



A complex, abstract background graphic is composed of numerous small, glowing blue and purple dots connected by thin lines, forming a three-dimensional structure that resembles a neural network or a circuit board. This graphic serves as a backdrop for the text elements.

Build. Unify. Scale.

WIFI SSID:Spark+AISummit | Password: UnifiedDataAnalytics

Internals of Speeding up PySpark with Arrow

Ruben Berenguel, Consultant

#UnifiedDataAnalytics #SparkAISummit



WHOAMI

- > RUBEN BERENGUEL (@BERENGUEL)
- > PHD IN MATHEMATICS
- > (BIG) DATA CONSULTANT
- > LEAD DATA ENGINEER USING PYTHON, GO AND SCALA
- > RIGHT NOW AT AFFECTV

WHAT IS PANDAS?

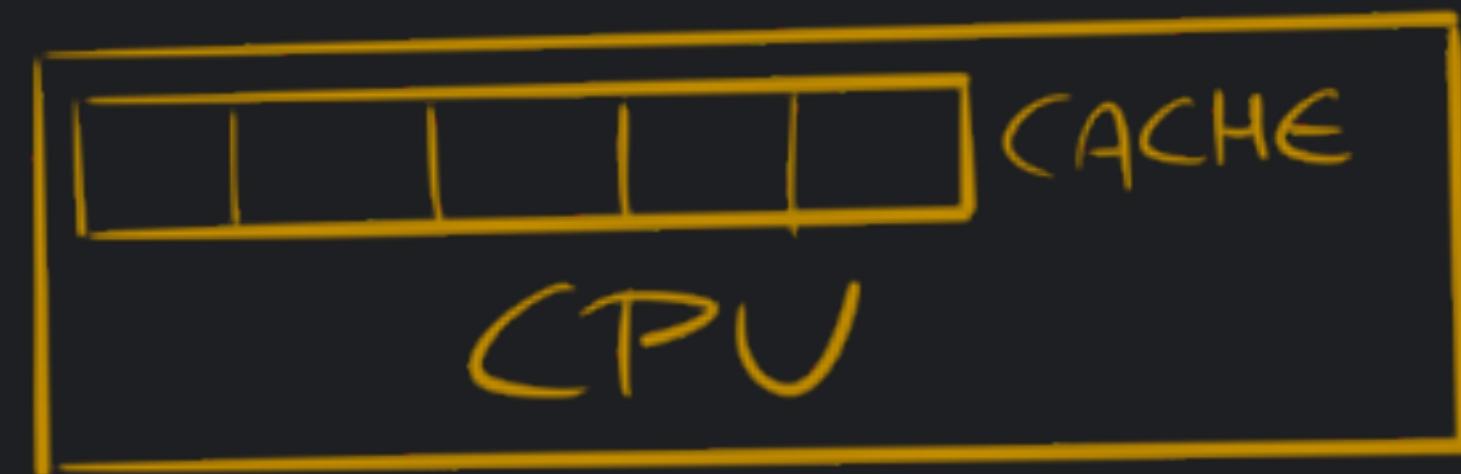
- > PYTHON DATA ANALYSIS LIBRARY
- > USED EVERYWHERE DATA AND PYTHON APPEAR IN JOB OFFERS
- > EFFICIENT (IS COLUMNAR AND HAS A C AND CYTHON BACKEND)

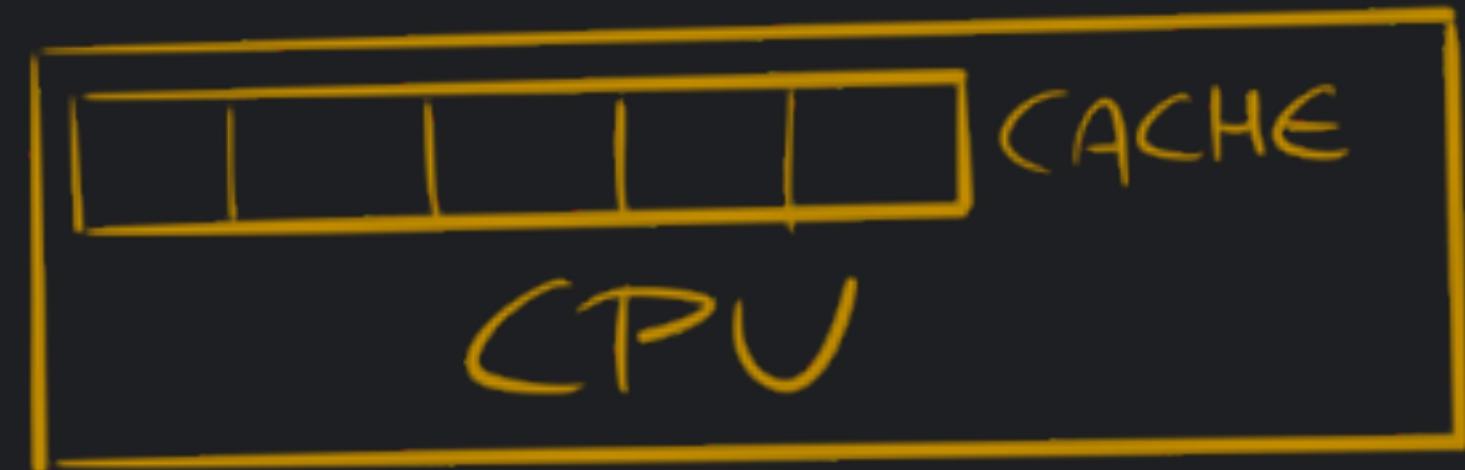
c_1	c_2	c_3	\dots	c_{42}	c_{43}	\dots
\vdots	\vdots	\vdots				

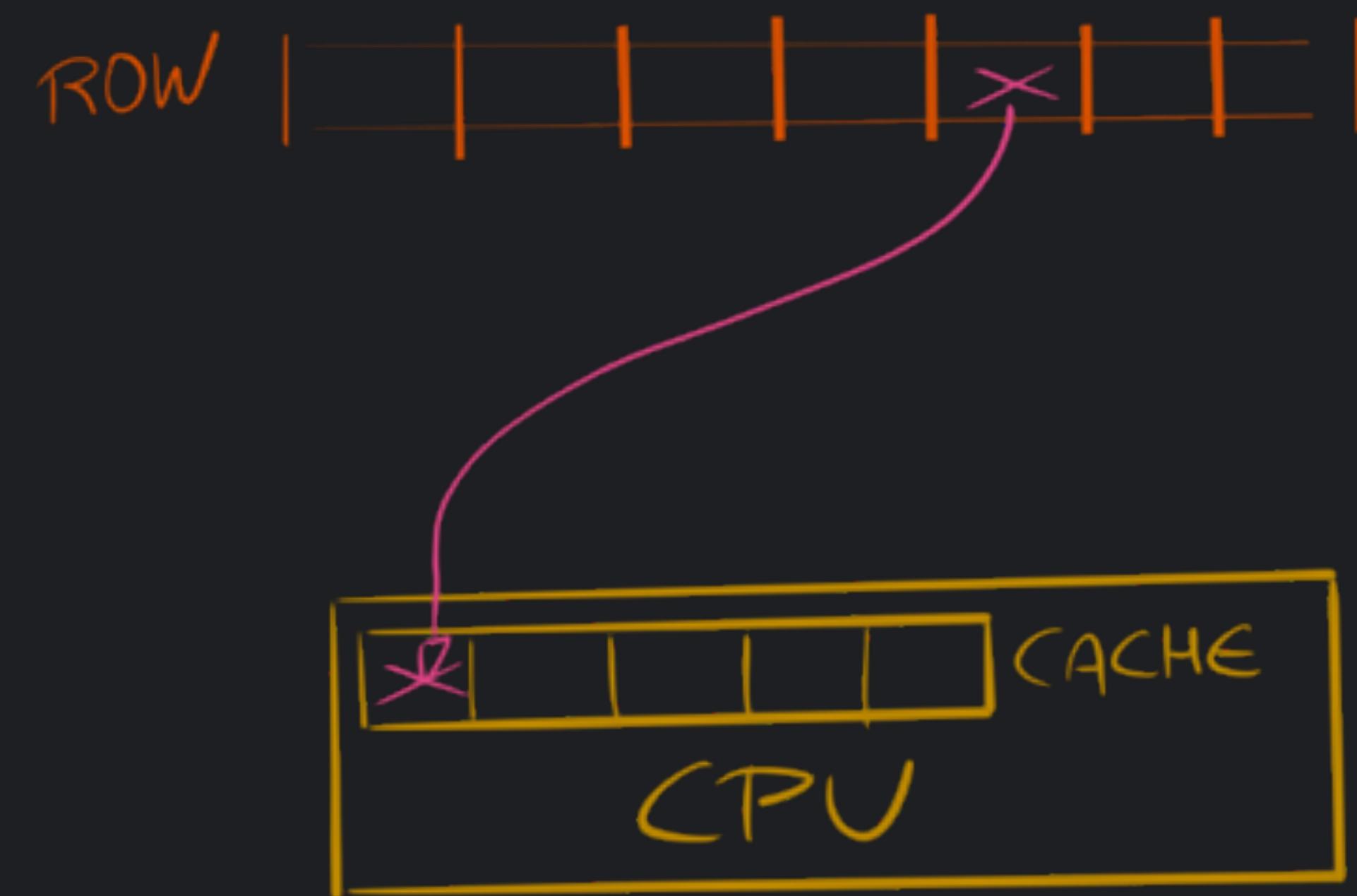
ROW

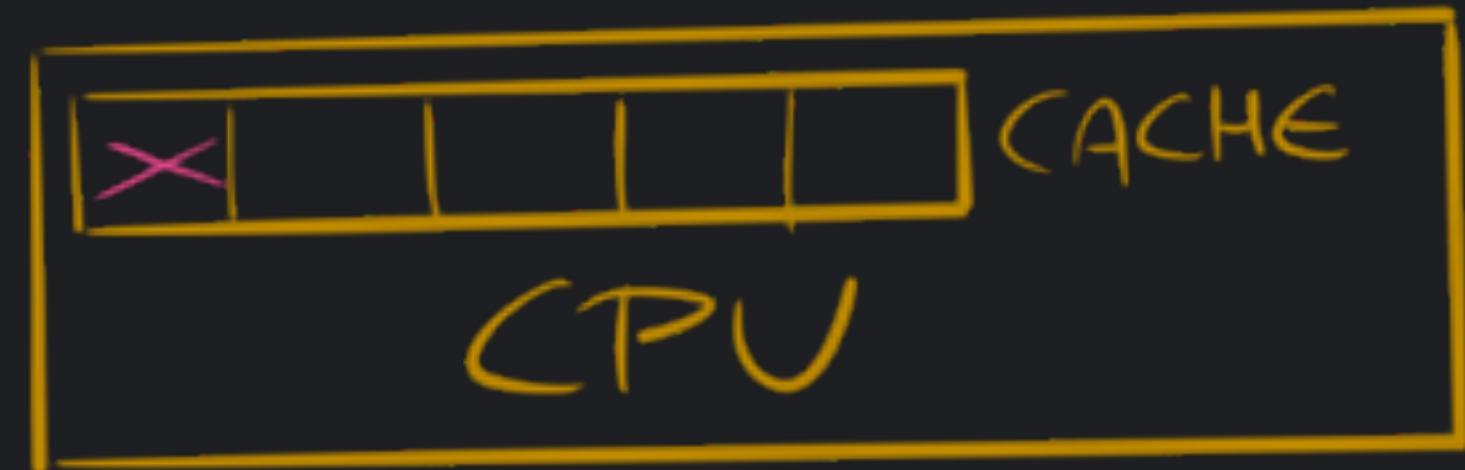


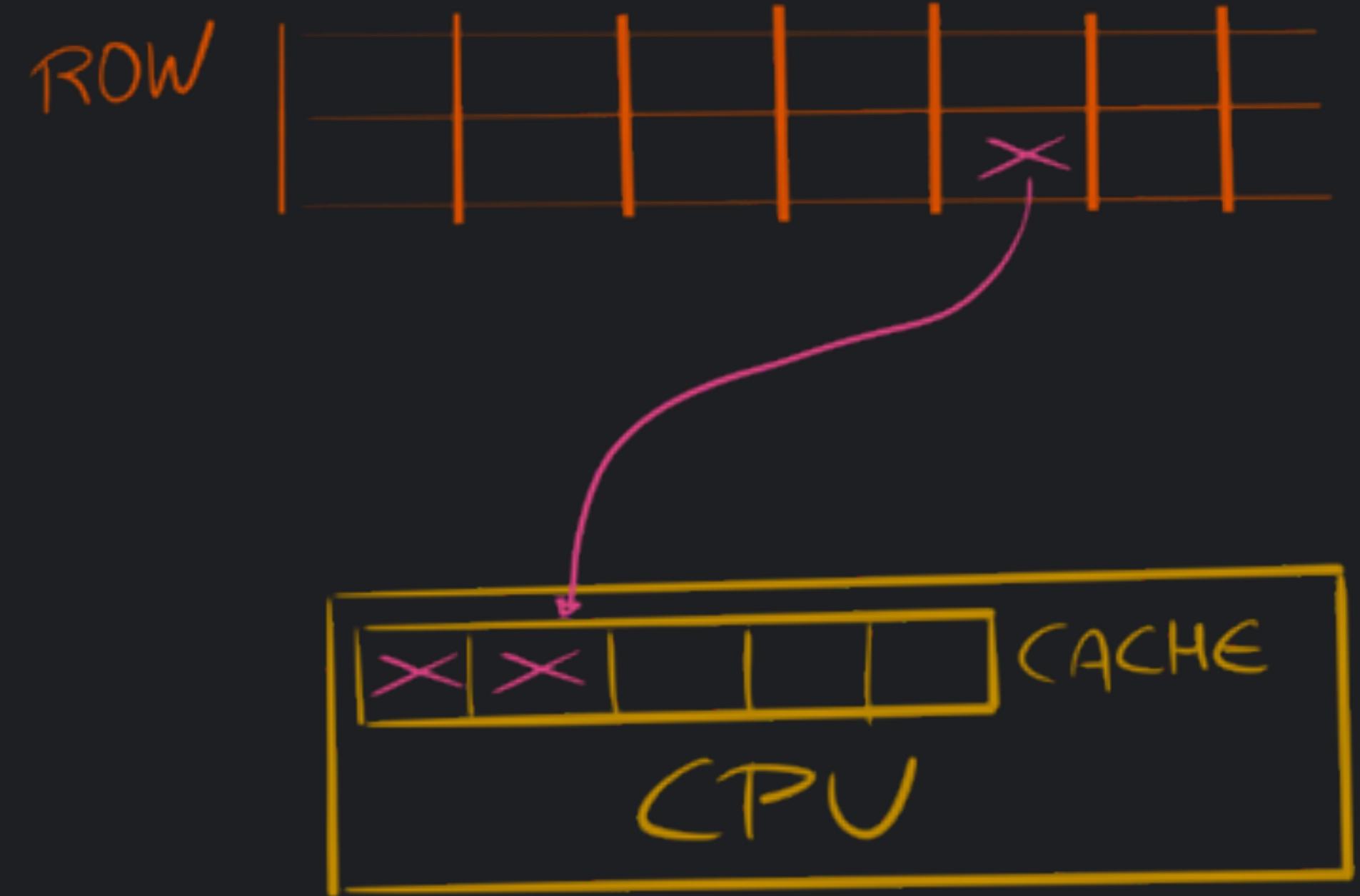
CPU

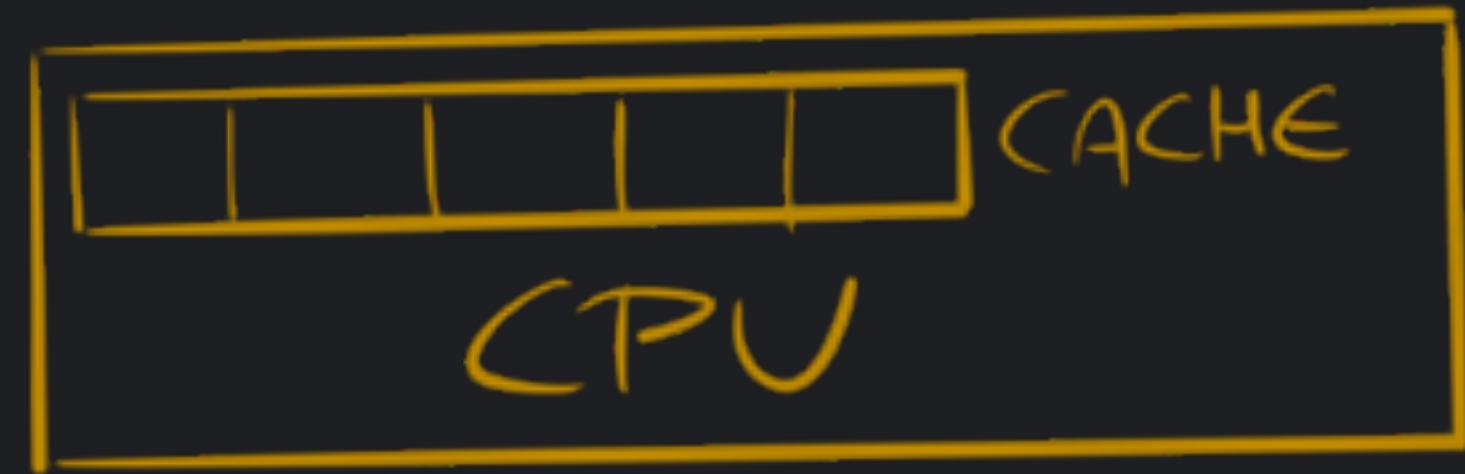


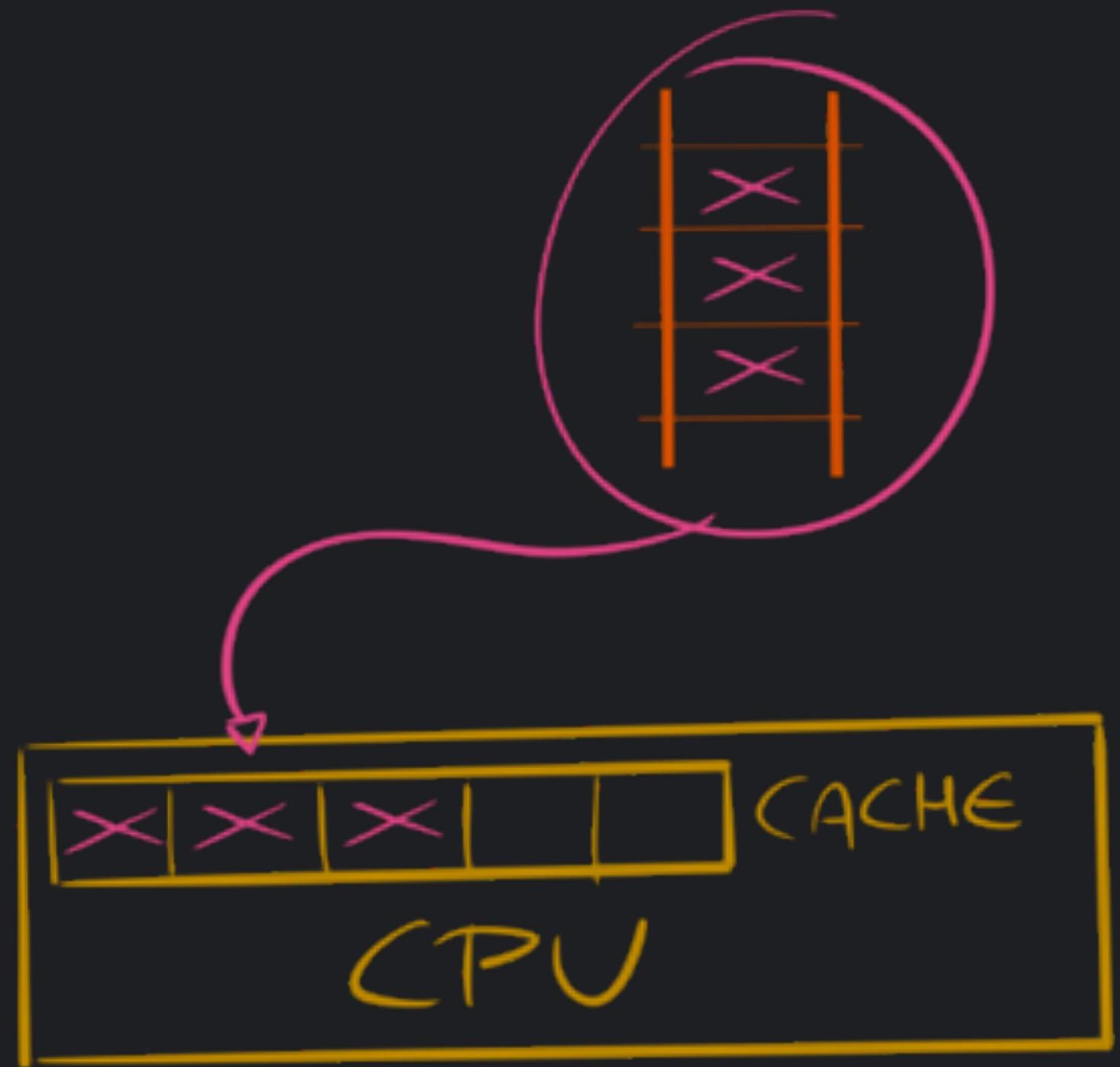












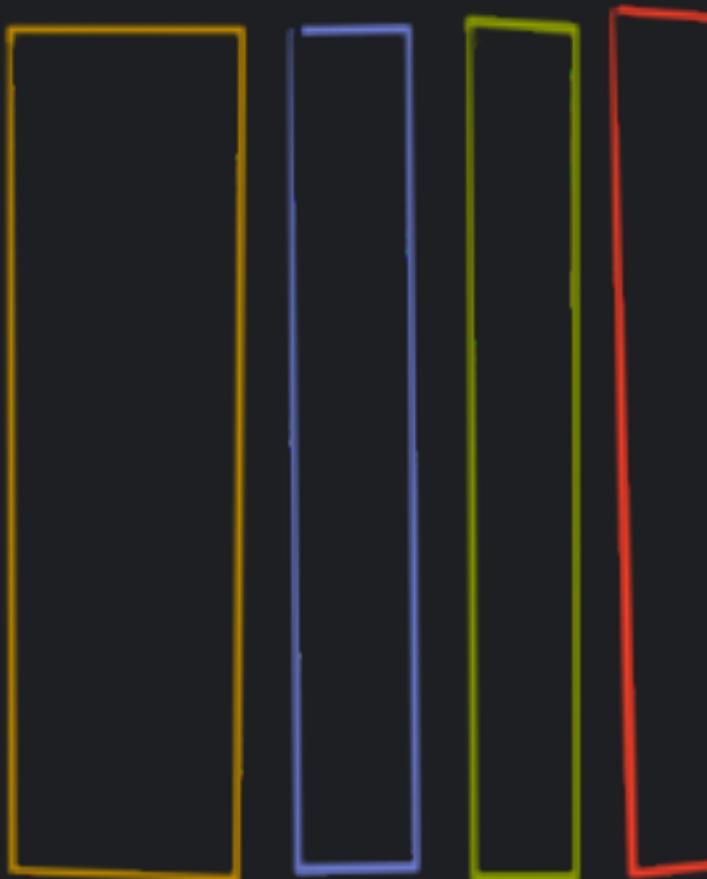
HOW DOES PANDAS MANAGE COLUMNAR DATA?

Pandas DataFrame

	A	B	C	D	
1	4.4	a	2	0.1	
2	3.1	b	42	0.4	
3	2	c	8	0.9	
4	8.3	d	15	0.3	
:					

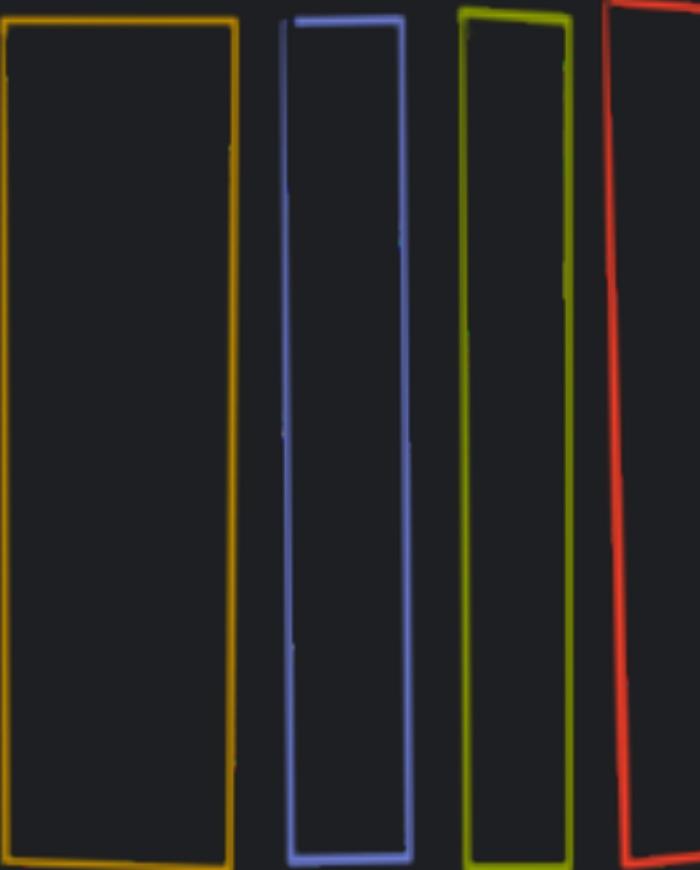
Pandas DataFrame

	A	B	C	D
1	4.4	a	2	0.1
2	3.1	b	42	0.4
3	2	c	8	0.9
4	8.3	d	15	0.3
:				



Pandas DataFrame

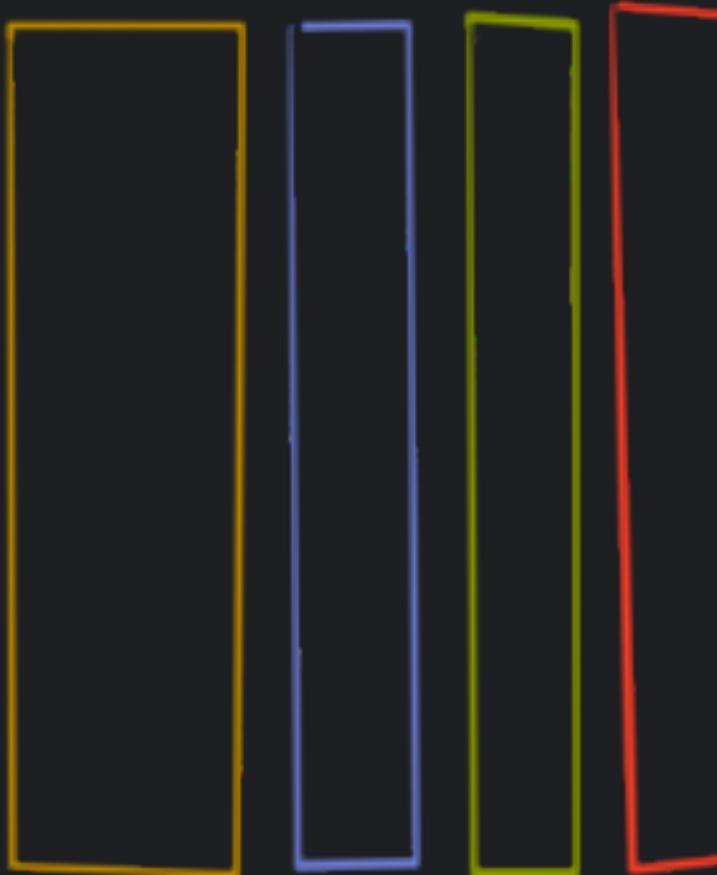
	A	B	C	D	
1	4.4	a	2	0.1	
2	3.1	b	42	0.4	
3	2	c	8	0.9	
4	8.3	d	15	0.3	
:					



INDEX
SECTION

Pandas DataFrame

	A	B	C	D	
1	4.4	a	2	0.1	
2	3.1	b	42	0.4	
3	2	c	8	0.9	
4	8.3	d	15	0.3	
:					

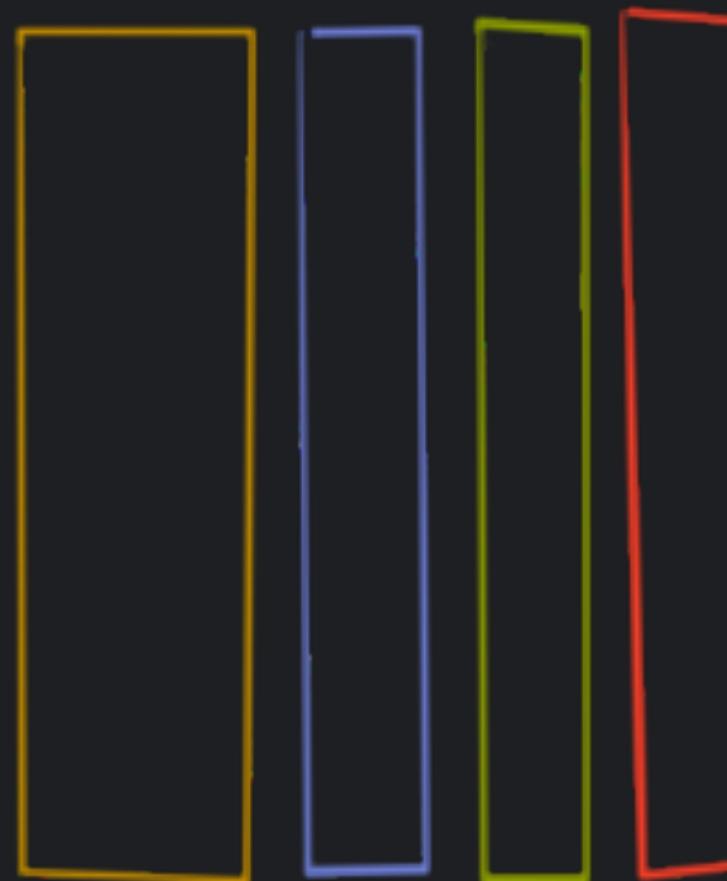


$2 \times N$ floatBlock

Pandas DataFrame

	A	B	C	D	
1	4.4	a	2	0.1	
2	3.1	b	42	0.4	
3	2	c	8	0.9	
4	8.3	d	15	0.3	
:					

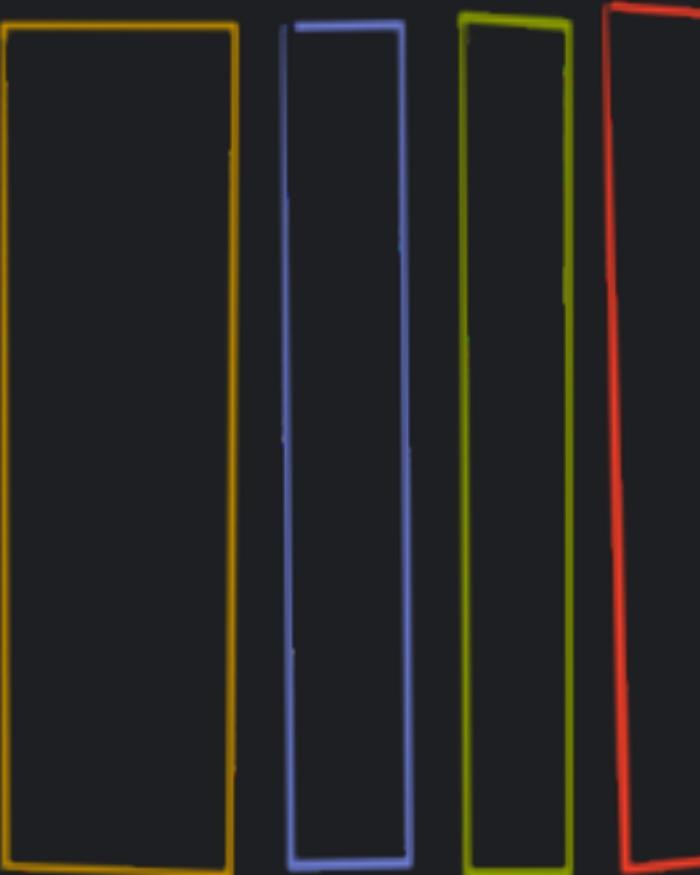
1×N ObjectBlock



Pandas DataFrame

	A	B	C	D
1	4.4	a	2	0.1
2	3.1	b	42	0.4
3	2	c	8	0.9
4	8.3	d	15	0.3
:				

1xN IntBlock



Pandas DataFrame

	A	B	C	D
1	4.4	a	2	0.1
2	3.1	b	42	0.4
3	2	c	8	0.9
4	8.3	d	15	0.3
:				

BlockManager



WHAT IS ARROW?

- > CROSS-LANGUAGE IN-MEMORY COLUMNAR FORMAT LIBRARY
- > OPTIMISED FOR EFFICIENCY ACROSS LANGUAGES
- > INTEGRATES SEAMLESSLY WITH PANDAS

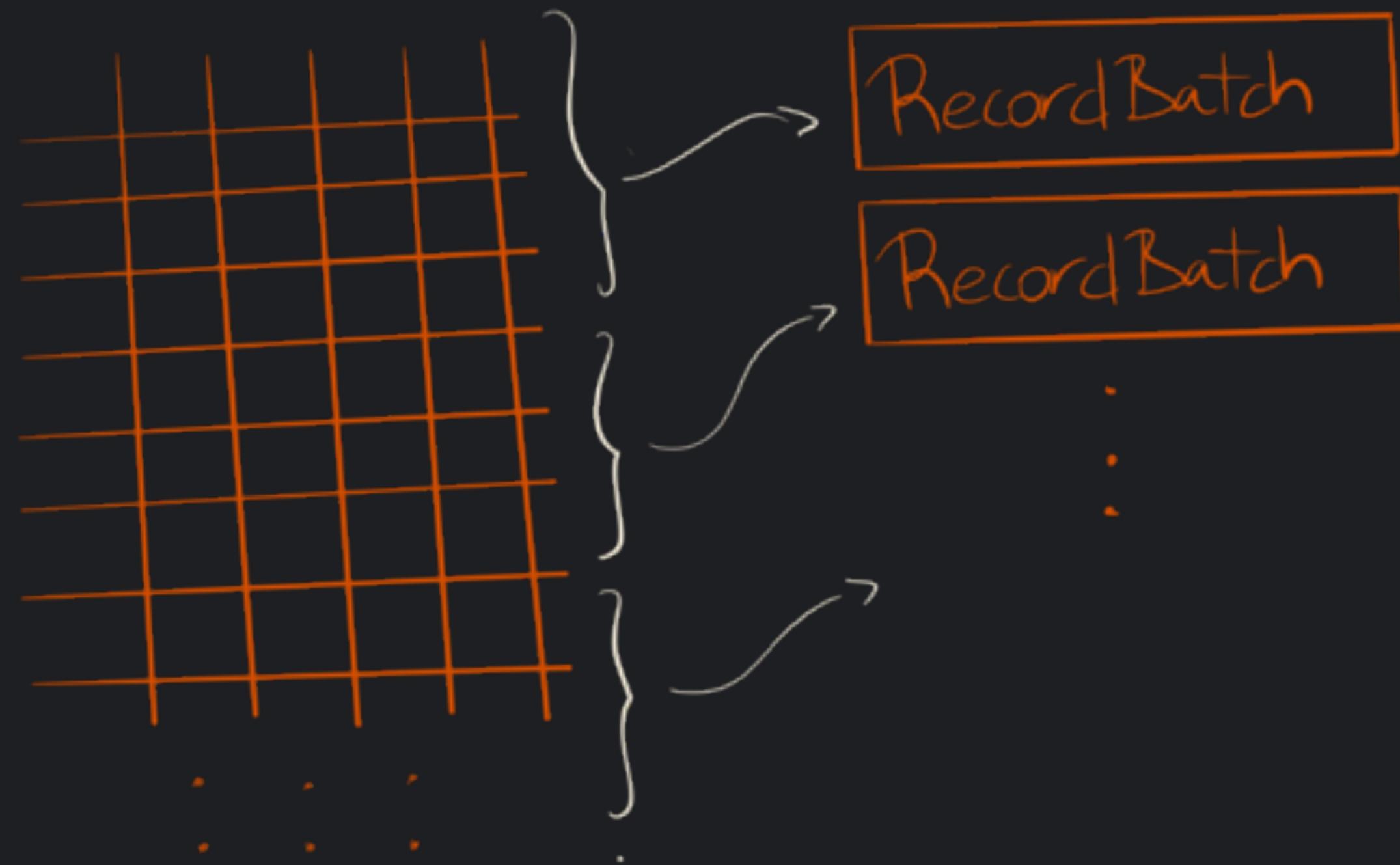
HOW DOES ARROW MANAGE COLUMNAR DATA?

Table

A 10x10 grid of orange lines on a black background. The grid consists of 9 horizontal rows and 9 vertical columns, creating a total of 81 square cells. The lines are thin and have a slightly irregular, hand-drawn appearance.

#UnifiedDataAnalytics #SparkAISummit

Table



Some rows

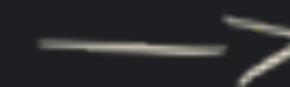
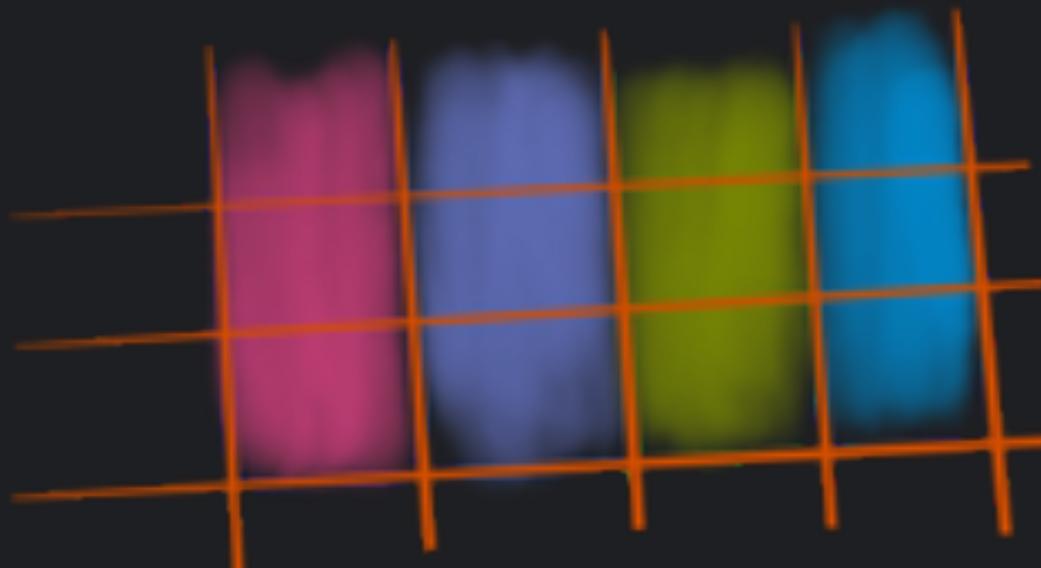


Some rows



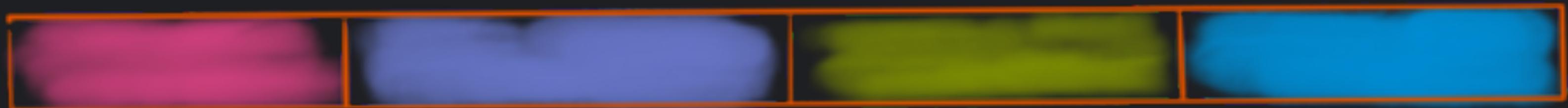
RecordBatch

Some rows



RecordBatch

(+ metadata)



Arrow Table

1	2	3	4
1	2	3	4
1	2	3	4
1	2	3	4
1	2	3	4

•

•



- > ARROW USES RecordBatches
- > PANDAS USES BLOCKS HANDLED BY A BlockManager
- > YOU CAN CONVERT AN ARROW Table INTO A PANDAS DataFrame EASILY

Arrow Table



Pandas BlockManager



WHAT IS SPARK?

- > DISTRIBUTED COMPUTATION FRAMEWORK
 - > OPEN SOURCE
 - > EASY TO USE
- > SCALES HORIZONTALLY AND VERTICALLY

HOW DOES SPARK WORK?

SPARK
USUALLY
RUNS ON TOP
OF A CLUSTER
MANAGER





AND A
DISTRIBUTED
STORAGE

Cluster Manager

Distributed Storage



A SPARK PROGRAM
RUNS IN THE DRIVER

THE DRIVER REQUESTS
RESOURCES FROM THE
CLUSTER MANAGER TO
RUN TASKS

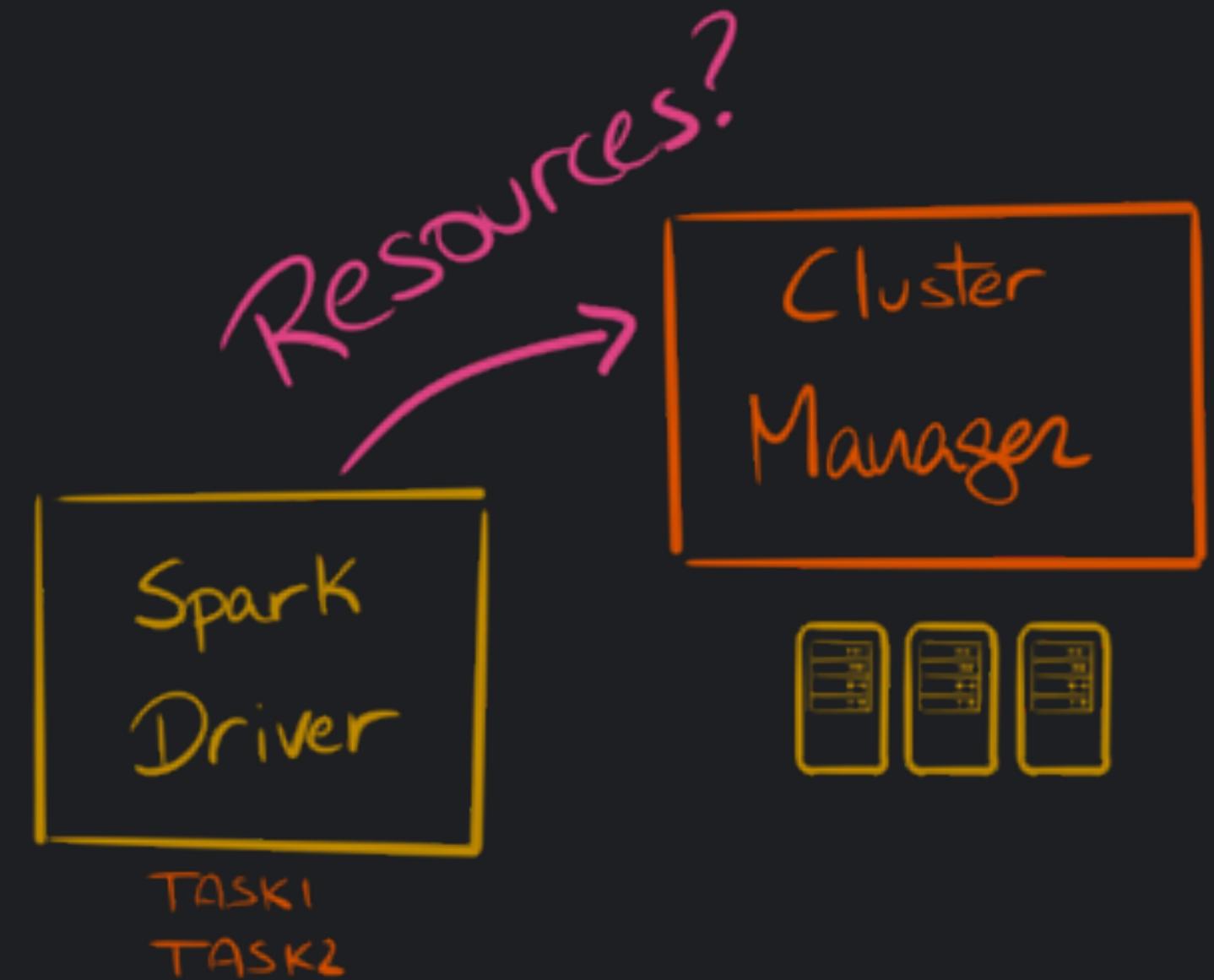
Spark
Driver



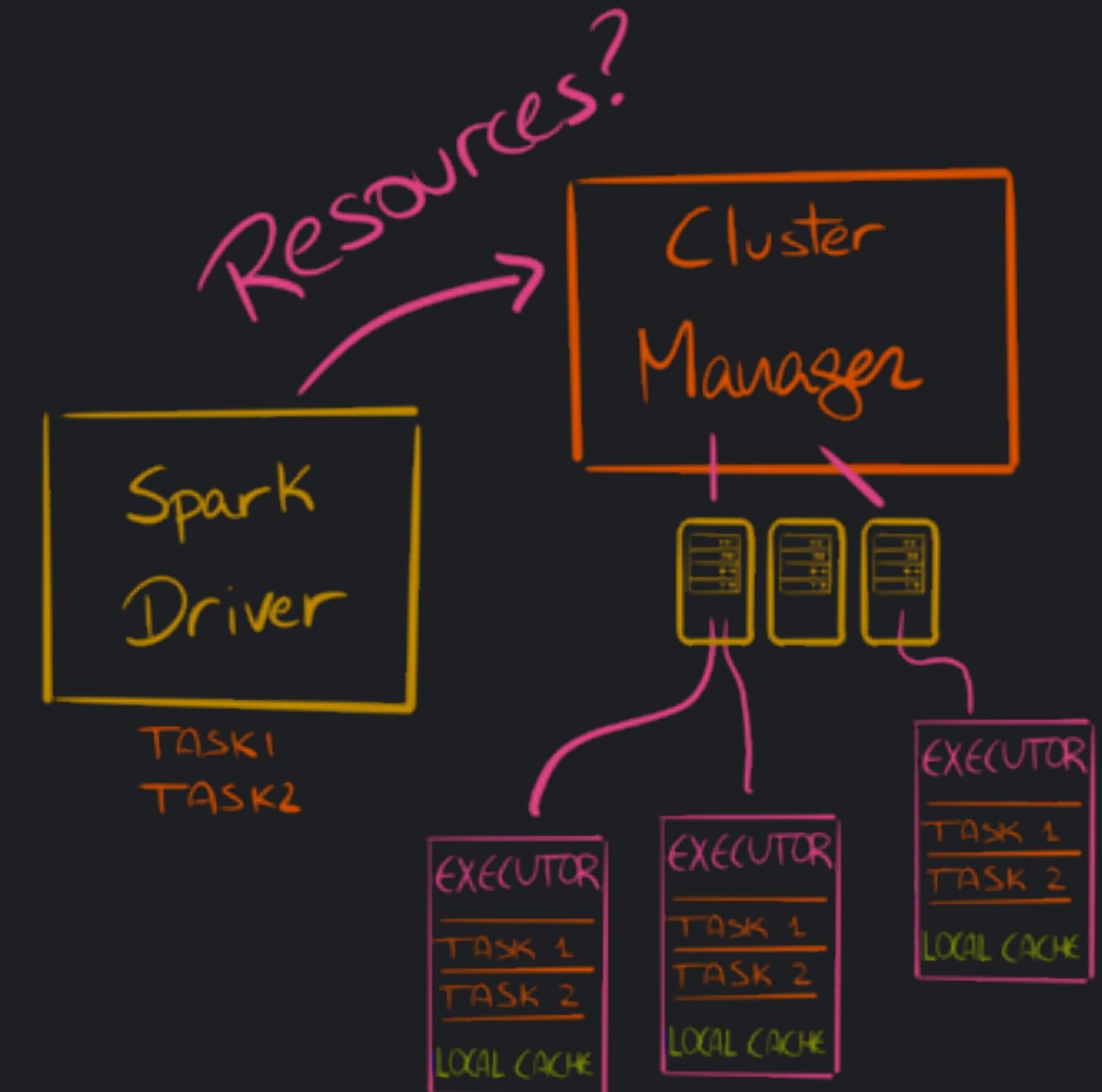
THE DRIVER REQUESTS RESOURCES FROM THE CLUSTER MANAGER TO RUN TASKS



THE DRIVER REQUESTS
RESOURCES FROM THE
CLUSTER MANAGER TO
RUN TASKS



THE DRIVER REQUESTS RESOURCES FROM THE CLUSTER MANAGER TO RUN TASKS



THE MAIN BUILDING BLOCK
IS THE RDD:
**RESILIENT DISTRIBUTED
DATASET**



RDD



RDD



Partitions

RDD



Partitions



RDD



RDD



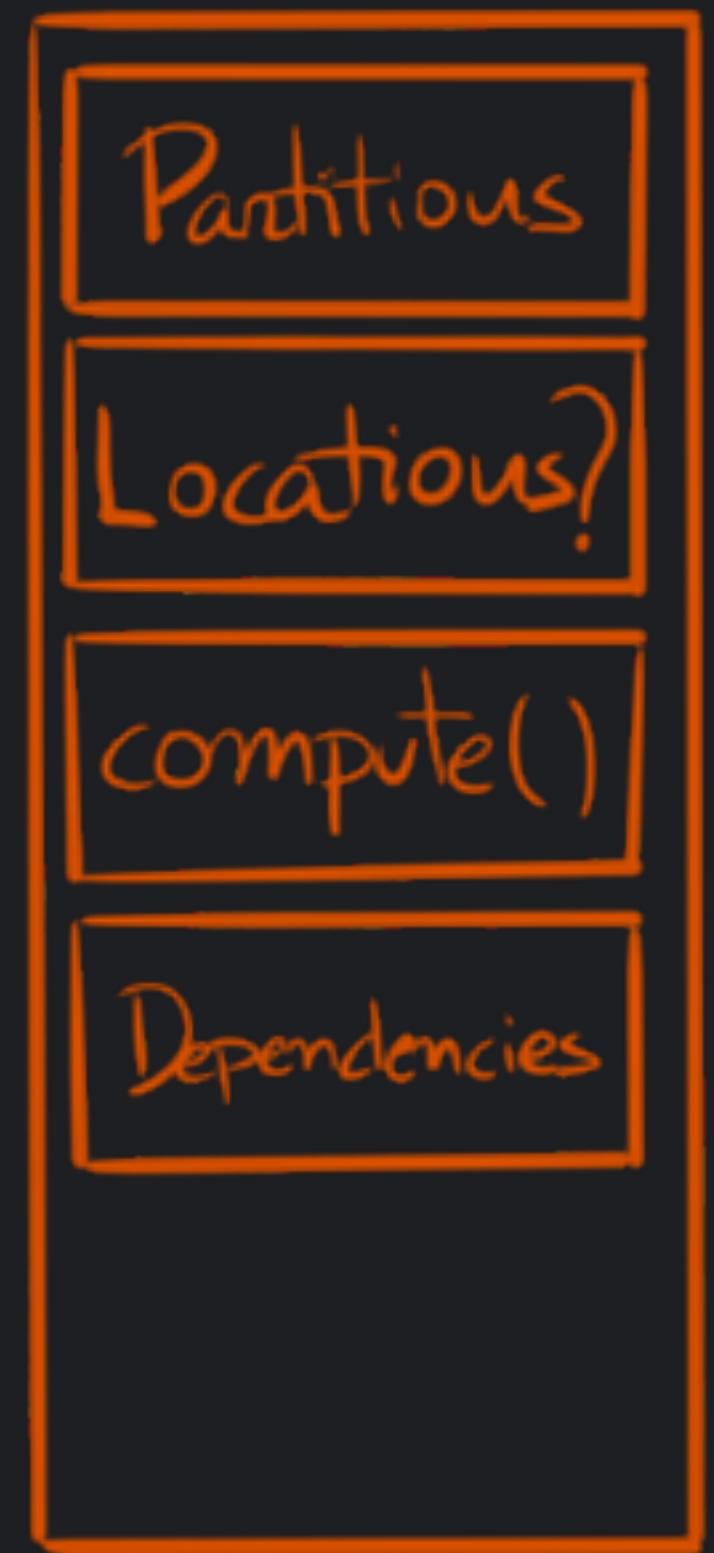
RDD



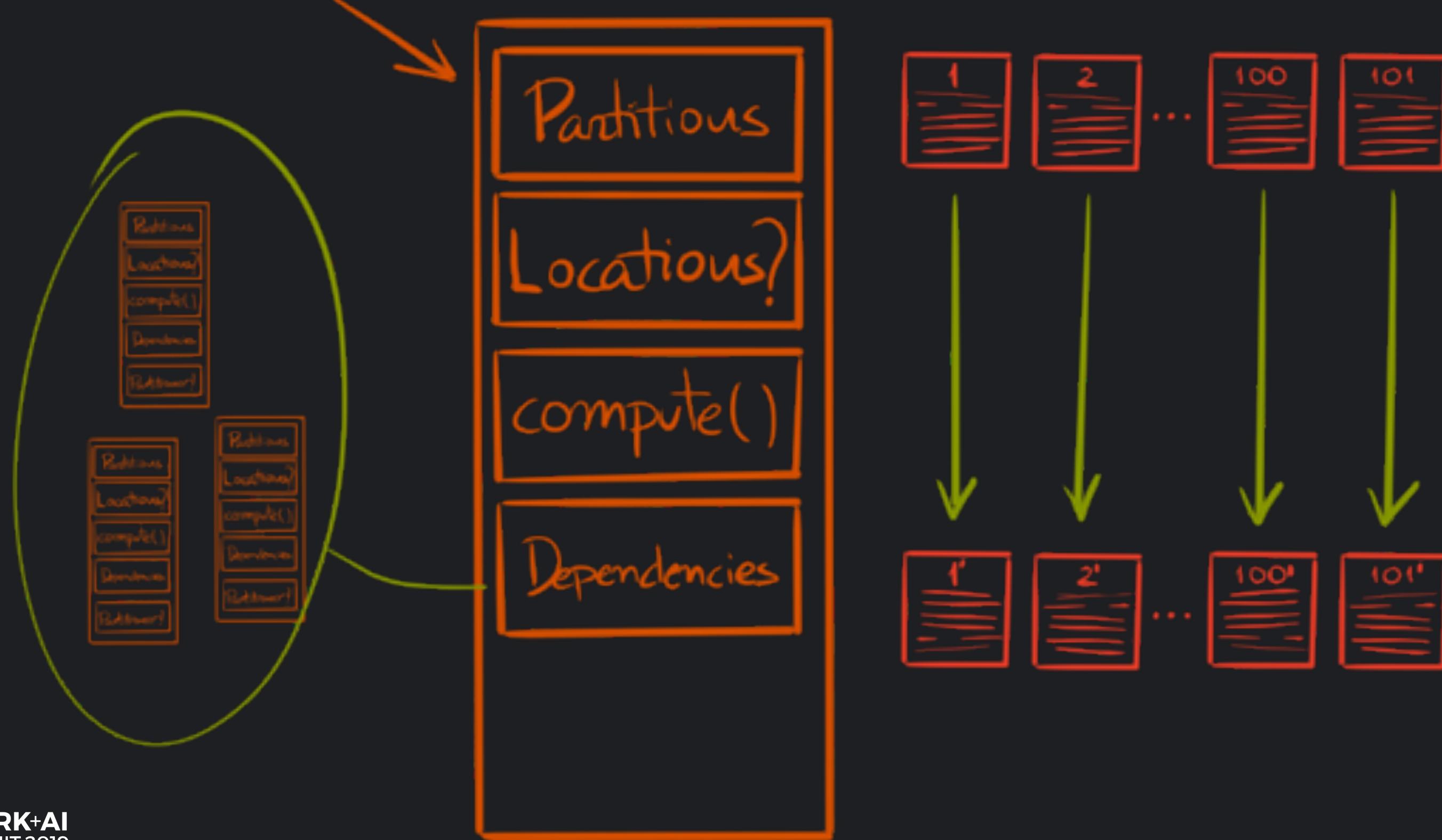
RDD



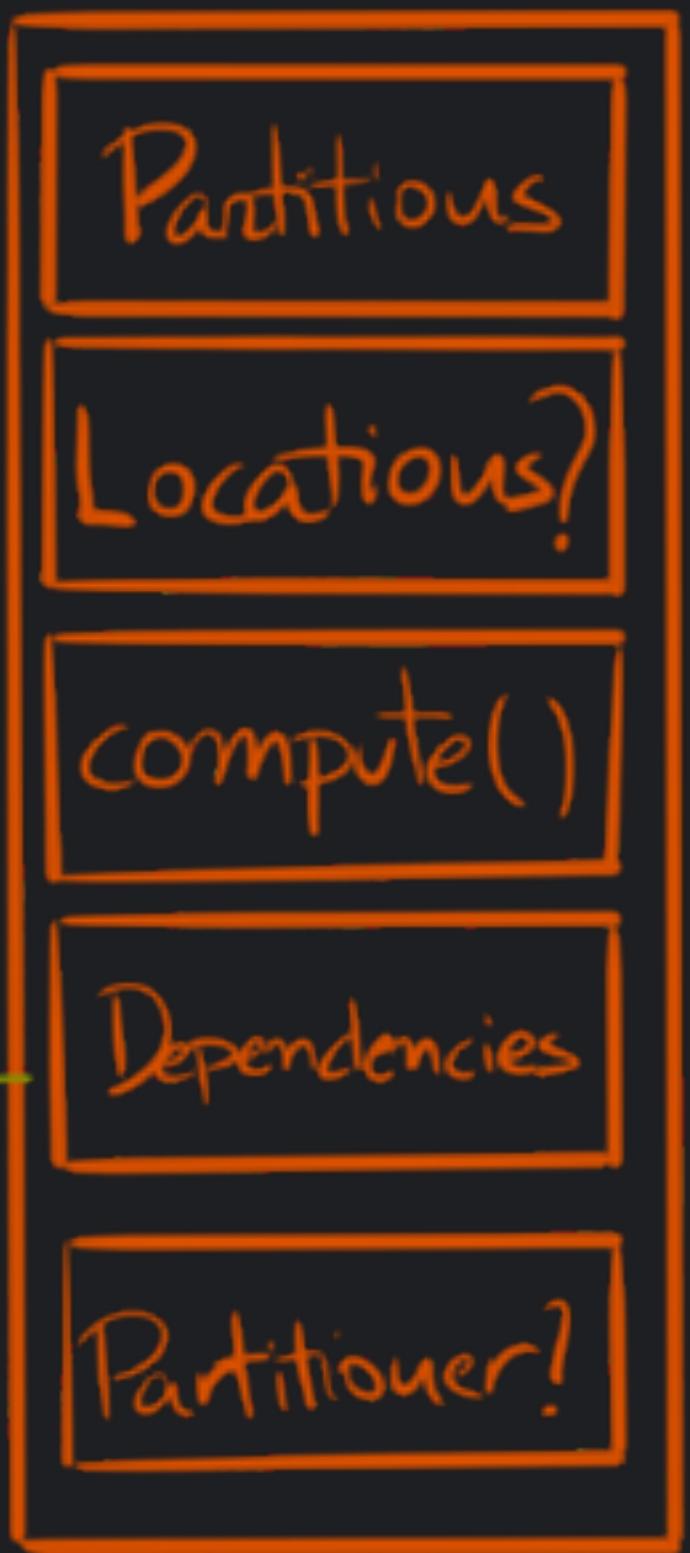
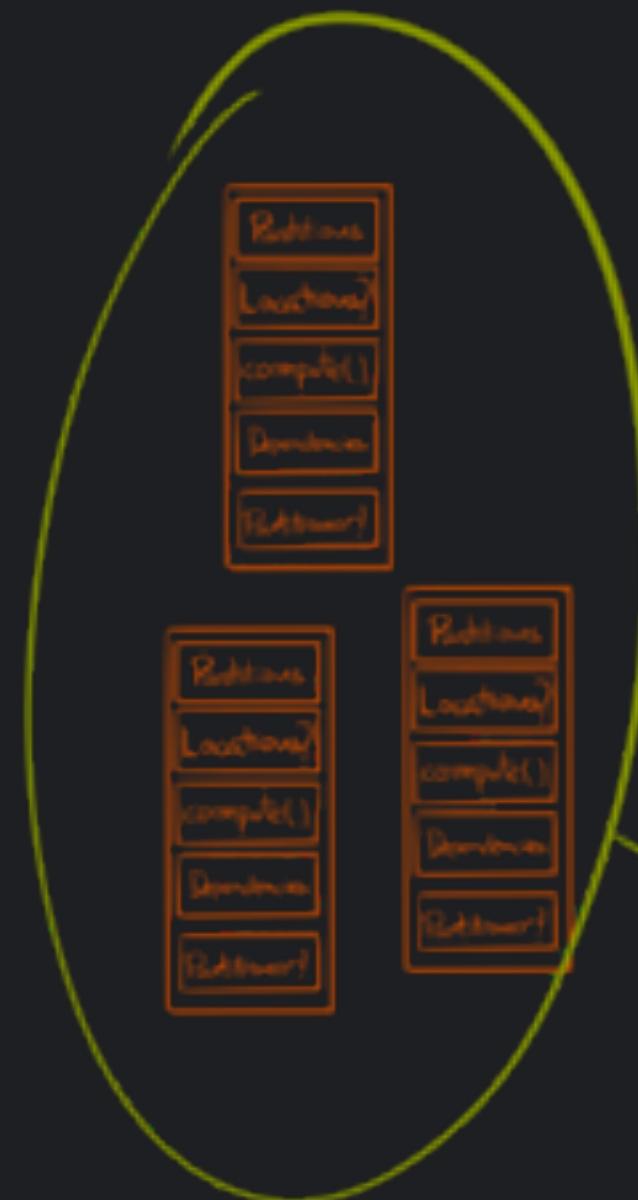
RDD



RDD



RDD



PYSPARK

PYSPARK OFFERS A
PYTHON API TO THE SCALA
CORE OF SPARK

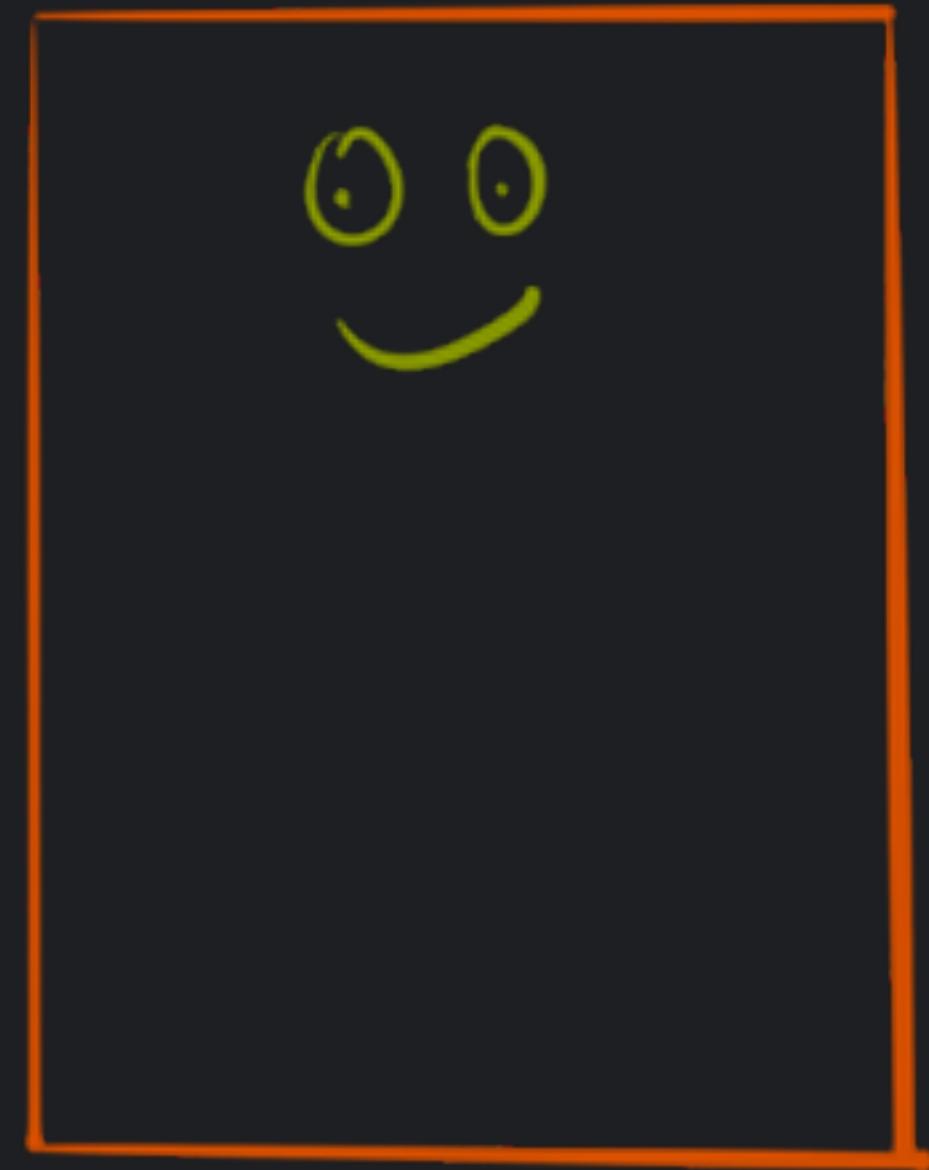
IT USES THE
PY4J BRIDGE

```
# Connect to the gateway
gateway = JavaGateway(
    gateway_parameters=GatewayParameters(
        port=gateway_port,
        auth_token=gateway_secret,
        auto_convert=True))

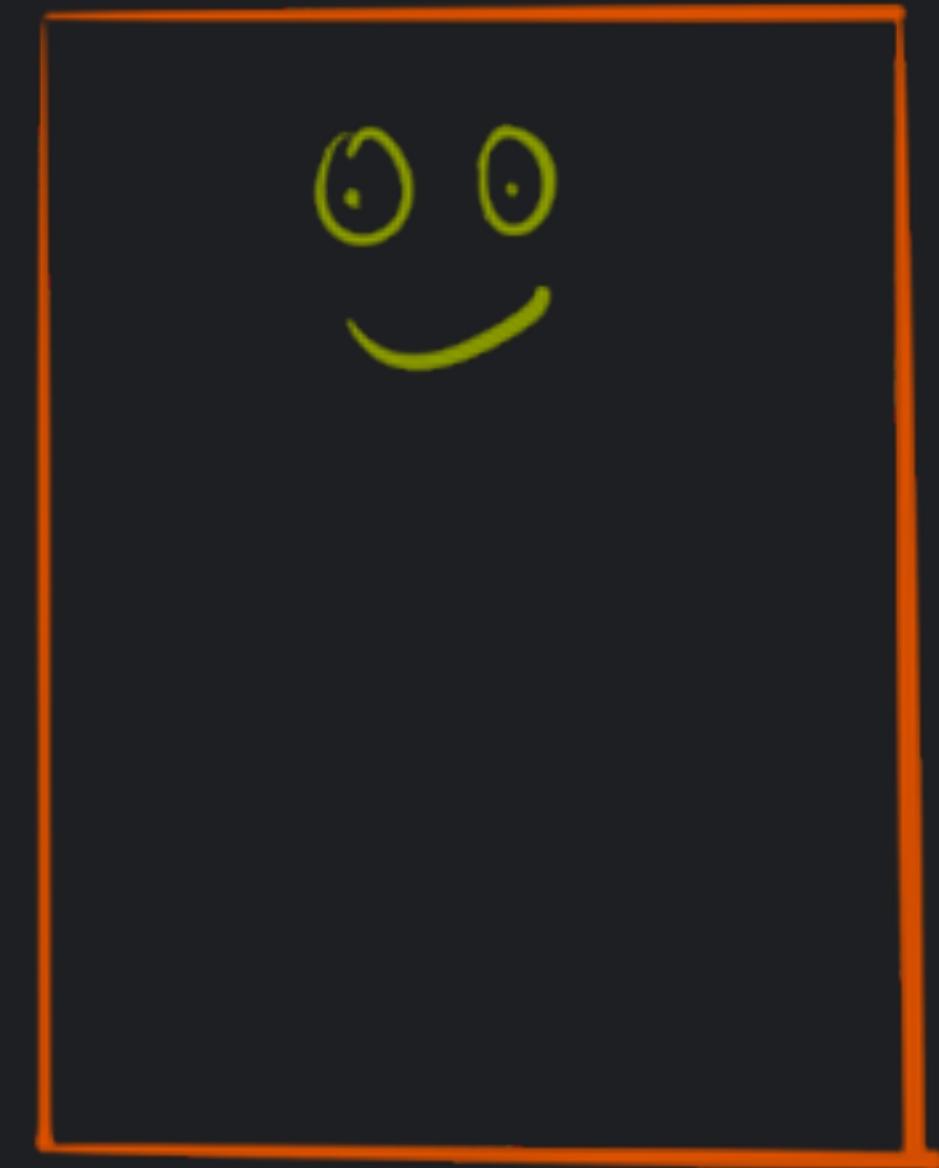
# Import the classes used by PySpark
java_import(gateway.jvm, "org.apache.spark.SparkConf")
java_import(gateway.jvm, "org.apache.spark.api.java.*")
java_import(gateway.jvm, "org.apache.spark.api.python.*")

.
.
.

return gateway
```



Python in
RDD Land



RDD in
Python Land

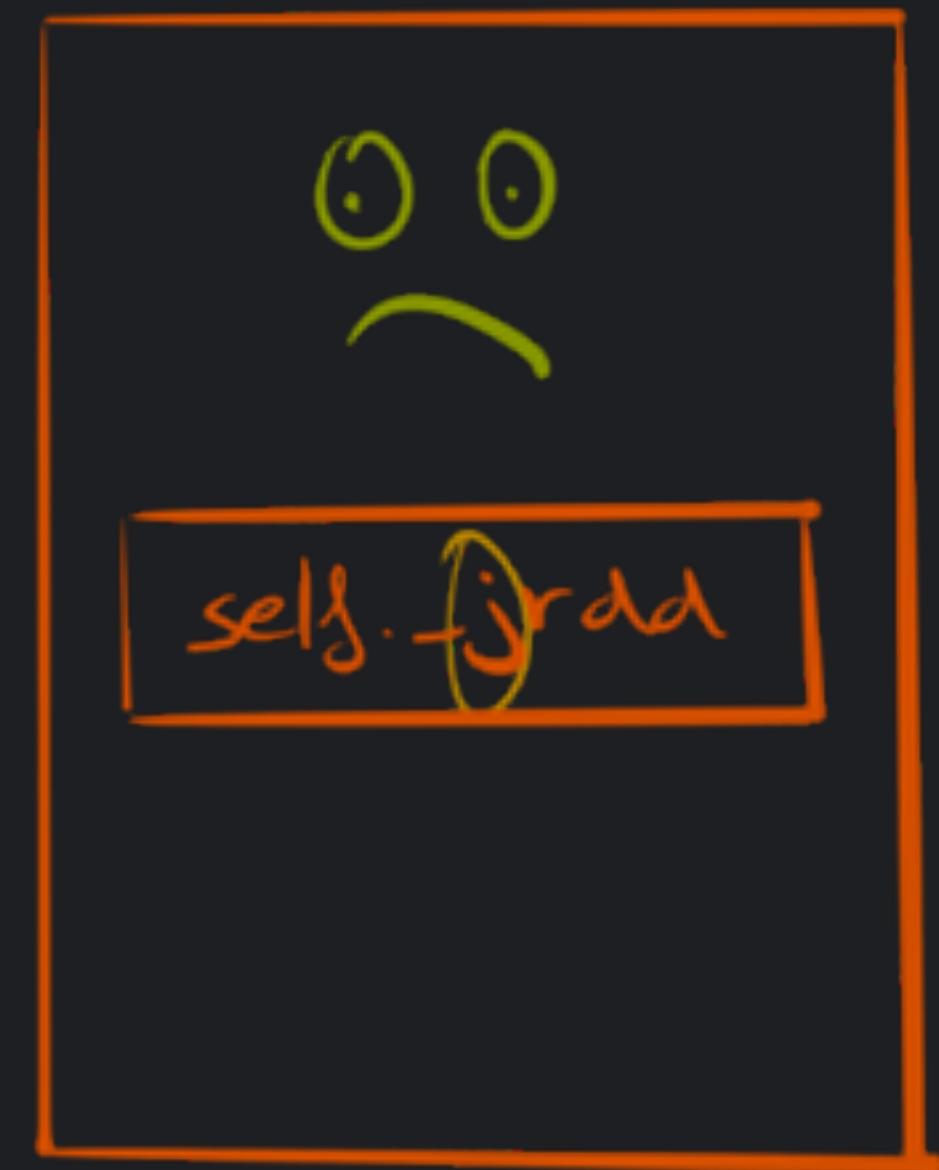


RDD in
Python Land

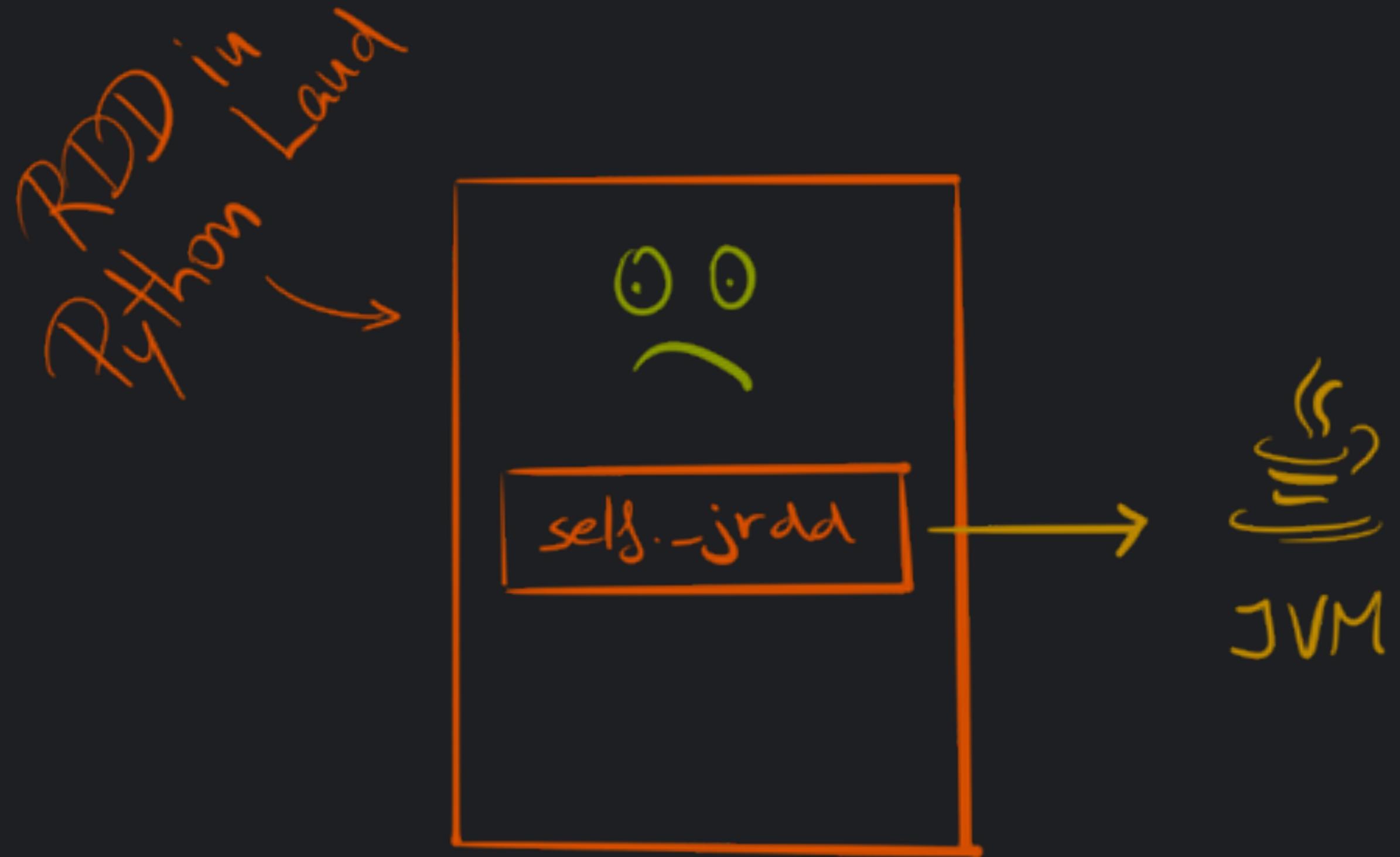


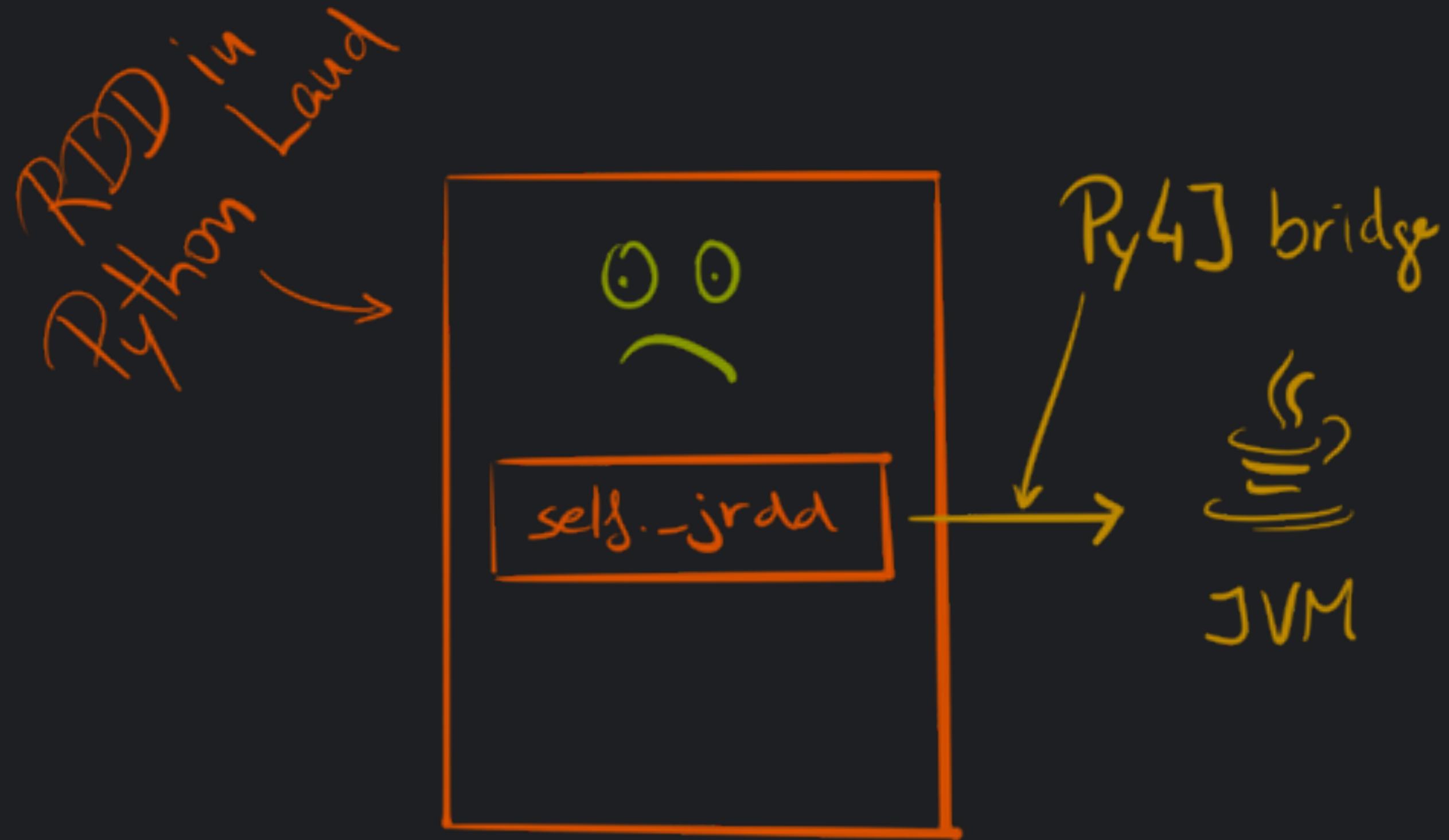
M
is for
murder

RDD in
Python Land



j is for
java





THE MAIN ENTRYPOINTS
ARE RDD AND
PipelinedRDD(RDD)

PipelinedRDD
BUILDS IN THE JVM A
PythonRDD





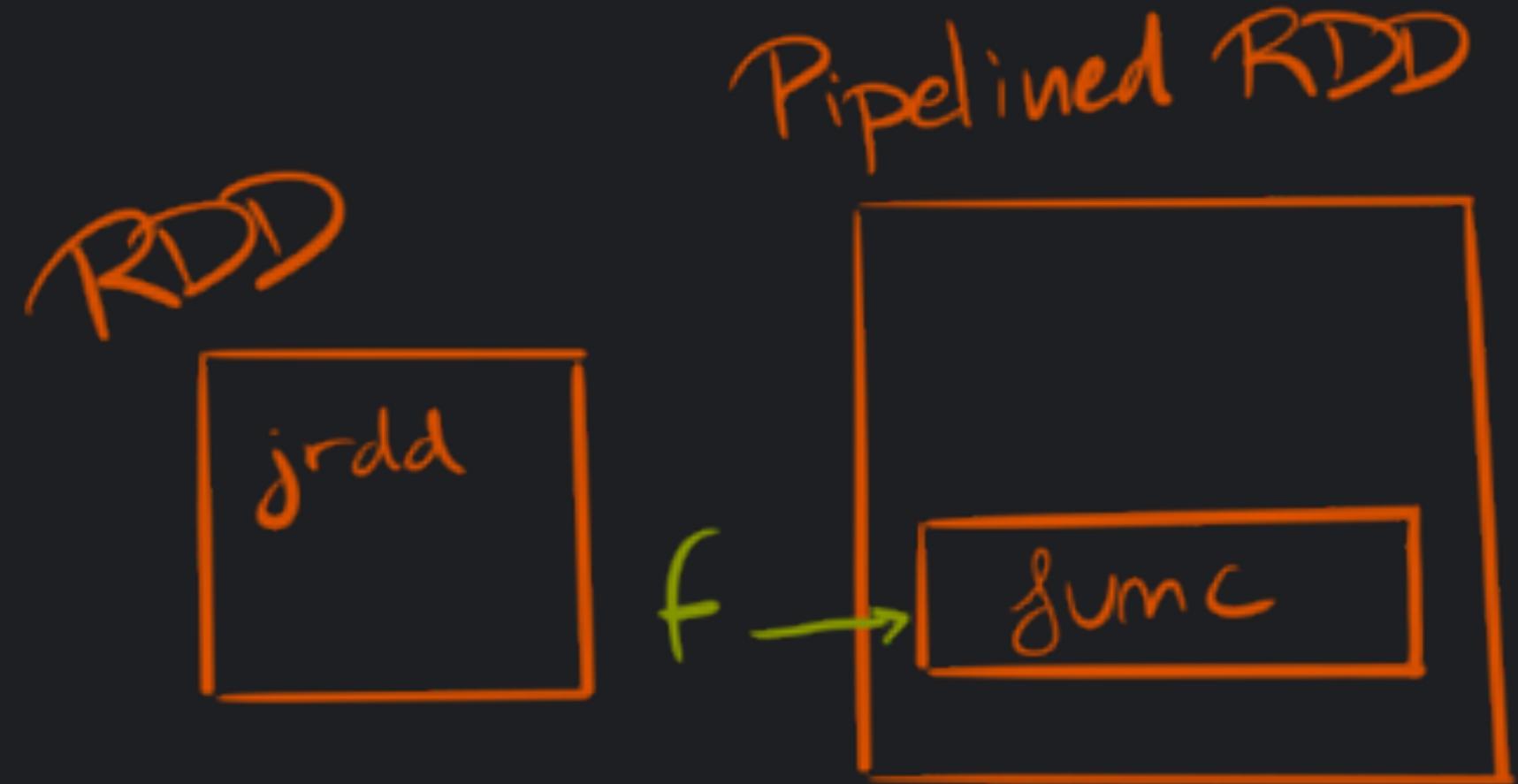


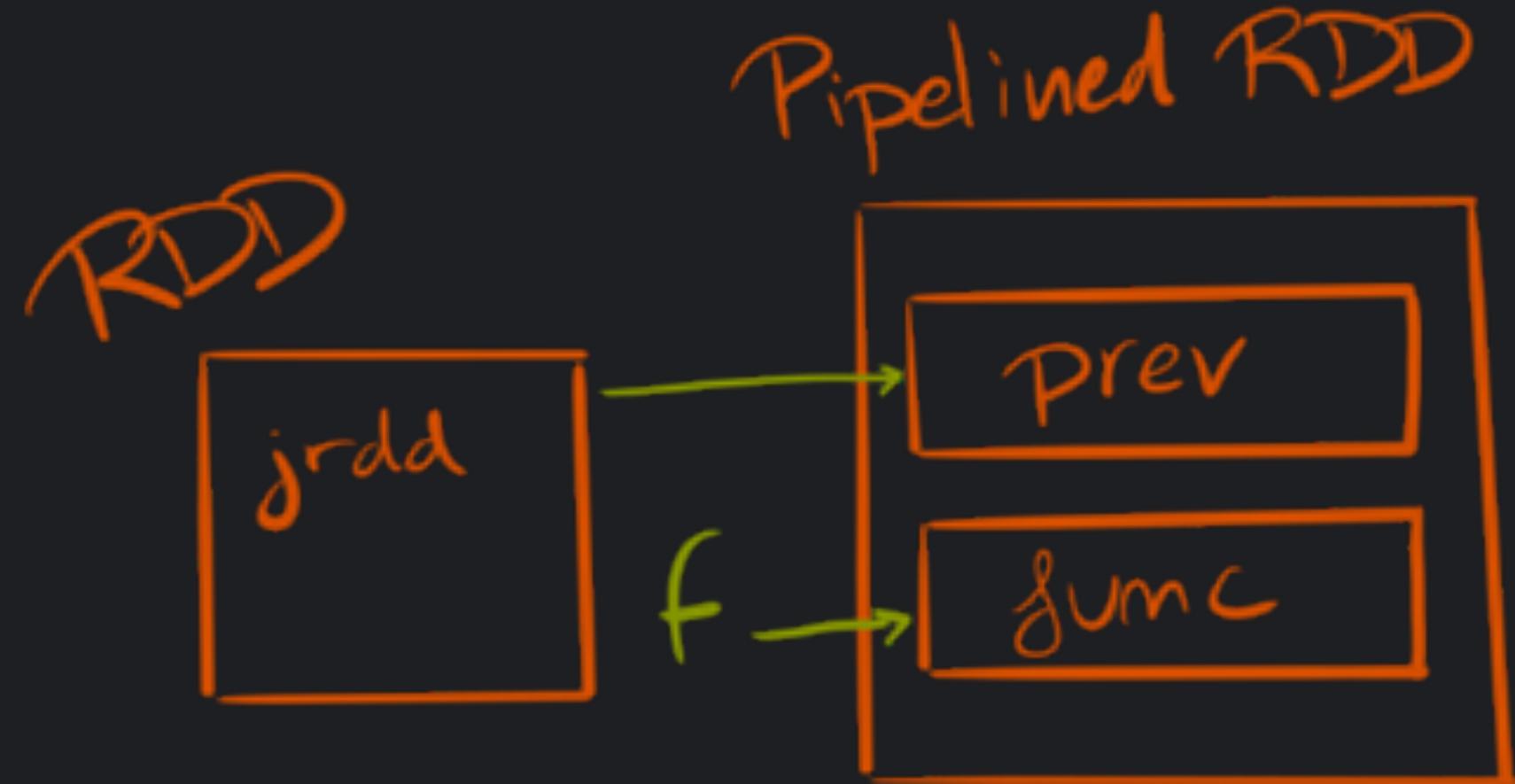
RDD

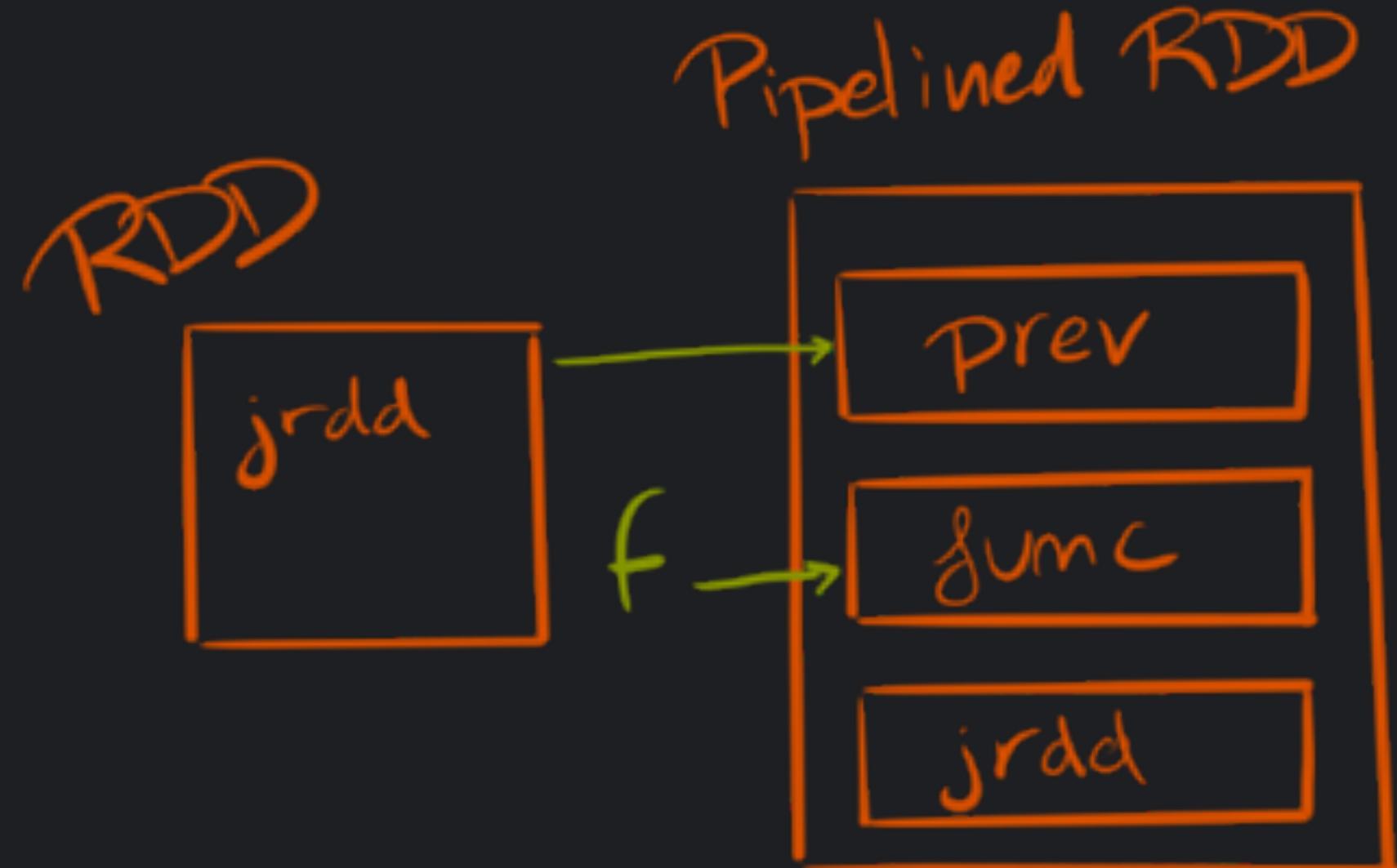
jRDD

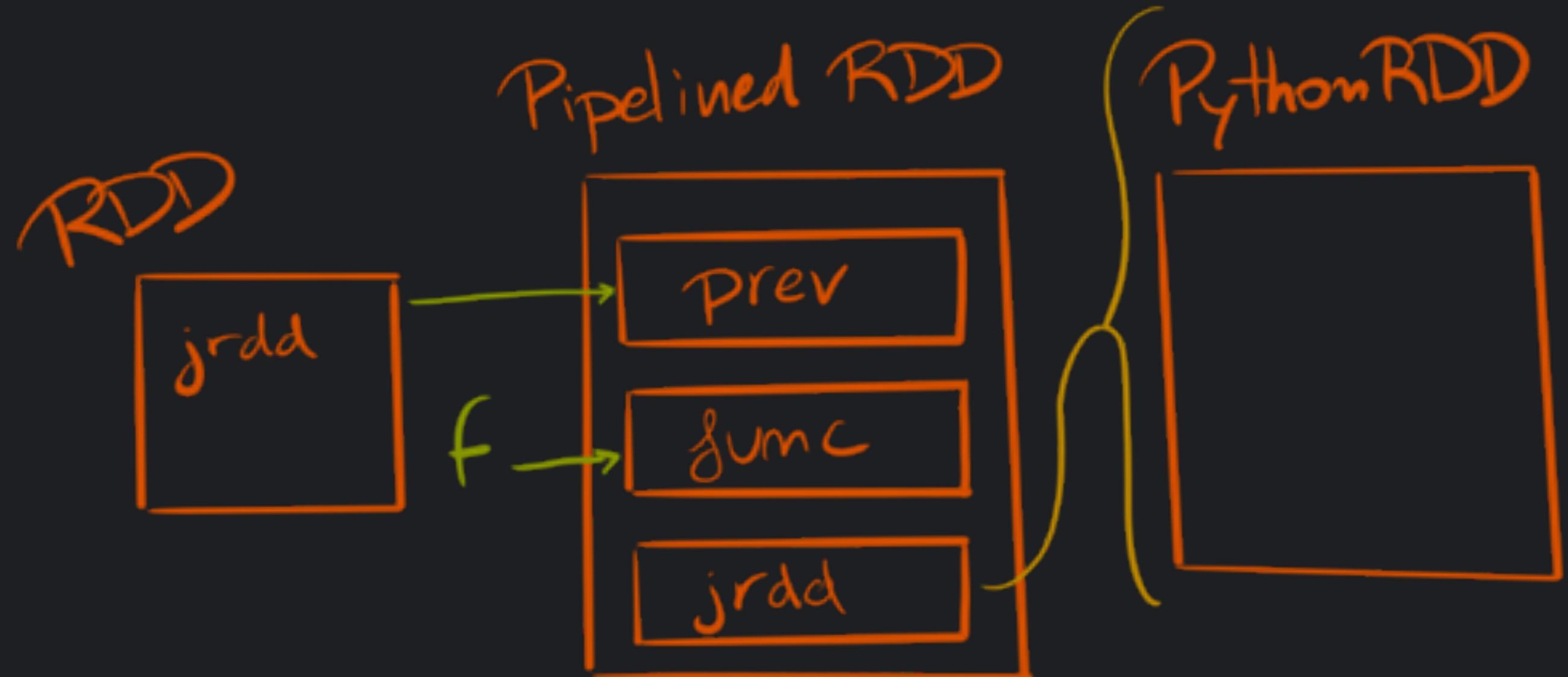
MAP(f)





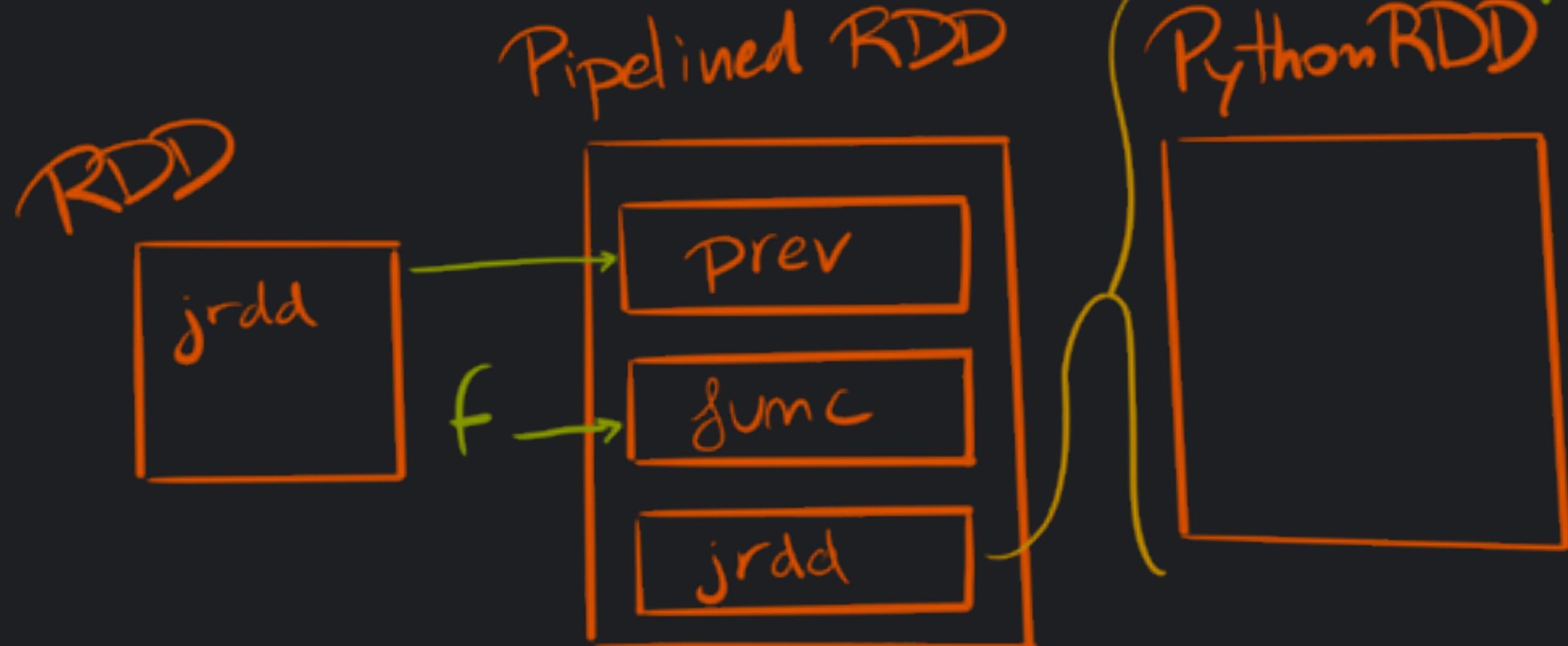






class in Scala

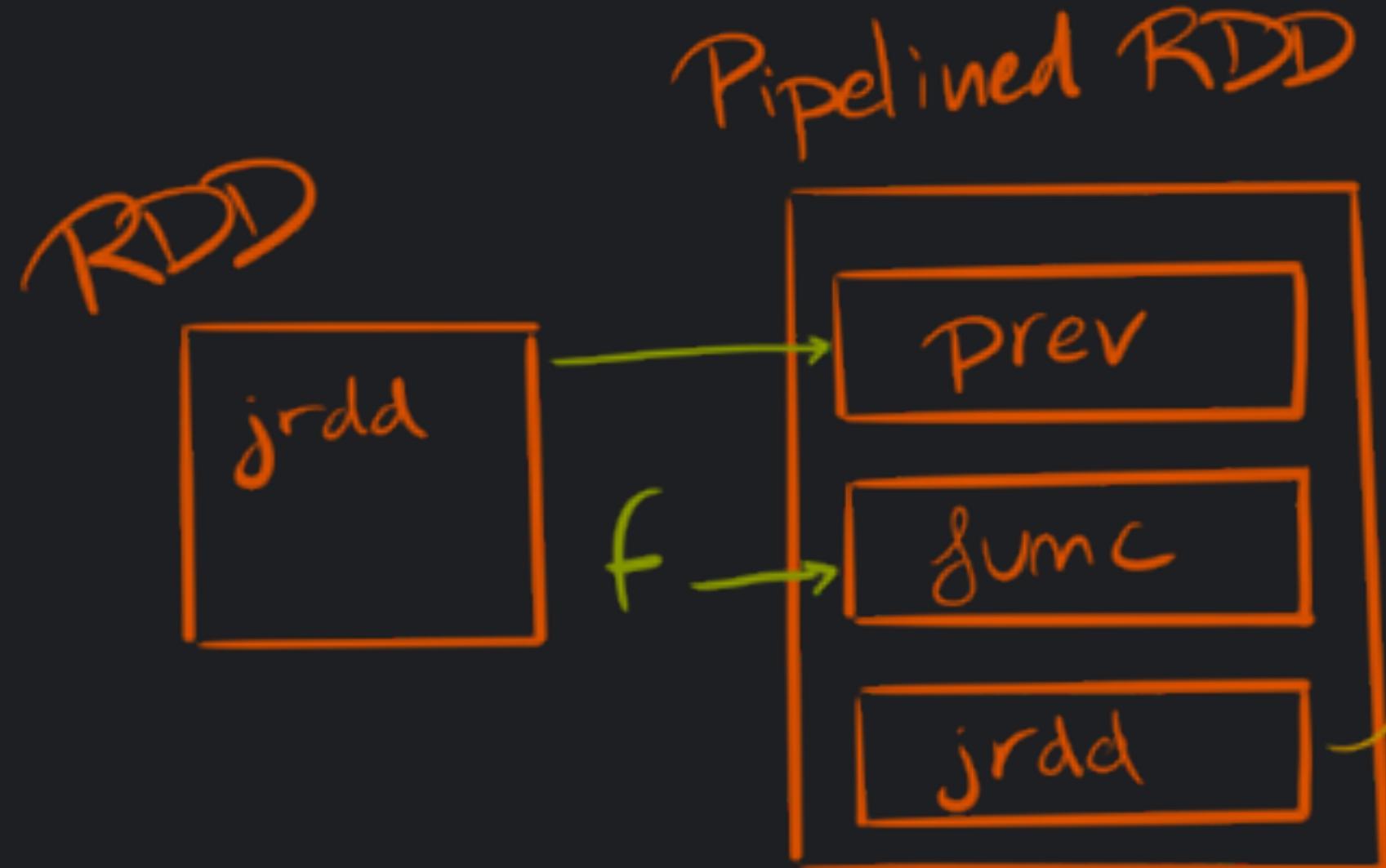
PythonRDD



class in Scala

PythonRDD

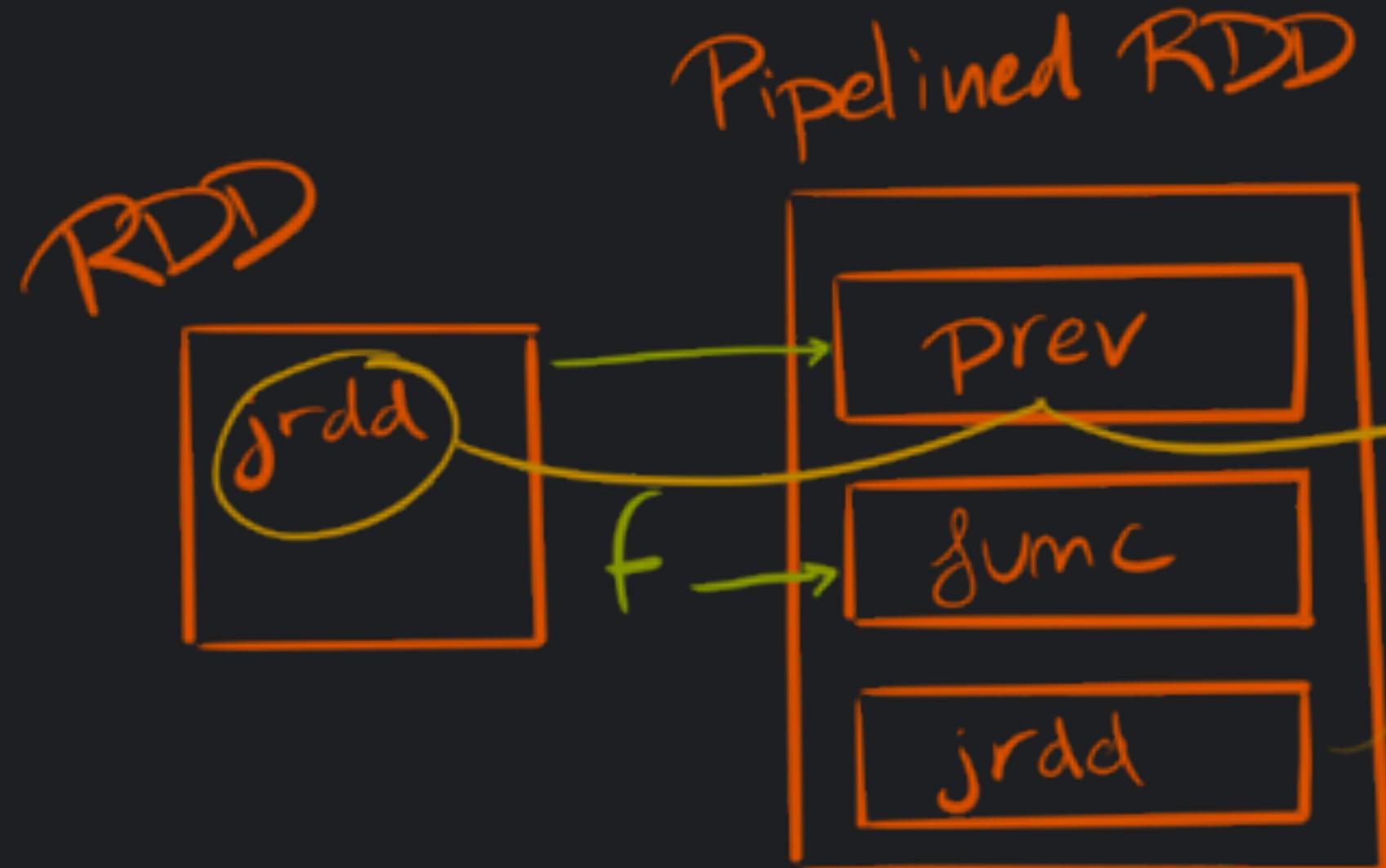
Dependencies



class in Scala

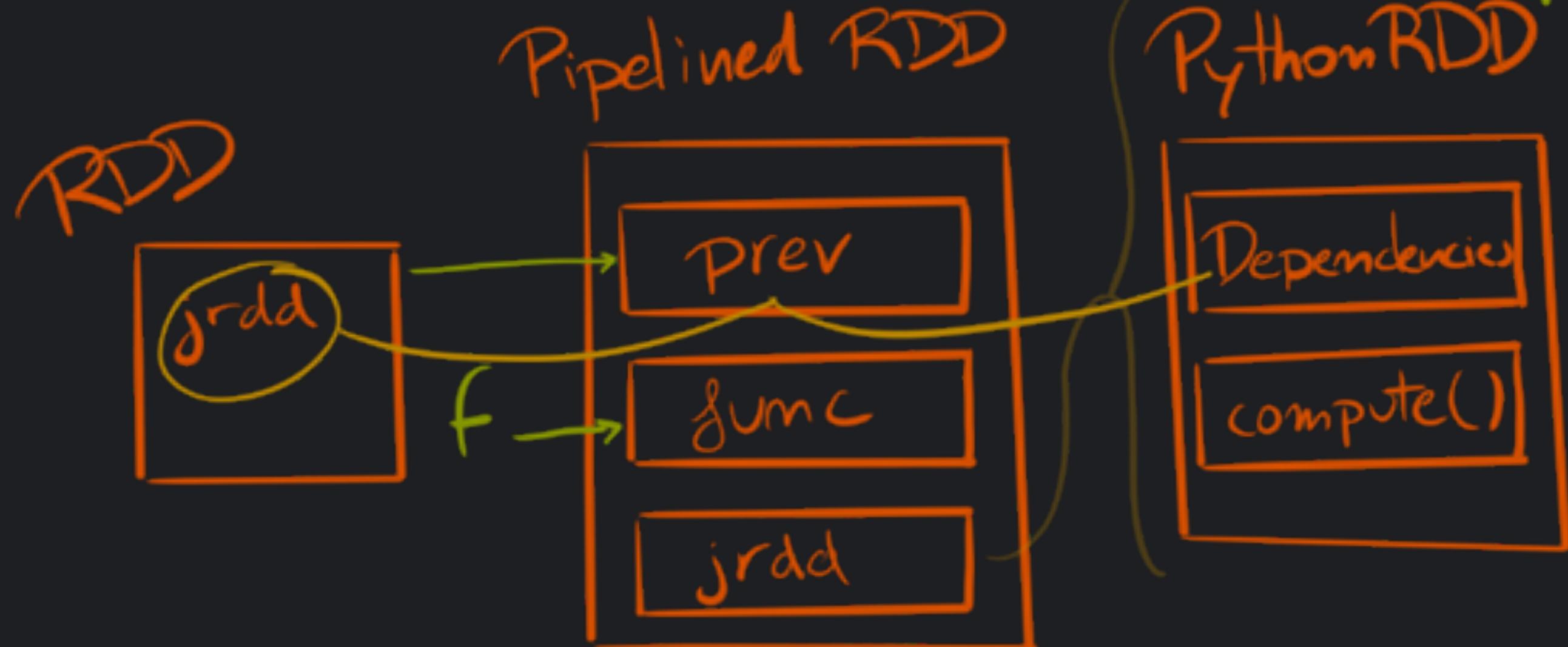
PythonRDD

Dependencies



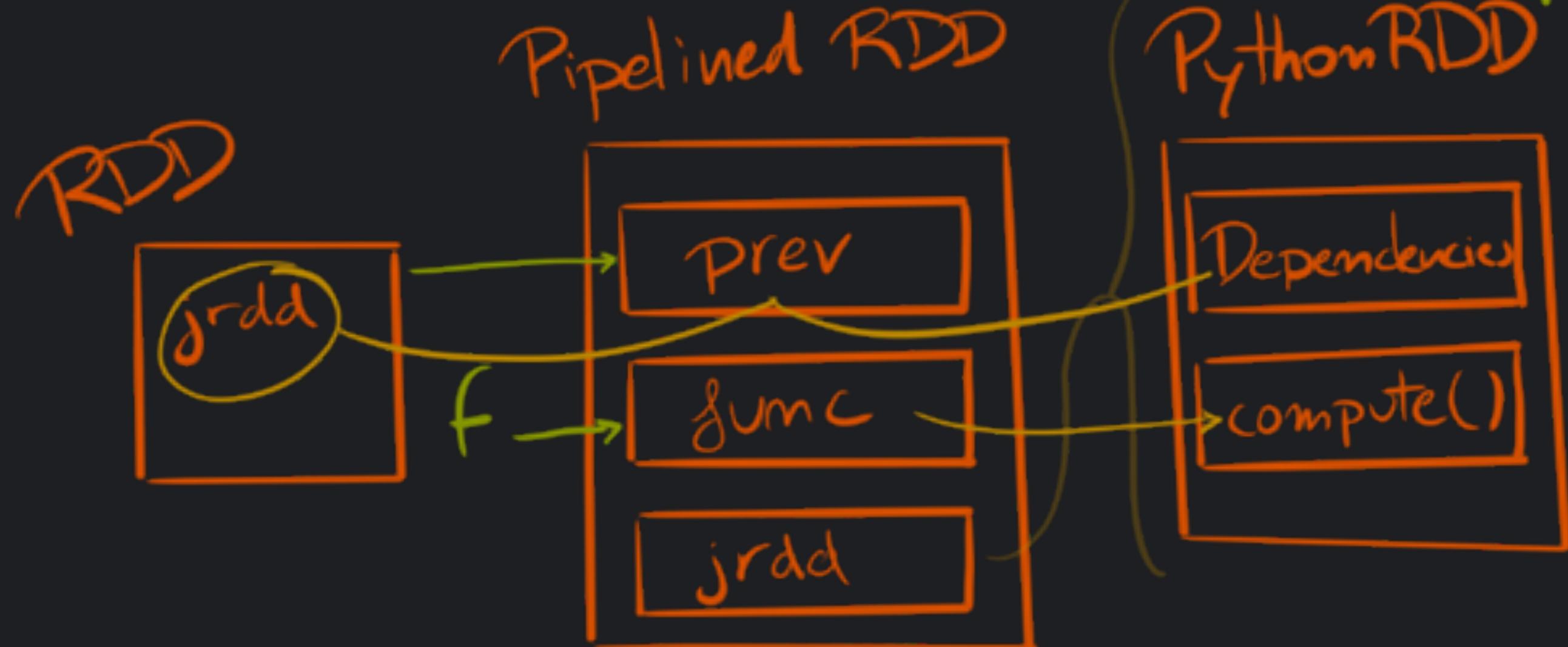
class in Scala

PythonRDD



class in Scala

PythonRDD





THE MAGIC IS
IN
compute

compute
IS RUN ON EACH
EXECUTOR AND STARTS
A PYTHON WORKER VIA
PythonRunner



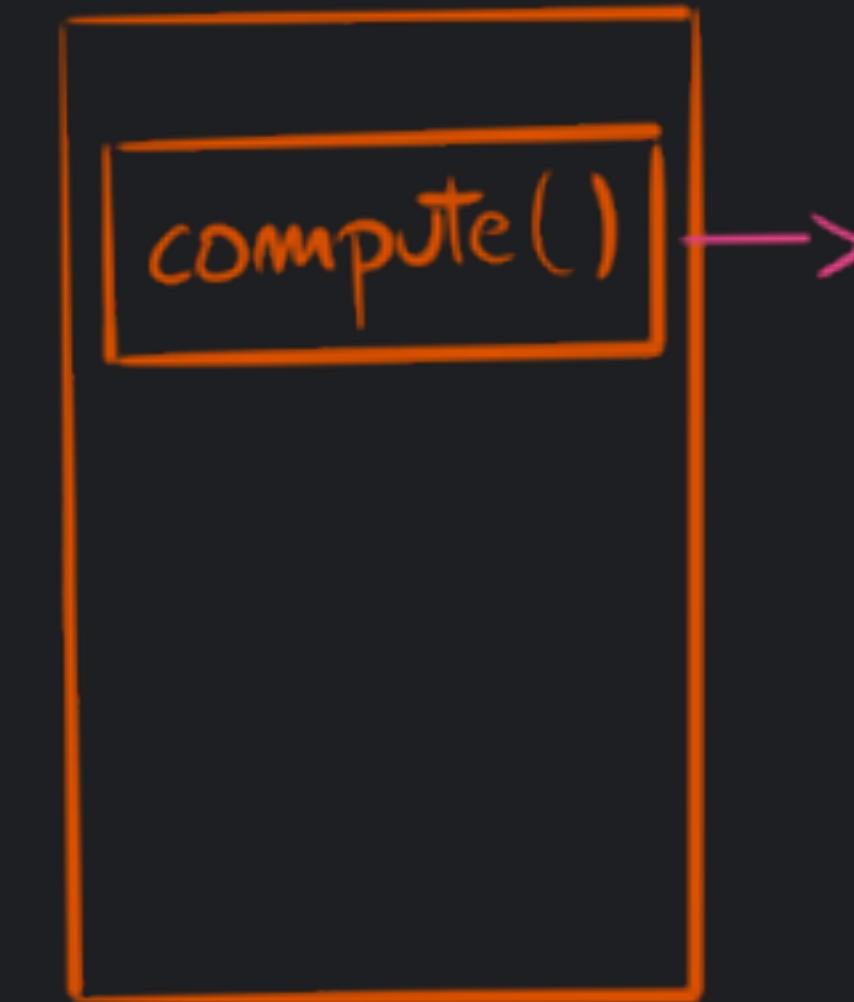


PythonRDD

compute()

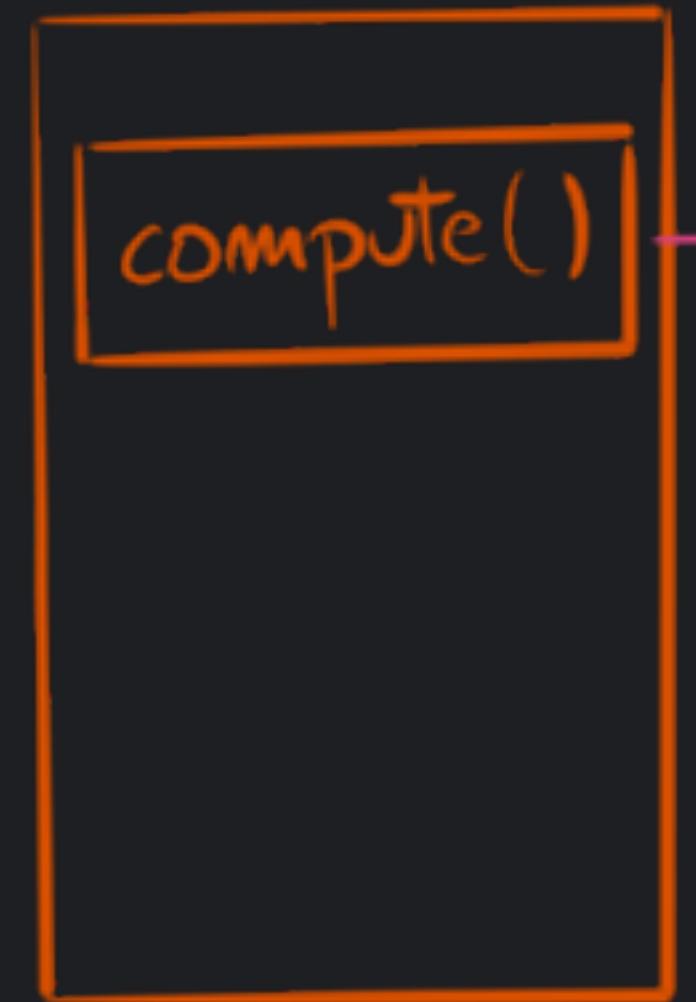
Scala!

PythonRDD



Scala!

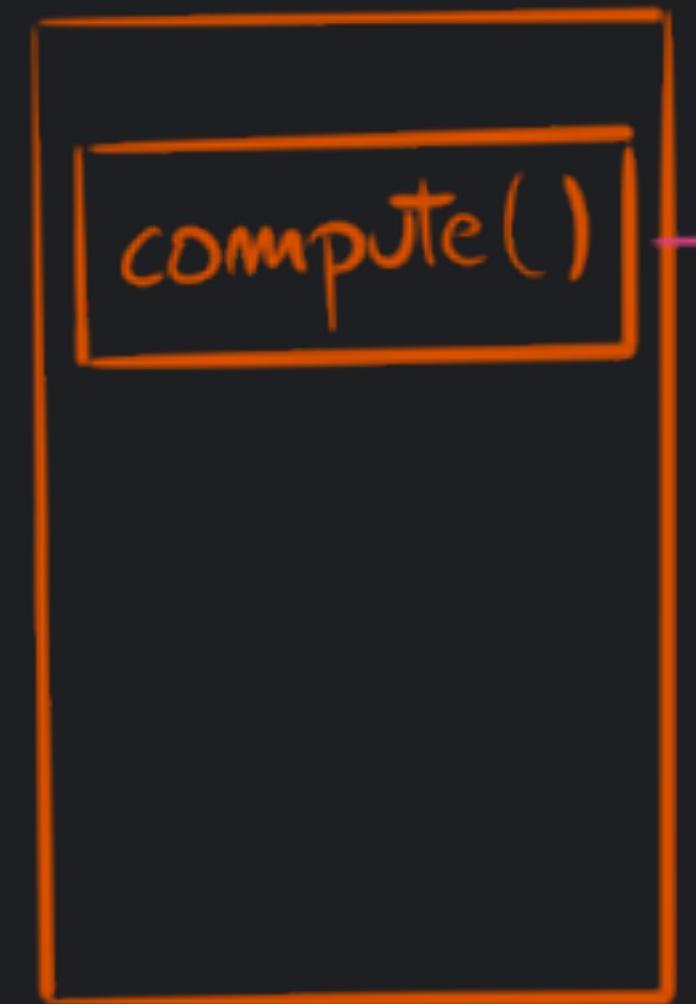
PythonRDD



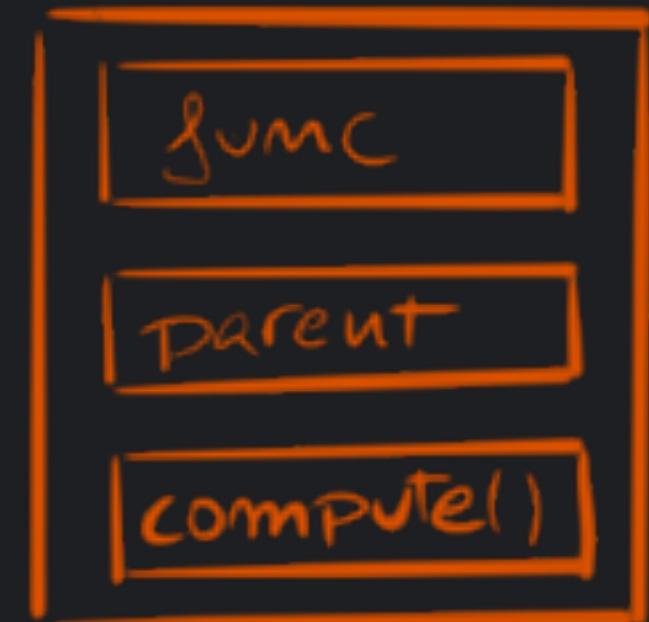
Scala!

← Python!

PythonRDD



PythonRunner



Scala!

← Python!

WORKERS ACT AS STANDALONE PROCESSORS OF STREAMS OF DATA

- > CONNECTS BACK TO THE JVM THAT STARTED IT
 - > LOAD INCLUDED PYTHON LIBRARIES
- > DESERIALIZES THE PICKLED FUNCTION COMING FROM THE STREAM
- > APPLIES THE FUNCTION TO THE DATA COMING FROM THE STREAM
 - > SENDS THE OUTPUT BACK



#UnifiedDataAnalytics #SparkAISummit

BUT . . . WASN'T SPARK
MAGICALLY OPTIMISING
EVERYTHING?

YES. FOR SPARK
Dataframe



SPARK WILL GENERATE
A PLAN
(A DIRECTED ACYCLIC GRAPH)
TO COMPUTE THE
RESULT

AND THE PLAN WILL BE
OPTIMISED USING
CATALYST



DEPENDING ON THE FUNCTION, THE
OPTIMISER WILL CHOOSE

PythonUDFRunner

OR

PythonArrowRunner

(BOTH EXTEND PythonRunner)



UDF RUNNER

SERIALIZED
BATCH

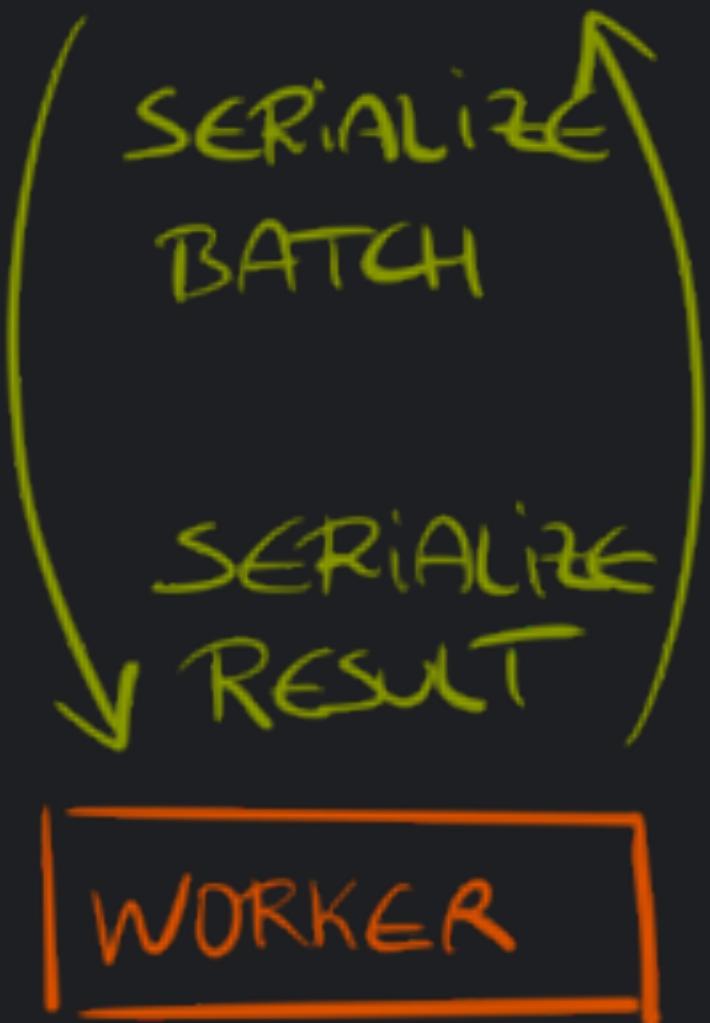


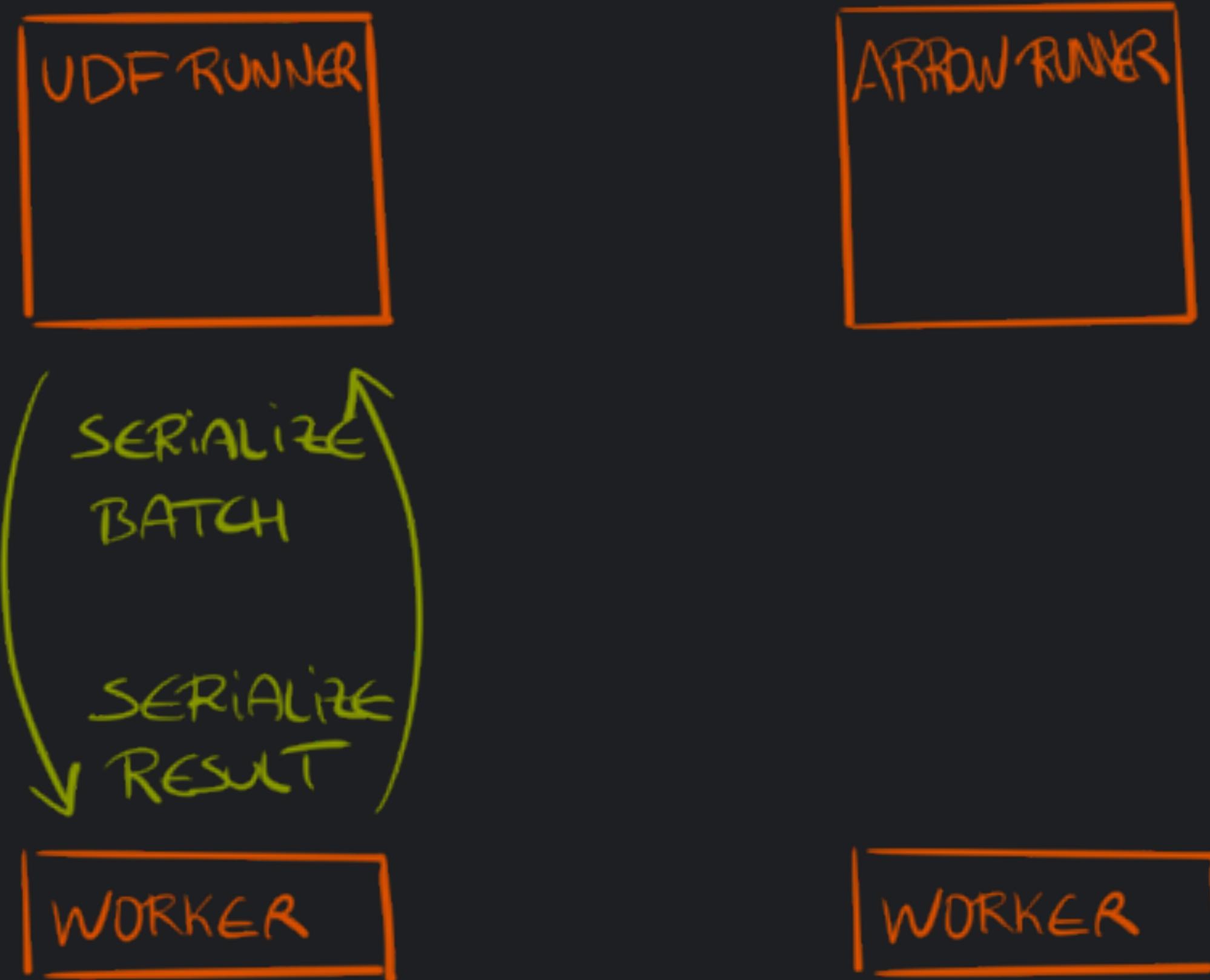
WORKER

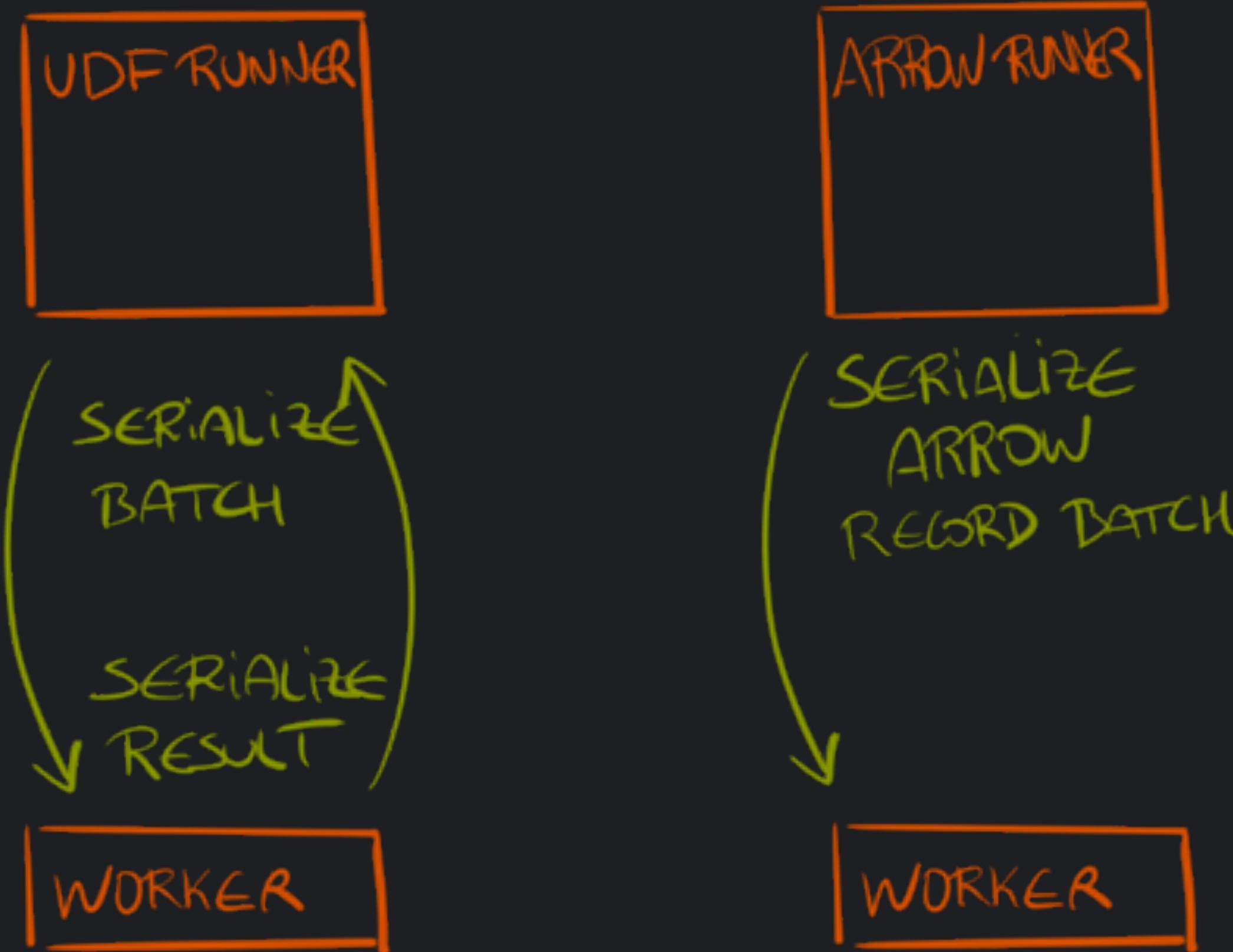


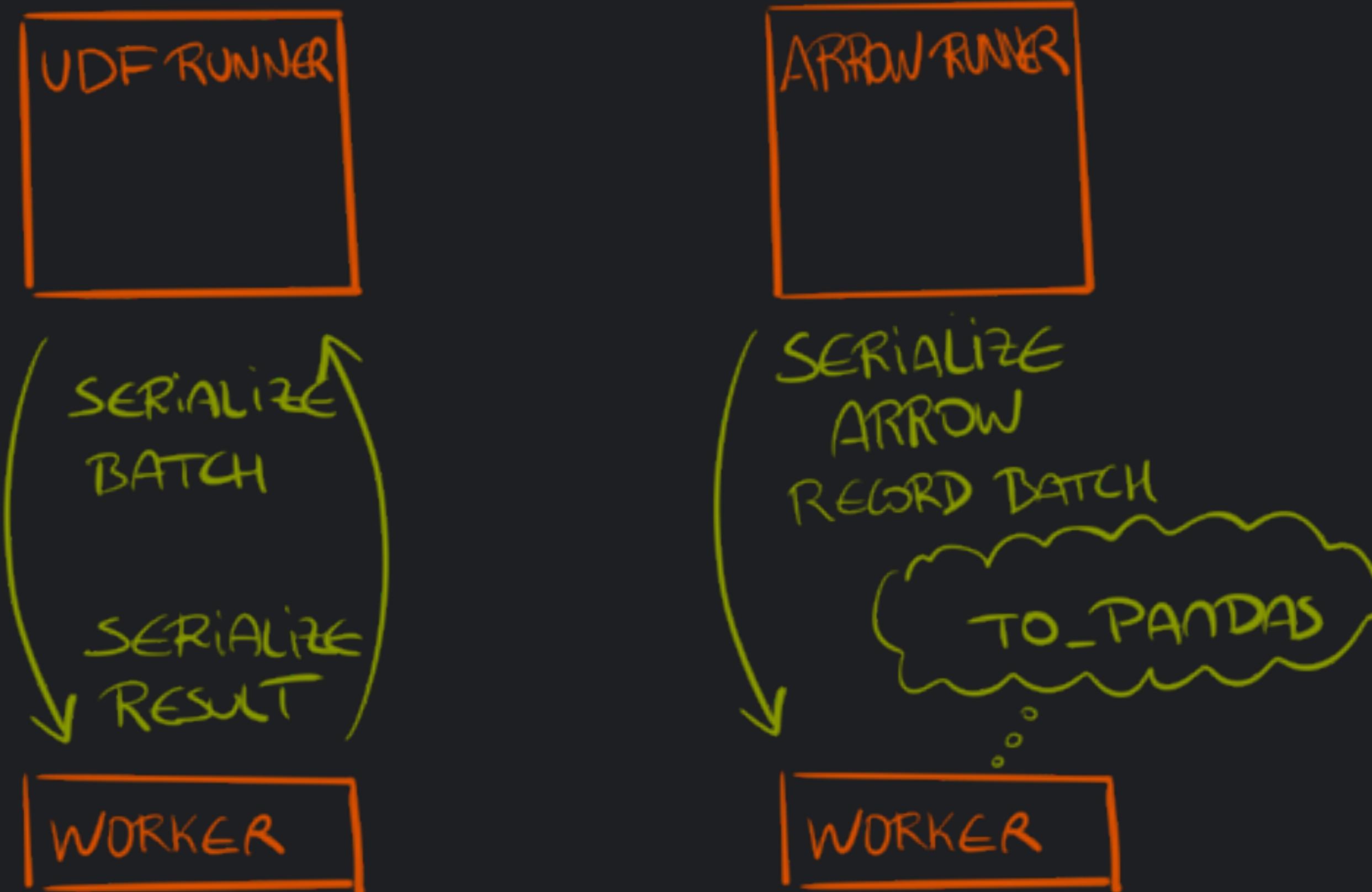


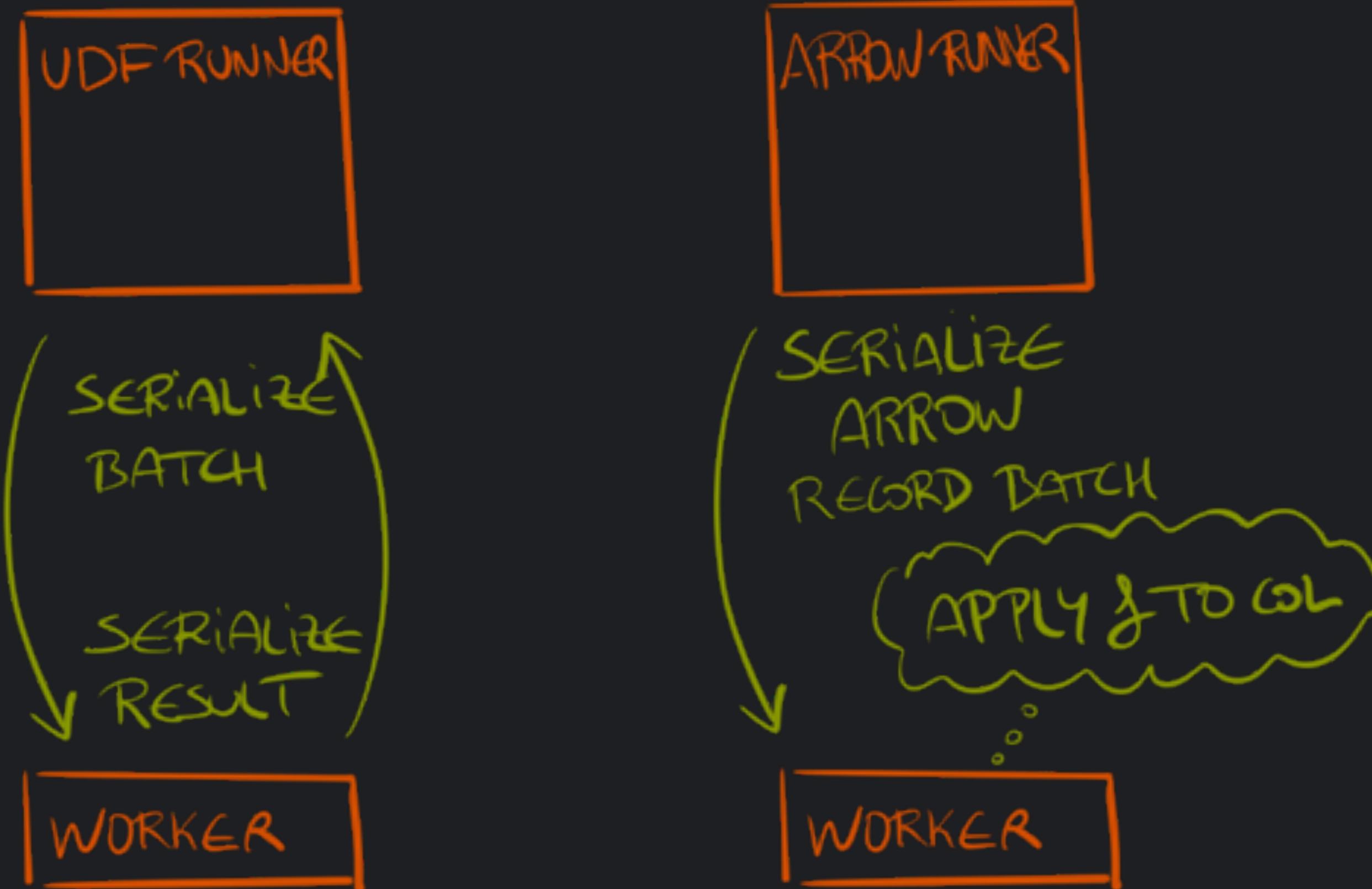
UDF RUNNER

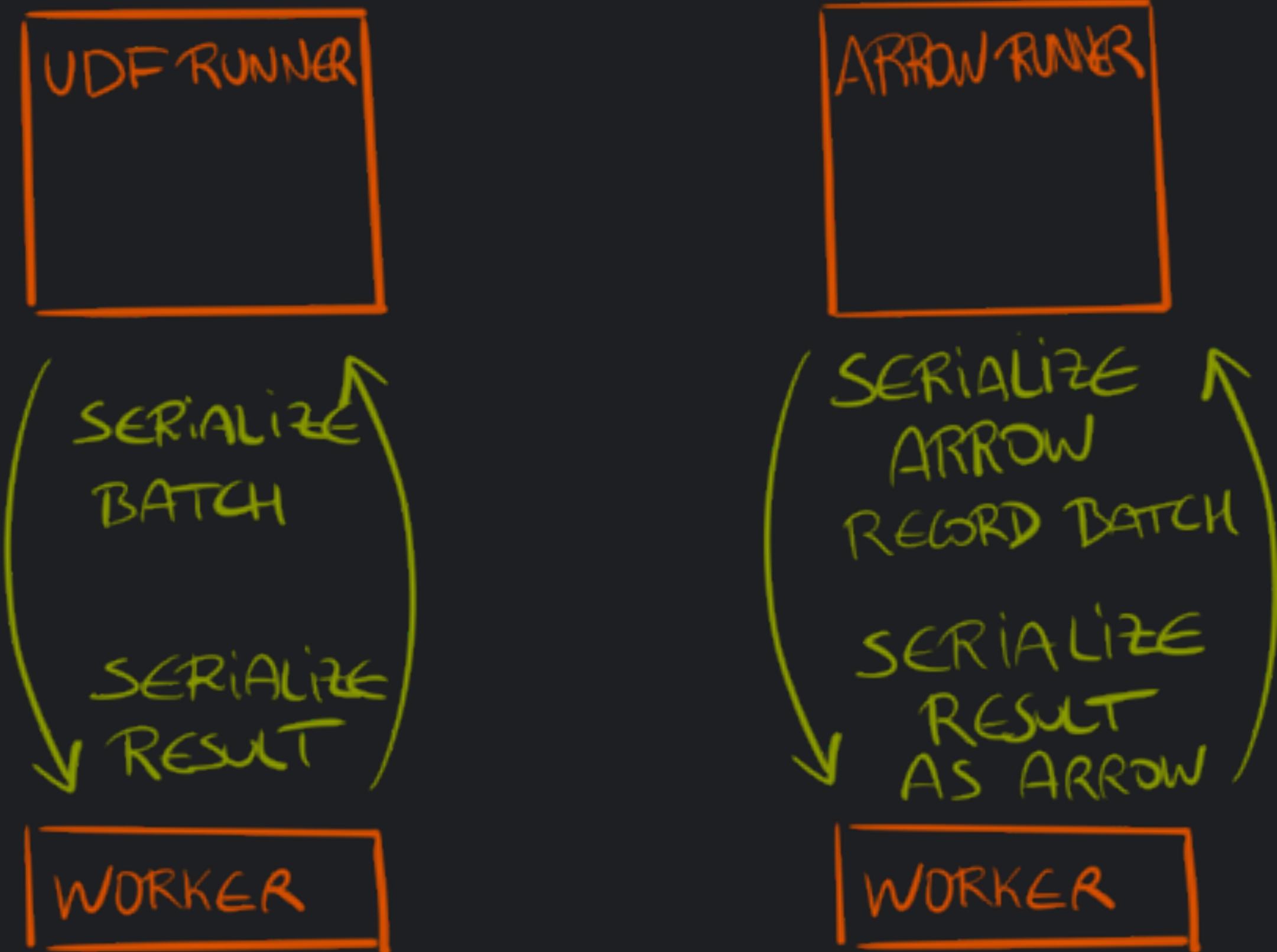












IF WE CAN DEFINE OUR FUNCTIONS
USING PANDAS Series
TRANSFORMATIONS WE CAN SPEED UP
PYSPARK CODE FROM 3X TO 100X!

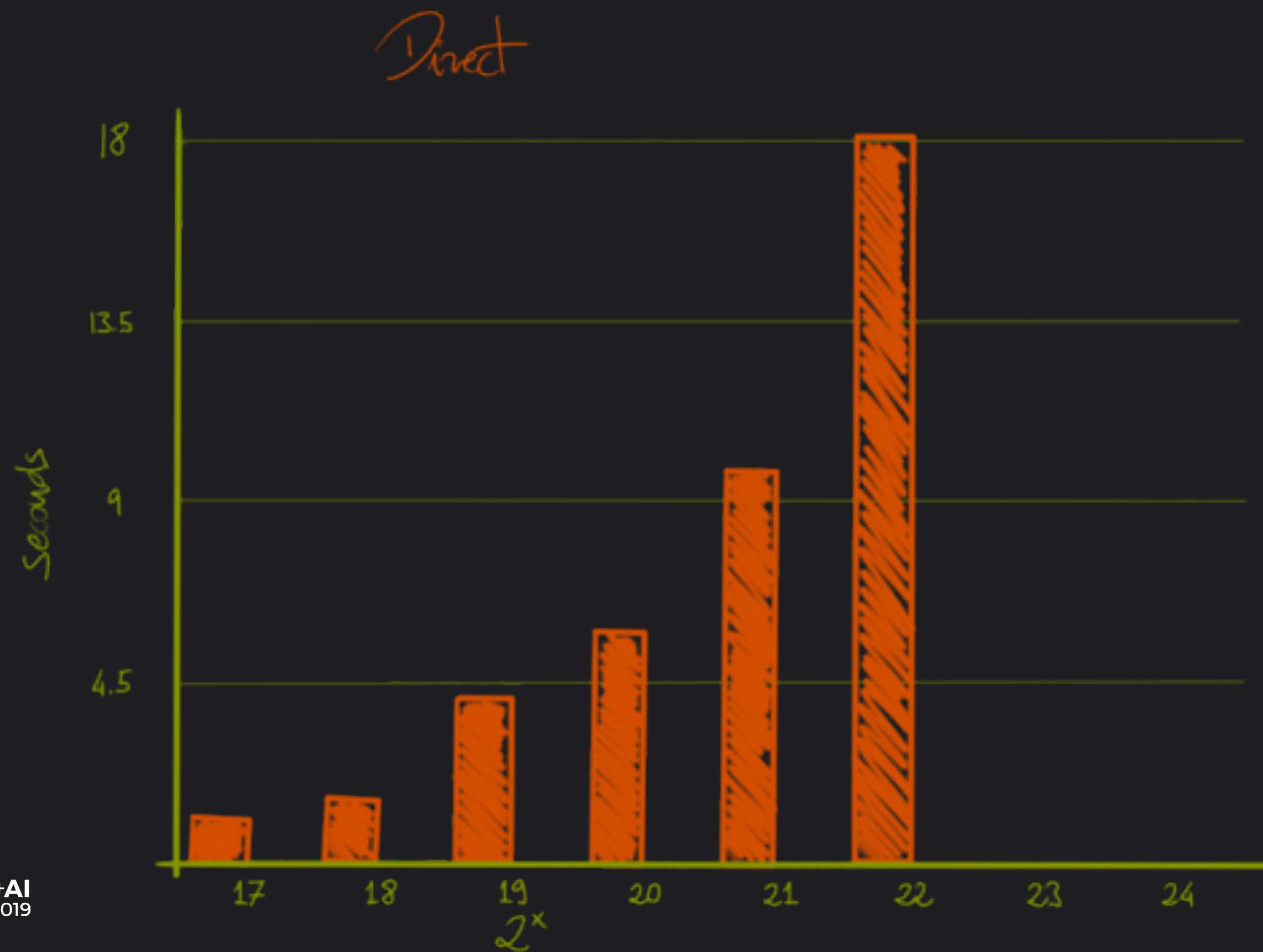
QUICK EXAMPLES

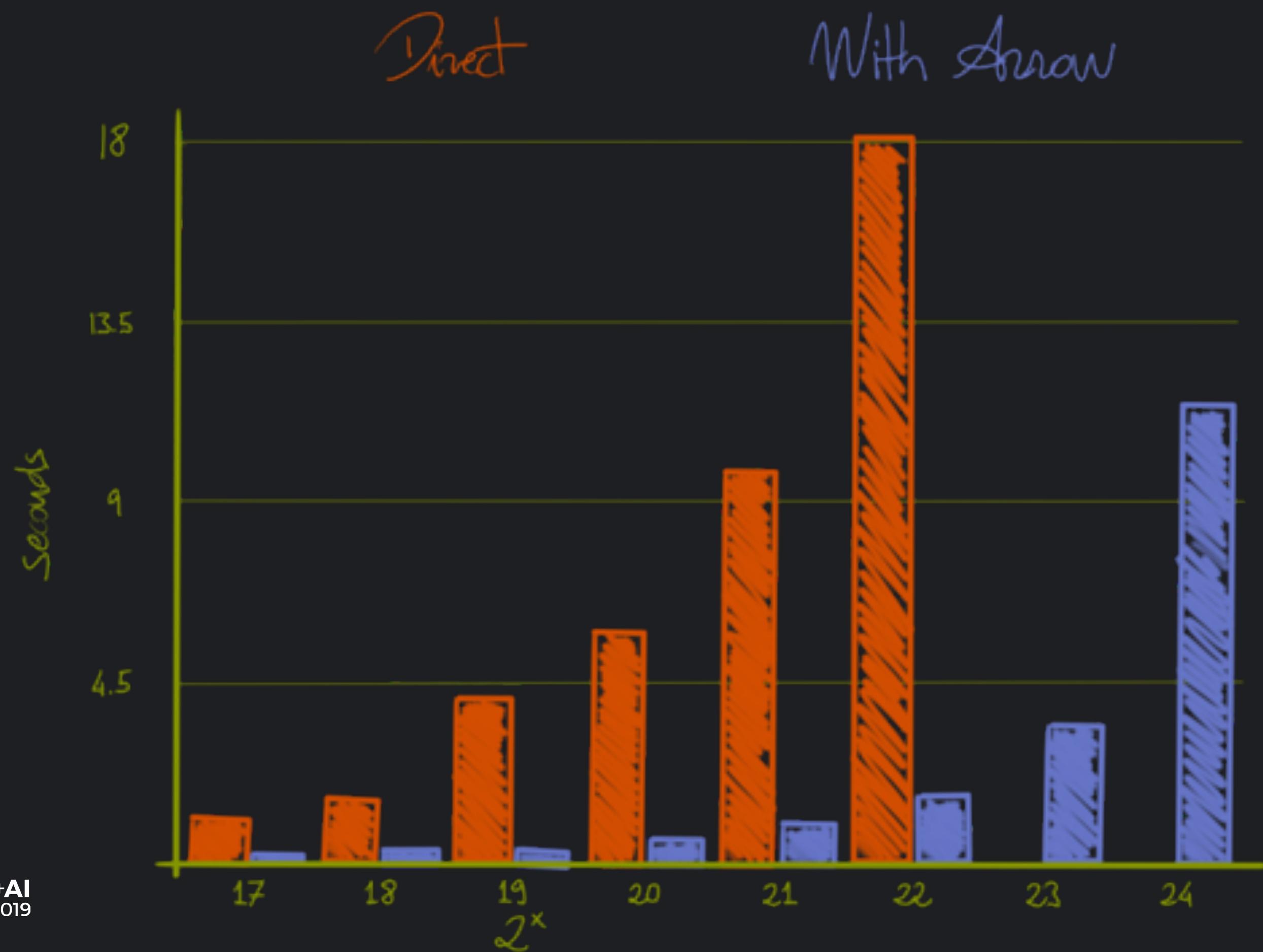
THE BASICS: toPandas

```
from pyspark.sql.functions import rand  
  
df = spark.range(1 << 20).toDF("id").withColumn("x", rand())  
  
spark.conf.set("spark.sql.execution.arrow.enabled", "false")  
pandas_df = df.toPandas() # we'll time this
```

```
from pyspark.sql.functions import rand  
  
df = spark.range(1 << 20).toDF("id").withColumn("x", rand())  
  
spark.conf.set("spark.sql.execution.arrow.enabled", "false")  
pandas_df = df.toPandas() # we'll time this
```

```
from pyspark.sql.functions import rand  
  
df = spark.range(1 << 20).toDF("id").withColumn("x", rand())  
  
spark.conf.set("spark.sql.execution.arrow.enabled", "true")  
pandas_df = df.toPandas() # we'll time this
```





THE FUN: .groupBy

```
from pyspark.sql.functions import rand, randn, floor
from pyspark.sql.functions import pandas_udf, PandasUDFType

df = spark.range(20000000).toDF("row").drop("row") \
    .withColumn("id", floor(rand()*10000)).withColumn("spent", (randn()+3)*100)

@pandas_udf("id long, spent double", PandasUDFType.GROUPED_MAP)
def subtract_mean(pdf):
    spent = pdf.spent
    return pdf.assign(spent=spent - spent.mean())

df_to_pandas_arrow = df.groupby("id").apply(subtract_mean).toPandas()
```

```
from pyspark.sql.functions import rand, randn, floor
from pyspark.sql.functions import pandas_udf, PandasUDFType

df = spark.range(20000000).toDF("row").drop("row") \
    .withColumn("id", floor(rand()*10000)).withColumn("spent", (randn()+3)*100)

@pandas_udf("id long, spent double", PandasUDFType.GROUPED_MAP)
def subtract_mean(pdf):
    spent = pdf.spent
    return pdf.assign(spent=spent - spent.mean())

df_to_pandas_arrow = df.groupby("id").apply(subtract_mean).toPandas()
```

Executors



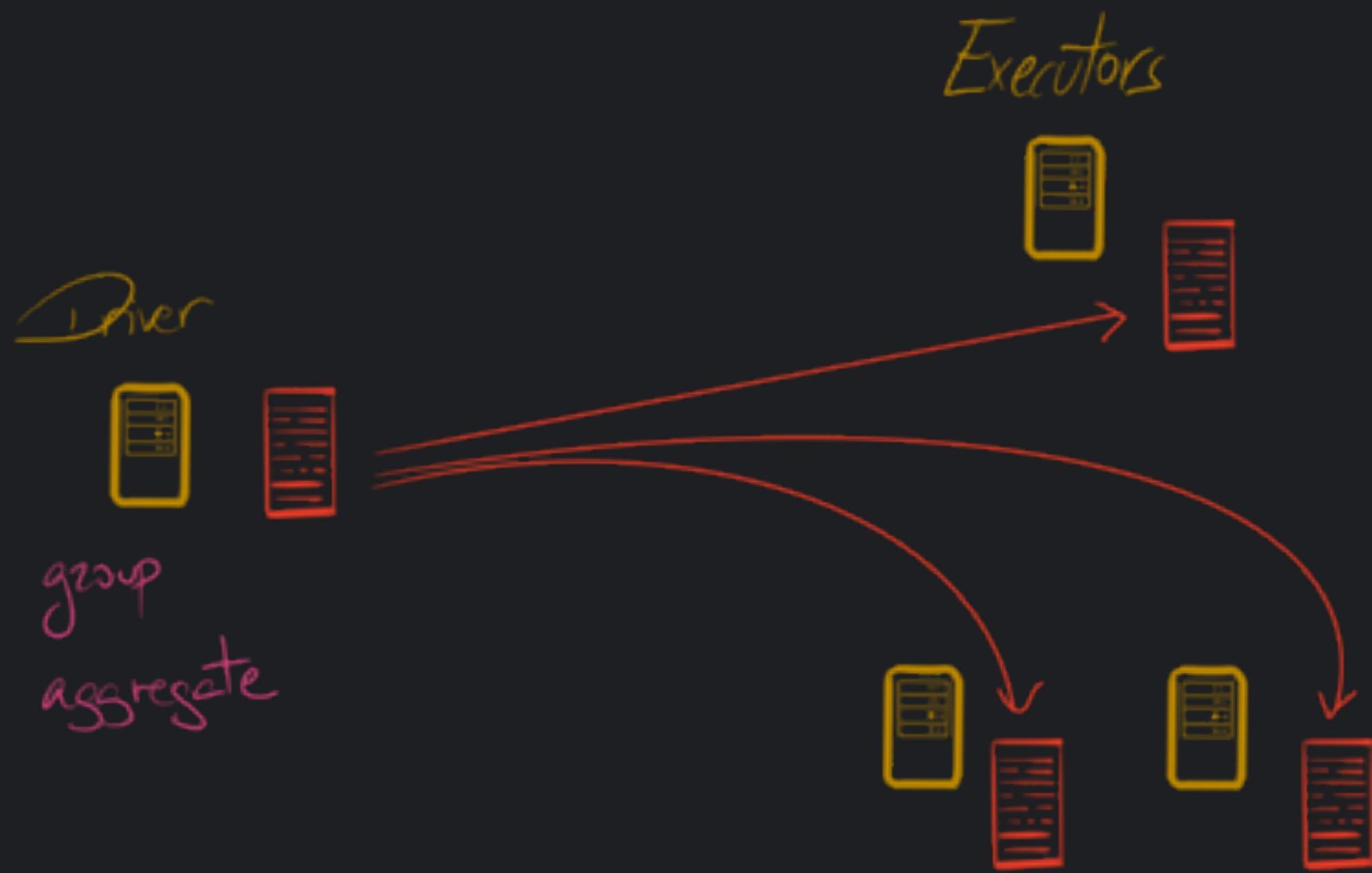
Driver



group

aggregate





Driver

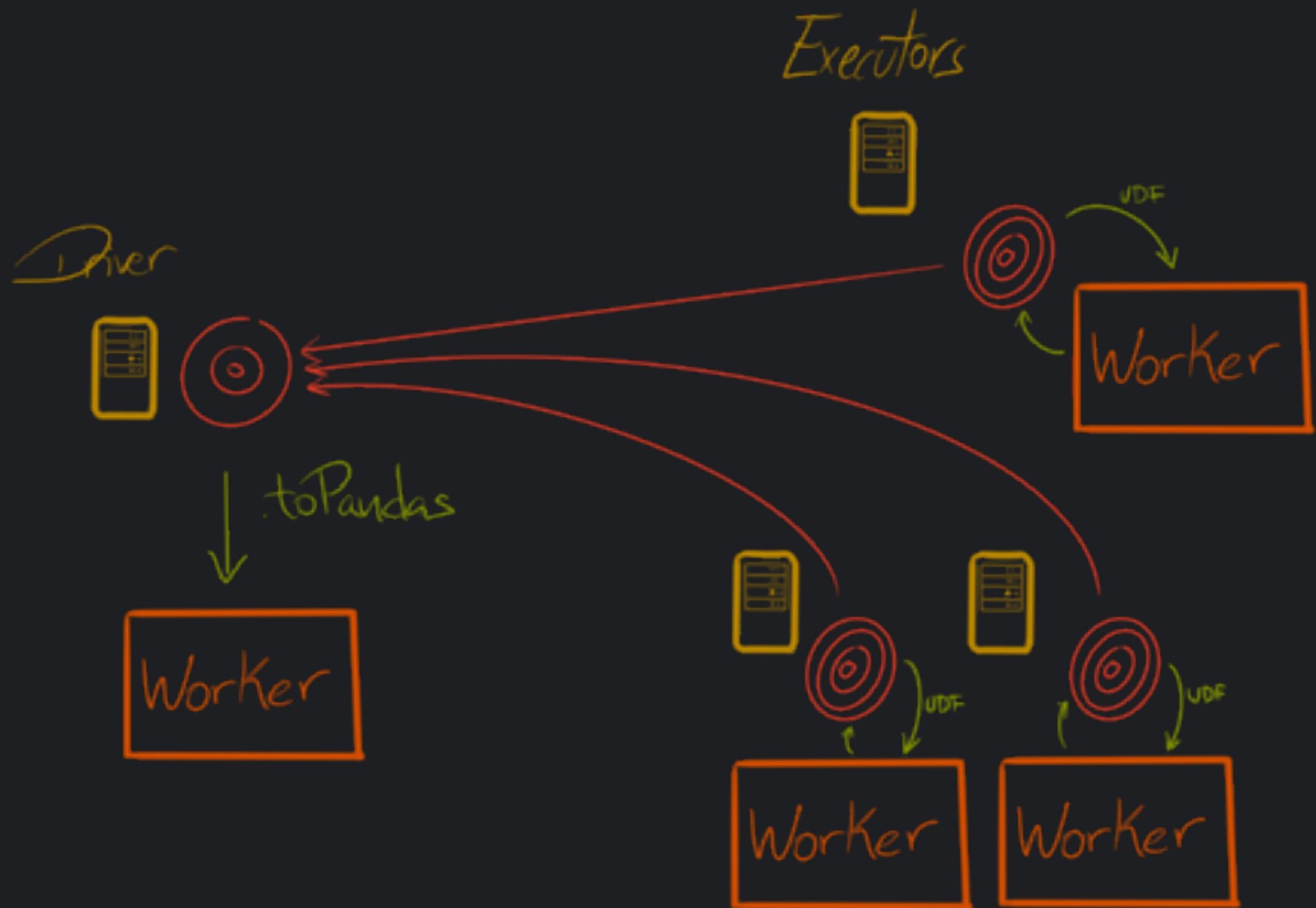


group

aggregate

Executors





BEFORE YOU MAY HAVE DONE SOMETHING LIKE..

```
import numpy as np
from pyspark.sql.functions import collect_list

grouped = df2.groupby("id").agg(collect_list('spent').alias("spent_list"))
as_pandas = grouped.toPandas()
as_pandas["mean"] = as_pandas["spent_list"].apply(np.mean)
as_pandas["subtracted"] = as_pandas["spent_list"].apply(np.array) - as_pandas["mean"]
df_to_pandas = as_pandas.drop(columns=["spent_list", "mean"]).explode("subtracted")
```

```
import numpy as np
from pyspark.sql.functions import collect_list

grouped = df2.groupby("id").agg(collect_list('spent').alias("spent_list"))
as_pandas = grouped.toPandas()
as_pandas["mean"] = as_pandas["spent_list"].apply(np.mean)
as_pandas["subtracted"] = as_pandas["spent_list"].apply(np.array) - as_pandas["mean"]
df_to_pandas = as_pandas.drop(columns=["spent_list", "mean"]).explode("subtracted")
```

Driver



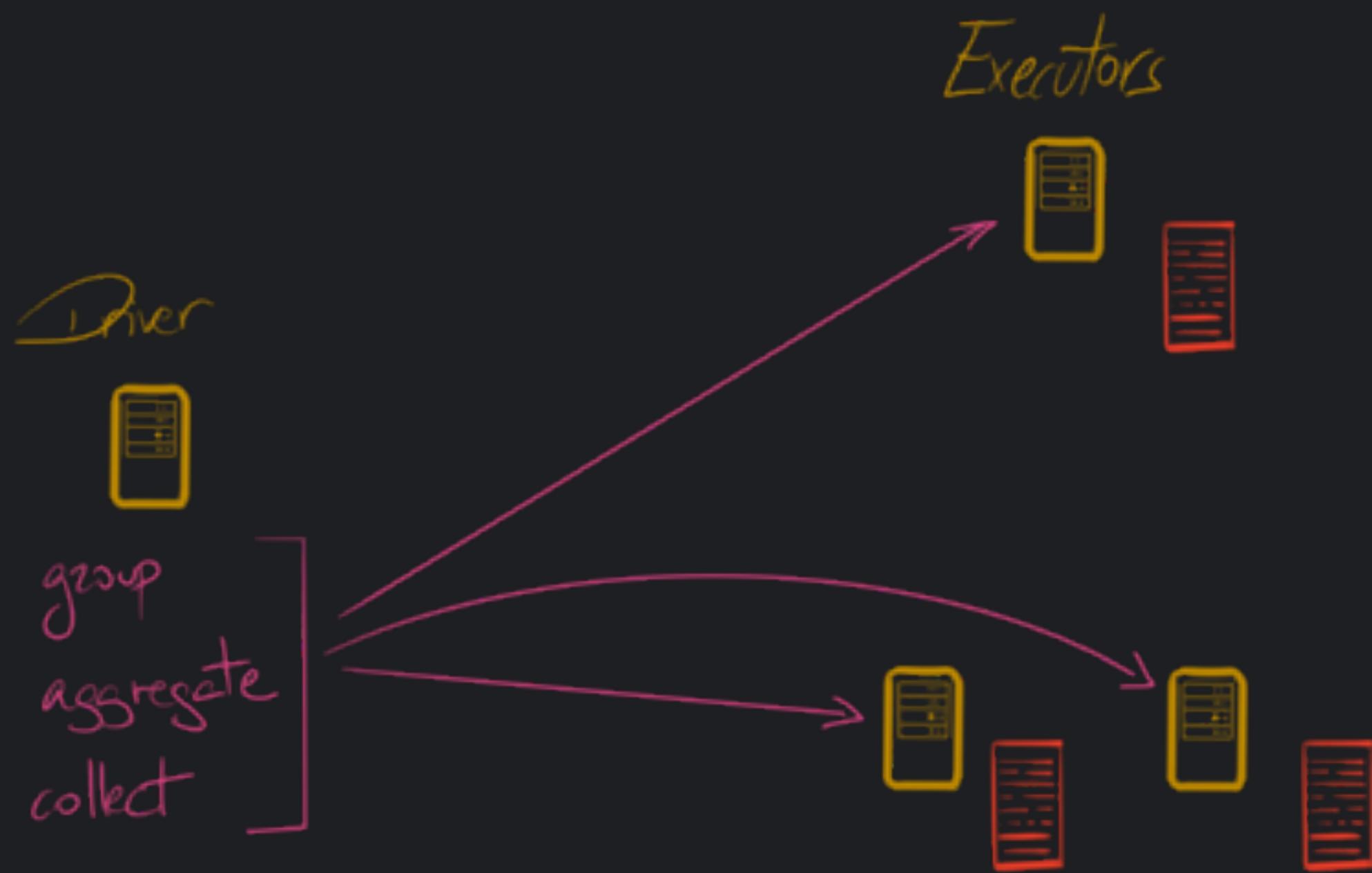
group

aggregate

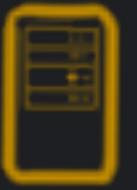
collect

Executors





Driver



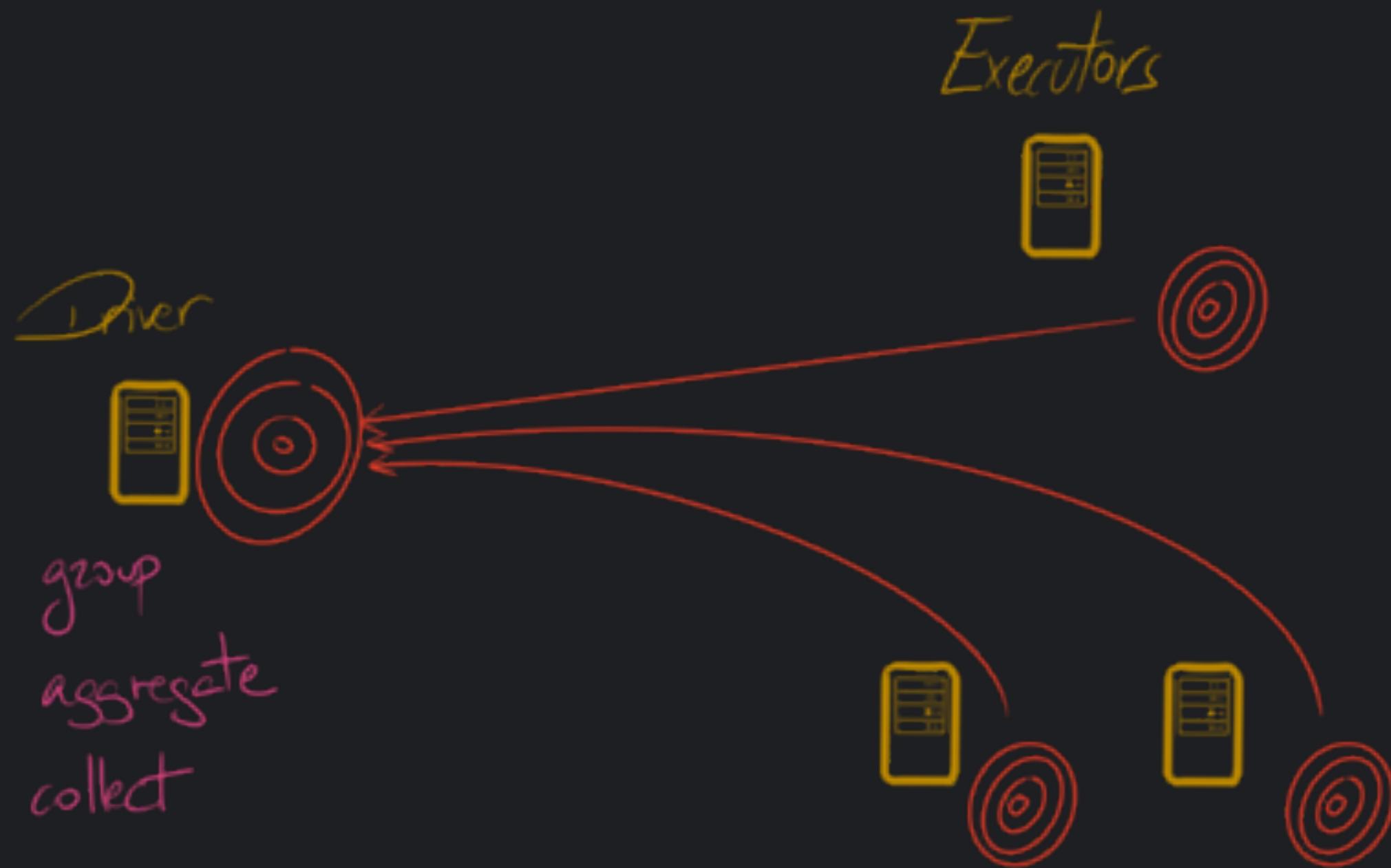
group

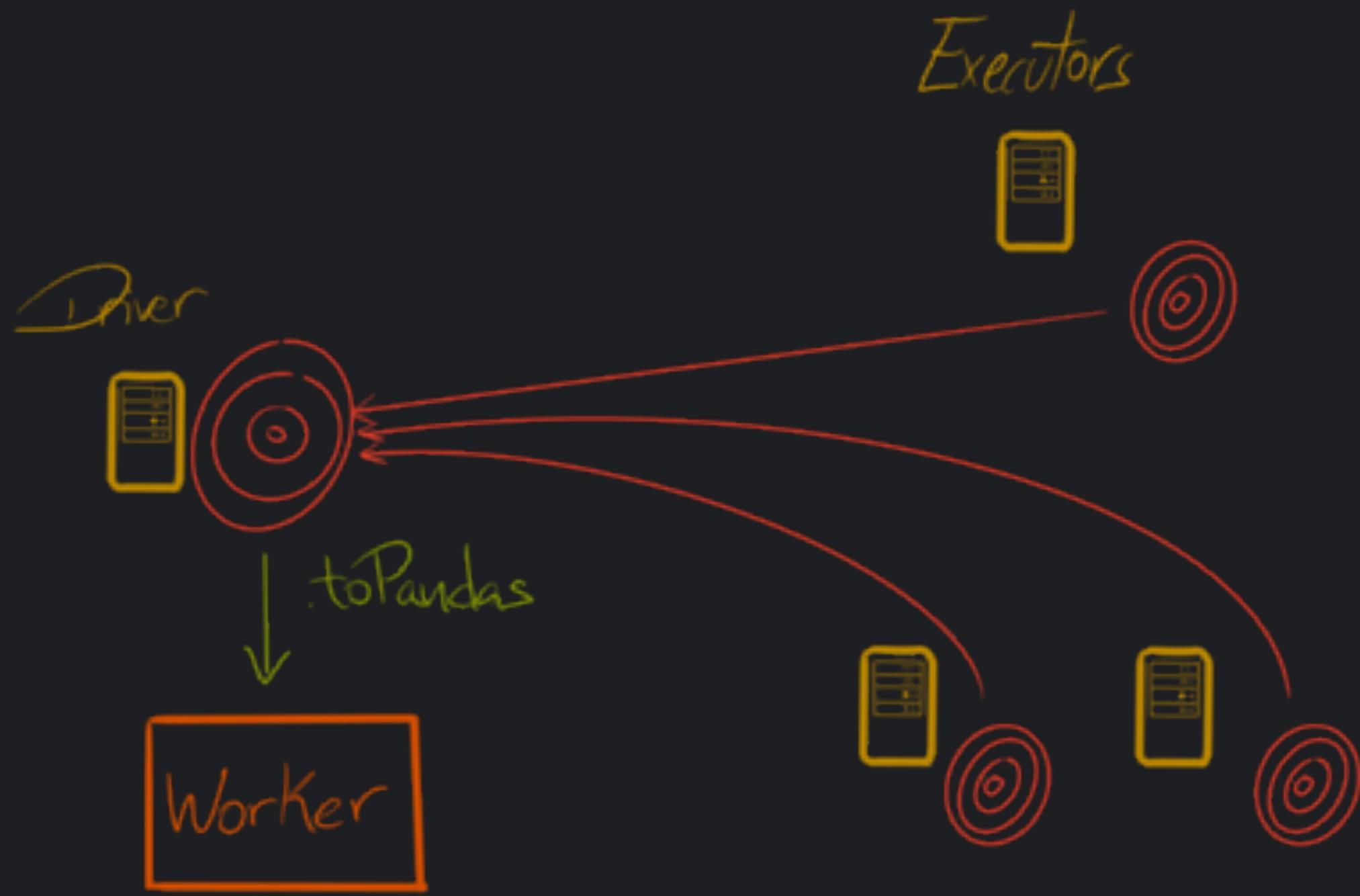
aggregate

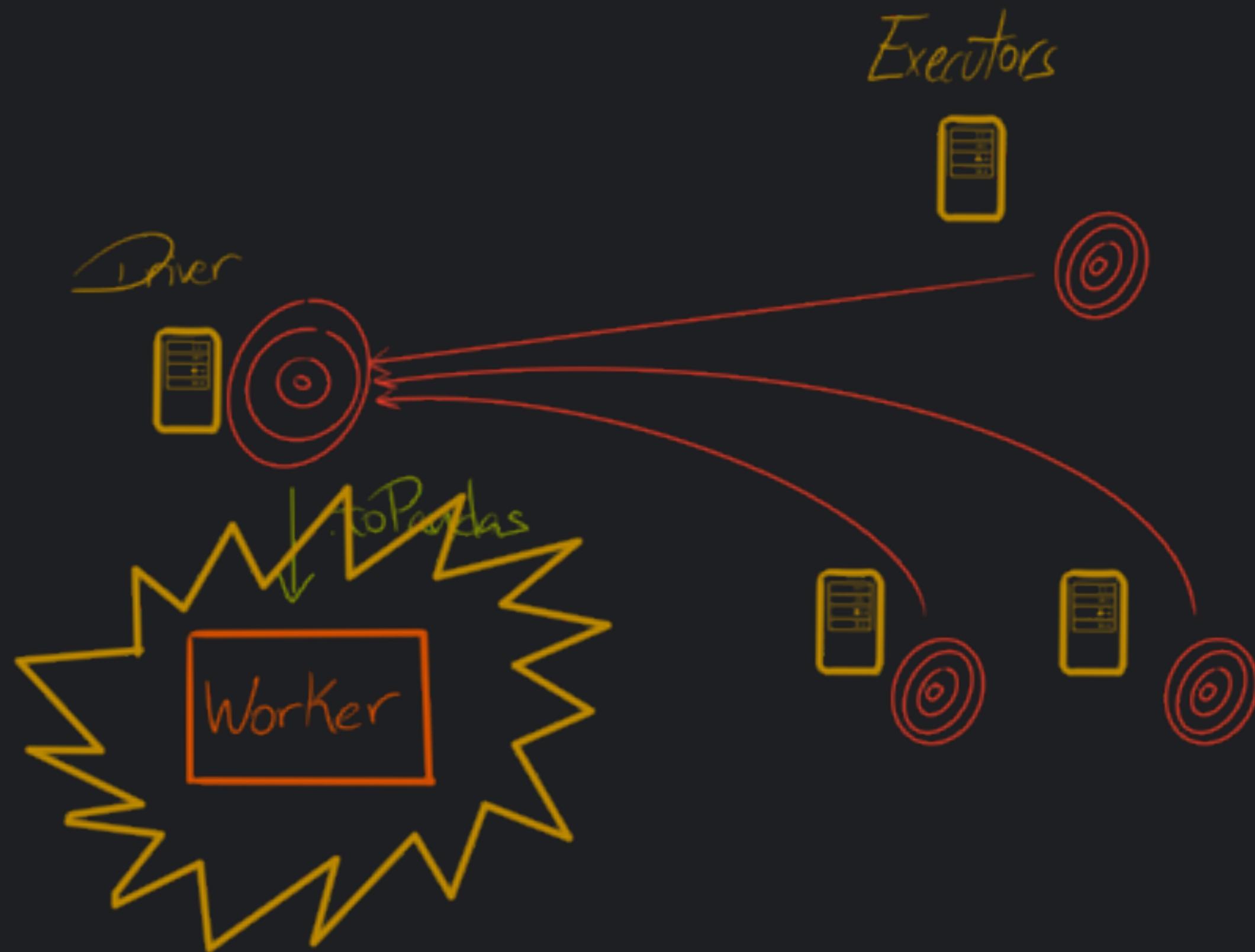
collect

Executors









TLDR: USE¹ ARROW AND PANDAS UDFS

¹IN PYSPARK

RESOURCES

- > [SPARK DOCUMENTATION](#)
- > [HIGH PERFORMANCE SPARK BY HOLDEN KARAU](#)
- > [THE INTERNALS OF APACHE SPARK 2.4.2 BY JACEK LASKOWSKI](#)
 - > [SPARK'S GITHUB](#)
 - > [BECOME A CONTRIBUTOR](#)

QUESTIONS?

A person is playing a double bass on stage. They are wearing a patterned shirt and dark pants. The stage is lit with red lights, creating a warm atmosphere. The person is looking down at their instrument.

THANKS!

GET THE SLIDES FROM MY GITHUB:

github.com/rberenguel/

THE REPOSITORY IS

pyspark-arrow-pandas





A complex, abstract background graphic is composed of numerous small, glowing blue and purple dots connected by thin lines, forming a three-dimensional grid-like structure that suggests a neural network or a data flow through a system of nodes.

Build. Unify. Scale.

WIFI SSID:Spark+AISummit | Password: UnifiedDataAnalytics

FURTHER REFERENCES

ARROW

[ARROW'S HOME](#)

[ARROW'S GITHUB](#)

[ARROW SPEED BENCHMARKS](#)

[ARROW TO PANDAS CONVERSION BENCHMARKS](#)

[POST: STREAMING COLUMNAR DATA WITH APACHE ARROW](#)

[POST: WHY PANDAS USERS SHOULD BE EXCITED BY APACHE ARROW](#)

[CODE: ARROW-PANDAS COMPATIBILITY LAYER CODE](#)

[CODE: ARROW TABLE CODE](#)

[PYARROW IN-MEMORY DATA MODEL](#)

[BALLISTA: A POC DISTRIBUTED COMPUTE PLATFORM \(RUST\)](#)

[PYJAVA: POC ON JAVA/SCALA AND PYTHON DATA INTERCHANGE WITH ARROW](#)

PANDAS

PANDAS' HOME

PANDAS' GITHUB

GUIDE: IDIOMATIC PANDAS

CODE: PANDAS INTERNALS

DESIGN: PANDAS INTERNALS

TALK: DEMYSTIFYING PANDAS' INTERNALS. BY MARC GARCIA

MEMORY LAYOUT OF MULTIDIMENSIONAL ARRAYS IN NUMPY

SPARK/PYSPARK

[CODE: PYSPARK SERIALIZERS](#)

[JIRA: FIRST STEPS TO USING ARROW \(ONLY IN THE PYSPARK DRIVER\)](#)

[POST: SPEEDING UP PYSPARK WITH APACHE ARROW](#)

[ORIGINAL JIRA ISSUE: VECTORIZED UDFS IN SPARK](#)

[INITIAL DOC DRAFT](#)

[POST BY BRYAN CUTLER \(LEADER FOR THE VEC UDFS PR\)](#)

[POST: INTRODUCING PANDAS UDF FOR PYSPARK](#)

[CODE: ORG.APACHE.SPARK.SQL.VECTORIZED](#)

[POST BY BRYAN CUTLER: SPARK TOPANDAS\(\) WITH ARROW. A DETAILED LOOK](#)

PY4J

PY4J'S HOME
PY4J'S GITHUB
CODE: REFLECTION ENGINE

TABLE FOR toPandas

2^X	DIRECT (S)	WITH ARROW (S)	FACTOR
17	1.08	0.18	5.97
18	1.69	0.26	6.45
19	4.16	0.30	13.87
20	5.76	0.61	9.44
21	9.73	0.96	10.14
22	17.90	1.64	10.91
23	(00M)	3.42	
24	(00M)	11.40	

EOF