

Shuffle: bring together data that is related but resides in different nodes. It happens when doing a wide transformation (JOIN, GROUP BY)

Q: What is the total sales per store? Perform an aggregated sum.

ie. Process data for 4 stores in n files, each containing Store_ID & Sales_volume data

❓ How does it work?

1 Spark reads the data from files into partitions 



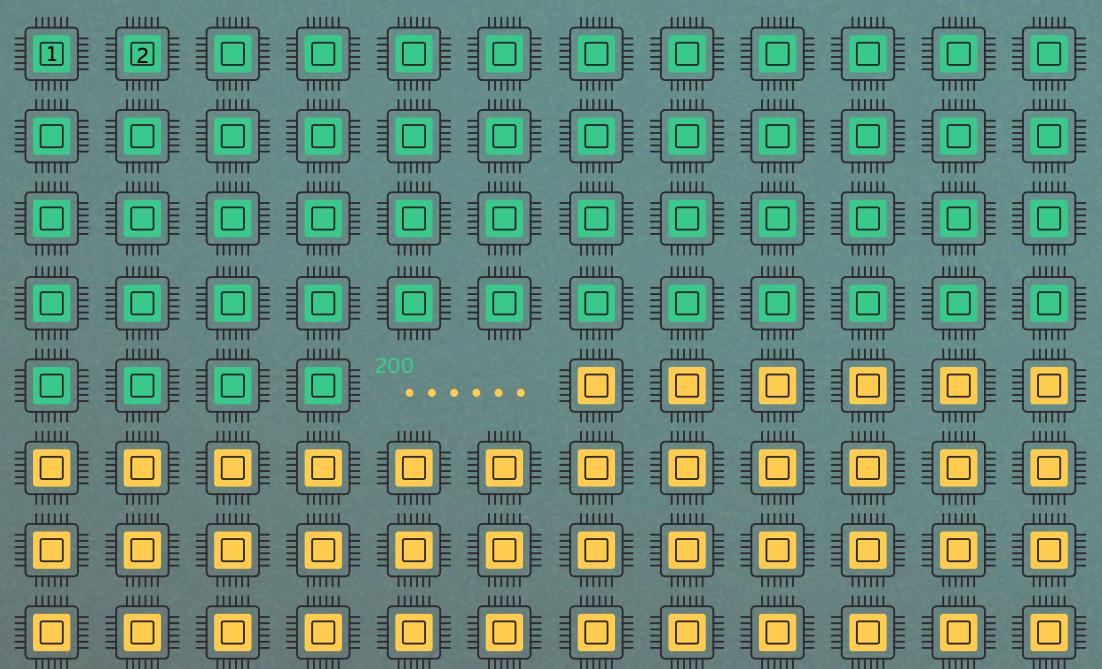
❓ Why shuffle partitions matter?

- Shuffle partitions run JOIN/ GROUP BY operations (Spark default: 200 partitions)
- Each partition ran by one core.
- Ideal processing per partition: ~1-200Mb

Consequences of not effective partition handling for your Spark job:

1. Slow completion time
 2. Cluster underutilization
- ie. a 1000 core cluster running with default setup: only **200** partitions work while the other **800** are idle.
- Larger processing per working core makes the whole cluster slower and underutilized

1000 core cluster w/ default partition setup



Working | Idle



Albert Campillo

  **Repost**

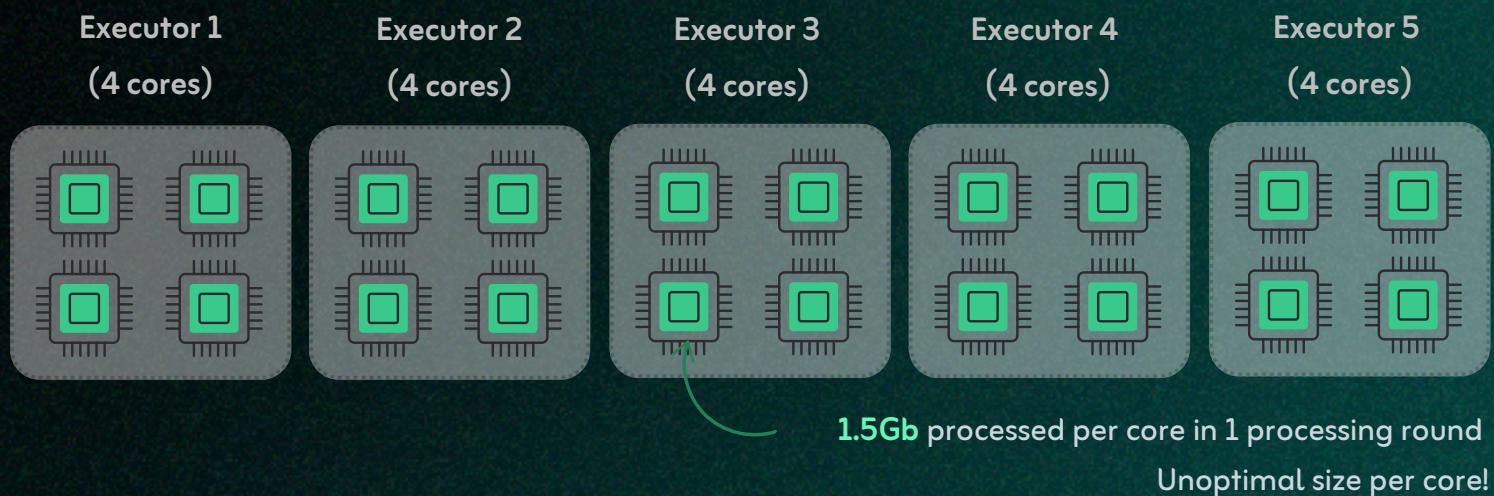
SCENARIO 1: The data partition per shuffle is very large (300 GB)

- 5 executors, 20 cores (4 cores per executor); default shuffle partitions (`spark.sql.shuffle.partition = 200`)
- Amount of data shuffled: 300Gb (`shuffle.write = 300`)

Case 1

No change in DSP (= 200)

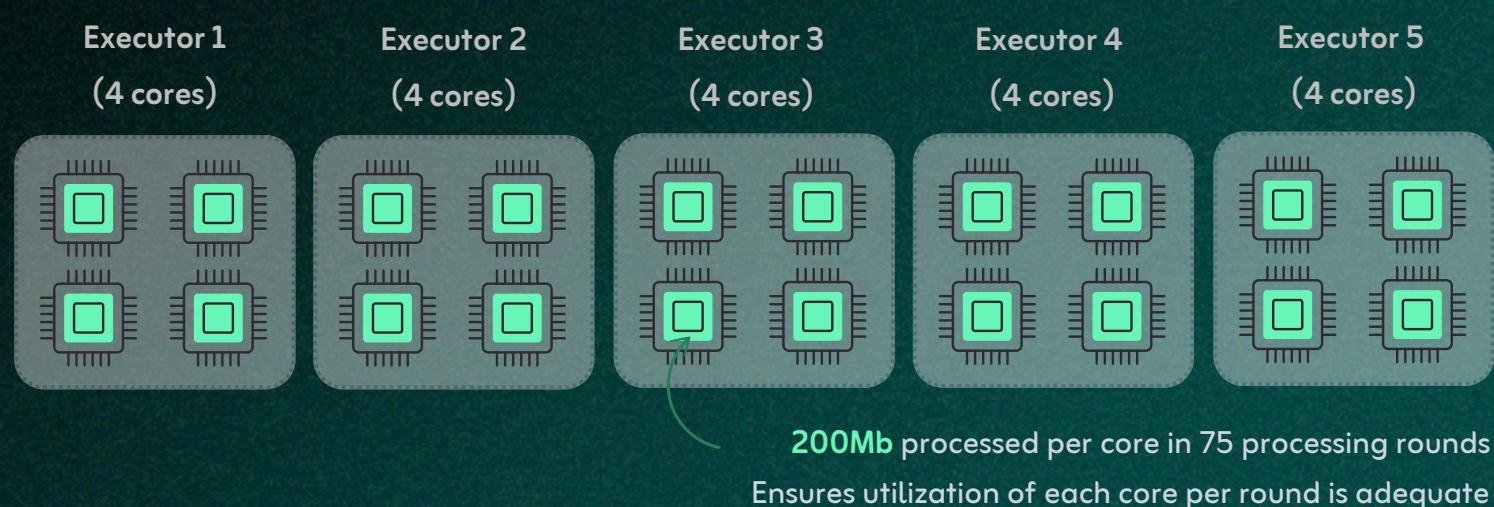
- Size per shuffle partition:
 $(300\text{Gb} / 200 \text{ DSP}) = 1.5\text{Gb/core}$
- 1 processing round
($1.5\text{Gb} \times 20 \text{ cores}$).
Unoptimal size per core!



Case 2

Change SP based on ideal processing size (ie. 200 Mb)

- $200\text{mb/SP} : (300\text{Gb} / X \text{ DSP}) = 1500 \text{ SP}$
- 75 processing rounds
($1500 \text{ SP} / 20 \text{ cores}$)



SCENARIO 2: The data partition per shuffle is very small (50 MB)

- 3 executors, 12 cores (4 cores per executor); default shuffle partitions (`spark.sql.shuffle.partition = 200`)
- Amount of data shuffled: 50Mb (`shuffle.write = 50`)

Case 1

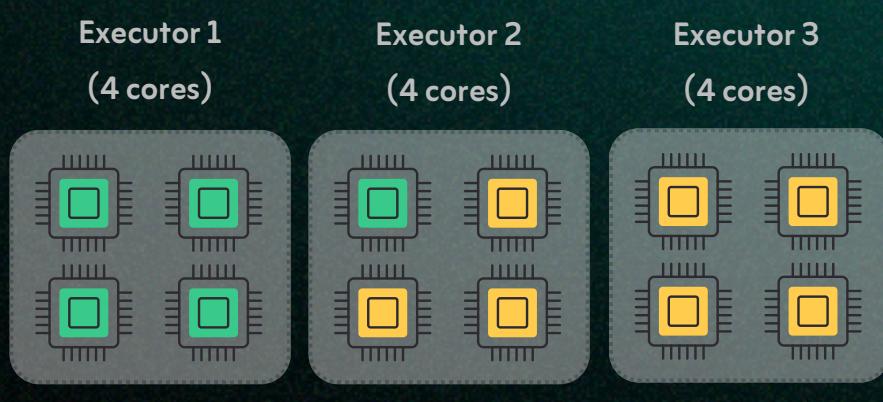
No change in DSP (= 200) --> Size per shuffle partition: $(50\text{Mb} / 200 \text{ DSP}) = 250\text{Kb/core}$ (very small vs. optimal ~200Mb)

Case 2

Adjust the SP to higher value (ie. 10Mb/ core)

- Number of SP = $50\text{Mb} / 10\text{Mb} = 5 \text{ SP}$

Underutilized, as 7 cores are **idle**

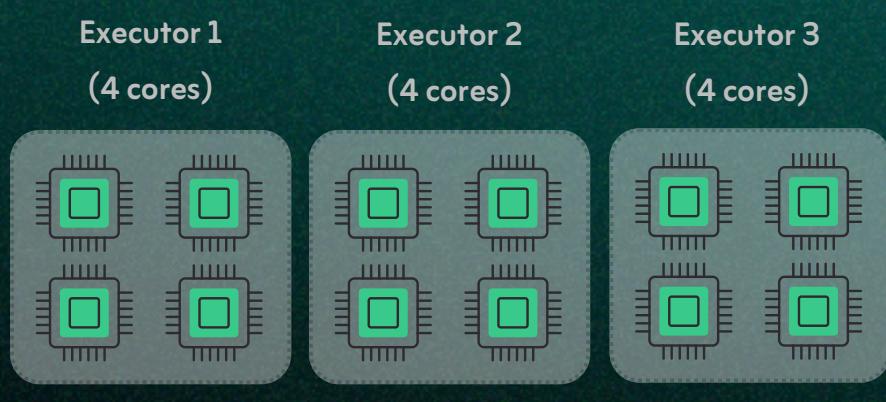


Case 3

Adjust the number of SP (all 12 cores)

- Size of SP = $50\text{Mb} / 12 \text{ cores} = \sim 4.16\text{Mb/core}$

Ensures job completion is faster (more cores working)



Albert Campillo

