

CS x55: DISTRIBUTED SYSTEMS [SPARK]

Transformations: Narrow and Wide

Though their numbers are few

Don't let them beguile you

Innocuous though

they may seem

The wrong invocation

Is all it takes

To amplify inefficiencies

And protract computations

Shrideep Pallickara
Computer Science
Colorado State University



Frequently asked questions from the previous class survey

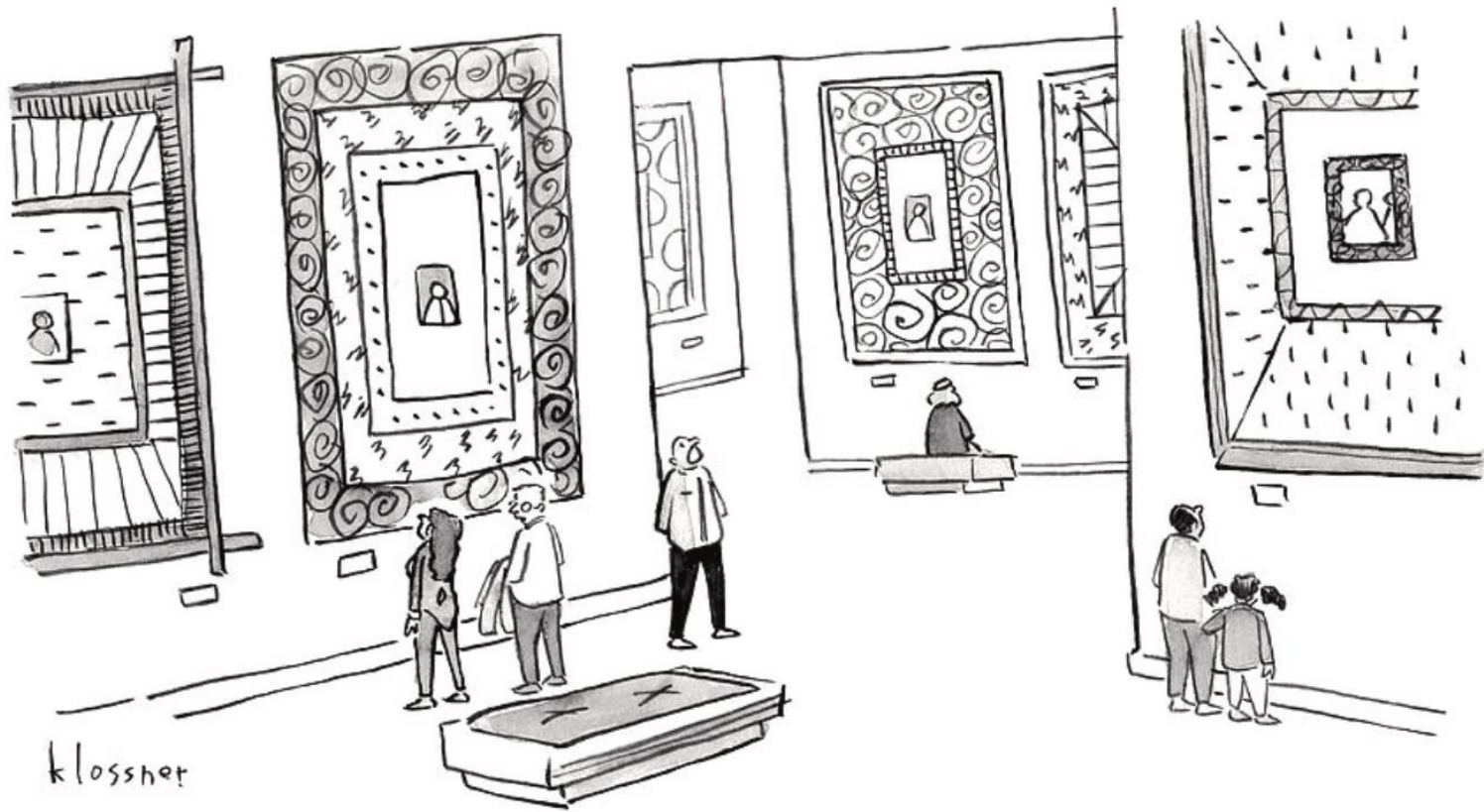
- What happens if you persist an RDD and then you don't use it?
- How do Spark and Hadoop make money?
- Data underpinning transformation operations ... where are they (on disk, in memory)?
- How do you choose the right persistence strategy?



Topics covered in this lecture

- Data Frames
 - ▣ Column manipulations
- Orchestration Plans





DATAFRAMES

"Painting was her love, but framing was her passion."

Paint was her love. John Klossner. New Yorker. April 2023.

Spark DataFrame

- DataFrames consist of
 - ▣ A series of **records** (like rows in a table) that are of type `Row`
 - ▣ A number of columns (like columns in a spreadsheet)
- Rows
 - ▣ You can create rows by manually instantiating a `Row` object with the values that belong in each column
- Columns
 - ▣ You can select, manipulate, and remove columns from `DataFrames` and these operations are represented as **expressions**



Schemas

- A **schema** defines the column names and types of a DataFrame
- You can let a data source define the schema (called schema-on-read) or define it explicitly
- Note that only DataFrames have schemas
 - ▣ Rows themselves **do not** have schemas
 - ▣ If you create a Row manually?
 - You must specify the values **in the same order** as the schema of the DataFrame to which they might be appended



We can create DataFrames from raw data sources

- Spark has six “**core**” data sources
 - ▣ CSV
 - ▣ JSON
 - ▣ Parquet
 - ▣ ORC: Apache Optimized Row Columnar (ORC) file format
 - ▣ JDBC/ODBC connections
 - ▣ Plain-text files
- Hundreds of external data sources written by the community
 - ▣ E.g.: Cassandra, HBase, MongoDB, AWS, Redshift, XML etc.



The foundation for reading data in Spark is the `DataFrameReader`

- We access this through the `SparkSession` via the `read` attribute: `spark.read`
- After we have a `DataFrame` reader, we specify several aspects:
 - ▣ The format: Input data source format
 - ▣ The schema
 - ▣ The read mode `{Permissive, DropMalformed, Failfast}`
 - ▣ A series of options
- **The format, options, and schema each return a `DataFrameReader`** that can undergo *further* transformations and are all optional



However, at a minimum, the DataFrameReader must have a **path** from which to read

```
spark.read.format("csv")  
  .option("mode", "FAILFAST")  
  .option("inferSchema", "true")  
  .option("path", "path/to/file(s)")  
  .schema(someSchema)  
  .load()
```



Writing data is quite similar to that of reading data

- Instead of the `DataFrameReader` , we have the `DataFrameWriter`
- We access the `DataFrameWriter` on a per-`DataFrame` basis via the `write` attribute:

```
dataFrame.write
```



Writing Data

- After we have a `DataFrameWriter`, we specify three values:
 - ▣ The format, a series of options, and the save mode (e.g., `append` or `overwrite`)
- **At a minimum**, you must supply a path
- Options may vary from data source to data source

```
dataframe.write.format( "csv" )  
                .option( "mode", "APPEND" )  
                .option( "dateFormat", "yyyy-MM-dd" )  
                .option( "path", "path/to/file(s)" )  
                .save ( )
```



You can make any DataFrame into a table or view

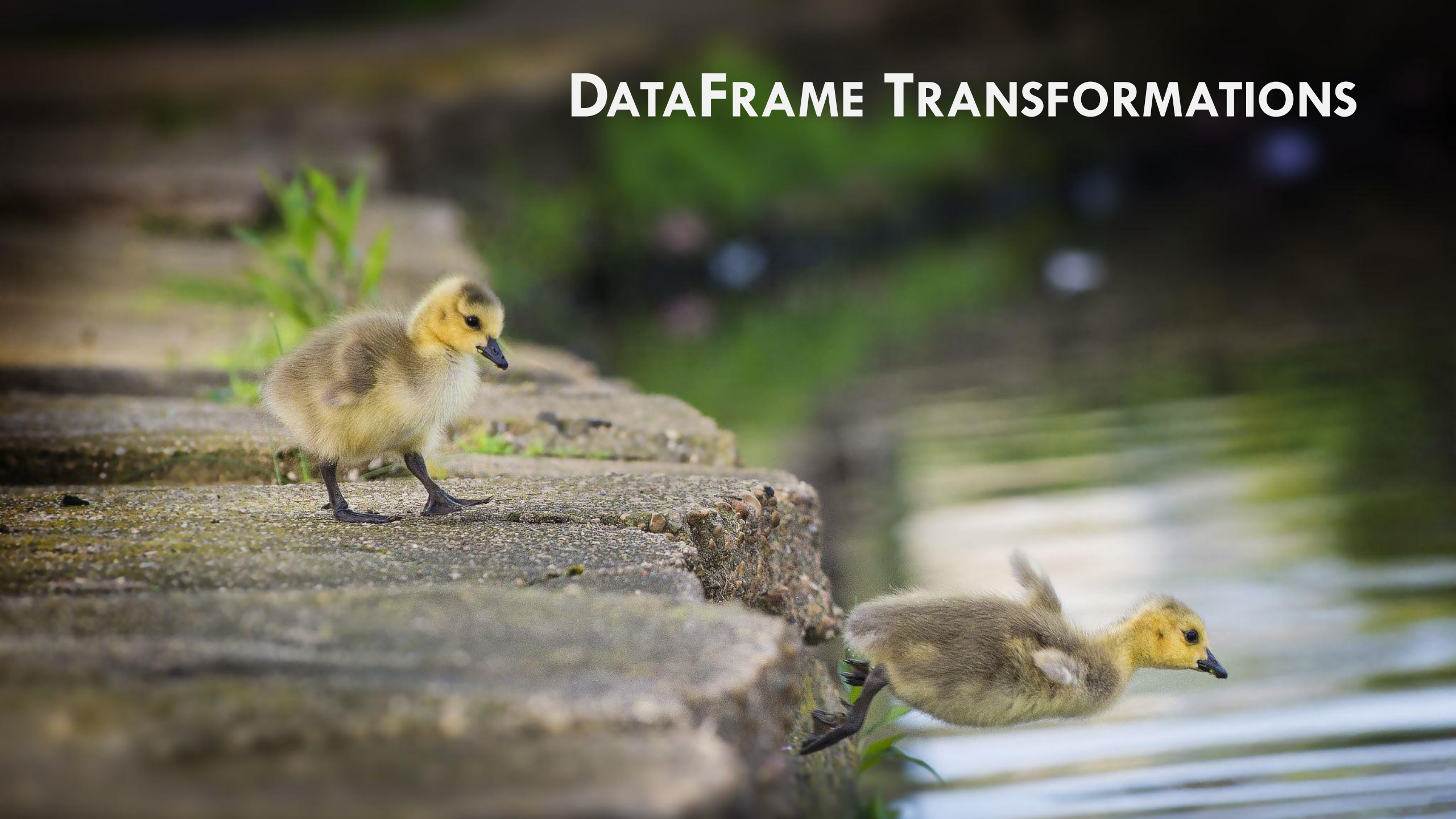
- Done via a simple method call: `createOrReplaceTempView`
- This then allows you to query the data using SQL

```
val df = spark.read
    .format("json" )
    .load("/data/flight-data/json/2022-summary.json")

df.createOrReplaceTempView("dfTable")
```



DATAFRAME TRANSFORMATIONS



DataFrame transformations

- Add rows or columns
- Remove rows or columns
- Transform a row into a column (or vice versa)
- Change the order of rows based on the values in columns



Adding Columns

- Use the **withColumn** method on the DataFrame
- For example, let's add a column that just adds the number "1" as a column:

```
df.withColumn("numberOne", lit(1))
```

`lit` is short for "literal"
`lit()` wraps a constant so Spark can treat it like a column in expressions



Renaming Columns

- Done using the `withColumnRenamed` method
- Will rename the column with the name of the string in the first argument to the string in the second argument:

```
df.withColumnRenamed ("DEST_COUNTRY_NAME", "dest")
```



Removing Columns

- Done using a method called **drop**

```
df.drop("ORIGIN_COUNTRY_NAME" )
```

- We can drop multiple columns by passing in multiple columns as arguments

```
dfWithLongColName.drop("ORIGIN_COUNTRY_NAME",  
                        "DEST_COUNTRY_NAME")
```



Filtering Rows

- To **filter** rows, we create an expression that evaluates to true or false
 - ▣ Those rows where the expression evaluates to `false` are filtered out

```
df.filter( col( "count" ) < 2)
```



Getting Unique Rows

- A very common use case is to extract the unique or **distinct** values in a DataFrame
 - ▣ These values can be in **one or more columns**
 - ▣ Done by using the **distinct** method on a DataFrame
 - Allows **deduplication** of any rows that are in that DataFrame.
 - ▣ Again, this is a transformation that will return a **new** DataFrame with only unique rows:

```
df.select("ORIGIN_COUNTRY_NAME", "DEST_COUNTRY_NAME")  
  .distinct()
```



Random Samples

- You might want to sample some random records from a `DataFrame`
- Done by using the **sample** method on a `DataFrame`
 - ▣ Specify a fraction of rows to extract from a `DataFrame` and whether the sample will be with or without replacement

```
val seed = 5
val withReplacement = false
val fraction = 0.5
df.sample(withReplacement, fraction, seed )
```



Random Splits

- ❑ Random splits are helpful when you need to break up a `DataFrame` into a random “splits” of the original `DataFrame`
- ❑ Often used with machine learning algorithms to create training, validation, and test sets

```
val dataFrames =
```

```
    df.randomSplit(Array (0.60, 0.20, 0.20 ), seed )
```

- ❑ E.g.: 60% for training, 20% for validation (tuning), 20% for test (reality check)
- ❑ `seed`: ensures the split is reproducible



□ **withColumn(columnName, func)**

- ▣ Return a Dataframe with the additional column

- ▣ Invocation: `df.withColumn("dogYears", df.age / 7)`

□ **dropColumn(columnName)**

- ▣ Return a Dataframe without the column

- ▣ Invocation: `df.dropColumn("age")`



□ **select(columnNames)**

- ▣ Return a DataFrame with the specified columns
- ▣ Invocation: `df.select("firstName", "age")`

□ **describe(columnName)**

- ▣ Compute summary statistics over DataFrame columns
- ▣ Invocation: `df.describe("age"), df.describe()`



Column Manipulations

[3/4]

```
val df = Seq(  
    ("Peterson", "Marcus", 54),  
    ("Batey", "Edward", 36),  
    ("Bruce", "Karen", 35)  
).toDF("lastName", "firstName", "age")  
  
df.withColumn("dogYears", df.age / 7.0)  
df.describe("age", "dogYears")
```



Column Manipulations

[4/4]

```
+-----+-----+-----+
|summary|      age| dogYears|
+-----+-----+-----+
|  count|        3|        3|
|   mean|  41.6667|  5.95238|
| stddev| 10.69268|  1.52753|
|    min|       35|        5|
|    max|       54|  7.714286|
+-----+-----+-----+
```



Dataframe joins

- `join(other, <columnComparison>, <joinType>)`
 - ▣ Performs a join between 2 Dataframes
 - ▣ Invocation: `df1.join(df2, Seq("id"))`
 - ▣ Spark matches rows where the value of `id` in `df1` equals the value of `id` in `df2`
 - ▣ The result is a new DataFrame containing columns from both, aligned on the matching keys



Join column comparison

- Supports a variety of criteria
 - ▣ Sequence of column names (e.g., `Seq("id", "age")`)
 - ▣ Elaborate comparison definitions (e.g., `df1("age") >= df2("age")`)



Join Type

- DataFrames may perform multiple styles of join
 - ▣ Inner: typical dataset join with key-to-key match
 - ▣ Outer, left-outer, right-outer: result contains all rows, filling in columns with 'null' values where data doesn't exist
 - ▣ Left-semi, right-semi: similar to outer join, but result only contains rows in specified source dataset



Example: Spark SQL

```
val df = Seq(  
    ("Peterson", "Marcus", 54),  
    ("Batey", "Edward", 36),  
    ("Bruce", "Karen", 35)  
).toDF("lastName", "firstName", "age")  
  
df.createOrReplaceTempView("people")  
spark.sql("SELECT firstName, age, age / 7.0 as dogYears  
FROM people where age < 50")
```





TUNING THE LEVEL OF PARALLELISM

Tuning the level of parallelism

- Every RDD has a **fixed number of partitions**
 - ▣ Determine the degree of parallelism when executing operations
- During aggregations or grouping operations, you can ask Spark to use a specific number of partitions
 - ▣ This will override defaults that Spark uses



Example: Tuning the level of parallelism

```
data = [("a", 3), ("b", 4), ("a", 1)]

sc.parallelize(data).
  reduceByKey(lambda x, y: x + y) #default

sc.parallelize(data).
  reduceByKey(lambda x, y: x + y, 10) #Custom
```



What if you want to tune parallelism outside of grouping and aggregation operations?

- There is `repartition()`
 - ▣ **Shuffles data across the network** to create a new set of partitions
 - ▣ Very **expensive operation!**
- There is the `coalesce()` operation
 - ▣ Allows avoiding data movement
 - But only if you are **decreasing** the number of partitions
 - ▣ Check `rdd.getNumPartitions()` and make sure you are coalescing to fewer partitions than current



DATASETS Vs DATAFRAMES



Datasets vs DataFrames

- In Spark's supported languages, Datasets make sense only in Java and Scala, whereas in Python and R only DataFrames make sense
- This is because Python and R are **not compile-time type-safe**
 - ▣ Types are dynamically *inferred* or assigned during execution, not during compile time
- The reverse is true in Scala and Java: types are bound to variables and objects at compile time



Scala **typing**: explicit when you must, inferred when you can

- You can tell Scala exactly what you mean

- ▣ `val number: Int = 25`

- Or let Scala figure it out for you

- ▣ `val inferredNumber = 25`

- Scala infers this as an `Int`

- `val` vs `var`

- ▣ `val` creates an **immutable** value (constant).

- ▣ `var` creates a **mutable** variable you can reassign

- `var anotherMutableInt = 5`

- `anotherMutableInt = 15`

Scala's type inference often lets you omit type declarations, which makes code leaner.

However, for clarity (and also when inference may mislead) explicit type declarations are valuable!





ORCHESTRATION PLANS

Executing Spark code in clusters: Overview

- Write DataFrame/Dataset/SQL Code
- If the code is valid, Spark converts this to a **Logical Plan**
- Spark transforms this Logical Plan to a **Physical Plan**, checking for optimizations along the way
- Spark then executes this Physical Plan (which involves RDD manipulations) on the cluster

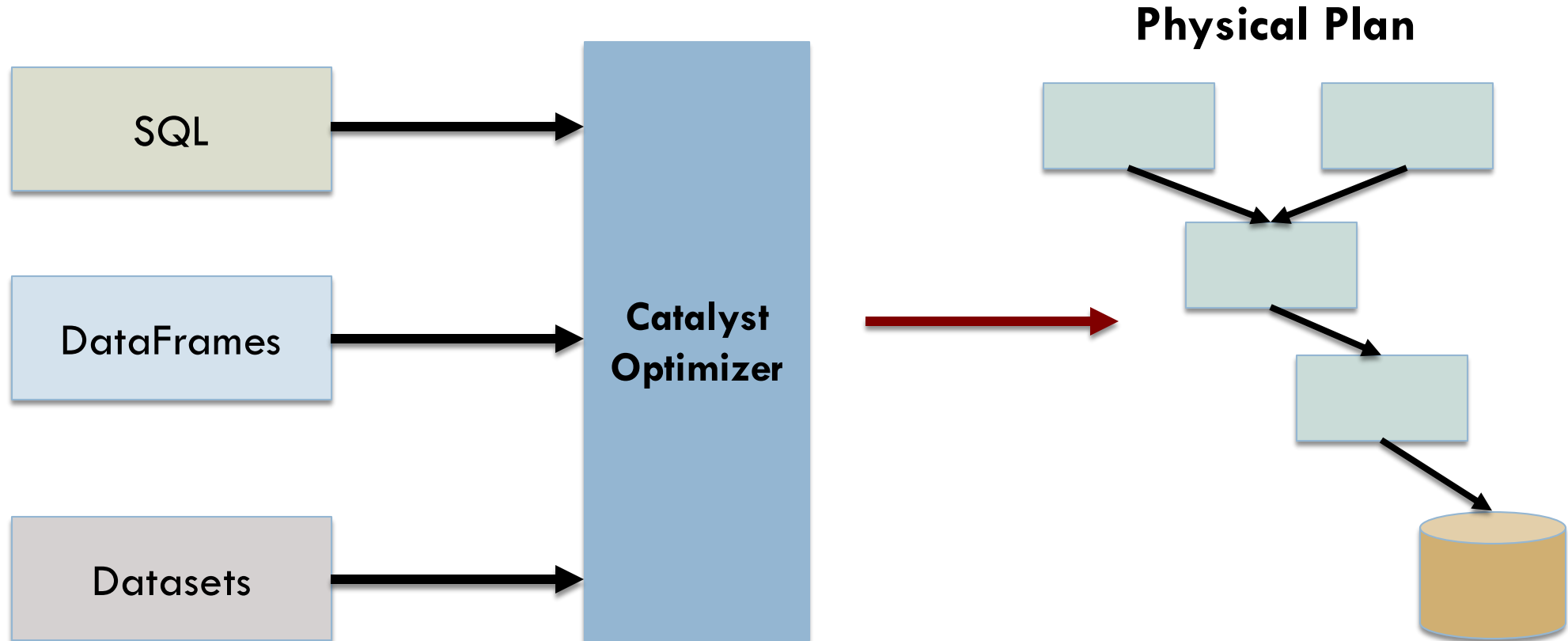


Once you have the code ready

- Code is submitted either through the console or via a submitted job
- This code passes through the **Catalyst Optimizer**
 - ▣ Decides *how* the code should be executed
 - ▣ Lays out a **plan** for doing so before, finally, the code is *run*
 - And the result returned to the user



The Catalyst Optimizer



Logical Planning

- The **logical plan** only represents a set of abstract transformations
 - ▣ Does not refer to executors or drivers
 - ▣ Simply converts the user's set of expressions into the most optimized version
- Converting user's code into an *unresolved* logical plan
 - ▣ This plan is unresolved because although your code might be valid, the tables or columns that it refers to might or might not exist

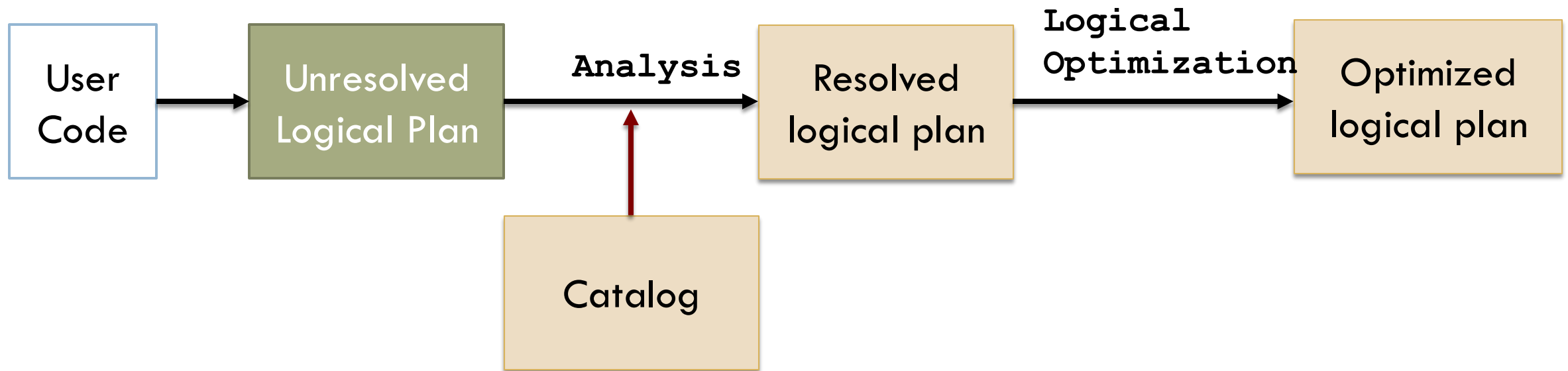


How are columns and tables resolved?

- ❑ Spark uses the **catalog**, a repository of all table and DataFrame information, to resolve columns and tables in the analyzer optimizations
- ❑ The analyzer might reject the unresolved logical plan if the required table or column name does not exist in the catalog
- ❑ If the analyzer can resolve it, the result is passed through the Catalyst Optimizer



The Structured API Logical Planning Process



Catalyst Optimizer

- A **collection of rules** that attempt to optimize the logical plan by pushing down predicates or selections
- Catalyst is **extensible**
 - ▣ Users can include their own rules for domain-specific optimizations



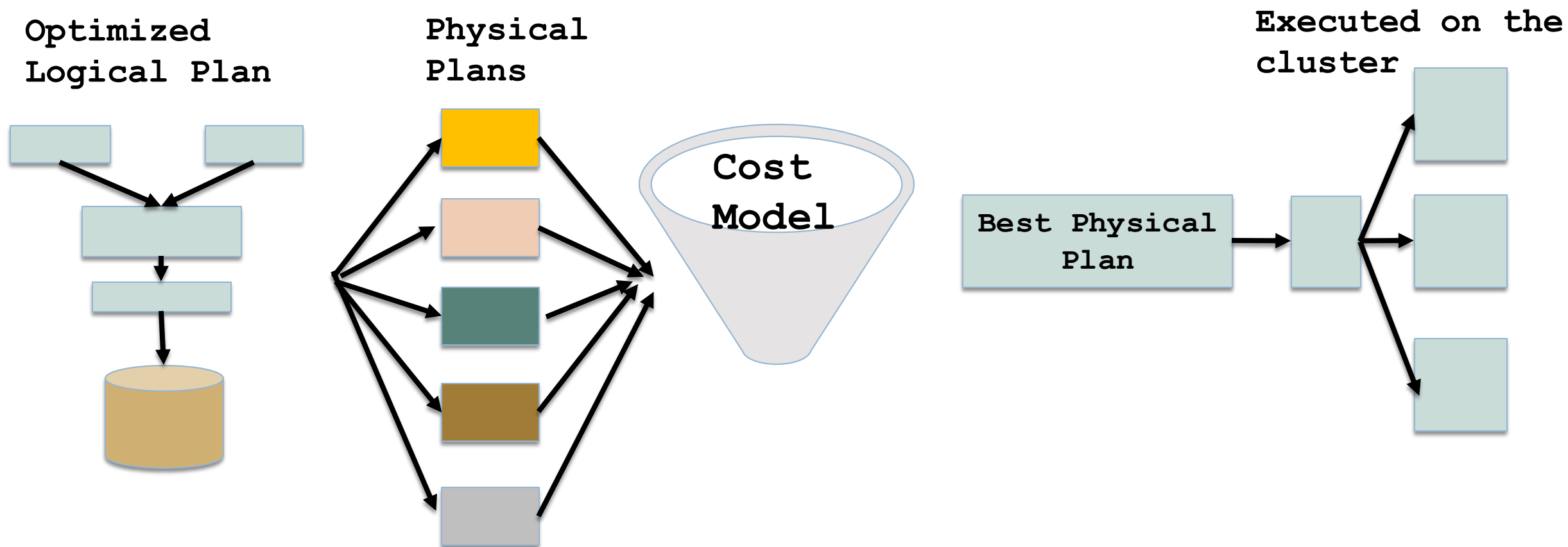
- The physical plan specifies how the logical plan will execute on the cluster
- Involves generating *different* physical execution strategies and comparing them through a **cost model**
- An example of the cost comparison might be choosing how to perform a given join by looking at the physical attributes of a given table
 - ▣ How big the table is or
 - ▣ How big its partitions are



- Physical planning results in a series of RDDs and transformations
- This is why Spark is also referred to as a compiler
 - ▣ Takes queries in DataFrames, Datasets, and SQL and compiles them into RDD transformations



The Physical Planning Process



Execution

- Spark performs further optimizations at runtime
- Generating native Java bytecode that can remove entire tasks or stages during execution
- Finally, the result is returned to the user





WIDE AND NARROW TRANSFORMATIONS

Transformations and Dependencies

- Two categories of **dependencies**
 - ▣ Narrow
 - Each partition of the parent RDD is used by **at most one partition of the child** RDD
 - ▣ Wide
 - Multiple child RDD partitions may depend on a single parent RDD partition
- The narrow versus wide distinction has significant implications for the way Spark evaluates a transformation and ...
 - ▣ consequently, for its performance

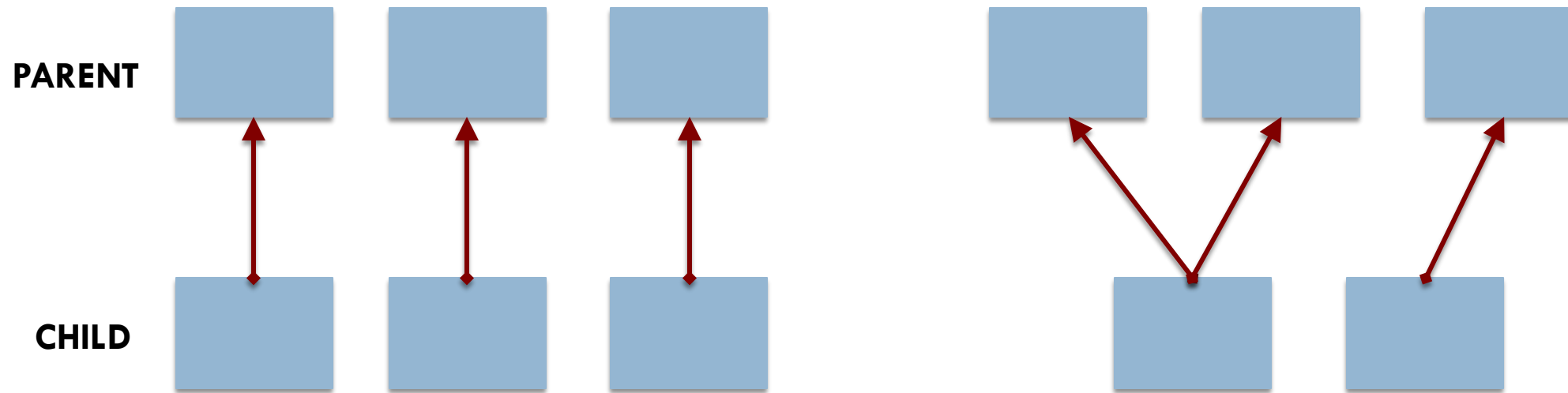


Narrow Transformations

- Narrow transformations are those in which each partition in the child RDD has simple, finite dependencies on partitions in the parent RDD
- Dependencies can be **determined at design time**, irrespective of the values of the records in the parent partitions
- Partitions in narrow transformations can either depend on:
 - ▣ One parent (such as in the `map` operator), or
 - ▣ A unique subset of the parent partitions that is known at design time (`coalesce`)
- Narrow transformations can be executed on an arbitrary subset of the data without any information about the other partitions



Dependencies between partitions for **narrow** transformations

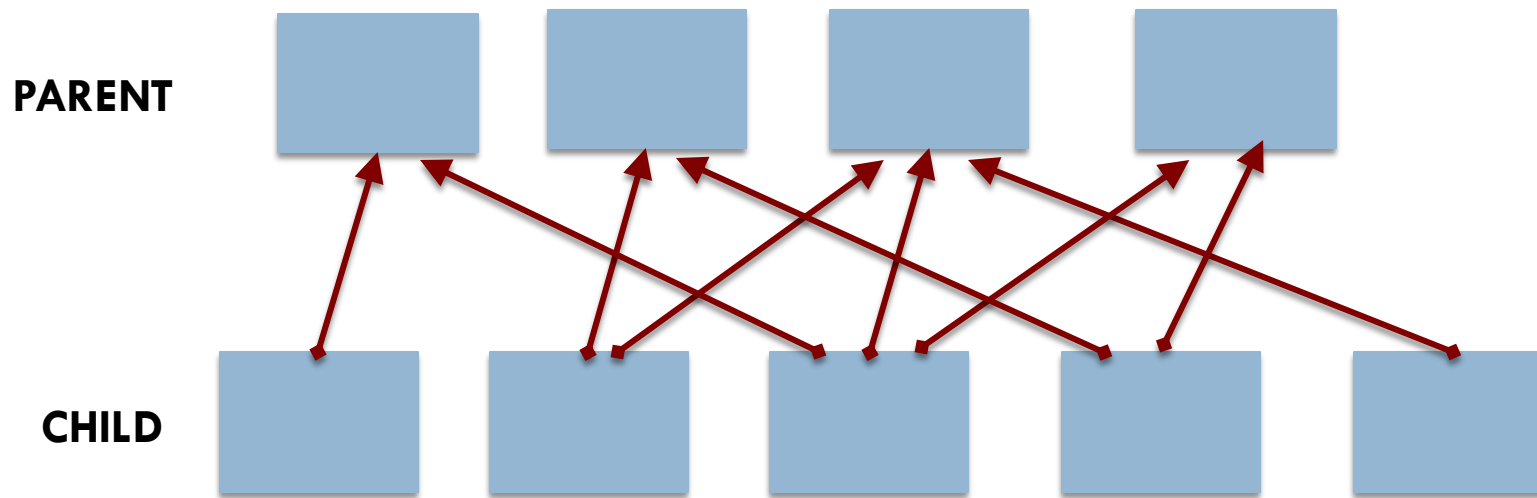


Wide Transformations

- Transformations with **wide dependencies** cannot be executed on arbitrary rows
- Require the data to be partitioned in a particular way, e.g., according to the **value** of their key
 - ▣ In sort, for example, records have to be partitioned so that keys in the same range are on the same partition
- Transformations with wide dependencies include `sort`, `reduceByKey`, `groupByKey`, `join`, and anything that calls the `rePartition` function



Dependencies between partitions for wide transformations



Wide dependencies **cannot be known fully before the data is evaluated**

The dependency graph for any operations that cause a **shuffle** (such as `groupByKey`, `reduceByKey`, `sort`, and `sortByKey`) follows this pattern



PAIR RDDs: WHAT TO WATCH FOR



Despite their utility, key/value operations can lead to a number of performance issues

- Most expensive operations in Spark fit into the key/value pair paradigm
 - ▣ Because **most wide transformations** are key/ value transformations,
 - And most require some fine tuning and care to be performant



In particular, operations on key/value pairs can cause ...

1. Out-of-memory errors in the driver
 2. Out-of-memory errors on the executor nodes
 3. Shuffle failures
 4. “Straggler tasks” or partitions, which are especially slow to compute
- The last three performance issues are all most often caused by **shuffles associated with the wide transformations**



Memory errors in the driver, are usually caused by actions

- ❑ Several key/value actions (including `countByKey`, `countByValue`, `lookUp`, and `collectAsMap`) return data to the driver
- ❑ In most instances they return *unbounded* data since the number of keys and the number of values are unknown
- ❑ In addition to number of records, the size of each record is an important factor in causing memory errors



Preventing out-of-memory errors with aggregation operations

[1 / 2]

- `combineByKey` and all of the aggregation operators built on top of it (`reduceByKey`, `foldLeft`, `foldRight`, `aggregateByKey`) may lead to memory errors if they cause the accumulator to become too large for one key
- What about `groupByKey`?
 - ▣ It is actually implemented using `combineByKey` where the accumulator is an iterator with all the data.



Preventing out-of-memory errors with aggregation operations

[2/2]

- Use functions that implement **map-side combinations**
 - ▣ Meaning that records with the same key are combined before they are shuffled
 - ▣ This can greatly reduce the shuffled read
- The following four functions are implemented to use map-side combinations
 - ▣ `reduceByKey`
 - ▣ `treeAggregate`: Use a tree pattern rather than a linear pattern to merge
 - ▣ `aggregateByKey`
 - ▣ `foldByKey` : Similar to `reduceByKey` but lets you use a starting value



The contents of this slide-set are based on the following references

- *Learning Spark: Lightning-Fast Big Data Analysis. 1st Edition. Holden Karau, Andy Konwinski, Patrick Wendell, and Matei Zaharia. O'Reilly. 2015. ISBN-13: 978-1449358624. [Chapters 1-4, 10]*
- Chambers, Bill, and Zaharia, Matei. *Spark: The Definitive Guide: Big Data Processing Made Simple. O'Reilly Media. ISBN-13: 978-1491912218. 2018. [Chapters 5 and 9].*
- Karau, Holden; Warren, Rachel. *High Performance Spark: Best Practices for Scaling and Optimizing Apache Spark. O'Reilly Media. 2017. ISBN-13: 978-1491943205. [Chapter 2]*

