



From 5 PB of logs to thousands of models - Story of a ML pipeline

Fabian Höring



ML related stuff
represents only
10% of your code



200 TB of data every day

A photograph of a server room with rows of server racks. The racks are illuminated with a blue light, and the floor is a light-colored tile with a grid pattern. The ceiling has recessed lighting fixtures. The text "Europe's largest Hadoop Cluster" is overlaid in the center of the image.

Europe's largest Hadoop Cluster

Given a description of a user and an item history build a model to predict how likely a user will be to click on the ad

YARN

APACHE
SparkTM

MLlib



HDFS

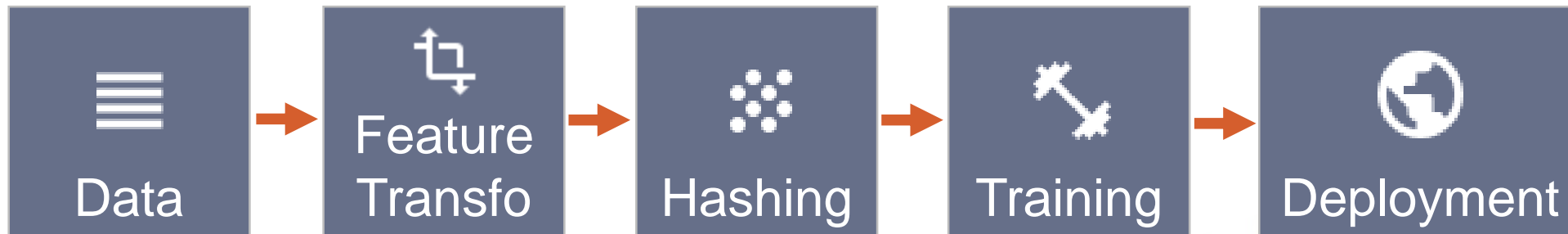


The image shows a close-up, low-angle view of several large, metallic industrial pipes. The pipes are arranged in a perspective that leads towards the background. The background is a bright, hazy sunset or sunrise, with the sun low on the horizon. In the distance, the silhouettes of industrial structures, including tall chimneys or towers, are visible against the bright sky. The overall color palette is dominated by warm, golden-yellow and orange tones. Two semi-transparent rectangular boxes are overlaid on the image. The first box, located in the upper left, is dark grey and contains the word "Offline" in white, bold, sans-serif font. The second box, located in the lower right, is dark blue and contains the word "Online" in white, bold, sans-serif font.

Offline

Online

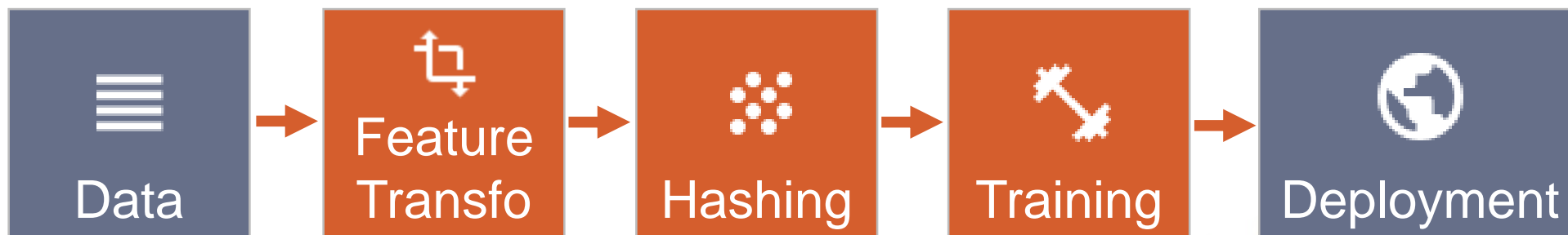
Offline Learning Workflow for 1 model



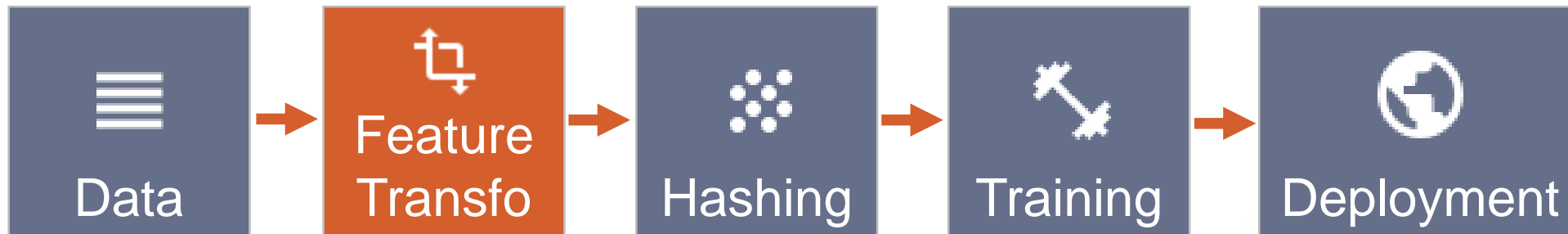


Preprocessing

Learning



Feature transformations



“The codepaths that actually generate input features may differ for training and inference time .. This is sometimes called “training/serving skew” and requires careful monitoring to detect and avoid”

— Google paper ML test score

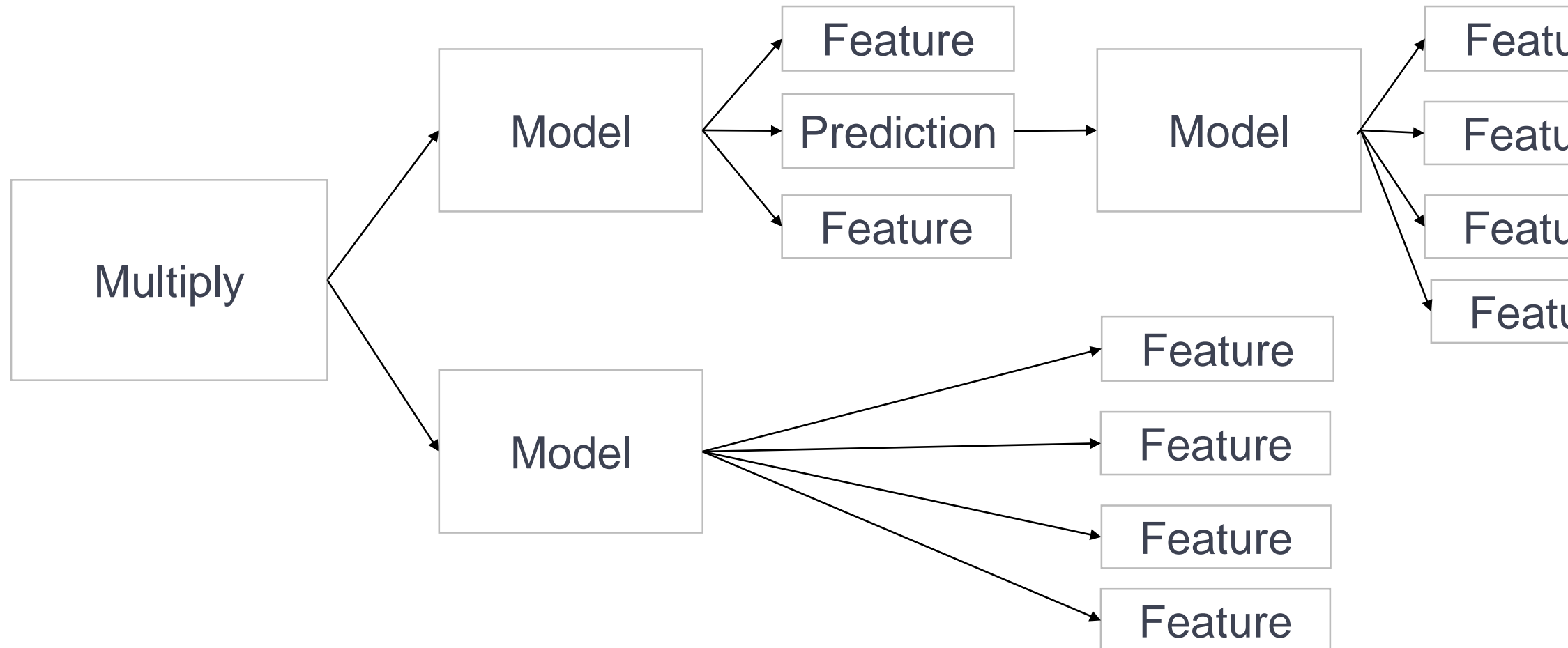
C# vs JVM



Invoking c# process with Spark pipe function

Feature Transformations

Everything is an ensemble



Hashing the data



Dataset is several TB big



Number of possible modalities
8 Billion



Encoding Categories

Feature User's Country

User's Country
France
Germany
Germany
..
Japan



Training examples

Encoding as Integer

Feature User's Country

User's Country
France = 1
Germany = 2
Germany = 2
..
Japan = 100

Encoding as Integer

Feature User's Country

User's Country
France = 1
Germany = 2
Germany = 2
..
Japan = 100



One hot encoding

Feature User's Country

User's Country
France
Germany
Germany
..
Japan



France	Germany	Japan	...
1	0	0	
0	1	0	
0	1	0	
	..		
0	0	1	



columns = number of distinct modalities

One hot encoding

Feature User's Country

France	Germany	Japan	Brasil	USA
1	0	0	0	0	..	0	0
0	1	0	0	0	..	0	0
0	1	0	0	0	..	0	0
	..						
0	0	1	0	0	..	0	0

190 columns

Hashing trick

Hash1	Hash2	Hash3	Hash4	Hash5	Hash6
1	0	0	1	0	0



France

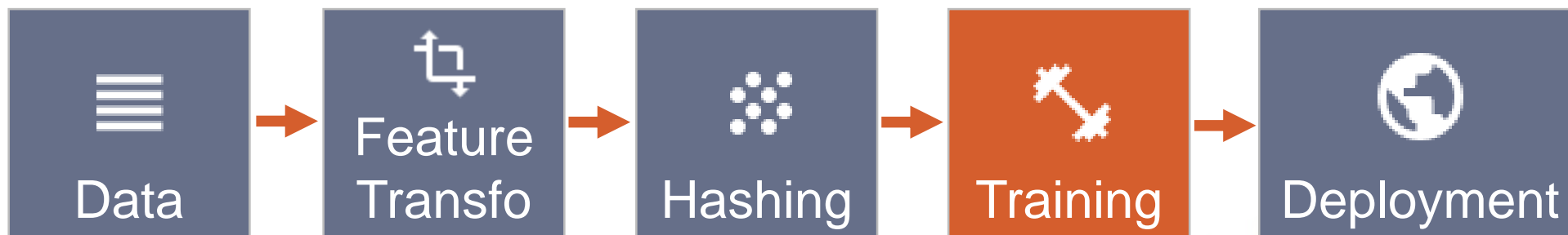
Chrome
Browser



Red Shoes

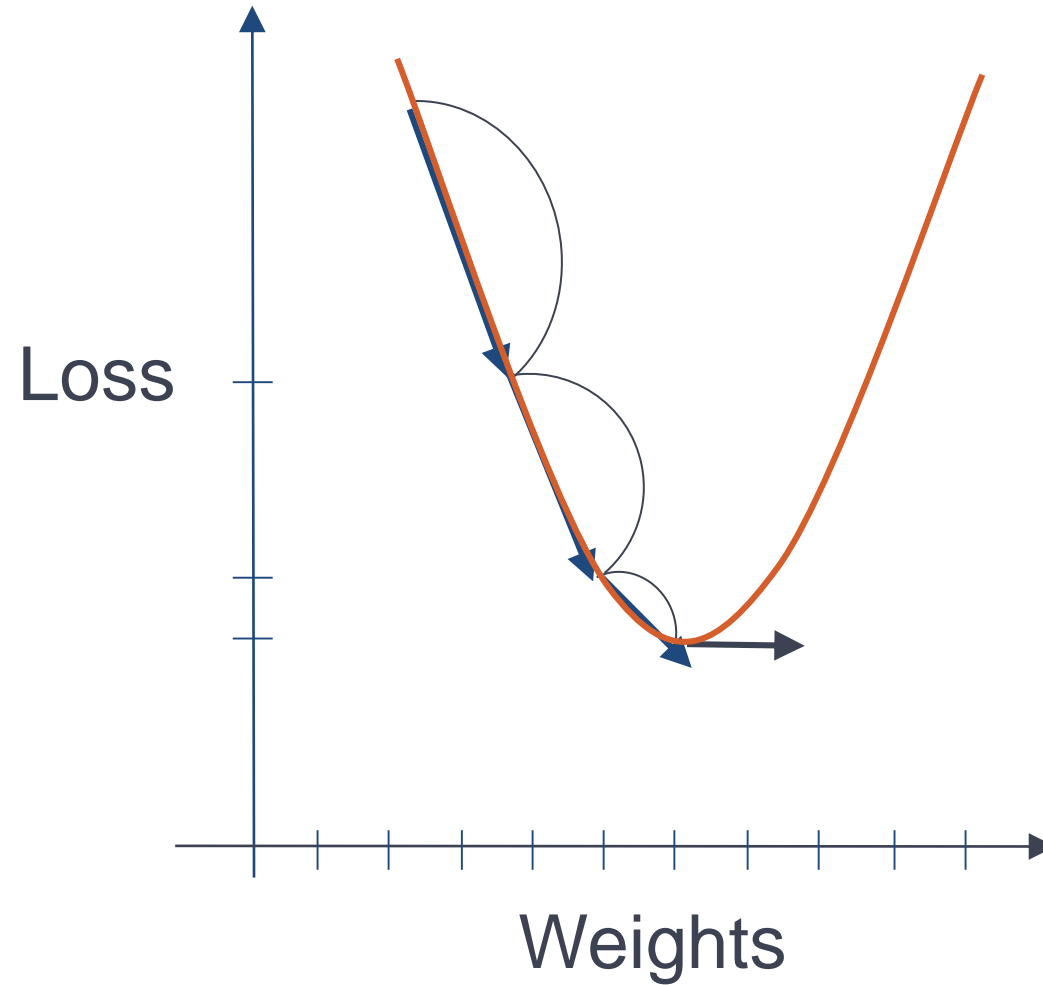
With the Hashing trick
one feature vector has
68 Mio values
& takes **256 MB** in
memory

APACHE
Spark[™]
Learning

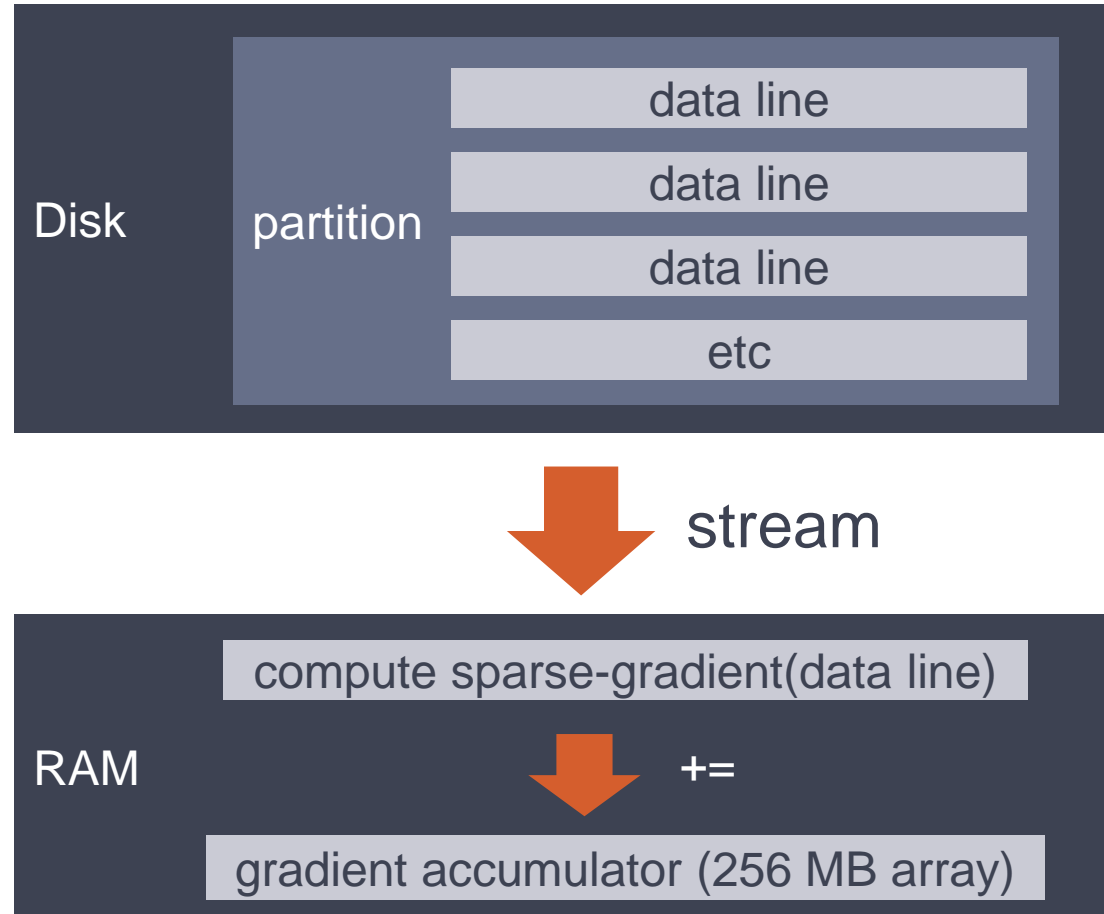


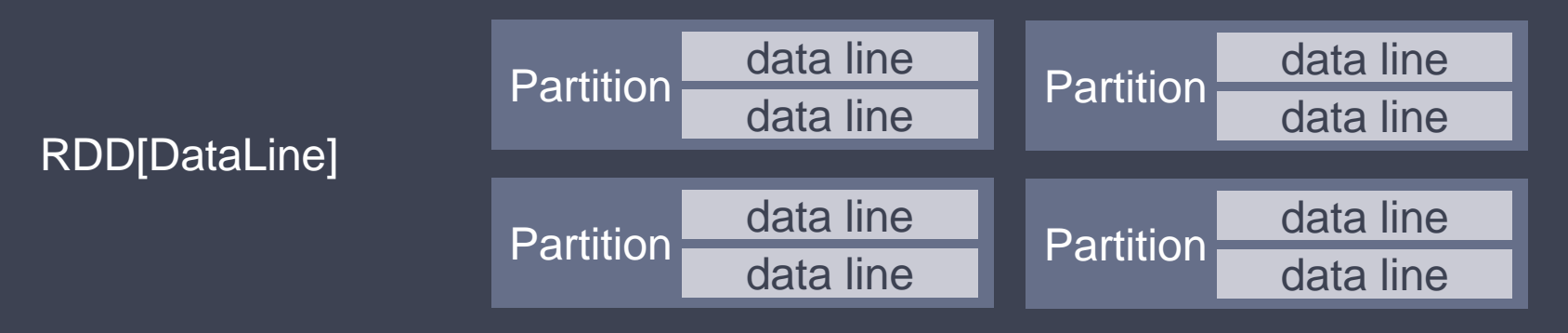
Logistic Regression

Gradient descent

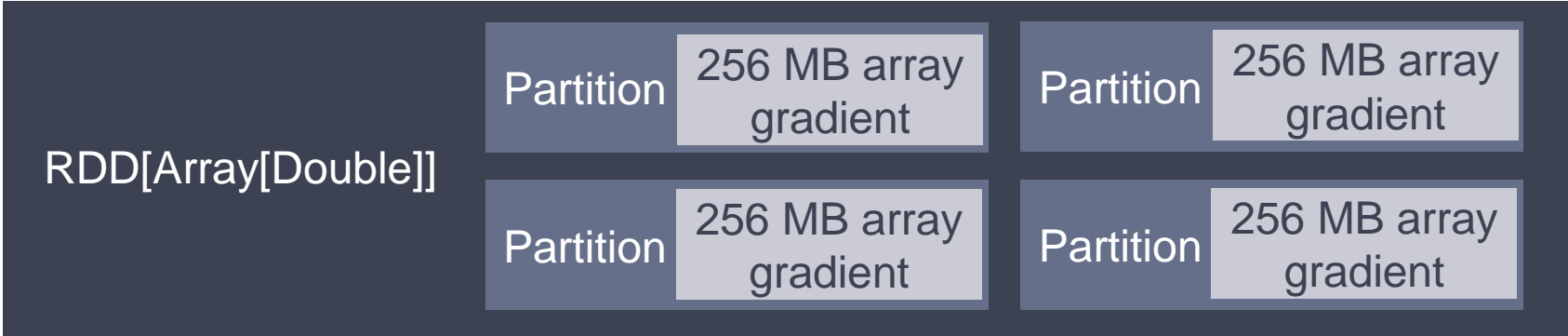
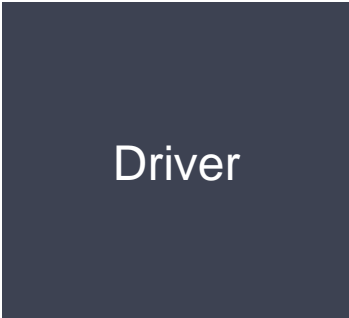


Computation for each partition



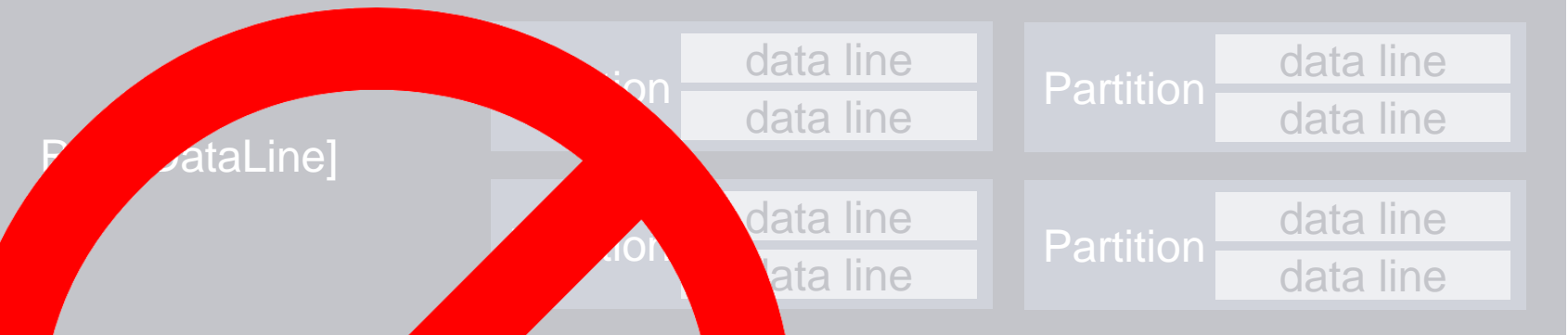


rdd.mapPartitions(...)

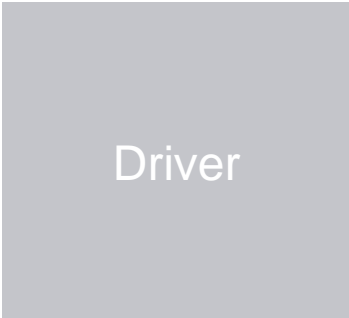
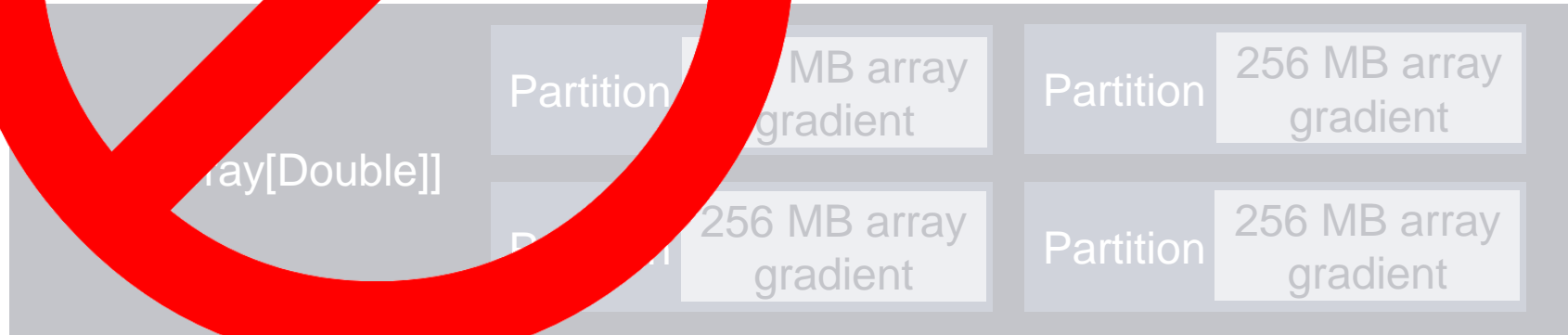


reduce(x,y => x + y)



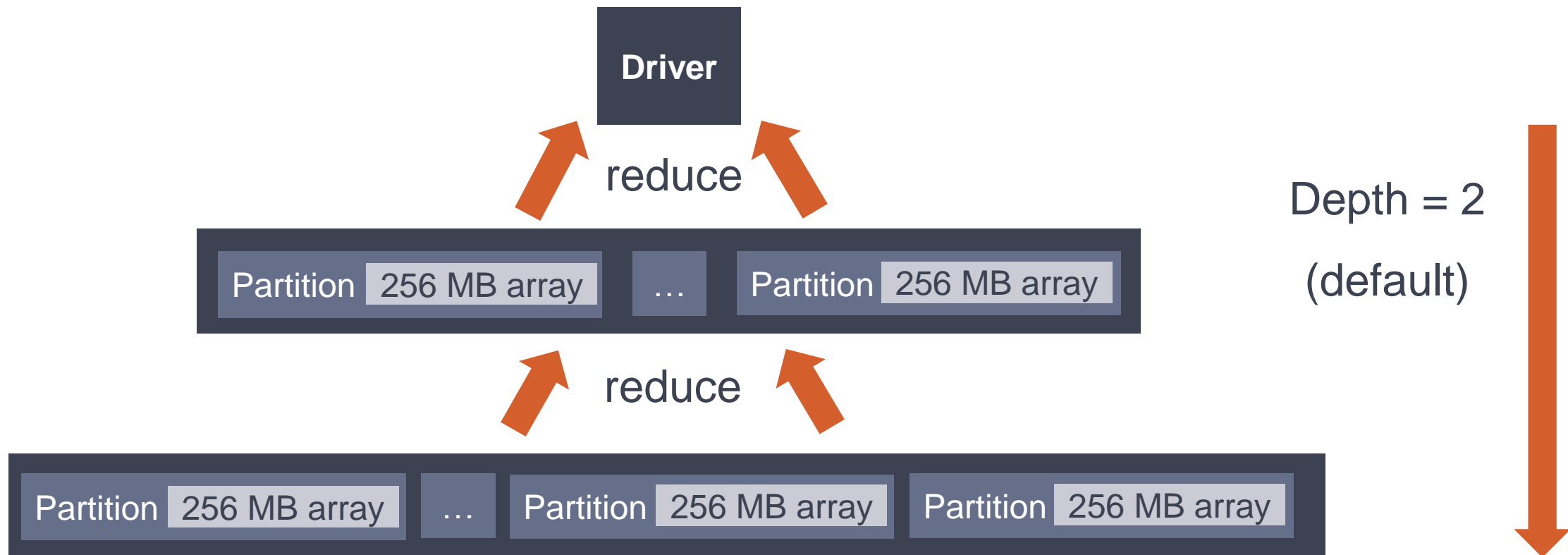


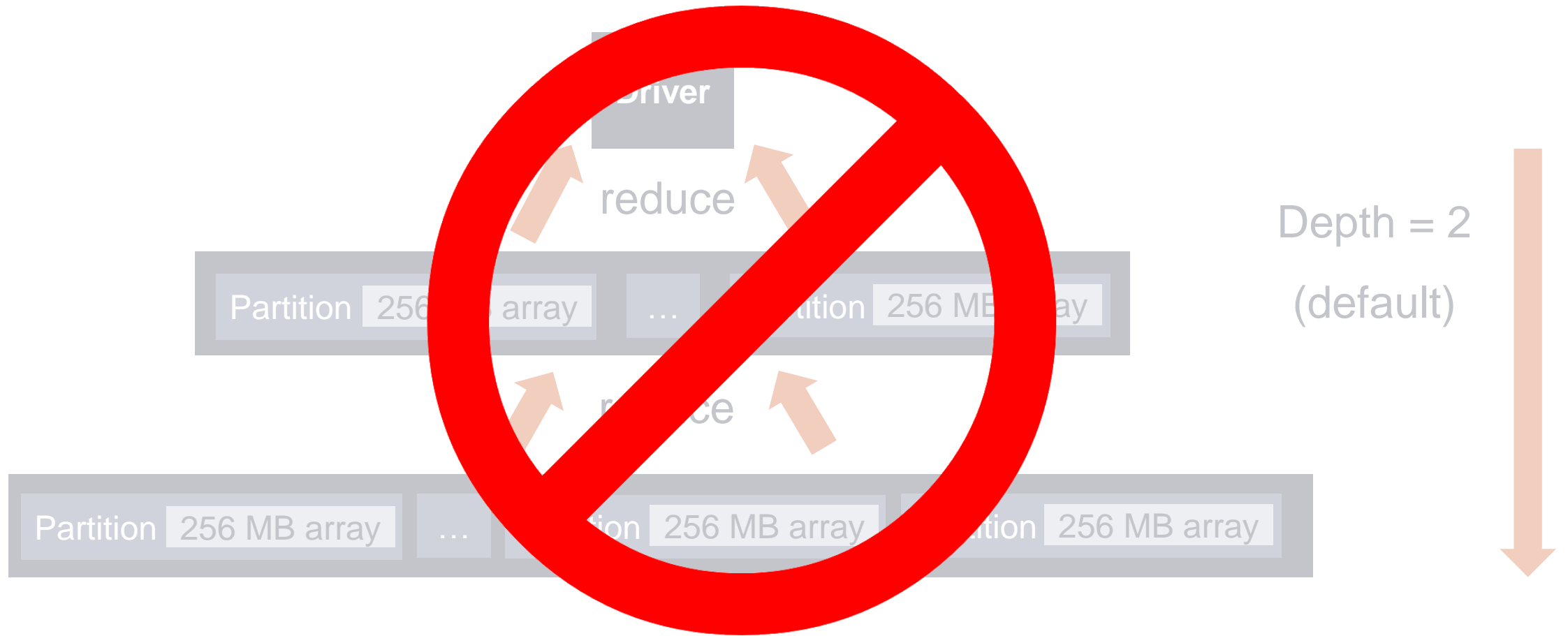
`sc.mapPartitions(...)`

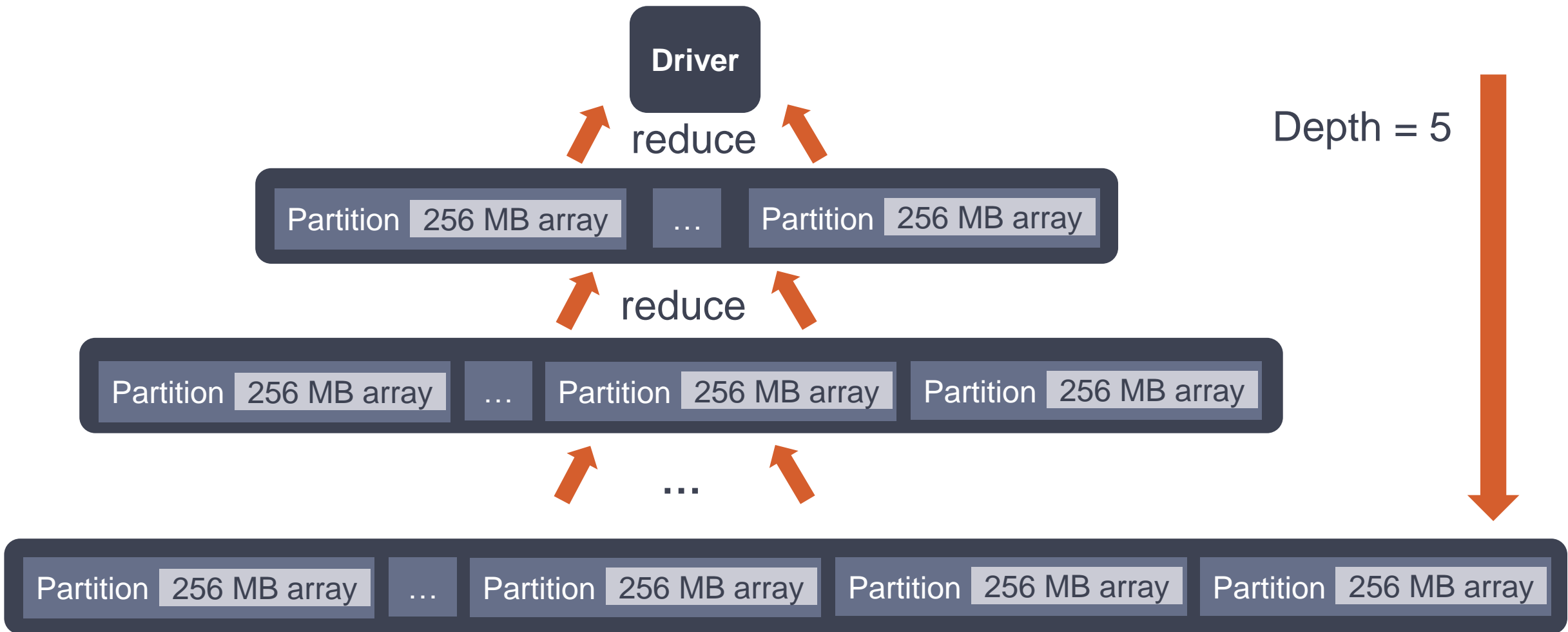


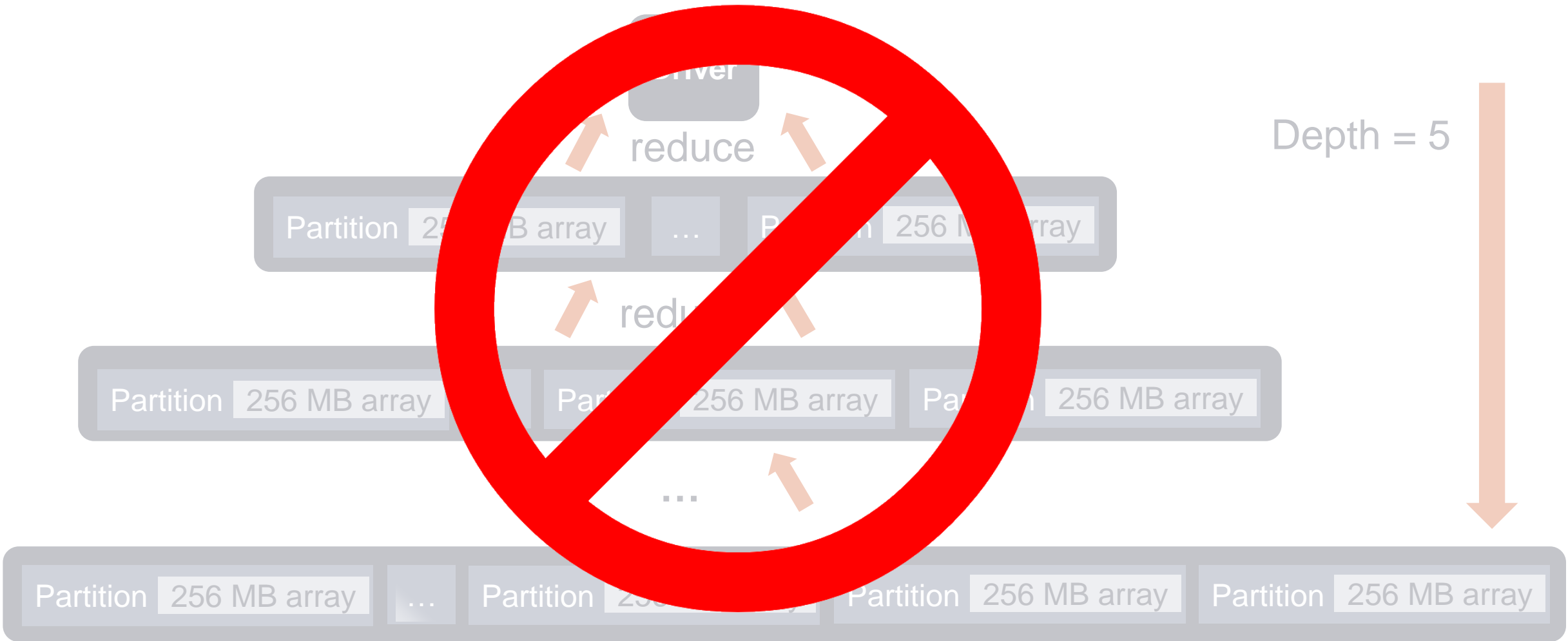
`reduce(x,y => x + y)`











Depth = 5

Partition 256 MB array

Partition 256 MB array

Partition 256 MB array

flatMap



Partition 0 4 MB 1 4 MB 2 4 MB

Partition 0 4 MB 1 4 MB 2 4 MB

Partition 0 4 MB 1 4 MB 2 4 MB

Partition 256 MB array

Partition 256 MB array

Partition 256 MB array

flatMap

→

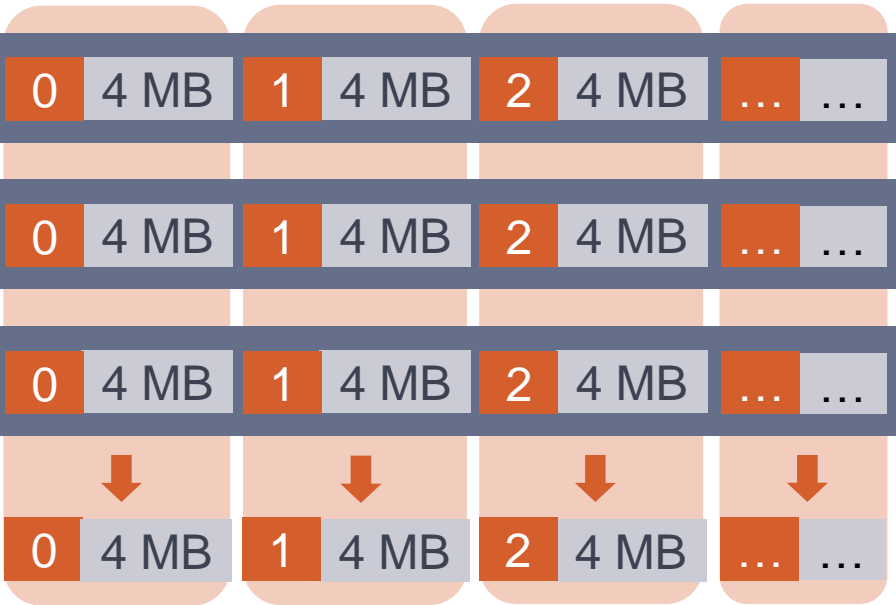
→

Partition 0 4 MB 1 4 MB 2 4 MB

Partition 0 4 MB 1 4 MB 2 4 MB

Partition 0 4 MB 1 4 MB 2 4 MB

reduce by key 0 4 MB 1 4 MB 2 4 MB



Partition 256 MB array

Partition 256 MB array

Partition 256 MB array

flatMap

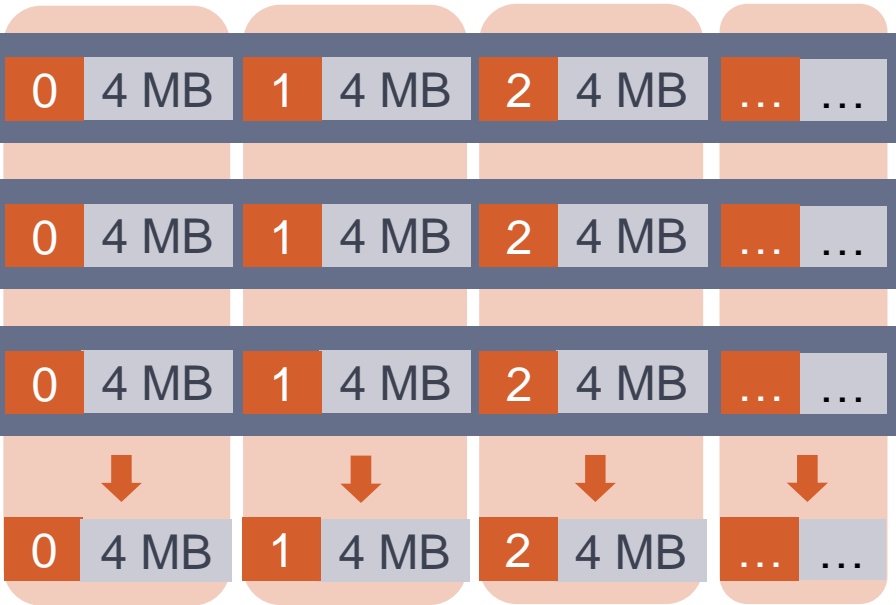
Partition 0 4 MB 1 4 MB 2 4 MB

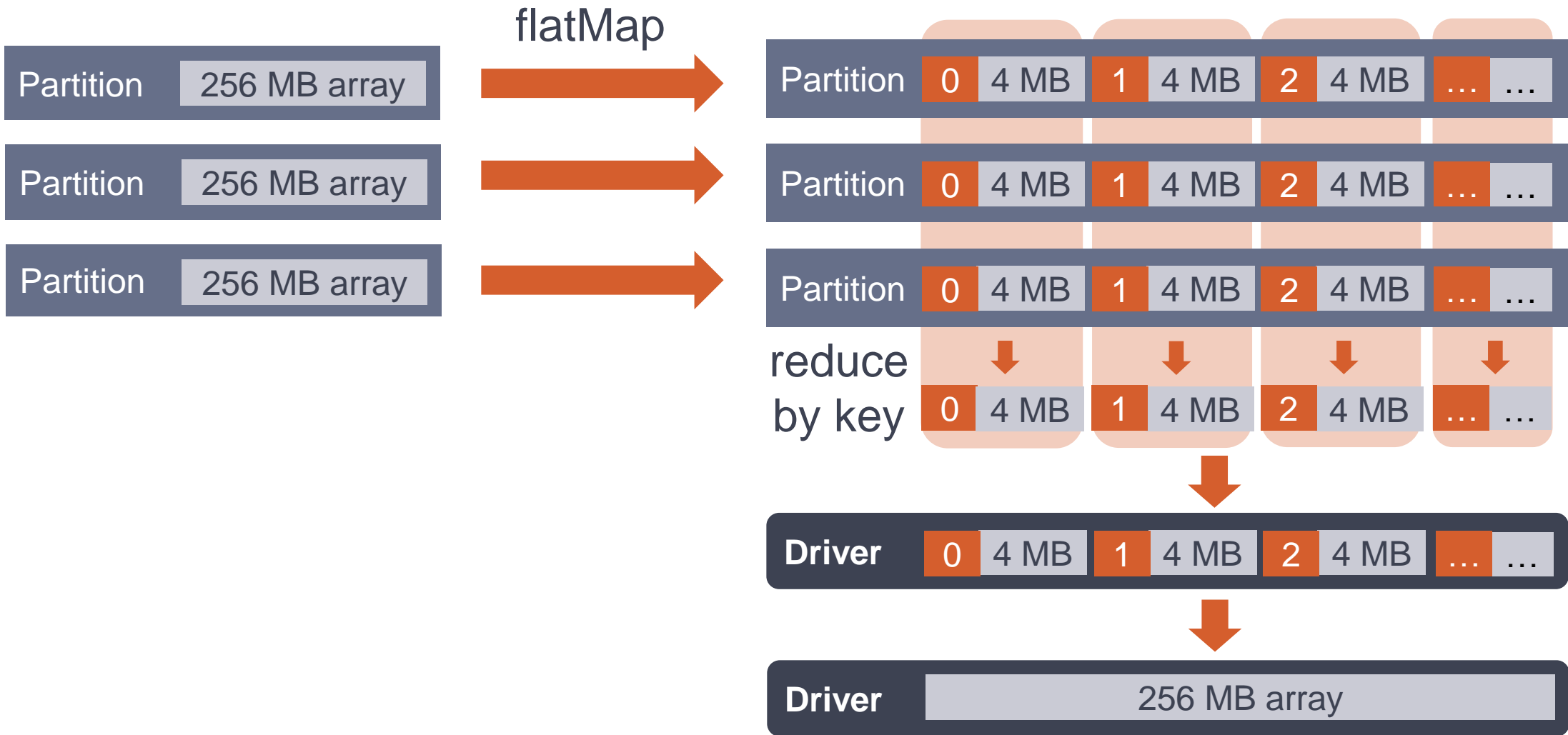
Partition 0 4 MB 1 4 MB 2 4 MB

Partition 0 4 MB 1 4 MB 2 4 MB

reduce by key 0 4 MB 1 4 MB 2 4 MB

Driver 0 4 MB 1 4 MB 2 4 MB

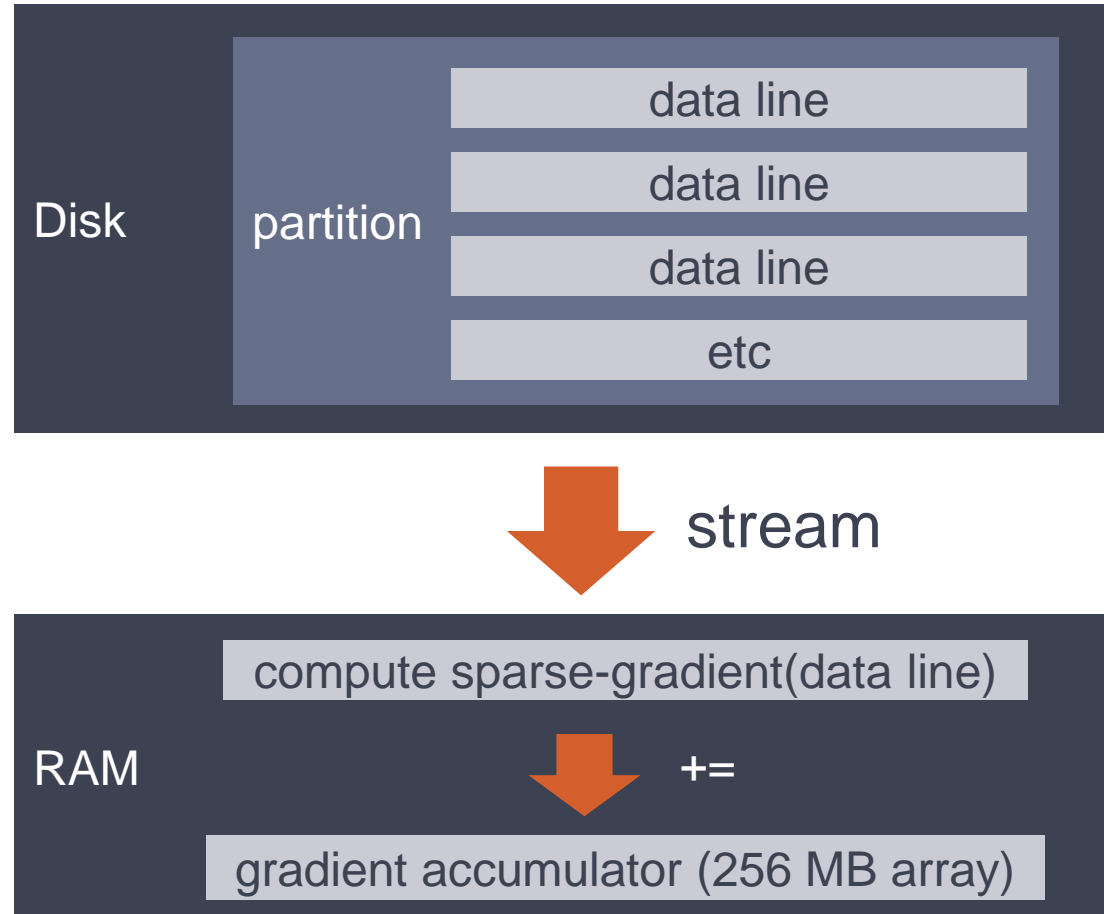




Out Of Memory — Java heap space



Computation for each partition



**Multiple weeks of
investigation**

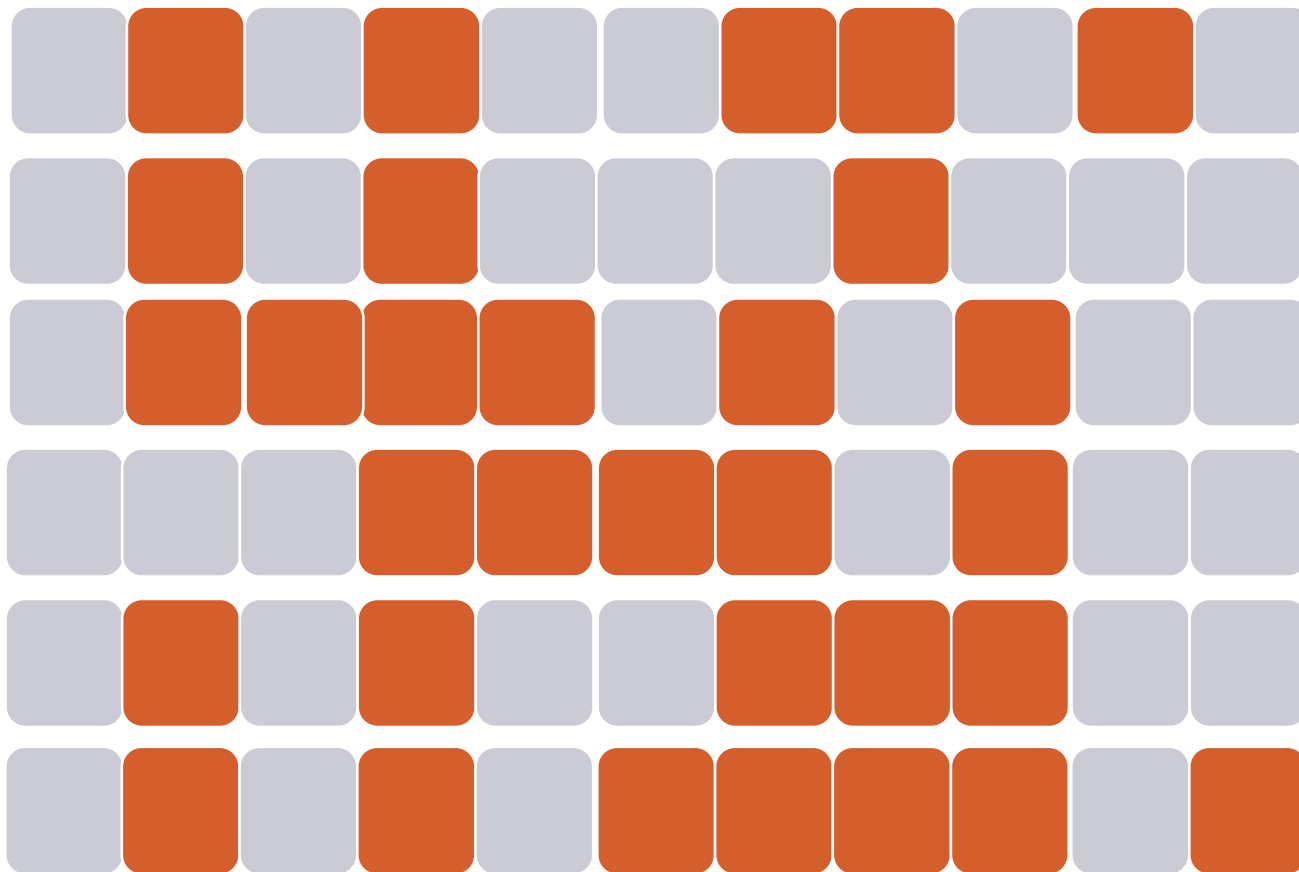
Activated G1 GC logs on Executor

- XX:+PrintGCDetails
- XX:+PrintGCTimeStamps
- XX:+PrintHeapAtGC
- XX:+PrintTenuringDistribution
- XX:+UnlockExperimentalVMOptions
- XX:+UnlockDiagnosticVMOptions
- XX:G1LogLevel=finest
- XX:+G1PrintRegionLivenessInfo

Humongous Object



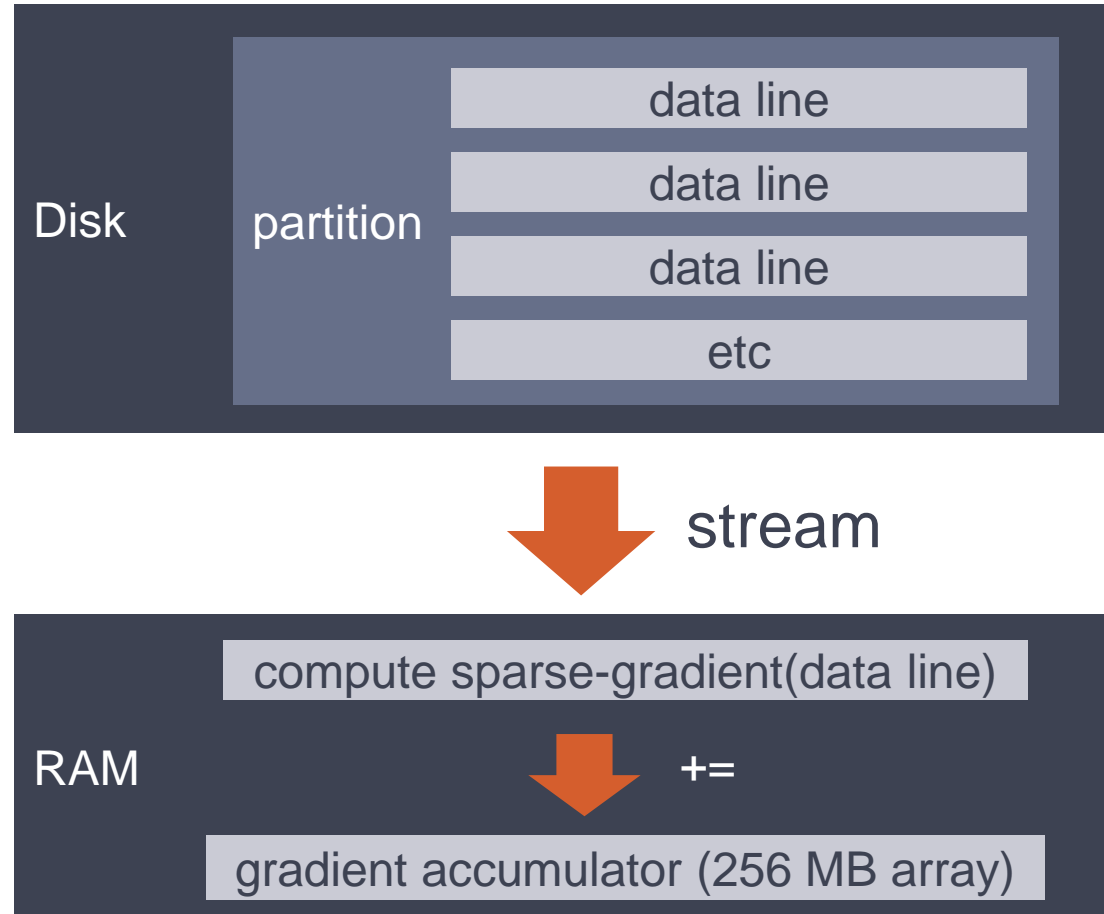
Fragmented Heap



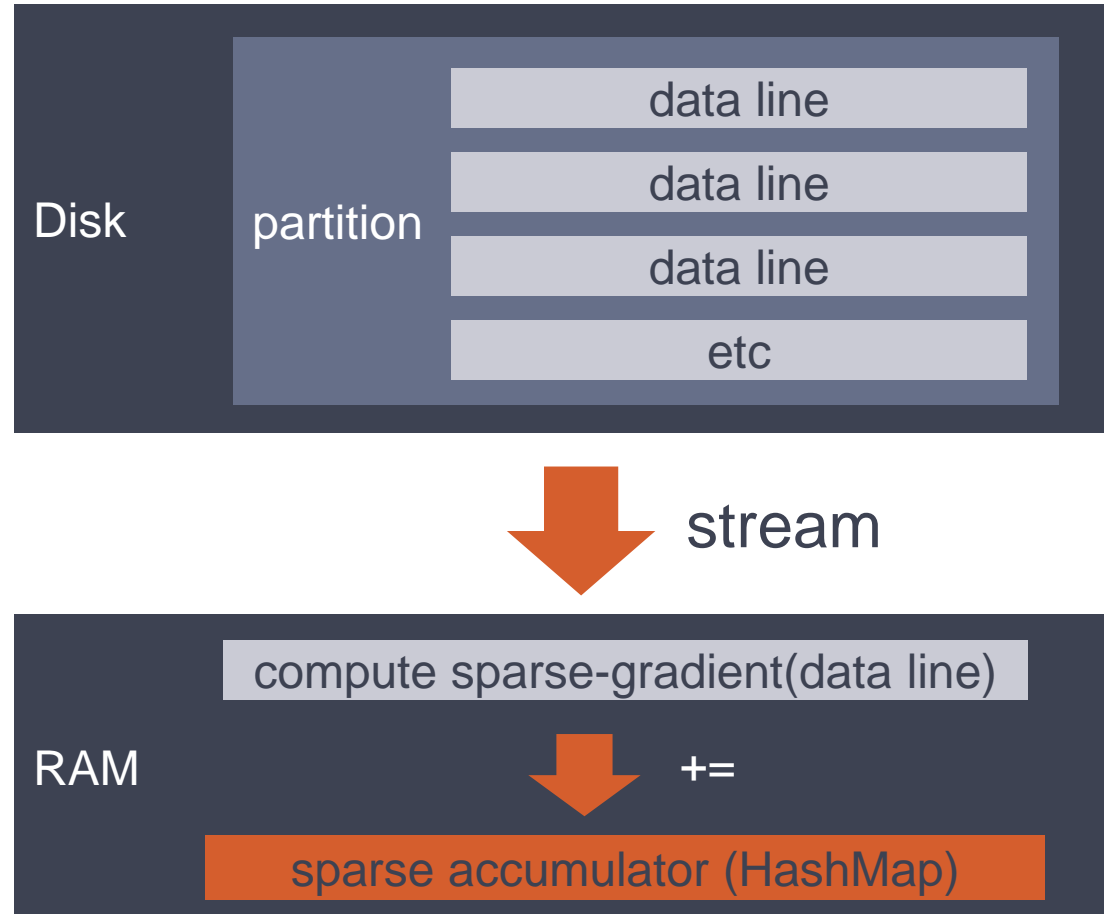
Taken

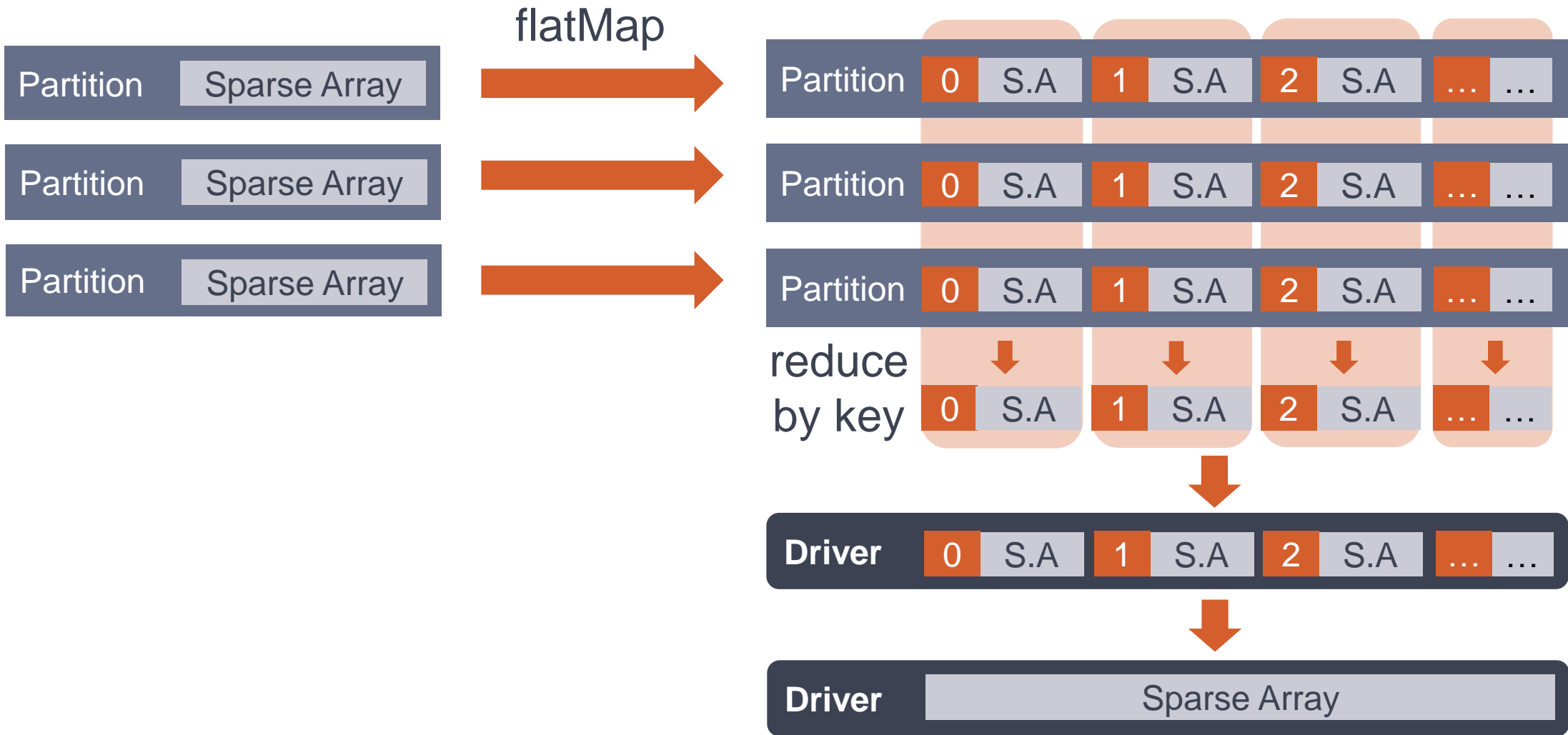
Free

Computation for each partition



Computation for each partition





Learning 1000s of models



**Several 1000s of models must be
learned every day**



**Yarn scheduling is based on
resource consumption only
(CPU & Memory)**



Scheduler



Training

Training

Training

Prepare
Config



Spark
Preprocess



Spark
Learning

Deploy

Scheduler

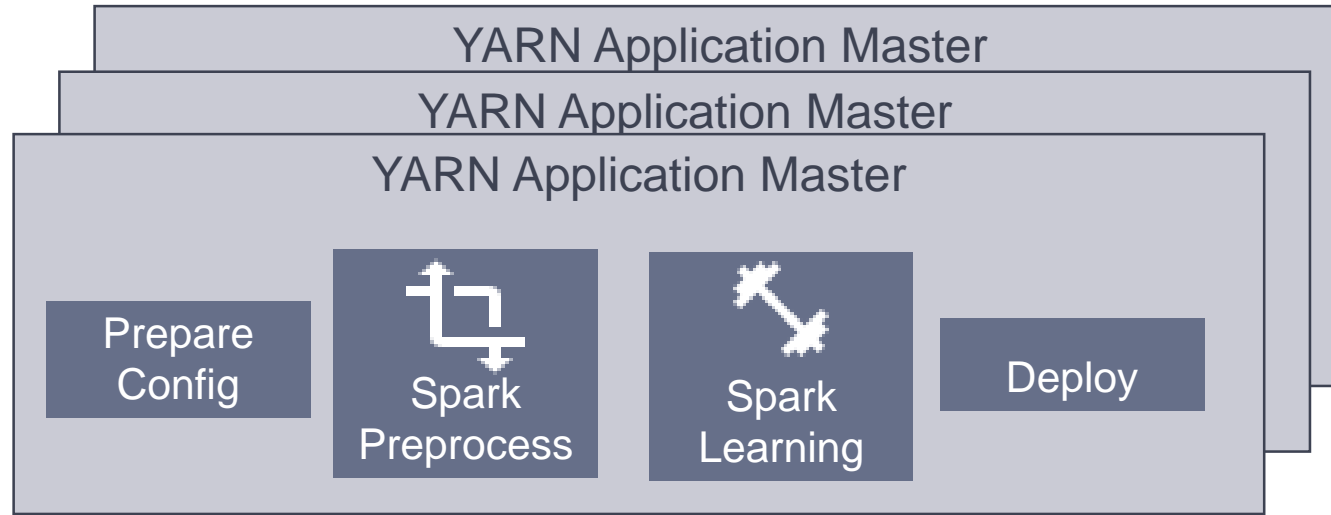
Prepare
Config

Spark
Preprocess

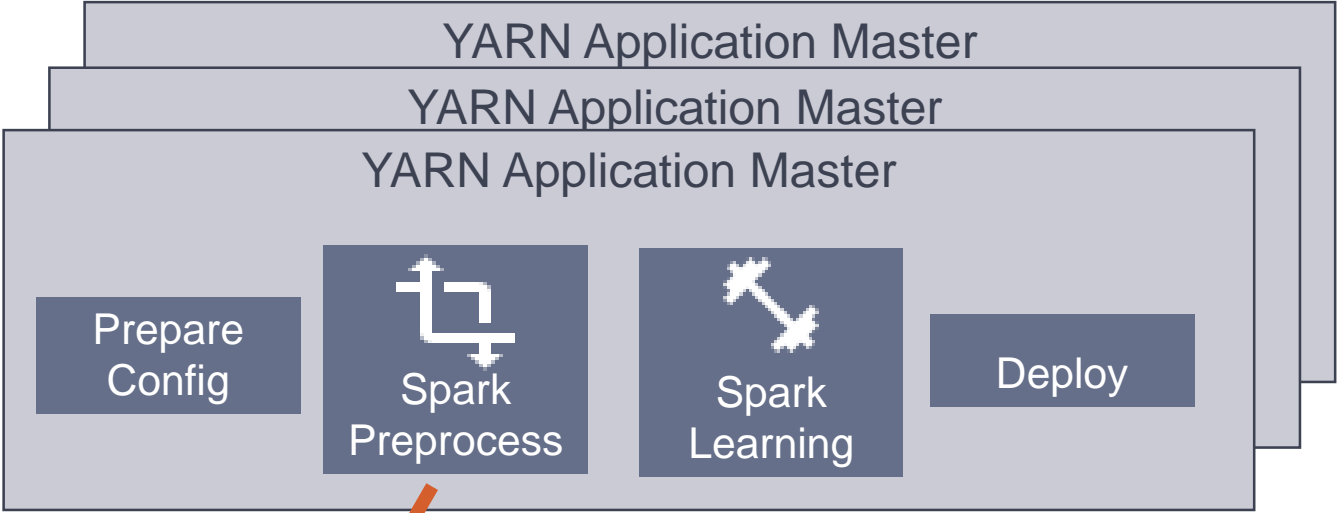
Deploy



Scheduler



Scheduler



Spark Application Master

Client



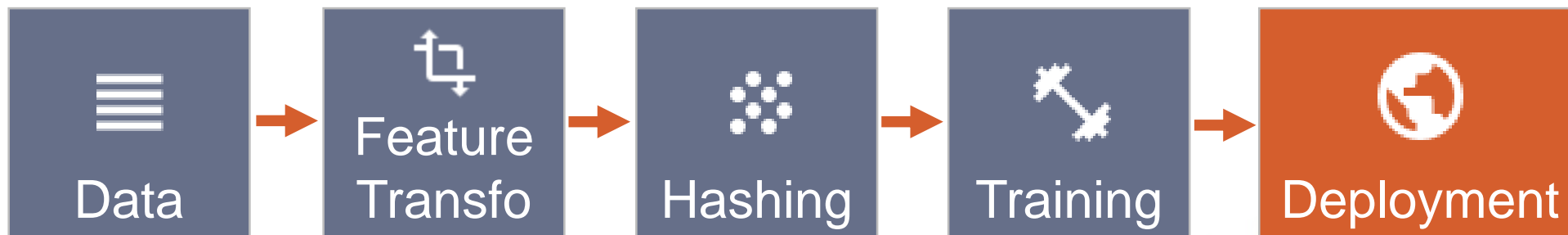
Spark Executor

Task

Spark Executor

Task

Model deployment





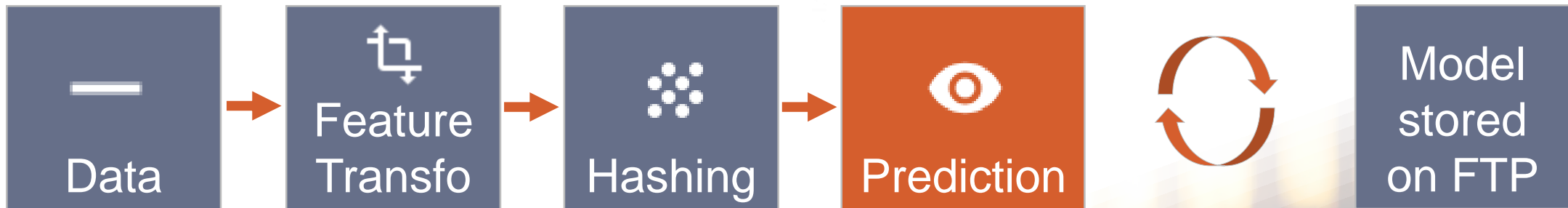
Online prediction



**Prediction must be made in less than
100 μ s**



Online Prediction Workflow



10 Million

predictions per second

Monitoring



Finatra Prometheus Graphite



Questions

f.horing@criteo.com
@f_hoering

