

# CS4225/CS5425 Big Data Systems for Data Science

## Spark II: Advanced Topics

Ai Xin  
School of Computing  
National University of Singapore  
[aixin@comp.nus.edu.sg](mailto:aixin@comp.nus.edu.sg)



# Mid-term Test Instructions

- Scope: all the lectures **before recess week**.
- Held in person; open book + notes, but no electronics usage
- Time: **13:45-15:30pm March 18**, Saturday (actual paper time: 1 hour and 15 minutes).
  - **Students are expected to be seated by 13:45.**
  - **We will start the test at 2pm sharp.**
  - **You are NOT allowed to take the test if you come after 2:30pm.**
- For both Grp L1 & L2.
- Venue: UTOWN AUDITORIUM 1/2.
- Seating plan: Canvas > CS4225/CS5425> Files > MidtermMatters.

# Recap: Demo\_I

```
1 df1 = spark.range(2, 100000000, 2)
2 df2 = spark.range(2, 100000000, 4)
3 df3 = df1.join(df2, ["id"])
4 df3.count()
```

## ▼ (4) Spark Jobs

### ▼ Job 0 [View](#) (Stages: 1/1)

Stage 1: 8/8 ⓘ

### ▼ Job 1 [View](#) (Stages: 1/1)

Stage 0: 8/8 ⓘ

### ▼ Job 2 [View](#) (Stages: 1/1, 2 skipped)

Stage 2: 0/8 ⓘ skipped

Stage 3: 0/8 ⓘ skipped

Stage 4: 8/8 ⓘ

### ▼ Job 3 [View](#) (Stages: 1/1, 3 skipped)


Stage 5: 0/8 ⓘ skipped

Stage 6: 0/8 ⓘ skipped

Stage 7: 0/8 ⓘ skipped

Stage 8: 1/1 ⓘ

▶  df1: pyspark.sql.dataframe.DataFrame = [id: long]

▶  df2: pyspark.sql.dataframe.DataFrame = [id: long]

▶  df3: pyspark.sql.dataframe.DataFrame = [id: long]

Out[1]: 2500000

```
1 df1.show(10)
```

## ▶ (1) Spark Jobs

+---+

| id|

+---+

| 2|

| 4|

| 6|

| 8|

| 10|

| 12|

| 14|

| 16|

| 18|

| 20|

+---+

```
1 df2.show(10)
```

## ▶ (1) Spark Jobs

+---+

| id|

+---+

| 2|

| 6|

| 10|

| 14|

| 18|

| 22|

| 26|

| 30|

| 34|

| 38|

+---+

```
1 df3.show(10)
```

## ▶ (3) Spark Jobs

+---+

| id|

+---+

| 22|

| 26|

| 34|

| 50|

| 54|

| 94|

| 110|

| 126|

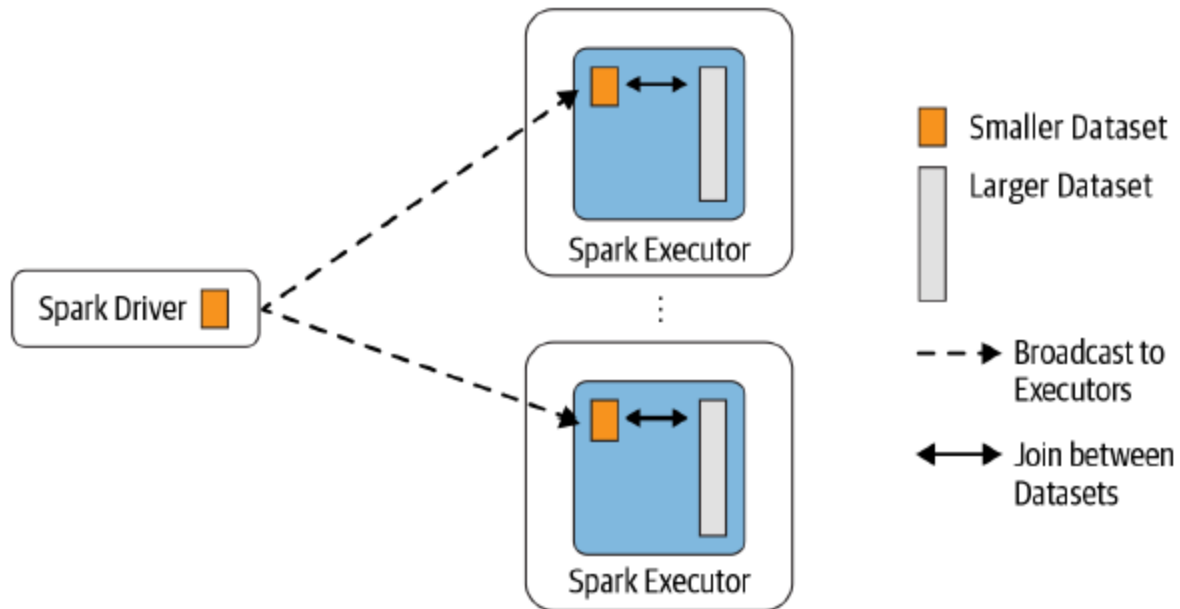
| 130|

| 190|

+---+

# Spark Join

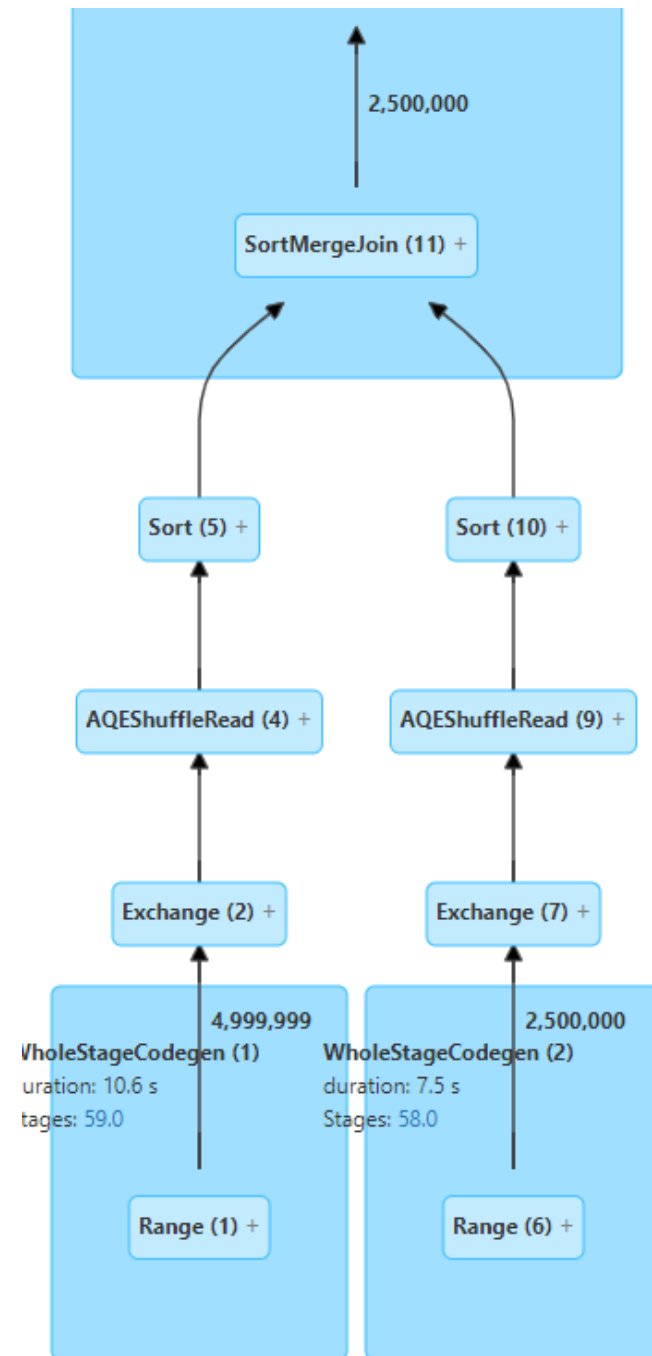
- Broadcast Hash Join (a.k.a. map-side-only join)
  - the smaller data set is broadcast to all executors



# Spark Join

- Shuffle Sort Merge Join
  - an efficient way to merge two large data sets over a common key that is sortable, unique, and can be assigned to or stored in the same partition
  - all rows within each data set with the same key are hashed on the same partition on the same executor

```
df1 = spark.range(2, 10000000, 2)
df2 = spark.range(2, 10000000, 4)
df3 = df1.join(df2, ["id"])
df3.count()
```

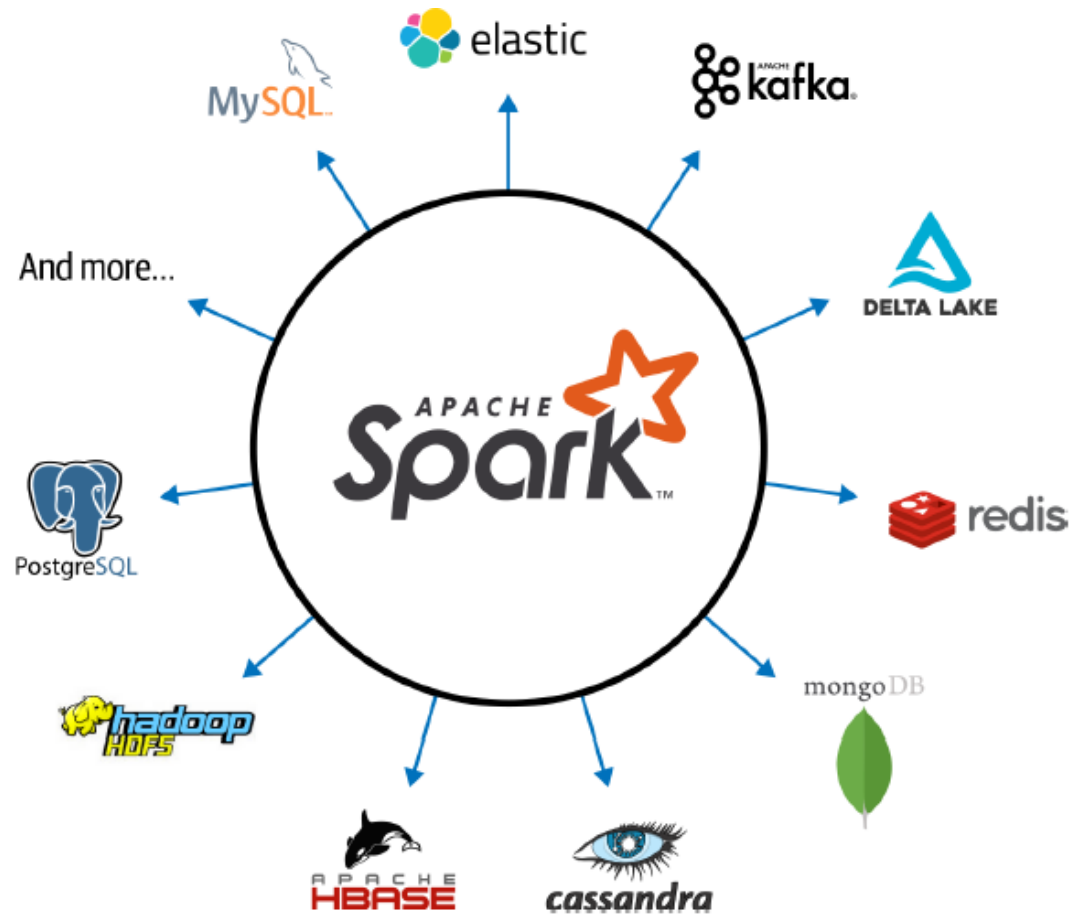


# Today's Plan

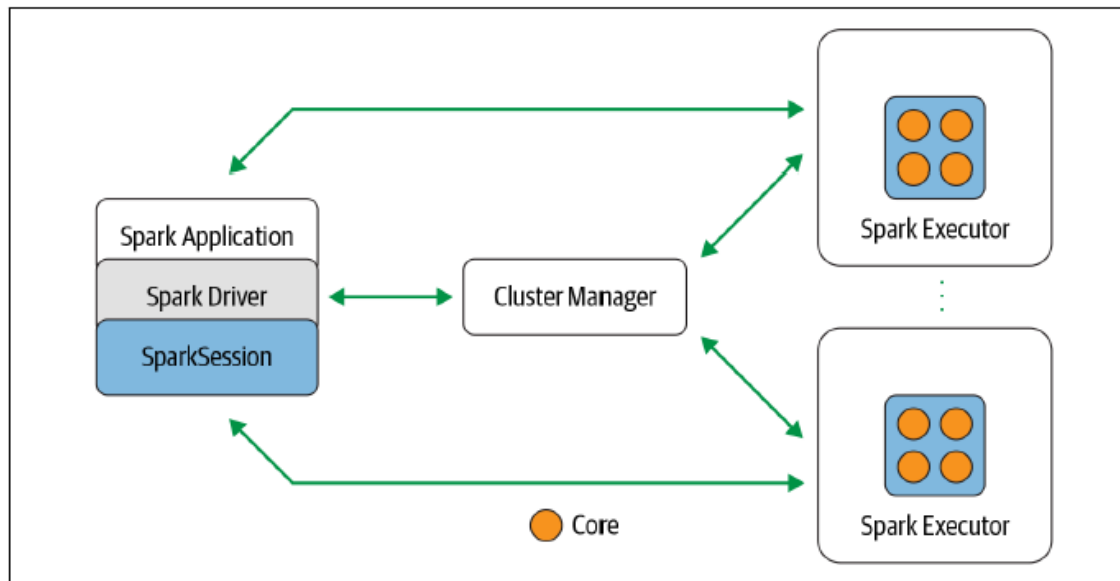
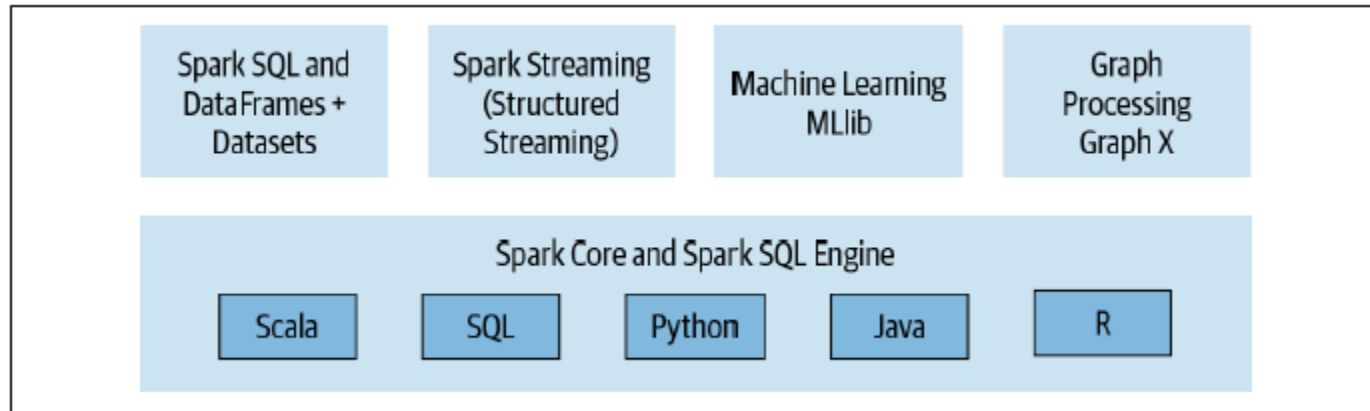
- **Spark SQL and Catalyst Optimizer**
- **Machine Learning with Mllib**
- **Structured Streaming**

# Spark Design Philosophy

- Speed
- Ease of use
- Modularity
- Extensibility



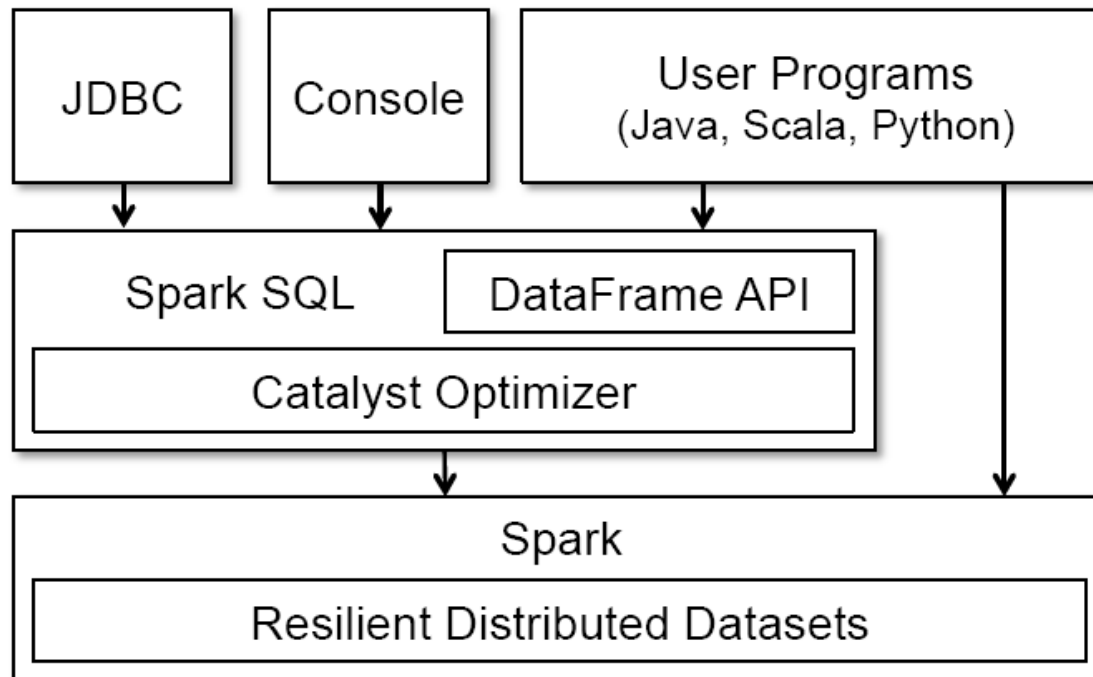
# Spark: a unified stack for distributed execution





# Spark SQL

- Unifies Spark components and permits abstraction to DataFrames/Datasets in Java, Scala, Python, and R
- Keep track of schema and support optimized relational operations



# RDD vs. DataFrame

- RDD

```
# Create an RDD of tuples (name, age)
dataRDD = sc.parallelize([("Brooke", 20), ("Denny", 31), ("Jules", 30),
    ("TD", 35), ("Brooke", 25)])
# Use map and reduceByKey transformations with their lambda
# expressions to aggregate and then compute average
```

```
agesRDD = (dataRDD
    .map(lambda x: (x[0], (x[1], 1))))
    .reduceByKey(lambda x, y: (x[0] + y[0], x[1] + y[1]))
    .map(lambda x: (x[0], x[1][0]/x[1][1])))
```

- DataFrame

```
# Create a DataFrame
data_df = spark.createDataFrame([("Brooke", 20), ("Denny", 31), ("Jules", 30),
    ("TD", 35), ("Brooke", 25)], ["name", "age"])
# Group the same names together, aggregate their ages, and compute an average
avg_df = data_df.groupBy("name").agg(avg("age"))
# Show the results of the final execution
avg_df.show()
```

```
+-----+-----+
|  name|avg(age)|
+-----+-----+
|Brooke|   22.5|
|  Jules|   30.0|
|    TD|   35.0|
|  Denny|   31.0|
+-----+-----+
```

# RDD vs. DataFrame

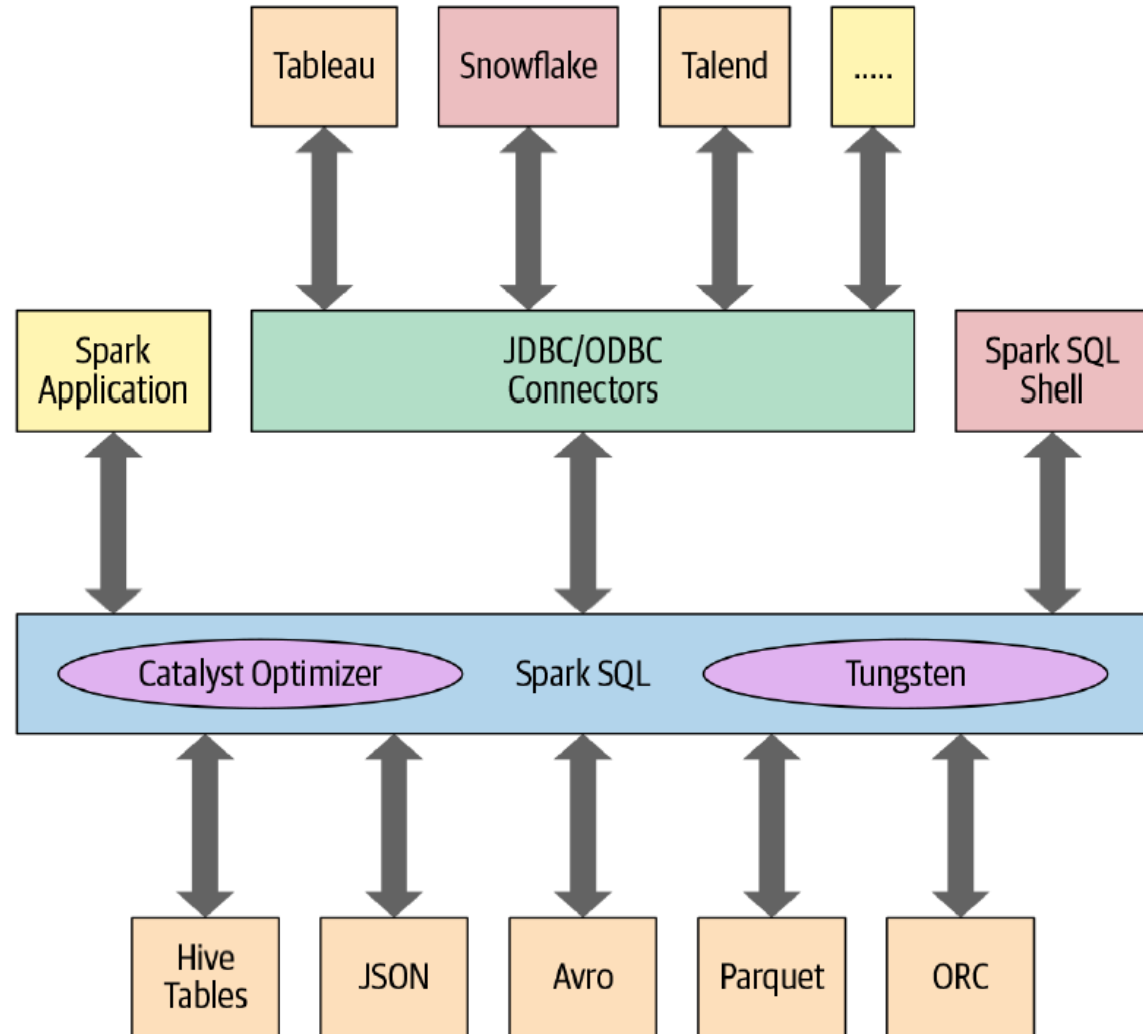
## ○ RDD

- Instruct Spark how to compute the query
- The intention is completely opaque to Spark
- Spark also does not understand the structure of the data in RDDs (which is arbitrary Python objects) or the semantics of user functions (which contain arbitrary code)

## ○ DataFrame

- Tell Spark what to do, instead of How to do
- The code is far more expressive as well as simpler
  - Using a domain specific language (DSL) similar to python pandas
  - Use high-level DSL operators to compose the query
- Spark can inspect or parse this query and understand our intention, it can then optimize or arrange the operations for efficient execution

# Spark SQL

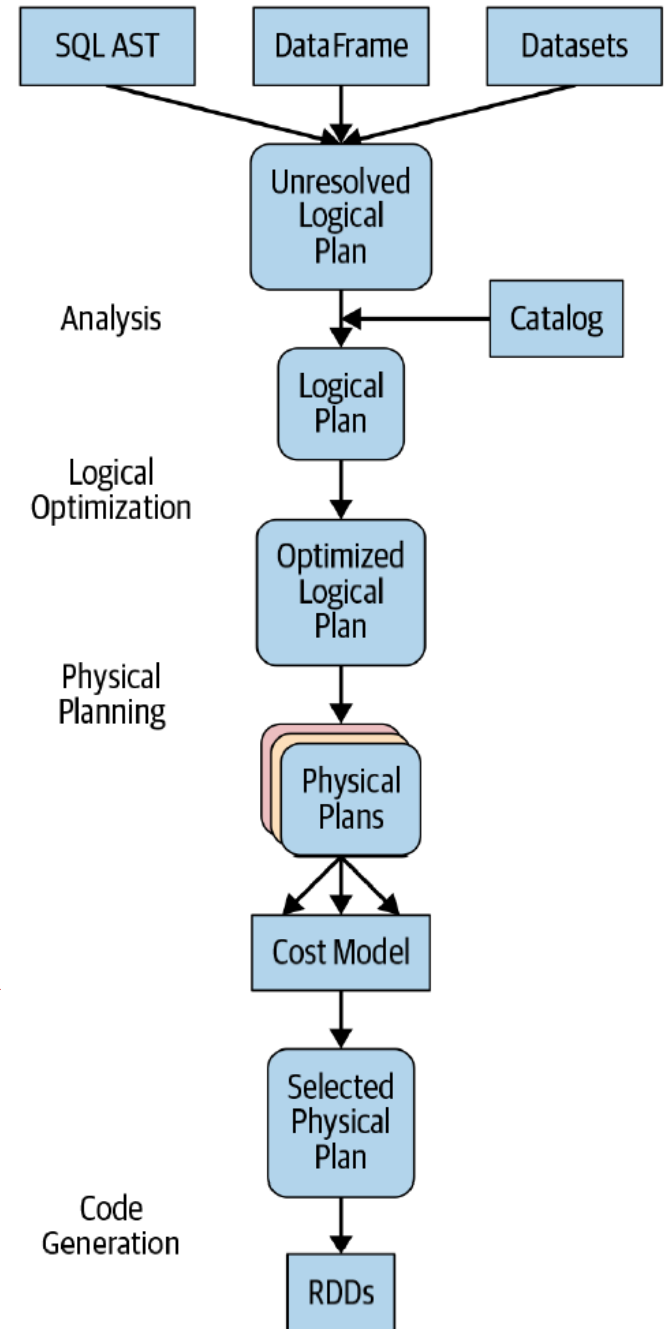


# The Catalyst Optimizer

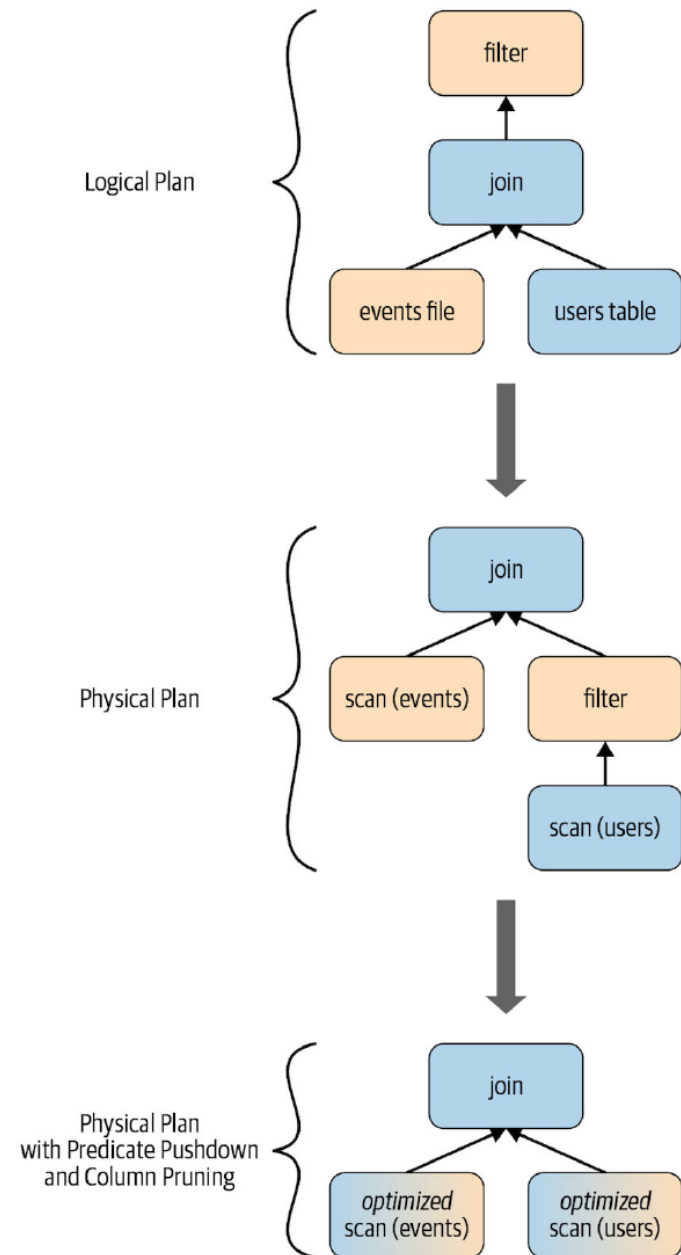
- Takes a computational query and converts it into an execution plan through four transformational phases:

1. Analysis
2. Logical optimization
3. Physical planning
4. Code generation

A Spark computation's four-phase journey

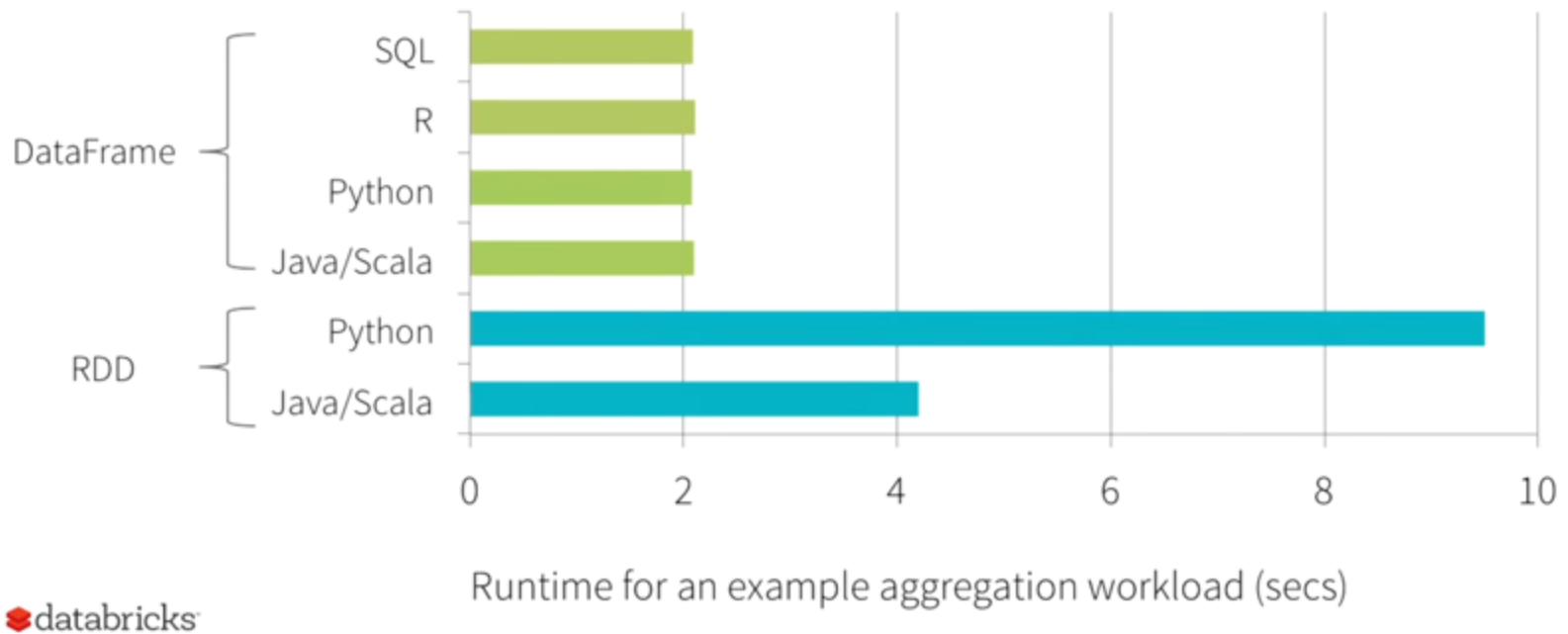


```
// In Scala
// Users DataFrame read from a Parquet table
val usersDF = ...
// Events DataFrame read from a Parquet table
val eventsDF = ...
// Join two DataFrames
val joinedDF = users
  .join(events, users("id") === events("uid"))
  .filter(events("date") > "2015-01-01")
```



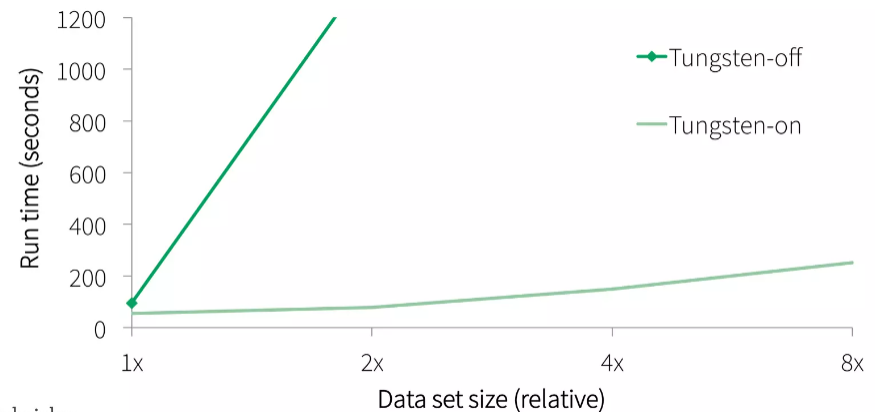
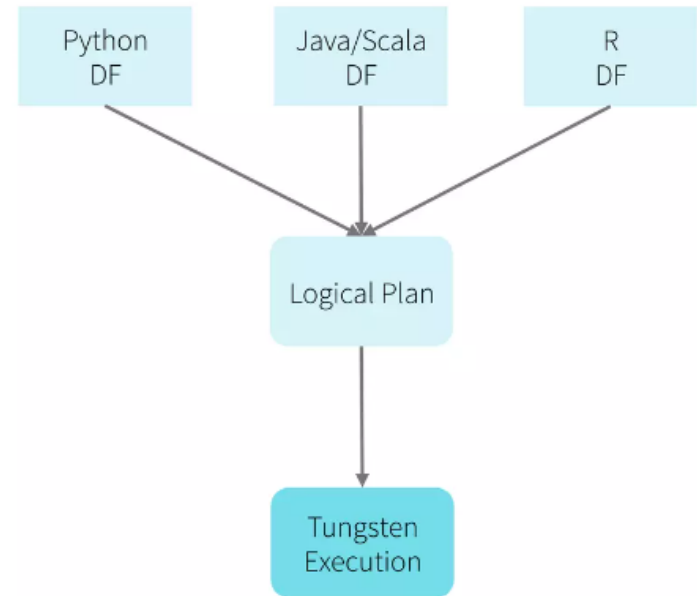
# Benefit of Logical Plan

- Performance Parity Across Languages



# Project Tungsten

- Objectives:
  - Substantially improve the memory and CPU efficiency of Spark applications
  - Push performance closer to the limits of modern hardware
- How?
  - Memory Management and Binary Processing
  - Cache-aware computation
  - Code generation

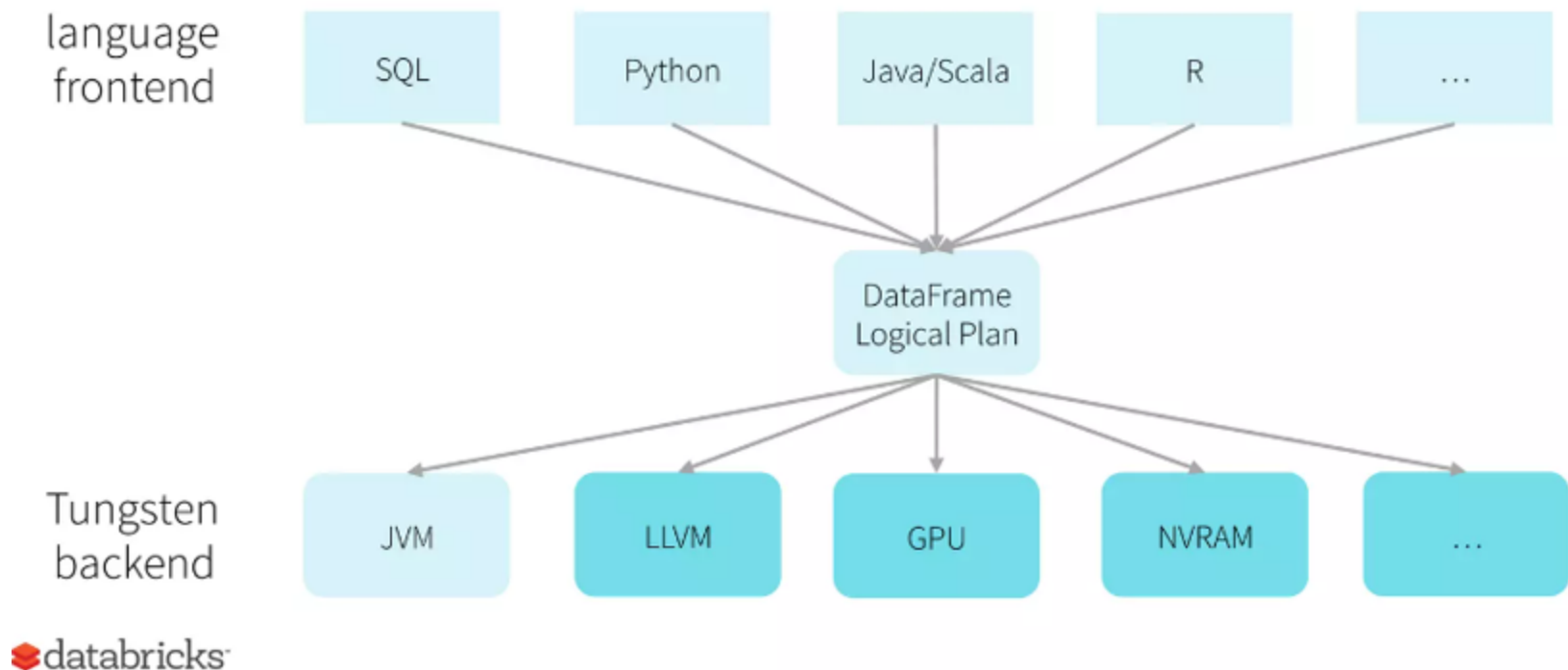


 databricks

Source: <https://youtu.be/VbSar607HM0>  
<https://youtu.be/5ajs8EIPWGI>

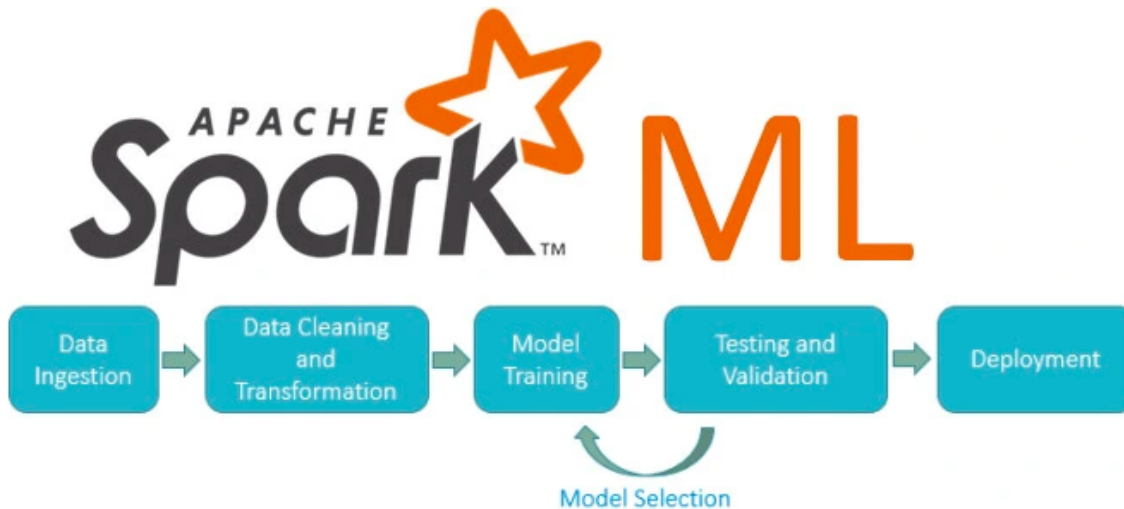


# Unified API, One Engine, Automatically Optimized

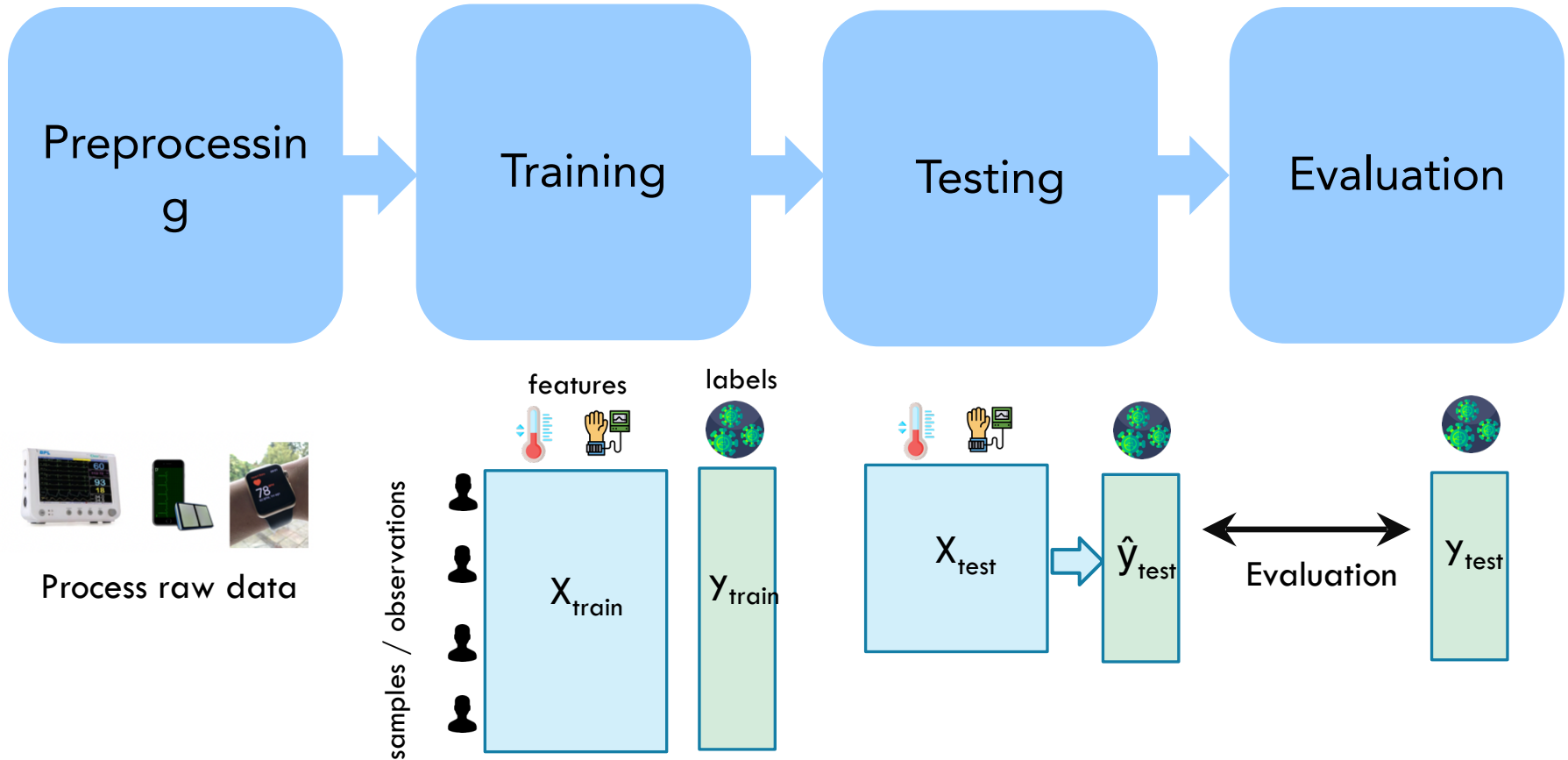


# Today's Plan

- Spark SQL
- Machine Learning with MLlib
- Structured Streaming



# Typical Machine Learning Pipeline



# Spark MLlib: Simple Logistic Regression Model

```
from pyspark.ml.classification import LogisticRegression

training = spark.read.format("libsvm").load("data/mllib/sample_libsvm_data.txt")

lr = LogisticRegression(maxIter=10)

lrModel = lr.fit(training)

print("Coefficients: " + str(lrModel.coefficients))
print("Intercept: " + str(lrModel.intercept))
```

# Pipelines

**Idea:** building complex pipeline out of simple building blocks  
(Note: scikit-learn pipelines are basically the same as Spark MLlib ones)

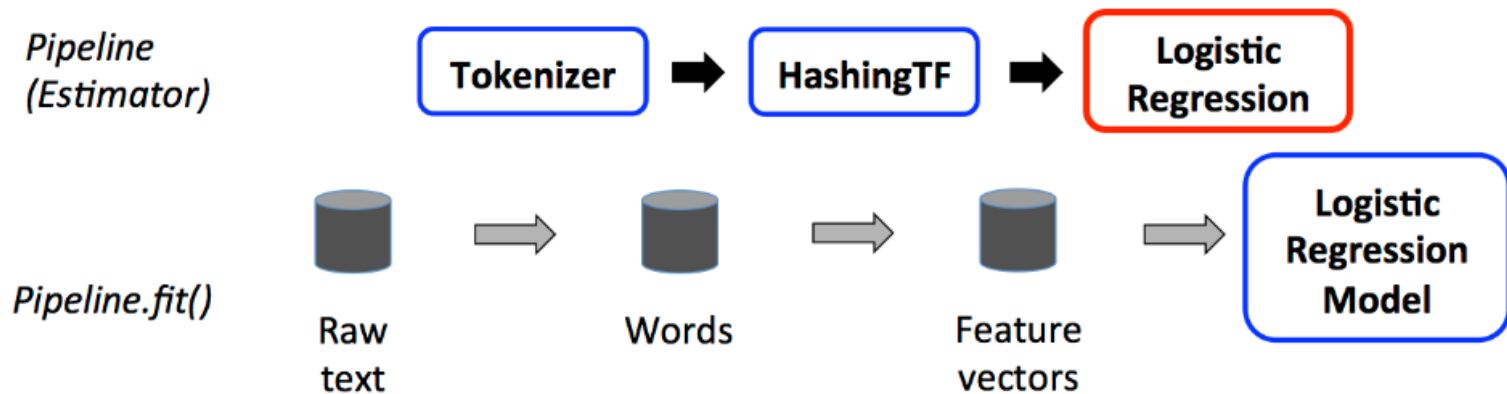


# Pipelines

**Idea:** building complex pipeline out of simple building blocks: e.g. normalization, feature transformation, model fitting.

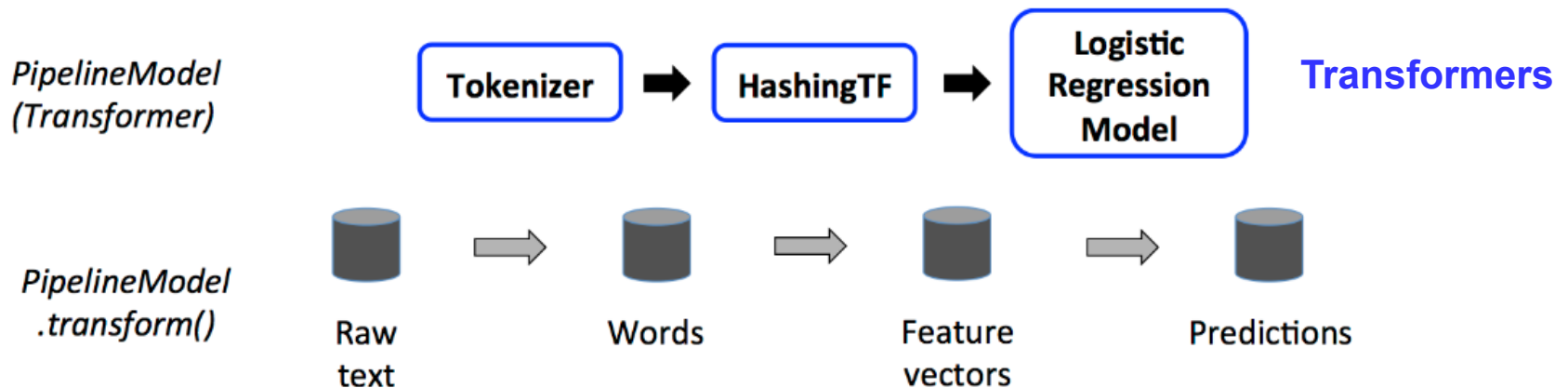
## Why?

- Better code reuse: without pipelines, we would repeat a lot of code, e.g. between the training and test pipelines, cross-validation, model variants, etc.
- Easier to perform cross validation, and hyperparameter tuning.



# Building Blocks: Transformers

- **Transformers** are for mapping DataFrames to DataFrames
  - Examples: one-hot encoding, tokenization
  - Specifically, a Transformer object has a `transform()` method, which performs its transformation
- Generally, these transformers output a new DataFrame which **append** their result to the original DataFrame.
  - Similarly, a fitted model (e.g. logistic regression) is a Transformer that transforms a DataFrame into one with the predictions appended.



# Building Blocks: Estimator

- **Estimator** is an algorithm which takes in data, and outputs a fitted model. For example, a learning algorithm (the LogisticRegression object) can be fit to data, producing the trained logistic regression model.
- They have a fit() method, which returns a Transformer.

```
from pyspark.ml.classification import LogisticRegression

training = spark.read.format("libsvm").load("data/mllib/sample_libsvm_data.txt")

lr = LogisticRegression(maxIter=10)

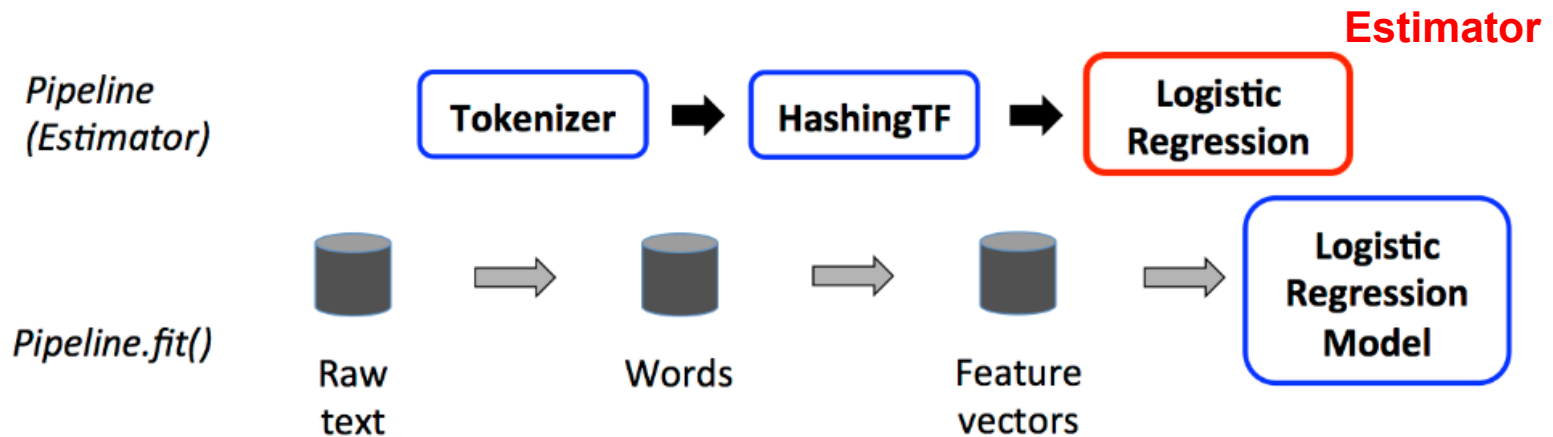
lrModel = lr.fit(training)

print("Coefficients: " + str(lrModel.coefficients))
print("Intercept: " + str(lrModel.intercept))
```



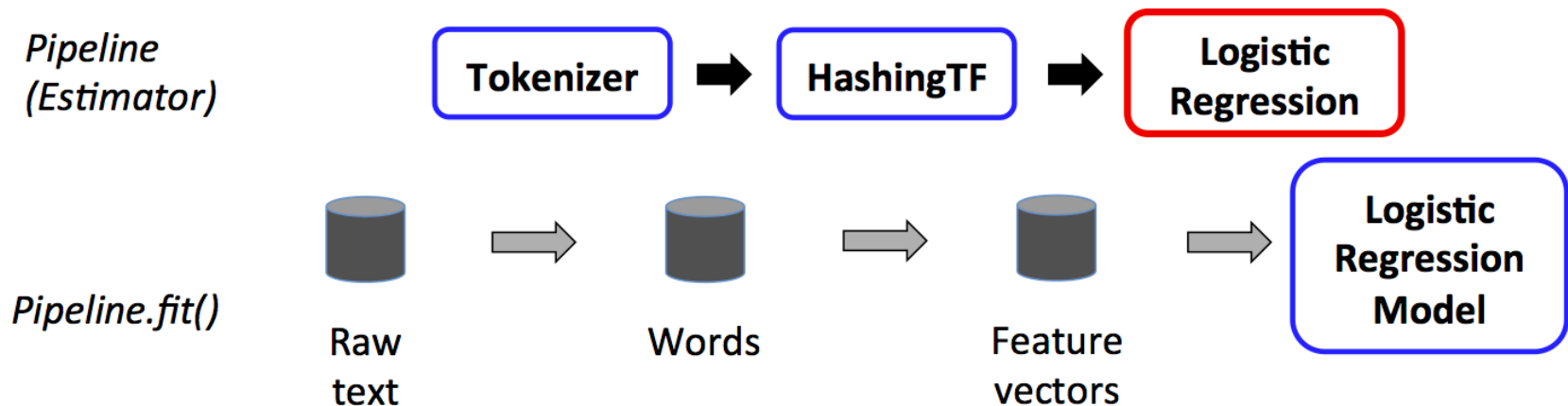
# Building Blocks: Estimator

- **Estimator** is an algorithm which takes in data, and outputs a fitted model. For example, a learning algorithm (the LogisticRegression object) can be fit to data, producing the trained logistic regression model.
- They have a fit() method, which returns a Transformer.



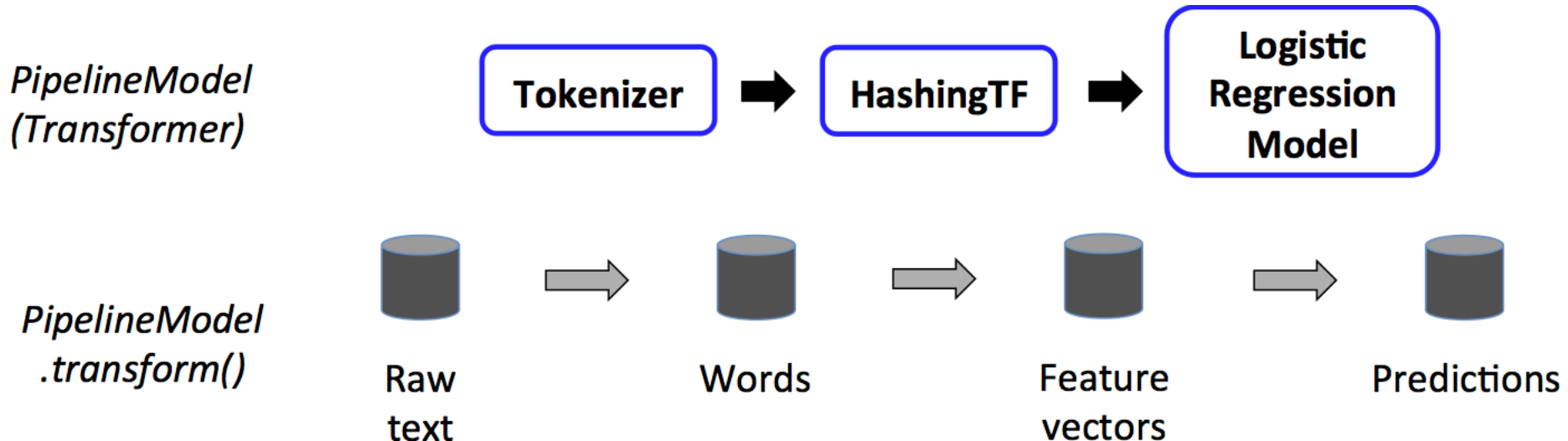
# Pipeline: Training Time

- A pipeline chains together multiple Transformers and Estimators to form an ML workflow.
- Pipeline is an Estimator. When `Pipeline.fit()` is called:
  - Starting from the beginning of the pipeline:
  - For Transformers, it calls `transform()`
  - For Estimators, it calls `fit()` to fit the data and returns a fitted model



# Pipeline: Test Time

- The output of `Pipeline.fit()` is the estimated pipeline model (of type `PipelineModel`).
  - It is a transformer, and consists of a series of Transformers.
  - When its `transform()` is called, each stage's `transform()` method is called.



# Demo\_3: Machine Learning Pipeline

```
# Prepare training documents from a list of (id, text, label) tuples.
training = spark.createDataFrame([
    (0, "a b c d e spark", 1.0),
    (1, "b d", 0.0),
    (2, "spark f g h", 1.0),
    (3, "hadoop mapreduce", 0.0)
], ["id", "text", "label"])
```

```
# Configure an ML pipeline, which consists of three stages: tokenizer, hashingTF, and lr.
tokenizer = Tokenizer(inputCol="text", outputCol="words")
hashingTF = HashingTF(inputCol=tokenizer.getOutputCol(), outputCol="features")
lr = LogisticRegression(maxIter=10, regParam=0.001)
pipeline = Pipeline(stages=[tokenizer, hashingTF, lr])
```

```
# Fit the pipeline to training documents.
model = pipeline.fit(training)
```

```

1 # Prepare test documents
2 test = spark.createDataFrame([
3     (4, "spark i j k", 1.0),
4     (5, "l m n", 0.0),
5     (6, "spark hadoop spark", 1.0),
6     (7, "apache hadoop", 0.0)
7 ], ["id", "text", "label"])

```

```

# Configure an ML pipeline, which consists of three stages: tokenizer, hashingTF, and lr.
tokenizer = Tokenizer(inputCol="text", outputCol="words")
hashingTF = HashingTF(inputCol=tokenizer.getOutputCol(), outputCol="features")
lr = LogisticRegression(maxIter=10, regParam=0.001)
pipeline = Pipeline(stages=[tokenizer, hashingTF, lr])

```

```

1 # Make predictions on test documents and print columns of interest.
2 pred_test = model.transform(test)
3 pred_test.show(truncate = False)

```

Python ▶ ▼ - ✕

▶ (3) Spark Jobs

▶  pred\_test: pyspark.sql.dataframe.DataFrame = [id: long, text: string ... 6 more fields]


id	text	label	words	features	rawPrediction	probability	prediction
4	spark i j k	1.0	[spark, i, j, k]	(262144, [19036, 68693, 173558, 213660], [1.0, 1.0, 1.0, 1.0])	[0.5288285522796805, -0.5288285522796805]	[0.6292098489668488, 0.37079015103315116]	0.0
5	l m n	0.0	[l, m, n]	(262144, [1303, 52644, 248090], [1.0, 1.0, 1.0])	[4.169141395340055, -4.169141395340055]	[0.984770006762304, 0.015229993237696027]	0.0
6	spark hadoop spark	1.0	[spark, hadoop, spark]	(262144, [173558, 198017], [2.0, 1.0])	[-1.8649814141188985, 1.8649814141188985]	[0.13412348342566147, 0.8658765165743385]	1.0
7	apache hadoop	0.0	[apache, hadoop]	(262144, [68303, 198017], [1.0, 1.0])	[5.415644272001849, -5.415644272001849]	[0.9955732114398529, 0.00442678856014711]	0.0

```

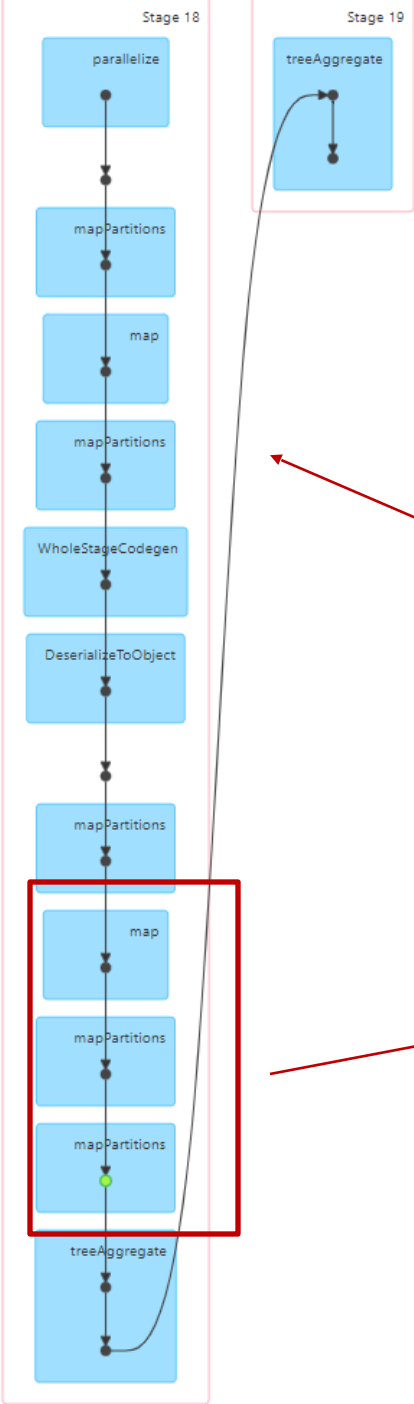
1 # compute accuracy on the test set
2 predictionAndLabels = pred_test.select("prediction", "label")
3 evaluator = MulticlassClassificationEvaluator(metricName="accuracy")
4 print("Test set accuracy = " + str(evaluator.evaluate(predictionAndLabels)))

```

▶ (1) Spark Jobs

▶  predictionAndLabels: pyspark.sql.dataframe.DataFrame = [prediction: double, label: double]

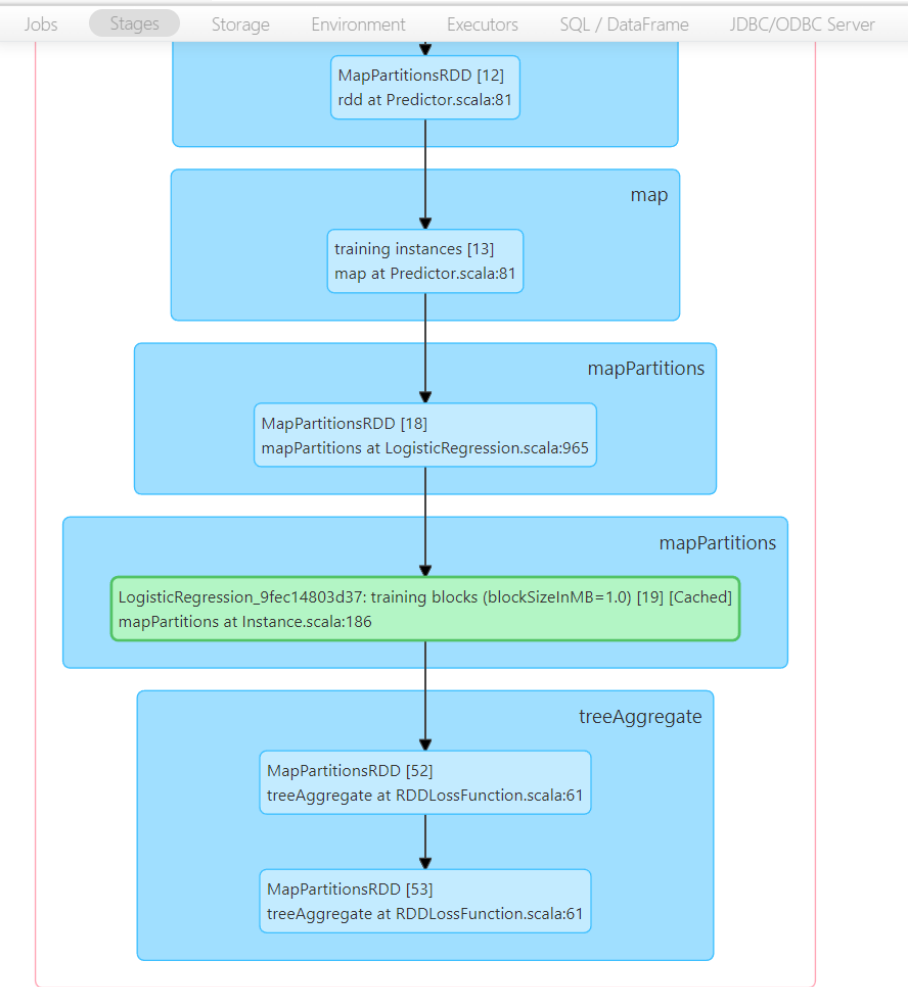
Test set accuracy = 0.75



```
1 # Fit the pipeline to training documents.
2 model = pipeline.fit(training)
```

- ▼ (12) Spark Jobs
- ▶ Job 0 [View](#) (Stages: 2/2)
  - ▶ Job 1 [View](#) (Stages: 2/2)
  - ▶ Job 2 [View](#) (Stages: 2/2)
  - ▶ Job 3 [View](#) (Stages: 2/2)
  - ▶ Job 4 [View](#) (Stages: 2/2)
  - ▶ Job 5 [View](#) (Stages: 2/2)
  - ▶ Job 6 [View](#) (Stages: 2/2)
  - ▶ Job 7 [View](#) (Stages: 2/2)
  - ▶ Job 8 [View](#) (Stages: 2/2)
  - ▼ Job 9 [View](#) (Stages: 2/2)
    - Stage 18: 8/8 [i](#)
    - Stage 19: 2/2 [i](#)
  - ▶ Job 10 [View](#) (Stages: 2/2)
  - ▶ Job 11 [View](#) (Stages: 2/2)

```
# Configure an ML pipeline, which consists of three stages: tokenizer, hashingTF, and lr.
tokenizer = Tokenizer(inputCol="text", outputCol="words")
hashingTF = HashingTF(inputCol=tokenizer.getOutputCol(), outputCol="features")
lr = LogisticRegression(maxIter=10, regParam=0.001)
pipeline = Pipeline(stages=[tokenizer, hashingTF, lr])
```

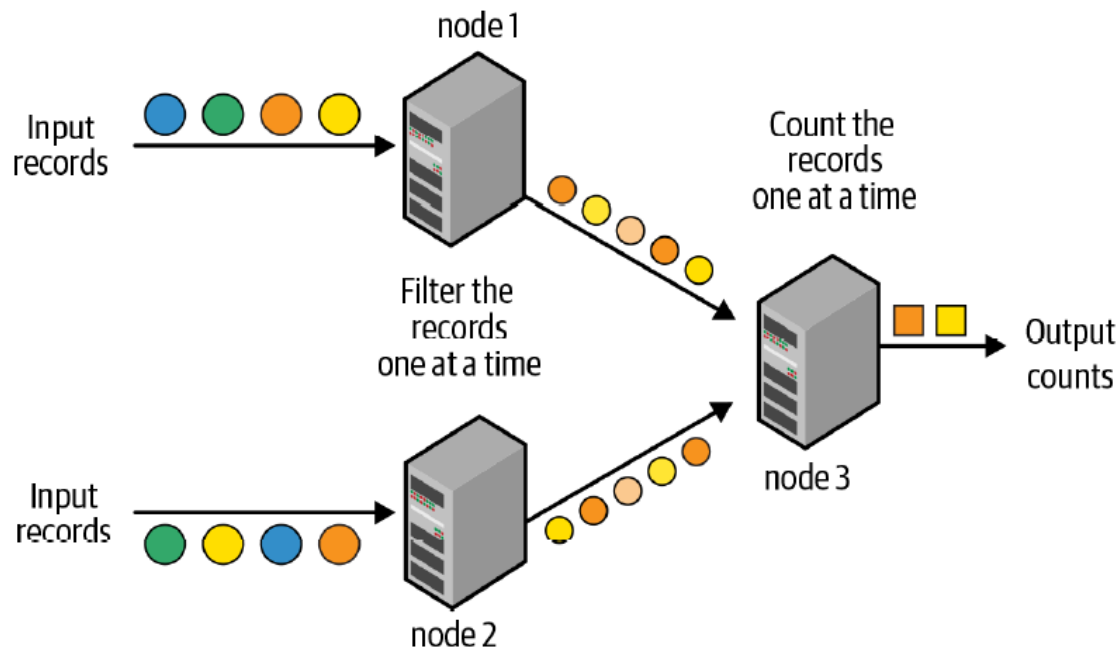


# Today's Plan

- **Spark SQL**
- **Machine Learning with MLlib**
- **Structured Streaming**

# Traditional Model

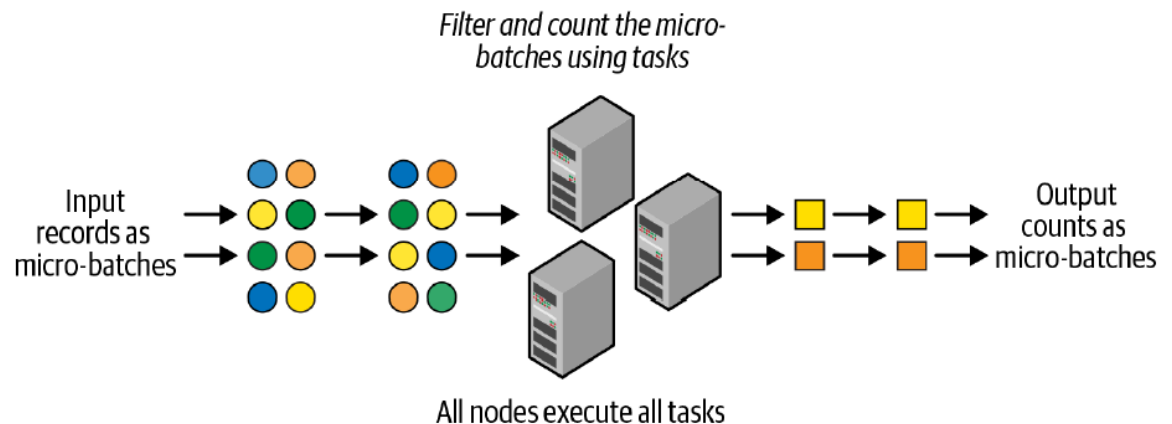
- Traditional record-at-a-time processing model
  - can achieve very low latencies (e.g. milliseconds)
  - not very efficient at recovering from
    - node failures
    - straggler nodes: nodes that are slower than others





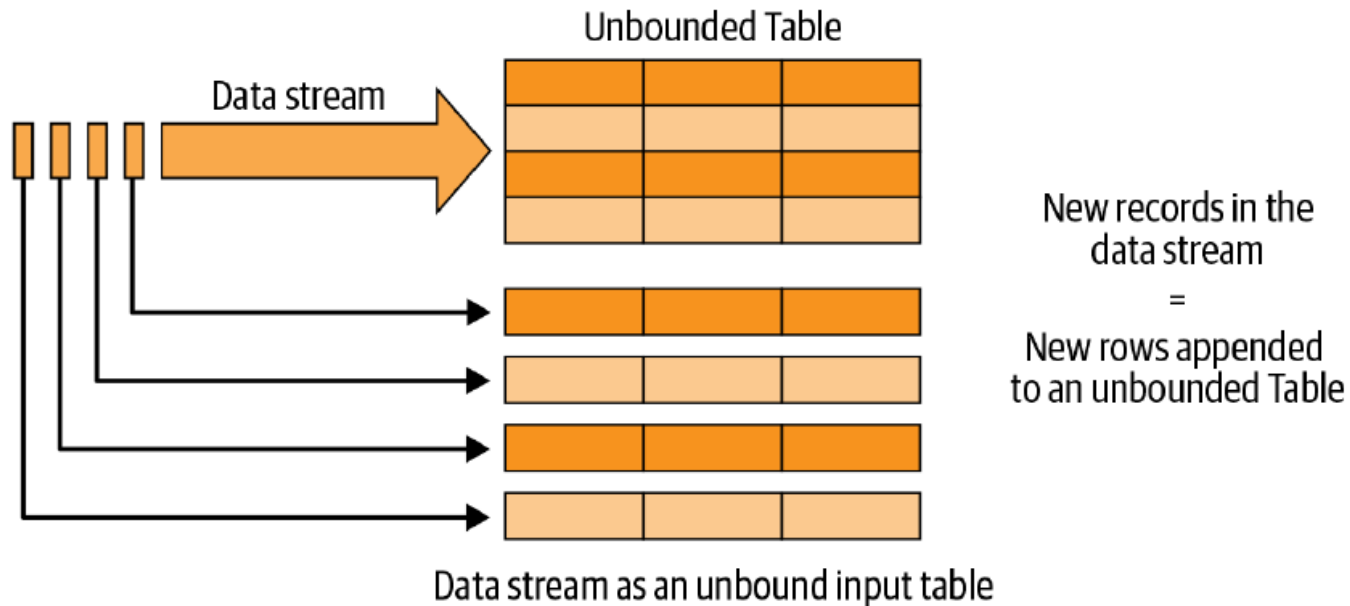
# Micro-Batch Stream Processing

- Structured Streaming uses a micro-batch processing model
  - divides the data from the input stream into micro batches
  - each batch is processed in the Spark cluster in a distributed manner
  - small deterministic tasks generate the output in micro-batches
- Advantages over traditional model
  - quickly and efficiently recover from failures and straggler executors
  - deterministic nature ensures end-to-end exactly-once processing guarantees
- Disadvantages: latencies of a few seconds
  - OK for many applications
  - Application may incur more than a few seconds delay in other parts of pipeline

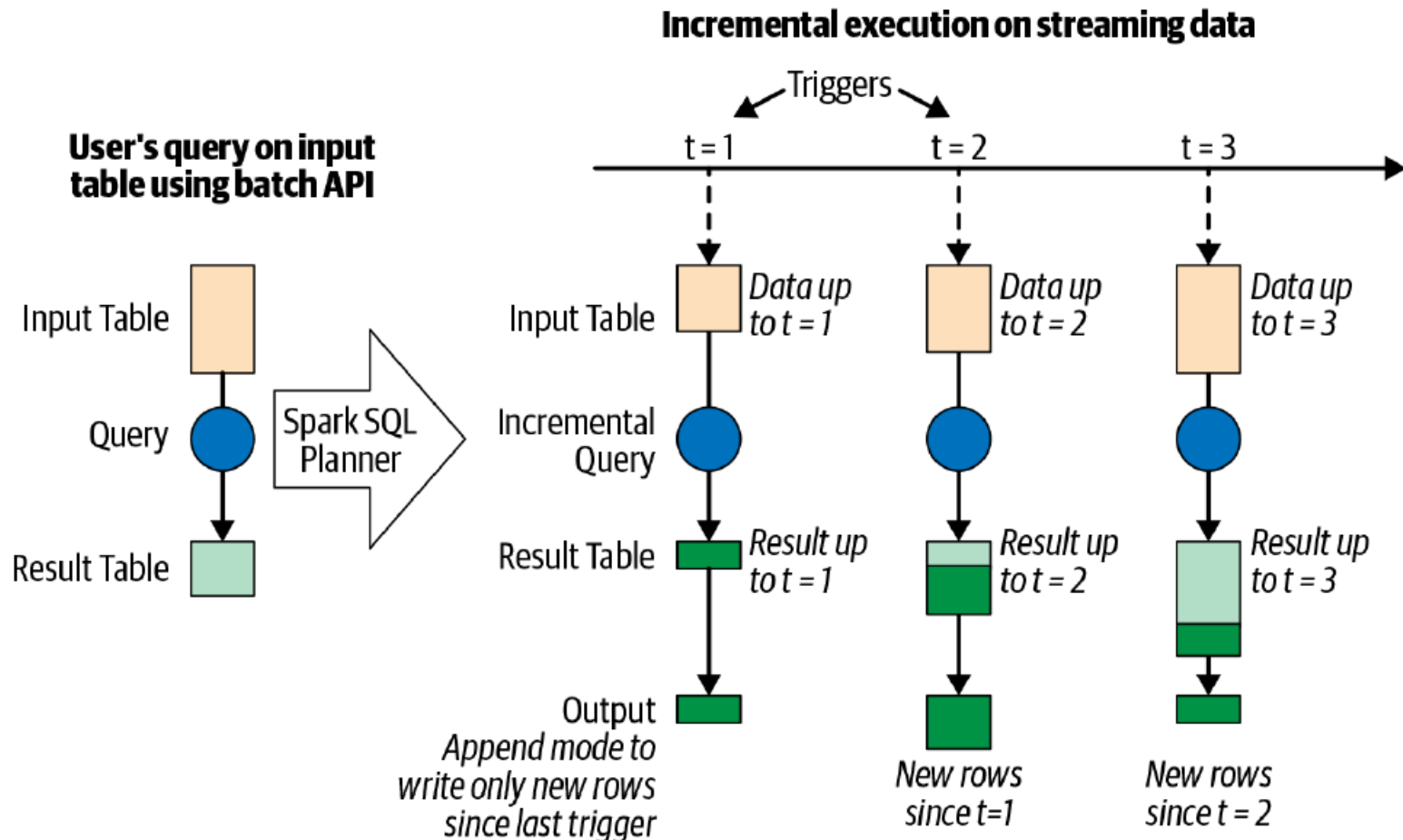


# The Philosophy of Structured Streaming

- For developers, writing stream processing pipelines should be as easy as writing batch pipelines.
  - A single, unified programming model and interface for batch and stream processing
  - A broader definition of stream processing
- The Structured Streaming programming model: data stream as an unbounded table



# The Structured Streaming processing model



Users express query on streaming data using a batch-like API and Structured Streaming incrementalizes them to run on streams.

# Five Steps to Define a Streaming Query

- Step 1: Define input sources
- Step 2: Transform data
- Step 3: Define output sink and output mode
  - Output writing details (where and how to write the output)
  - Processing details (how to process data and how to recover from failures)
- Step 4: Specify processing details
  - Triggering details: when to trigger the discovery and processing of newly available streaming data.
  - Checkpoint Location: store the streaming query process info for failure recovery
- Step 5: Start the query

# Practical\_3: a simple streaming example

Practical\_3 Python ▾

File Edit View Run Help [Last edit was 8 minutes ago](#) [Give feedback](#)

---

Cmd 1

```
1 spark.conf.set("spark.sql.shuffle.partitions", 5)
```

Command took 0.13 seconds -- by aixin@comp.nus.edu.sg at 2/13/2023, 3:35:13 PM on Test

>

Cmd 2

```
1 static = spark.read.json("/databricks-datasets/definitive-guide/data/activity-data/")
2 dataSchema = static.schema
3
```

▶ (3) Spark Jobs

▶ static: pyspark.sql.dataframe.DataFrame = [Arrival\_Time: long, Creation\_Time: long ... 8 more fields]

Command took 38.98 seconds -- by aixin@comp.nus.edu.sg at 2/13/2023, 3:35:17 PM on Test

---

Cmd 3

```
1 streaming = spark.readStream.schema(dataSchema).option("maxFilesPerTrigger", 1)\
2 .json("/databricks-datasets/definitive-guide/data/activity-data")
3
```

▶ streaming: pyspark.sql.dataframe.DataFrame = [Arrival\_Time: long, Creation\_Time: long ... 8 more fields]

Command took 0.36 seconds -- by aixin@comp.nus.edu.sg at 2/13/2023, 3:26:19 PM on Test

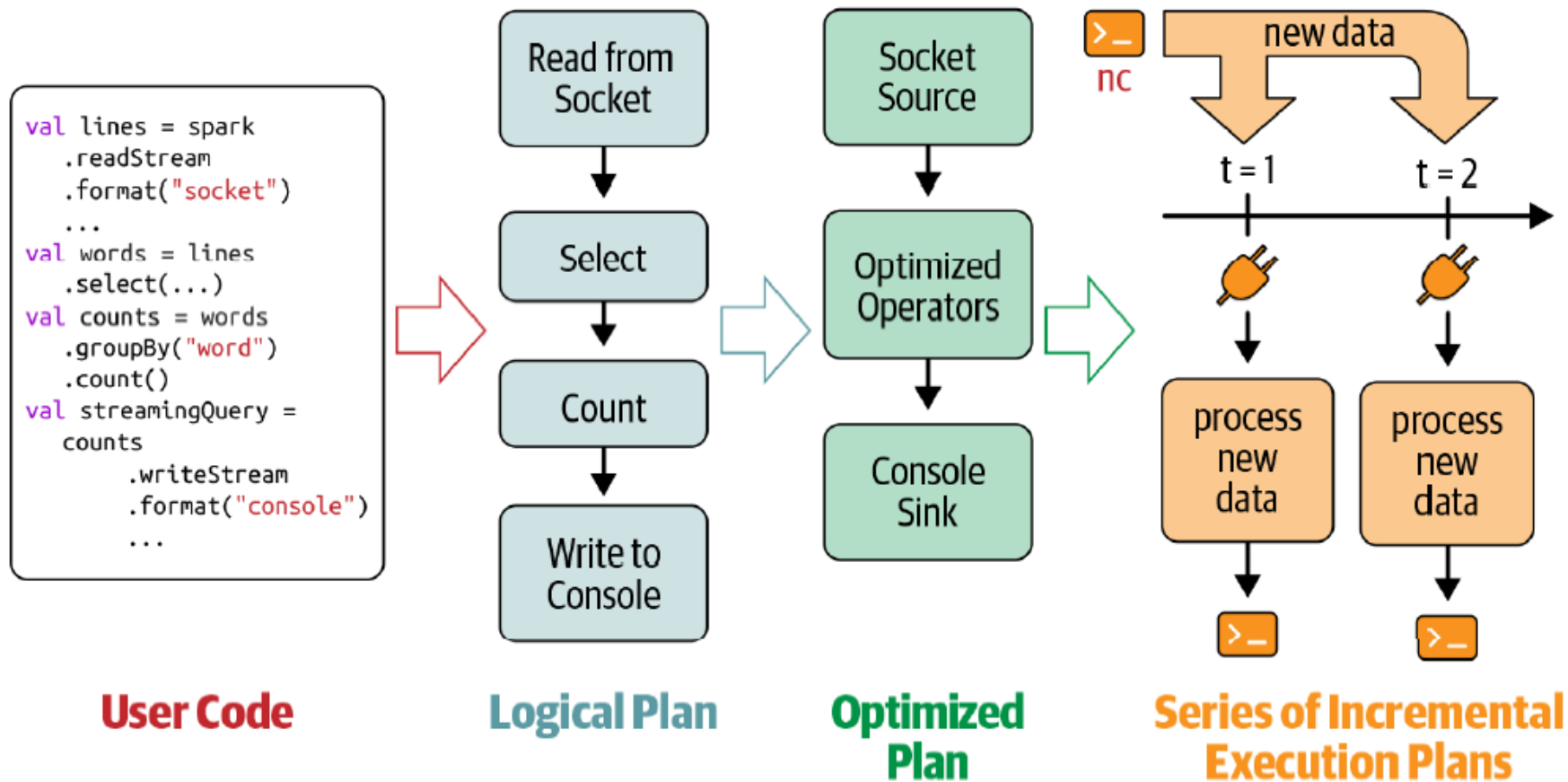
---

Cmd 4

```
1 activityCounts = streaming.groupBy("gt").count()
2
```

Source: <https://github.com/databricks/Spark-The-Definitive-Guide>

# Incremental execution of streaming queries



# Data Transformation

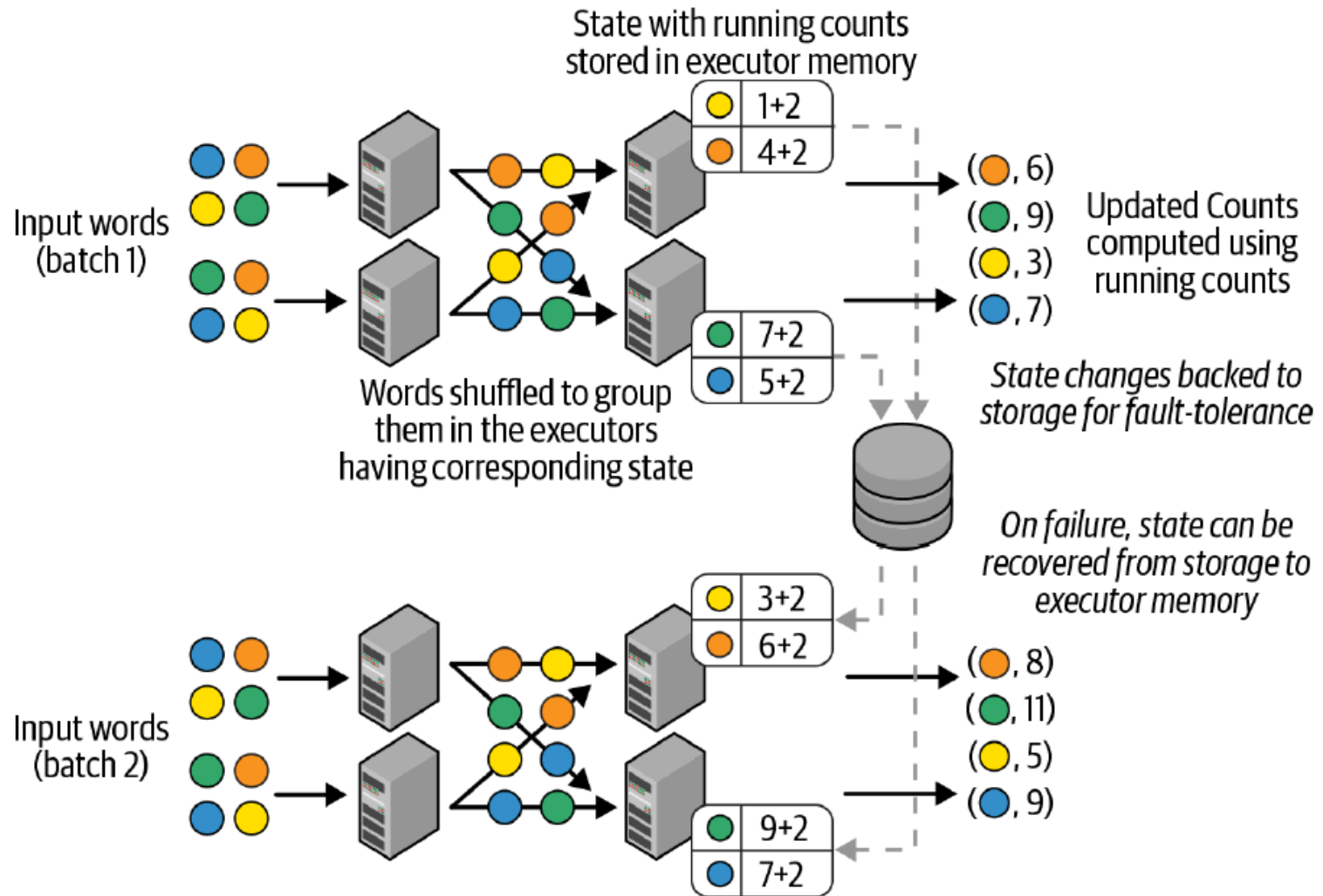
- Stateless Transformation

- Process each row individually without needing any information from previous rows
- Projection operations: `select()`, `explode()`, `map()`, `flatMap()`
- Selection operations: `filter()`, `where()`

- Stateful Transformation

- A simple example: `DataFrame.groupBy().count()`
- In every micro-batch, the incremental plan adds the count of new records to the previous count generated by the previous micro-batch
- The partial count communicated between plans is the state
- The state is maintained in the memory of the Spark executors and is checkpointed to the configured location to tolerate failures.

# Distributed state management in Structured Streaming





# Stateful Streaming Aggregations

- Aggregations Not Based on Time

- Global aggregations

```
runningCount = sensorReadings.groupBy().count()
```

- Grouped aggregations

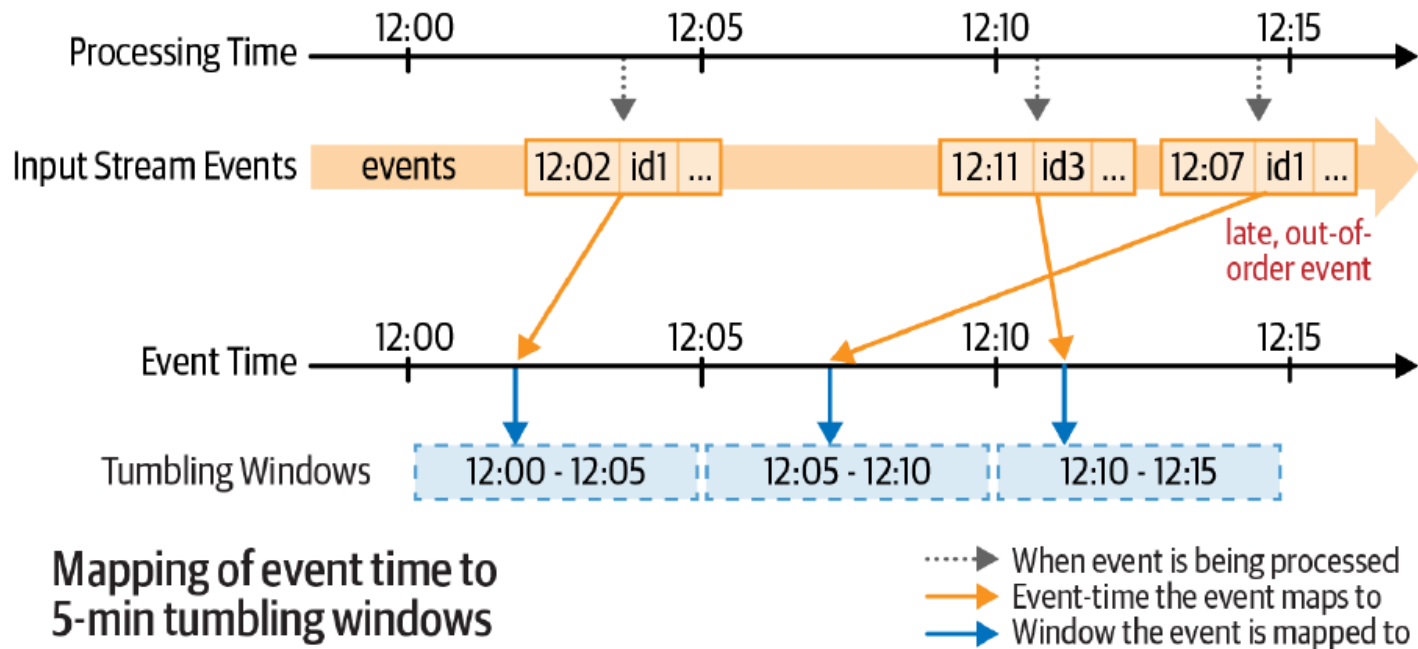
```
baselineValues = sensorReadings.groupBy("sensorId").mean("value")
```

- All built-in aggregation functions in DataFrames are supported
  - `sum()`, `mean()`, `stddev()`, `countDistinct()`, `collect_set()`, `approx_count_distinct()`, and etc.

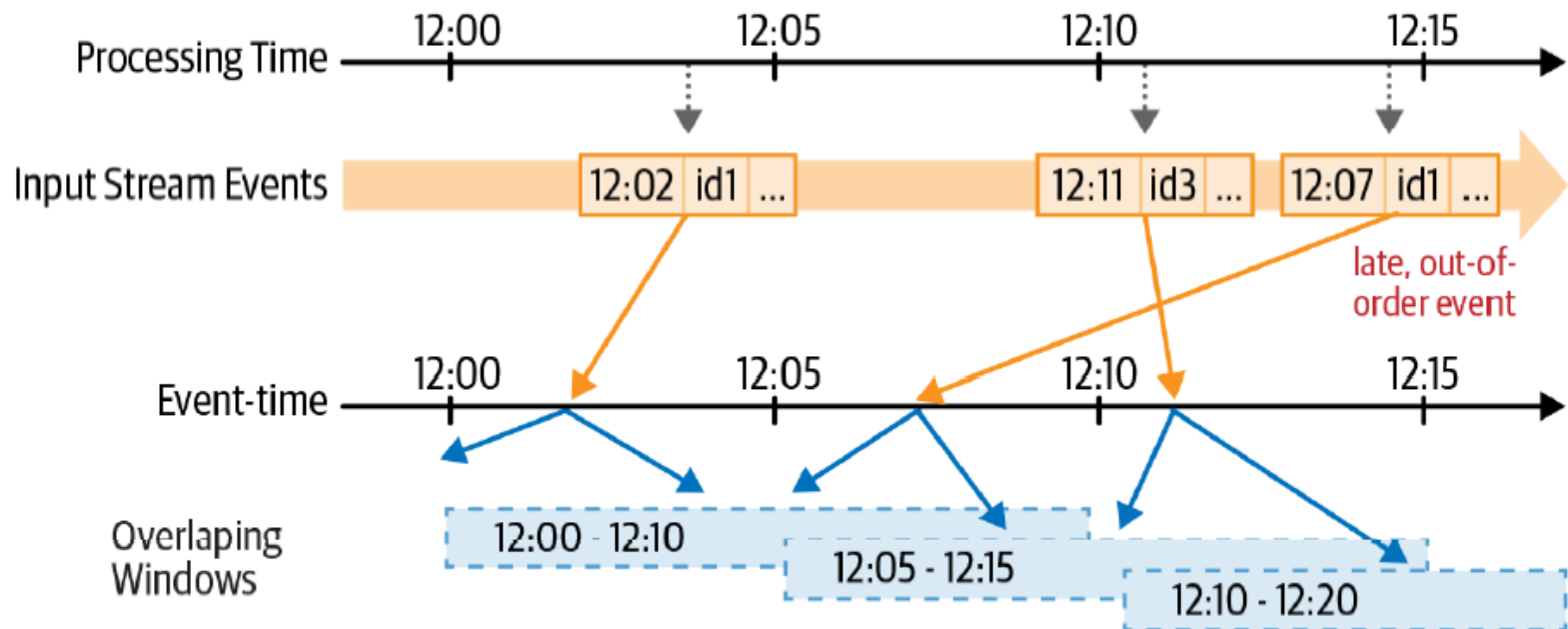
# Stateful Streaming Aggregations

- Aggregations with Event-Time Windows

```
(sensorReadings  
  .groupBy("sensorId", window("eventTime", "5 minute"))  
  .count())
```



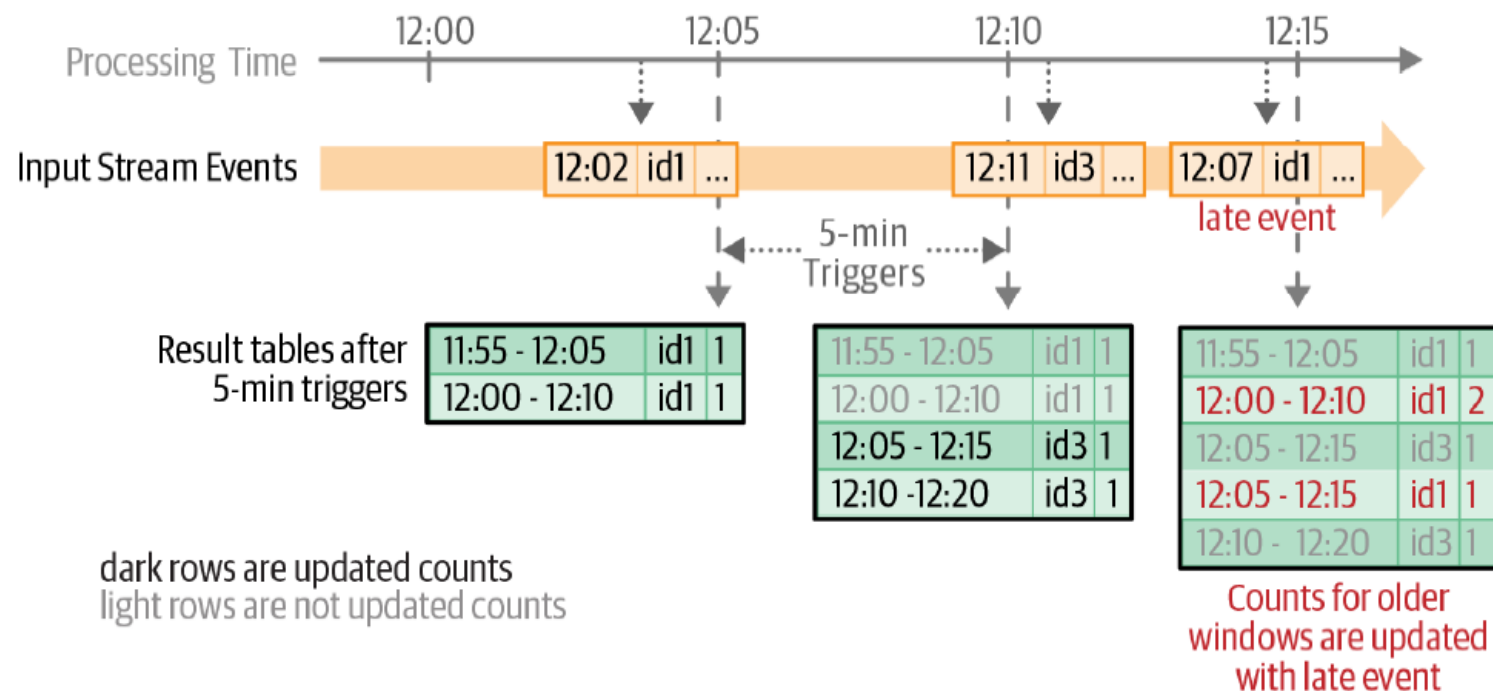
```
(sensorReadings
  .groupBy("sensorId", window("eventTime", "10 minute", "5 minute"))
  .count())
```



**Mapping of event time to overlapping windows of length 10 mins and sliding interval 5 mins**

- .....➤ When event is being processed
- Event-time the event maps to
- Window the event is mapped to

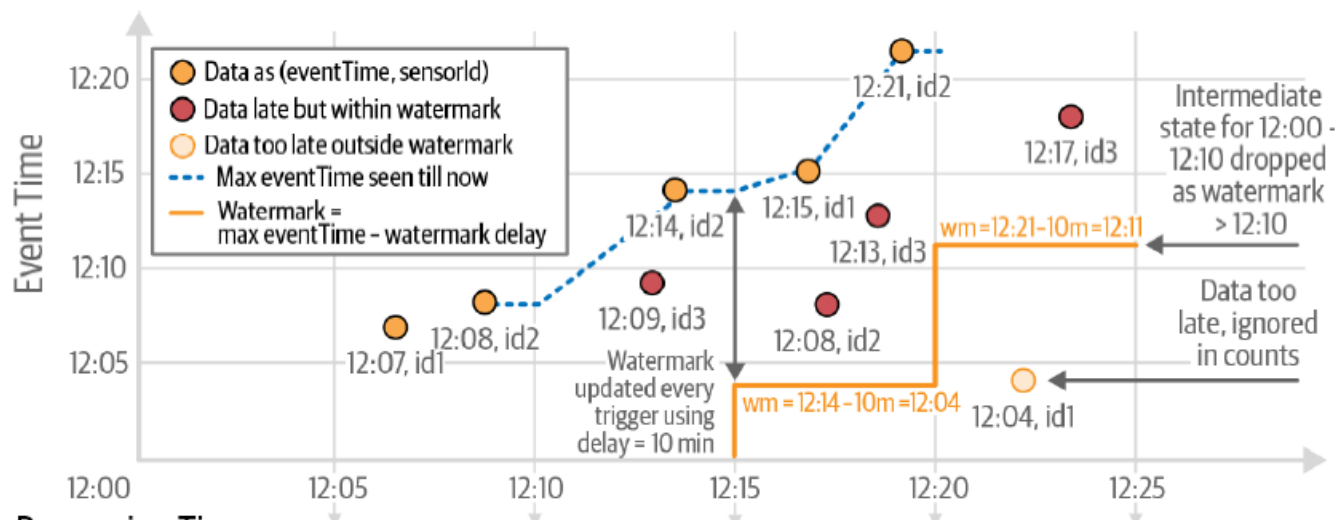
- Updated counts in the result table after each five-minute trigger



# Handling Late Data with Watermarks

(sensorReadings

```
.withWatermark("eventTime", "10 minutes")
.groupBy("sensorId", window("eventTime", "10 minutes", "5 minutes"))
.count())
```



Processing Time  
with 5 min triggers

12:00 - 12:10	id1	1
12:00 - 12:10	id2	1
12:05 - 12:15	id1	1
12:05 - 12:15	id2	1

Result Tables after each trigger

12:00 - 12:10	id1	1
12:00 - 12:10	id2	1
12:00 - 12:10	id3	1
12:05 - 12:15	id1	1
12:05 - 12:15	id2	2
12:05 - 12:15	id3	1
12:10 - 12:20	id2	1

dark rows  
are updated  
counts

12:00 - 12:10	id1	1
12:00 - 12:10	id2	2
12:00 - 12:10	id3	1
12:05 - 12:15	id1	2
12:05 - 12:15	id2	3
12:05 - 12:15	id3	2
12:10 - 12:20	id2	1
12:10 - 12:20	id1	1
12:10 - 12:20	id3	1
...		

12:00 - 12:10	id1	1
12:00 - 12:10	id2	2
12:00 - 12:10	id3	1
12:05 - 12:15	id1	2
12:05 - 12:15	id2	3
12:05 - 12:15	id3	2
12:10 - 12:20	id2	1
12:10 - 12:20	id1	1
12:10 - 12:20	id3	2
...		

Table *not*  
updated with  
too late data  
(12:04, id1)

Table updated  
with late data  
(12:17, id3)

Watermarking in  
Windowed Grouped Counts

# Performance Tuning

- Besides tuning Spark SQL engine, a few other considerations
  - Cluster resource provisioning appropriately to run 24/7
  - Number of partitions for shuffles to be set much lower than batch queries
  - Setting source rate limits for stability
  - Multiple streaming queries in the same Spark application

# Acknowledgements

- CS4225 slides by He Bingsheng and Bryan Hooi
- Jules S. Damji, Brooke Wenig, Tathagata Das & Denny Lee, “Learning Spark: Lightning-Fast Data Analytics”
- Bill Chambers, Matei Zaharia, “Spark: The Definitive Guide”
- Spark SQL: Relational Data Processing in Spark, SIGMOD’15
- <https://spark.apache.org/docs/latest/ml-pipeline.html>