Oscar Mike Claure Cabrera

# Content

Overview of relevant concepts:
- Introduction
- Core components
- Unified framework
- RDDs
- Lazy vs eager evaluation
- Catalyst optimizer
- Shuffling
- Partitioning

Performance evaluation
- Explore query plans
- Spark UI
- Type of joins

Introduction



What is Apache Spark?

A distributed data processing engine designed for big data and large-scale computation
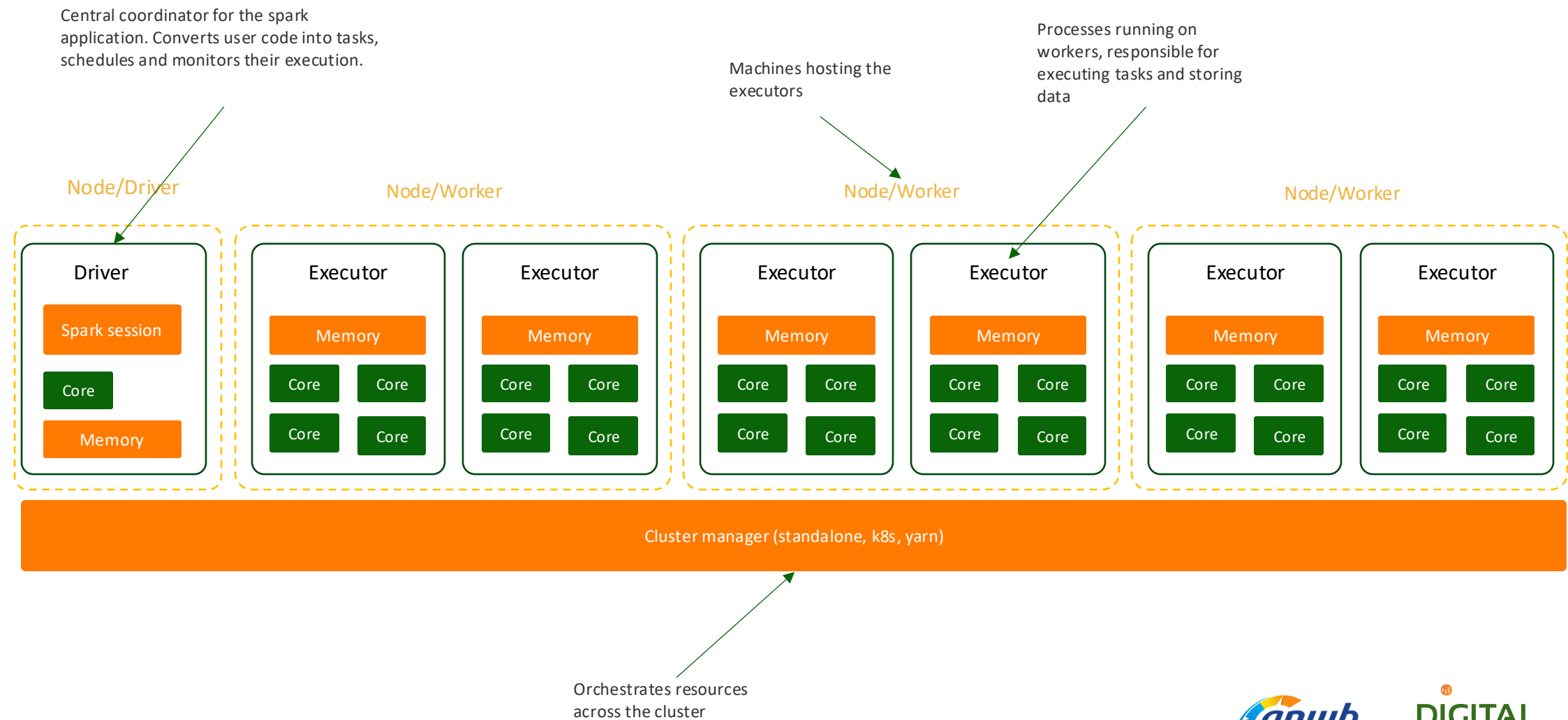
Spark's key features

- Speed: In-memory processing for faster computation
- Scalability: Handles petabytes of data across large clusters
- Versatility: Supports multiple programming languages (API, Java, Scala, R, SQL)
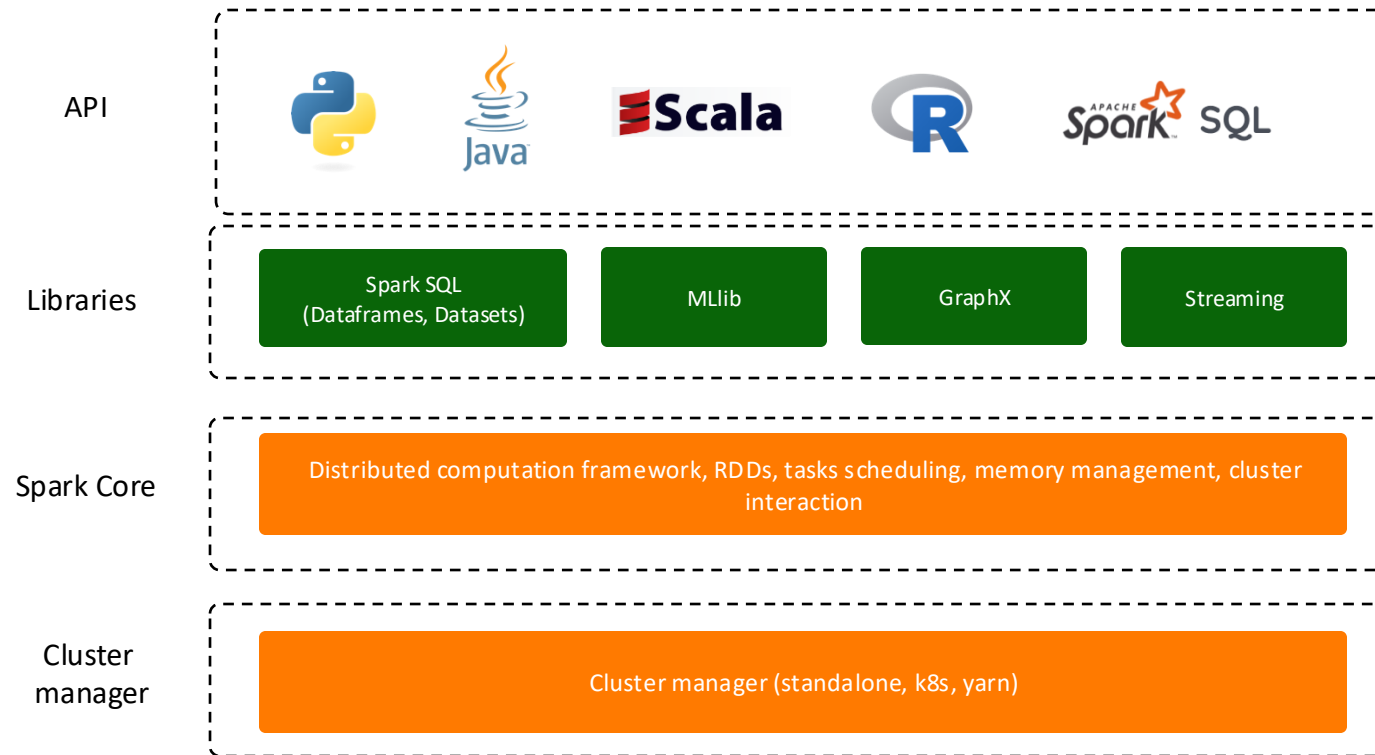
Common use cases:

- ETL pipelines
- Data Analytics
- Machine Learning workflows
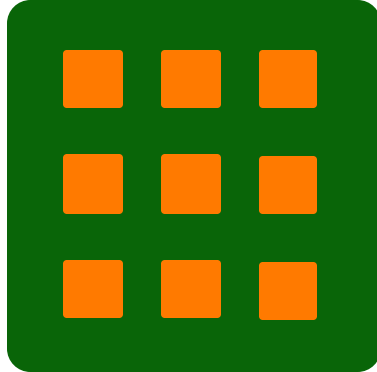
# Spark's core components

Central coordinator for the spark application. Converts user code into tasks, schedules and monitors their execution.

Machines hosting the executors

Processes running on workers, responsible for executing tasks and storing data

**Node/Driver**

**Node/Worker**

**Node/Worker**

**Node/Worker**

| Driver |
|---|
| Spark session |
| Core |
| Memory |

| Executor |
|---|
| Memory |
| Core / Core |
| Core / Core |

| Executor |
|---|
| Memory |
| Core / Core |
| Core / Core |

| Executor |
|---|
| Memory |
| Core / Core |
| Core / Core |

| Executor |
|---|
| Memory |
| Core / Core |
| Core / Core |

| Executor |
|---|
| Memory |
| Core / Core |
| Core / Core |

| Executor |
|---|
| Memory |
| Core / Core |
| Core / Core |

**Cluster manager (standalone, k8s, yarn)**

Orchestrates resources across the cluster

& DIGITAL POWER
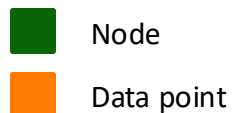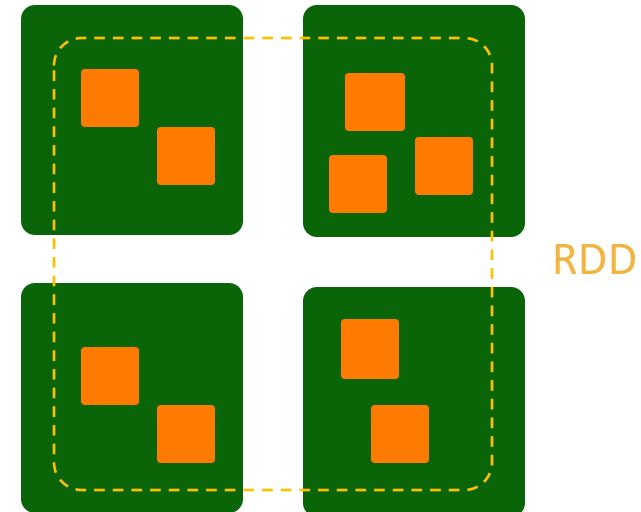
# Spark's unified framework

# RDDs

Resilient distributed dataset (RDD), is a fault-tolerant, immutable collection of elements that can be operated in parallel across a cluster of machines.

my_var = [1, 2, 3, 4, 5, …, N]

my_rdd = sc.parallelize([1, 2, 3, 4, 5, …, N])

RDD

■ Node

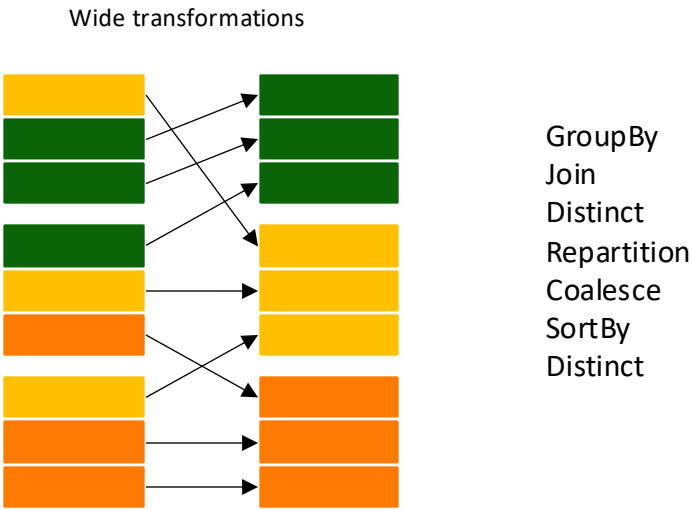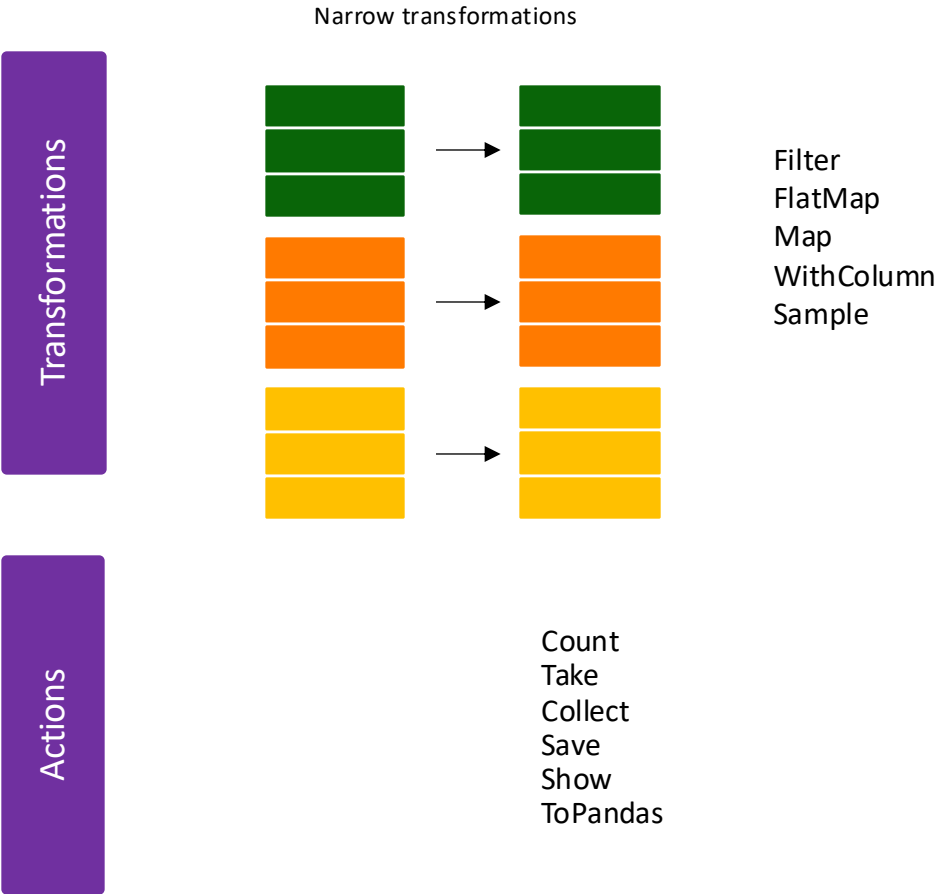■ Data point

# Lazy vs eager evaluation

## Lazy

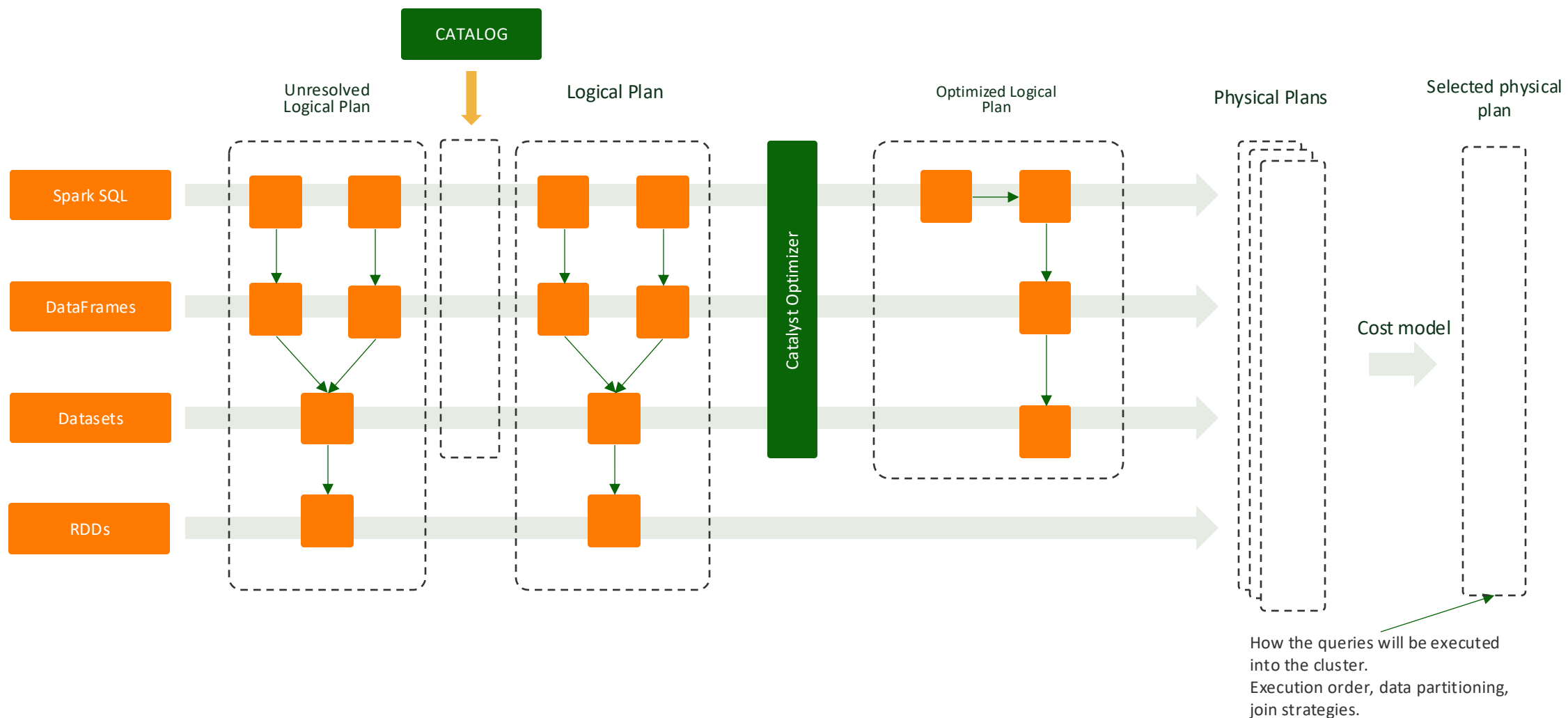Evaluation of expressions is delayed until
their results are needed

## Eager

Evaluation of expressions occur every time
a new expression is declared.
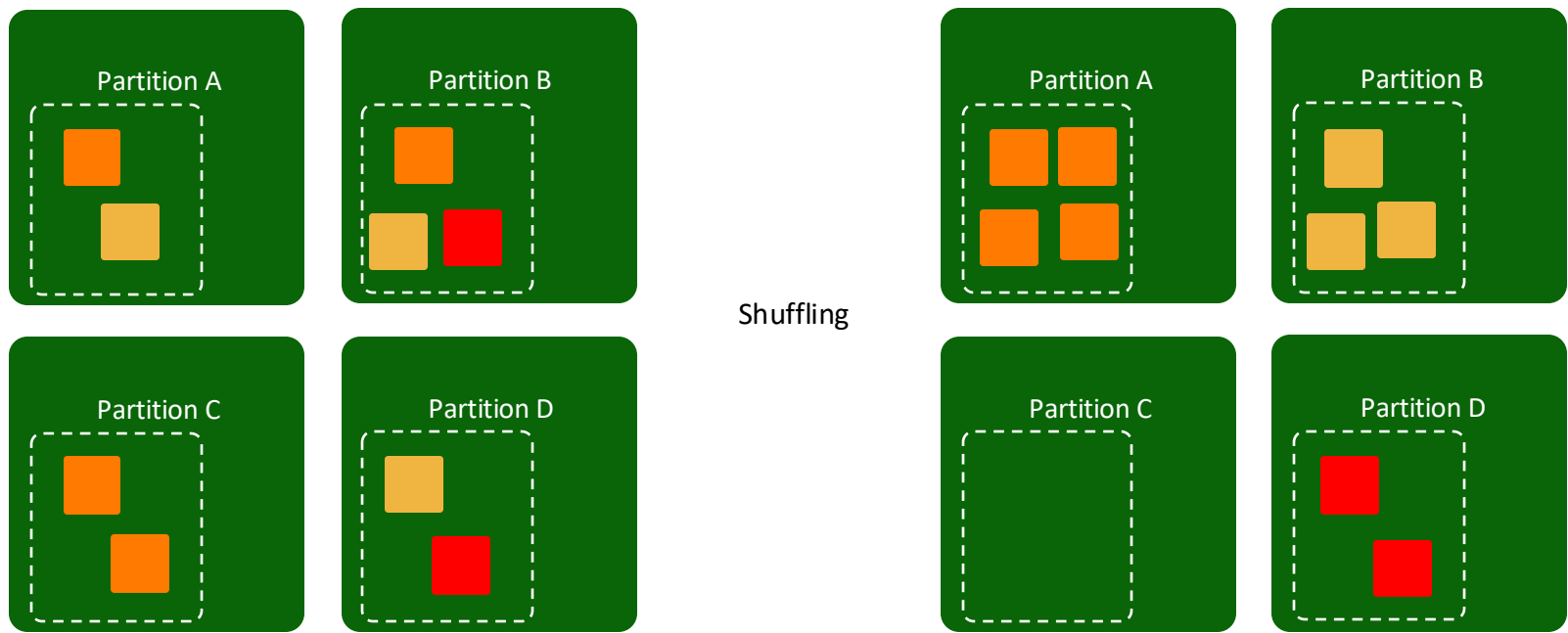
Transformations and actions

Narrow transformations

Wide transformations

**Transformations**

Filter
FlatMap
Map
WithColumn
Sample

GroupBy
Join
Distinct
Repartition
Coalesce
SortBy
Distinct

**Actions**

Count
Take
Collect
Save
Show
ToPandas

# Catalyst optimizer and Tungsten engine



CATALOG

Unresolved Logical Plan

Logical Plan

Optimized Logical Plan

Physical Plans

Selected physical plan

Spark SQL

DataFrames

Datasets

RDDs

Catalyst Optimizer

Cost model

How the queries will be executed into the cluster.
Execution order, data partitioning, join strategies.

Source: Key topics in Apache Spark, Catalyst Optimizer

anwb & DIGITAL POWER

# Shuffling and partitioning



Shuffling

Node

Data point

# Memory management



Executor Container

JVM

On-heap memory
`spark.executor.memory`

Spark.memory.fraction

Storage memory
RDD, DataFrames caching,
Broadcast variables

Execution memory
GroupBy,Joins, Sorts,
Shuffles

User memory
User defined Data Structures:
Variable, Objects, metadata

Reserved memory
300MB

Off-heap memory
`spark.executor.offHeap.enabled`
`spark.executor.offHeap.size`

Overhead
`spark.executor.memoryOverheadsize`
~10-20% of `spark.executor.memory`

Source: Spark's memory management overview