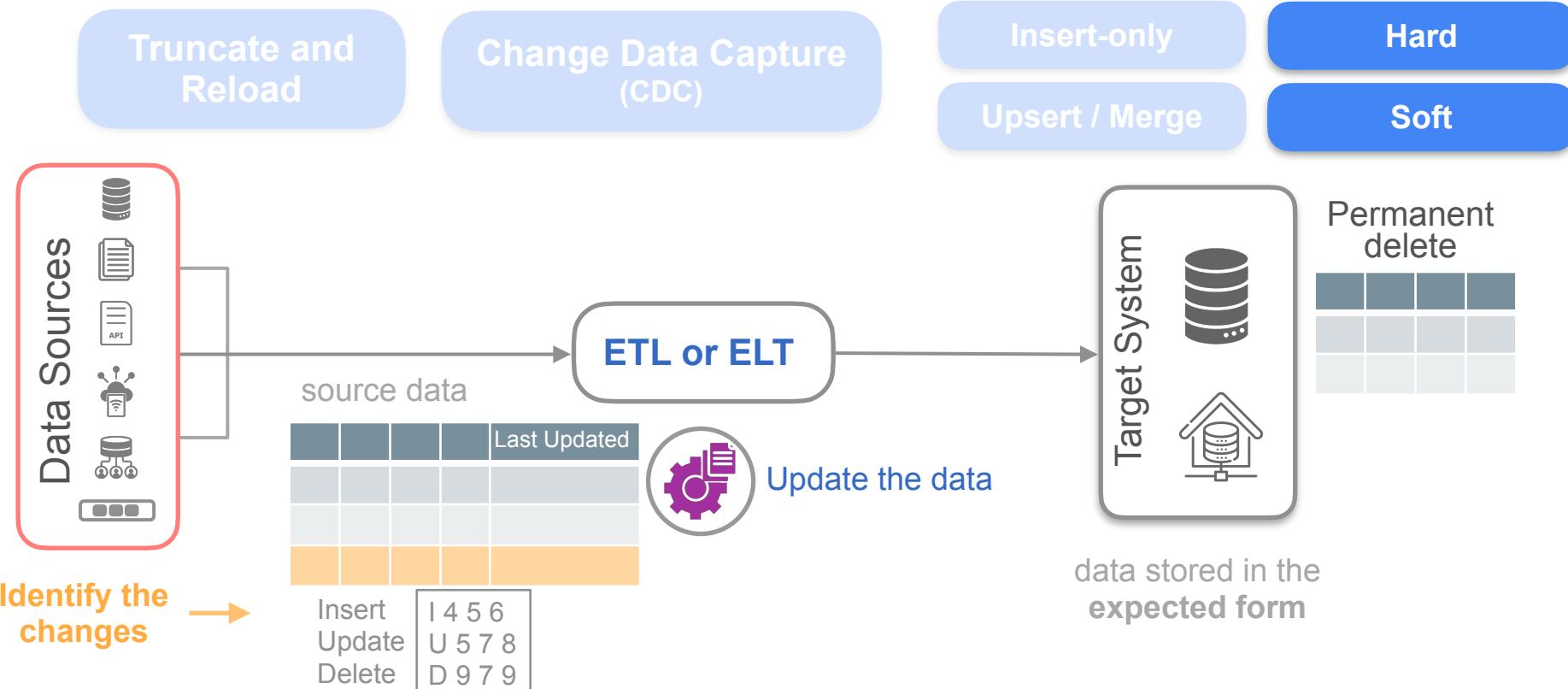
# Transformations for Data Updating



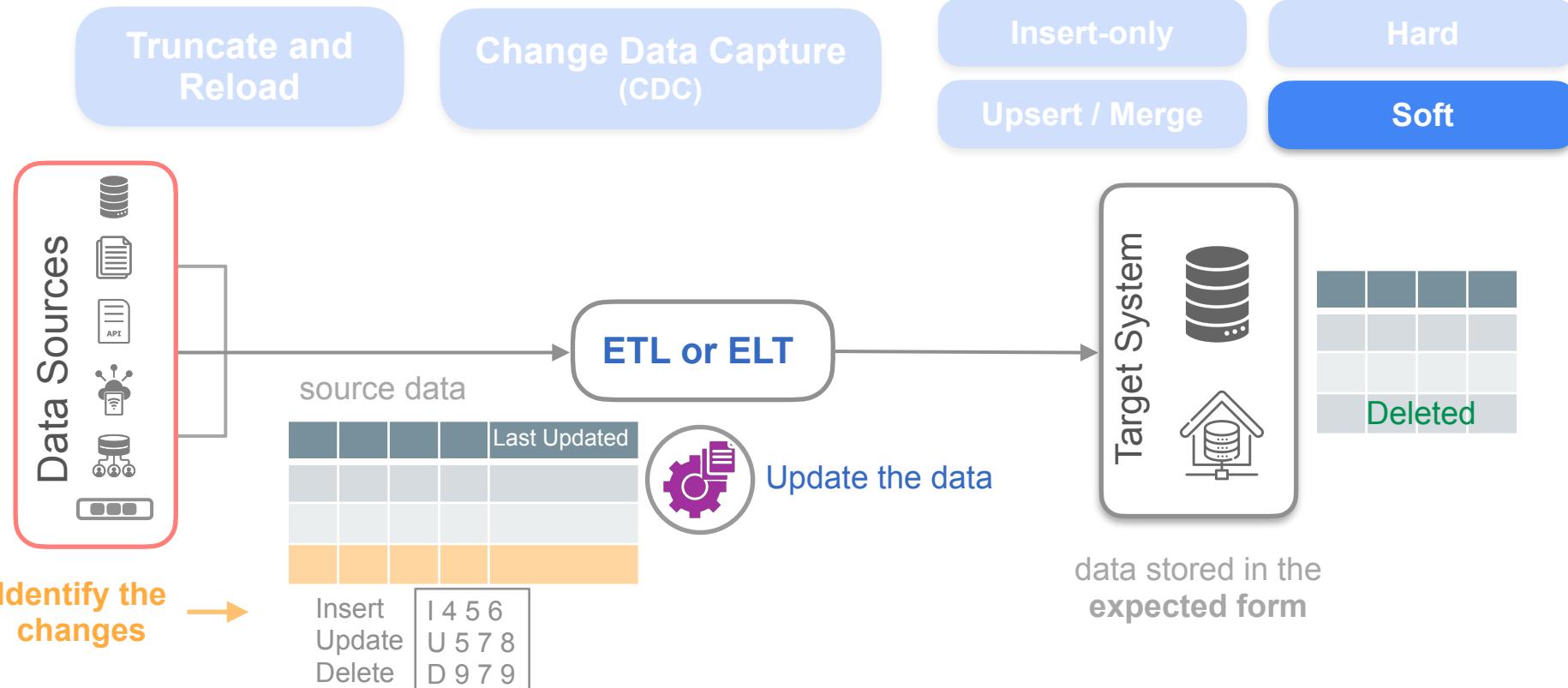
# Transformations for Data Updating



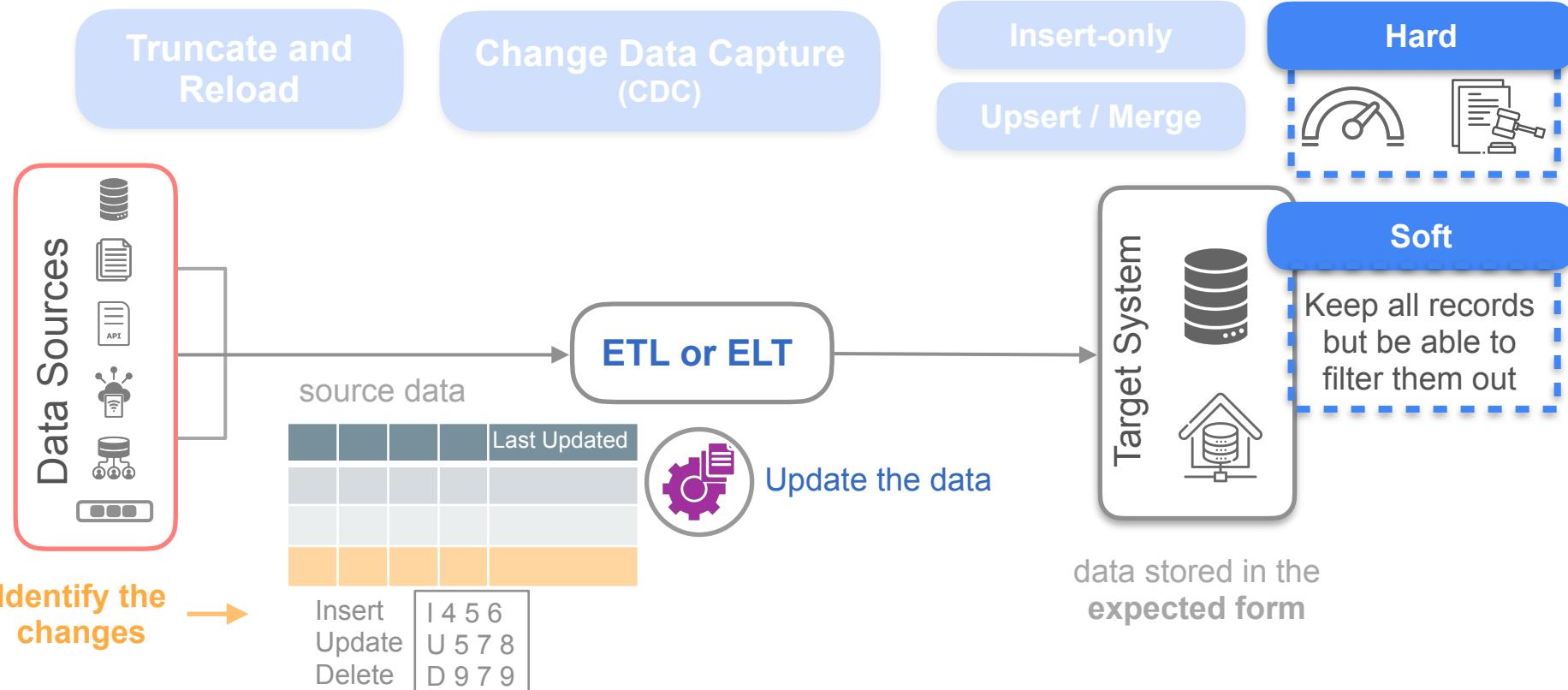
# Transformations for Data Updating



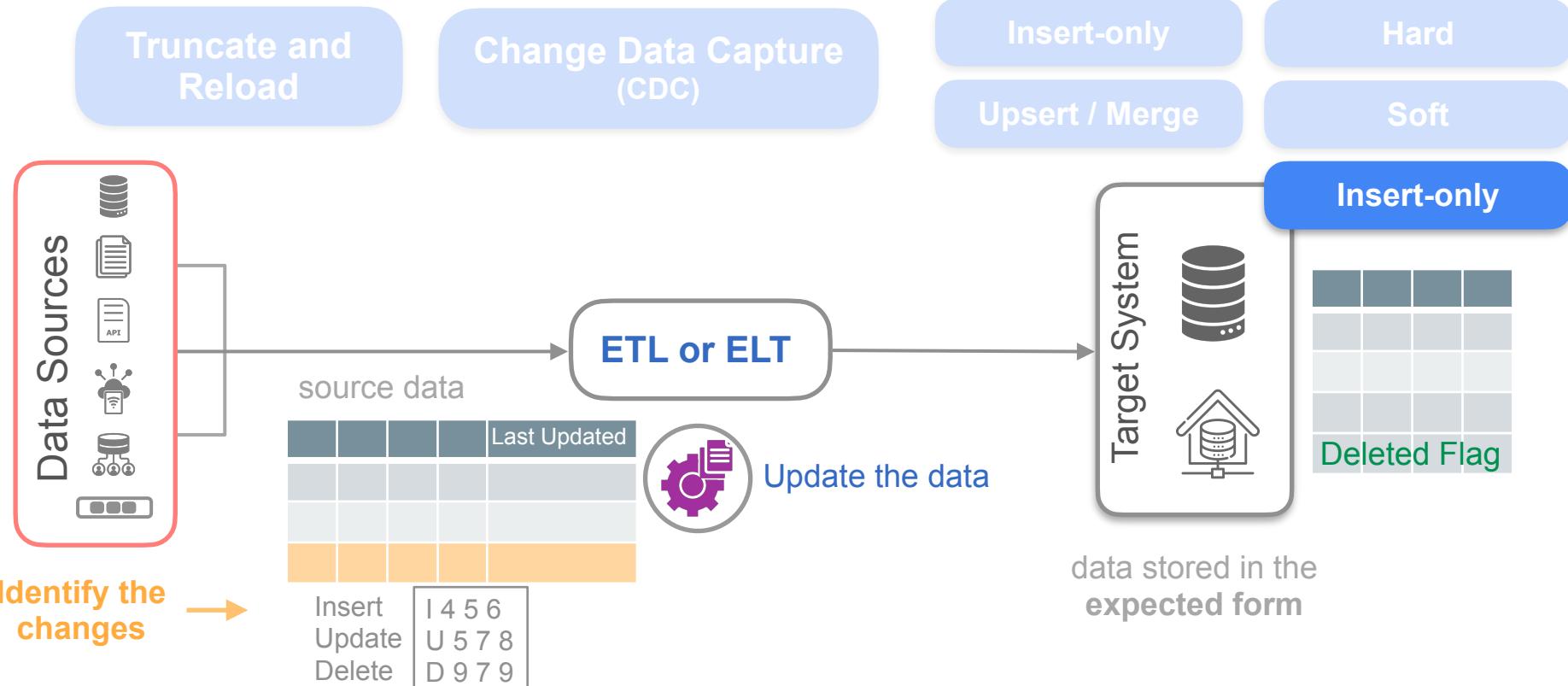
# Transformations for Data Updating



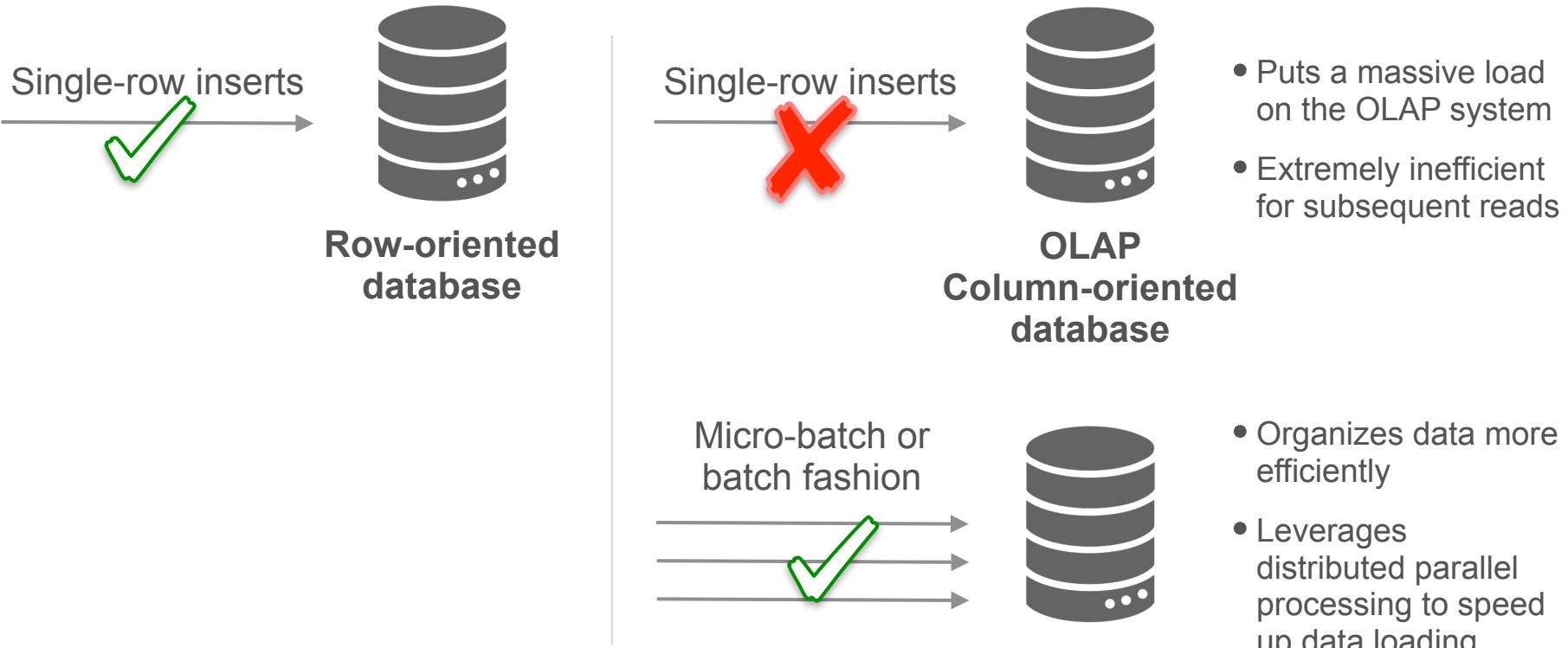
# Transformations for Data Updating



# Transformations for Data Updating



# Transformations for Data Updating





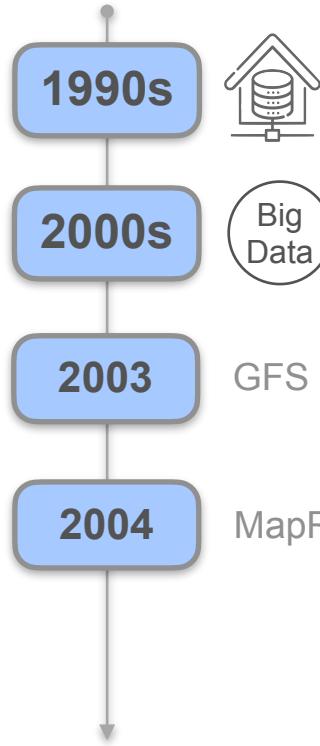
DeepLearning.AI

## Batch Transformations

---

**Distributed Processing  
Framework —  
Hadoop MapReduce**

# Big Data Era



- Traditional monolithic data warehouses were not able to handle massive amounts of data effectively
- Commodity hardware became cheap and ubiquitous
- Several innovations in large-scale distributed storage and computing:

2003  
Google

The Google File System

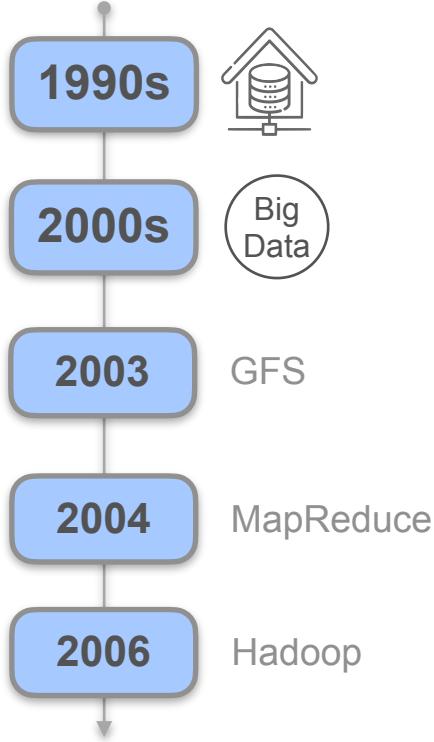
GFS

2004  
Google

MapReduce

Parallel programming paradigm for processing data distributed over GFS

# Big Data Era



2006



- **Hadoop Distributed File System (HDFS)**  
Key ingredients of current big data engines (EMR, Spark)
- **Hadoop MapReduce**  
Influences many of today's distributed systems

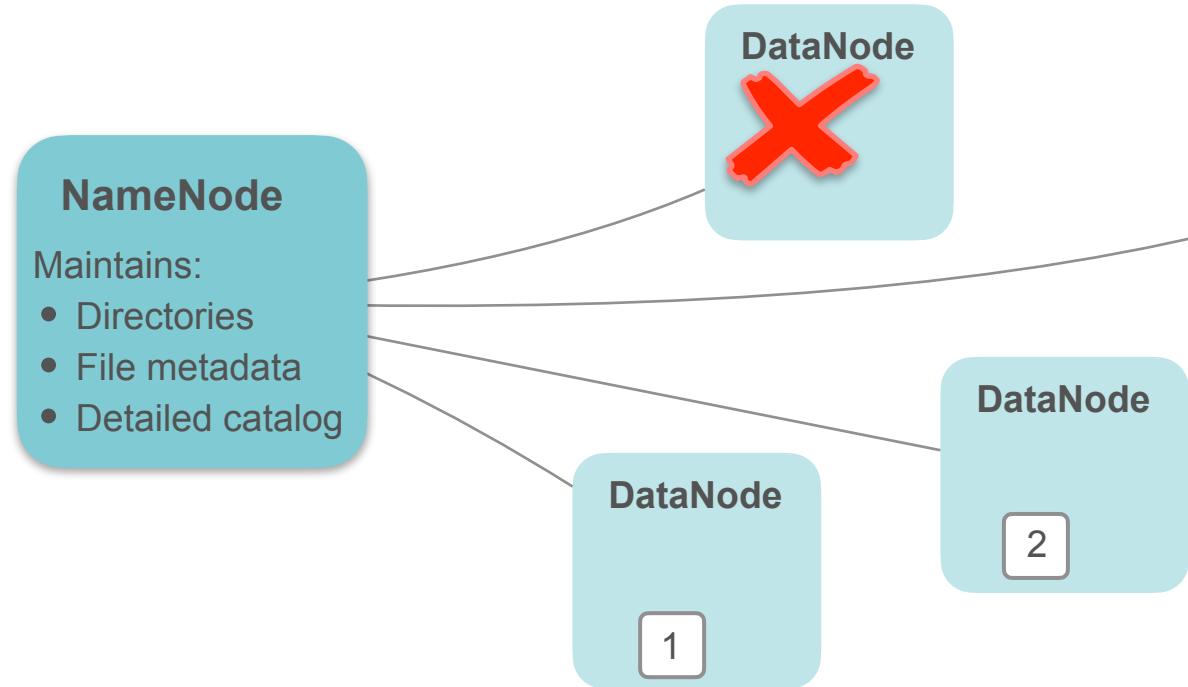
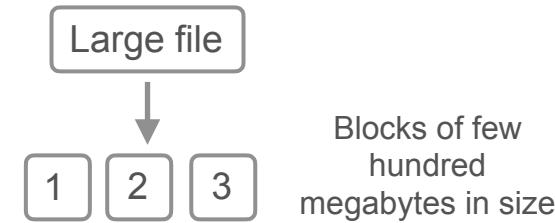
# Hadoop Distributed File System

**HDFS:** combines compute and storage on the same nodes

**Object Storage:** limited compute support for internal processing

# Hadoop Distributed File System

**HDFS:** combines compute and storage on the same nodes



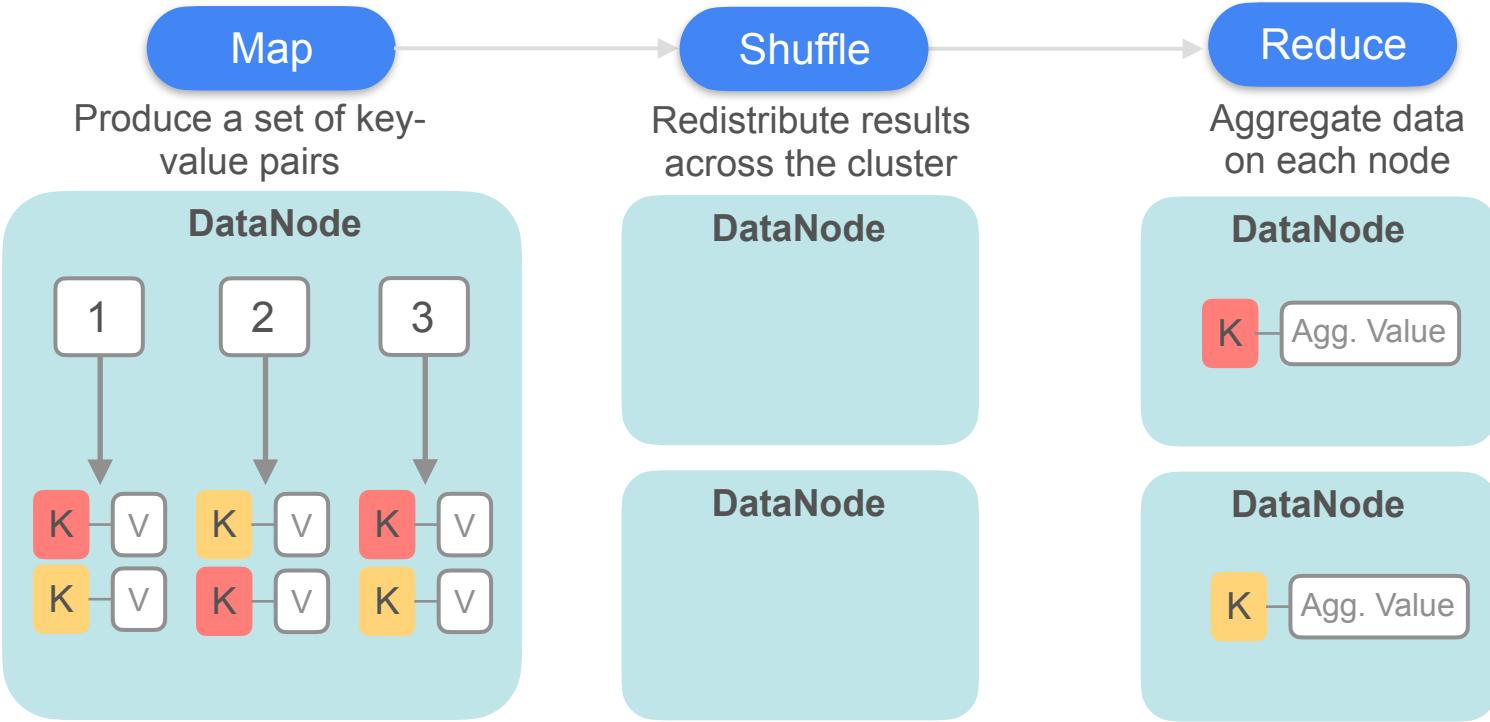
**DataNode**

Combining compute and storage allows in-place data processing

**Replication**  
increases durability and availability of the data

# Hadoop MapReduce

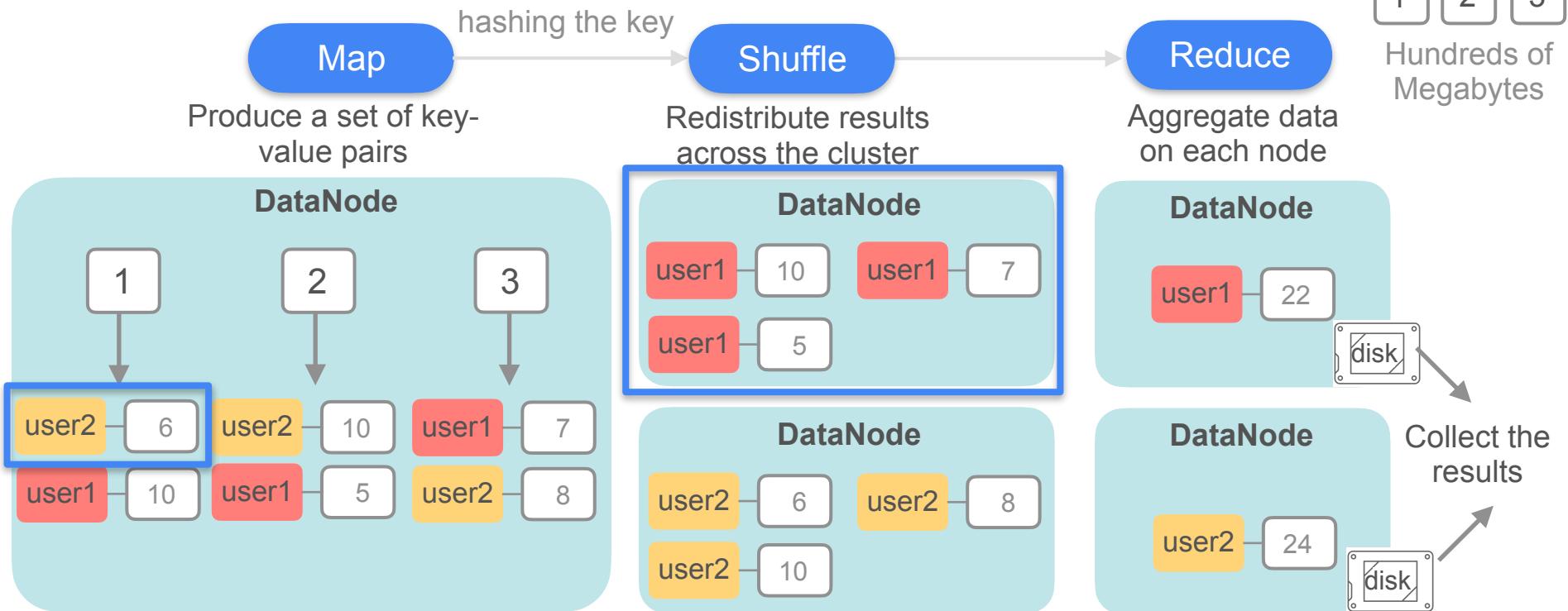
**MapReduce:** send computation code to the nodes that contain the data



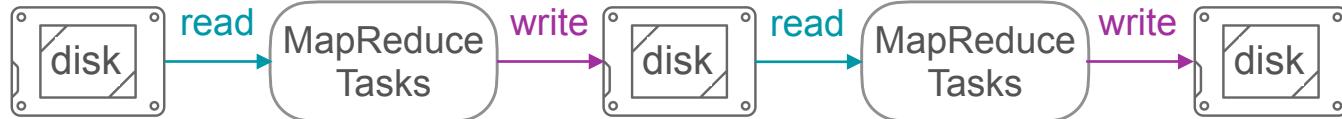
# Hadoop MapReduce

```
SELECT user_id, COUNT(*)  
FROM user_events  
GROUP BY user_id;
```

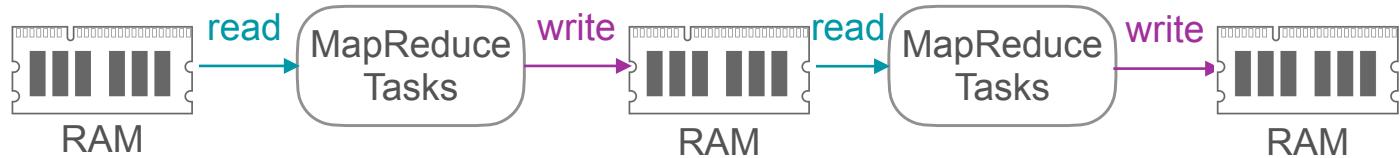
## Petabytes



# Hadoop MapReduce - Shortcomings



Pros	Cons
<ul style="list-style-type: none"><li>• Simplifies state and workflow management</li><li>• Minimizes memory consumption</li></ul>	<ul style="list-style-type: none"><li>• Drives high-disk bandwidth utilization</li><li>• Increases processing time</li></ul>



- In-memory processing speeds up data processing tasks
- Disk is treated as a second-class data storage layer



DeepLearning.AI

## Batch Transformations

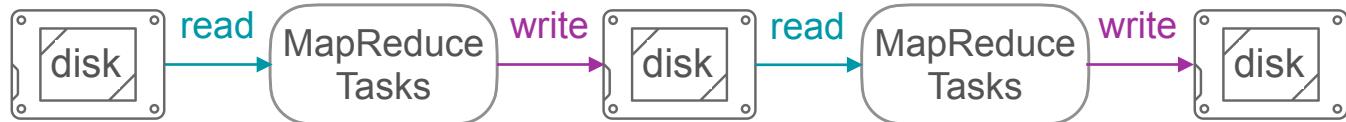
---

**Distributed Processing  
Framework —  
Spark**

# Distributed processing framework

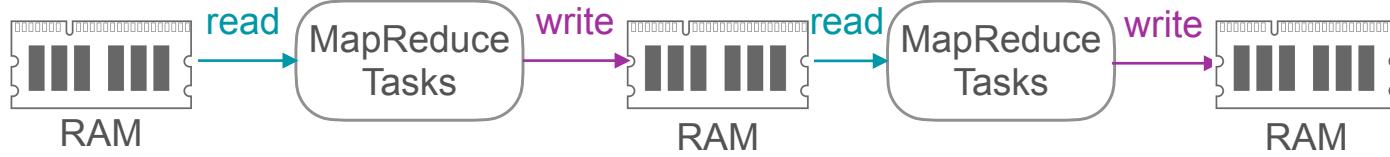


Integrate compute and persistent storage via the HDFS.



Computing engine designed for processing large datasets.

- Retains intermediate results in memory
- Limits disk I/O interactions, enabling faster computation



## Separate persistent storage systems

Locally  
Data centers  
Cloud



# Spark

Unified  
Platform

You can run different types of analytical workloads:

- Perform SQL queries
- Train and test ML algorithms
- Apply streaming transformations

Spark  
SQL

MLlib

Spark  
Streaming

GraphX

Spark core APIs



Python



Java



Scala

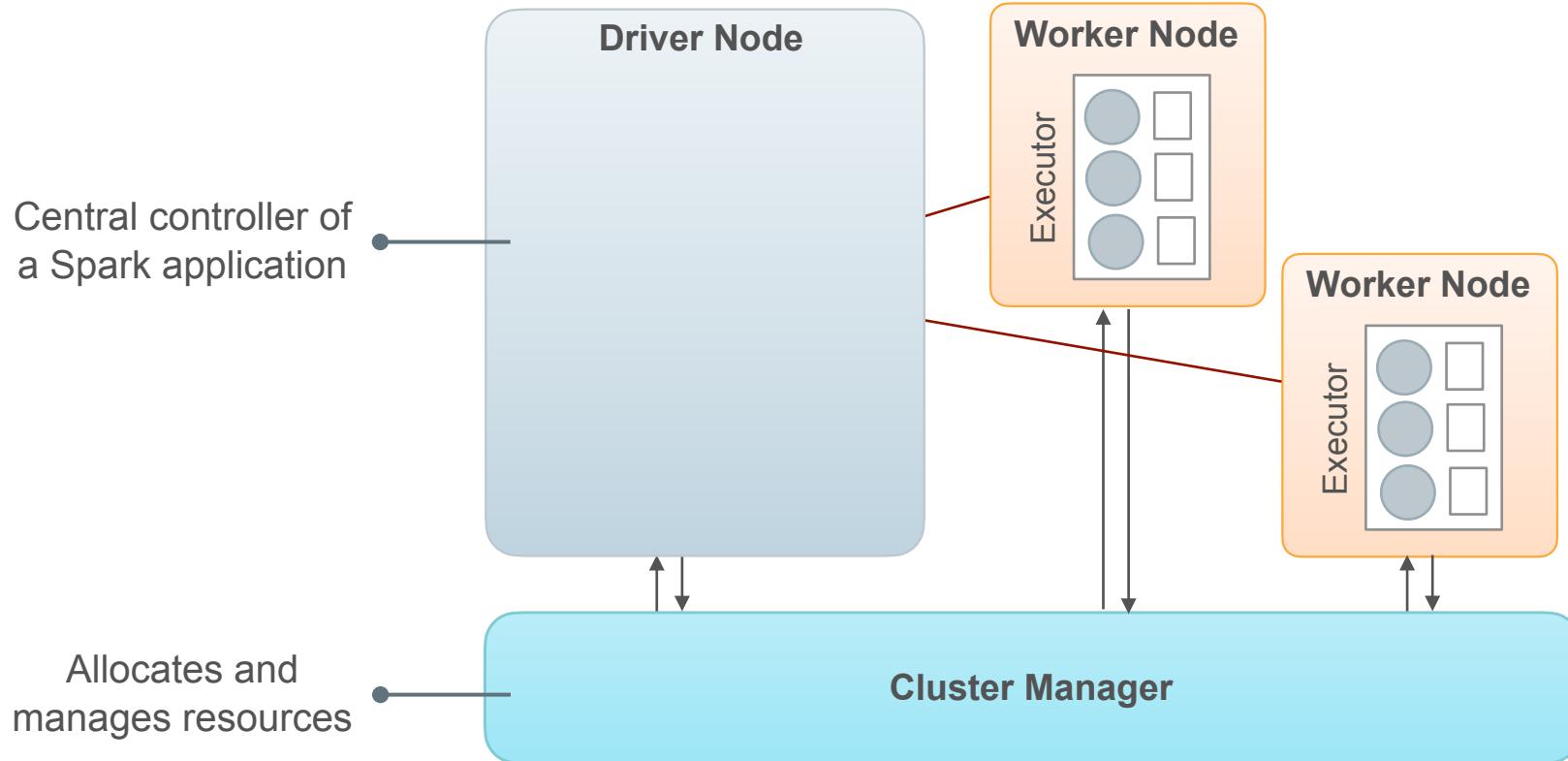


R

External third-party libraries

# Spark Application

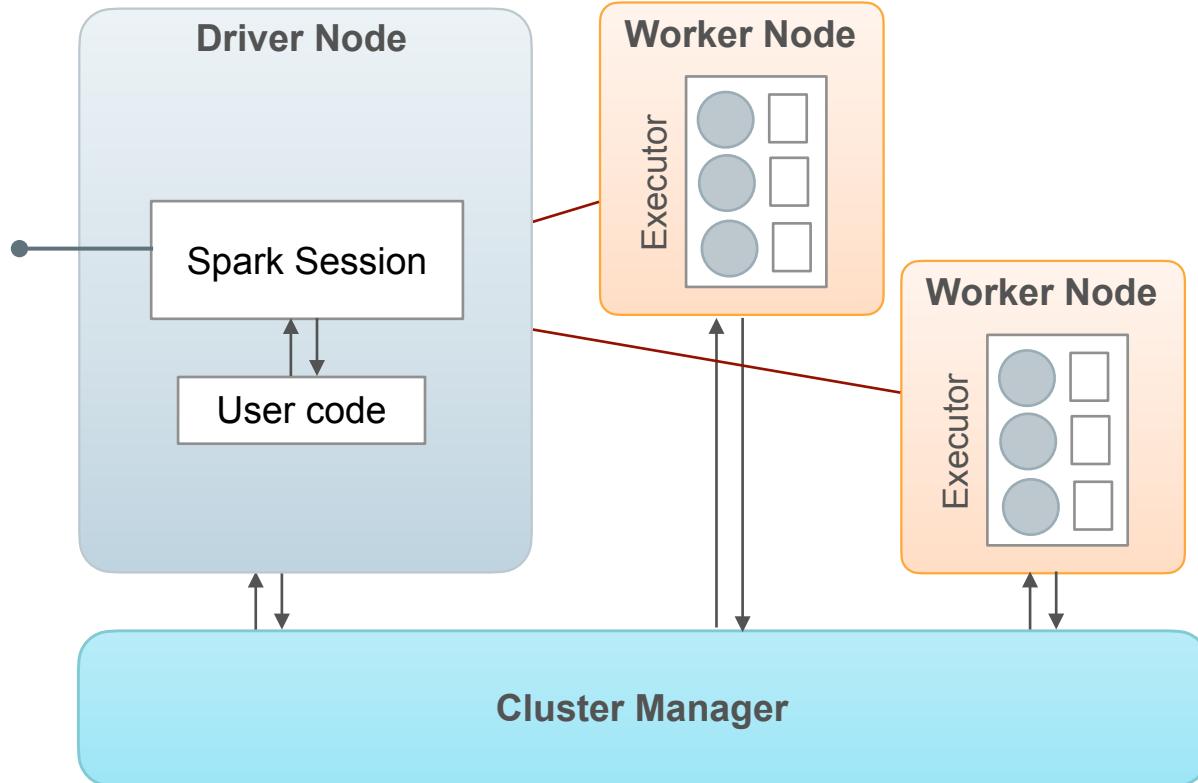
CPU core



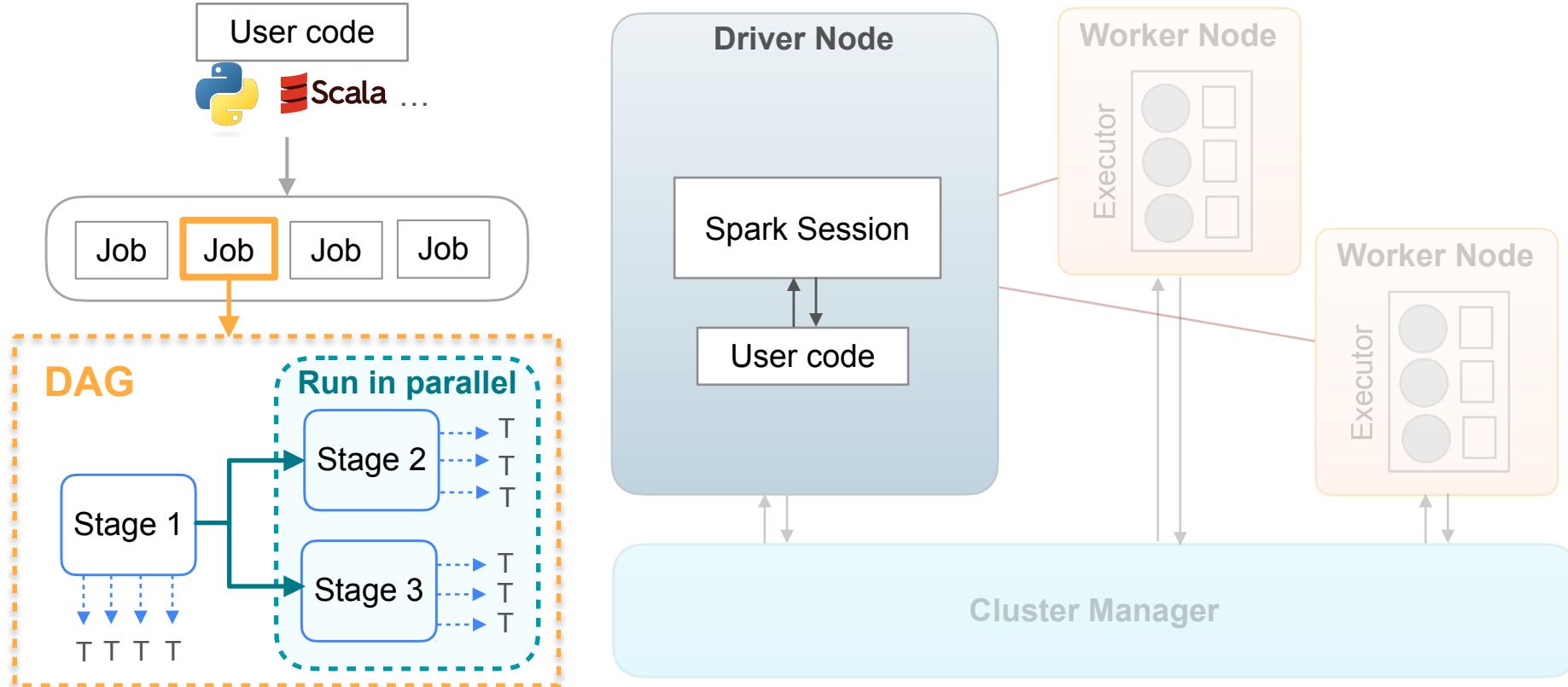
# Spark Application

Single unified entry point to Spark's functionality:

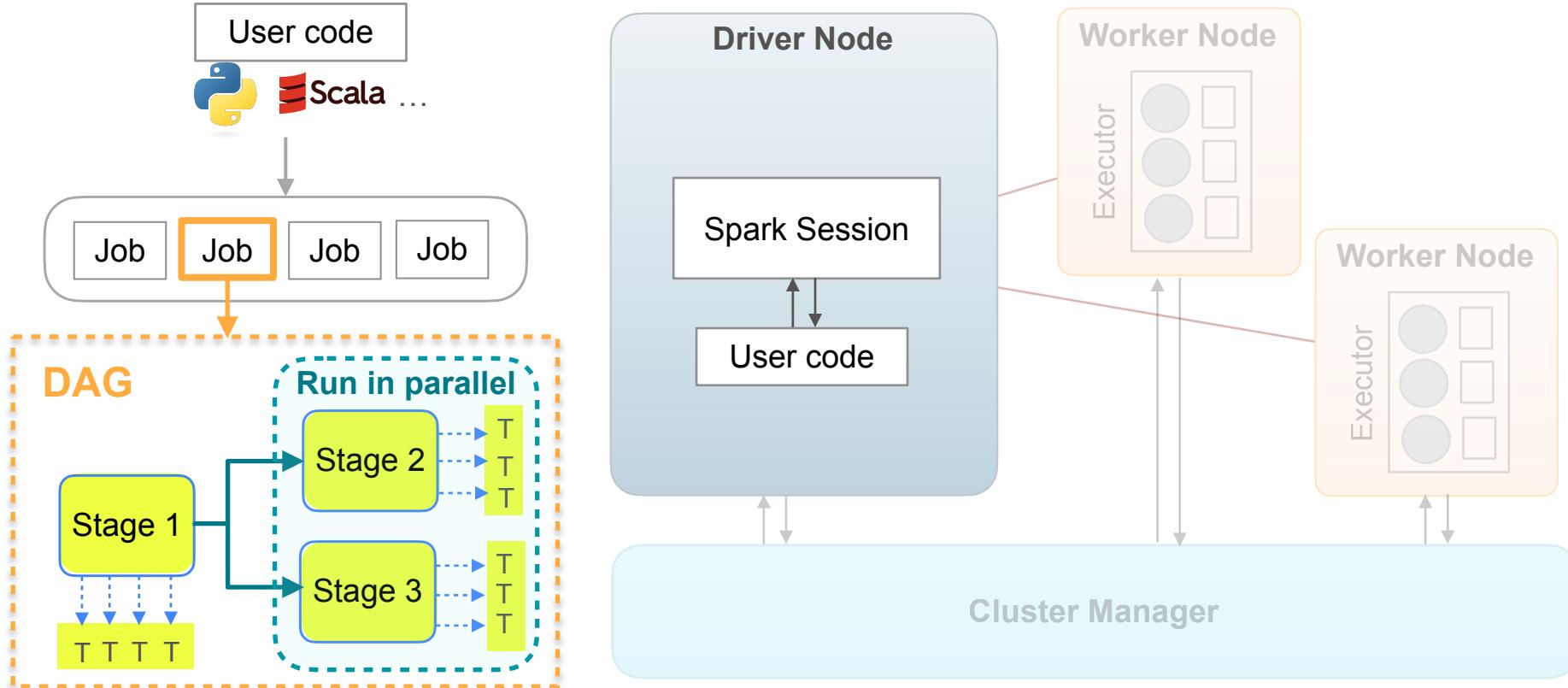
- Define Dataframes
- Read data from sources
- Perform SQL queries



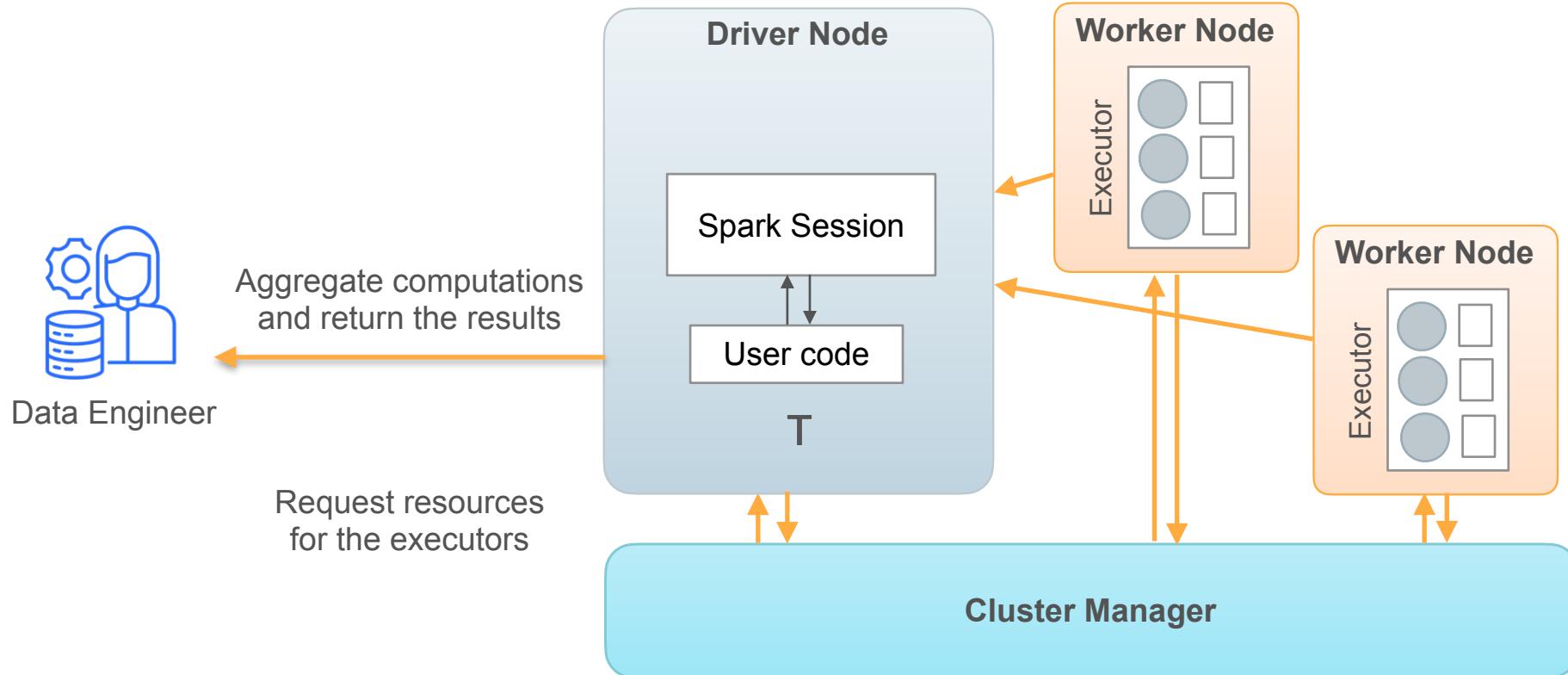
# Spark Application



# Spark Application



# Spark Application





DeepLearning.AI

## Batch Transformations

---

## Spark DataFrames

# Spark Dataframes

## High-level Structured API

Spark DataFrames

Use simpler and more expressive operations:  
Filtering, selecting, counting, aggregating and grouping

## Immutable data structures

## Low-level API

Resilient Distributed Dataset (RDD)  
actual partitioned collection of records

Need to manually define and optimize operations

# Spark Dataframes

## High-level Structured API

Spark DataFrames

Use simpler and more expressive operations:  
Filtering, selecting, counting, aggregating and grouping

### Transformation

Filtering, selecting, joining and grouping

### Action

Counting, showing, and saving

Immutable  
DataFrame

Lazy evaluation

DataFrame

DataFrame

- Transformations: recorded as a lineage, executed when an action is invoked
- Allows Spark to optimize the execution plan
- Lineage & immutability: fault-tolerance of DataFrames



DeepLearning.AI

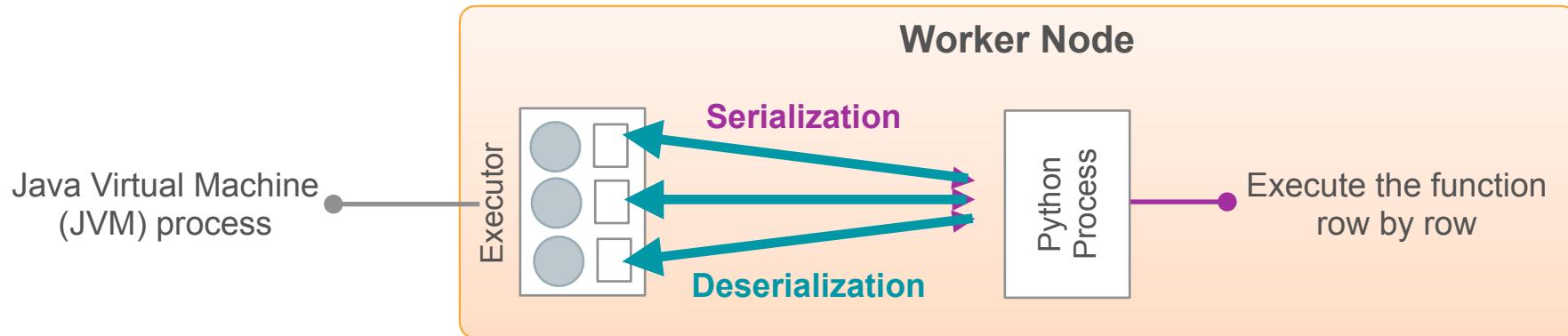
## Batch Transformations

---

**Demo: Working with Spark  
DataFrames Using Python**

# Python UDFs

Spark is native in Scala



- The JVM and Python processes will compete for memory resources.
- Serialization and deserialization are expensive.
- Consider writing UDFs in Scala or Java: they run directly within the JVM.



DeepLearning.AI

## Batch Transformations

---

**Demo: Working with Spark SQL**



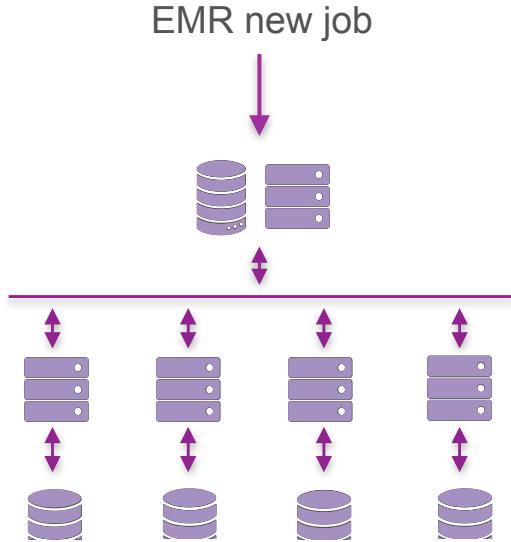
DeepLearning.AI

# Batch Transformations

---

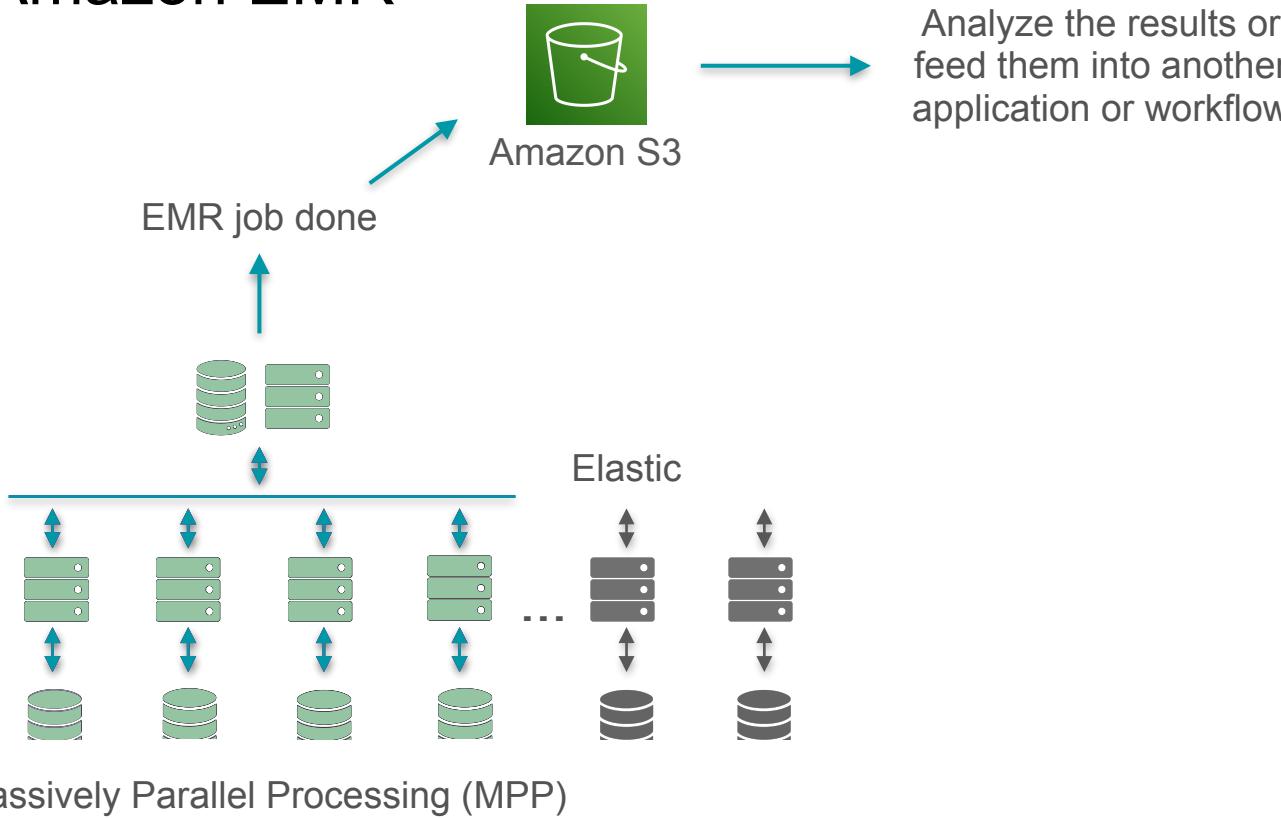
**EMR**

# Amazon EMR

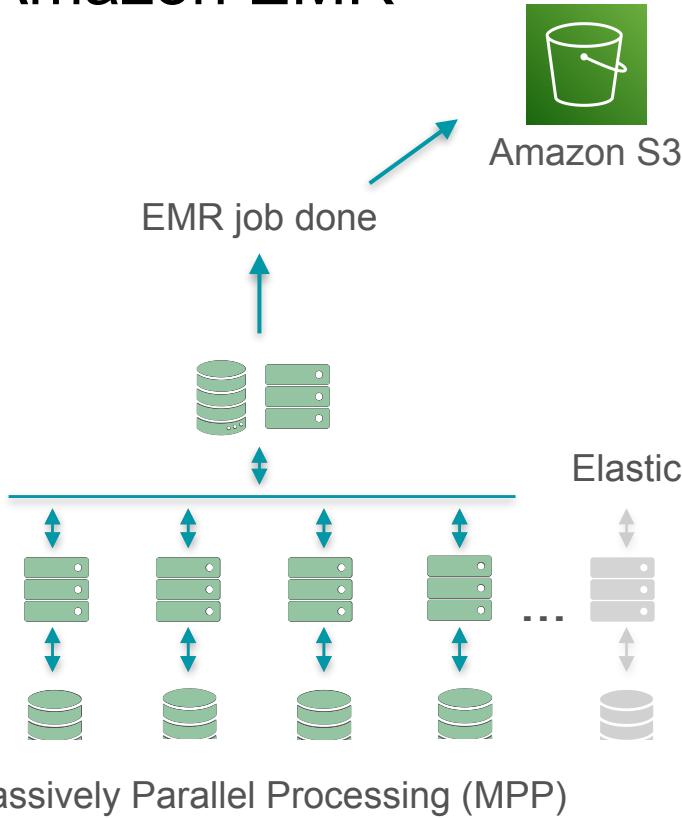


Massively Parallel Processing (MPP)

# Amazon EMR



# Amazon EMR

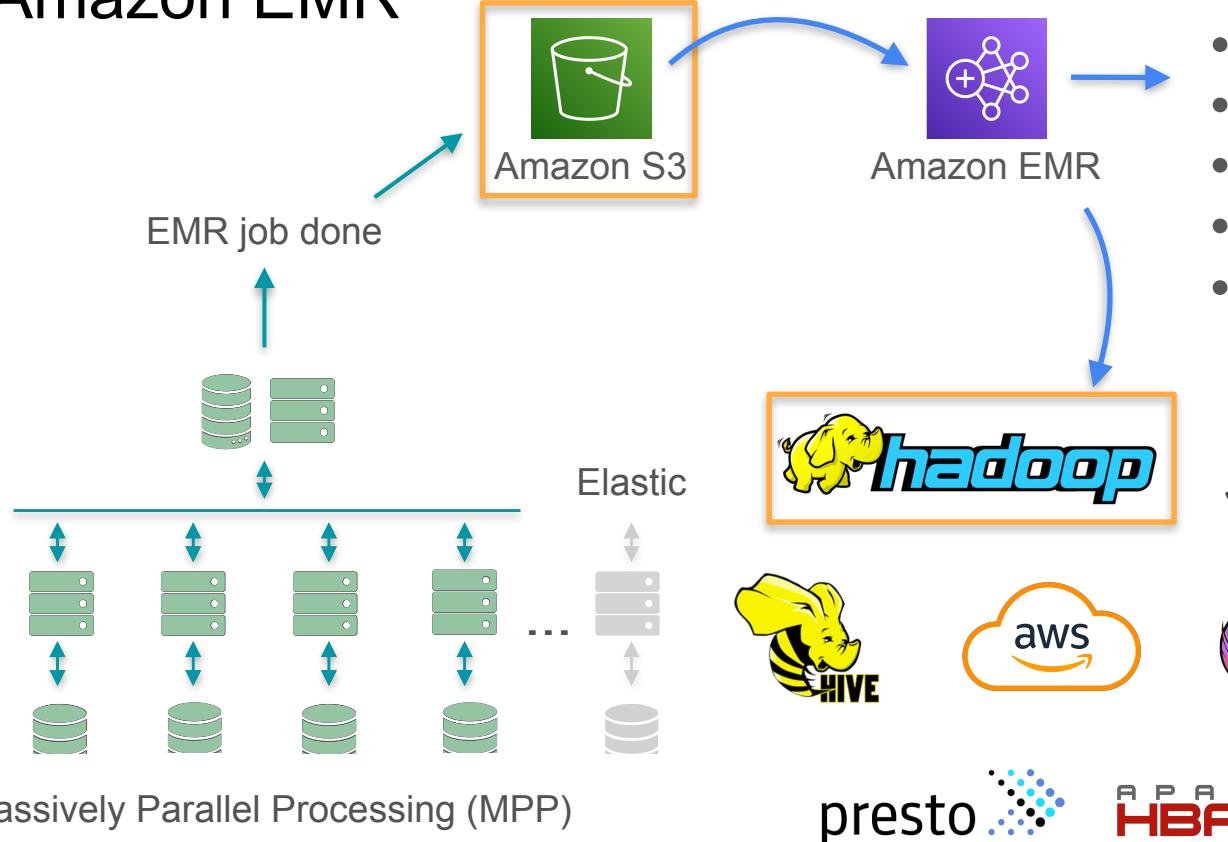


Management Environment

- Simplifies setup and scaling
- Natively integrates with AWS
- Focus on data workflows



# Amazon EMR

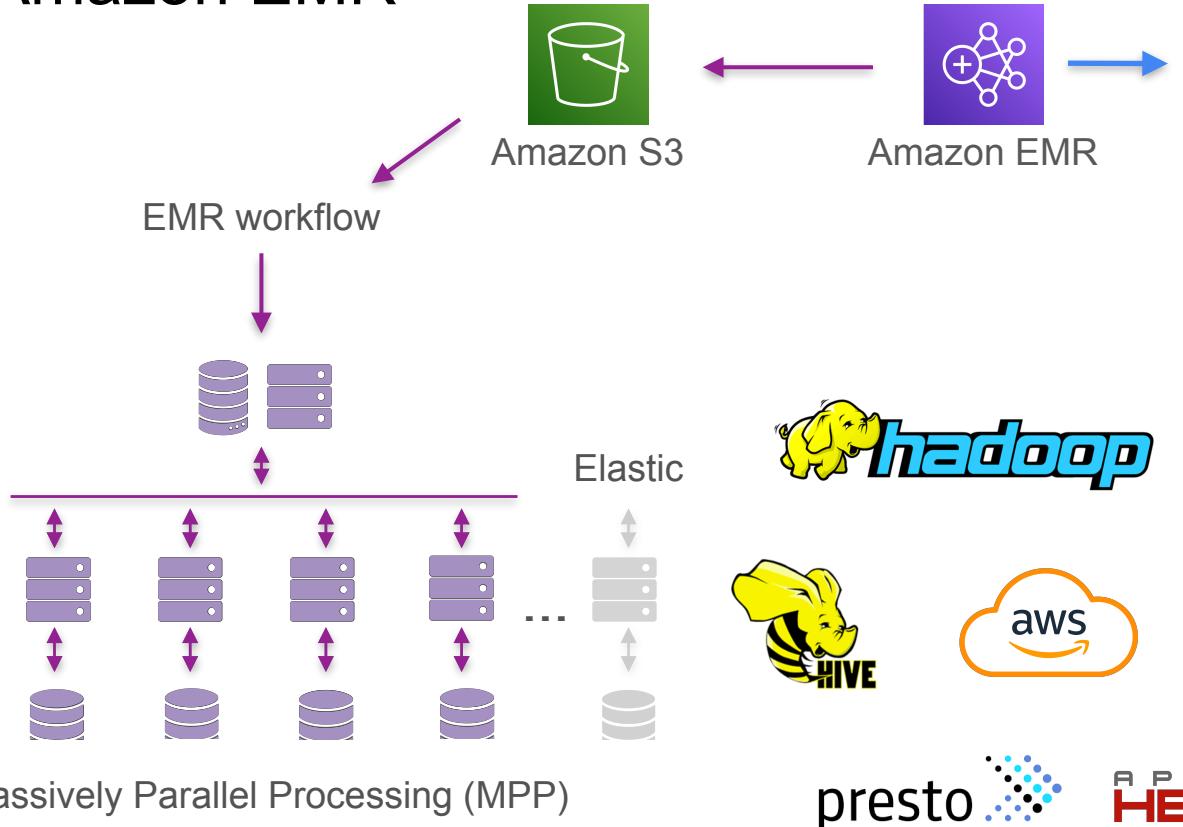


Management Environment

- Simplifies setup and scaling
- Natively integrates with AWS
- Focus on data workflows
- Decouple compute and storage
- Analyze large amount of data



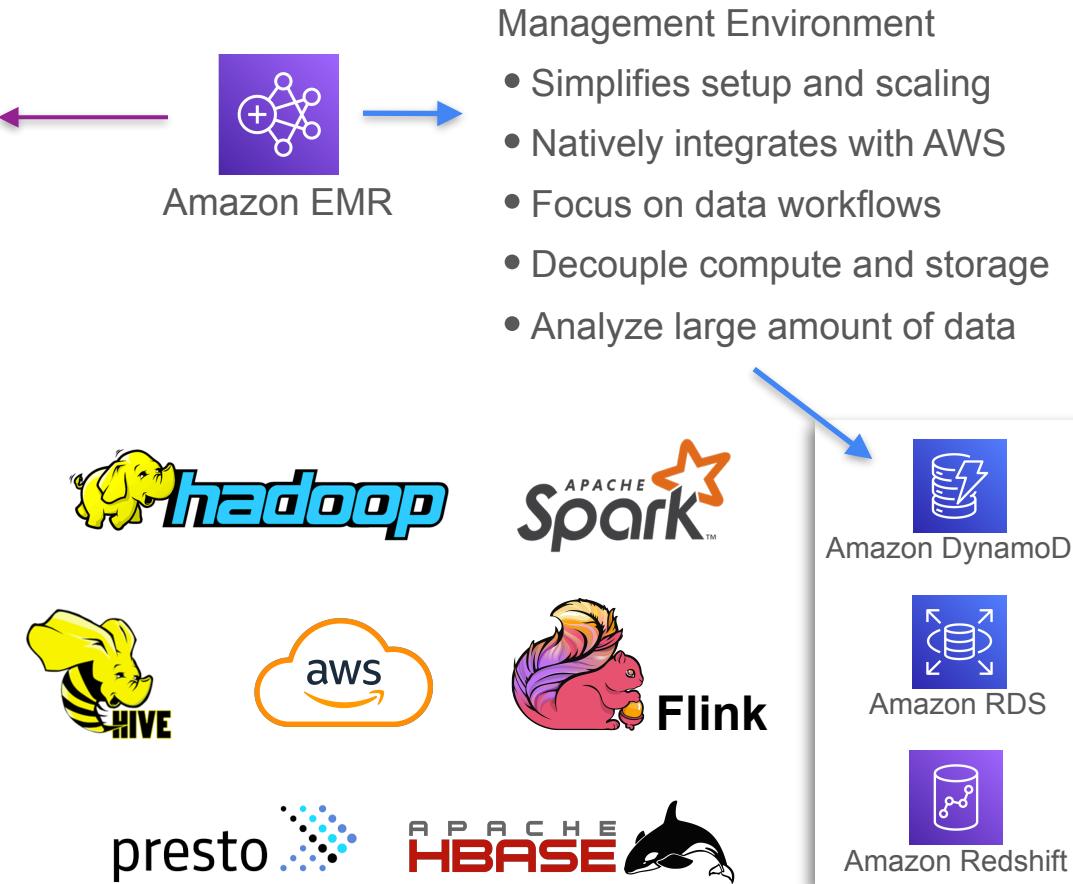
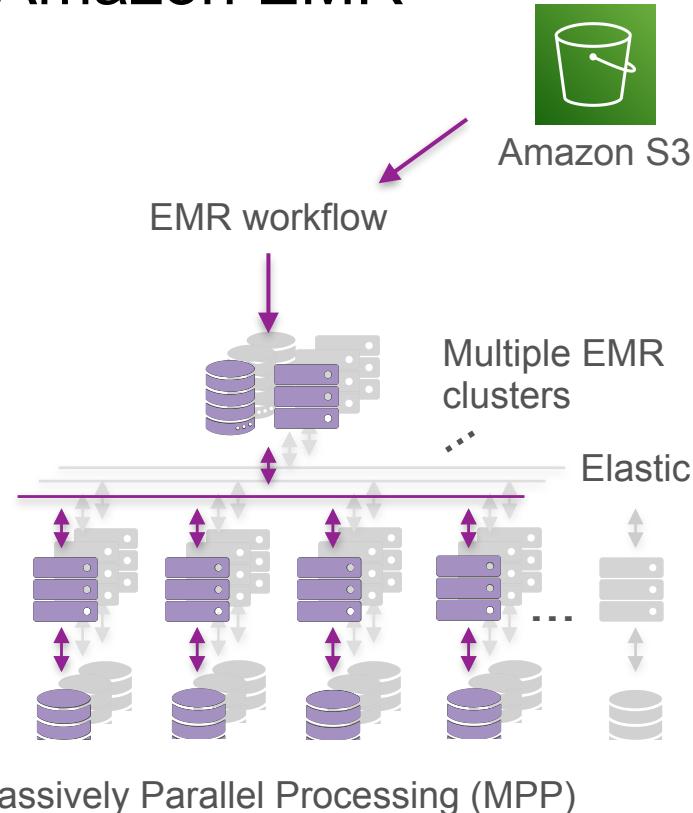
# Amazon EMR



- Management Environment
- Simplifies setup and scaling
  - Natively integrates with AWS
  - Focus on data workflows
  - Decouple compute and storage
  - Analyze large amount of data



# Amazon EMR





DeepLearning.AI

## Batch Transformations

---

## Technical Considerations

Transformation Complexity	Code Reusability and Testability	Team's skills and background
<p><b>Simple transformation</b></p> <p>Python and SQL approaches have comparable performance:</p> <ul style="list-style-type: none"><li>• translated to the same execution plan</li><li>• executed by the same engine</li></ul>		
<p><b>Complex transformation</b></p> <p>Cannot be implemented in SQL / not in a straightforward way</p> <ul style="list-style-type: none"><li>• transpose operation:<ul style="list-style-type: none"><li>• DataFrame: df.T</li><li>• SQL: Not supported</li></ul></li><li>• normalizing and cleaning data:<ul style="list-style-type: none"><li>• Python may require more code than SQL</li></ul></li></ul>		

Transformation Complexity	Code Reusability and Testability	Team's skills and background
<p><b>Simple transformation</b></p> <p>Python and SQL approaches have comparable performance:</p> <ul style="list-style-type: none"> <li>• translated to the same execution plan</li> <li>• executed by the same engine</li> </ul>	<ul style="list-style-type: none"> <li>• DataFrames: more testable, maintainable, modular/reusable code.</li> <li>• SQL: no good notion of reusability for more complex query component</li> </ul>	
<p><b>Complex transformation</b></p> <p>Cannot be implemented in SQL / not in a straightforward way</p> <ul style="list-style-type: none"> <li>• transpose operation: <ul style="list-style-type: none"> <li>• DataFrame: df.T</li> <li>• SQL: Not supported</li> </ul> </li> <li>• normalizing and cleaning data: <ul style="list-style-type: none"> <li>• Python may require more code than SQL</li> </ul> </li> </ul>		

Transformation Complexity	Code Reusability and Testability	Team's skills and background
<p><b>Simple transformation</b></p> <p>Python and SQL approaches have comparable performance:</p> <ul style="list-style-type: none"> <li>• translated to the same execution plan</li> <li>• executed by the same engine</li> </ul>	<ul style="list-style-type: none"> <li>• DataFrames: more testable, maintainable, modular/reusable code.</li> <li>• SQL: no good notion of reusability for more complex query component</li> </ul>	<ul style="list-style-type: none"> <li>• SQL: writing SQL queries might be simpler and easier than working with DataFrames</li> </ul>
<p><b>Complex transformation</b></p> <p>Cannot be implemented in SQL / not in a straightforward way</p> <ul style="list-style-type: none"> <li>• transpose operation:           <ul style="list-style-type: none"> <li>• DataFrame: df.T</li> <li>• SQL: Not supported</li> </ul> </li> <li>• normalizing and cleaning data:           <ul style="list-style-type: none"> <li>• Python may require more code than SQL</li> </ul> </li> </ul>		

# When to use Spark DataFrames?



**Distributed Framework:** load the entire data into a cluster of nodes

*Use Spark:* if data doesn't fit entirely into memory or you want to leverage distributed computations

**Non-distributed Framework:** load the entire data into the memory

*Use pandas:* if data can fit into the memory

## Best Practices

- Extract only the data you need from the source
- Apply transformation inside the source database to reduce the size of the ingested data



DeepLearning.AI

# Streaming Transformations

---

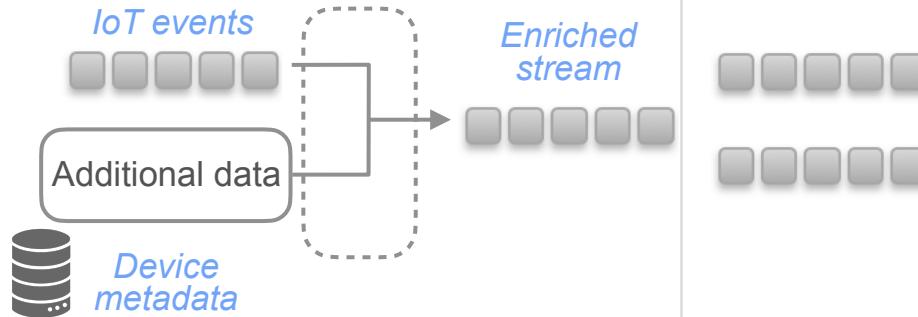
## Stream Processing

## Streaming Transformations

Prepares data for downstream consumption by converting a stream of events into another stream

### Enrich a Stream

Processor



### Join Two Streams



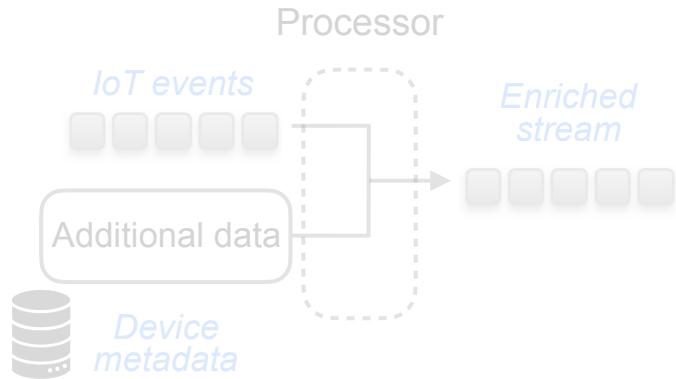
## C2W2 Streaming Ingestion Lab



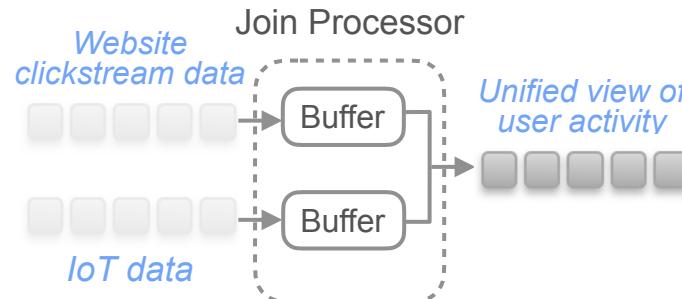
## Streaming Transformations

Prepares data for downstream consumption by converting a stream of events into another stream

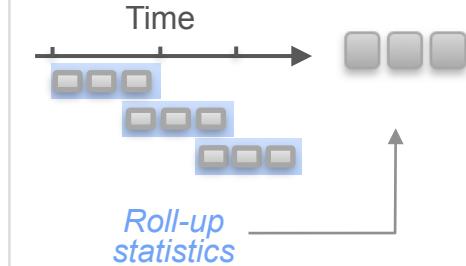
### Enrich a Stream



### Join Two Streams



### Windowed queries



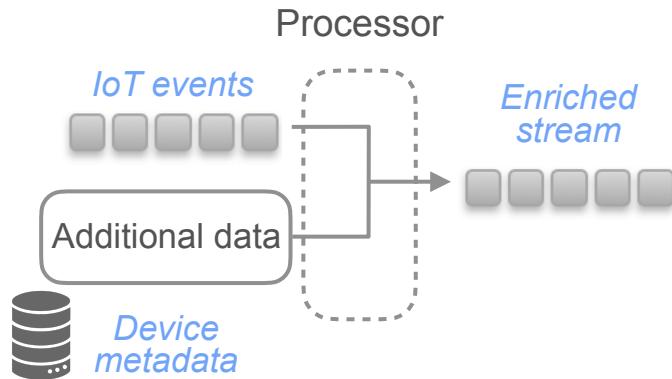
## C2W2 Streaming Ingestion Lab



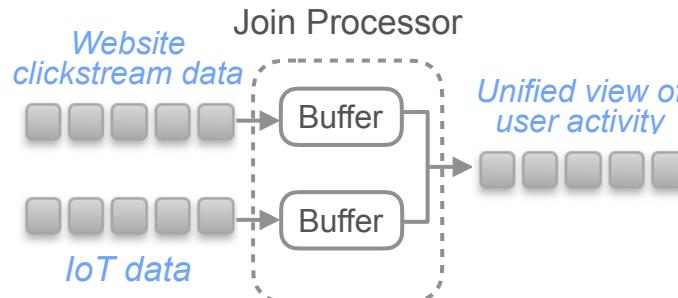
## Streaming Transformations

Prepares data for downstream consumption by converting a stream of events into another stream

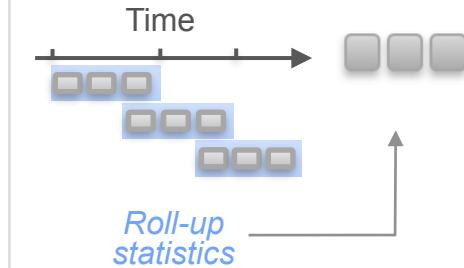
### Enrich a Stream



### Join Two Streams



### Windowed queries



### Streaming Platform



Amazon Kinesis  
Data Streams



### Stream processor



Apache Flink

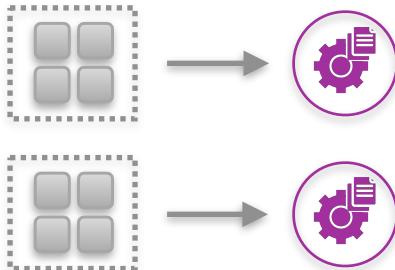
- Distributed stream processing tools
- Open-source
- Python or SQL queries

# Considerations for Choosing a Stream Processing Tool

- Your use case
- Latency requirements
- Performance capabilities of the stream processing framework

## Microbatch stream processing

Process each batch in parallel



Up to 2 minutes

## Spark Streaming

### Example:



Sales metrics

## True stream processing



Apache Flink

- Each node continuously listens to messages and updates its dependent nodes
- Delivers processed events at very low latency
- Comes with significant overhead

### Example:



Security metrics



DeepLearning.AI

# Data Transformations

---

## Week 3 Summary

# Batch Transformations

## Considerations:

- Size of the data you need to transform
- Memory specification of the machine processing the code

## Distributed processing frameworks

- Process large data that doesn't fit into memory
- Leverages parallel computation



# Batch Transformations

## Considerations:

- Size of the data you need to transform
- Memory specification of the machine processing the code

## Distributed processing frameworks

- Process large data that doesn't fit into memory
- Leverages parallel computation



**hadoop**  
MapReduce



complex transformations

- Spark DataFrames
- Spark SQL

simple transformations

# Batch Transformations

## Considerations:

- Size of the data you need to transform
- Memory specification of the machine processing the code

## Distributed processing frameworks

- Process large data that doesn't fit into memory
- Leverages parallel computation



complex transformations

- Spark DataFrames
- Spark SQL

simple transformations

# Streaming Transformations

Understand your use case and performance requirements



# Batch Transformations

## Considerations:

- Size of the data you need to transform
- Memory specification of the machine processing the code

## Distributed processing frameworks

- Process large data that doesn't fit into memory
- Leverages parallel computation



# Streaming Transformations

Understand your use case and performance requirements



## Lab

Implement a CDC pipeline

