

Session 10

APACHE SPARK

Jan 29, 2024



* Why can't we use Pandas for TBs of data??

- ① Memory problem (load everything in Ram)
- ② Single core

100 Mb csv file \rightarrow 1GB - 2GB

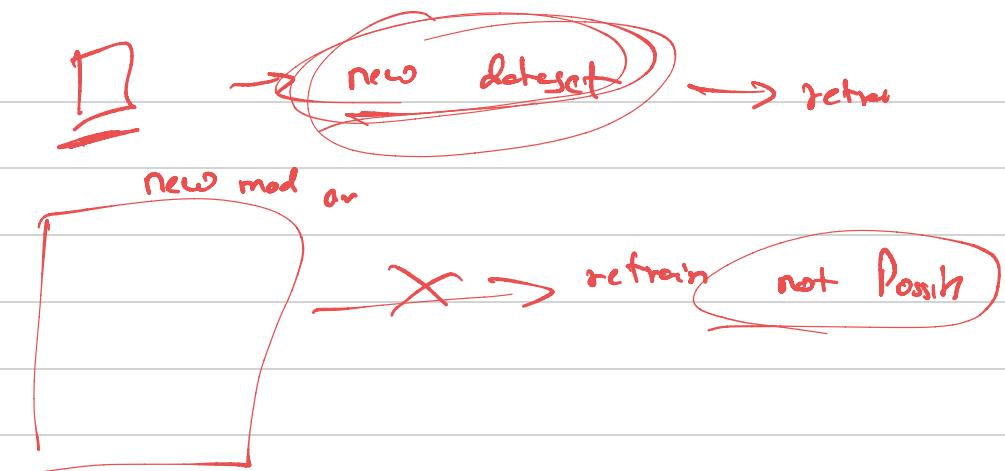
1GB CS-Vile \rightarrow 5GB of Ram space

A100 instance \rightarrow \$3.4

GPT-4 \rightarrow Open ai use

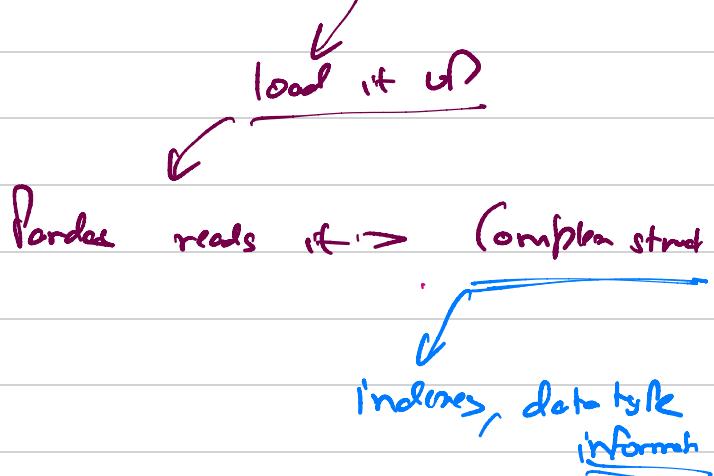
25,000 A100 instances
1 - 2 month
\$ million
\$125 million

512 GB Ram
80 GB VRAM



→ Why does it take more memory??

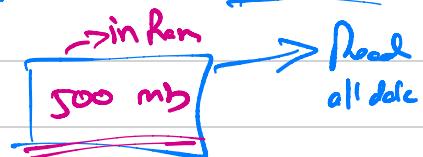
① Plain-text (CSV)



500 mb → CSV



→ Read it in Pandas



$$1\text{ GB} + 500\text{ mb} \rightarrow 1500\text{ mb}$$

of
Spec

→ What can we do to optimize this.

① Read in chunks

② Convert int64 → int32

③ Categorical data → objects

category dtype

→ consumes
lot less

Swap Spec

RAM

memory

Swap Spec

4GB

→ HDD

memory

Reading / writing → binary format

>> more efficient
than CSV

→ Dask → why do we need it??

→ Open source python source

(i) Chunk - bare processing.

- * Process data in chunks → sequentially
→ Parallel,

a months → 1 baby
→ 2 women → 4-5 months

(ii) Lazy Evaluation.

(iii) Distributed Computation → distribute task
within a computer
all the cores
d
4 cores

Pandas → >> ① Setup
② functionalities

③ Easier to Use

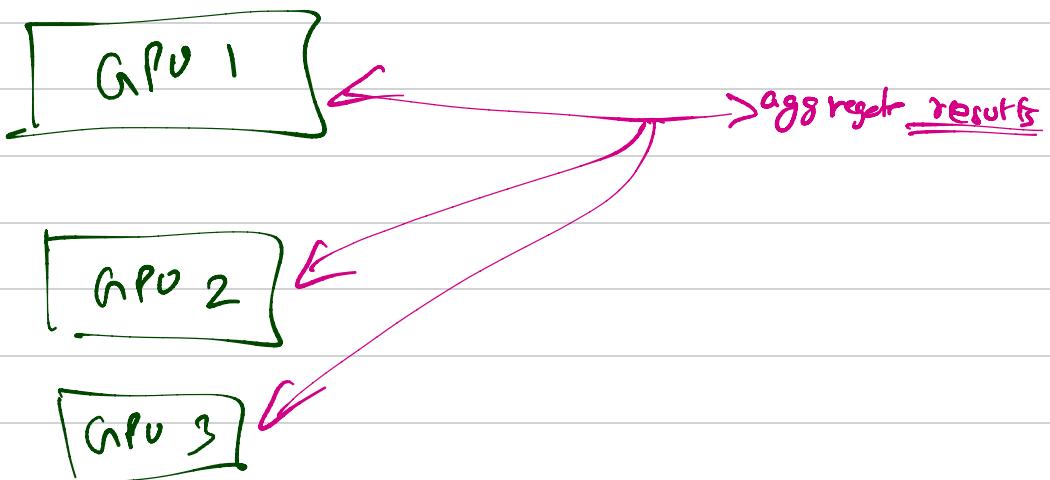
Not very Common

- ① Small data → Pandas
- ② Spark

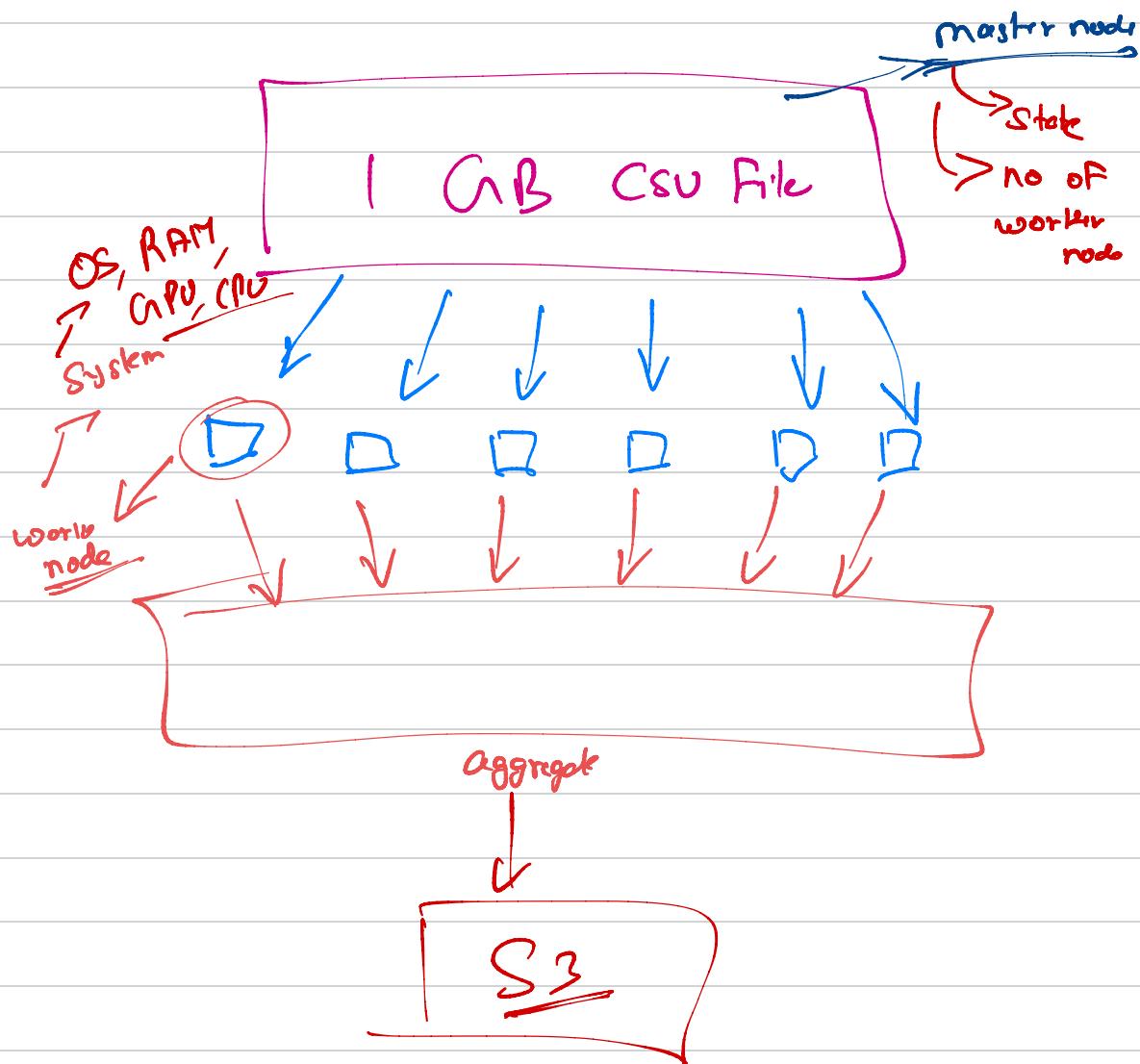


Task = 2

Sensor Flow



Spark = Why we need ??



Why we use Spark

- ① Distributed Computing
- ② Open Source
- ③ Lazy Evaluation
- ④ Not in memory.

$mr. 16n \text{ lone} \rightarrow 100 \text{ worker num}$
 $P_{2.8n \text{ lots}}$

range(1, 1000000)

for i in range ... ?

$i = \underline{4} \rightarrow \underline{\text{in my ran}}$

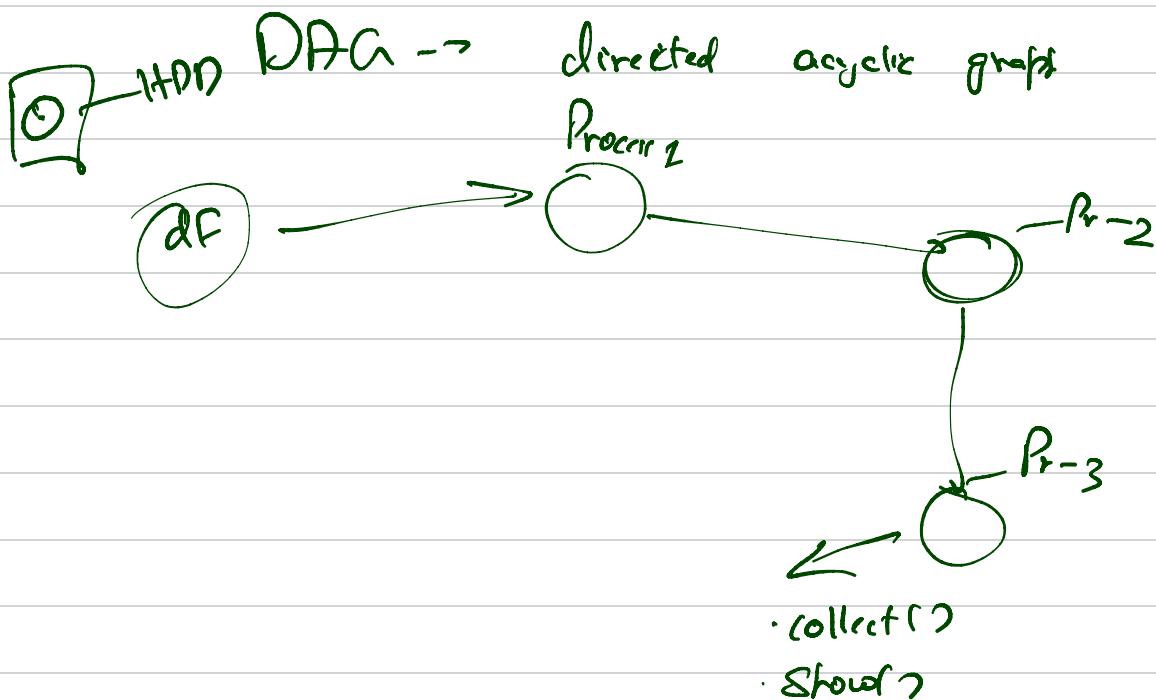
$dF \rightarrow \text{width, height, length}$

width

yield (generator function)

Why we need lazy evaluation:

① Efficient in computation



Spark \rightarrow SQL - spark

(> MLlib

↳ mgboost
↳ random forest

