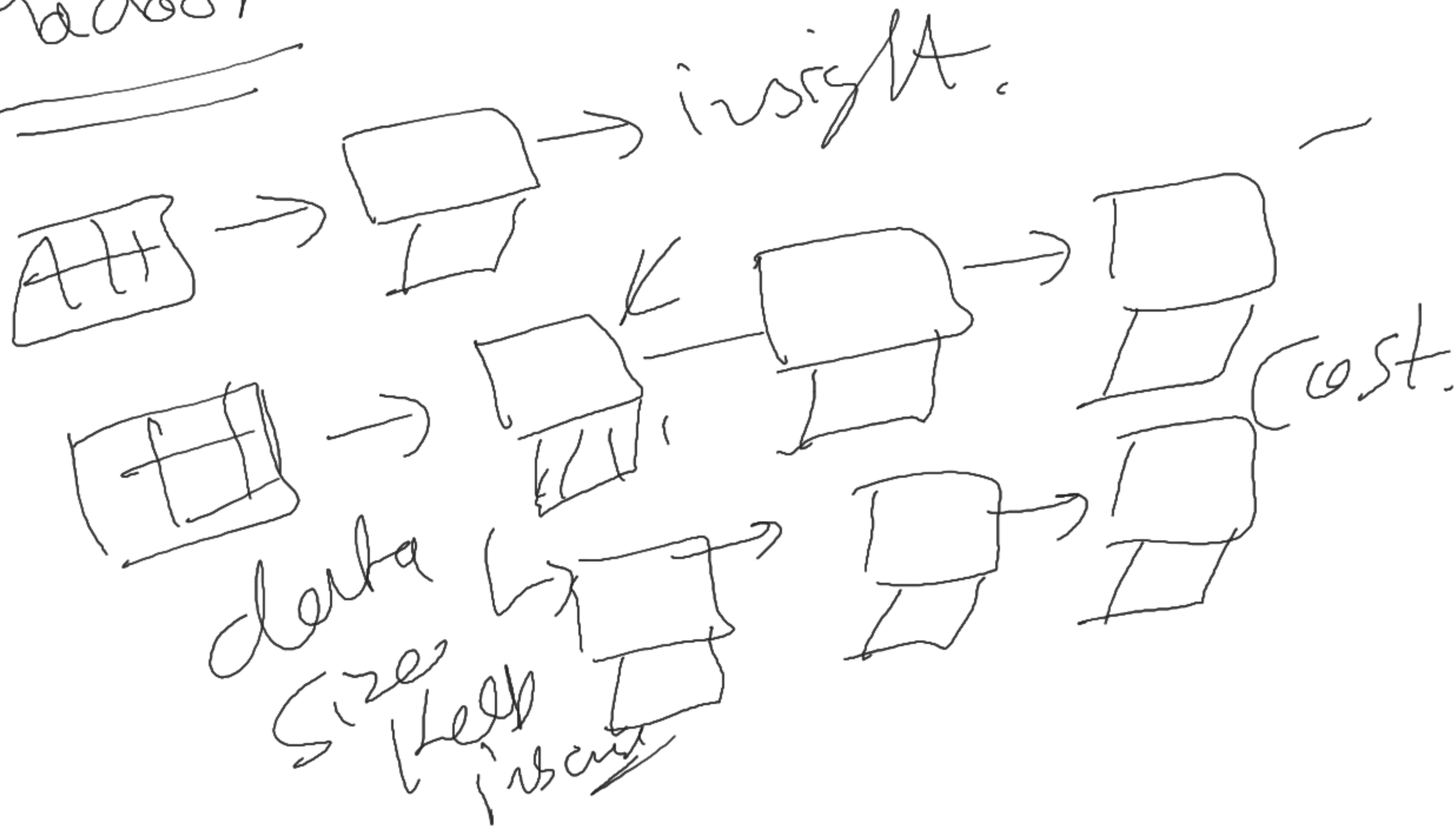


Hadoop



Google
AWS

Azure

inside

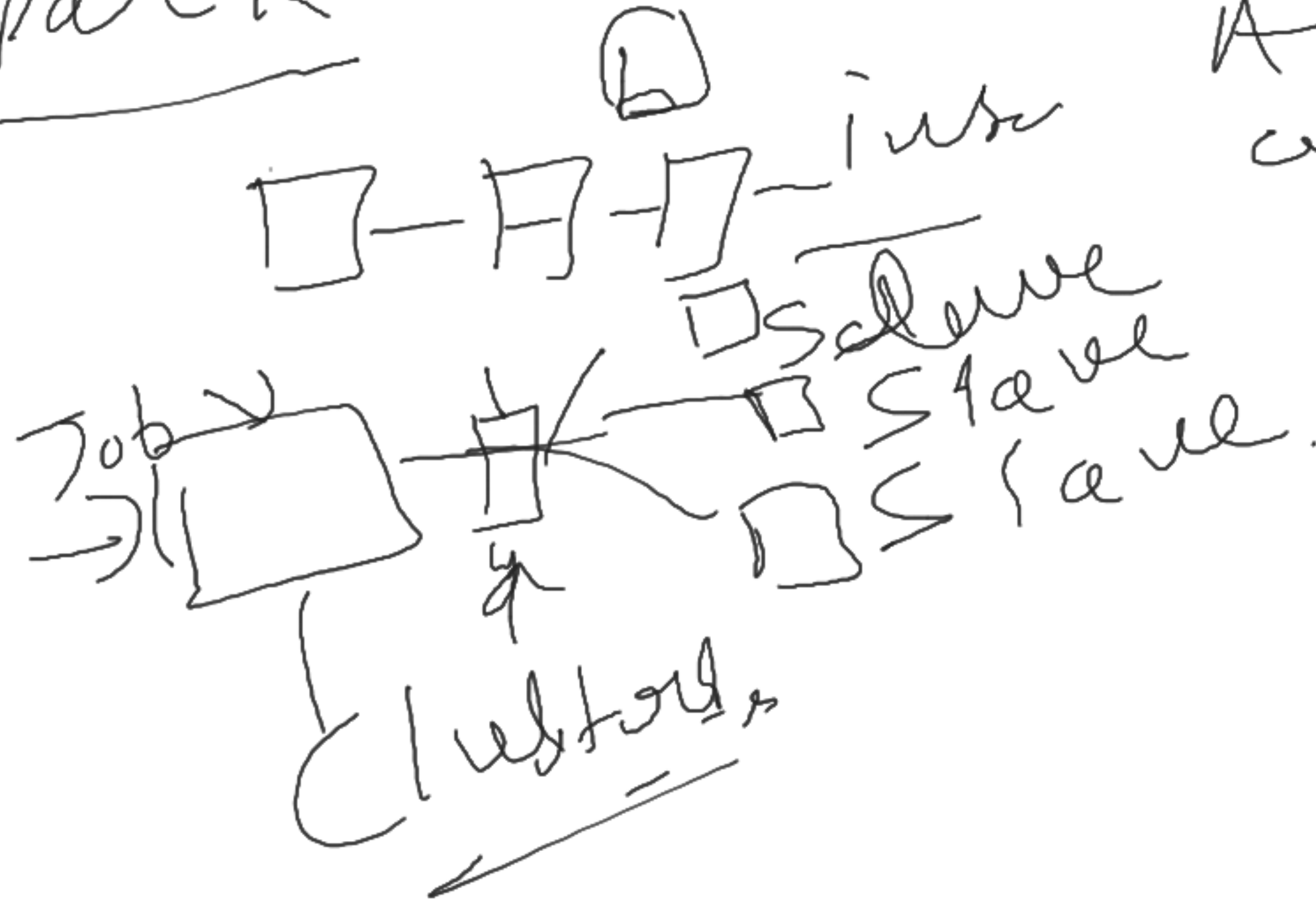
Harddisk @

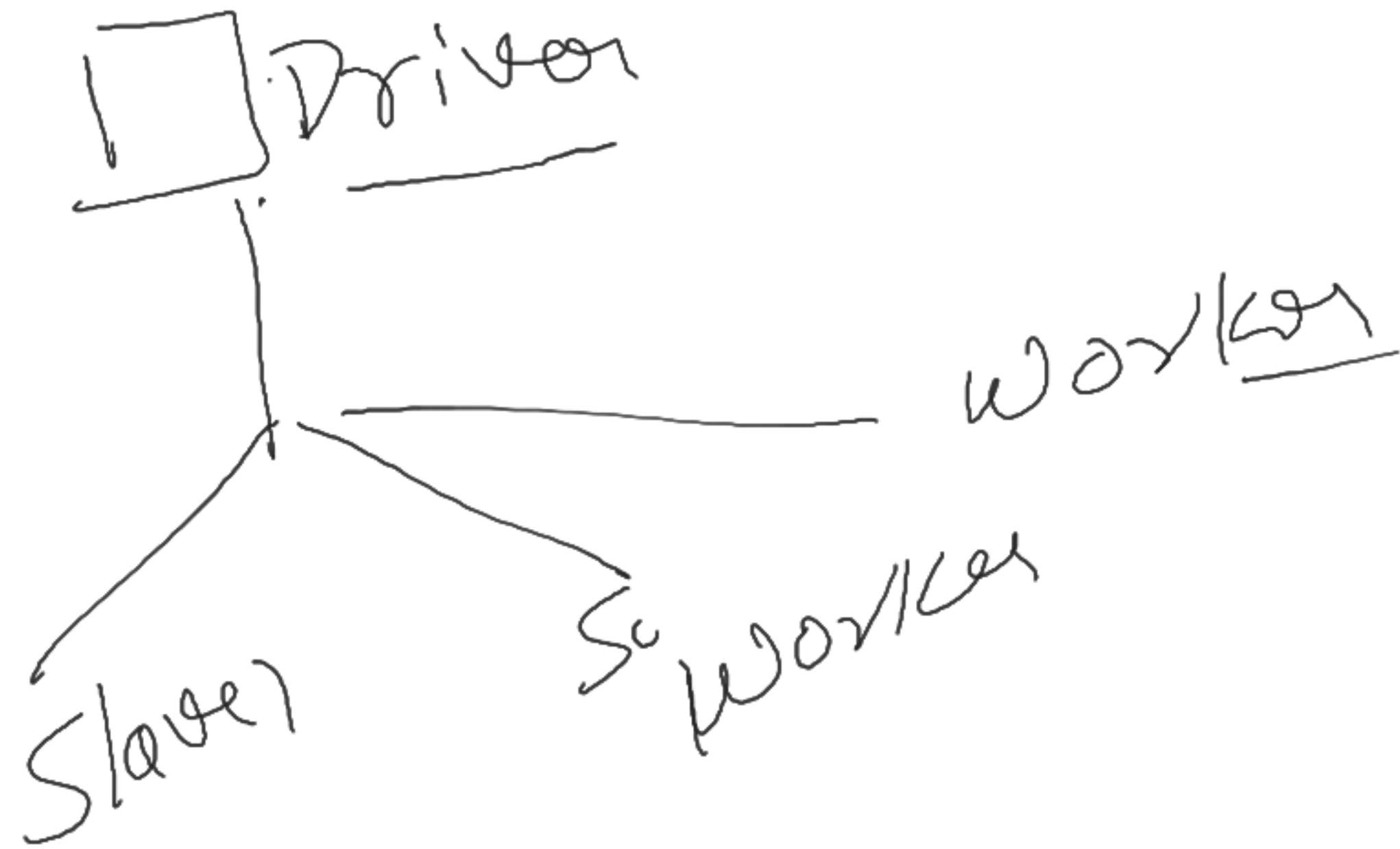
Spark \rightarrow 100 times

in memory process

Spark

AWS
cPP





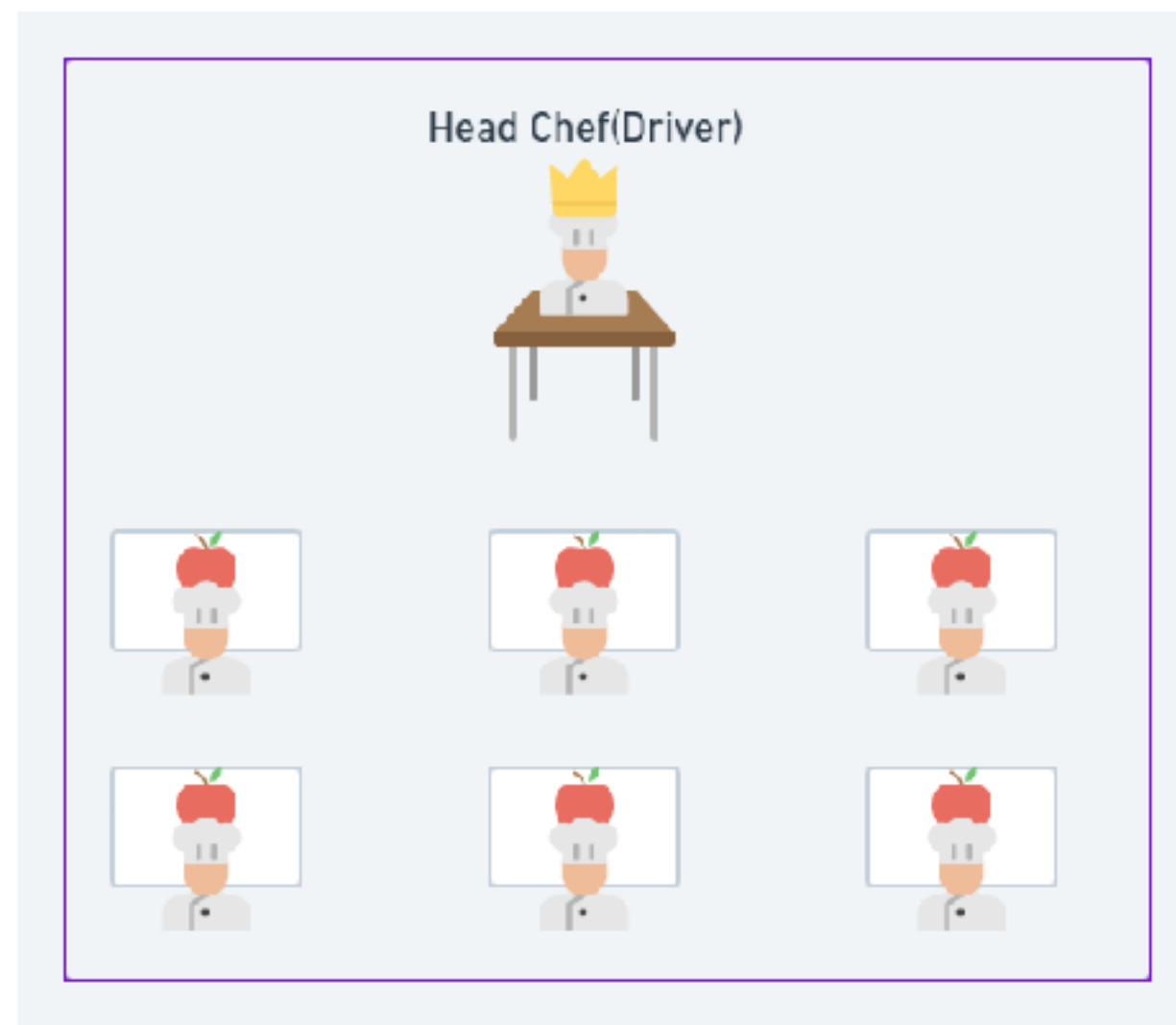
Restr → Fruit Salad

1 Apple → ~~Person~~ Chef

1000 → ~~Person~~ Chef

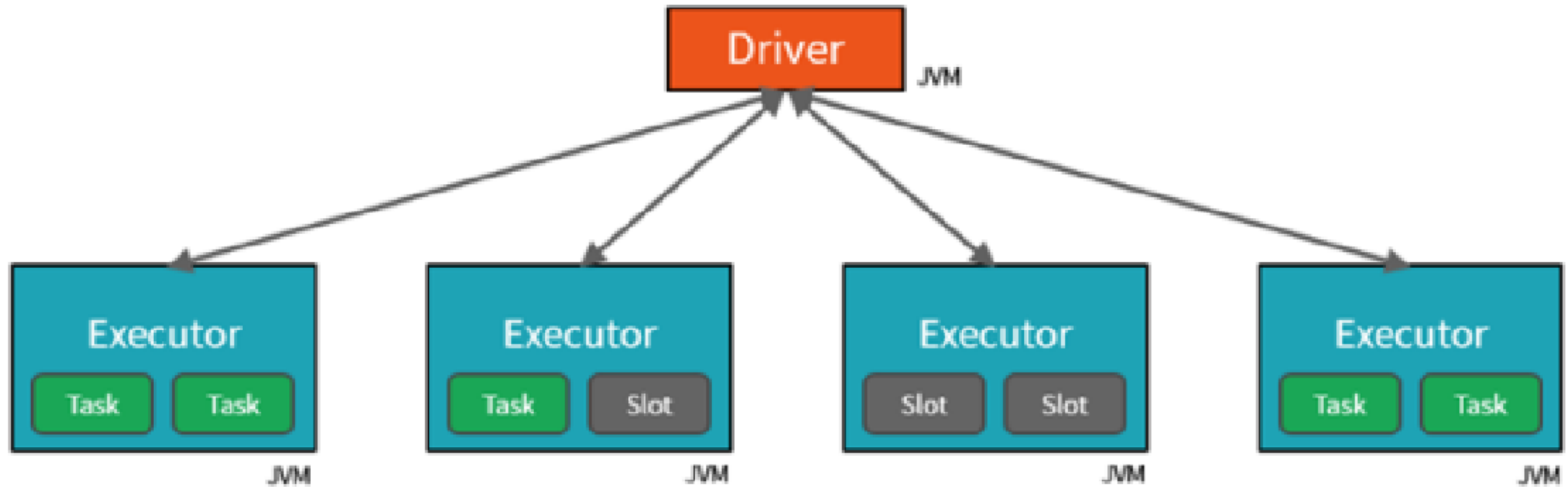
too much



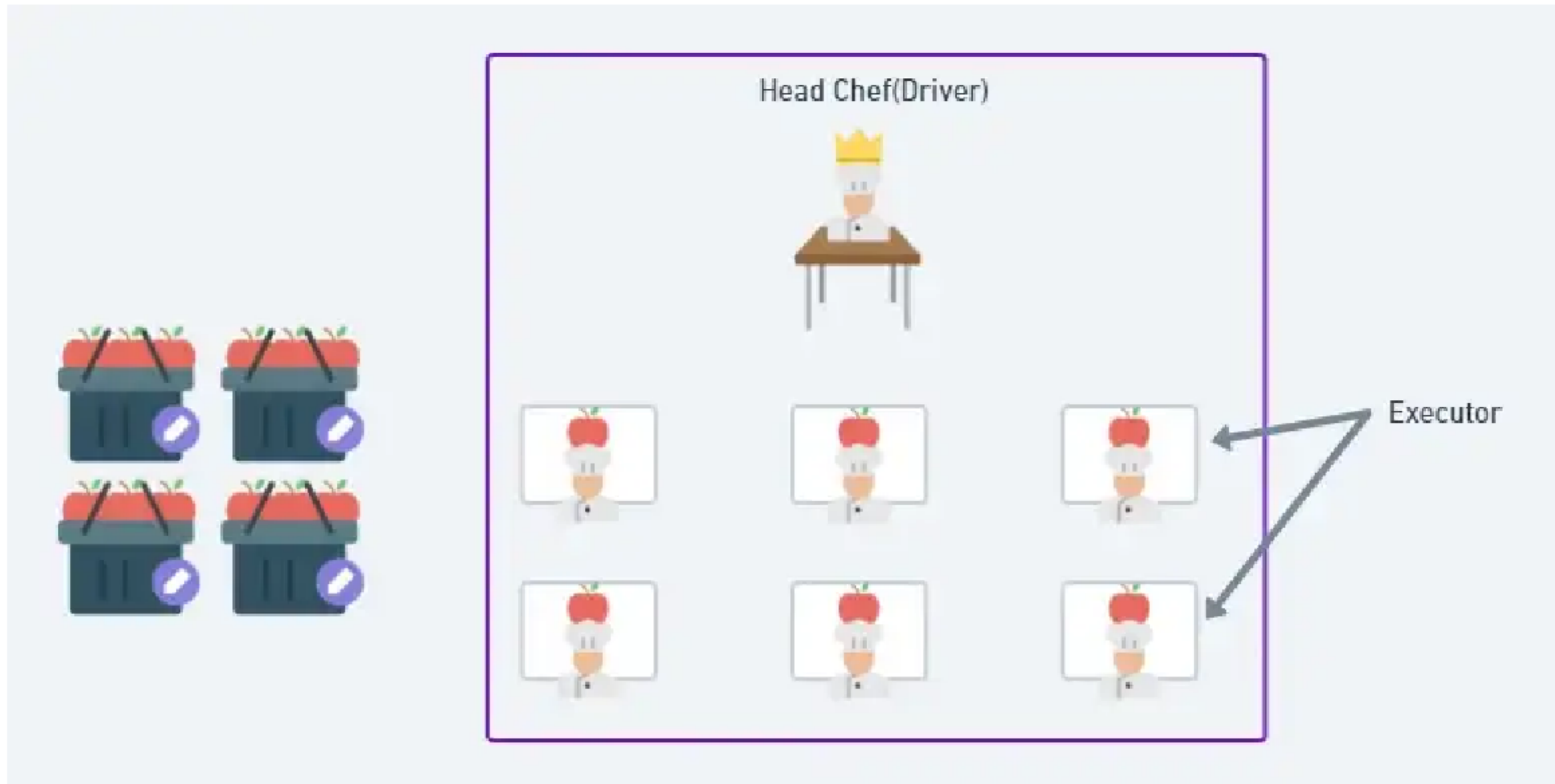


Driver: Consider Driver as a head chef/manager of the restaurant who distributes the work among the other chef.

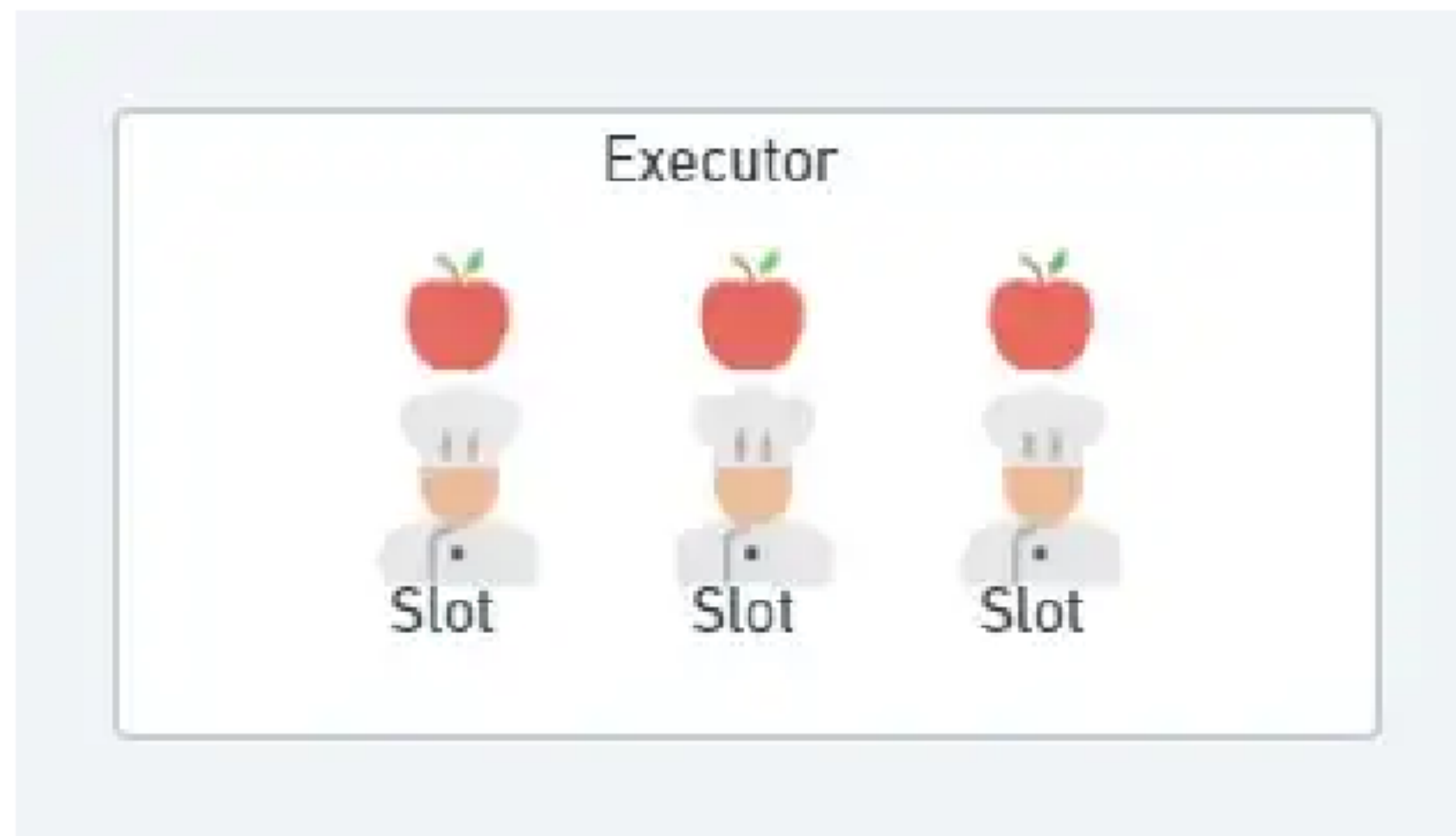
In spark, the Driver is responsible for assigning Tasks to the executors. From a Data engineer's perspective if we want to filter a table consisting of 1000 rows these rows will be divided and distributed across the executors by Driver. Note that the driver will not be allowed to view the data.



Executors: Consider executors as a workstation/chopping board where the chef performs the task(cutting apples).



In spark, Cores/Slots are the CPU that performs the actual task. The total number of cores present is responsible for executing more tasks parallelly.



In spark, the executors are the Java Virtual Machine where your task is being performed. The executors are responsible for running the individual task and sending the result to the driver.

Slots/Cores: Slots or Cores are referred to the chefs performing the tasks of cutting apples. An executor consists of Cores/Slots. An executor(workstation) can consist of one or more cores/slots.

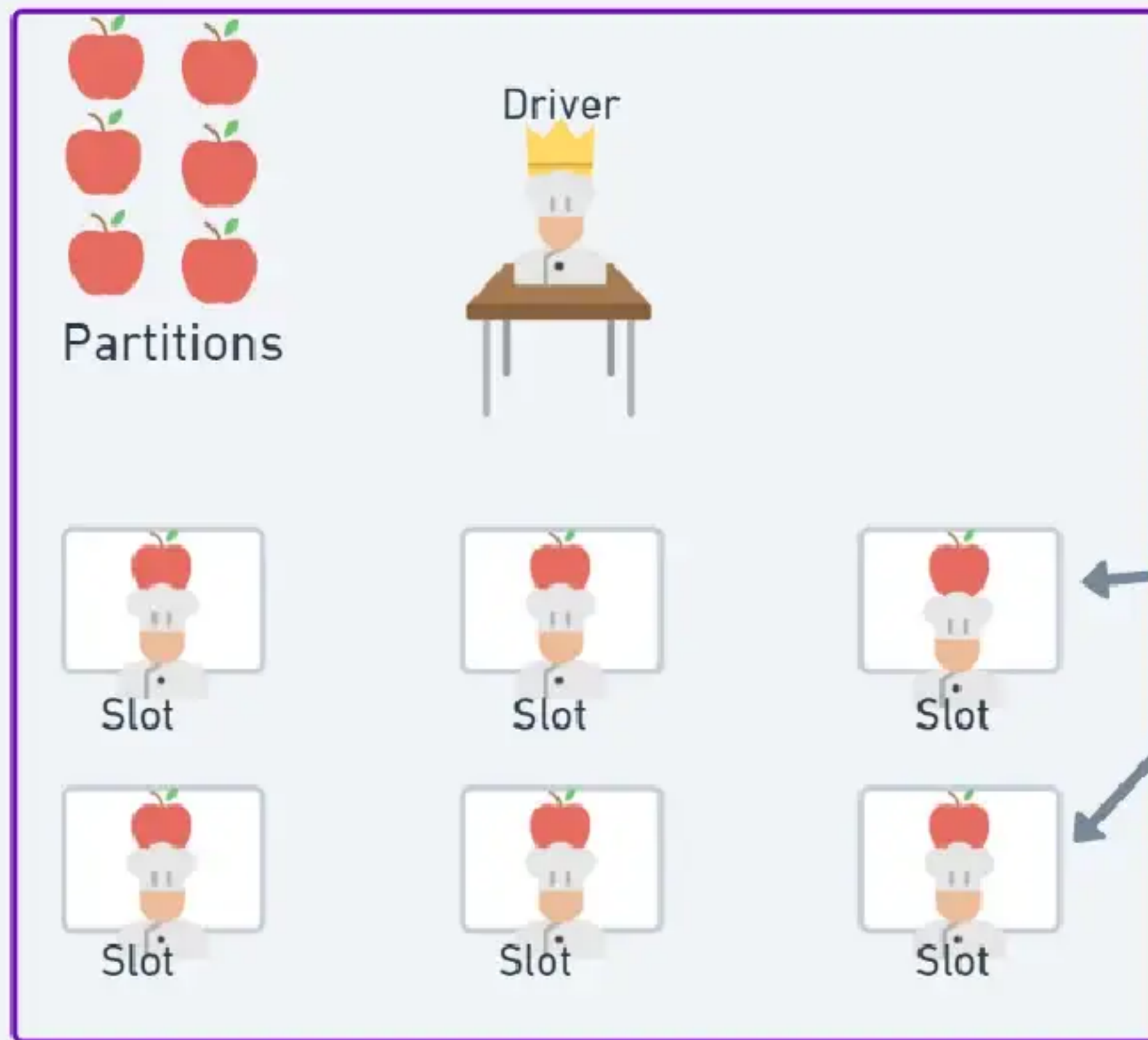
Jobs, Tasks, Dataset, Partitions:

Dataset is a collection of apples in a basket.

Partitioning is the process of splitting the dataset into smaller parts. Partitioning is done by the driver.

Tasks are usually the total number of partitions that needs to process for example if we have 4 baskets of 3 apples in each then we have 12 apples(Dataset) to be cut. This is divided into 12 partitions by driver hence 12 tasks have to be performed in total. Each task is performed by each core. Hence if we have 3 executors of 2 cores in each executor then we have 6 cores in total therefore there can be only a maximum of 6 tasks that can be performed parallelly.

Job is the overall instruction that the driver receives. In this case, we need to cut 4 baskets of 3 apples in each.



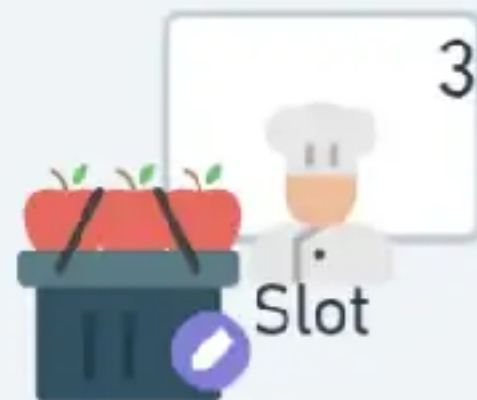
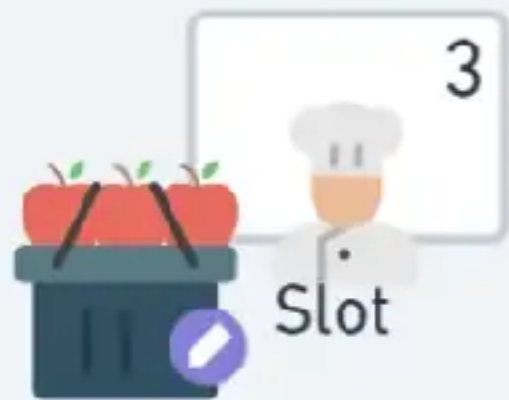
Cluster

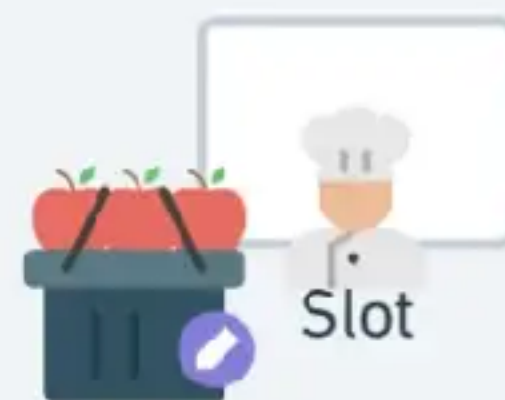
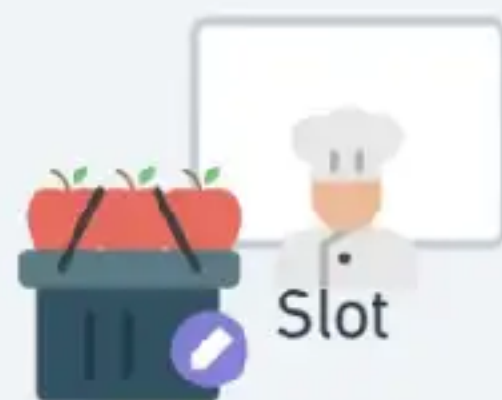
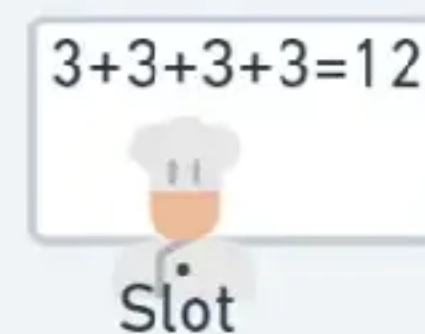
Executor

Stages:

Stages in Spark represent groups of tasks that can be executed together to compute the same operation on multiple machines. As we know already that the driver is not allowed to see the dataset then have you wondered how the driver knows the total no of rows in the dataset or apples in our example? Spark tackles this issue by distributing the dataset across executors and asking them for the count. Let's understand the entire process by going through two stages.

Stage 1: In this stage, the driver divides the dataset and distributes them to a set of executors. The driver now asks the executors to count the number of rows in the dataset. In this stage, the work of the executors is to find the no of apples in each basket. Now the count is stored within each executor on the disk to use in later stages this is also called shuffle data.







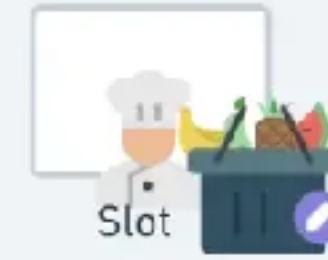
	Apple
	Strawberry
	Grape
	Orange



	Banana
	melon
	Pineapple
	Pear



	Apple
	Strawberry
	Grape
	Orange



	Banana
	melon
	Pineapple
	Pear



	Banana
	melon
	Pineapple
	Pear
	Apple
	Strawberry
	Grape
	Orange

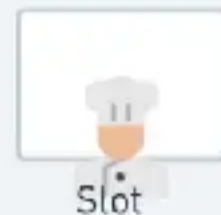
	Apple
	Strawberry



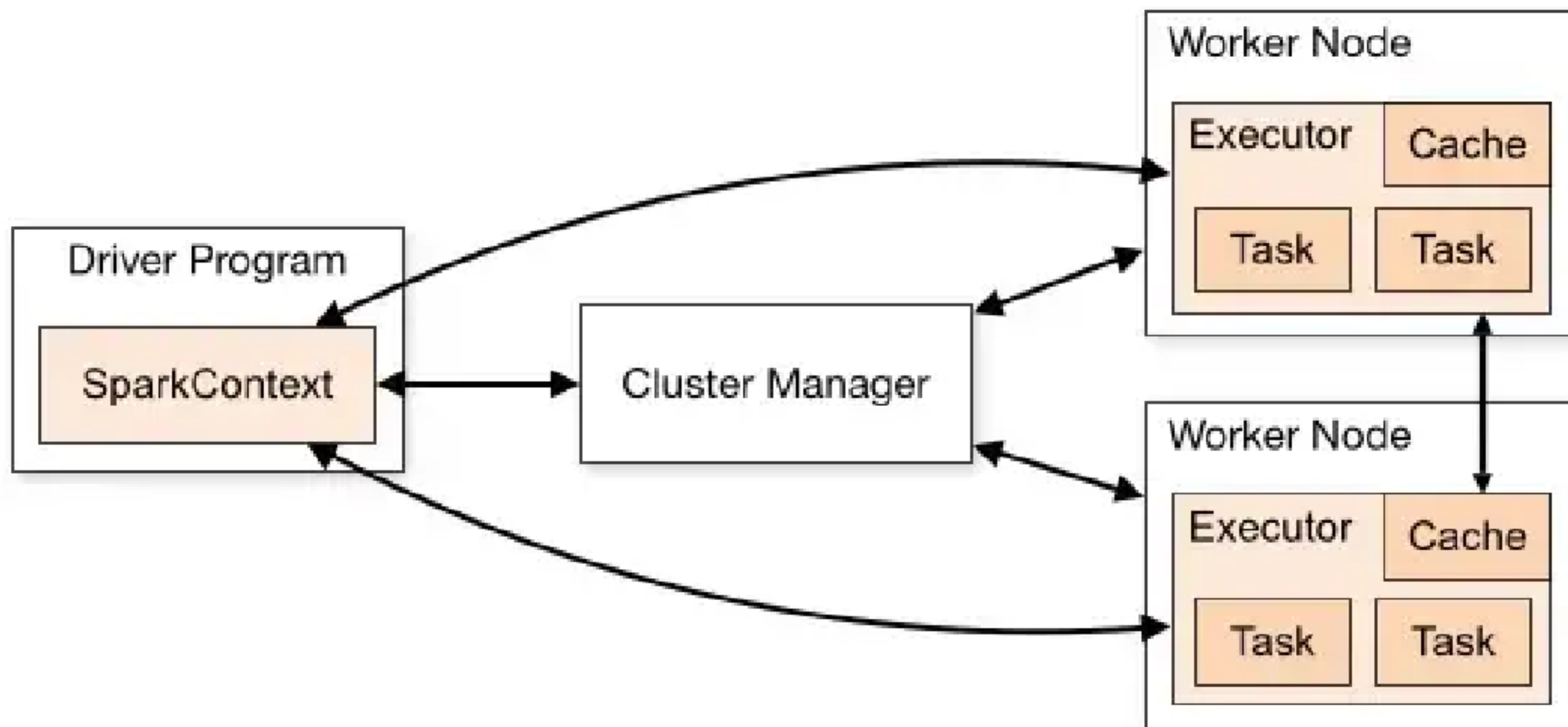
	Banana
	Pineapple



	Grape
	Orange



	Pineapple
	Pear



```
orders_df=spark.read.csv('dbfs:/FileStore/tables/Data/orders/
part_00000',schema="""order_id INT,order_date
DATE,order_customer_id INT,order_status STRING""")
```

```
1 orders_df.count()
```

Python ▶ ▾ - ✕

▼ (2) Spark Jobs

- ▶ Job 2 [View](#) (Stages: 1/1)
- ▶ Job 3 [View](#) (Stages: 1/1, 1 skipped)

Out[7]: 68883