

SPEEDING UP PYSPARK WITH ARROW



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?

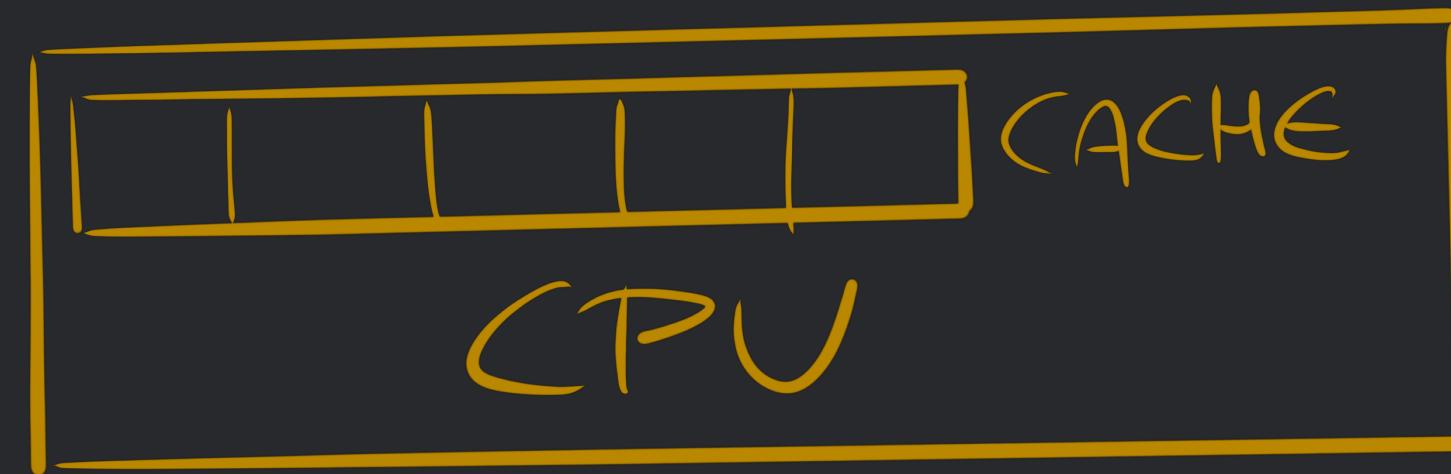
- > PYTHON DATA ANALYSIS LIBRARY
- > USED EVERYWHERE DATA AND PYTHON APPEAR IN JOB OFFERS
- > EFFICIENT (IS COLUMNAR AND HAS A C AND CYTHON BACKEND)

ROW

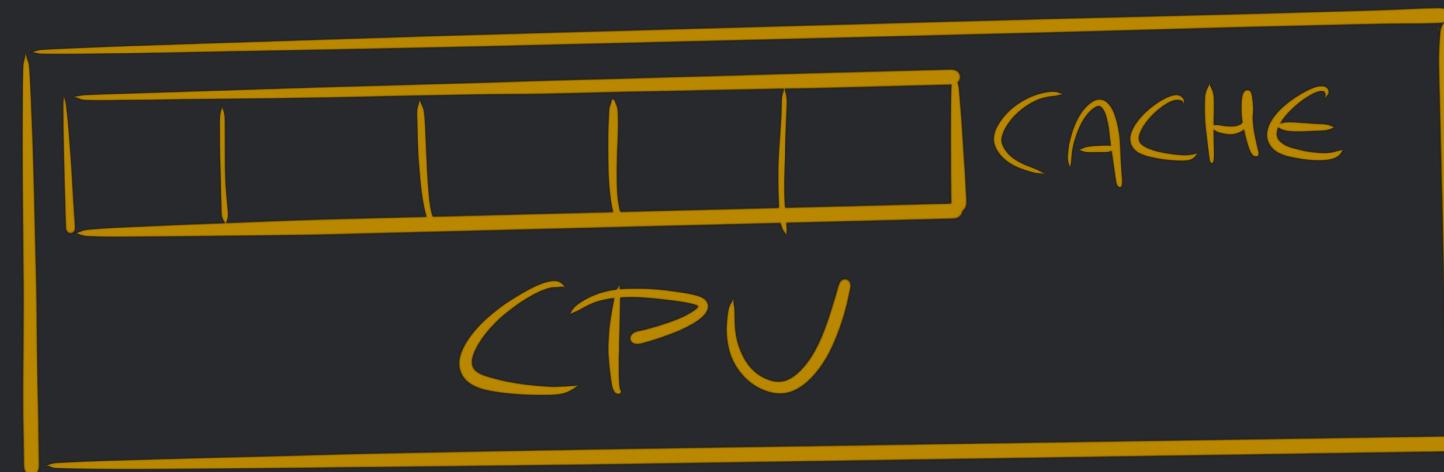
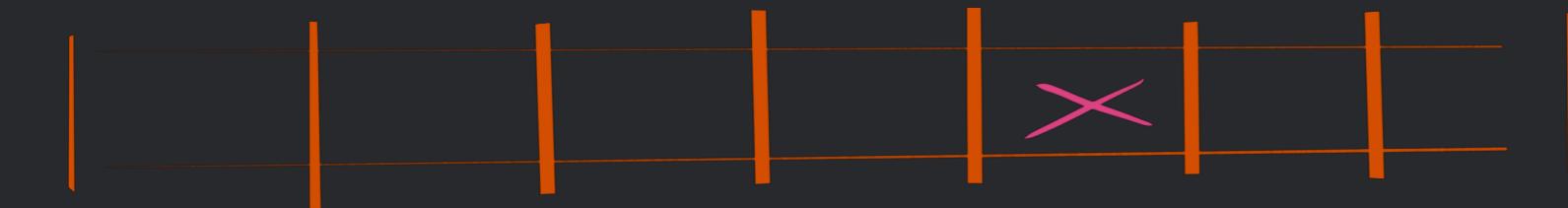


CPU

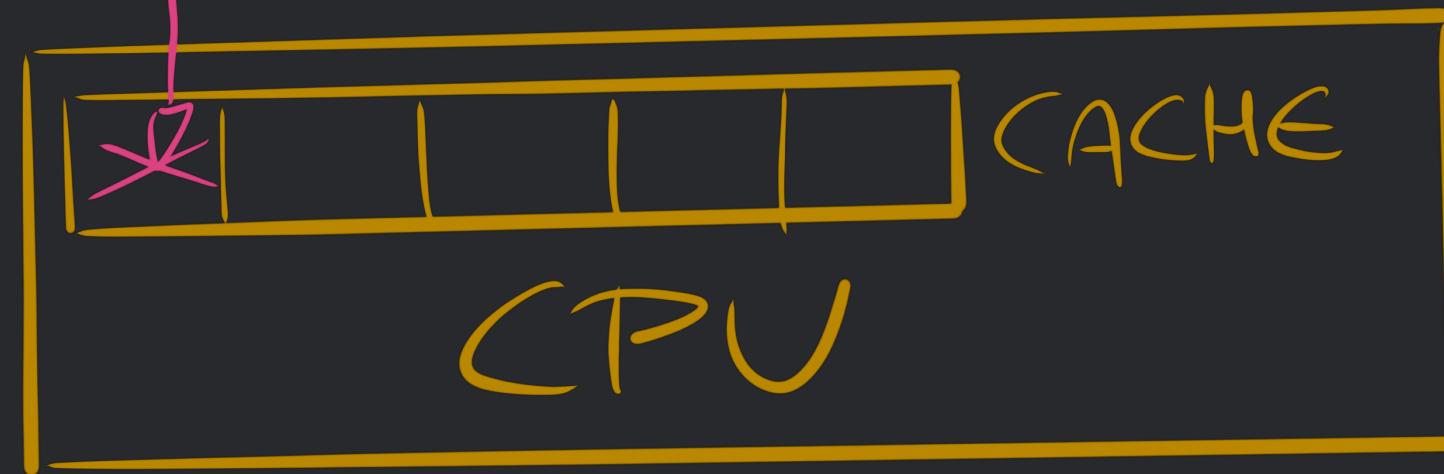
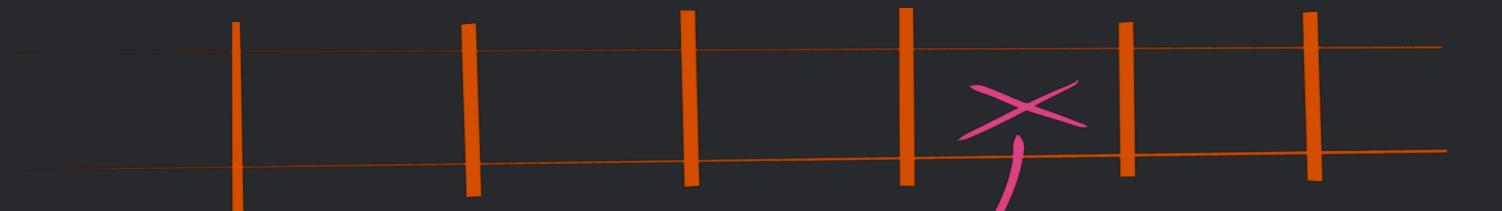
ROW



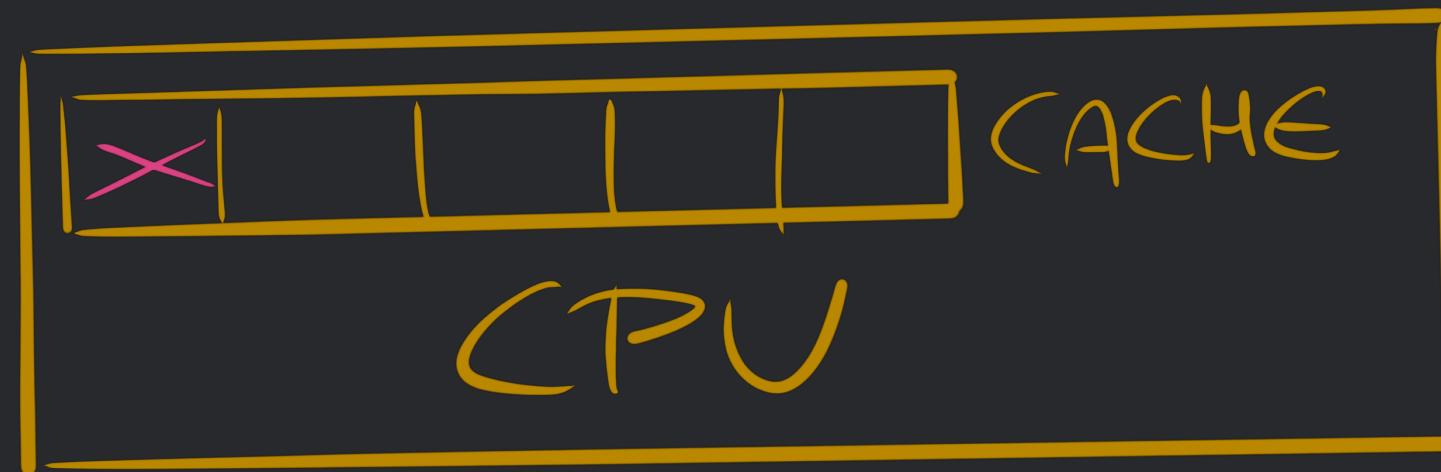
ROW



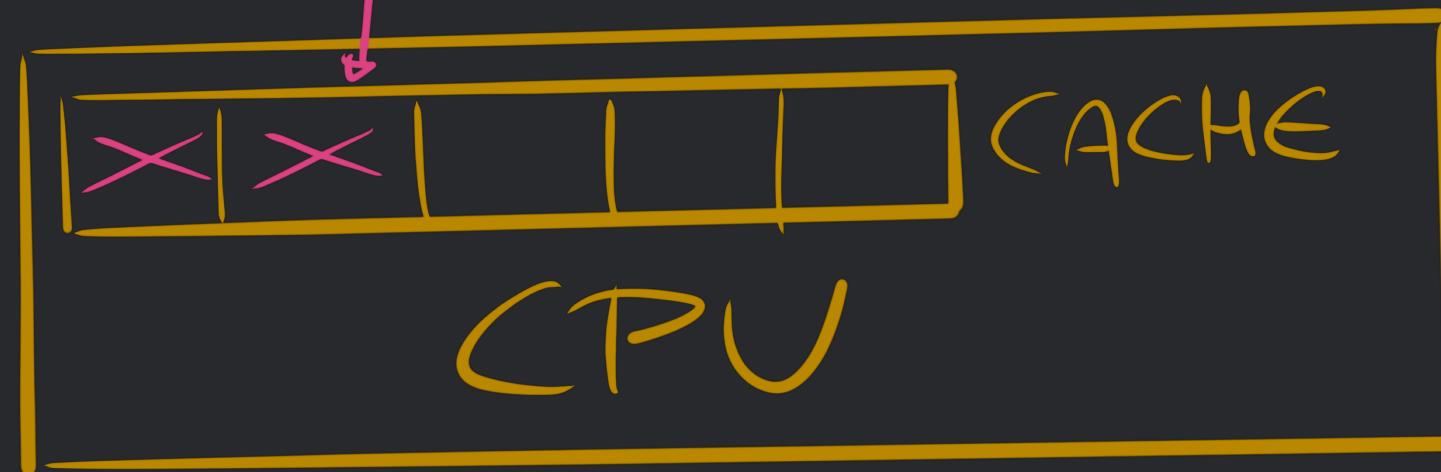
ROW

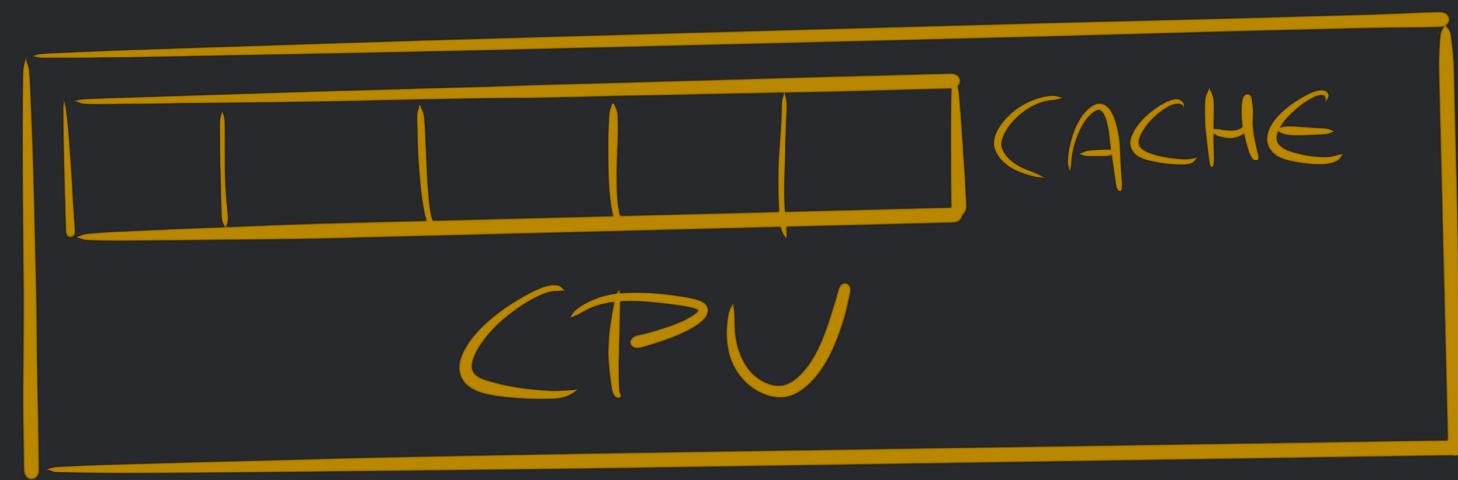
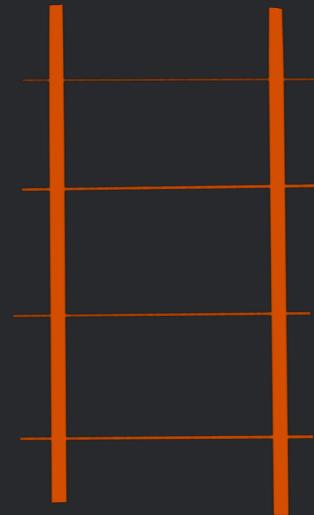


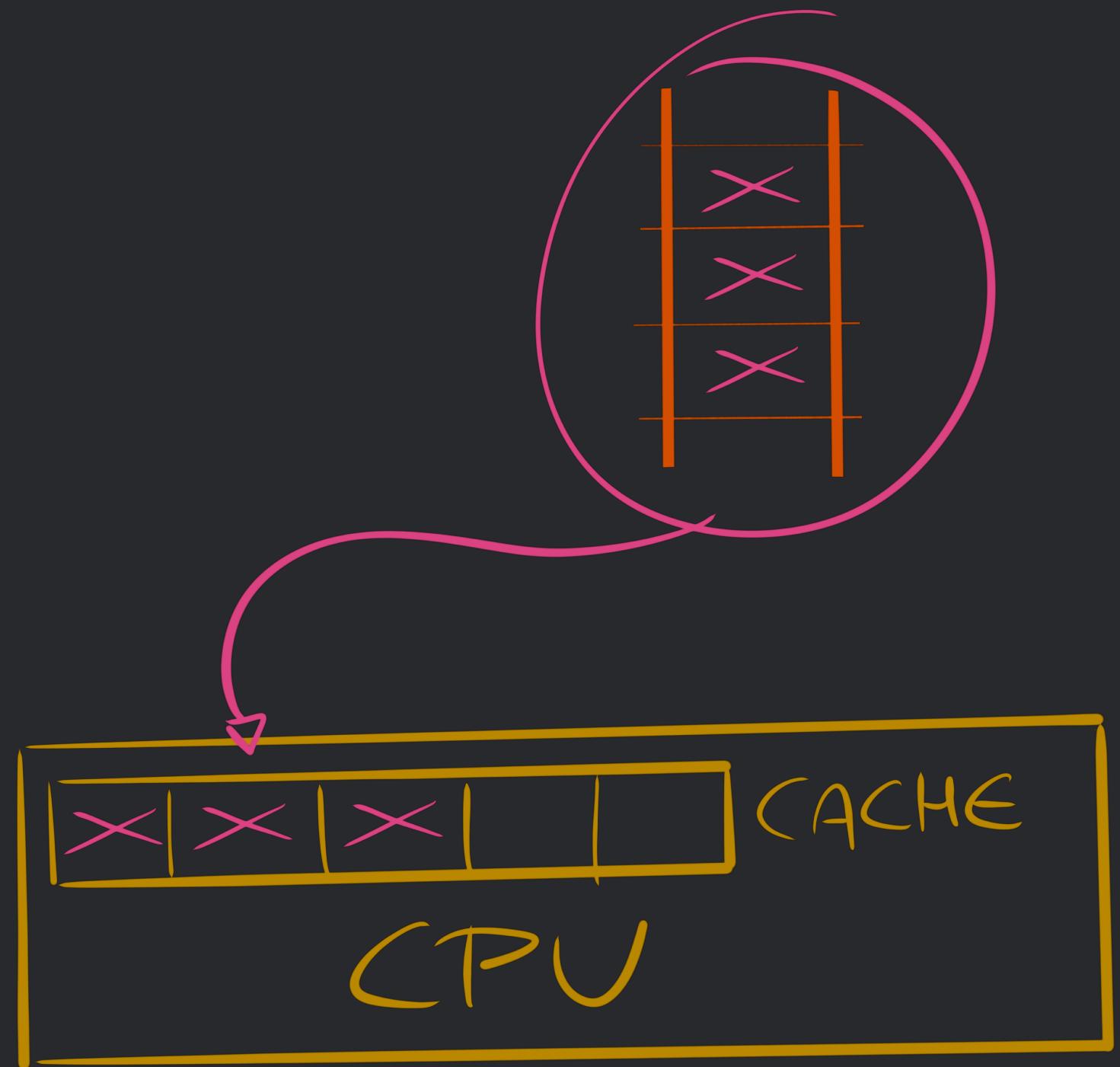
ROW



ROW







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
:				

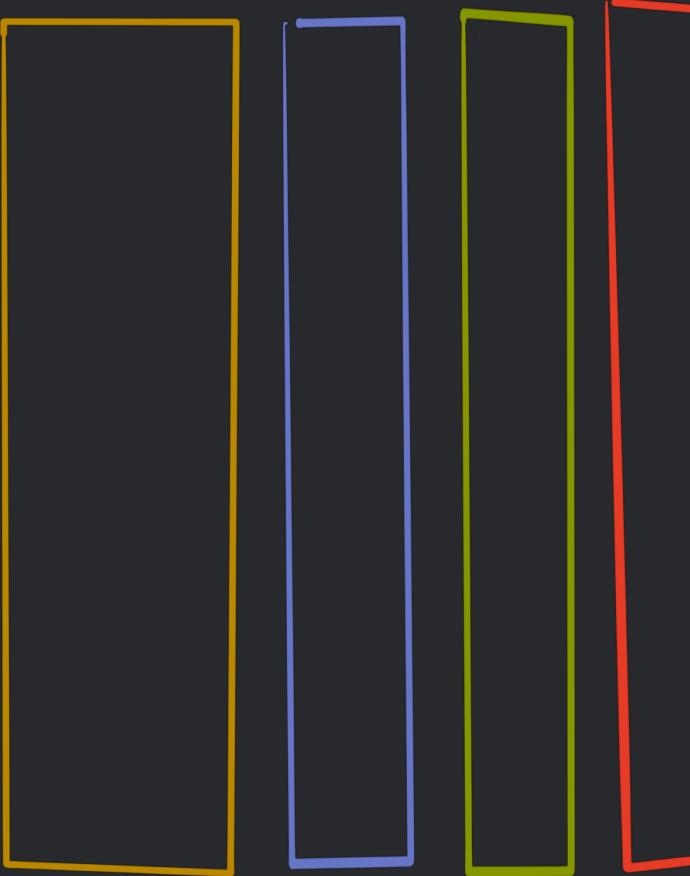
Pandas DataFrame

	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

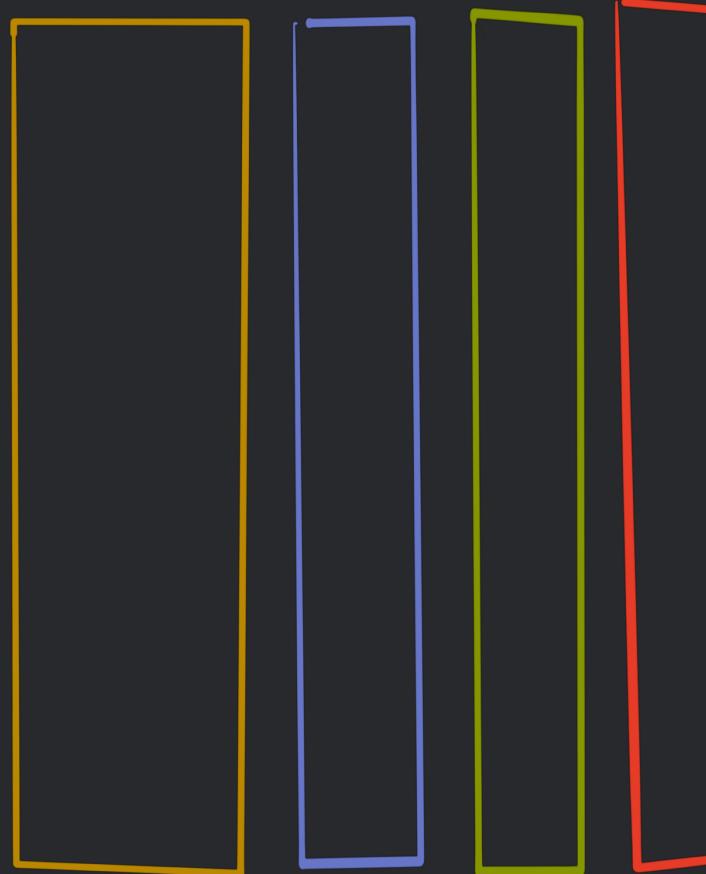
	A	B	C	D
1	4.4	a	2	0.1
2	3.1	b	42	0.4
3	2	c	8	0.9
4	8.3	d	15	0.3
:				



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
:				

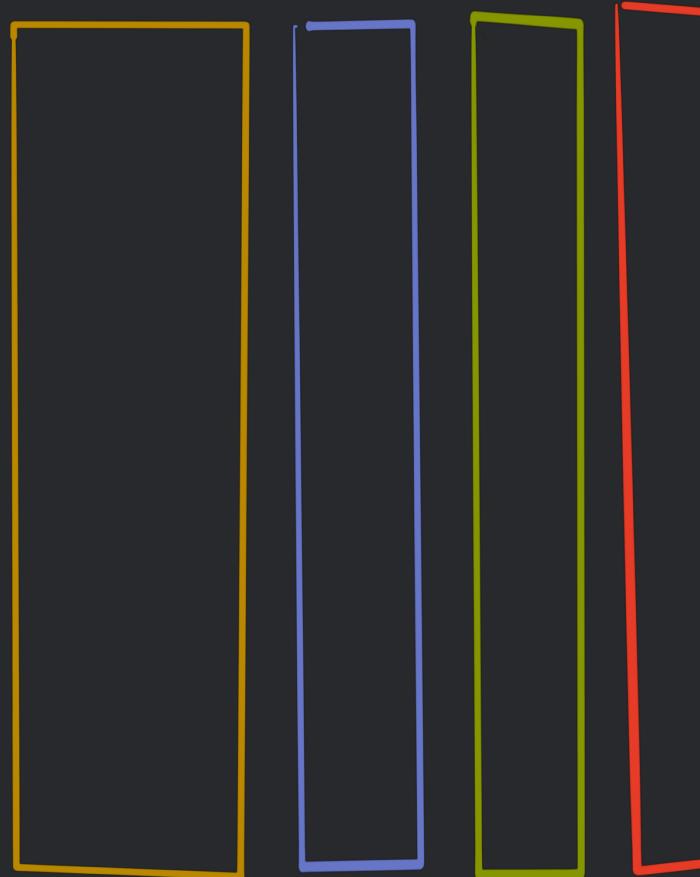


$2 \times N$ FloatBlock

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
⋮				

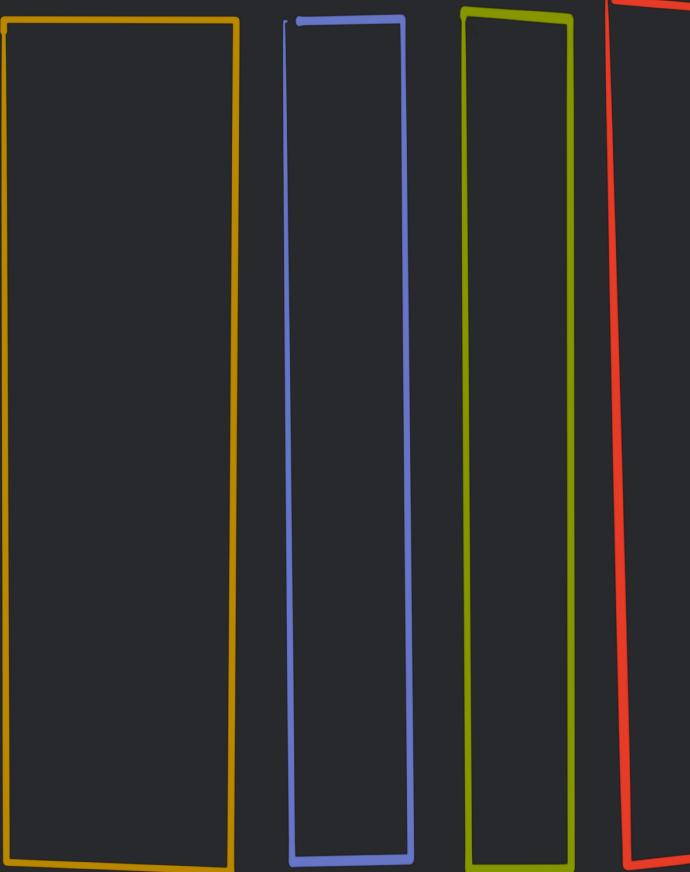
1 × N ObjectBlock



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
⋮				

$1 \times N$ 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?

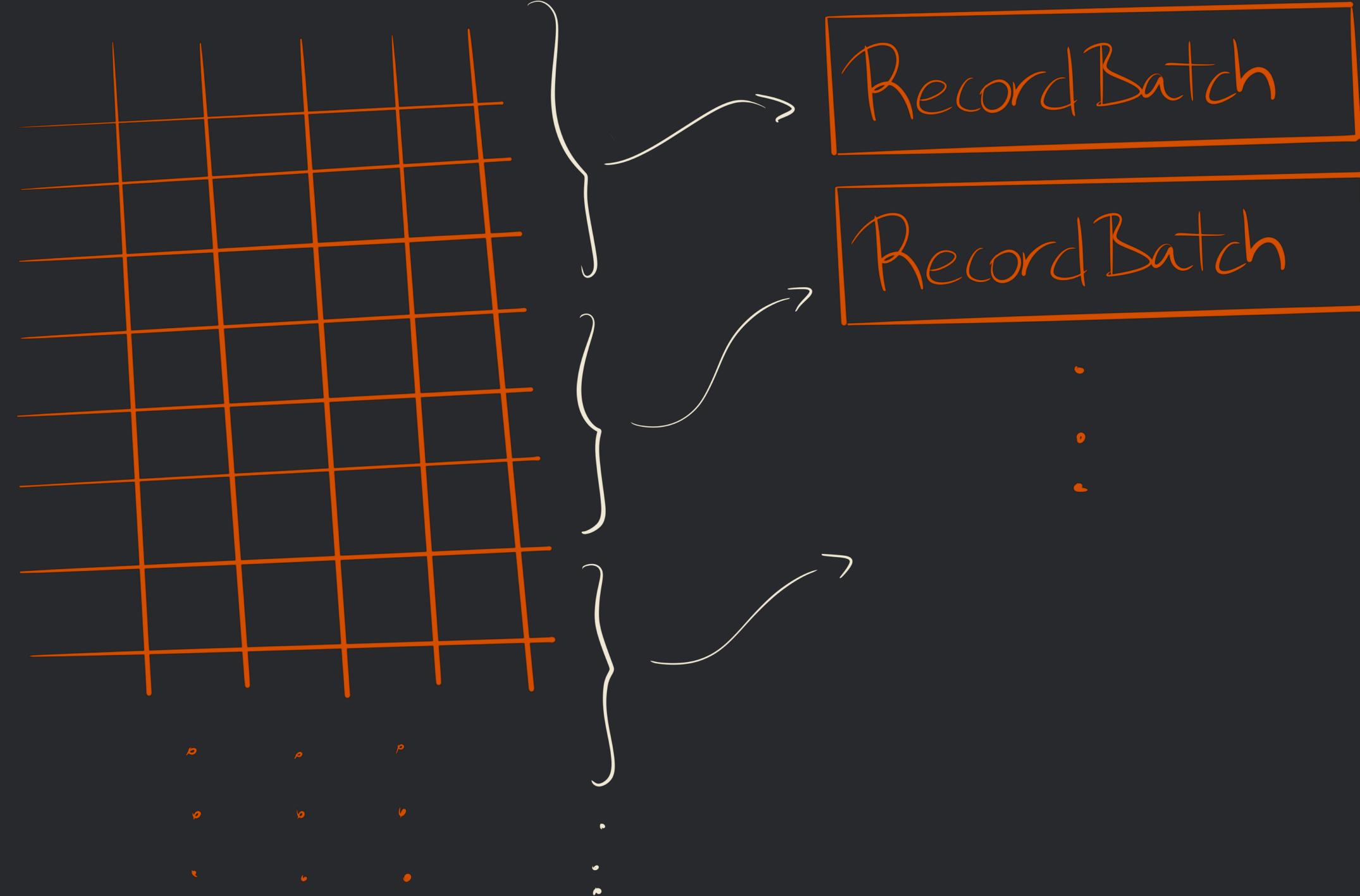
Table

P P P

P P P

P P P

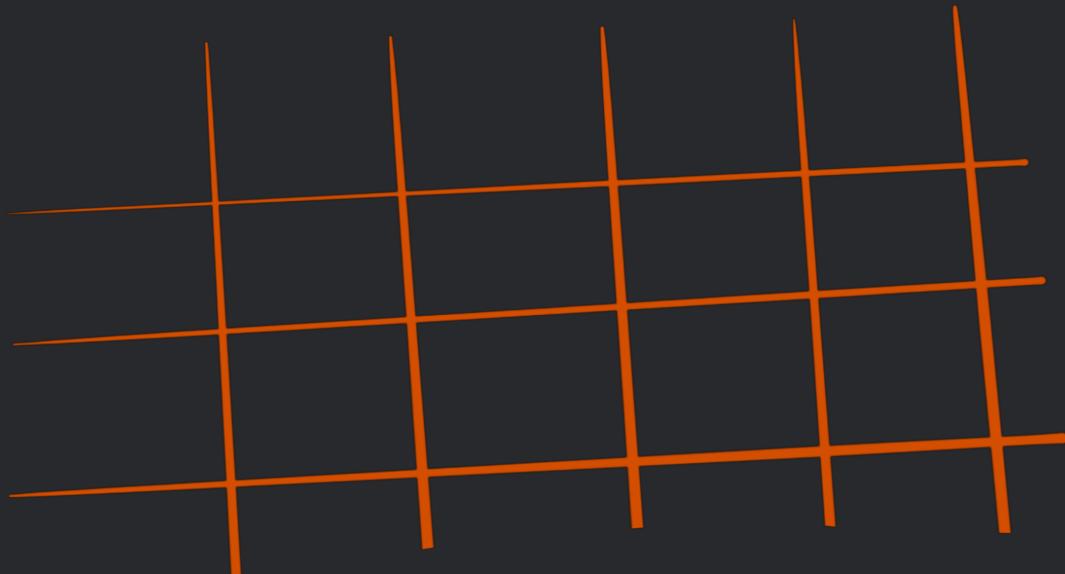
Table



Some rows

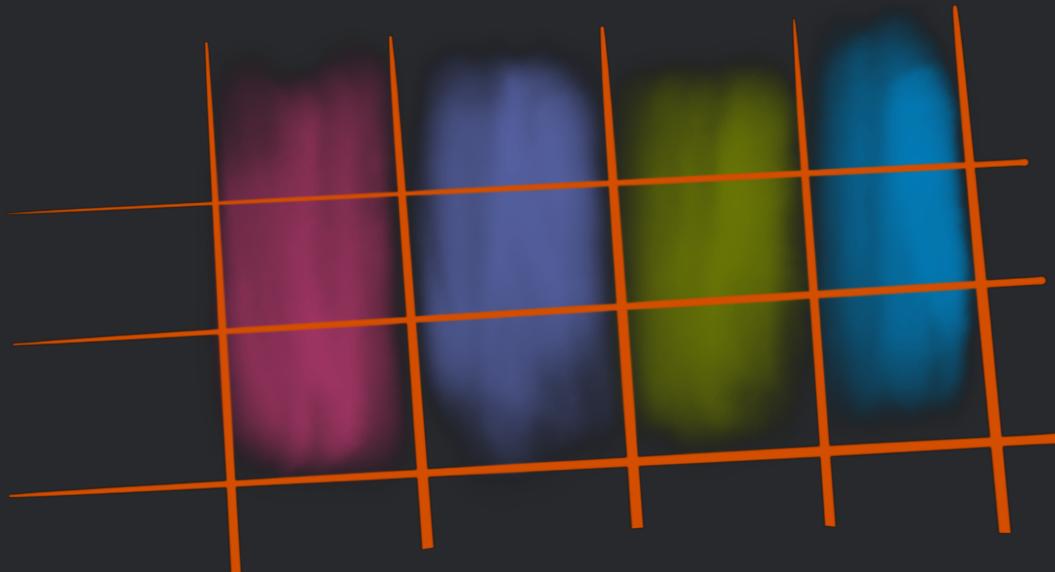


Some rows



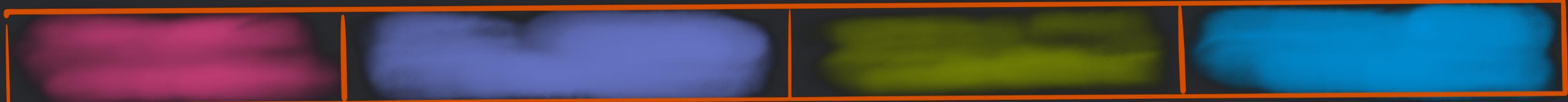
RecordBatch

Some rows



RecordBatch

(+ metadata)

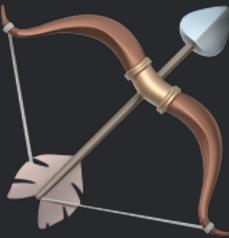


Arrow Table

•

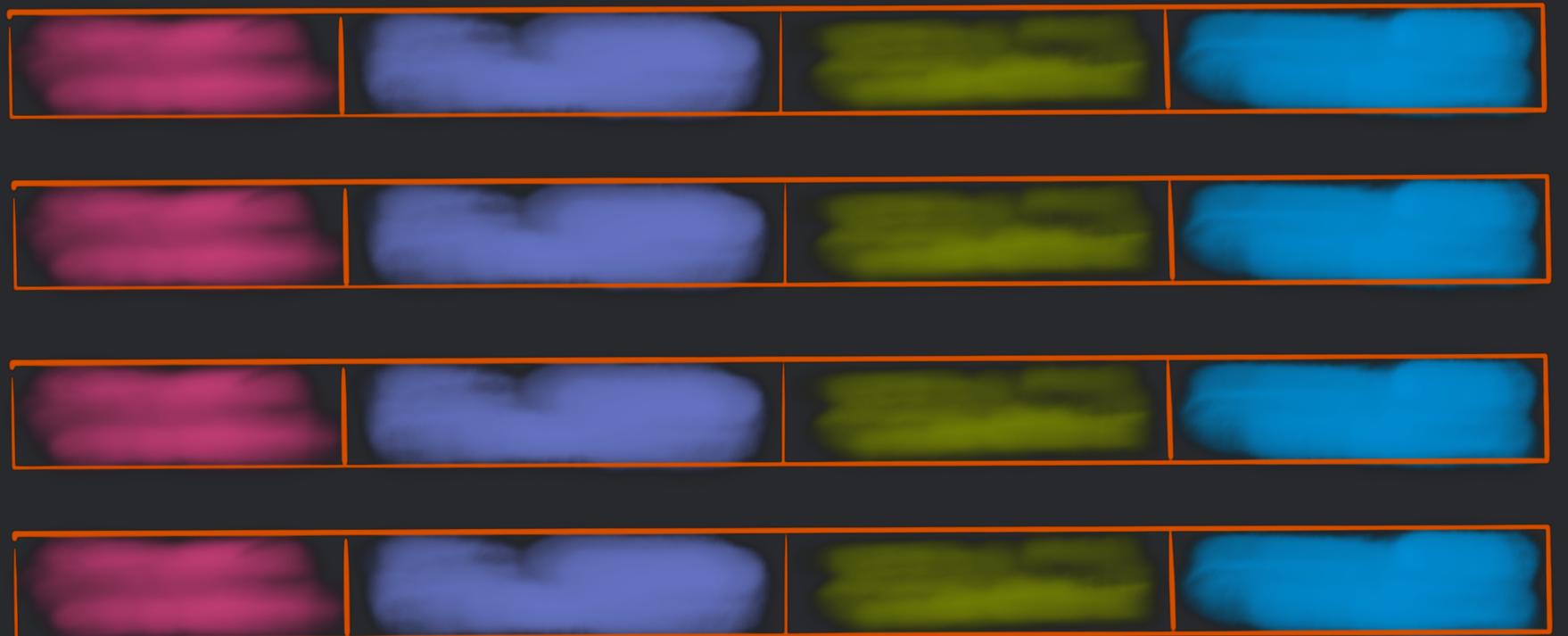
•

•

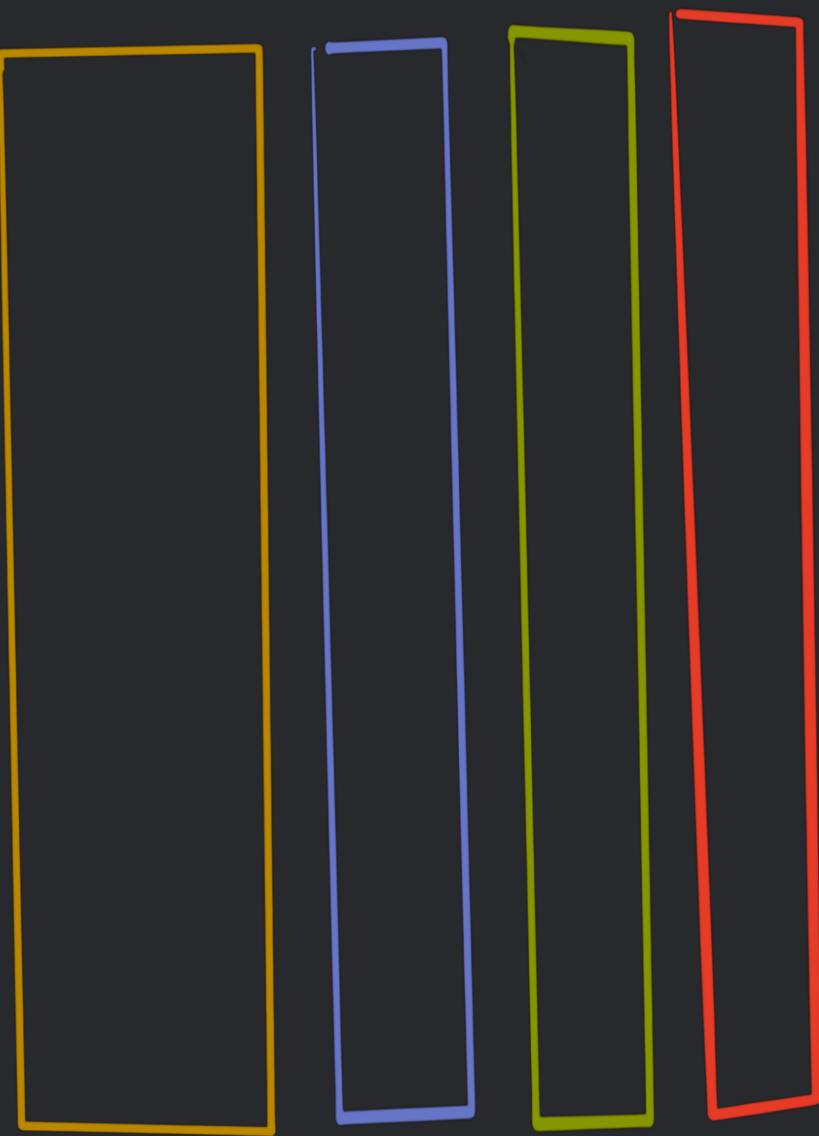


- > ARROW USES RecordBatches
- > PANDAS USES BLOCKS HANDLED BY A BlockManager
- > YOU CAN CONVERT AN ARROW Table INTO A PANDAS DataFrame EASILY

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 SITS
ON TOP OF A
**CLUSTER
MANAGER**



Cluster Manager



AND A
**DISTRIBUTED
STORAGE**

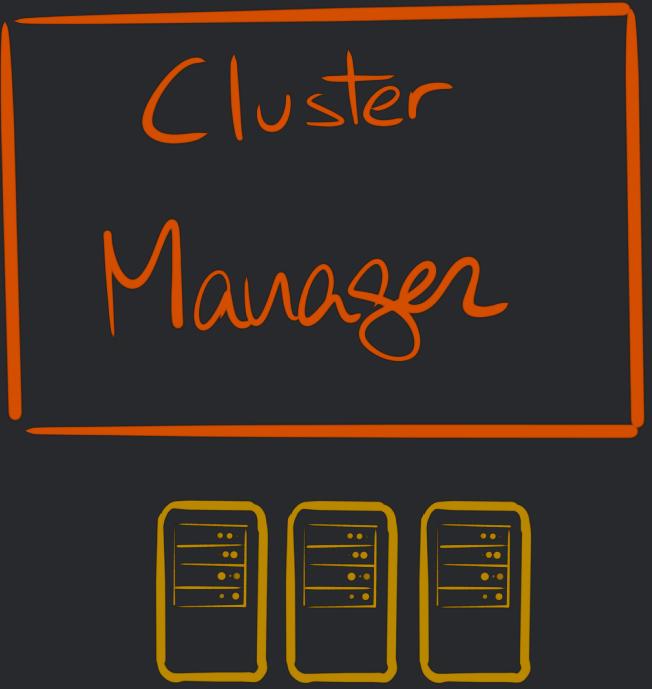
Cluster Manager

Distributed Storage

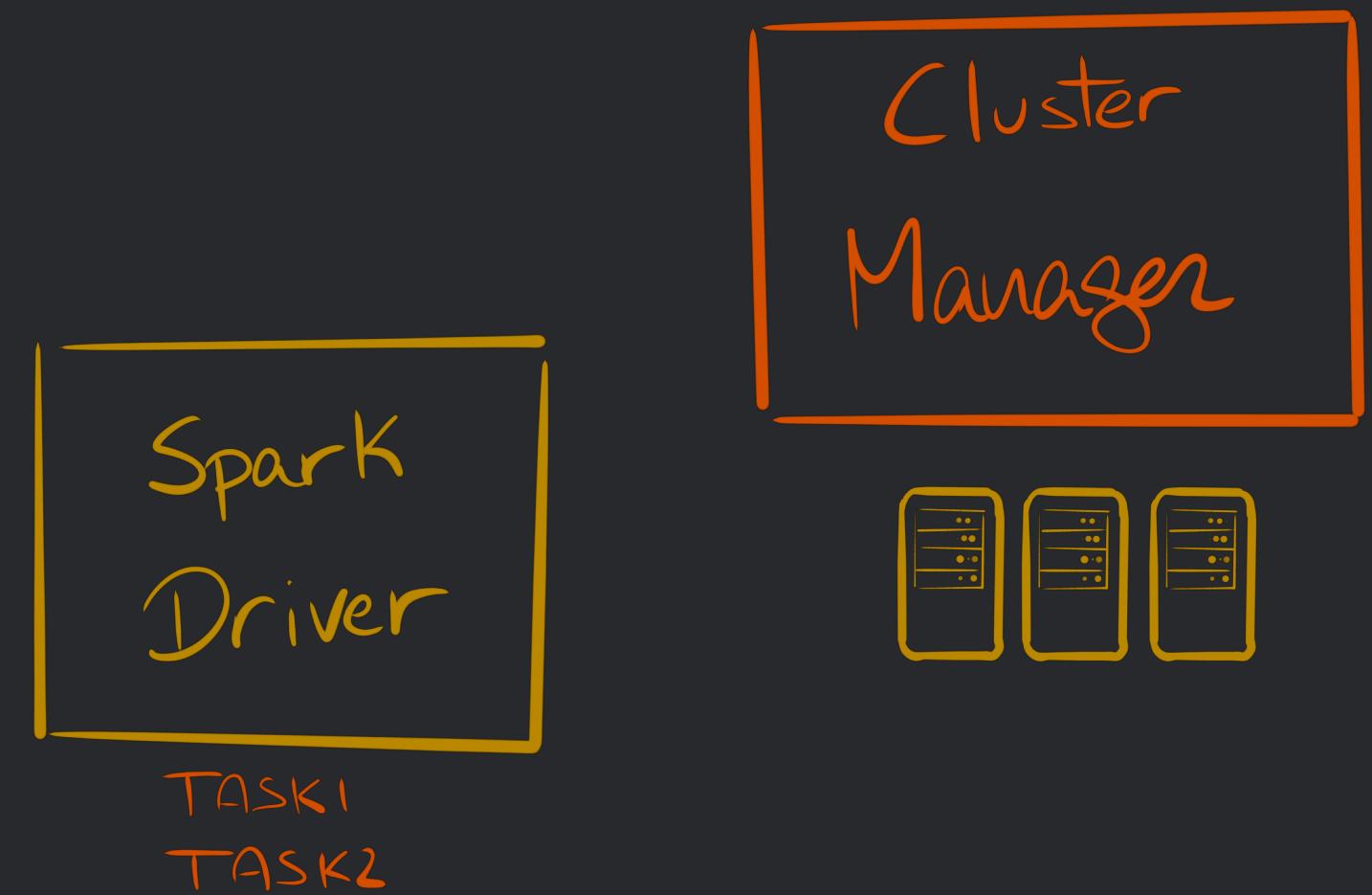


A SPARK PROGRAM
RUNS IN THE DRIVER

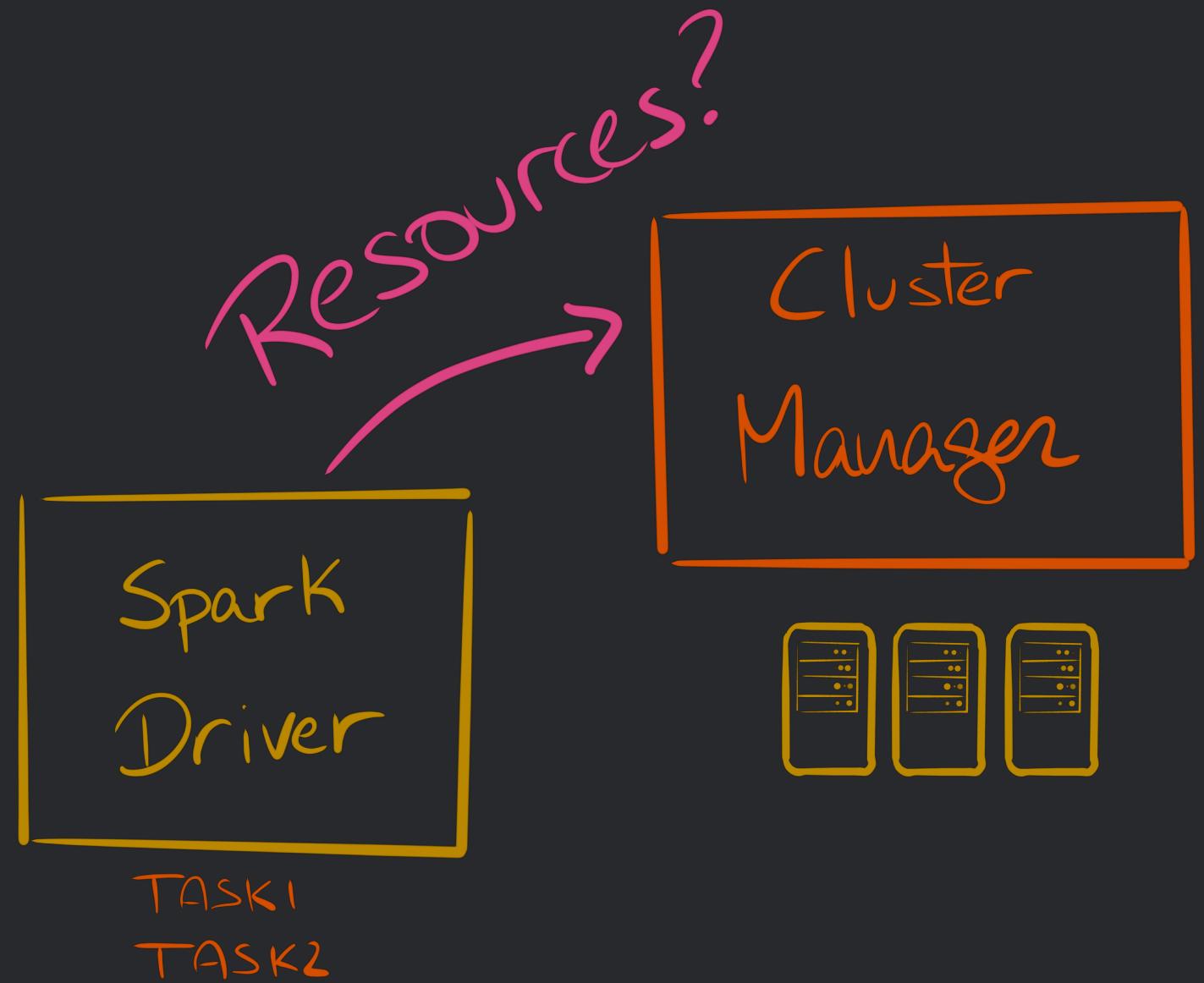
THE DRIVER REQUESTS
RESOURCES FROM THE
CLUSTER MANAGER TO
RUN TASKS



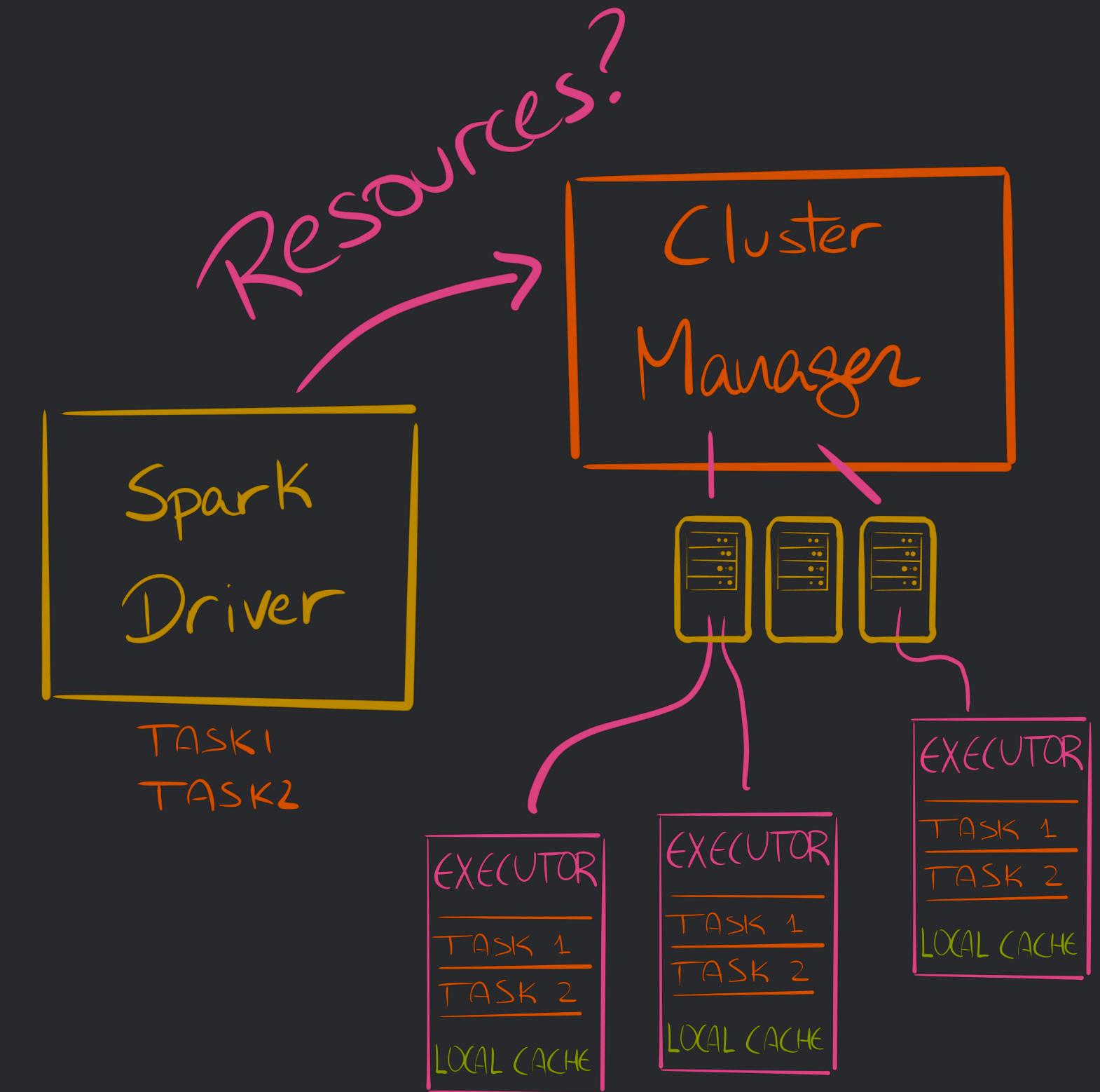
THE DRIVER REQUESTS
RESOURCES FROM THE
CLUSTER MANAGER TO
RUN TASKS



THE DRIVER REQUESTS
RESOURCES FROM THE
CLUSTER MANAGER TO
RUN TASKS



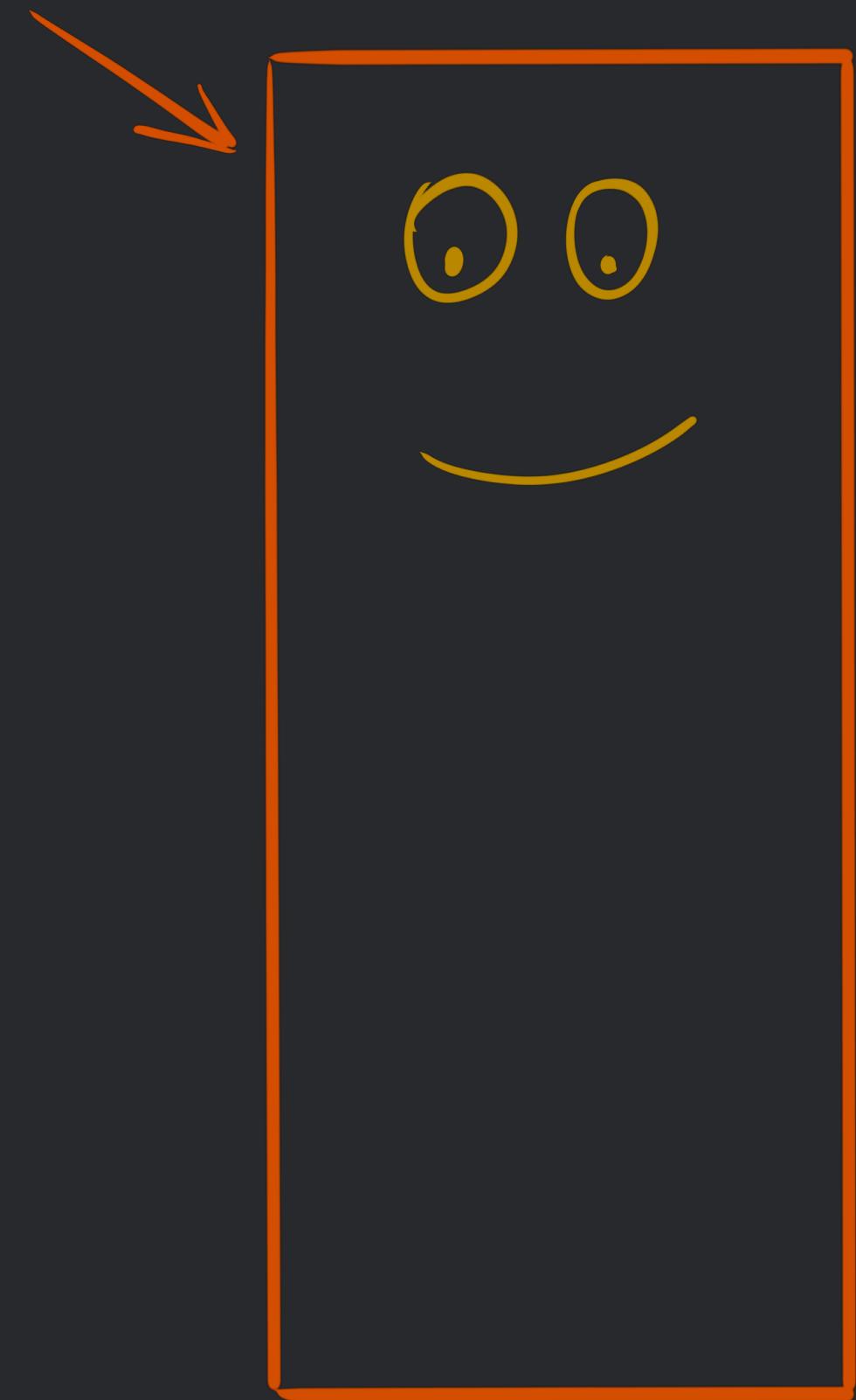
THE DRIVER REQUESTS
RESOURCES FROM THE
CLUSTER MANAGER TO
RUN TASKS



THE MAIN BUILDING BLOCK
IS THE RDD:
RESILIENT DISTRIBUTED
DATASET



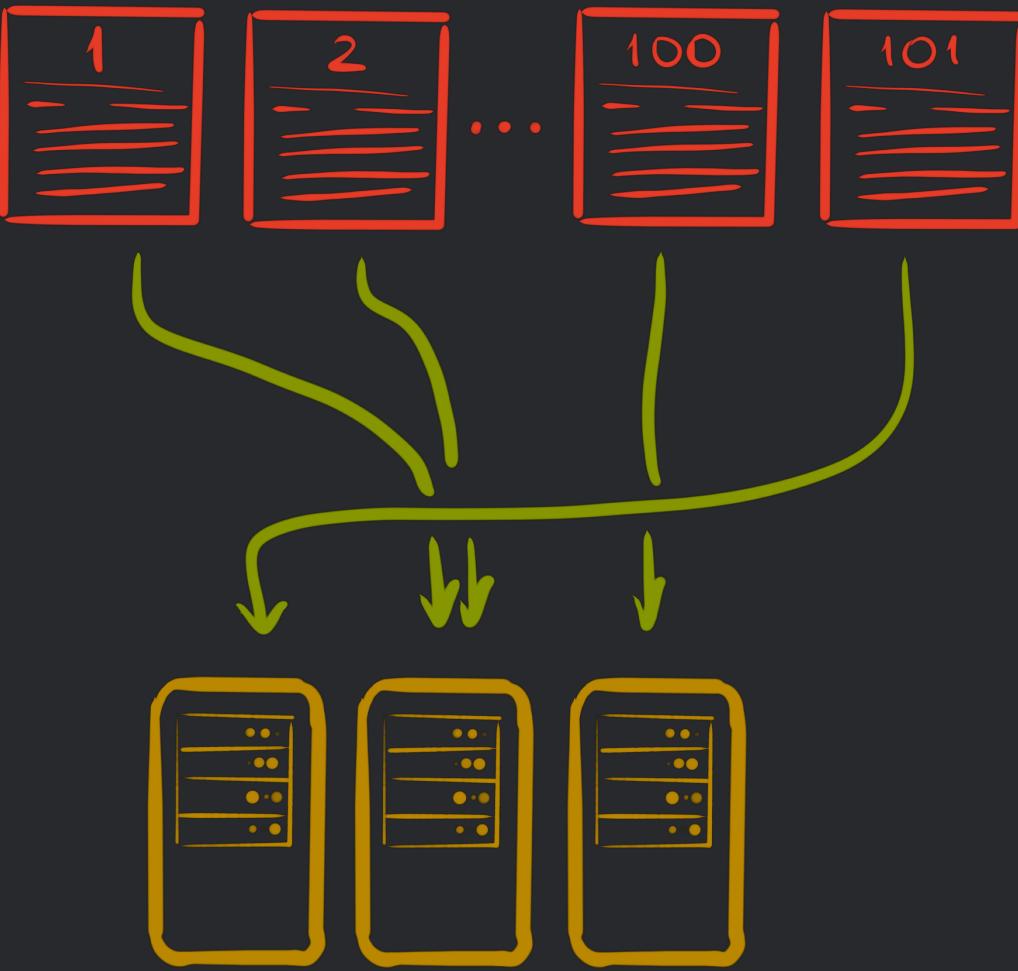
RDD



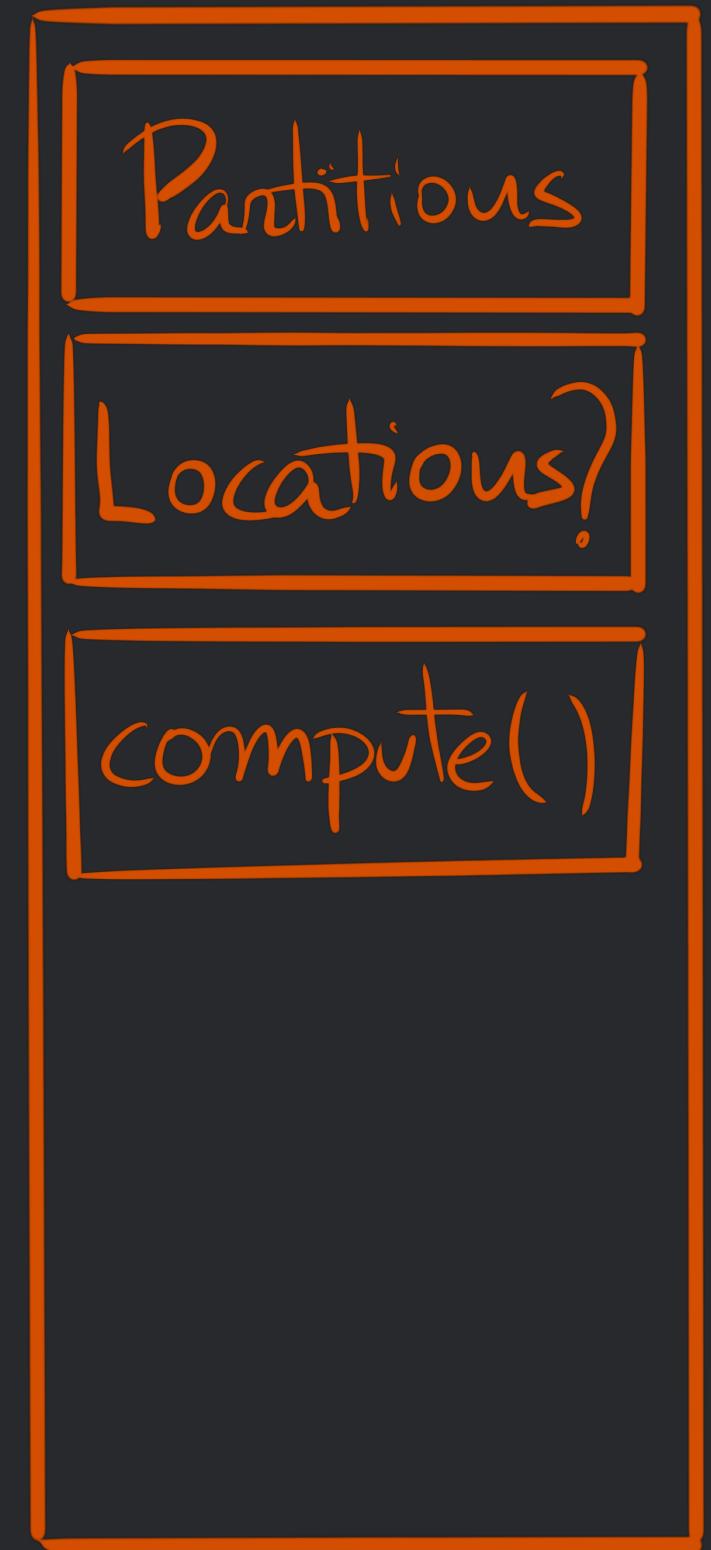
RDD



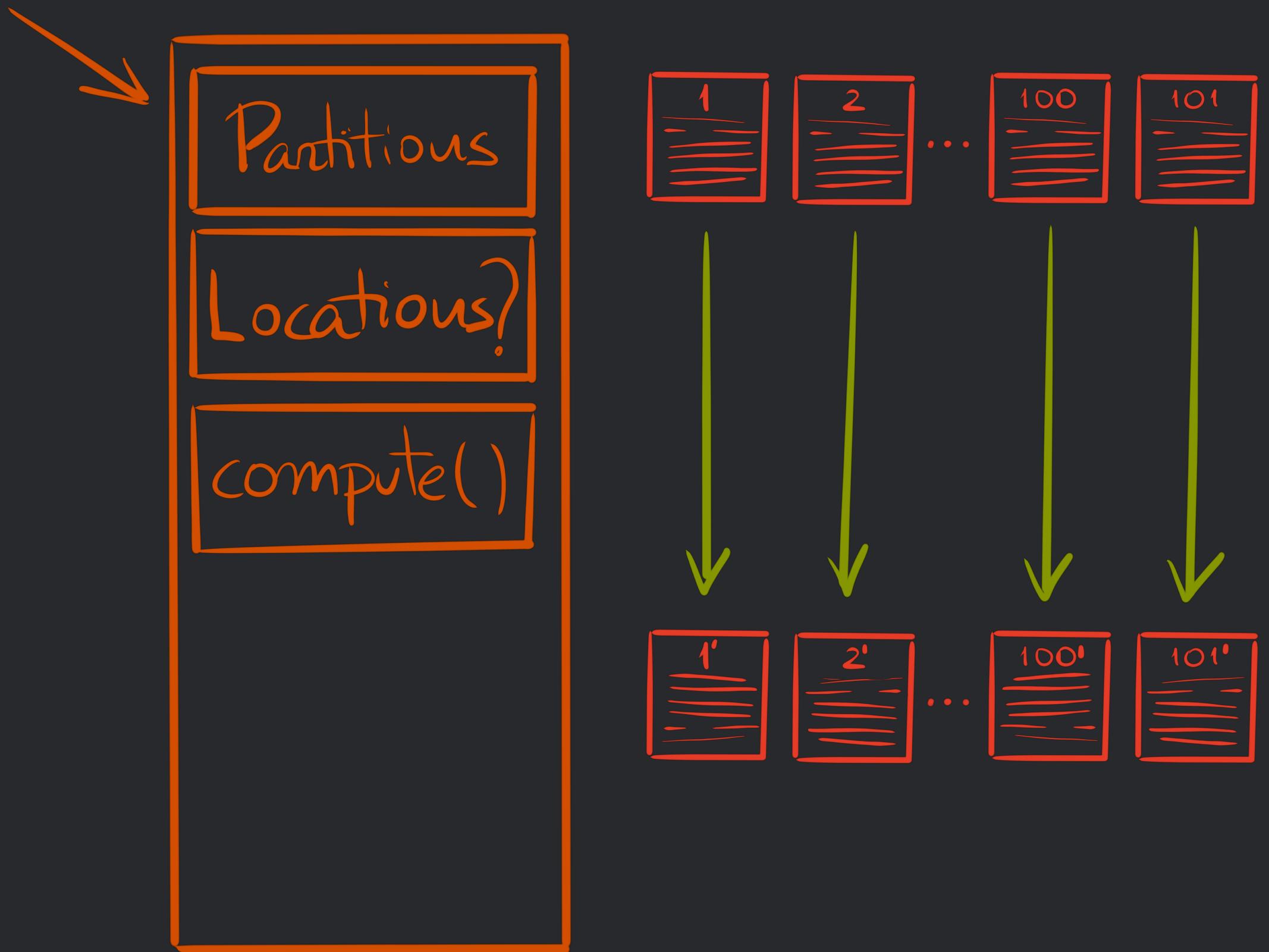
RDD



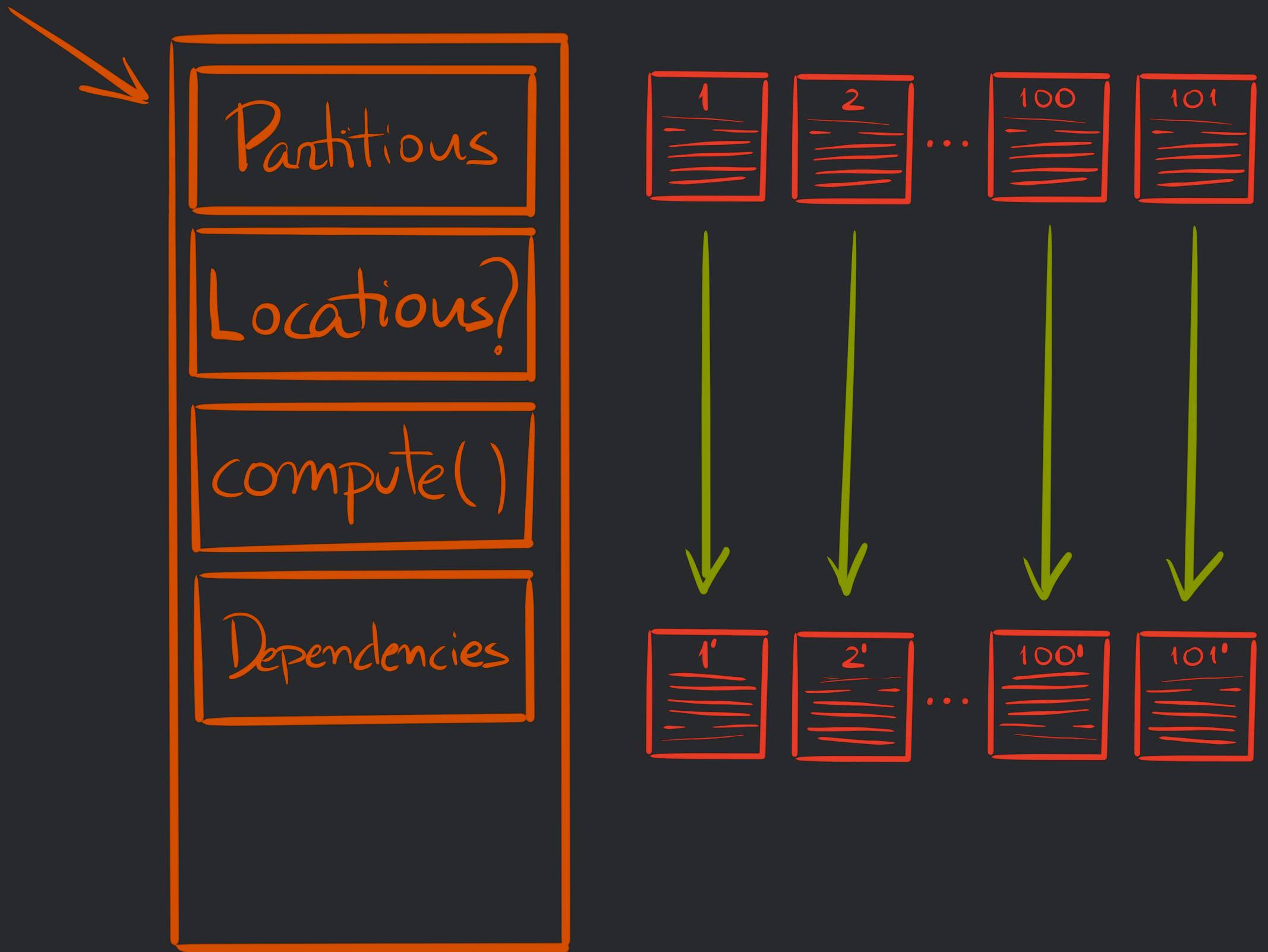
RDD



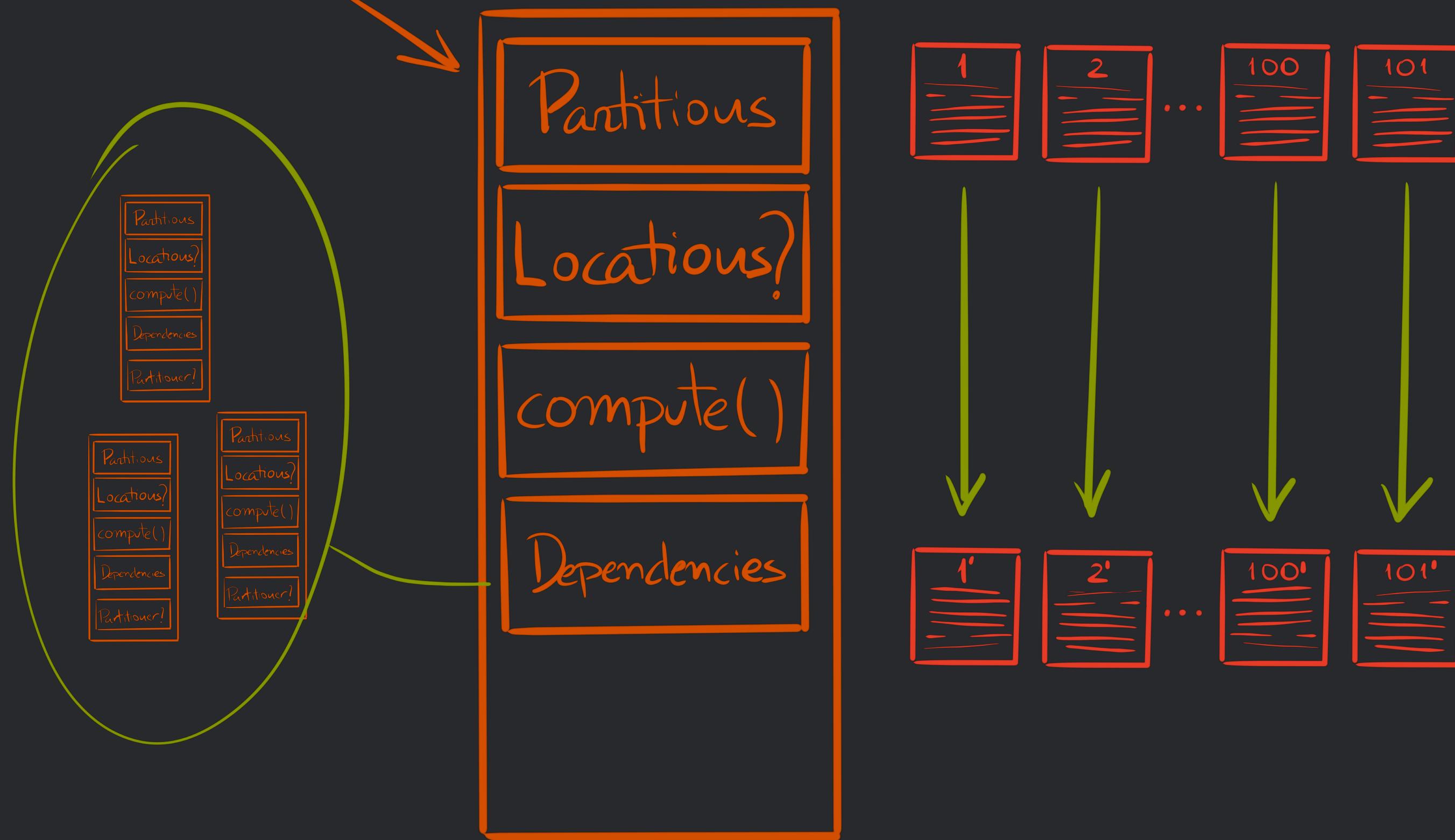
RDD



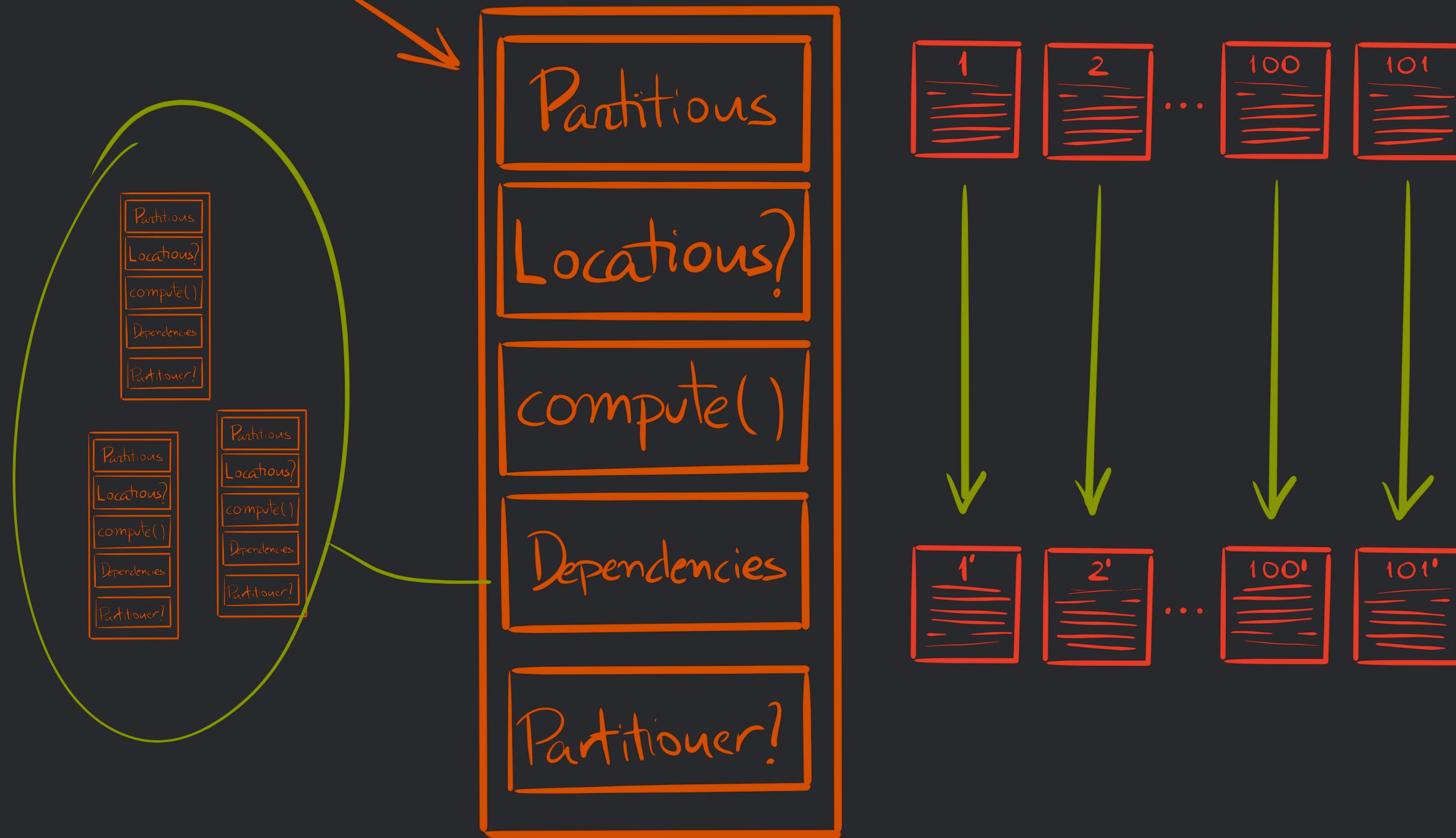
RDD



RDD



RDD



PYSPARK

PYSPARK OFFERS A
PYTHON API TO THE SCALA
CORE OF SPARK

IT USES THE
PY4J BRIDGE

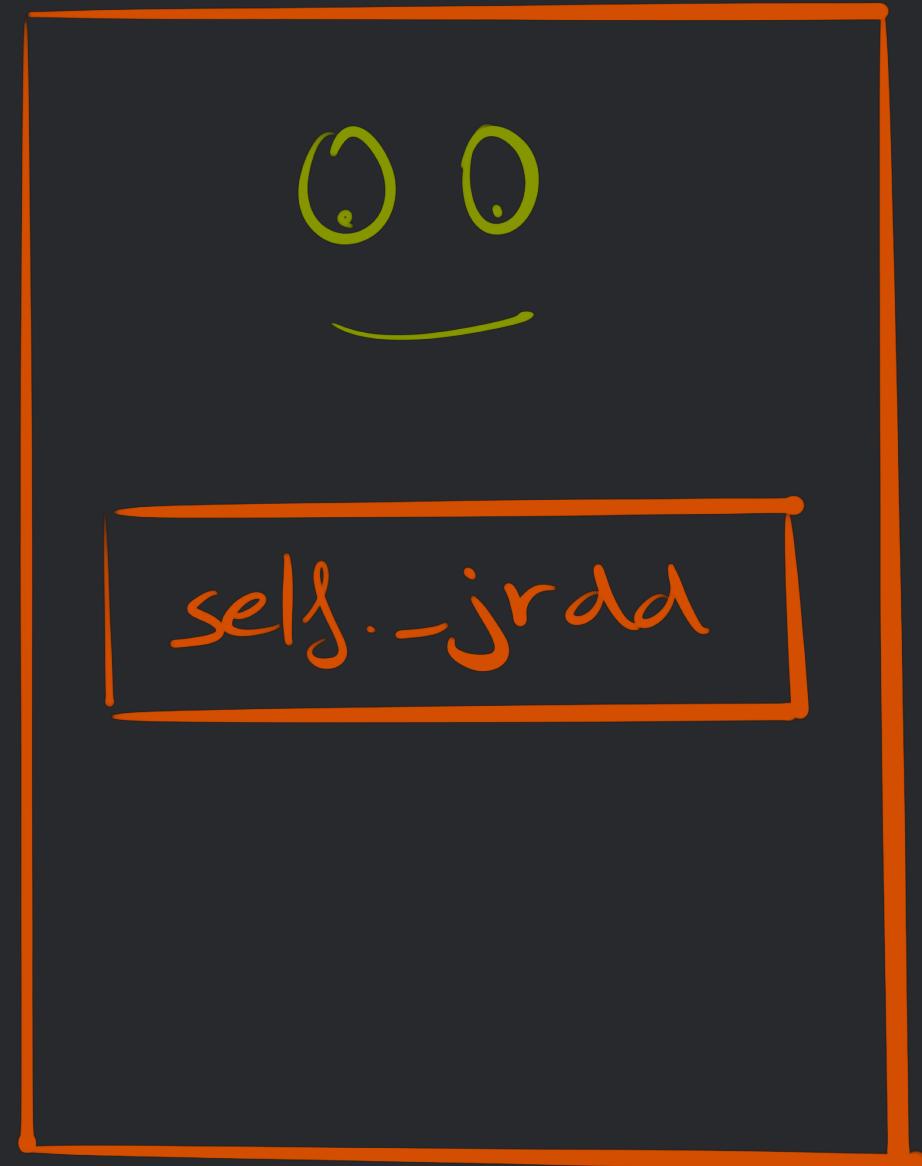
```
# Connect to the gateway
gateway = JavaGateway(
    gateway_parameters=GatewayParameters(
        port=gateway_port,
        auth_token=gateway_secret,
        auto_convert=True))

# 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.*")

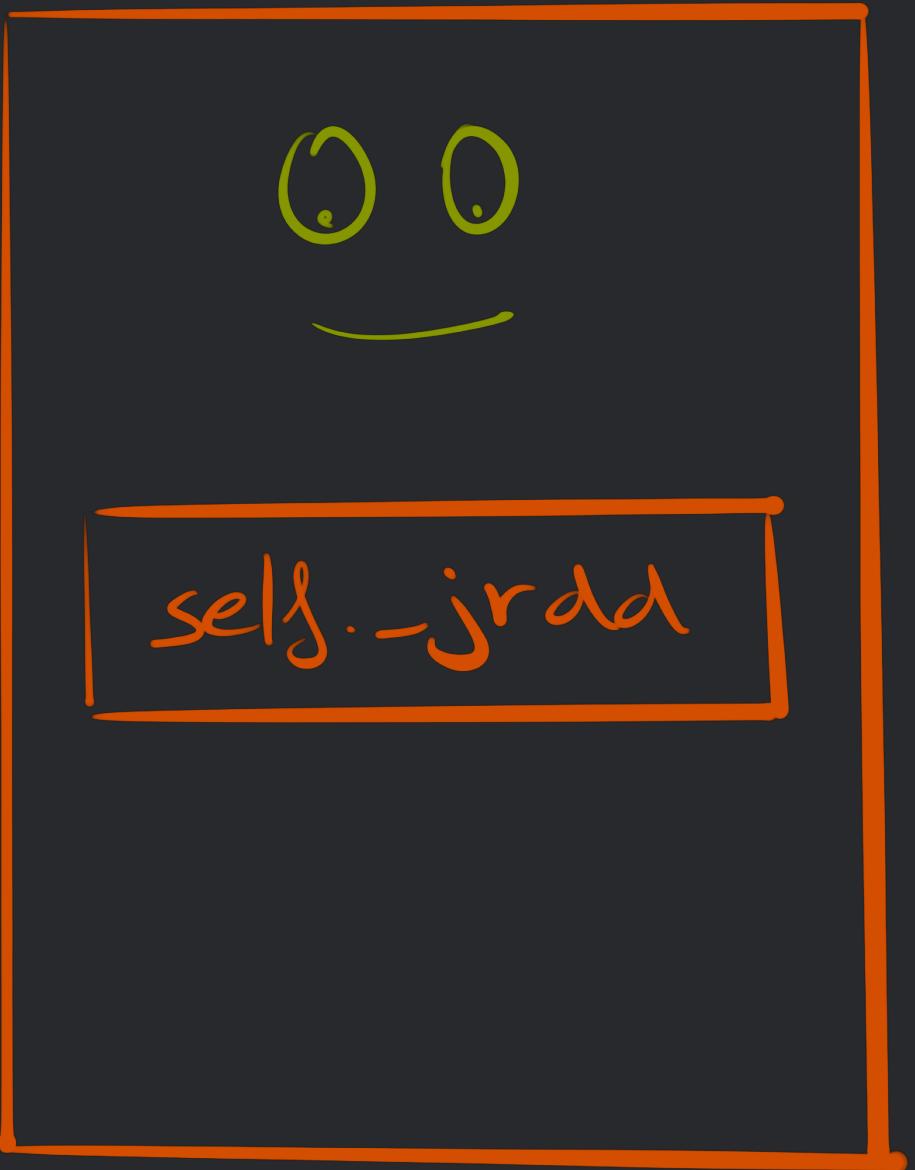
.
.
.

return gateway
```


RDD in
Python Land

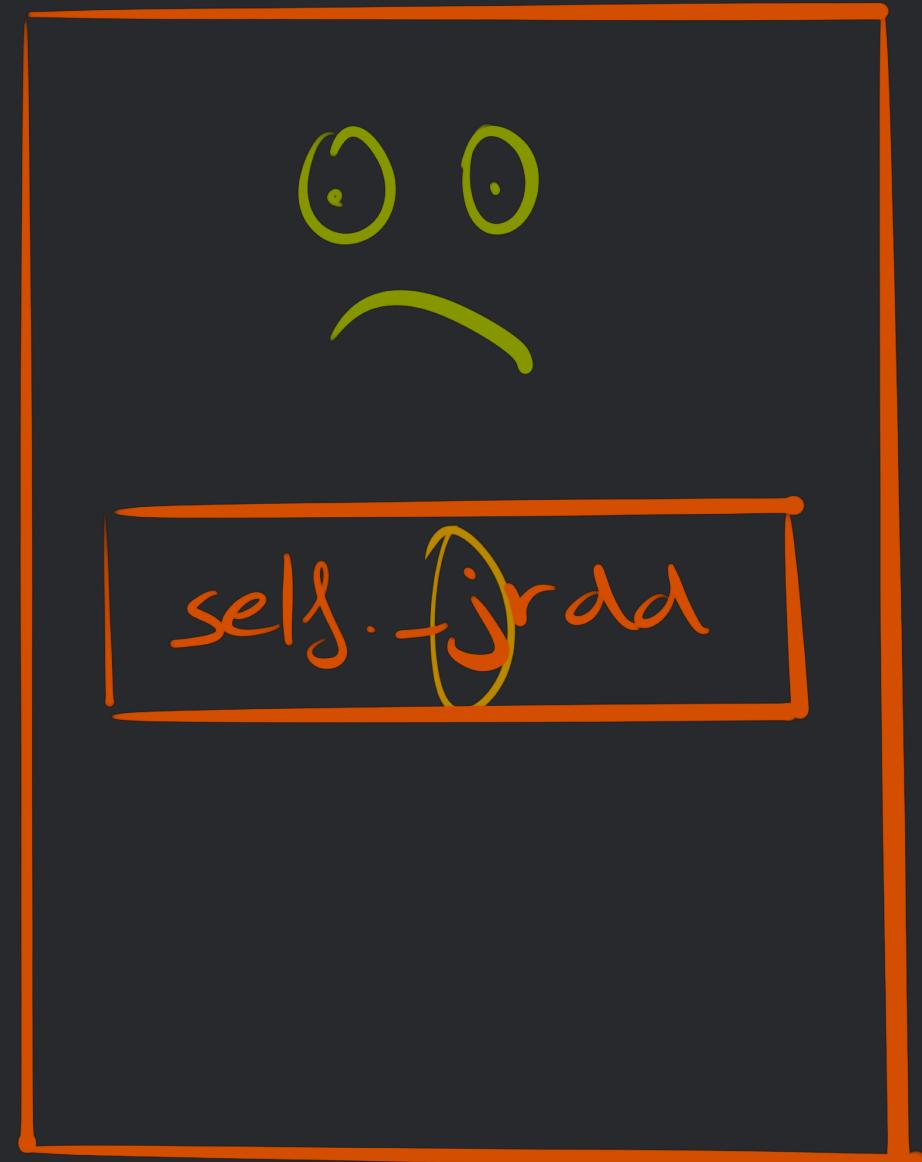


RDD in
Python Land



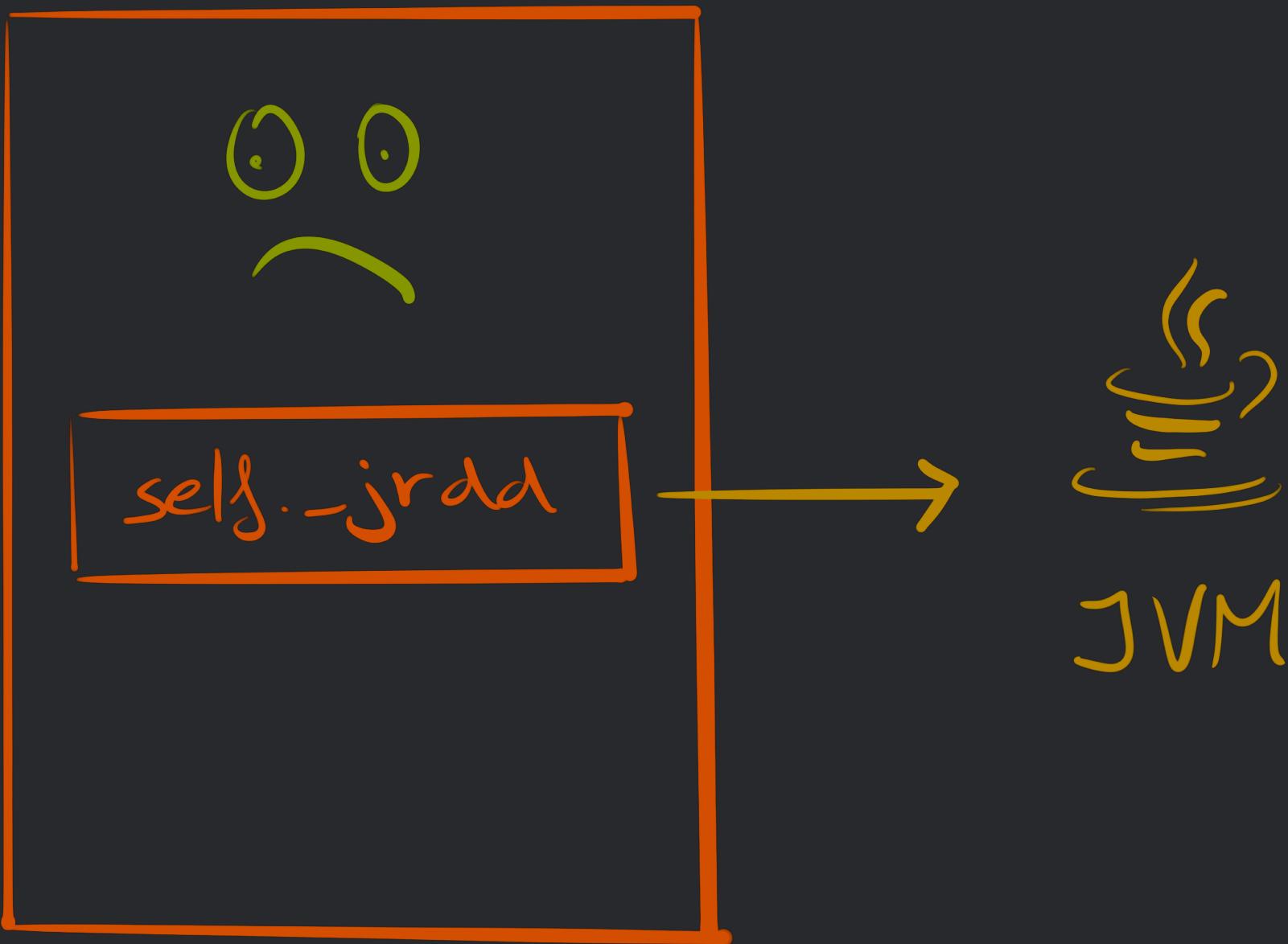
M
is for
murder

RDD in
Python Land



j
is for
java

RDD in
Python Land



RDD in
Python Land



Py4J bridge
JVM

THE MAIN ENTRYPOINTS
ARE RDD AND
PipelinedRDD(RDD)

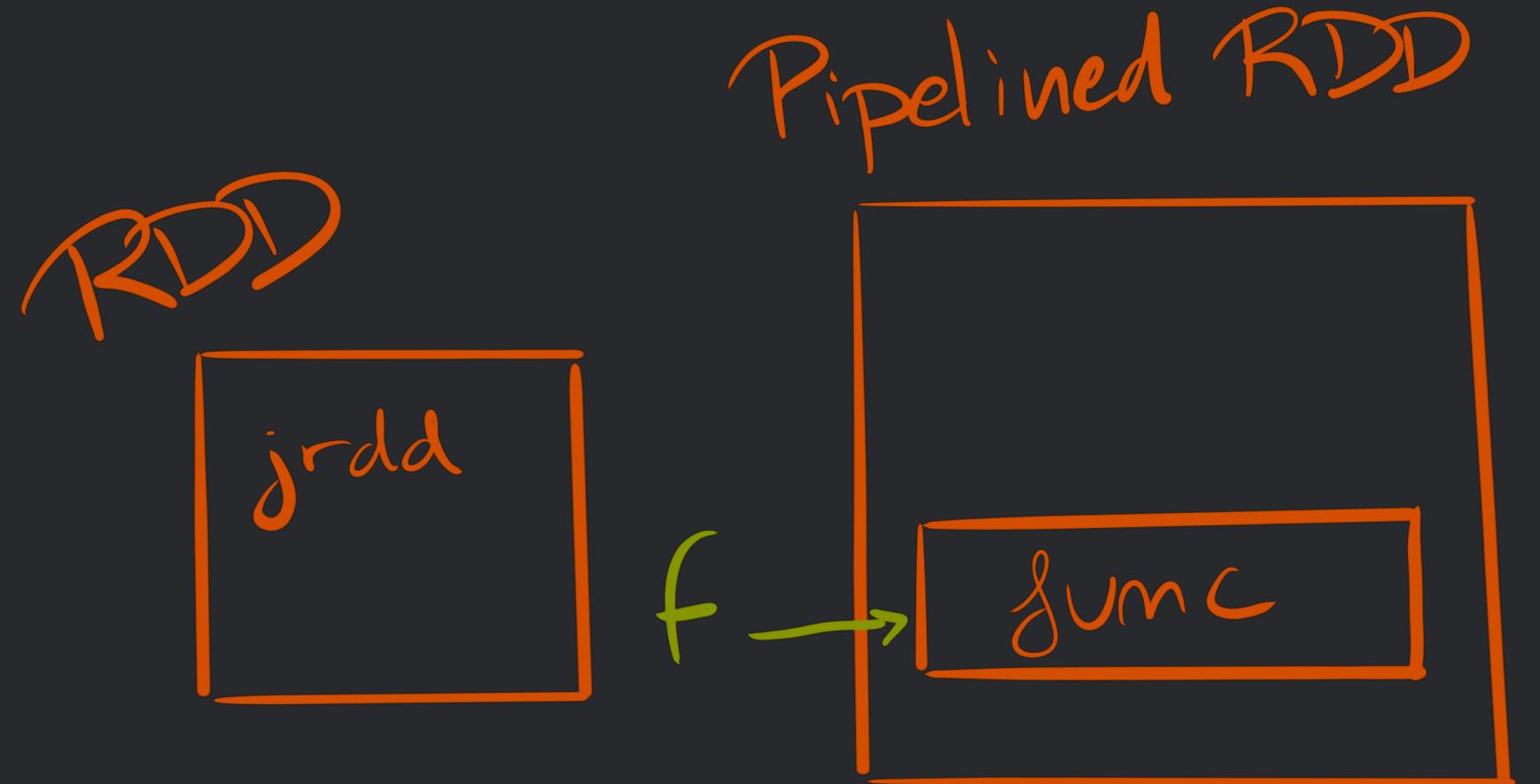
PipelinedRDD
BUILDS IN THE JVM A
PythonRDD

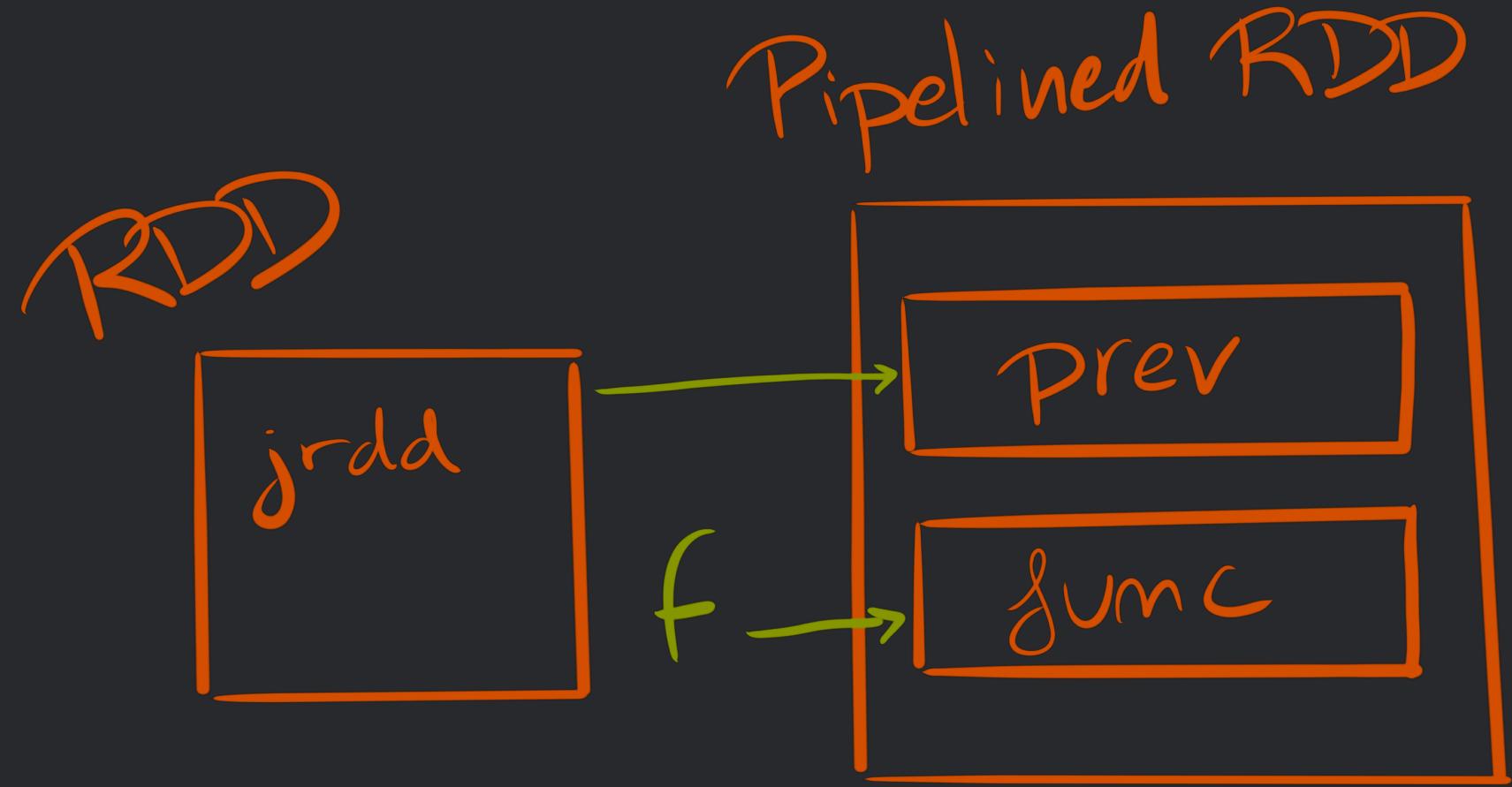


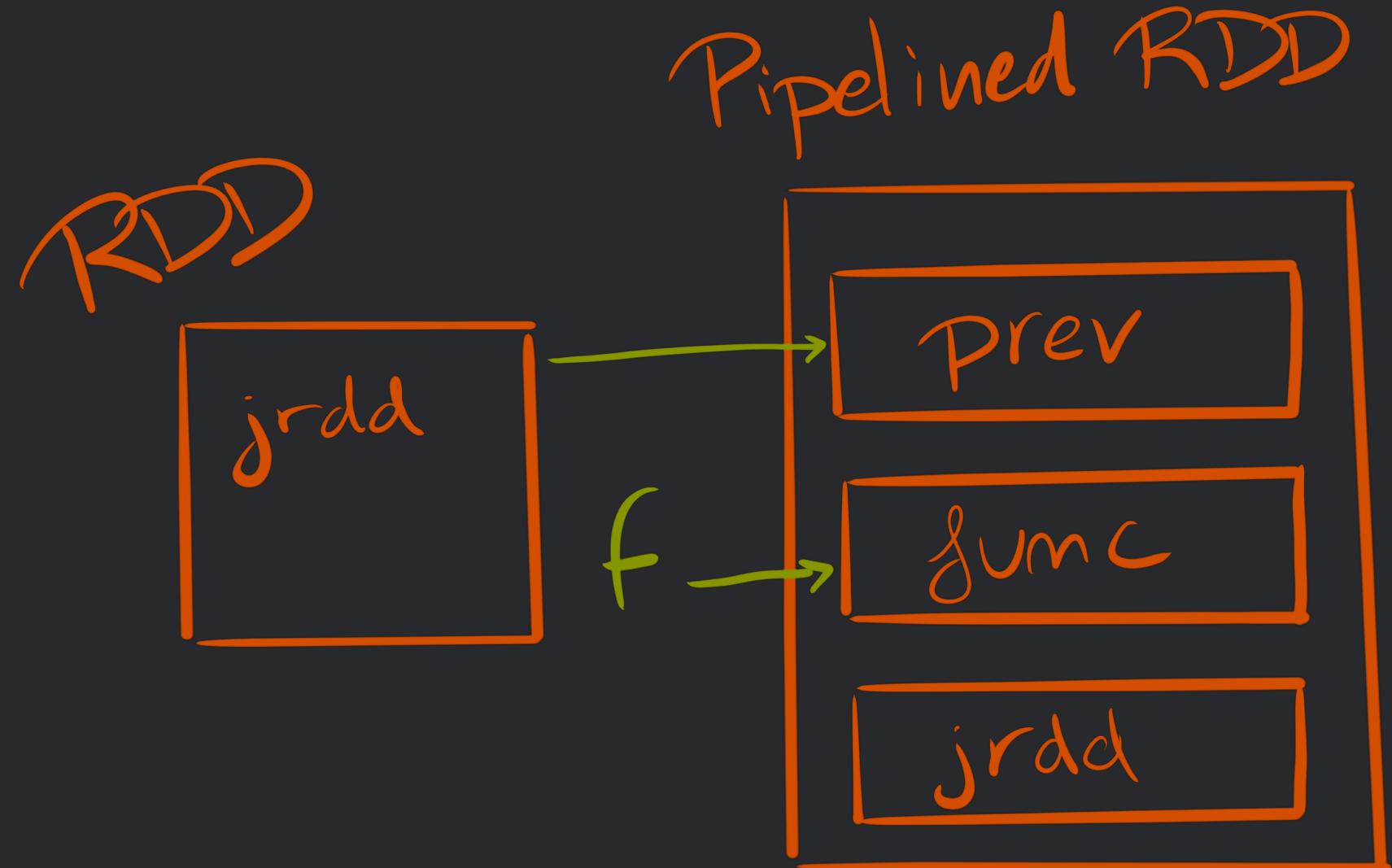
RDD

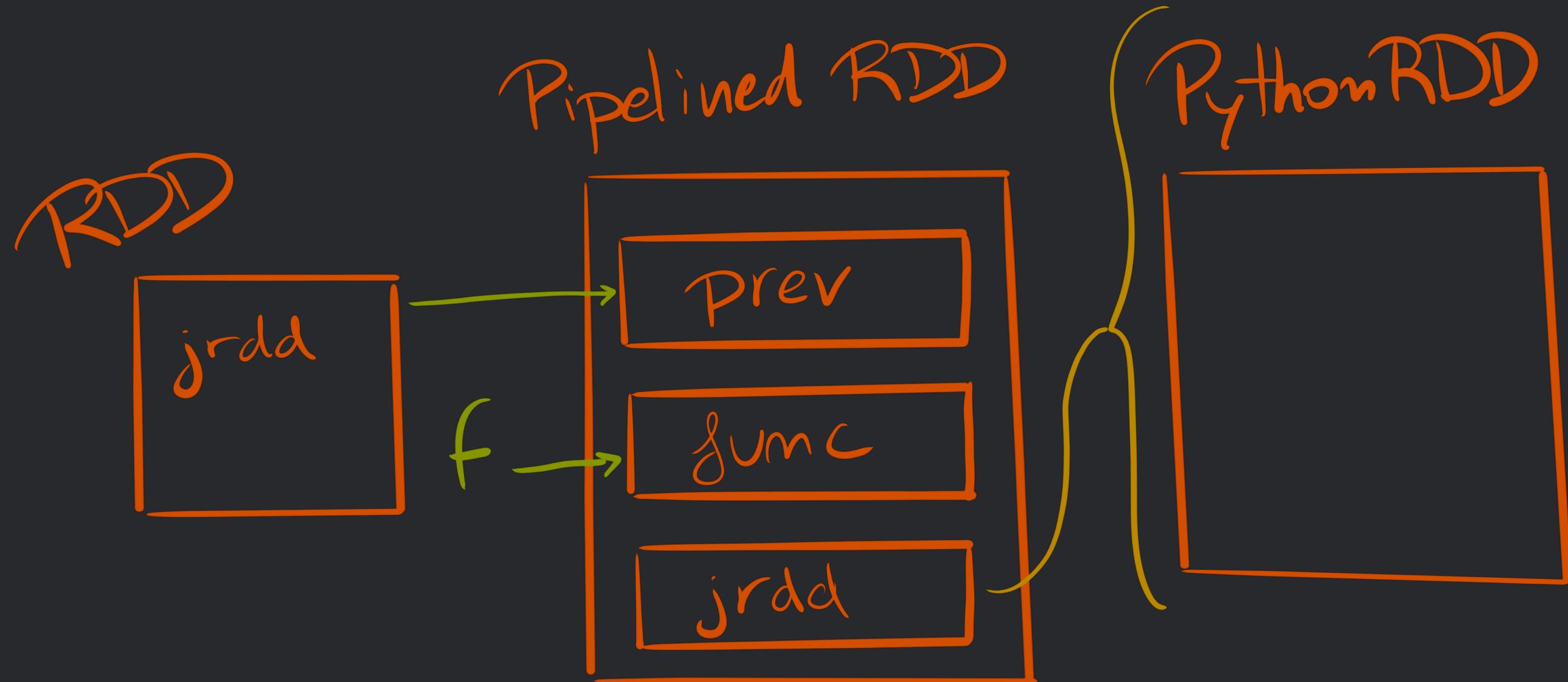
jRDD

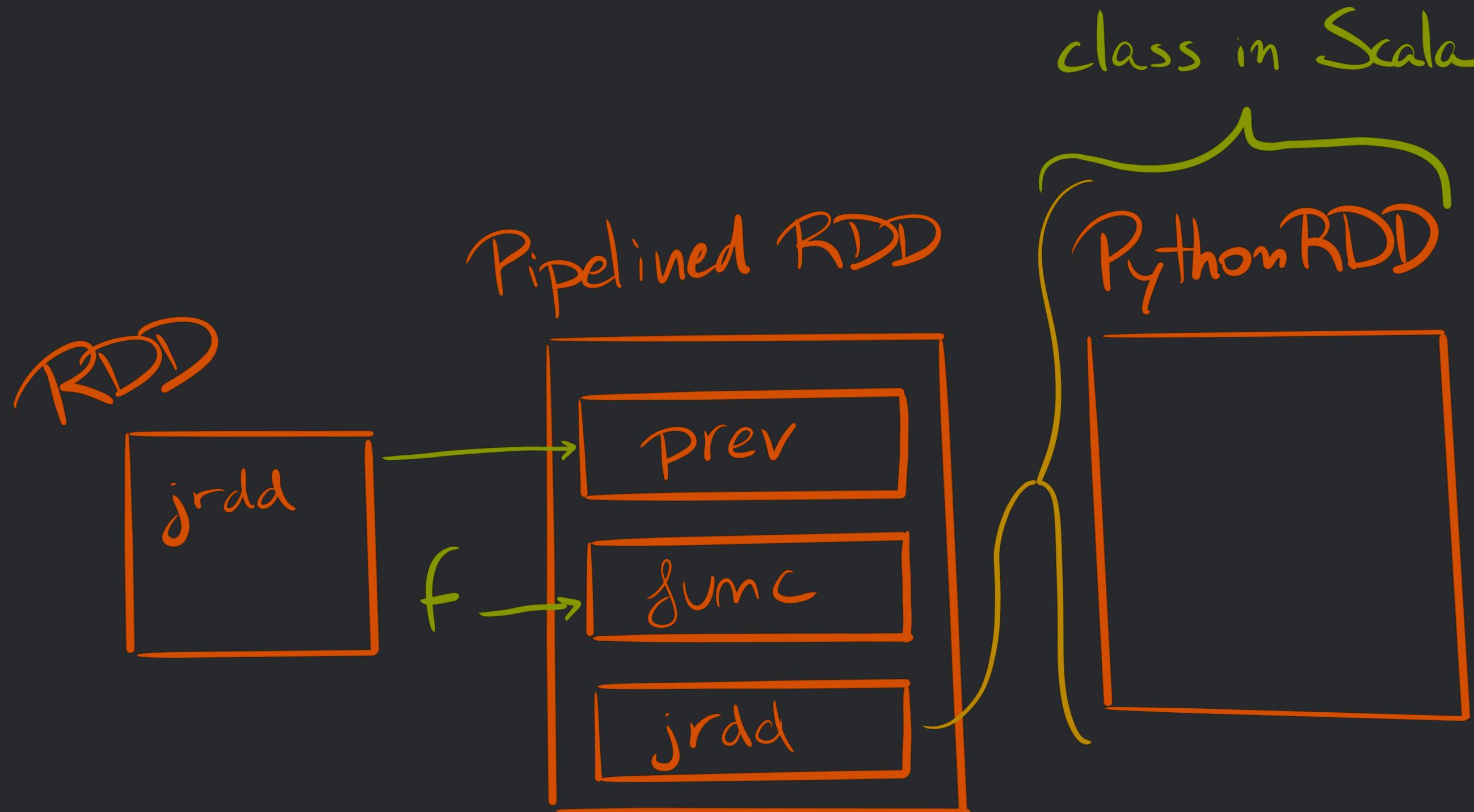
MAP(f)











class in Scala

PythonRDD

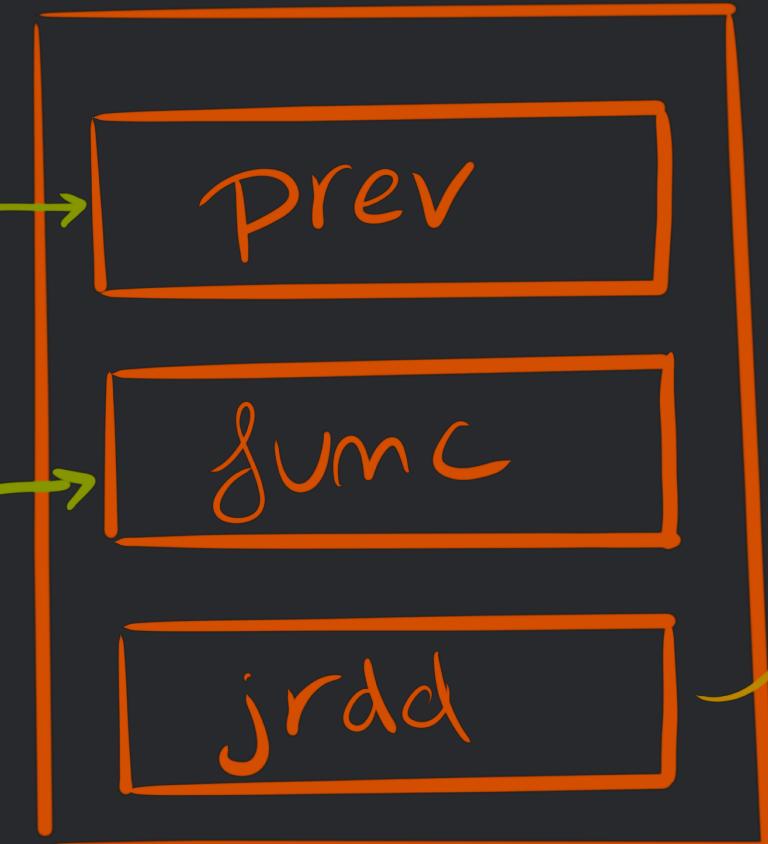
Dependencies

Pipelined RDD

RDD

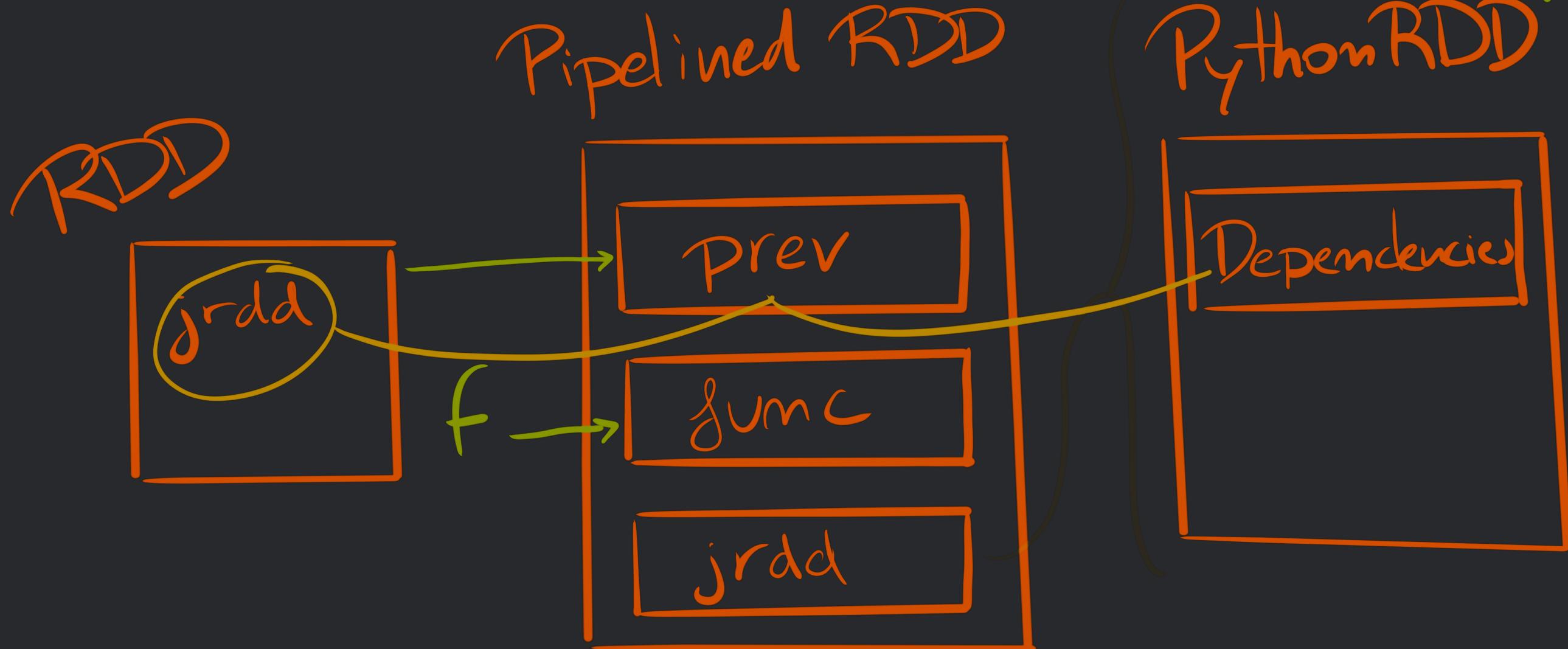
jrrdd

f



class in Scala

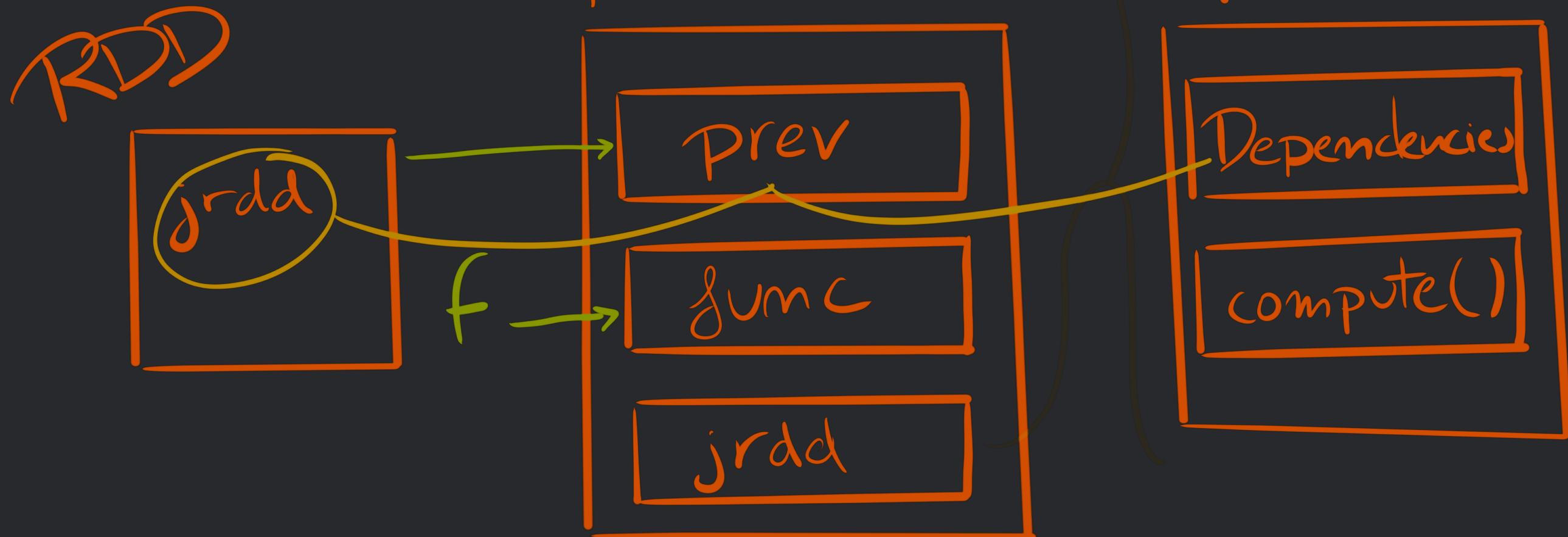
PythonRDD



class in Scala

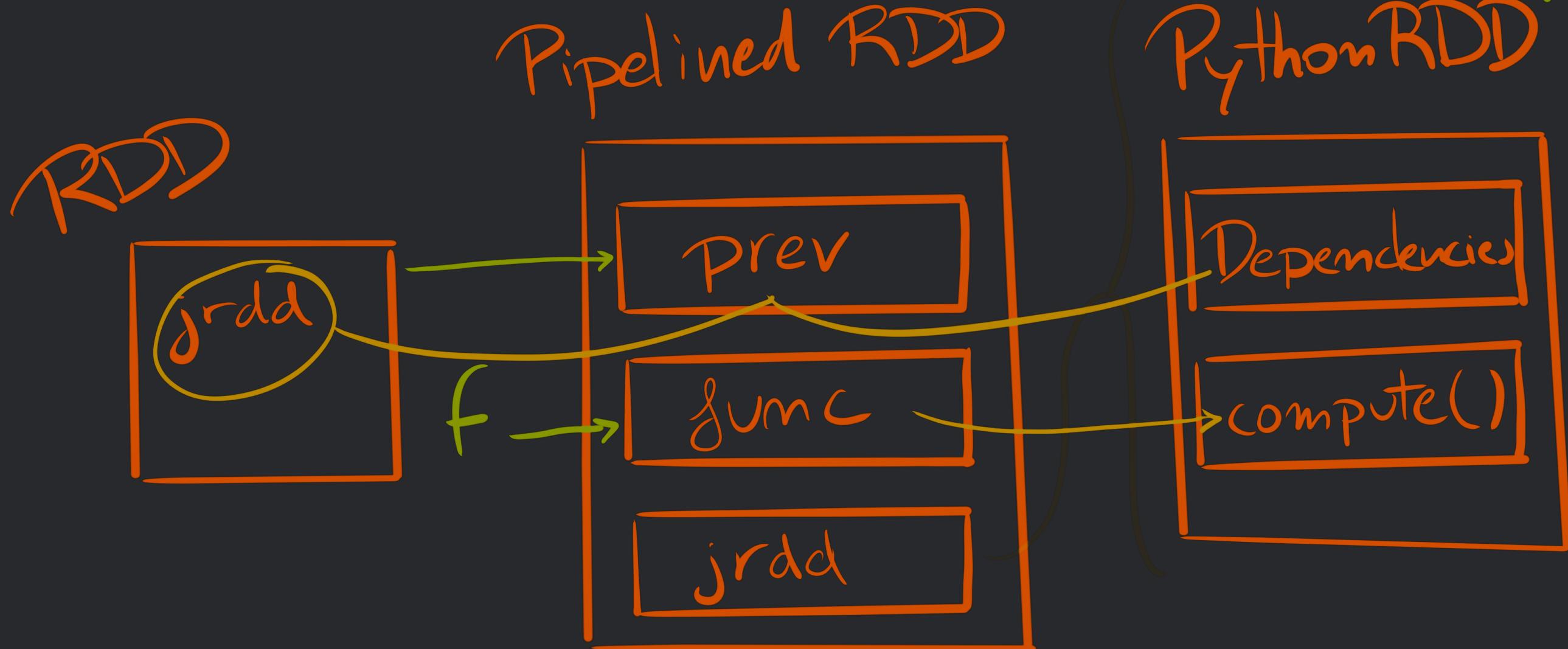
PythonRDD

Pipelined RDD



class in Scala

PythonRDD



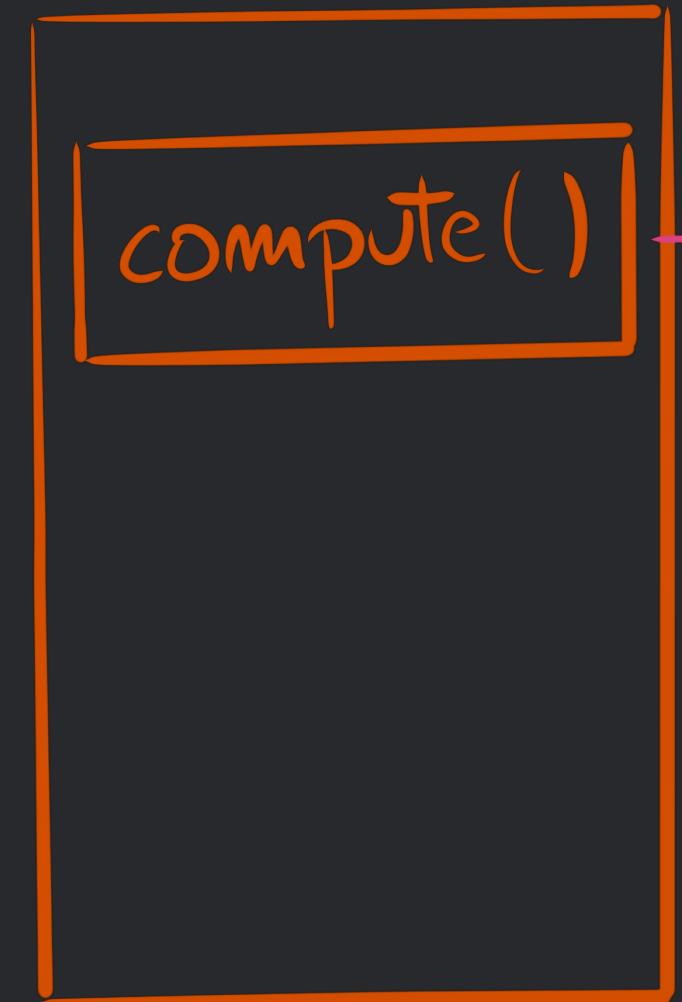


THE MAGIC IS
IN
compute

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



PythonRDD

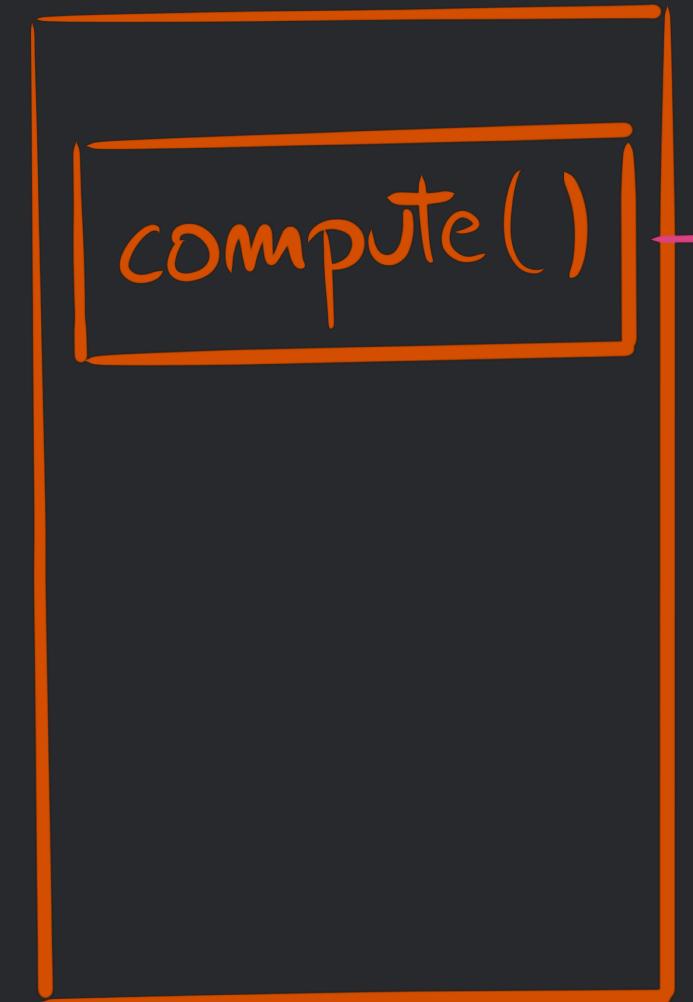


PythonRunner

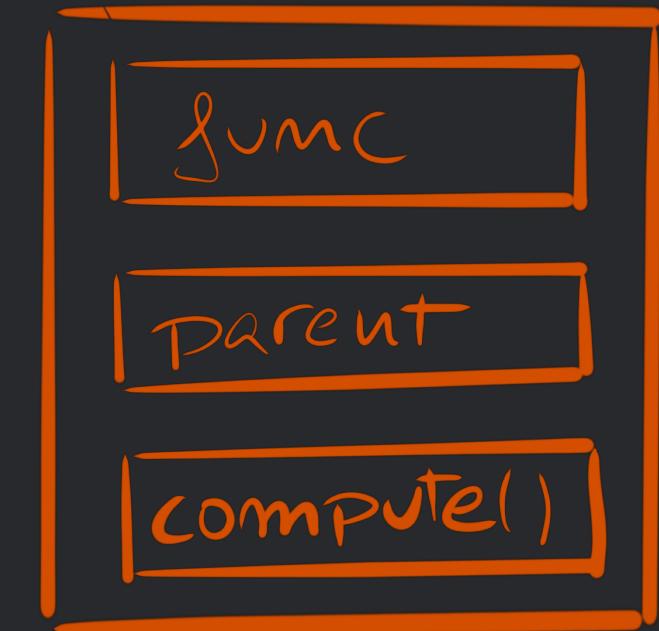


Scala!

PythonRDD



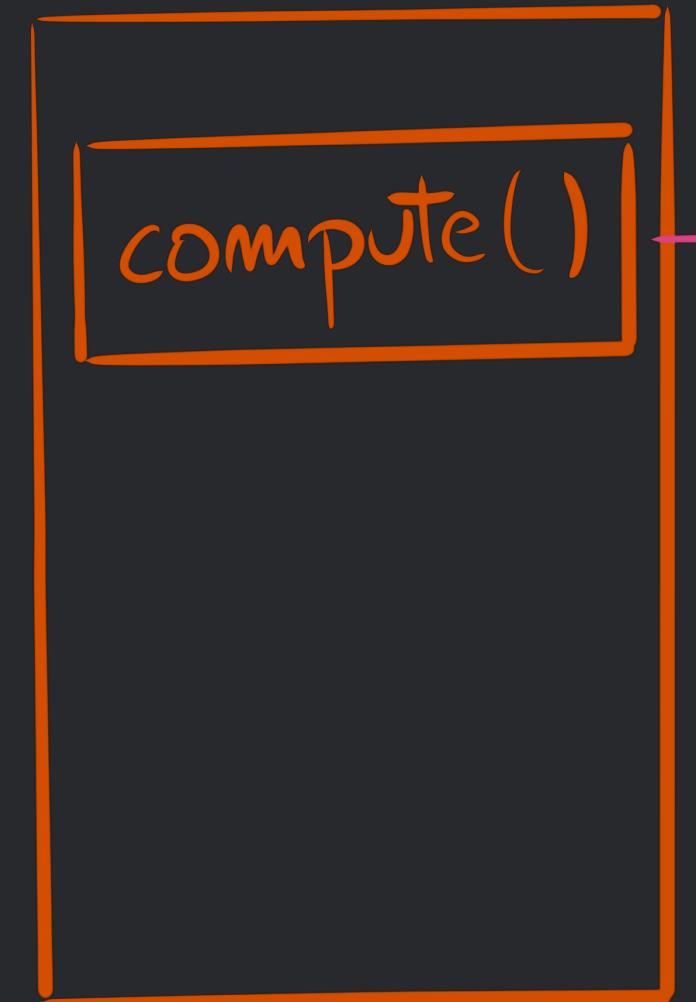
PythonRunner



Scala!

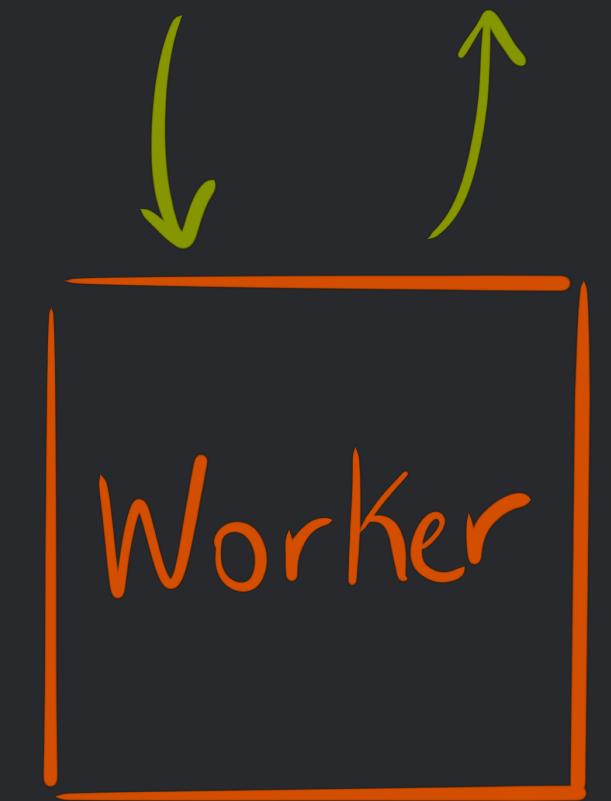
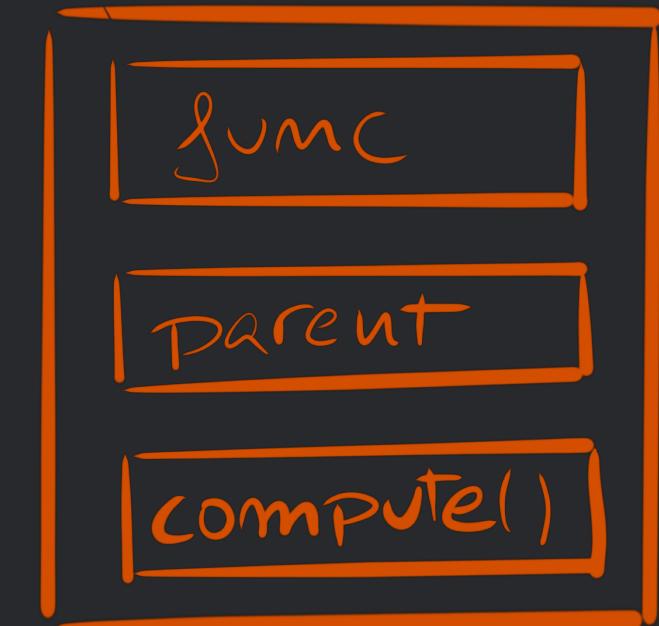
← Python!

PythonRDD



Scala!

PythonRunner



← 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

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

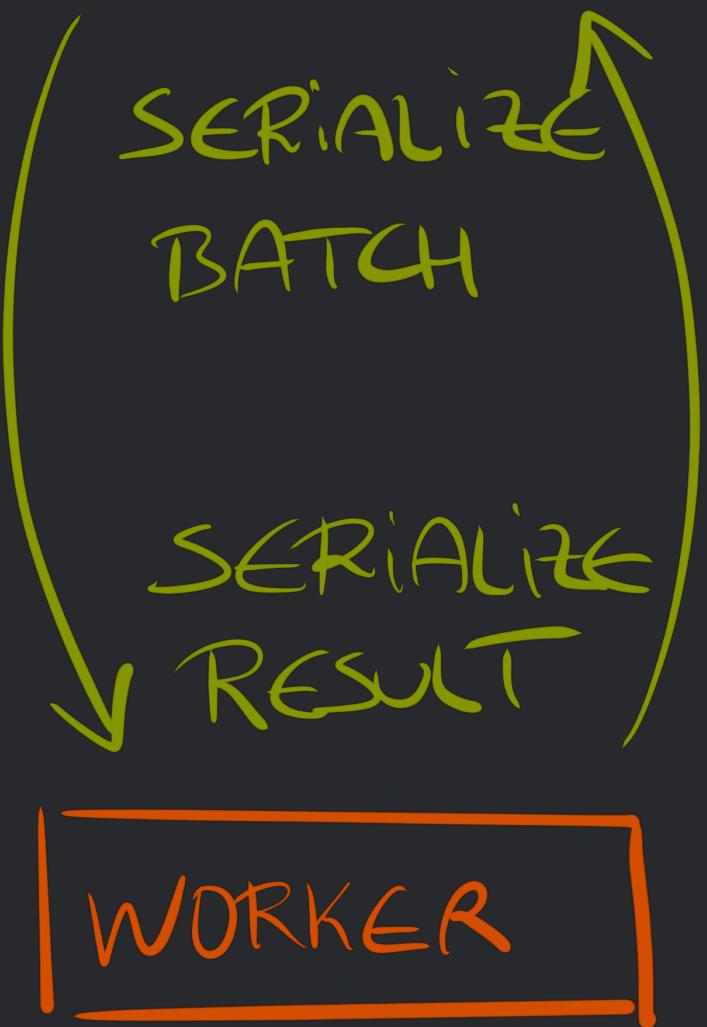
SERIALIZED
BATCH

SERIALIZED
RESULT

WORKER

UDF RUNNER

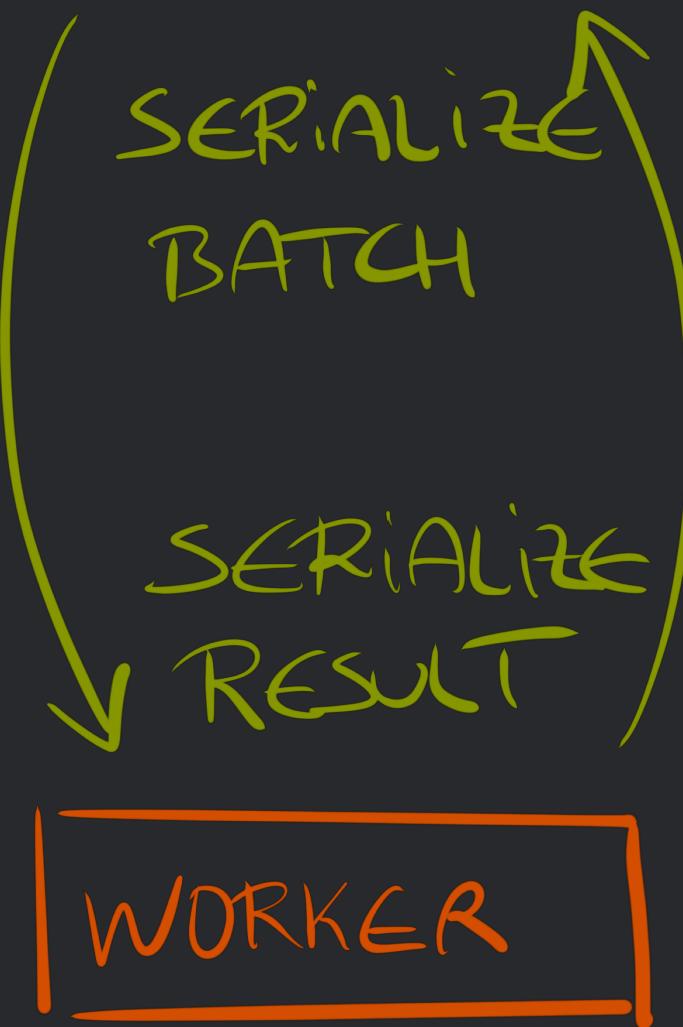
ARROW RUNNER



WORKER

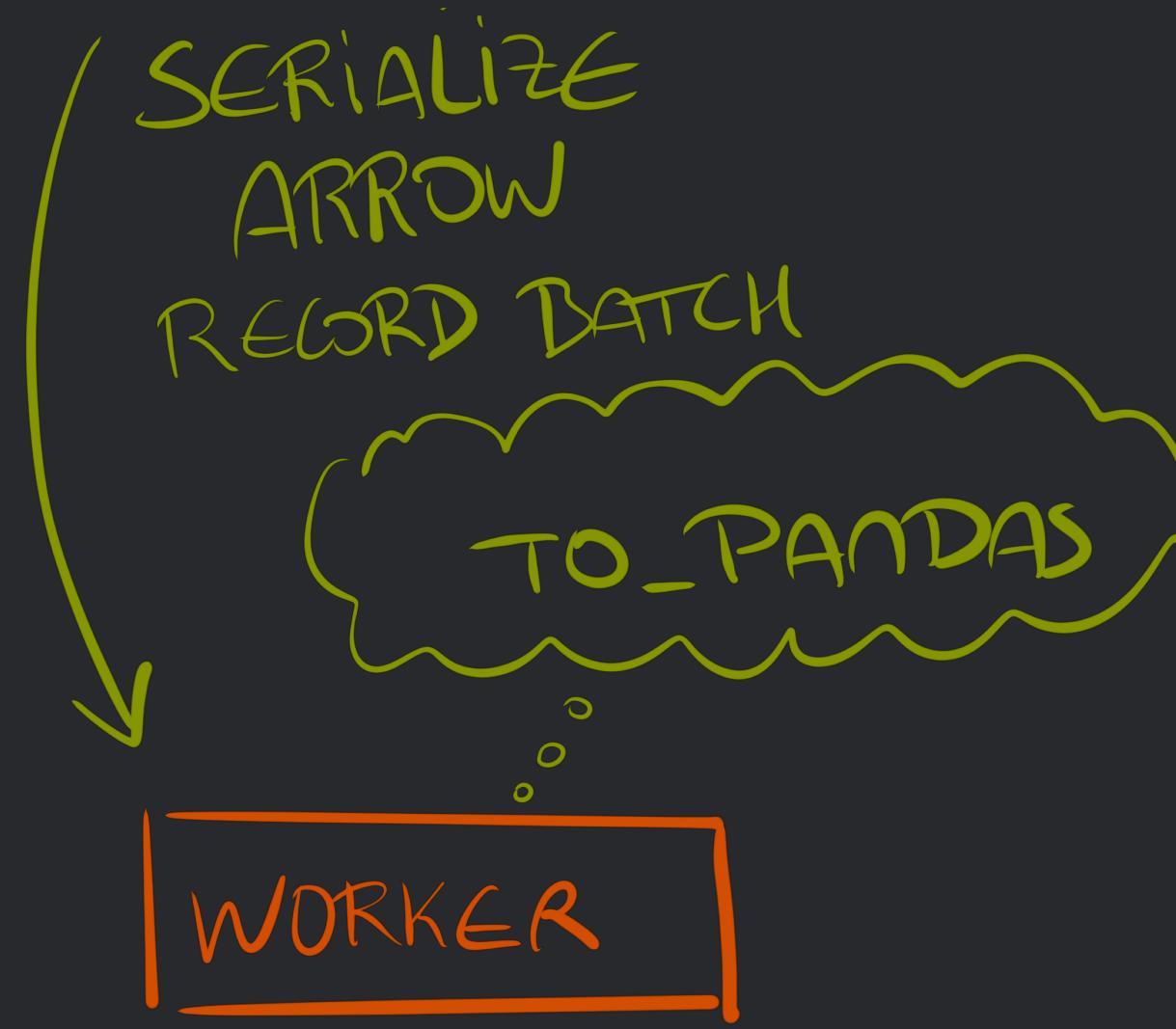
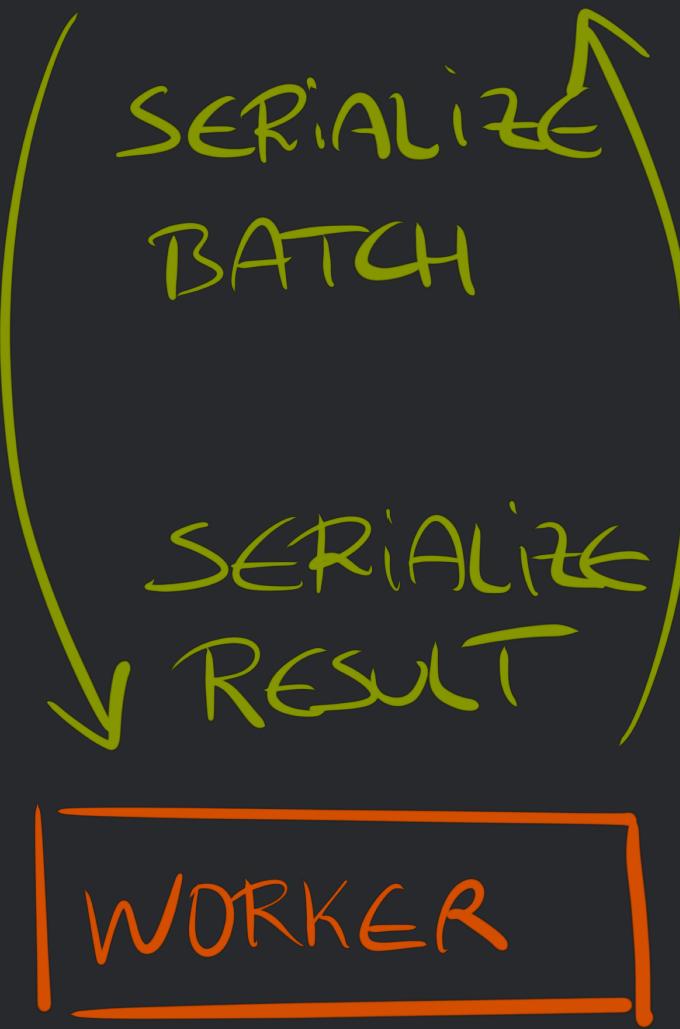
UDF RUNNER

ARROW RUNNER



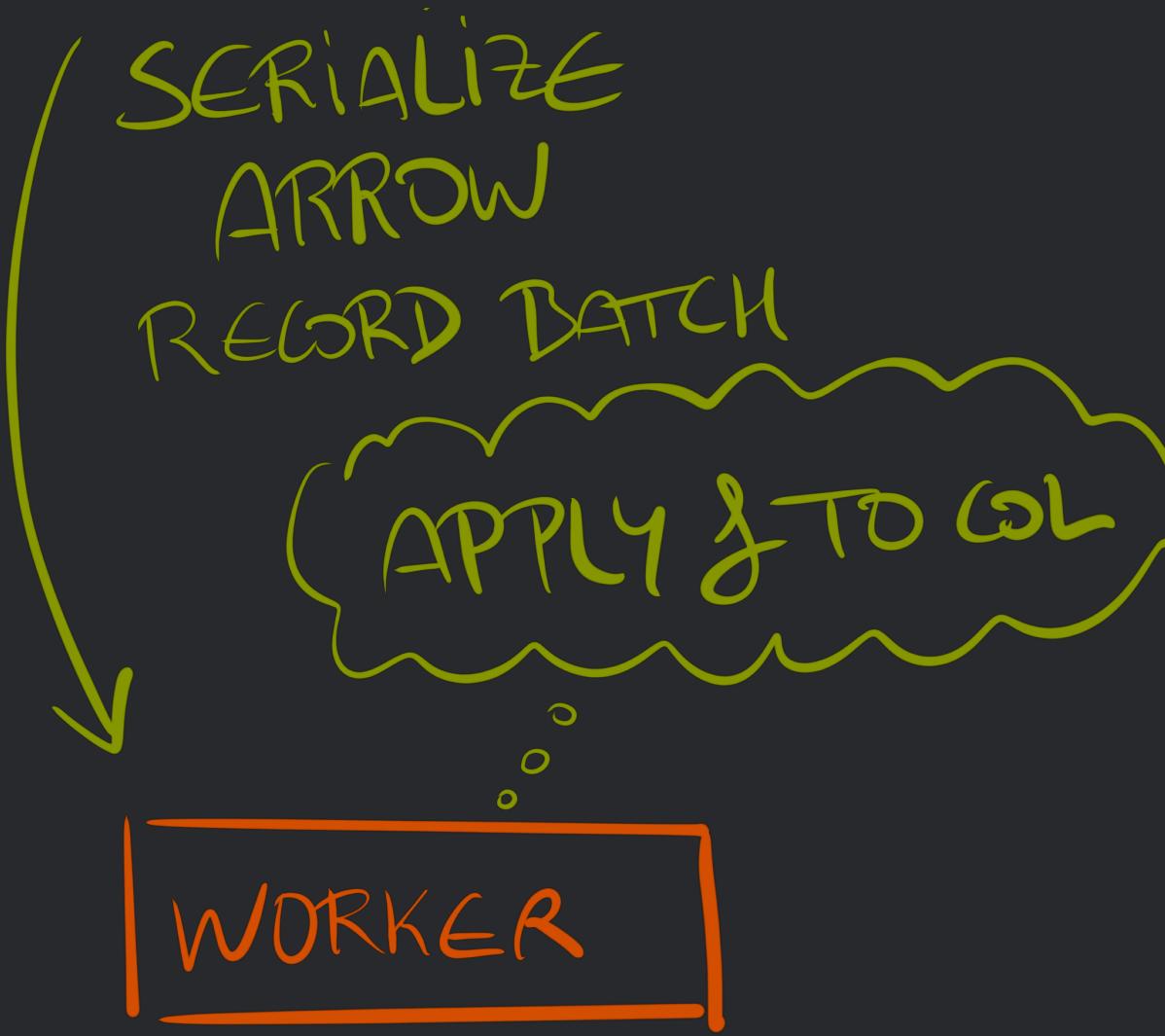
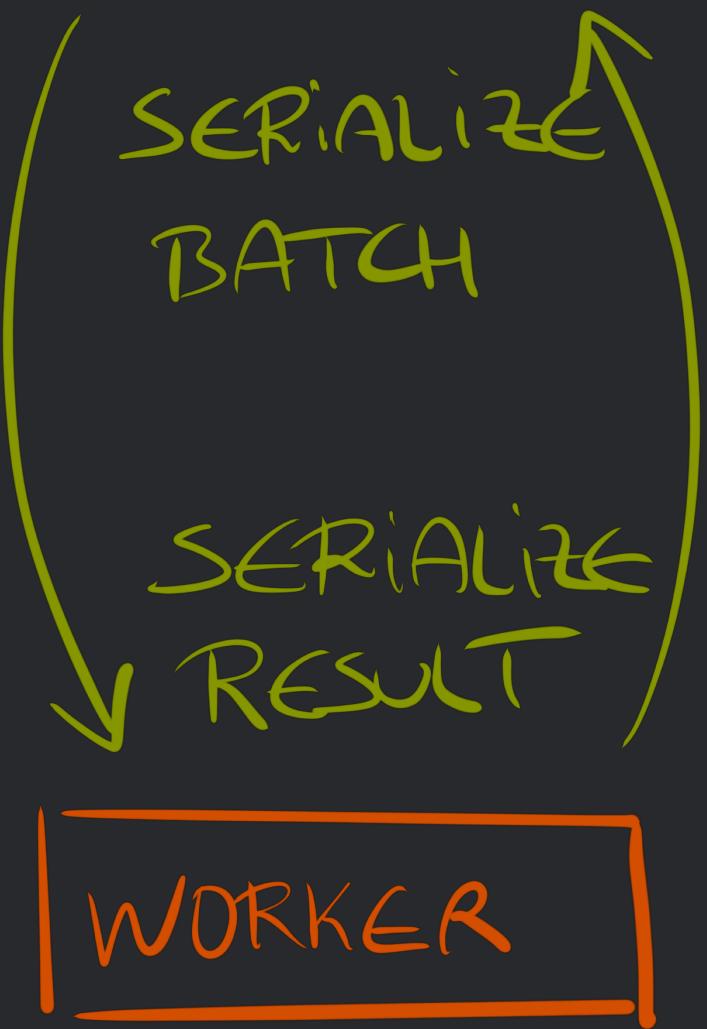
UDF RUNNER

ARROW RUNNER



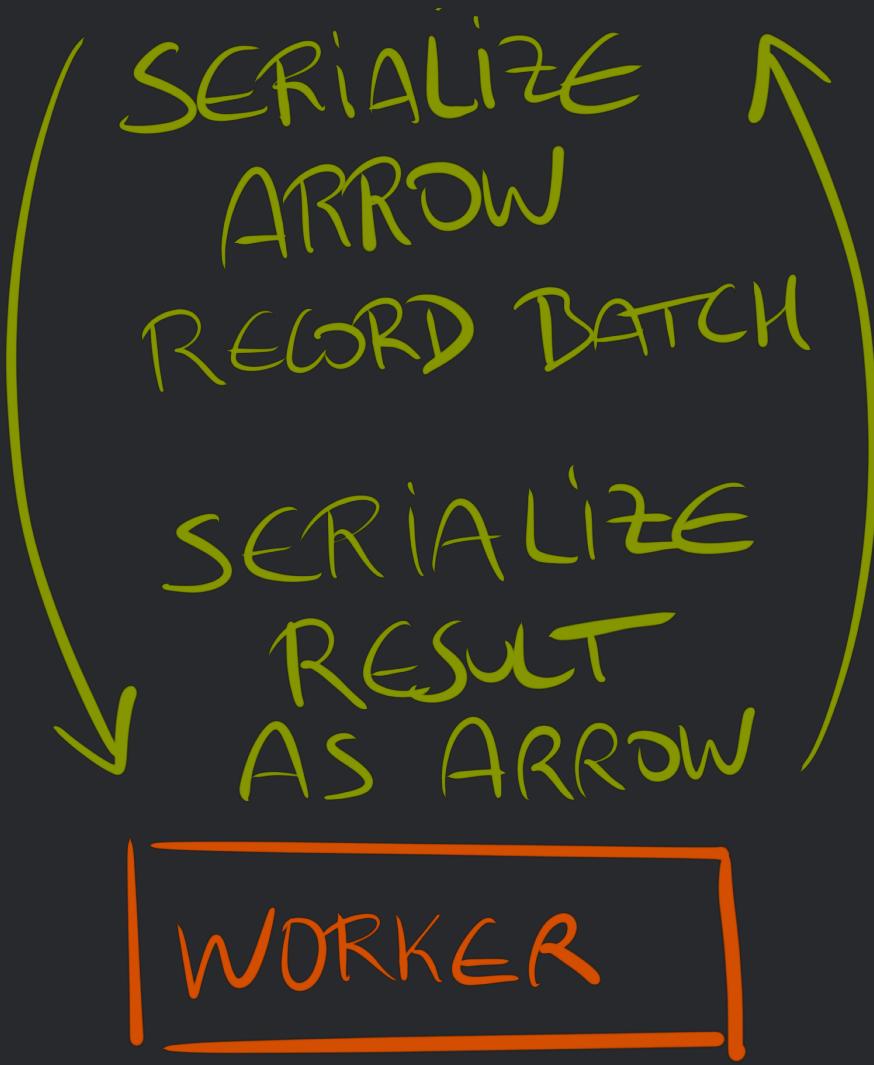
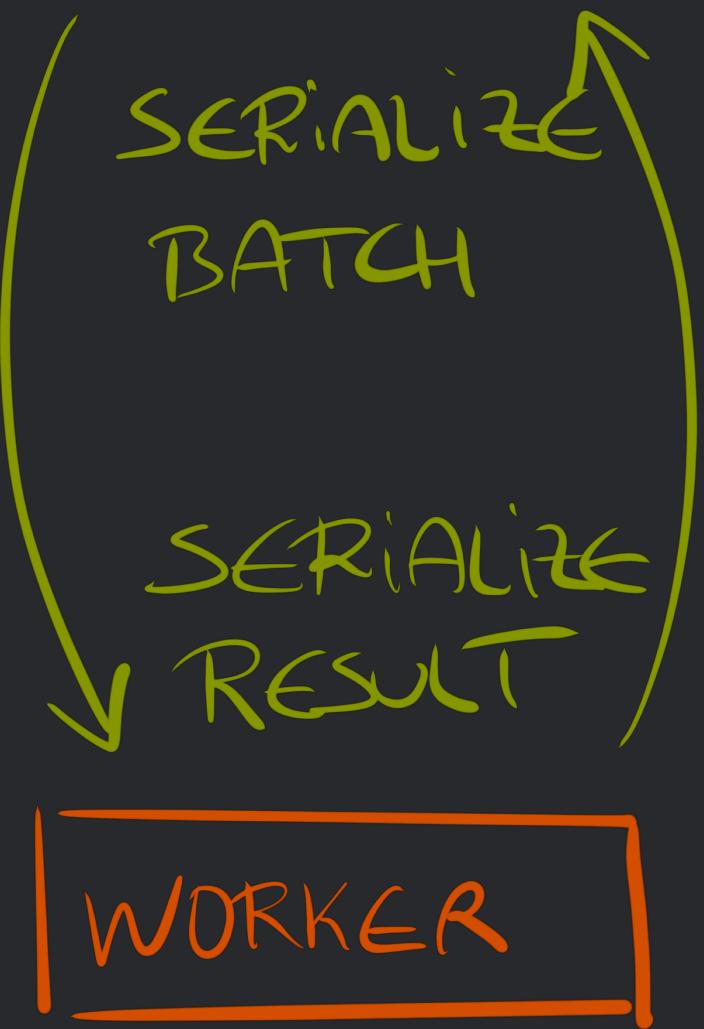
UDF RUNNER

ARROW RUNNER



UDF RUNNER

ARROW RUNNER



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

RESOURCES

- > [SPARK DOCUMENTATION](#)
- > [HIGH PERFORMANCE SPARK BY HOLDEN KARAU](#)
- > [MASTERING APACHE SPARK 2.3 BY JACEK LASKOWSKI](#)
 - > [SPARK'S GITHUB](#)
 - > [BECOME A CONTRIBUTOR](#)

QUESTIONS?



THANKS!

FURTHER REFERENCES

ARROW

ARROW'S HOME

ARROW'S GITHUB

ARROW SPEED TESTS

ARROW TO PANDAS CONVERSION SPEED

STREAMING COLUMNAR DATA WITH APACHE ARROW

WHY PANDAS USERS SHOULD BE EXCITED BY APACHE ARROW

ARROW-PANDAS COMPATIBILITY LAYER CODE

ARROW TABLE CODE

PYARROW IN-MEMORY DATA MODEL

PANDAS

PANDAS' HOME

PANDAS' GITHUB

IDIOMATIC PANDAS GUIDE

PANDAS INTERNALS CODE

PANDAS INTERNALS DESIGN

DEMYSTIFYING PANDAS' INTERNALS (TALK BY MARC GARCIA)

MEMORY LAYOUT OF MULTIDIMENSTIONAL ARRAYS (NUMPY)

SPARK/PYSPARK

PYSPARK SERIALIZERS CODE

FIRST STEPS TO USING ARROW (ONLY IN THE PYSPARK DRIVER)

SPEEDING UP PYSPARK WITH APACHE ARROW

ORIGINAL JIRA ISSUE: VECTORIZED UDFS IN SPARK

INITIAL DOC DRAFT

BLOG POST BY BRYAN CUTLER (LEADER FOR THE VEC UDFS PR)

INTRODUCING PANDAS UDF FOR PYSPARK

ORG.APACHE.SPARK.SQL.VECTORIZED

PY4J

PY4J'S HOME
PY4J'S GITHUB
REFLECTION ENGINE

EOF