

imgflip.com

# Issue WITH Pandas

TBs

- ① It tries to load everything in RAM
- ② Single Core
- ③ Does not use distributed Computation

id, name, age, salary  
1, Alice, 29, 100000  
2, Bob, 31, 95000  
...  
} every character → 1 byte

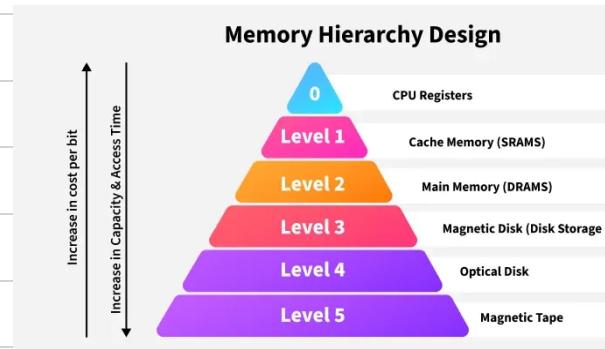
UTF-8

Operational + structural  
overhead

float 64  
float 32  
-  
-  
-

low space

## Swap Space



using the `chunksize` parameter.

### Using `chunksize` parameter in `read_csv()`

For instance, suppose you have a large [CSV file](#) that is too large to fit into memory. The file contains 1,000,000 ( 10 Lakh ) rows so instead we can load it in chunks of 10,000 ( 10 Thousand) rows- 100 times rows i.e. You will process the file in 100 chunks, where each chunk contains 10,000 rows using Pandas like this:

```
import pandas as pd

# Load a large CSV file in chunks of 10,000 rows
for chunk in pd.read_csv('large_file.csv', chunksize=10000):
    print(chunk.shape) # process the shape of each chunk
```

Output:

```
import pandas as pd

# Load a large CSV file in chunks of 10,000 rows
for chunk in pd.read_csv('large_file.csv', chunksize=10000):
    print(chunk.shape)

[1]: (10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
(10000, 9)
```

Load a Massive File as small chunks in Pandas

<https://www.geeksforgeeks.org/how-to-load-a-massive-file-as-small-chunks-in-pandas/>

For optimize → chunk size  
→ Quantization

Float 64 → Float 16

Model Variant	Parameters	VRAM (FP16)	VRAM (4-bit Quantization)
DeepSeek-LLM 7B	7 billion	16 GB	4 GB
DeepSeek-LLM 67B	67 billion	154 GB	38 GB
DeepSeek V2 16B	16 billion	37 GB	9 GB
DeepSeek V2 236B	236 billion	543 GB	136 GB
DeepSeek V2.5 236B	236 billion	543 GB	136 GB
DeepSeek V3 671B	671 billion	1,543 GB	386 GB

<https://api.proxpc.com/media/uploads/2025/01/28/vram-deepseek-table-1.png>

DASK

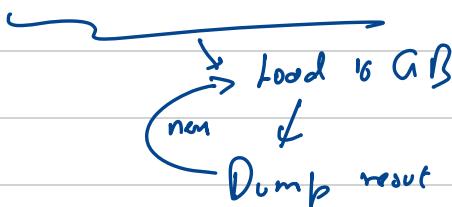


70% of API / function  
is similar to

- ① Does chunk wise processing

Pandas

100 GB - file

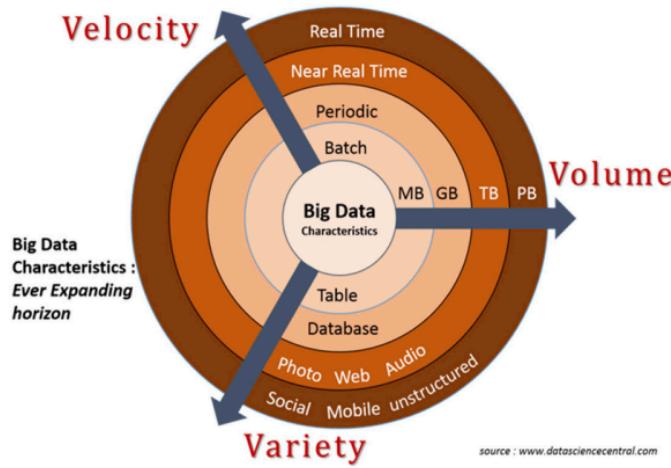


- ② Lazy Evaluation

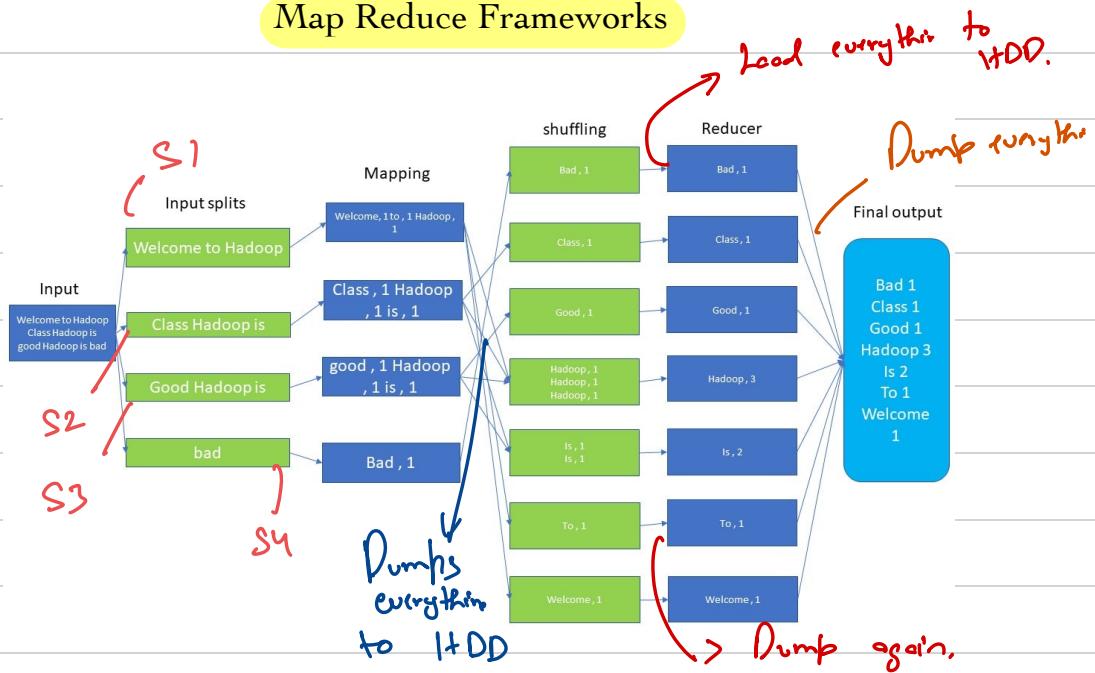
③

Distributed Computing. (not very stable)

## Big Data



## Map Reduce Frameworks



Inverted Indexing → Google Score.

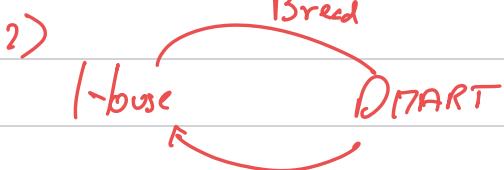
## Apache Spark

- ① Keeps Everything in RAM
- ② Open Source
- ③ Distributed Computing
- ④ Lazy Evaluation

Groceries:

- ① Milk
- ② Bread
- ③ 2 kg potatoes
- ④ 1 kg basmati Rice

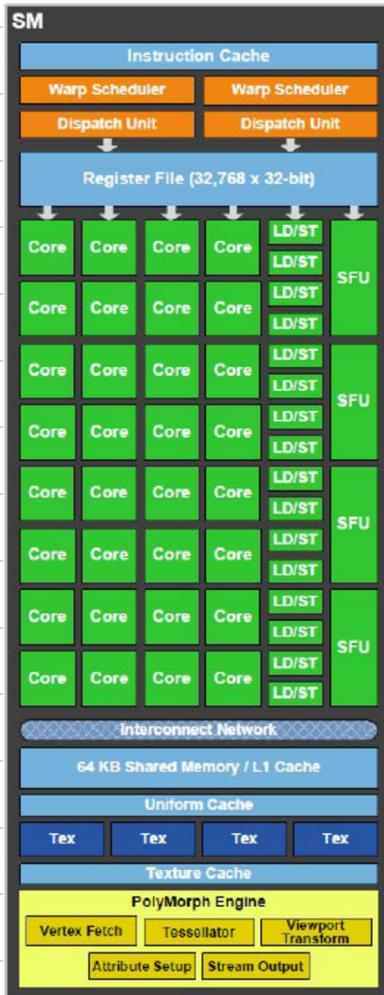
Eager Evaluation



Lazy Evaluation

→ In End, I buy  
these stuffs

Java, C/C++



The NVIDIA RTX 4090 (AD102, Ada Lovelace) features:

- 128 Streaming Multiprocessors (SMs) [techpowerup.com +12](#)

- 16,384 CUDA cores (derived from 128 SMs × 128 CUDA cores per SM) [techpowerup.com +10](#)

So in summary: 128 SMs, packing a total of 16,384 CUDA cores.

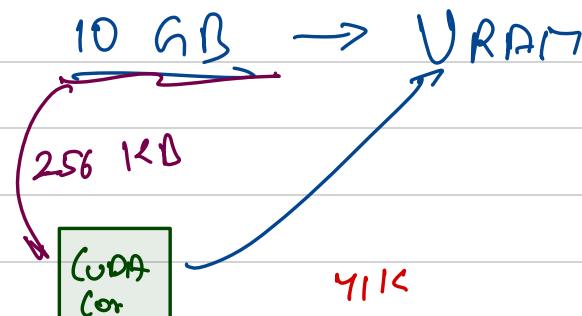
The RTX 5090 (based on the GB202 die, part of NVIDIA's Blackwell architecture) features:

- 170 Streaming Multiprocessors (SMs) — each GPC has TPCs, and in total the GPU includes 170 SM units [images.nvidia.com +14](#).
- 21,760 CUDA cores — that's 170 SMs × 128 CUDA cores per SM [kitguru.net +6](#).

So inside an RTX 5090 you get 170 SMs and 21,760 CUDA cores.

- ① Multiplication  
② Component

→ Dump entire data in VRAM



## What is a key limitation of using pandas for Big Data processing that Apache Spark addresses more effectively?

0 users have participated

- A Pandas has more complex functions for simple data analysis tasks. 0%
- B Pandas operates primarily on in-memory data and may struggle with datasets that exceed available RAM. 0%
- C Pandas has a more extensive library for machine learning algorithms compared to Spark. 0%
- D Pandas can only process data in structured formats like CSV or Excel. 0%

[End Quiz Now](#)



Santhosh K

2/3

95.47



Renu Singh

1/1

95.80



Sagar Deka

1/1

93.17

1n 4 89.37

1n 4 89.20

1n 4 89.20

1n 4 89.00

1n 4 86.93

1n 4 84.73

1n 4 82.33

## What is the primary advantage of using Apache Spark over traditional data processing tools for handling Big Data?

0 users have participated

- A Spark can only handle structured data types 0%
- B Spark requires a deep understanding of Java programming. 0%
- C Spark restricts the use of machine learning libraries. 0%
- D Spark provides fast data processing for large, complex datasets. 0%

[End Quiz Now](#)



Santhosh K

2/2

184.56



Renu Singh

2/2

190.53



Sagar Deka

2/2

182.06

2/2 4 181.23

2/2 4 179.20

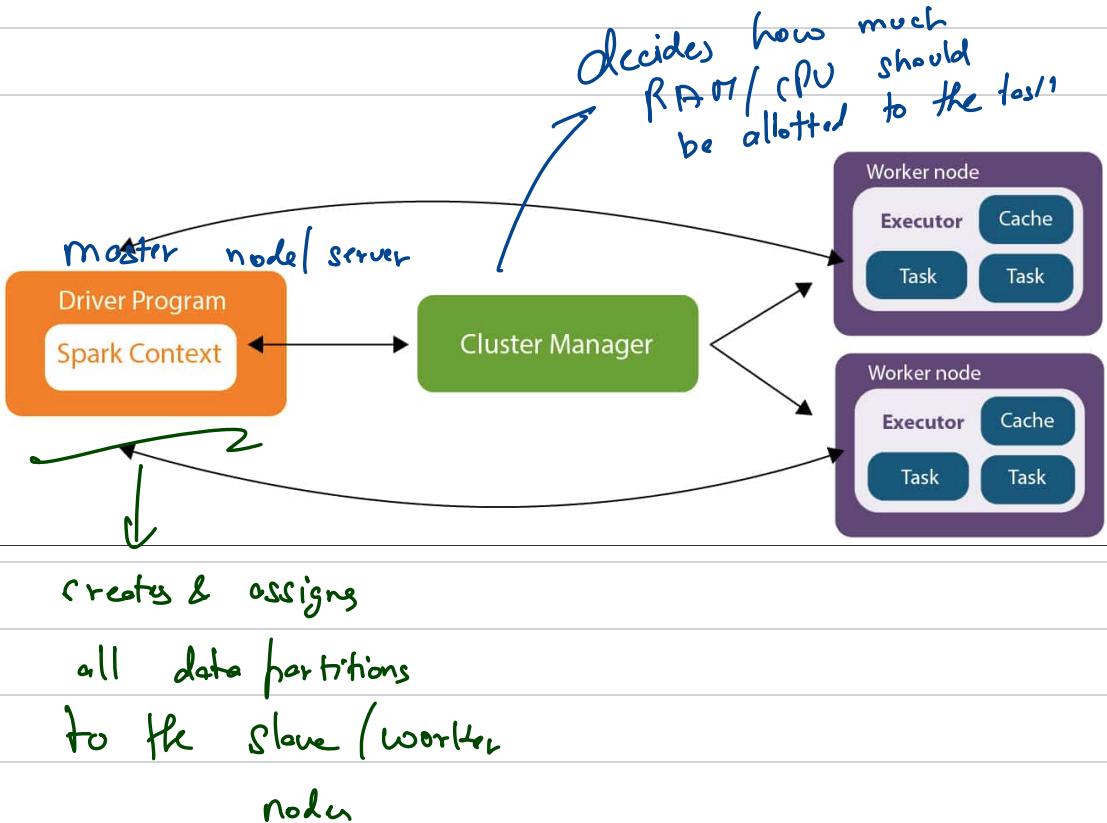
2/2 4 178.67

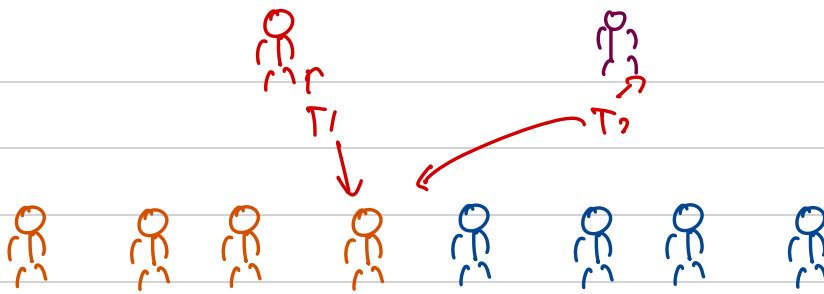
2/2 4 178.50

2/2 4 177.40

2/2 4 175.00

2/2 4 170.06





### How does PySpark handle data partitioning by default when creating an RDD?

0 users have participated

- A PySpark keeps all data on a single partition to enhance data security. 0%
- B PySpark does not support partitioning; all data must be managed manually. 0%
- C By automatically determining the number of partitions based on available resources and datasets. 0%
- D PySpark randomly scatters data across the network without regard for resource optimization. 0%

[End Quiz Now](#)

