

# SPARK SUMMARY DOCUMENT

# TABLE OF CONTENTS

<b>SPARK CORE</b>	<b>7</b>
Spark	7
Spark Components	7
Client	8
Driver	8
Executors	9
Worker	9
Cluster Manager	9
Slots	10
Spark session	10
Dataframes	11
Partitions	11
Glom	11
Transformations	11
Narrow transformations	12
Wide transformation	12
Action	12
Lazy evaluation	12
Datasets	13
DAG	13
Spark Lifecycle (Explain Plan)	13
Catalyst optimizer	14
Catalyst optimizer trees	14
Logical plan	15
Physical plan	15
Types of transformations in optimizer	15
Dynamic partition pruning	15
Adaptive Query Execution	16
AQE Scenarios	16
Spark submit	16
Spark submit master Url	17
Dag scheduler	18
Task scheduler	18
Jobs	19
Stages	19
Tasks	19
Spark config	19

Additional config properties	20
Dynamic allocation	20
Execution modes	21
Local Mode	21
Client mode	21
Cluster mode	22
Serialization	23
Shuffle	23
Caching	24
Uncache Table	24
Persist Options	24
RDD	25
Types of RDD	25
Broadcast Variables	26
Accumulators	26
Repartition	27
Coalesce	27
Repartition vs Coalesce	28
Catalog	28
<b>SPARK READ/WRITE</b>	<b>29</b>
Generic File options	29
Read	29
Write	30
CSV	30
Json	31
Text	31
Parquet	31
Orc	31
Database	32
<b>TABLES AND VIEWS</b>	<b>33</b>
Global Managed Table	33
Global Unmanaged/External Table	33
Local Table (a.k.a) Temporary Table (a.k.a) Temporary View	33
Global Temporary View	34
Global Permanent View	34
<b>DATAFRAMES</b>	<b>35</b>
CreateDataFrame	35
printSchema	36
createOrReplaceTempView	36
Schema	36

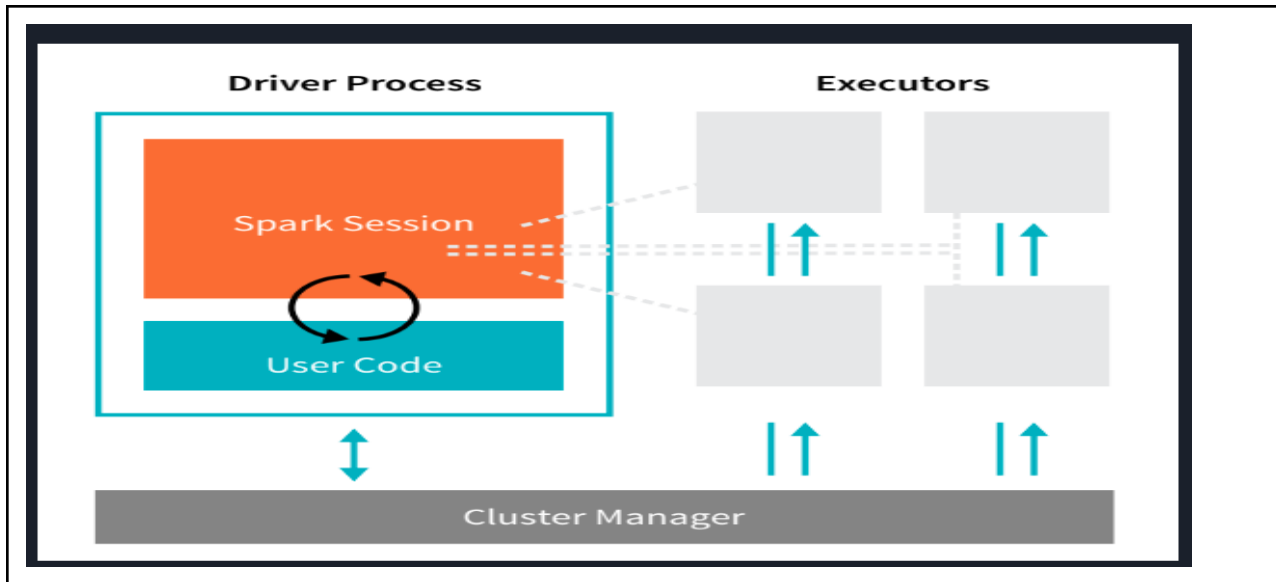
col	36
row	37
Null in spark	37
expr	37
select	37
selectExpr	38
withColumn	38
withColumnRenamed	39
alias	39
drop	40
dropDuplicates	40
dropna	40
where/filter	42
distinct	42
sample	42
split	43
sort/orderBy	44
sortwithinpartitions	45
asc_nulls_first	45
desc_nulls_first	46
asc_nulls_last	46
desc_nulls_last	46
limit	46
first	47
take	47
collect	47
show	47
from_unixtime	48
to_date	48
datediff	48
explode	48
agg	48
avg	49
isin	49
range	49
<b>AGGREGATE FUNCTIONS:</b>	<b>52</b>
agg	52
avg/mean	52
count	53
max	53

min	53
pivot	53
sum	54
<b>JOINS:</b>	<b>55</b>
join	55
Inner Join	56
Full Outer Join	56
Right Outer join	57
Left Semi Join	57
Left Anti Join	58
Self Join	58
Join Algorithms in Spark	58
Broadcast Hash join	58
Sort Merge Join	59
Shuffle Hash join	60
<b>PERFORMANCE TUNING CONFIG OPTIONS</b>	<b>61</b>
Performance Tuning considerations	61
Improve spark performance	61
spark.sql.inMemoryColumnarStorage.compressed	61
spark.sql.inMemoryColumnarStorage.batchSize	61
spark.sql.files.maxPartitionBytes	61
spark.sql.files.openCostInBytes	61
spark.sql.files.minPartitionNum	62
spark.sql.broadcastTimeout	62
spark.sql.autoBroadcastJoinThreshold	62
spark.sql.shuffle.partitions	62
spark.sql.sources.parallelPartitionDiscovery.threshold	63
spark.sql.sources.parallelPartitionDiscovery.parallelism	63
<b>TUNING IN SPARK:</b>	<b>65</b>
Need of Tuning	65
Areas of tuning	65
Data serialization	65
Memory tuning	65
Memory management in spark.	65
Determining memory usage	66
Measuring impact of GC	66
Components of memory in JVM	66
Goal Of GC	66
Steps to avoid Full GC	67
Data locality	67

<b>References And Important Links for Spark:</b>	<b>68</b>
Advanced Apache Spark Training - Sameer Farooqui (Databricks)	68
Preparation material :	68
Useful Links:	68
<b>Interesting White papers:</b>	<b>68</b>
Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing	69
Spark SQL: Relational Data Processing in Spark	69
Discretized Streams: Fault-Tolerant Streaming Computation at Scale	69
Discretized Streams: An Efficient and Fault-Tolerant Model for Stream Processing on Large Clusters	69
Shark: Fast Data Analysis Using Coarse-grained Distributed Memory	69
Spark: Cluster Computing with Working Sets	69

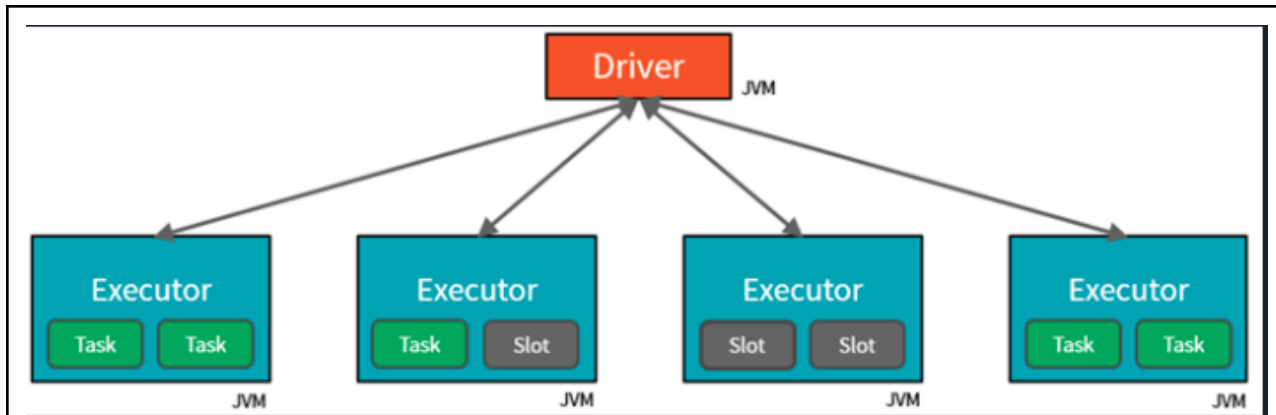
## SPARK CORE

Concept	Definition	Example/Code
Spark	Unified computing engine and a set of libraries for big data.	
Spark Components	<ul style="list-style-type: none"><li>• Driver</li><li>• Cluster Manager</li><li>• Executors</li><li>• Client</li></ul>	



Client	<ul style="list-style-type: none"> <li>• Responsible for starting the driver program.</li> <li>• Can be done via spark-submit or a spark-shell script or a custom application using spark api.</li> <li>• Prepares classpath and all application configuration for spark application.</li> <li>• Also passes required arguments</li> </ul>	
Driver	<ul style="list-style-type: none"> <li>• Orchestrates and monitors spark execution.</li> <li>• We have 1 driver per spark application.</li> <li>• Driver is responsible for requesting resources from the cluster manager.</li> <li>• Responsible for breaking application logic into stages and tasks.</li> <li>• Responsible for sending tasks to executors and Collecting results back.</li> <li>• Not responsible for scheduling tasks on executors</li> </ul>	





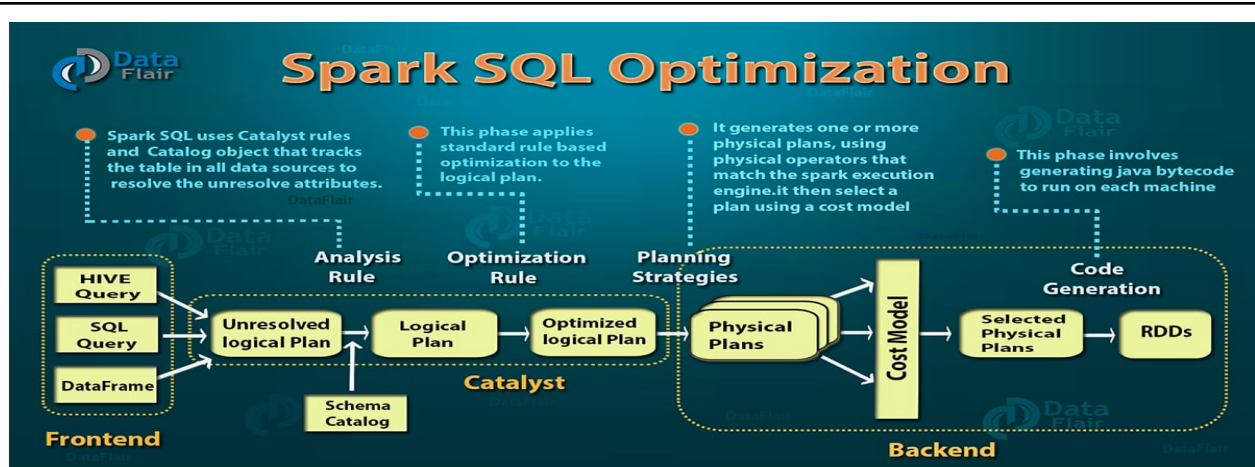
Executors	<ul style="list-style-type: none"> <li>• Accepts tasks from the driver, executes them and returns the results.</li> <li>• Each executor has several task slots or CPU cores(not physical).</li> <li>• Slots indicate the number of tasks that can be run in parallel.</li> <li>• We can set the task slots 2-3 times the number of actual physical CPU cores.</li> <li>• This is set by the <code>spark.executor.cores</code>.</li> </ul>	
Worker	<ul style="list-style-type: none"> <li>• Any node that can run application code in the cluster.</li> <li>• Has a finite or fixed number of Executors allocated at any point in time.</li> </ul>	
Cluster Manager	<ul style="list-style-type: none"> <li>• Keeps track of resources available.</li> <li>• Spark runs in local mode or inside a cluster.</li> <li>• Spark has below cluster managers: <ul style="list-style-type: none"> <li>○ Standalone <ul style="list-style-type: none"> <li>■ Most basic and default</li> <li>■ Can run only spark applications</li> </ul> </li> </ul> </li> </ul>	

	<ul style="list-style-type: none"> <li>○ Yarn <ul style="list-style-type: none"> <li>■ Hadoop's resource manager</li> <li>■ Can run spark and other java applications</li> </ul> </li> <li>○ Mesos <ul style="list-style-type: none"> <li>■ Scalable and fault tolerant.</li> <li>■ Can run spark as well several python,java and other applications</li> </ul> </li> <li>○ Kubernetes</li> </ul>	
Slots	<ul style="list-style-type: none"> <li>● Present inside an executor</li> <li>● Slots indicate the number of tasks that can be run in parallel.</li> </ul>	
Spark session	<ul style="list-style-type: none"> <li>● Sparksession instance is the handle through which spark executes spark user defined manipulations across cluster</li> <li>● There is 1-1 between spark session and spark application</li> <li>● Gives access to low level configs and contexts</li> <li>● subsume previous entry points to the Spark like the SparkContext, SQLContext, HiveContext, SparkConf, and StreamingContext</li> <li>● Spark Session allows you to create JVM runtime parameters, define DataFrames and Datasets, read from data sources, access catalog metadata, and issue Spark SQL queries.</li> </ul>	

Dataframes	<ul style="list-style-type: none"> <li>• Most common structured api in spark and represents a table of rows and columns in spark.</li> <li>• It can be considered as a spreadsheet with data across multiple computers</li> </ul>	
Partitions	<ul style="list-style-type: none"> <li>• A Partition is a logical chunk of your DataFrame</li> <li>• Data is split into Partitions so that each Executor core/thread can operate on a single part, enabling parallelization.</li> <li>• Spark by default creates 1 partition for every 128 MB of the file.</li> </ul>	<p>Eg: If you have 4 data partitions and you have 4 executor cores/threads, you can process everything in parallel, in a single pass.</p> <p>Syntax:</p> <pre>rdd = sc.parallelize(range(1,11)) rdd.getNumPartitions()  sc.defaultParallelism</pre>
Glom	<ul style="list-style-type: none"> <li>• returns an RDD having the elements within each partition in a separate list</li> </ul>	<code>rdd.glom().collect()</code>
Transformations	<ul style="list-style-type: none"> <li>• Dataframes are immutable.</li> <li>• Transformations are operations where we instruct Spark as to how to go about making changes to our dataframe.</li> <li>• Transformations are lazily evaluated.</li> <li>• They eventually return another DataFrame</li> <li>• There are two types of transformations: narrow and wide transformation</li> </ul>	

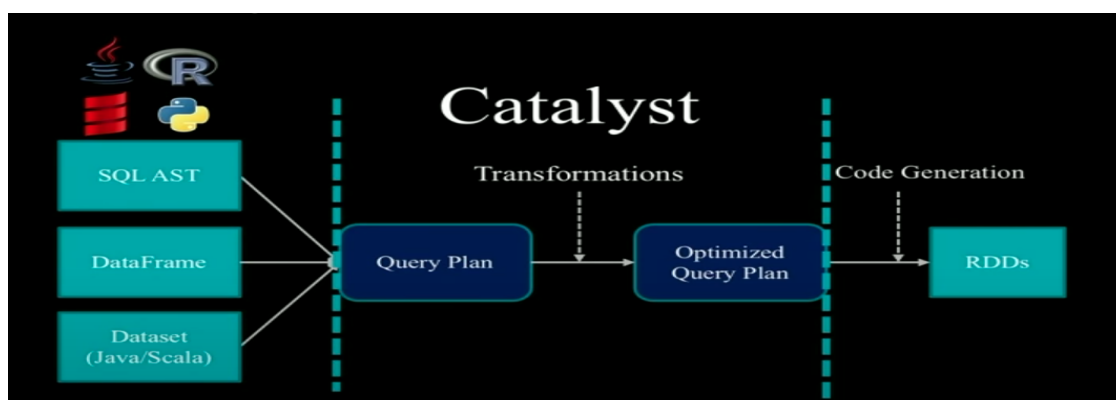
Narrow transformations	<ul style="list-style-type: none"> <li>• In this transformation, we do not perform a shuffle on the parent rdd to get a new rdd.</li> <li>• The data required to compute the records in a single partition in child rdd resides at most one partition of the parent RDD.</li> <li>• narrow dependencies can be grouped into one stage.</li> <li>• If we partition buckets for a column on which we want to perform joins, we can avoid shuffle</li> </ul>	<ul style="list-style-type: none"> <li>• Map</li> <li>• Mapvalues</li> <li>• Filter</li> <li>• Union</li> <li>• Flatmap</li> <li>• Mappartitions</li> <li>• Join</li> </ul>
Wide transformation	<ul style="list-style-type: none"> <li>• In this transformation, we perform a shuffle on the parent rdd to get new rdd</li> <li>• The data required to compute the records in a single partition in child rdd resides at more than one partition of the parent RDD.</li> </ul>	<ul style="list-style-type: none"> <li>• groupBy</li> <li>• leftouterjoin</li> <li>• Groupbykey</li> <li>• Reducebykey</li> <li>• Distinct</li> <li>• Intersection</li> <li>• repartition</li> </ul>
Action	<ul style="list-style-type: none"> <li>• An action instructs spark to compute a result from a series of transformations.</li> <li>• 3 kinds of actions: <ul style="list-style-type: none"> <li>○ Action to view data on console.</li> <li>○ Action to collect data to native objects.</li> <li>○ Action to write to an output data source.</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• count</li> <li>• collect</li> <li>• take(n)</li> <li>• top()</li> <li>• countByValue()</li> <li>• reduce()</li> <li>• fold()</li> <li>• aggregate()</li> <li>• foreach()</li> </ul>
Lazy evaluation	<ul style="list-style-type: none"> <li>• Spark pipelines all the transformations and performs the actual execution at the end.</li> <li>• This is called lazy evaluation.</li> <li>• Done to improve performance and</li> </ul>	

	optimize execution	
Datasets	<ul style="list-style-type: none"> <li>Statistically typed code available in java and python.</li> </ul>	
DAG	<ul style="list-style-type: none"> <li>DAG stands for directed acyclic graph.</li> <li>Has no directed circles.</li> <li>Used to construct execution plans for spark code, form different stages and optimize on the same.</li> </ul>	
Spark Lifecycle (Explain Plan)	<ul style="list-style-type: none"> <li>User submits application to spark cluster.</li> <li>Driver takes spark code and identifies transformations and actions,</li> <li>Logical plan(consisting of tasks and stages) is constructed here which is passed through an optimizer and a physical plan is constructed.</li> <li>In physical plan stages are combined based on the nature of transformations.</li> <li>Wide transformations which require a shuffle to be performed , and any reads from external sources mark the start of a new stage.</li> </ul>	Analysis->logical optimization->physical planning->code generation



Catalyst optimizer

- Interface between high level structured API's and low level RDD's.
- Converts initial query plan to optimized version.
- The user programs are defined as trees and what catalyst does is convert one tree to another.
- There are 2 kinds of query plans:-
  - Logical plan
  - Physical plan.



Catalyst optimizer trees

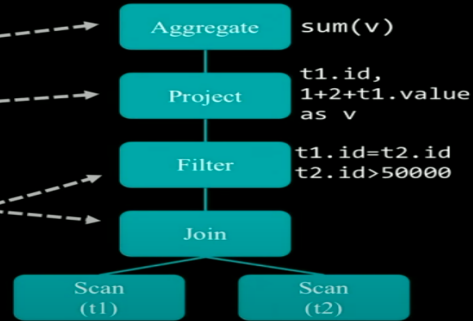
- Abstraction of a user program.

# Trees: Abstractions of Users' Programs

## Query Plan

```

SELECT sum(v)
FROM (
  SELECT
    t1.id,
    1 + 2 + t1.value AS v
  FROM t1 JOIN t2
  WHERE
    t1.id = t2.id AND
    t2.id > 50000) tmp
    
```



Logical plan	<ul style="list-style-type: none"> <li>Abstraction of user programs without execution details.</li> <li>Has list of output columns and property called constraints</li> </ul>	
Physical plan	<ul style="list-style-type: none"> <li>Describes types of computation</li> </ul>	
Types of transformations in optimizer	<ul style="list-style-type: none"> <li>Expression to expression: Techniques used are:                             <ul style="list-style-type: none"> <li>Predicate pushdown(pushing filter downstream so that less data is read)</li> <li>Column pruning - get only required columns from DB.</li> </ul> </li> <li>Logical plan to logical plan</li> <li>Logical plan to physical plan</li> </ul>	
Dynamic partition pruning	<ul style="list-style-type: none"> <li>Intended to skip over data you do not need in the results of a query</li> <li>Done at Logical planning level to find the dimensional filter and propagate across the join to the other side of the scan.</li> <li>Done at Physical level to wire it together in a way that this filter executes only once on the dimension side.</li> </ul>	spark.sql.optimizer. dynamicPartitionPruning .enabled

	<ul style="list-style-type: none"> <li>Then the results of the filter gets reused directly in the scan of the table. And with this two fold approach we can achieve significant speed ups in many queries in Spark.</li> <li>Best use case is for joining fact and dimension tables</li> </ul>	
Adaptive Query Execution	<ul style="list-style-type: none"> <li>Cost optimization is done based on initial statistics which is different from runtime statistics.</li> <li>Aqe aims to optimize on the same.</li> <li>This is done at stage boundaries.</li> <li>Not applicable for streaming queries</li> </ul>	spark.sql.adaptive.enabled=true needs to be set
AQE Scenarios	<ul style="list-style-type: none"> <li>Data Skew - partitions are uneven which leads to longer execution time, aqe has a skew partitioner which breaks down larger partitions.</li> <li>Dynamically changing partitions - if aqe feels less partitions are needed than what is allocated it does a coalesce.</li> <li>Introducing broadcast joins instead of sort merge joins - if we are joining 2 tables and 1 is smaller than 10mb we can convert the join to a broadcast thereby reducing shuffles.</li> </ul>	<ul style="list-style-type: none"> <li>spark.sql.adaptive.coalescePartitions.enabled</li> <li>spark.sql.adaptive.coalescePartitions.minPartitionNum</li> <li>spark.sql.adaptive.coalescePartitions.initialPartitionNum</li> <li>spark.sql.adaptive.advisoryPartitionSizeInBytes</li> <li>spark.sql.adaptive.localShuffleReader.enabled - convert sort merge join to broadcast join</li> </ul>
Spark submit	<ul style="list-style-type: none"> <li>Spark submit command is used to submit a job to a spark cluster once an application is packaged.</li> <li>Arguments of spark submit are</li> </ul>	<p>Syntax:</p> <pre>/bin/spark-submit \ --class &lt;main-class&gt; \</pre>



	<p>below.</p> <ul style="list-style-type: none"> <li>• <code>--class</code>: The entry point for your application (e.g. <code>org.apache.spark.examples.SparkPi</code>)</li> <li>• <code>--master</code>: The master url for the cluster (e.g. <code>spark://23.195.26.187:7077</code>)</li> <li>• <code>--deploy-mode</code>: Whether to deploy your driver on the worker nodes (cluster) or locally as an external client (client) (default: client) †</li> <li>• <code>--conf</code>: Arbitrary Spark configuration property in <code>key=value</code> format. For values that contain spaces, wrap “key=value” in quotes (as shown).</li> <li>• <code>application-jar</code>: Path to a bundled jar including your application and all dependencies. The URL must be globally visible inside of your cluster, for instance, an <code>hdfs://</code> path or a <code>file://</code> path that is present on all nodes.</li> <li>• <code>application-arguments</code>: Arguments passed to the main method of your main class, if any.</li> </ul>	<pre>--master &lt;master-url&gt; \ --deploy-mode &lt;deploy-mode&gt; \ --conf &lt;key&gt;=&lt;value&gt; \  ... # other options &lt;application-jar&gt; \ [application-arguments]</pre>
Spark submit master Url	<p>Spark master can have the below options:</p> <ul style="list-style-type: none"> <li>• <b>Local</b> - Run Spark locally with one worker thread (i.e. no parallelism at all)</li> <li>• <b>local[K]</b>- Run Spark locally with K worker threads (ideally, set this to the number of cores on your machine).</li> <li>• <b>local[K,F]</b> - Run Spark locally with K worker threads and F maxFailures</li> </ul>	

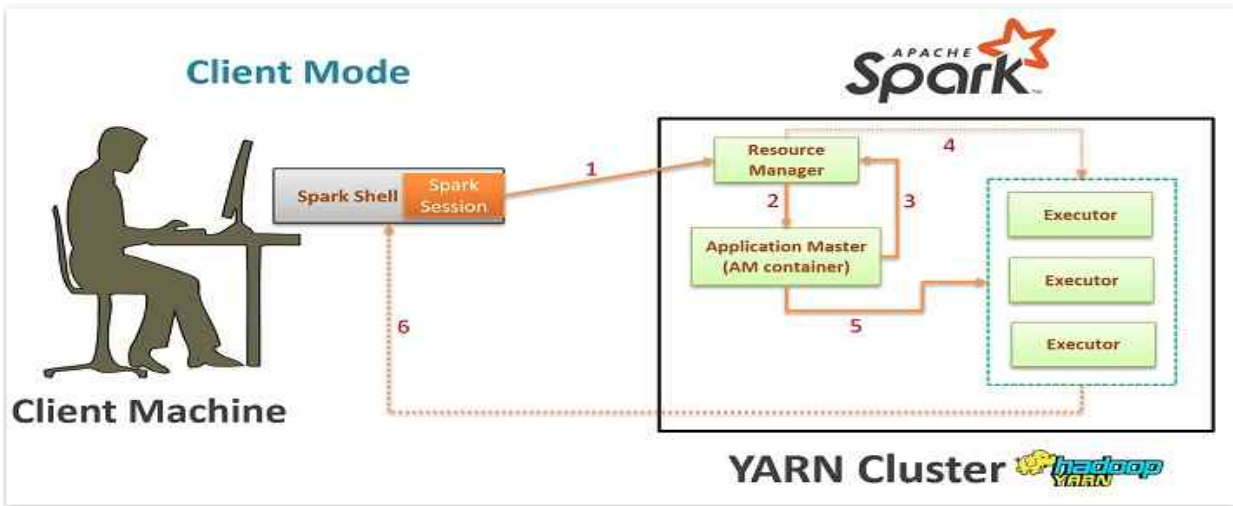
	<p>(see spark.task.maxFailures for an explanation of this variable)</p> <ul style="list-style-type: none"> <li>● <b>local[*]</b> - Run Spark locally with as many worker threads as logical cores on your machine.</li> <li>● <b>local[*,F]</b> - Run Spark locally with as many worker threads as logical cores on your machine and F maxFailures.</li> <li>● <b>spark://HOST:PORT</b> - Connect to the given Spark standalone cluster master. The port must be whichever one your master is configured to use, which is 7077 by default.</li> <li>● <b>spark://HOST1:PORT1,HOST2:PORT2</b> - Connect to the given Spark standalone cluster with standby masters with Zookeeper. Port :7077 by default</li> <li>● <b>mesos://HOST:PORT</b>- Connect to the given Mesos cluster.</li> <li>● <b>Yarn</b> - Connect to a YARN cluster in client or cluster mode depending on the value of --deploy-mode.</li> </ul>	
Dag scheduler	<ul style="list-style-type: none"> <li>● DAGScheduler is the scheduling layer of Apache Spark that implements stage-oriented scheduling.</li> <li>● It transforms a logical execution plan to a physical execution plan (using stages).</li> <li>● Uses a fifo scheduler</li> <li>● Fair scheduler runs jobs in round robin fashion</li> <li>● To enable fair scheduler set spark.scheduler.mode= Fair</li> </ul>	
Task scheduler	<ul style="list-style-type: none"> <li>● These schedulers get sets of tasks submitted to them from the DAGScheduler for each stage, and are responsible for sending the tasks to the cluster, running them,</li> </ul>	

	<p>retrying if there are failures, and mitigating stragglers.</p> <ul style="list-style-type: none"> <li>• They return events to the DAGScheduler.</li> </ul>	
Jobs	<ul style="list-style-type: none"> <li>• A Job is a sequence of stages, triggered by an action such as count(), collect(), read() or write().</li> <li>• Each parallelized action is referred to as a Job.</li> <li>• The results of each Job (parallelized/distributed action) are returned to the Driver from the Executor.</li> <li>• Depending on the work required, multiple Jobs will be required.</li> </ul>	
Stages	<ul style="list-style-type: none"> <li>• Each job that gets divided into smaller sets of tasks is a stage.</li> <li>• Stage is a step in a physical execution plan</li> <li>• A stage is a set of parallel tasks - one task per partition</li> <li>• Each stage contains a sequence of narrow transformations</li> </ul>	
Tasks	<ul style="list-style-type: none"> <li>• A task is a unit of work that is sent to the executor.</li> <li>• Each stage has some tasks, one task per partition.</li> <li>• The same task is done over different partitions of the RDD.</li> </ul>	
Spark config	<ul style="list-style-type: none"> <li>• spark.executor.instances: Number of executors for the spark application.</li> <li>• spark.executor.memory: Amount of memory to use for each executor</li> </ul>	

	<p>that runs the task.</p> <ul style="list-style-type: none"> <li>• spark.executor.cores: Number of concurrent tasks an executor can run.</li> <li>• spark.driver.memory: Amount of memory to use for the driver.</li> <li>• spark.driver.cores: Number of virtual cores to use for the driver process.</li> <li>• spark.sql.shuffle.partitions: Number of partitions to use when shuffling data for joins or aggregations.</li> <li>• spark.default.parallelism: Default number of partitions in resilient distributed datasets (RDDs) returned by transformations like join and aggregations.</li> <li>• Yarn has different configuration settings which can be looked at as well.</li> </ul>	
Additional config properties	<ul style="list-style-type: none"> <li>• spark.speculation - relaunch one or more tasks if they are running slowly in a stage.</li> <li>• spark.memory.fraction - <ul style="list-style-type: none"> <li>○ percentage of memory used for computation in shuffles, joins, sort and aggregations</li> <li>○ If gc is invoked multiple times before task completed, decrease this so that amount of memory used for caching is reduced</li> </ul> </li> </ul>	
Dynamic allocation	<ul style="list-style-type: none"> <li>• scales the number of executors registered with this application up and down based on the workload.</li> </ul>	<ul style="list-style-type: none"> <li>• spark.dynamicAllocation.enabled</li> <li>• spark.dynamicAllocation.minExecutors</li> </ul>

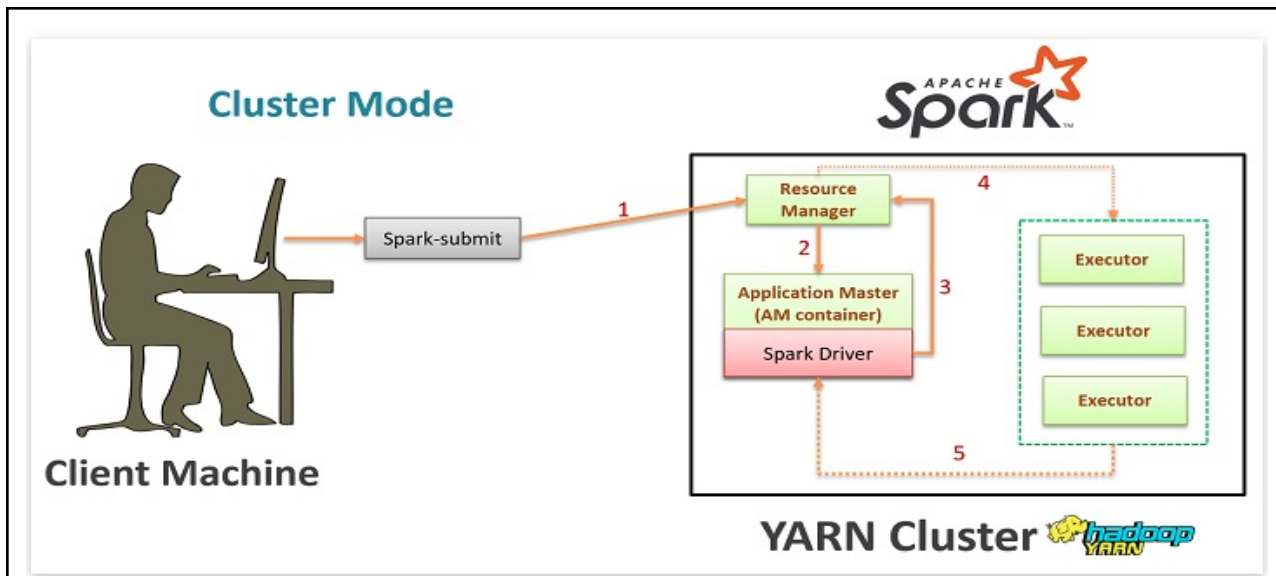
		<ul style="list-style-type: none"> <li>• spark.dynamicAllocation.maxExecutors</li> <li>• spark.dynamicAllocation.initialExecutors</li> </ul>
Execution modes	<p>Spark execution can be in one of 3 modes:</p> <ul style="list-style-type: none"> <li>• Local mode.</li> <li>• Client Mode</li> <li>• Cluster mode</li> </ul>	
Local Mode	<ul style="list-style-type: none"> <li>• This is like running a program on someone's laptop or desktop using a single JVM.</li> <li>• You should have defined &amp; used spark context objects, imported spark libraries and processed data from your local system files.</li> <li>• Everything runs LOCALLY and there is no concept of NODE involved and nothing runs in DISTRIBUTED mode.</li> <li>• The driver &amp; the executor are created inside a single JVM process.</li> </ul>	
Client mode	<ul style="list-style-type: none"> <li>• In Client mode, the Driver is started in the Local machine i.e. Driver is outside of the Cluster.</li> <li>• The Executors will be running inside the Cluster.</li> <li>• Entire application is dependent on the Local machine since the Driver resides in here.</li> <li>• In case of any issue with the local machine , the driver will go off . Subsequently the entire application will go off.</li> </ul>	

- Good for Debugging or Testing.



Cluster mode

- In Cluster Mode, the Driver & Executor both run inside the Cluster.
- You submit the spark job from your local machine to a Cluster machine inside the Cluster
- Such machines are usually called Edge Nodes.
- This is the approach used in Production use cases.



Serialization	<ul style="list-style-type: none"> <li>There are 2 types of serialization supported in spark: <ul style="list-style-type: none"> <li>Java serialization: default serialization, easy but slow.</li> <li>Kryo Serialization - <ul style="list-style-type: none"> <li>10 times faster than default serialization</li> <li>Requires custom registration of classes.</li> <li>Best in case of network intensive applications.</li> </ul> </li> </ul> </li> <li>All data saved to disk in spark is serialized by default</li> </ul>	<p>Kryo:</p> <pre>conf.set("Spark.serializer", "org.apache.spark.serializer.kryoserializer")</pre>
Shuffle	<ul style="list-style-type: none"> <li>A Shuffle refers to an operation where data is re-partitioned across a Cluster - i.e. when data needs to move between executors(comparison done).</li> <li>Join and any operation that ends with ByKey(wide transformation)</li> </ul>	

	<p>will trigger a Shuffle.</p> <ul style="list-style-type: none"> <li>It is a costly operation because a lot of data can be sent via the network.</li> <li>Types of shuffle implementations in spark are: <ul style="list-style-type: none"> <li>Bypassmergesortshufflewriter</li> <li>Sortshufflewriter</li> <li>Unshuffle writer</li> </ul> </li> </ul>	
Caching	<ul style="list-style-type: none"> <li>Shuffle files are by definition temporary files and will eventually be removed.</li> <li>However, we can cache data explicitly to accomplish the same thing that happens inadvertently with shuffle files.</li> <li>This is used for optimizing the spark query.</li> <li>LRU cache is used by default.</li> <li>Cache is eager by default, we can change this by using keyword lazy.</li> </ul>	<pre>spark.catalog.cacheTable("tableName")</pre> <pre>dataFrame.cache()</pre> <p>Unpersist() - used to forcibly remove rdd from head</p> <pre>spark.sql("cache lazy table tableName")</pre>
Uncache Table	<ul style="list-style-type: none"> <li>Removes the entries and associated data from the in-memory and/or on-disk cache for a given table or view.</li> <li>The underlying entries should already have been brought to cache by the previous CACHE TABLE operation.</li> <li>The UNCACHE TABLE on a non-existent table throws an Exception if IF EXISTS is not specified.</li> </ul>	<p><b>Syntax</b></p> <pre>UNCACHE TABLE [ IF EXISTS ] table_name</pre> <p><b>Parameters</b></p> <ul style="list-style-type: none"> <li>var: table_name</li> <li>def: The name of the table or view to be uncached.</li> </ul> <p><b>Example:</b></p> <pre>UNCACHE TABLE t1;</pre>
Persist Options	<ul style="list-style-type: none"> <li>MEMORY_ONLY-</li> </ul>	



	<ul style="list-style-type: none"> <li>○ rdd stored in memory</li> <li>● MEMORY_ONLY_SER - <ul style="list-style-type: none"> <li>○ rdd stored in memory</li> <li>○ Cost of serialization/deserialization incurred</li> </ul> </li> <li>● MEMORY_AND_DISK <ul style="list-style-type: none"> <li>○ Default storage option for cache</li> <li>○ First we store in memory, if space is not available we push the oldest partition to disk</li> </ul> </li> <li>● MEMORY_AND_DISK_SER <ul style="list-style-type: none"> <li>○ Above with data serialized</li> </ul> </li> <li>● MEMORY_AND_DISK_DESER <ul style="list-style-type: none"> <li>○ Default storage level of spark since 3.1.1</li> </ul> </li> <li>● MEMORY_ONLY_2 <ul style="list-style-type: none"> <li>○ 2 replicas created of the rdd</li> <li>○ Do only if rdd is critical</li> </ul> </li> <li>● DISK_ONLY <ul style="list-style-type: none"> <li>○ Keep replication factor to 1</li> </ul> </li> </ul>	
RDD	<ul style="list-style-type: none"> <li>● Rdd represents an immutable partitioned collection of records that can be operated on in parallel.</li> <li>● RDDs are java,python or scala objects.</li> <li>● Dataframes and datasets are internally converted into rdd's</li> </ul>	
Types of RDD	<ul style="list-style-type: none"> <li>● You have rdd interfaces and multiple implementations of the</li> </ul>	

	<p>same.</p> <ul style="list-style-type: none"> <li>• Things that we specify in an rdd implementation are <ul style="list-style-type: none"> <li>◦ Set of partitions</li> <li>◦ List of dependencies on parent RDD's</li> <li>◦ Function to compute a partition given to a parent.</li> <li>◦ Optional preferred locations.</li> <li>◦ Optional partition information.</li> <li>◦ Broadly there are two types of rdd's:- <ul style="list-style-type: none"> <li>■ Generic RDD -</li> <li>■ Key Value RDD custom partitions by keys are done.</li> </ul> </li> </ul> </li> </ul>	
Broadcast Variables	<ul style="list-style-type: none"> <li>• Shared immutable variable cached on every machine on the cluster</li> <li>• Present in the driver and sent across clusters.</li> <li>• Broadcastjoin is used when we are joining 2 tables and want to broadcast a complete table(size limit 10mb)</li> <li>• No shuffle incurred</li> </ul>	
Accumulators	<ul style="list-style-type: none"> <li>• Accumulators are used to write data across a cluster.</li> <li>• Eg. We are reading a file in chunks and we want to count the number of blank lines.</li> <li>• We can use an accumulator to do this.</li> </ul>	

	<ul style="list-style-type: none"> <li>Executors just update the value.</li> </ul>	
Repartition	<ul style="list-style-type: none"> <li>Repartition is a method in spark which is used to perform a full shuffle on the data present and create partitions based on the user's input. <ul style="list-style-type: none"> <li>Method 1 : Repartition using Column Name</li> <li>Method 2 : Repartition using integer value</li> </ul> </li> <li>By default, 200 partitions are created if the number is not specified in the repartition clause.</li> <li>Use Repartition only when you want to increase the number of partitions (or) if you want to perform a full shuffle on the data.</li> </ul>	<ul style="list-style-type: none"> <li>Repartition using Column Name  <code>df = df.repartition('_1')</code> </li> <li>Repartition using integer value  <code>df = df.repartition(3)</code> </li> </ul>
Coalesce	<ul style="list-style-type: none"> <li>Used to reduce the number of partitions in a dataframe.</li> <li>We cannot increase the number of partitions using coalesce</li> <li>Unlike repartition, coalesce doesn't perform a shuffle to create the partitions.</li> <li>Suppose there are 100 partitions with 10 records in each partition, and if the partition size is reduced to 50, it would retain 50 partitions and append the other values to these existing partitions thereby having 20 records in each partition.</li> </ul>	<code>df3 = df.coalesce(2)</code>

Repartition vs Coalesce	<ul style="list-style-type: none"> <li>• Repartition can be used under these scenarios: <ul style="list-style-type: none"> <li>○ when you want your output partitions to be of equally distributed chunks.</li> <li>○ increase the number of partitions.</li> <li>○ perform a shuffle of the data and create partitions.</li> </ul> </li> <li>• Coalesce can be used under the following scenarios, <ul style="list-style-type: none"> <li>○ when you want to decrease the number of partitions.</li> <li>○ in order to avoid shuffle when partitions are decreasing.</li> </ul> </li> </ul>	
Catalog	<ul style="list-style-type: none"> <li>• Catalog is the interface for managing a metastore (aka metadata catalog) of relational entities (e.g. database(s), tables, functions, table columns and temporary views).</li> <li>• Catalog is available using <code>SparkSession.catalog</code> property.</li> </ul>	

# SPARK READ/WRITE

Concept	Definition	Example/Code
Generic File options	<ul style="list-style-type: none"> <li>• Ignore Corrupt Files</li> <li>• Ignore Missing Files</li> <li>• pathGlobFilter- include files with file names matching the pattern.</li> <li>• recursiveFileLookup - used to recursively load files and it disables partition inferring. Its default value is false</li> </ul>	<pre>spark.sql ("set spark.sql.files.ignoreCorruptFiles=true")  spark.sql.files.ignoreMissingFiles  spark.read. load("examples/src/main/resources/dir1", format="parquet", pathGlobFilter="*.parquet")  recursive_loaded_df = spark.read.format("parquet") .option("recursiveFileLookup", "true") .load("examples/src/main/resources/dir1")</pre>
Read	<ul style="list-style-type: none"> <li>• The read modes are: <ul style="list-style-type: none"> <li>○ permissive — All fields are set to null and corrupted records are placed in a string column called <code>_corrupt_record</code></li> <li>○ dropMalformed — Drops all rows containing corrupt records.</li> </ul> </li> </ul>	<p>Syntax:</p> <pre>DataFrameReader.format(...).option("key", "value").schema(...).load()</pre> <p><b>FOSL</b></p>

	<ul style="list-style-type: none"> <li>○ failFast — Fails when corrupt records are encountered.</li> <li>● If you want to control the maximum partition size while reading files. <ul style="list-style-type: none"> <li>○ spark.files.maxpartitionBytes</li> </ul> </li> </ul>	
Write	<ul style="list-style-type: none"> <li>● DataFrame writer default format is the parquet file</li> </ul> <p>Save modes are:</p> <ul style="list-style-type: none"> <li>● append — appends output data to files that already exist</li> <li>● overwrite — completely overwrites any data present at the destination</li> <li>● errorIfExists — Spark throws an error if data already exists at the destination</li> <li>● ignore — if data exists do nothing with the dataframe</li> </ul>	<p>Syntax:</p> <pre>DataFrameWriter.format(...).option(...).partitionBy(...).bucketBy(...).sortBy( ...).save()</pre> <p>FOPBSS</p>
CSV		<p><b>Read:</b></p> <pre>csvFile = spark.read.format("json") .option("mode", "FAILFAST") .option("inferSchema", "true") .load("/data/flight-data/json/2010-summary.json") .show(5)</pre> <p><b>Write:</b></p> <pre>csvFile.write.format("csv")</pre>

		<pre>.mode("overwrite") .option("sep", "\t"). .save("/tmp/my-tsv-file.tsv")</pre>
Json	<ul style="list-style-type: none"> <li>Spark SQL can automatically infer the schema of a JSON dataset and load it as a Dataset[Row].</li> <li>Each line must contain a separate, self-contained valid JSON object</li> <li>For a regular multi-line JSON file, set the multiLine option to true.</li> <li>Default date format in json is yyyy-MM-dd</li> </ul>	<p><b>Read:</b></p> <pre>spark.read.format("json") .option("mode", "FAILFAST") .option("inferSchema", "true") .load("/data/flight-data/json/2010-summary.json").show(5)</pre> <p><b>Write:</b></p> <pre>csvFile.write.format("json") .mode("overwrite") .save("/tmp/my-json-file.json")</pre>
Text	<ul style="list-style-type: none"> <li>When you write a text file, you need to be sure to have only one string column; otherwise, the write will fail</li> </ul>	<p><b>Read:</b></p> <pre>lines = sc.textFile("/home/deepak/test1.txt")</pre> <p><b>Write:</b></p> <pre>lines.coalesce(1).write.format("text") .option("header", "false") .mode("append").save("output.txt")</pre>
Parquet	<ul style="list-style-type: none"> <li>Parquet is a columnar format that is supported by many other data processing systems</li> <li>When reading Parquet files, all columns are automatically converted to be nullable for compatibility reasons.</li> </ul>	<p><b>Read:</b></p> <pre>spark.read.format("parquet") .load("/data/flight-data/parquet/summary.parquet") .show(5)</pre> <p><b>Write:</b></p> <pre>csvFile.write.format("parquet") .mode("overwrite") .save("/tmp/my-parquet-file.parquet")</pre>

Orc		<p><b>Read:</b></p> <pre>spark.read.format("orc") .load("/data/flight-data/orc/2010-summary.orc") .show(5)</pre> <p><b>Write:</b></p> <pre>csvFile.write.format("orc") .mode("overwrite") .save("/tmp/my-json-file.orc")</pre>
Database		<pre>driver = "org.sqlite.JDBC" path = "/data/flight-data/jdbc/my-sqlite.db" url = "jdbc:sqlite:" + path tablename = "flight_info"</pre> <p><b>Read:</b></p> <pre>dbDataFrame = spark.read.format("jdbc") .option("url", url) .option("dbtable", tablename) .option("driver", driver).load()</pre> <p><b>Write:</b></p> <pre>newPath = "jdbc:sqlite://tmp/my-sqlite.db" csvFile.write .jdbc(newPath, tablename, mode="overwrite", properties=props)</pre>



# TABLES AND VIEWS

Global Managed Table	<ul style="list-style-type: none"> <li>• A managed table is a Spark SQL table for which Spark manages both the data and the metadata.</li> <li>• A global managed table is available across all clusters.</li> <li>• When you drop the table both data and metadata get dropped.</li> </ul>	<code>dataframe.write.saveAsTable("my_table")</code>
Global Unmanaged/External Table	<ul style="list-style-type: none"> <li>• Spark manages the metadata, while you control the data location. As soon as you add the 'path' option in the dataframe writer it will be treated as a global external/unmanaged table.</li> <li>• When you drop a table only metadata gets dropped.</li> <li>• A global unmanaged/external table is available across all clusters.</li> </ul>	<code>Dataframe.write.option('path', "&lt;your-storage-path&gt;").saveAsTable("my_table")</code>
Local Table (a.k.a) Temporary Table (a.k.a) Temporary View	<ul style="list-style-type: none"> <li>• Spark session scoped.</li> <li>• A local table is not accessible from other clusters (or if using a databricks notebook not in other notebooks as well).</li> <li>• It is not registered in</li> </ul>	<code>Dataframe.createOrReplaceTempView()</code>

	the metastore.	
Global Temporary View	<ul style="list-style-type: none"> <li>Spark application scoped global temporary views are tied to a system preserved temporary database <code>global_temp</code>.</li> <li>This view can be shared across different spark sessions (or if using databricks notebooks, then shared across notebooks).</li> </ul>	<pre>dataframe.createOrReplaceGlobalTempView("my_global_view")</pre> <p>can be accessed as,</p> <ul style="list-style-type: none"> <li><code>spark.read.table("global_temp.my_global_view")</code></li> </ul>
Global Permanent View	<ul style="list-style-type: none"> <li>Persist a data frame as a permanent view.</li> <li>The view definition is recorded in the underlying metastore.</li> <li>You can only create a permanent view on a global managed table or global unmanaged table. Not allowed to create a permanent view on top of any temporary views or dataframe.</li> <li>Note: Permanent views are only available in SQL API — not available in dataframe API</li> </ul>	<pre>spark.sql("CREATE VIEW permanent_view AS SELECT * FROM table")</pre>

# DATAFRAMES

Function	Use	Example
CreateDataFrame	<ul style="list-style-type: none"> <li>Creates a DataFrame from an RDD, a list or a pandas.DataFrame</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.session</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>createDataFrame(data, schema=None,samplingRatio=None, verifySchema=True)</pre> <ul style="list-style-type: none"> <li>data: <ul style="list-style-type: none"> <li>type:dataRDD or iterable</li> <li>def: an RDD of any kind of SQL data representation</li> </ul> </li> <li>schema: <ul style="list-style-type: none"> <li>type: pyspark.sql.types.DataType, str or list, optional</li> <li>def:</li> </ul> </li> <li>samplingRatio: <ul style="list-style-type: none"> <li>type:float, optional</li> <li>def:the sample ratio of rows used for inferring</li> </ul> </li> <li>verifySchema <ul style="list-style-type: none"> <li>type:bool, optional</li> <li>def:verify data types of every row against schema. Enabled by default.</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>l = [('Alice', 1)] spark.createDataFrame(l).collect()</pre>

printSchema	<ul style="list-style-type: none"> <li>Prints out the schema in the tree format.</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.DataFrame</li> </ul> </li> </ul>	<b>Syntax:</b> <pre>DataFrame.printSchema()</pre> <b>Example:</b> <pre>df.printSchema()</pre>
createOrReplaceTempView	<ul style="list-style-type: none"> <li>Creates or replaces a local temporary view with this DataFrame.</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.DataFrame</li> </ul> </li> </ul>	<b>Syntax:</b> <pre>DataFrame.createOrReplaceTempView(name)</pre> <b>Example:</b> <pre>df.createOrReplaceTempView("people")</pre>
Schema	<ul style="list-style-type: none"> <li>Returns the schema of this DataFrame as a pyspark.sql.types.StructType.</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.DataFrame</li> </ul> </li> </ul>	<b>Syntax:</b> <pre>DataFrame.schema</pre> <b>Example:</b> <pre>df.schema</pre>
col	<ul style="list-style-type: none"> <li>A column in a DataFrame.</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.functions</li> </ul> </li> </ul>	<b>Syntax:</b> <b>Example:</b> <ul style="list-style-type: none"> <li>Select a column out of a DataFrame <pre>df.colName</pre> <pre>df["colName"]</pre> </li> <li>Create from an expression <pre>df.colName + 1</pre> </li> </ul>

		<code>1 / df.colName</code>
row	<ul style="list-style-type: none"> <li>• A row in DataFrame.</li> <li>• The fields in it can be accessed: <ul style="list-style-type: none"> <li>◦ like attributes (row.key)</li> <li>◦ like dictionary values (row[key])</li> </ul> </li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <p><b>Example:</b></p> <pre>row = Row(name="Alice", age=11) row.name, row.age</pre>
Null in spark	Use None	
expr	<ul style="list-style-type: none"> <li>• Parses the expression string into the column that it represents</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.functions</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>expr(str)</pre> <p><b>Example:</b></p> <pre>df.select(expr("length(name)")).collect()</pre>
select	<ul style="list-style-type: none"> <li>• Projects a set of expressions and returns a new DataFrame.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.select(*cols)</pre> <ul style="list-style-type: none"> <li>• cols <ul style="list-style-type: none"> <li>◦ type:str, Column, or list</li> <li>◦ def:</li> </ul> </li> </ul>

	<p><b>DataFrame</b></p>	<ul style="list-style-type: none"> <li>■ column names (string) or expressions (Column).</li> <li>■ If one of the column names is '*', that column is expanded to include all columns in the current DataFrame.</li> </ul> <p><b>Example:</b></p> <pre>df.select('name', 'age').collect()  df.select(df.name, (df.age + 10).alias('age')).collect()</pre>
selectExpr	<ul style="list-style-type: none"> <li>• Projects a set of SQL expressions and returns a new DataFrame.</li> <li>• This is a variant of select() that accepts SQL expressions.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.DataFrame</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.selectExpr(*expr)</pre> <p><b>Example:</b></p> <pre>df.selectExpr("age * 2", "abs(age)").collect()</pre>
withColumn	<ul style="list-style-type: none"> <li>• Returns a new DataFrame by adding a column or replacing the existing column that has the same name.</li> <li>• The column expression must be an expression over</li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.withColumn(colName, col)</pre> <ul style="list-style-type: none"> <li>• colName <ul style="list-style-type: none"> <li>◦ Type:str</li> <li>◦ def: string, name of the new column.</li> </ul> </li> <li>• col <ul style="list-style-type: none"> <li>◦ type:Column</li> </ul> </li> </ul>

	<p>this DataFrame; attempting to add a column from some other DataFrame will raise an error.</p> <ul style="list-style-type: none"> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.DataFrame</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>def: a Column expression for the new column.</li> </ul> <p><b>Example:</b></p> <pre>df.withColumn('age2', df.age + 2).collect()</pre>
withColumnRenamed	<ul style="list-style-type: none"> <li>Returns a new DataFrame by renaming an existing column. This is a no-op if the schema doesn't contain the given column name.</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.DataFrame</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.withColumnRenamed(existing, new)</pre> <ul style="list-style-type: none"> <li>existing: <ul style="list-style-type: none"> <li>Type:str</li> <li>def:string, name of the existing column to rename.</li> </ul> </li> <li>new: <ul style="list-style-type: none"> <li>Type:str</li> <li>def:string, new name of the column.</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.withColumnRenamed('age', 'age2').collect()</pre>
alias	<ul style="list-style-type: none"> <li>Returns this column aliased with a new name or names.</li> <li>In the case of expressions that return more than one column, such as explode).</li> <li>package: <ul style="list-style-type: none"> <li>pyspark.sql.</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>Column.alias(*alias, **kwargs)</pre> <pre>alias type: str def: desired column names (collects all positional arguments passed)</pre> <p><b>Example:</b></p>

	column	<pre>df.select(df.age.alias("age2")).collect()</pre>
cast	<ul style="list-style-type: none"> <li>Convert the column into a different dataType.</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.Column</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>Column.cast(dataType)</pre> <p><b>Example:</b></p> <pre>df.select(df.age.cast("string").alias('ages')).collect()</pre>
drop	<ul style="list-style-type: none"> <li>Returns a new DataFrame that drops the specified column.</li> <li>This is a no-op if the schema doesn't contain the given column name(s).</li> <li>package: <ul style="list-style-type: none"> <li>pyspark.sql.DataFrame</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.drop(*cols)</pre> <ul style="list-style-type: none"> <li>cols: <ul style="list-style-type: none"> <li>type:str or :class:`Column`</li> <li>def: a name of the column, or the Column to drop</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.drop('age').collect()</pre> <pre>df.drop(df.age).collect()</pre>
dropDuplicates	<ul style="list-style-type: none"> <li>Return a new DataFrame with duplicate rows removed, optionally only considering certain columns.</li> <li>For a static batch DataFrame, it just drops duplicate rows.</li> <li>For a streaming DataFrame, it will keep all data</li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.dropDuplicates(subset=None)</pre> <p><b>Example:</b></p> <pre>df.dropDuplicates(['name', 'height']).show()</pre>



	<p>across triggers as an intermediate state to drop duplicate rows.</p> <ul style="list-style-type: none"> <li>You can use <code>withWatermark()</code> to limit how late the duplicate data can be and the system will accordingly limit the state.</li> <li><code>drop_duplicates()</code> is an alias for <code>dropDuplicates()</code></li> <li>Package: <ul style="list-style-type: none"> <li><code>pyspark.sql.DataFrame</code></li> </ul> </li> </ul>	
dropna	<ul style="list-style-type: none"> <li>Returns a new <code>DataFrame</code> omitting rows with null values.</li> <li><code>DataFrame.dropna()</code> and <code>DataFrameNaFunctions.drop()</code> are aliases of each other.</li> <li>Package: <ul style="list-style-type: none"> <li><code>pyspark.sql.DataFrame</code></li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.dropna(how='any', thresh=None, subset=None)</pre> <ul style="list-style-type: none"> <li><code>how</code> <ul style="list-style-type: none"> <li>type: <code>str</code>, optional</li> <li>def: <code>'any'</code> or <code>'all'</code>. If <code>'any'</code>, drop a row if it contains any nulls. If <code>'all'</code>, drop a row only if all its values are null.</li> </ul> </li> <li><code>thresh</code>: <ul style="list-style-type: none"> <li>type: <code>int</code>, optional</li> <li>def: default <code>None</code> If specified, drop rows that have less than <code>thresh</code> non-null values. This overwrites the <code>how</code> parameter.</li> </ul> </li> <li><code>subset</code> <ul style="list-style-type: none"> <li>type: <code>str</code>, tuple or list, optional</li> </ul> </li> </ul>

		<ul style="list-style-type: none"> <li>◦ def: optional list of column names to consider.</li> </ul> <p><b>Example:</b></p> <pre>df4.na.drop().show()</pre>
where/filter	<ul style="list-style-type: none"> <li>• Filters rows using the given condition.</li> <li>• where() is an alias for filter().</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.DataFrame</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.filter(condition)</pre> <ul style="list-style-type: none"> <li>• Condition: <ul style="list-style-type: none"> <li>◦ type:Column or str</li> <li>◦ def: a Column of types.BooleanType or a string of SQL expression.</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.filter(df.age &gt; 3).collect() df.where(df.age == 2).collect() df.filter("age &gt; 3").collect() df.where("age = 2").collect()</pre>
distinct	<ul style="list-style-type: none"> <li>• Returns a new DataFrame containing the distinct rows in this DataFrame.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.DataFrame</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.distinct()</pre> <p><b>Example:</b></p> <pre>df.distinct().count()</pre>
sample	<ul style="list-style-type: none"> <li>• Returns a sampled subset of this DataFrame.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.DataFrame</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.sample(withReplacement=N one, fraction=None, seed=None)</pre> <ul style="list-style-type: none"> <li>• withReplacement <ul style="list-style-type: none"> <li>◦ type:bool, optional</li> <li>◦ def: Sample with</li> </ul> </li> </ul>

		<p>replacement or not (default False).</p> <ul style="list-style-type: none"> <li>fraction <ul style="list-style-type: none"> <li>type:float, optional</li> <li>def: Fraction of rows to generate, range [0.0, 1.0].</li> </ul> </li> <li>seed <ul style="list-style-type: none"> <li>type:int, optional</li> <li>def: Seed for sampling (default a random seed).</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.sample(withReplacement=True, fraction=0.5, seed=3).count()  df.sample(False, fraction=1.0)</pre>
split	<ul style="list-style-type: none"> <li>Splits str around matches of the given pattern.</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.functions</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>split(str, pattern, limit=- 1)</pre> <ul style="list-style-type: none"> <li>str <ul style="list-style-type: none"> <li>type: Column or str</li> <li>def:a string expression to split</li> </ul> </li> <li>pattern <ul style="list-style-type: none"> <li>Type:str</li> <li>def: a string representing a regular expression. The regex string should be a Java regular expression.</li> </ul> </li> <li>limit <ul style="list-style-type: none"> <li>type:int, optional</li> <li>def: an integer which controls the number of times a pattern is applied.</li> <li>limit &gt; 0:The resulting array's length will not be more than limit, and the resulting array's</li> </ul> </li> </ul>

		<p>last entry will contain all input beyond the last matched pattern.</p> <ul style="list-style-type: none"> <li>◦ <code>limit &lt;= 0</code>: pattern will be applied as many times as possible, and the resulting</li> <li>◦ array can be of any size.</li> </ul> <p><b>Example:</b></p> <pre>df.select(split(df.s, '[ABC]', -1).alias('s')).collect()</pre>
sort/orderBy	<ul style="list-style-type: none"> <li>• Returns a new DataFrame sorted by the specified column(s).</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ <code>pyspark.sql.DataFrame</code></li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.sort(*cols, **kwargs)[source]</pre> <ul style="list-style-type: none"> <li>• <code>cols</code> <ul style="list-style-type: none"> <li>◦ type: str, list, or Column, optional</li> <li>◦ def: list of Column or column names to sort by.</li> </ul> </li> <li>• <code>ascending</code> <ul style="list-style-type: none"> <li>◦ type: bool or list, optional</li> <li>◦ def: boolean or list of boolean (default True). Sort ascending vs. descending. Specify list for multiple sort orders. If a list is specified, the length of the list must equal the length of the cols.</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.sort(df.age.desc()).collect()  df.sort("age", ascending=False).collect()  df.sort(asc("age")).collect()</pre>

		<pre>df.orderBy(desc("age"), "name").collect()</pre> <pre>df.orderBy(["age", "name"], ascending=[0, 1]).collect()</pre>
sortwithinpartitions	<ul style="list-style-type: none"> <li>• Returns a new DataFrame with each partition sorted by the specified column(s).</li> <li>• package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.DataFrame</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.sortWithinPartitions(*cols, **kwargs)[source]</pre> <ul style="list-style-type: none"> <li>• cols <ul style="list-style-type: none"> <li>◦ type: str, list, or Column, optional</li> <li>◦ def: list of Column or column names to sort by.</li> </ul> </li> <li>• ascending <ul style="list-style-type: none"> <li>◦ type: bool or list, optional</li> <li>◦ def: boolean or list of boolean (default True). Sort ascending vs. descending. Specify list for multiple sort orders. If a list is specified, the length of the list must equal the length of the cols.</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.sortWithinPartitions("age", ascending=False).show()</pre>
asc_nulls_first	<ul style="list-style-type: none"> <li>• Returns a sort expression based on the ascending order of the column, and null values return before non-null values.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.Column</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>Column.asc_nulls_first()</pre> <p><b>Example:</b></p> <pre>df.select(df.name).orderBy(df.name .asc_nulls_first()).collect()</pre>

desc_nulls_first	<ul style="list-style-type: none"> <li>• Returns a sort expression based on the descending order of the column, and null values appear before non-null values.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.Column</li> </ul> </li> </ul>	<p><b>Syntax:</b> Column.desc_nulls_first()</p> <p><b>Example:</b></p> <pre>df.select(df.name).orderBy(df.name.desc_nulls_first()).collect()</pre>
asc_nulls_last	<ul style="list-style-type: none"> <li>• Returns a sort expression based on the ascending order of the column, and null values appear after non-null values.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.Column</li> </ul> </li> </ul>	<p><b>Syntax:</b> Column.asc_nulls_last()</p> <p><b>Example:</b></p> <pre>df.select(df.name).orderBy(df.name.asc_nulls_last()).collect()</pre>
desc_nulls_last	<ul style="list-style-type: none"> <li>• Returns a sort expression based on the descending order of the column, and null values appear after non-null values.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.sql.Column</li> </ul> </li> </ul>	<p><b>Syntax:</b> Column.desc_nulls_last()</p> <p><b>Example:</b></p> <pre>df.select(df.name).orderBy(df.name.desc_nulls_last()).collect()</pre>
limit	<ul style="list-style-type: none"> <li>• Limits the result count to the number specified.</li> <li>• Package:</li> </ul>	<p><b>Syntax:</b> DataFrame.limit(num)</p> <p><b>Example:</b></p>

	<ul style="list-style-type: none"> <li>○ pyspark.sql.DataFrame</li> </ul>	df.limit(1).collect()
first		
take		
collect		
show	<ul style="list-style-type: none"> <li>● Prints the first n rows to the console.</li> <li>● Package: <ul style="list-style-type: none"> <li>○ pyspark.sql.DataFrame</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.show(n=20, truncate=True, vertical=False)[source]</pre> <ul style="list-style-type: none"> <li>● N <ul style="list-style-type: none"> <li>○ type: int, optional</li> <li>○ def: Number of rows to show.</li> </ul> </li> <li>● truncate <ul style="list-style-type: none"> <li>○ type: bool, optional</li> <li>○ def: If set to True, truncate strings longer than 20 chars by default. If set to a number greater than one, truncates long strings to length, and aligns cells right.</li> </ul> </li> <li>● vertical <ul style="list-style-type: none"> <li>○ type: bool, optional</li> <li>○ def: If set to True, print output rows vertically (one line per column value).</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.show(truncate=3) df.show(vertical=True)</pre>

from_unixtime	<ul style="list-style-type: none"> <li>Converts the number of seconds from the unix epoch (1970-01-01 00:00:00 UTC) to a string representing the timestamp of that moment in the current system time zone in the given format.</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.functions</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>from_unixtime(timestamp, format='yyyy-MM-dd HH:mm:ss')[source]</pre> <p><b>Example:</b></p> <pre>time_df = spark.createDataFrame([(1428476400,)], ['unix_time']) time_df.select(from_unixtime('unix_time').alias('ts')).collect()</pre>
to_date		
datediff	<ul style="list-style-type: none"> <li>Returns number of data from start to end</li> <li>Package: <ul style="list-style-type: none"> <li>pyspark.sql.functions</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>datediff(end, start)</pre> <p><b>Example:</b></p> <pre>df = spark.createDataFrame([('2015-04-08', '2015-05-10')], ['d1', 'd2']) df.select(datediff(df.d2, df.d1).alias('diff')).collect()</pre>
explode	<ul style="list-style-type: none"> <li>Returns a new row for each element in the given array or map.</li> <li>Uses the default column name for elements in the array and key and value for elements in the map unless specified otherwise.</li> <li>package:</li> </ul>	<p><b>Syntax:</b></p> <pre>explode(col)</pre> <p><b>Example:</b></p> <pre>eDF = spark.createDataFrame([Row(a=1, intlist=[1,2,3], mapfield={"a": "b"})]) eDF.select(explode(eDF.intlist).alias("anInt")).collect()</pre>



	<ul style="list-style-type: none"> <li>○ pyspark.sql.functions</li> </ul>	
agg	<ul style="list-style-type: none"> <li>● Aggregate on the entire DataFrame without groups</li> <li>● package:pyspark.sql.DataFrame</li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.agg(*exprs)</pre> <p><b>Example:</b></p> <pre>df.groupBy().agg() df.agg({"age": "max"}).collect() df.agg(F.min(df.age)).collect()</pre>
avg	<ul style="list-style-type: none"> <li>● Aggregate function: returns the average of the values in a group.</li> <li>● Package: <ul style="list-style-type: none"> <li>○ pyspark.sql.functions</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>avg(col)</pre> <p><b>Example:</b></p>
isin	<ul style="list-style-type: none"> <li>● A boolean expression that is evaluated to be true if the value of this expression is contained by the evaluated values of the arguments.</li> <li>● Package: <ul style="list-style-type: none"> <li>○ pyspark.sql.Column</li> </ul> </li> </ul>	<p><b>Syntax:</b></p> <pre>Column.isin(*cols)</pre> <p><b>Example:</b></p> <pre>df[df.name.isin("Bob", "Mike")].collect() df[df.age.isin([1, 2, 3])].collect()</pre>
range	<ul style="list-style-type: none"> <li>● Create a new RDD of int containing</li> </ul>	<p><b>Syntax:</b></p> <pre>range(start, end=None, step=1,</pre>

	<p>elements from start to end (exclusive), increased by every step.</p> <ul style="list-style-type: none"> <li>• Can be called the same way as python's built-in range() function.</li> <li>• If called with a single argument, the argument is interpreted as end, and start is set to 0.</li> <li>• Package: <ul style="list-style-type: none"> <li>◦ pyspark.SparkContext</li> </ul> </li> </ul>	<pre>numSlices=None)</pre> <ul style="list-style-type: none"> <li>• start <ul style="list-style-type: none"> <li>◦ Type:int</li> <li>◦ def:the start value</li> </ul> </li> <li>• end <ul style="list-style-type: none"> <li>◦ type:int, optional</li> <li>◦ def:the end value (exclusive)</li> </ul> </li> <li>• step <ul style="list-style-type: none"> <li>◦ type:int, optional</li> <li>◦ def:the incremental step (default: 1)</li> </ul> </li> <li>• numSlices <ul style="list-style-type: none"> <li>◦ type:int, optional</li> <li>◦ def:the number of partitions of the new RDD</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>sc.range(5).collect() sc.range(2, 4).collect() sc.range(1, 7, 2).collect()</pre>
--	---	---

<b>Actions</b>	<b>Wide transformation</b>	<b>Narrow transformations</b>
Count	Repartition	Coalesce
foreach	groupby	Filter
head	Join	map
Collect	distinct	sample
take		union
top		flatMap
reduce		drop
agg		selectexpr
printschema		withcolumnrenamed
show		
take		

## AGGREGATE FUNCTIONS:

Function	Use	Example
agg	<ul style="list-style-type: none"> <li>• Compute aggregates and returns the result as a DataFrame.</li> <li>• The available aggregate functions can be: <ul style="list-style-type: none"> <li>◦ built-in aggregation functions, such as avg, max, min, sum, count</li> <li>◦ group aggregate pandas UDFs, created with <code>pyspark.sql.functions.pandas_udf()</code></li> </ul> </li> <li>• package: <code>pyspark.sql.GroupedData</code></li> </ul>	<p><b>Syntax:</b></p> <pre>GroupedData.agg(*exprs</pre> <ul style="list-style-type: none"> <li>• <code>exprs:</code> <ul style="list-style-type: none"> <li>◦ type: dict</li> <li>◦ def: a dict mapping from column name (string) to aggregate functions (string), or a list of columns.</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>gdf = df.groupBy(df.name) sorted(gdf.agg({"*": "count"}).collect())  sorted(gdf.agg(F.min(df.age)).collect())</pre>
avg/mean	<ul style="list-style-type: none"> <li>• Computes average values for each numeric column for each group.</li> <li>• <code>mean()</code> is an alias for <code>avg()</code>.</li> <li>• package: <code>pyspark.sql.GroupedData</code></li> </ul>	<p><b>Syntax:</b></p> <pre>GroupedData.avg(*cols)</pre> <ul style="list-style-type: none"> <li>• <code>cols</code> <ul style="list-style-type: none"> <li>◦ type: str</li> <li>◦ def: column names. Non-numeric columns are ignored.</li> </ul> </li> </ul> <p><b>Example:</b></p>

		<pre>df.groupBy().avg('age').collect()</pre> <pre>df3.groupBy().avg('age', 'height').collect()</pre>
count	<ul style="list-style-type: none"> <li>Counts the number of records for each group</li> <li>package: pyspark.sql. GroupedData.</li> </ul>	<p><b>Syntax:</b></p> <pre>GroupedData.count()</pre> <p><b>Example:</b></p> <pre>sorted(df.groupBy(df.age).count().collect())</pre>
max	<ul style="list-style-type: none"> <li>Computes the max value for each numeric column for each group.</li> <li>package: pyspark.sql. GroupedData</li> </ul>	<p><b>Syntax:</b></p> <pre>GroupedData.max(*cols)[source]</pre> <p><b>Example:</b></p> <pre>df.groupBy().max('age').collect()</pre> <pre>df3.groupBy().max('age', 'height').collect()</pre>
min	<ul style="list-style-type: none"> <li>Computes the min value for each numeric column for each group.</li> <li>package: pyspark.sql. GroupedData</li> </ul>	<p><b>Syntax</b></p> <pre>GroupedData.min(*cols)</pre> <ul style="list-style-type: none"> <li>cols <ul style="list-style-type: none"> <li>type:str</li> <li>def: column names. Non-numeric columns are ignored.</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.groupBy().min('age').collect()</pre>
pivot	<ul style="list-style-type: none"> <li>Pivots a column of the current DataFrame and performs the</li> </ul>	<p><b>Syntax:</b></p> <pre>GroupedData.pivot(pivot_col, values=None)</pre>

	<p>specified aggregation.</p> <ul style="list-style-type: none"> <li>There are two versions of the pivot function: <ul style="list-style-type: none"> <li>one that requires the caller to specify the list of distinct values to pivot on</li> <li>one that does not.</li> </ul> </li> <li>package: <code>pyspark.sql.GroupedData</code></li> </ul>	<ul style="list-style-type: none"> <li><code>pivot_col</code> <ul style="list-style-type: none"> <li>type: <code>str</code></li> <li>def: Name of the column to pivot.</li> </ul> </li> <li><code>values :</code> <ul style="list-style-type: none"> <li>def: List of values that will be translated to columns in the output DataFrame.</li> </ul> </li> </ul> <p><b>Examples:</b></p> <ul style="list-style-type: none"> <li>Compute the sum of earnings for each year by course with each course as a separate column</li> </ul> <pre>df4.groupBy("year").pivot("course", ["dotNET", "Java"]).sum("earnings").collect()</pre> <ul style="list-style-type: none"> <li>Or without specifying column values (less efficient)</li> </ul> <pre>df4.groupBy("year").pivot("course").sum("earnings").collect()</pre>
sum	<ul style="list-style-type: none"> <li>Compute the sum for each numeric column for each group.</li> <li>package: <code>pyspark.sql.GroupedData</code></li> </ul>	<p><b>Syntax:</b></p> <pre>GroupedData.sum(*cols)</pre> <ul style="list-style-type: none"> <li><code>cols</code> <ul style="list-style-type: none"> <li>type :<code>str</code></li> <li>def: column names. Non-numeric columns are ignored.</li> </ul> </li> </ul> <p><b>Example:</b></p> <pre>df.groupBy().sum('age').collect()</pre> <pre>df3.groupBy().sum('age', 'height').collect()</pre>

--	--	--

## JOINS:

Function	Usage	Example
join	<ul style="list-style-type: none"> <li>Joins with another DataFrame, using the given join expression.</li> <li>package: pyspark.sql.DataFrame.join</li> </ul>	<p><b>Syntax:</b></p> <pre>DataFrame.join(other, on=None, how=None)</pre> <ul style="list-style-type: none"> <li>other <ul style="list-style-type: none"> <li>type: DataFrame</li> <li>def: Right side of the join</li> </ul> </li> <li>On <ul style="list-style-type: none"> <li>type: str, list or Column, optional</li> <li>def: a string for the join column name, a list of column names, a join expression (Column), or a list of Columns. If there is a string or a list of strings indicating the name of the join column(s), the column(s) must exist on both sides, and this performs an equi-join.</li> </ul> </li> </ul>

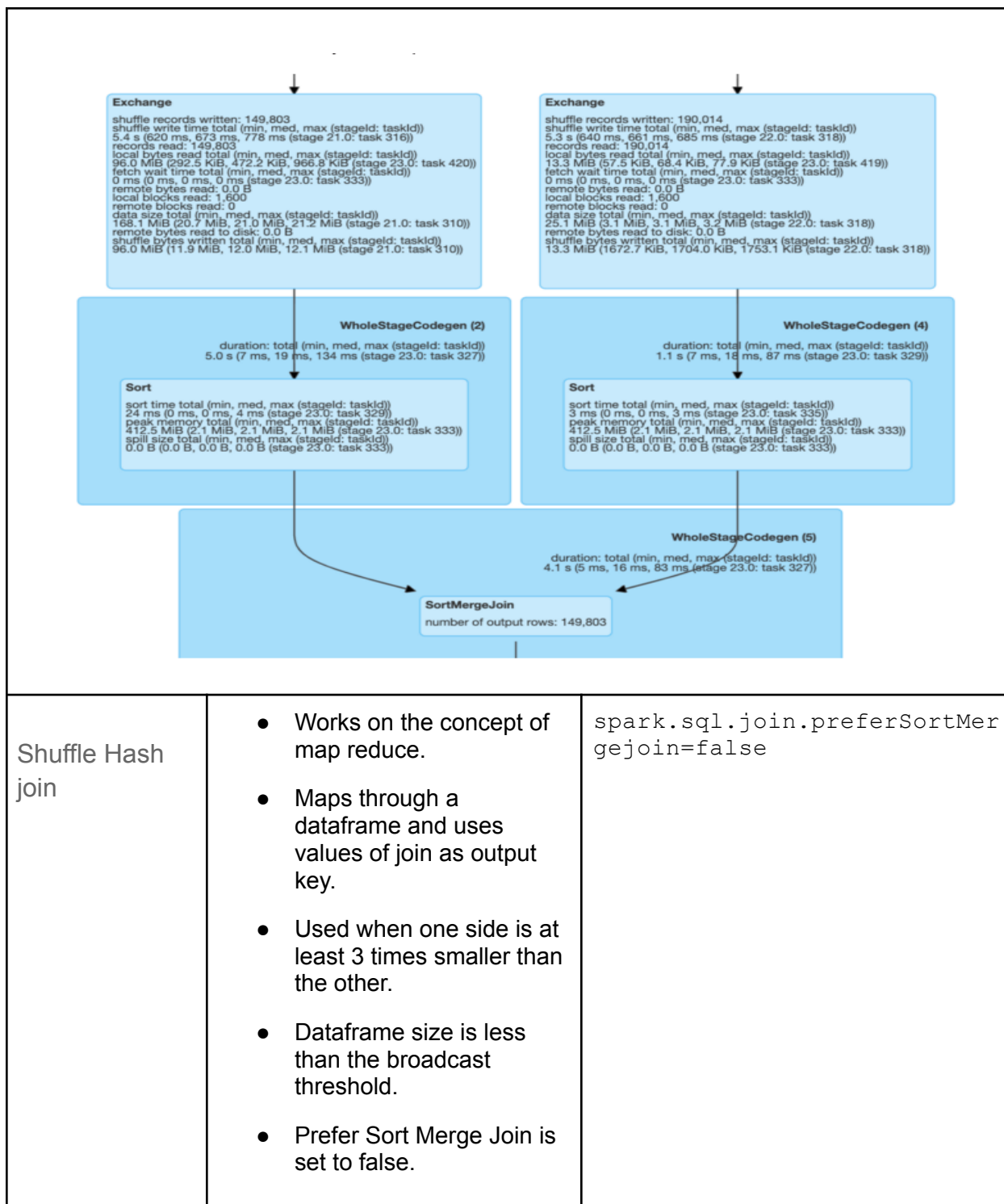
		<ul style="list-style-type: none"> <li>• how <ul style="list-style-type: none"> <li>◦ type: str, optional</li> <li>◦ def: <ul style="list-style-type: none"> <li>■ default inner.</li> <li>■ Must be one of: inner, cross, outer, full, fullouter, full_outer, left, leftouter, left_outer, right, rightouter, right_outer, semi, leftsemi, left_semi, anti, leftanti and left_anti.</li> </ul> </li> </ul> </li> </ul>
Inner Join	<ul style="list-style-type: none"> <li>• Spark Inner join is the default join and it's mostly used.</li> <li>• It is used to join two DataFrames/Datasets on key columns, and where keys don't match the rows get dropped from both datasets</li> </ul>	<pre>df.join(df2, df.name == df2.name, 'outer') .select(df.name, df2.height)</pre>
Full Outer Join	<ul style="list-style-type: none"> <li>• fullouter join returns all rows from both Spark DataFrame/Datasets,</li> </ul>	



	<p>where join expression doesn't match it returns null on respective record columns.</p>	
Right Outer join	<ul style="list-style-type: none"> <li>• is the opposite of left join.</li> <li>• here it returns all rows from the right DataFrame/Dataset regardless of match found on the left dataset, when join expression doesn't match, it assigns null for that record and drops records from left where match not found.</li> </ul>	
Left Semi Join	<ul style="list-style-type: none"> <li>• The Spark Left Semi join is similar to the inner join.</li> <li>• The difference being leftsemi join returns all columns from the left DataFrame/Dataset and ignores all columns from the right dataset.</li> <li>• In other words, this join returns columns from the only left dataset for the records matched in the right dataset on join expression; records not matched on join expression are ignored from both left and right datasets.</li> <li>• The same result can be achieved using select on the result of the inner join; however, using this join would be efficient.</li> </ul>	

Left Anti Join	<ul style="list-style-type: none"> <li>Left Anti join does the exact opposite of the Spark leftsemi join, leftanti join returns only columns from the left DataFrame/Dataset for non-matched records.</li> </ul>	
Self Join	<ul style="list-style-type: none"> <li>Though there is no self-join type available, we can use any of the above-explained join types to join DataFrame to itself. below example use inner self join</li> </ul>	
Join Algorithms in Spark	<ul style="list-style-type: none"> <li>The algorithm chosen depends on the size of the tables joined and the hints provided in join</li> <li>The commonly used joins in spark are: <ul style="list-style-type: none"> <li>Broadcast Hash join</li> <li>Sort Merge Join</li> <li>Shuffle Hash Join</li> </ul> </li> </ul>	
Broadcast Hash join	<ul style="list-style-type: none"> <li>Used when dataframe size is less than <code>spark.sql.autoBroadcastJoinThreshold</code>.</li> <li>Does not perform shuffle while doing a join.</li> <li>Broadcasting huge datasets results in a</li> </ul>	

	timeout	
	<pre> graph TD     In1[ ] --&gt; Filter[Filter number of output rows: 149,803]     Filter --&gt; P1[Project]     P1 --&gt; BHJ[BroadcastHashJoin number of output rows: 149,803]     In2[ ] --&gt; P2[Project]     P2 --&gt; BE[BroadcastExchange data size: 33.9 MiB time to collect: 266 ms time to build: 73 ms time to broadcast: 79 ms]     BE --&gt; BHJ   </pre>	
Sort Merge Join	<ul style="list-style-type: none"> <li>• If dataframes cannot be broadcasted, spark used sort merge join.</li> <li>• Uses node-to node communication strategy.</li> <li>• Two steps: <ul style="list-style-type: none"> <li>◦ First step exchanges and sorts datasets.</li> <li>◦ Second step merges the dataset.</li> </ul> </li> <li>• Better to use bucketing for optimization.</li> <li>• Default since spark 2.3</li> </ul>	<code>spark.sql.join.preferSortMergeJoin=true</code>



## PERFORMANCE TUNING CONFIG OPTIONS

Property	Definition	Default
Performance Tuning considerations	<ul style="list-style-type: none"> <li>Overhead of garbage collection</li> <li>Amount of memory used by objects</li> <li>Cost of accessing these objects</li> </ul>	
Improve spark performance	<ul style="list-style-type: none"> <li>Increase spark.default.parallelism and spark.sql.shuffle.partitions to improve spark performance</li> </ul>	
spark.sql.inMemoryColumnarStorage.compressed	<ul style="list-style-type: none"> <li>When set to true, Spark SQL will automatically select a compression codec for each column based on statistics of the data.</li> </ul>	true
spark.sql.inMemoryColumnarStorage.batchSize	<ul style="list-style-type: none"> <li>Controls the size of batches for columnar caching.</li> <li>Larger batch sizes can improve memory utilization and compression, but risk OOMs when caching data.</li> </ul>	10000
spark.sql.files.maxPartitionBytes	<ul style="list-style-type: none"> <li>The maximum number of bytes to pack into a single partition when reading files.</li> <li>This configuration is effective only when using file-based sources such as Parquet, JSON and ORC.</li> </ul>	134217728 (128 MB)
spark.sql.files.openCostInBytes	<ul style="list-style-type: none"> <li>The estimated cost to open a file, measured by the number of bytes, could be scanned at the</li> </ul>	4194304 (4 MB)

	<p>same time.</p> <ul style="list-style-type: none"> <li>• This is used when putting multiple files into a partition.</li> <li>• It is better to over-estimated, then the partitions with small files will be faster than partitions with bigger files (which is scheduled first).</li> <li>• This configuration is effective only when using file-based sources such as Parquet, JSON and ORC.</li> </ul>	
spark.sql.files.minPartitionNum	<ul style="list-style-type: none"> <li>• The suggested (not guaranteed) minimum number of split file partitions.</li> <li>• If not set, the default value is `spark.default.parallelism`.</li> <li>• This configuration is effective only when using file-based sources such as Parquet, JSON and ORC.</li> </ul>	Default Parallelism
spark.sql.broadcastTimeout	<ul style="list-style-type: none"> <li>• Timeout in seconds for the broadcast wait time in broadcast joins</li> </ul>	300
spark.sql.autoBroadcastJoinThreshold	<ul style="list-style-type: none"> <li>• Configures the maximum size in bytes for a table that will be broadcast to all worker nodes when performing a join.</li> <li>• By setting this value to -1 broadcasting can be disabled.</li> <li>• Note that currently statistics are only supported for Hive Metastore tables where the command <code>ANALYZE TABLE &lt;tableName&gt; COMPUTE STATISTICS noscan</code> has been run.</li> </ul>	10485760 (10 MB)
spark.sql.shuffle.partitions	<ul style="list-style-type: none"> <li>• Configures the number of partitions to use when shuffling</li> </ul>	200

	data for joins or aggregations.	
spark.sql.sources.parallelPartitionDiscovery.threshold	<ul style="list-style-type: none"> <li>Configures the threshold to enable parallel listing for job input paths.</li> <li>If the number of input paths is larger than this threshold, Spark will list the files by using a Spark distributed job.</li> <li>Otherwise, it will fallback to sequential listing.</li> <li>This configuration is only effective when using file-based data sources such as Parquet, ORC and JSON.</li> </ul>	32
spark.sql.sources.parallelPartitionDiscovery.parallelism	<ul style="list-style-type: none"> <li>Configures the maximum listing parallelism for job input paths.</li> <li>In case the number of input paths is larger than this value, it will be throttled down to use this value.</li> <li>Effective when using file-based data sources such as Parquet, ORC and JSON.</li> </ul>	10000





# TUNING IN SPARK:

Topic	Definition
Need of Tuning	<ul style="list-style-type: none"><li>• Since spark is of in-memory nature we can have bottlenecks due to resource, memory or network bandwidth.</li></ul>
Areas of tuning	<ul style="list-style-type: none"><li>• Data serialization-crucial for good network performance and reduce memory usage.</li><li>• Memory tuning</li></ul>
Data serialization	Covered earlier
Memory tuning	<ul style="list-style-type: none"><li>• Java objects occupy a lot of memory so we need to mind memory usage and also use appropriate objects</li></ul>
Memory management in spark.	<ul style="list-style-type: none"><li>• Spark has two categories of memory usage:<ul style="list-style-type: none"><li>◦ Storage</li><li>◦ Execution</li></ul></li><li>• Execution refers to memory used for computation in shuffle joins etc.</li><li>• Storage is memory used for caching and propagating internal data structures.</li><li>• When no execution is being done, storage can occupy all the memory.</li><li>• There is a subregion in memory which is always reserved for storage.</li><li>• <code>spark.memory.fraction</code> refers to the size of memory M</li><li>• <code>spark.memory.storageFraction</code> refers to size of storage fraction as</li></ul>

	a fraction of M
Determining memory usage	<ul style="list-style-type: none"> <li>• Put rdd in the cache and look at the storage tab in spark UI.</li> <li>• To estimate the memory of a particular object, use sizeEstimator's estimate method.</li> </ul>
Measuring impact of GC	<ul style="list-style-type: none"> <li>• This can be done by adding  <code>-verbose:gc</code>  <code>-XX:+PrintGCDetails</code>  <code>-XX:+PrintGCTimeStamps</code> to the Java options</li> </ul>
Components of memory in JVM	<ul style="list-style-type: none"> <li>• Java Heap space is divided into two regions: Young and Old. The Young generation is meant to hold short-lived objects while the Old generation is intended for objects with longer lifetimes.</li> <li>• The Young generation is further divided into three regions [Eden, Survivor1, Survivor2].</li> <li>• A simplified description of the garbage collection procedure: <ul style="list-style-type: none"> <li>○ When Eden is full, a minor GC is run on Eden and objects that are alive from Eden and Survivor1 are copied to Survivor2.</li> <li>○ The Survivor regions are swapped. If an object is old enough or Survivor2 is full, it is moved to Old. Finally, when Old is close to full, a full GC is invoked.</li> </ul> </li> </ul>
Goal Of GC	<ul style="list-style-type: none"> <li>• The goal of GC tuning in Spark is to ensure that only long-lived RDDs are stored in the Old generation and that the Young generation is sufficiently</li> </ul>

	sized to store short-lived objects.
Steps to avoid Full GC	<ul style="list-style-type: none"> <li>• Check if there are too many garbage collections by collecting GC stats. If a full GC is invoked multiple times before a task completes, it means that there isn't enough memory available for executing tasks.</li> <li>• If there are too many minor collections but not many major GCs, allocating more memory for Eden would help.</li> <li>• In the GC stats that are printed, if the OldGen is close to being full, reduce the amount of memory used for caching by lowering <code>spark.memory.fraction</code>;</li> </ul>
Data locality	<ul style="list-style-type: none"> <li>• Data locality is how close data is to the processing it.</li> <li>• This affects speed of execution.</li> <li>• The levels of data locality are: <ul style="list-style-type: none"> <li>◦ <code>PROCESS_LOCAL</code> data is in the same JVM as the running code. This is the best locality possible</li> <li>◦ <code>NODE_LOCAL</code> data is on the same node. Examples might be in HDFS on the same node, or in another executor on the same node. This is a little slower than <code>PROCESS_LOCAL</code> because the data has to travel between processes</li> <li>◦ <code>NO_PREF</code> data is accessed equally quickly from anywhere and has no locality preference</li> <li>◦ <code>RACK_LOCAL</code> data is on the same rack of servers. Data is on a different server on the</li> </ul> </li> </ul>

	<p>same rack so needs to be sent over the network, typically through a single switch</p> <ul style="list-style-type: none"> <li>○ ANY data is elsewhere on the network and not in the same rack</li> </ul>
--	--

## References And Important Links for Spark:

Advanced Apache Spark Training - Sameer Farooqui (Databricks)

- <https://www.youtube.com/watch?v=7ooZ4S7Ay6Y>

A Deep Dive into Spark SQL's Catalyst Optimizer with Yin Huai

- <https://www.youtube.com/watch?v=RmUn5vHlevc&t=671s>

Preparation material :

- Spark the definitive guide
- <http://spark.apache.org/docs/latest/index.html>

Useful Links:

- <https://towardsdatascience.com/spark-essentials-how-to-read-and-write-data-with-pyspark-5c45e29227cd?gi=9c79a01f751a>
- <https://sparkbyexamples.com/pyspark/pyspark-sql-types-datatype-with-examples/>

## Interesting White papers:

- Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing
  - [http://people.csail.mit.edu/matei/papers/2012/nsdi\\_spark.pdf](http://people.csail.mit.edu/matei/papers/2012/nsdi_spark.pdf)
- Spark SQL: Relational Data Processing in Spark
  - [http://people.csail.mit.edu/matei/papers/2015/sigmod\\_spark\\_sql.pdf](http://people.csail.mit.edu/matei/papers/2015/sigmod_spark_sql.pdf)
- Discretized Streams: Fault-Tolerant Streaming Computation at Scale
  - [http://people.csail.mit.edu/matei/papers/2013/sosp\\_spark\\_streaming.pdf](http://people.csail.mit.edu/matei/papers/2013/sosp_spark_streaming.pdf)
- Discretized Streams: An Efficient and Fault-Tolerant Model for Stream Processing on Large Clusters
  - [http://people.csail.mit.edu/matei/papers/2012/hotcloud\\_spark\\_streaming.p  
df](http://people.csail.mit.edu/matei/papers/2012/hotcloud_spark_streaming.pdf)
- Shark: Fast Data Analysis Using Coarse-grained Distributed Memory
  - [http://people.csail.mit.edu/matei/papers/2012/sigmod\\_shark\\_demo.pdf](http://people.csail.mit.edu/matei/papers/2012/sigmod_shark_demo.pdf)
- Spark: Cluster Computing with Working Sets
  - [http://people.csail.mit.edu/matei/papers/2010/hotcloud\\_spark.pdf](http://people.csail.mit.edu/matei/papers/2010/hotcloud_spark.pdf)