

To process 25 GB data in spark

1) How many CPU cores are required?

2) How many executors are required?

3) How much each executor memory is required?

4) What is total memory for all executor is required?

$$① \quad 25 \text{ GB} \rightarrow 25 \times 1024 \text{ MB} = 25600 \text{ MB}$$

1 core = 1 Task

spark - maxBytesPerPartition  $\rightarrow 128 \text{ MB}$

$$\text{No. of Partitions} = 25600 / 128 = 200 \rightarrow 200 \text{ Task}$$

$$\text{No. of CPU cores} = \text{No. of Partition} = 200$$

→ by default spark creates each partition for each block (128m)

We can increase the Number of Bytes per partition.

② How many executors will be required to process 25 GB of data?

Based on research's ideal number of ~~cores~~ <sup>maximum core</sup> for each executor is

Avg CPU cores for each executor = 4

2 - 5 → 4

Executor → 2 - 5

Total no of executors =  $200 / 4$

cores

③ How much each executor memory is required to process 25 GB of data.

CPU Core for each executor = 4

Expected memory for each ~~core~~ =

Memory for each executor =  $4 \times 512 \text{ MB}$

= 2 GB

$4 \times \text{default Partition size}$

$4 \times 128 \text{ MB}$

= 512 MB

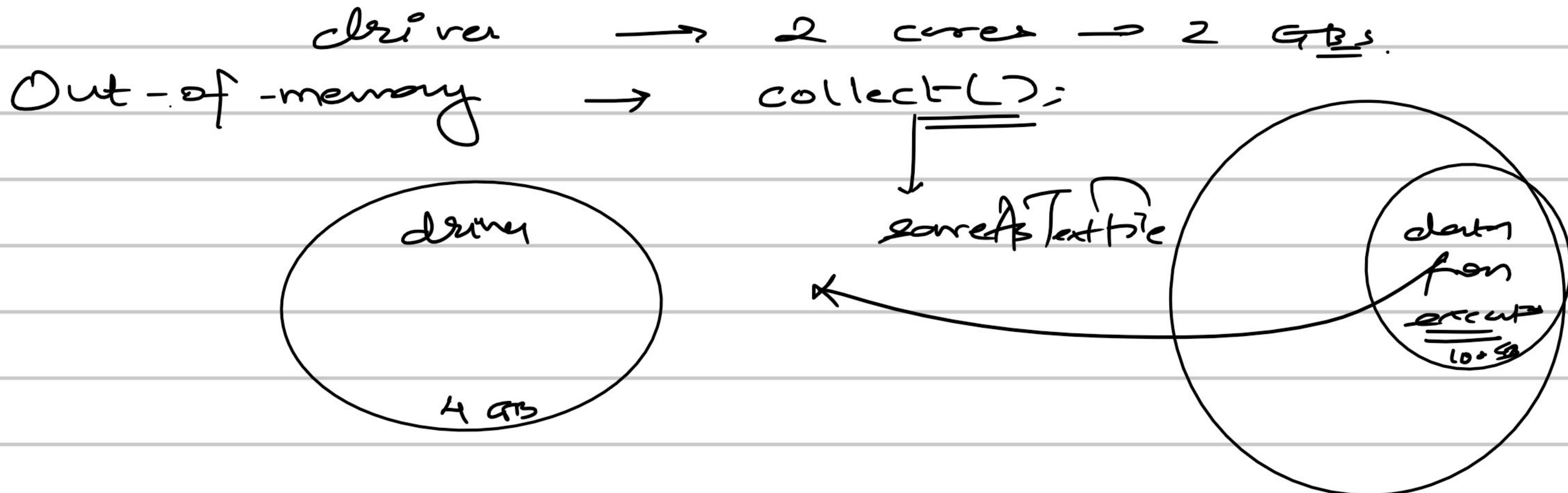
④ What is the total memory required to process 25 GB?

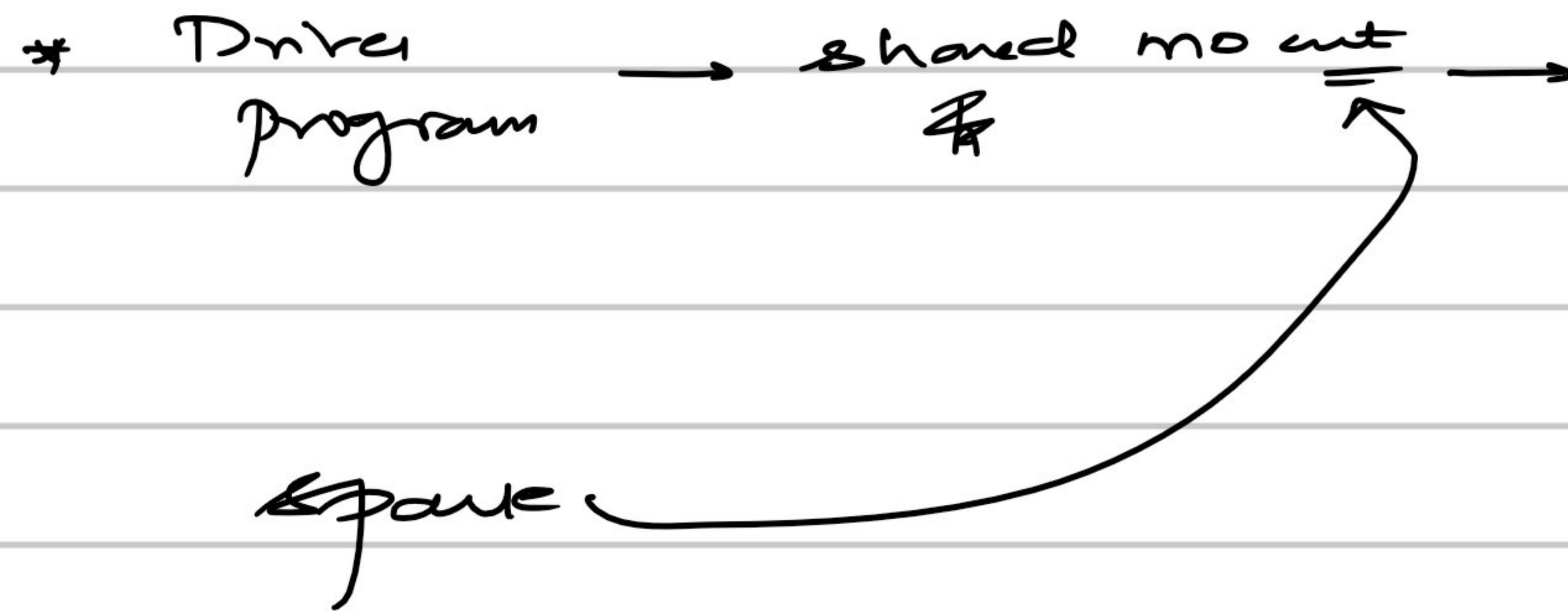
Total No. of executors = 50  
Memory for each = 2 GB.

Total memory for all executors = 100 GB

Executor memory is not less than 1.5 times of Spark reserved memory - Single Core executor memory should not be less than  $1.5 \times 300 \text{ MB} = 450 \text{ MB}$

⑤ How much memory is one to be allocated for driver?





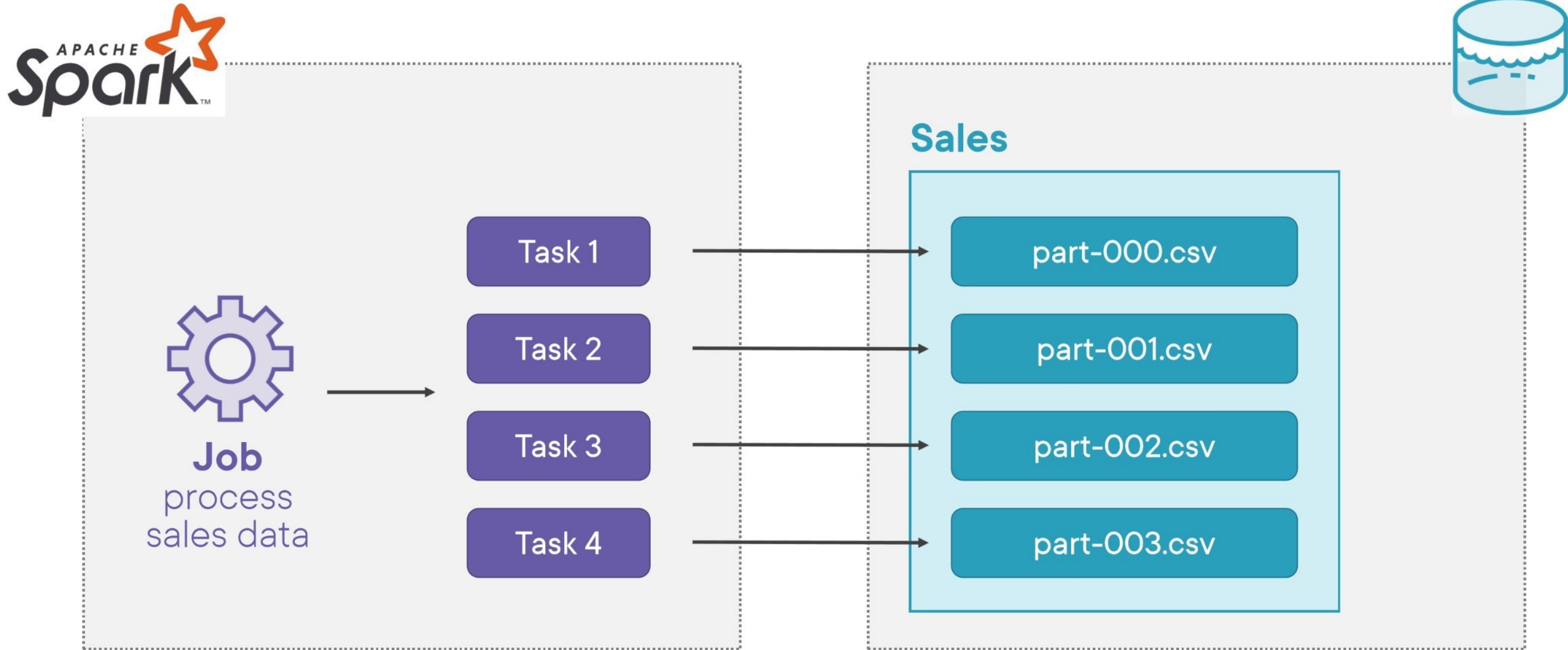
app → spark context → driver —  
executor

# Need for Delta Lake with Spark

---

**Data Lake** is a central repository to  
store all types of data at any scale,  
in the form of files

# How Spark Writes to Data Lake?



# ACID

Guarantees on Transactions

Each transaction in a Database provides ACID guarantees

# ACID

## Guarantees on Transactions

### Atomicity

All or no changes  
are written

Id	Name	Salary
1	A	10000
2	B	17000
3	C	12000

4	D	11000
5	E	14000



# ACID

## Guarantees on Transactions

### Consistency

Data always remain  
in valid state

Withdraw 200  
where Customer = 3

Customer	Balance
1	500
2	30
X 3	100

Constraint  
 $\text{Balance} \geq 0$

# ACID

## Guarantees on Transactions

### Isolation

Transaction must run isolated from other processes



**Writer**

4	D
5	E

Id	Name
1	A
2	B
3	C



**Reader**

# ACID

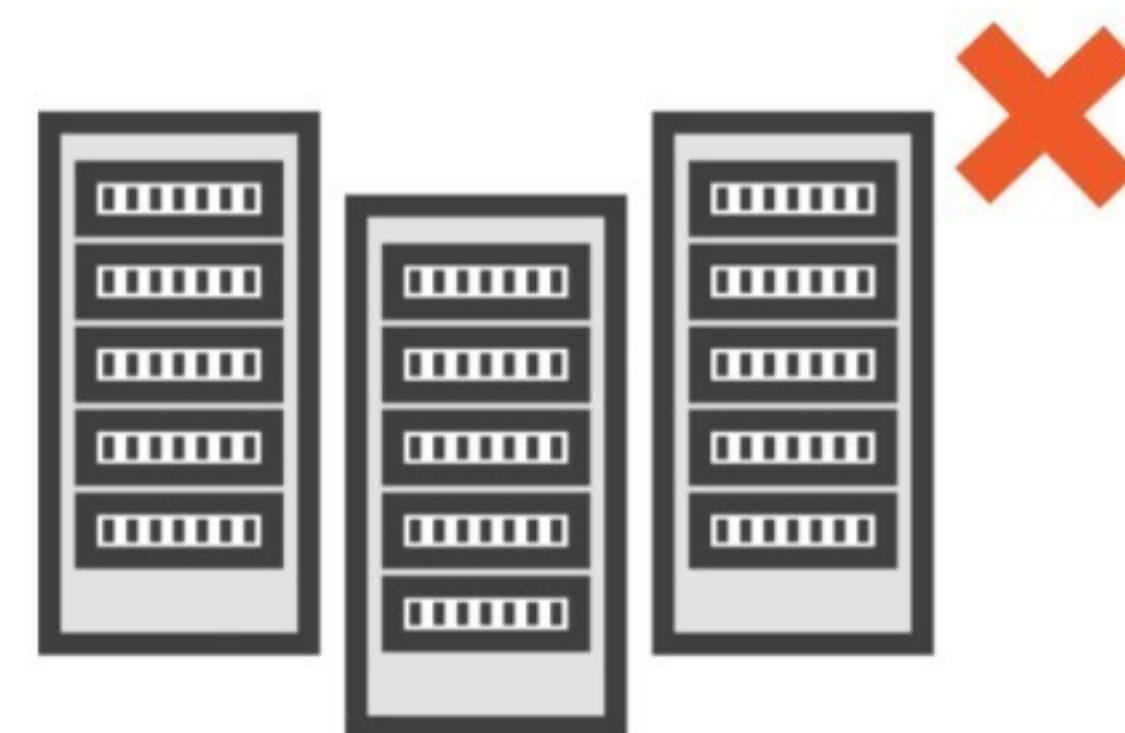
## Guarantees on Transactions

### Durability

Once committed,  
data persists even  
if system fails



Id	Name
1	A
2	B
3	C



Data Lake does **NOT** provide  
ACID guarantees

# 1 – Job Failure in Appending Data



**Writer**

Append 2 more part files  
using 2 tasks

Job failed with runtime  
error – Part file 4 could not  
be written

Part File 4



**Reader**

Read the folder  
(reads part files 1, 2 & 3)

Reads inconsistent data

**Breaks Atomicity and Consistency!**

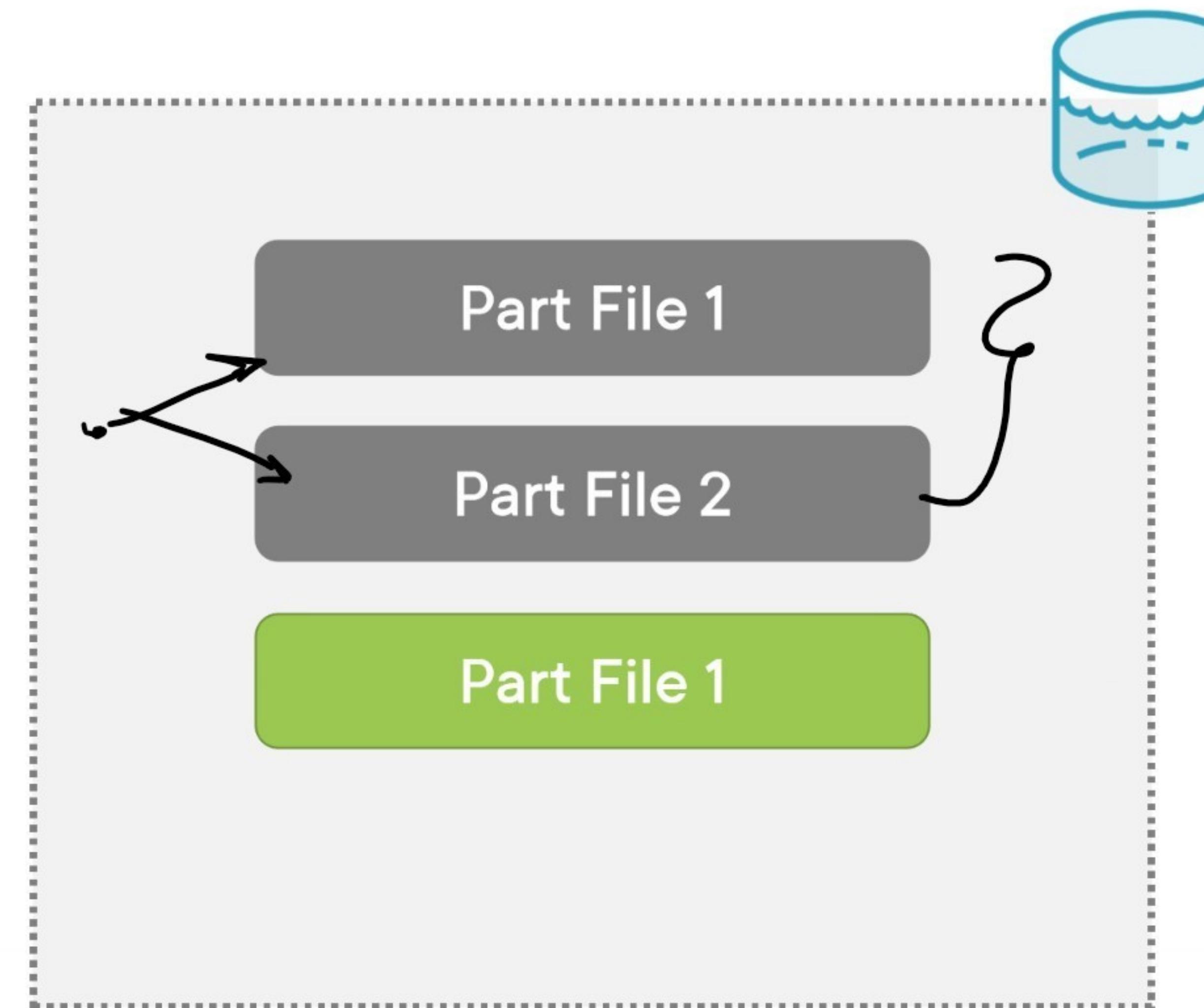
## 2 – Job Failure in Overwriting Data



**Writer**

Overwrite existing =  
Delete existing  
+ Write 2 new part files

Part File 2 



**Job failed with runtime error – New part file 2 could not be written**

Read the folder  
(reads part file 1 only)

**Reads inconsistent data + Previous data is lost**

**Breaks Atomicity, Consistency and Durability!**

# 3 – Simultaneous Read / Write



**Writer**

Write 2 more part files  
using 2 tasks

Only Part file 3 is written.  
Part 4 is still getting  
processed

Part File 4



**Reader**

Read the folder while  
writing is in progress  
(reads part files 1, 2 & 3)

Reads inconsistent data

**Breaks Consistency and Isolation!**

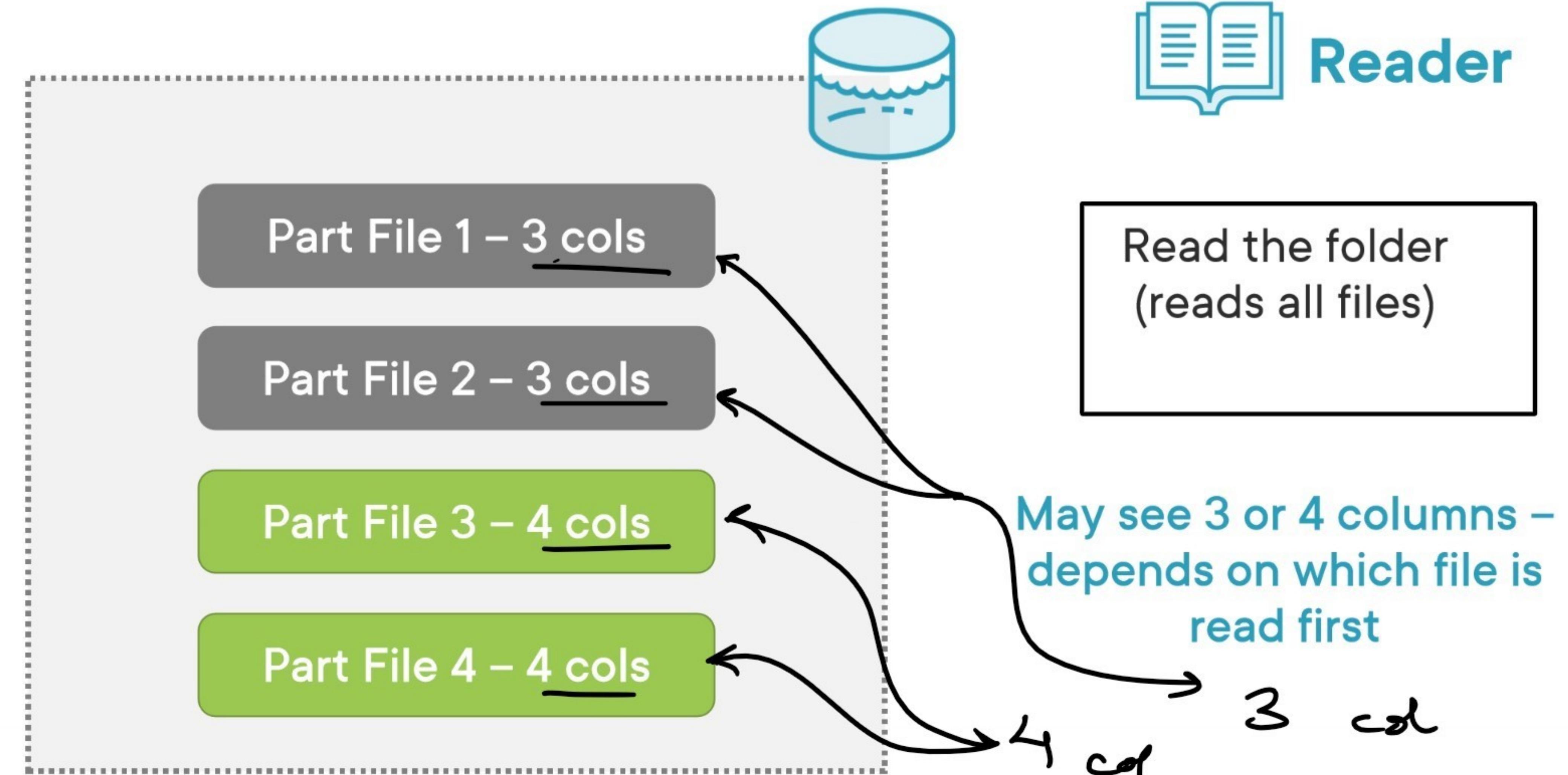
# 4 – Appending Data with New Schema



## Writer

Write 2 more part files  
with 4 output columns

New output files have 4  
columns. No schema  
validation before writing



**Breaks Consistency!**

# Challenges with Data Lake

## Data reliability issues

- Data corruption because of failures – no rollback!
- No data validation
- Consistency issues while reading data

## Handling Batch and Streaming data together is tough

### No updates / deletes / merge on files

- Difficult to implement GDPR / CCPA compliance

## Data quality issues

- Schema isn't verified before writing
- Cannot apply checks on data

## Query performance issues

## Difficult to maintain historical versions of data

Delta Lake can help us  
solve these challenges!

But is Delta Lake the only option?

Apache Iceberg  
Apache Hudi

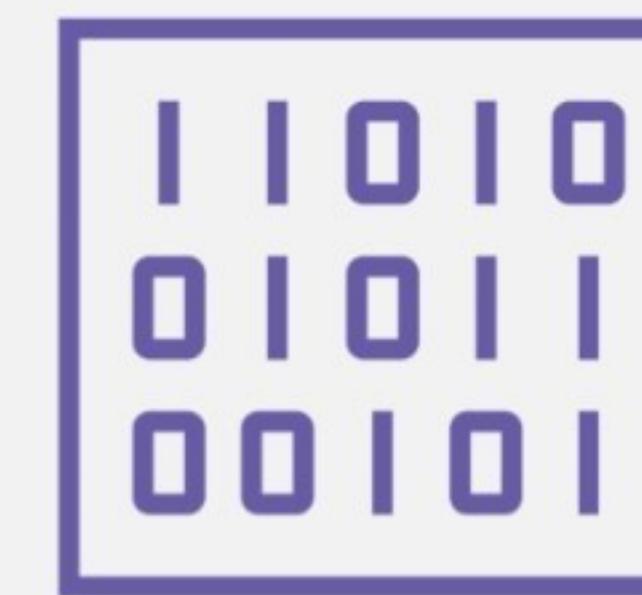
# How Delta Lake Works?

---

**Delta Lake is an open-source  
storage layer that brings reliability  
to Data Lakes**

# Writing Data in Parquet Format

`-format ("parquet");`

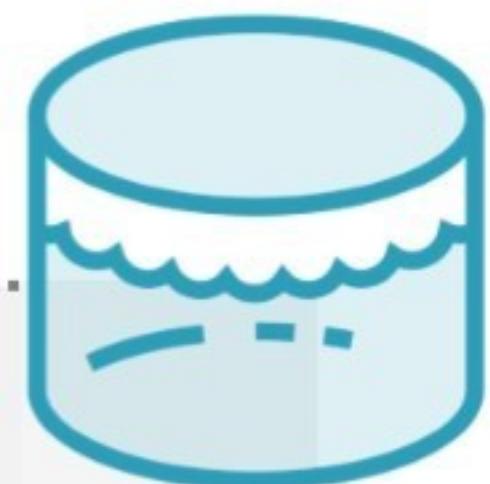


DataFrame

parquet  
format

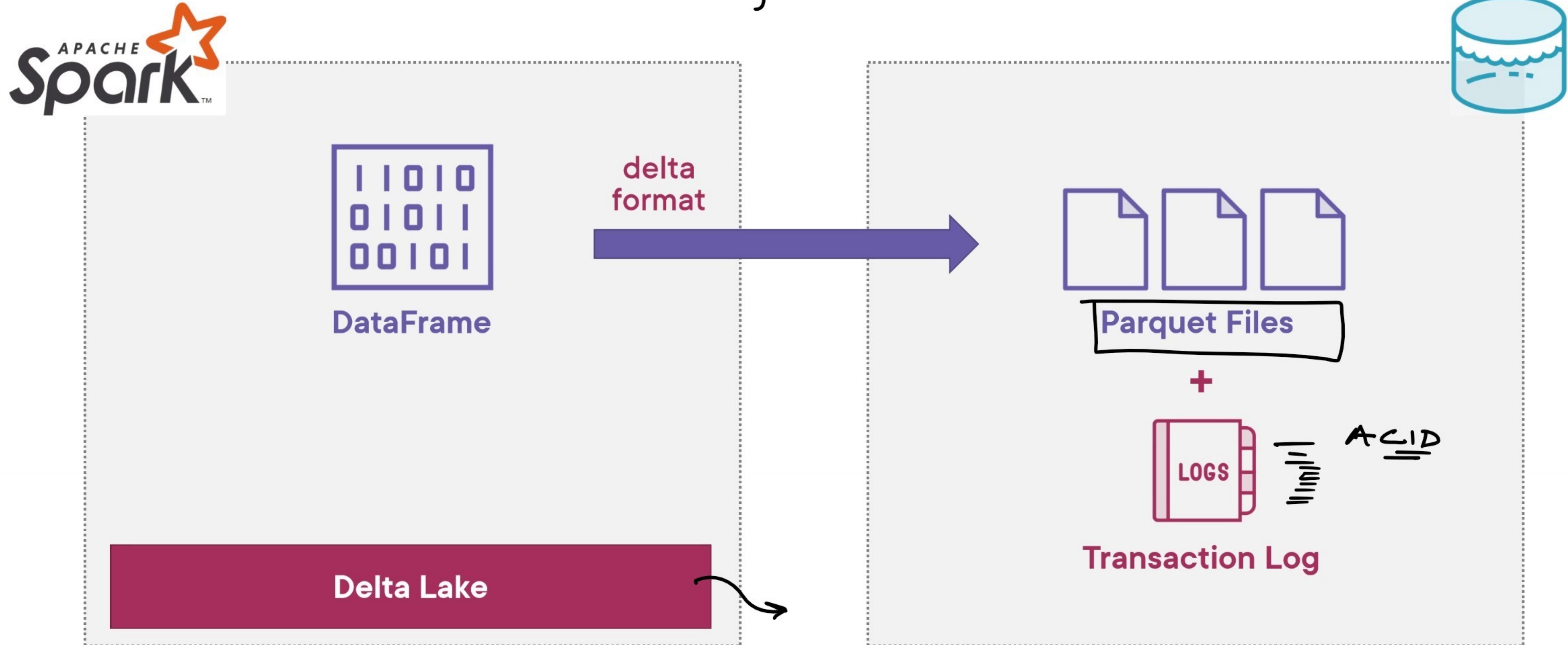


Parquet Files



# Writing Data in Delta Format

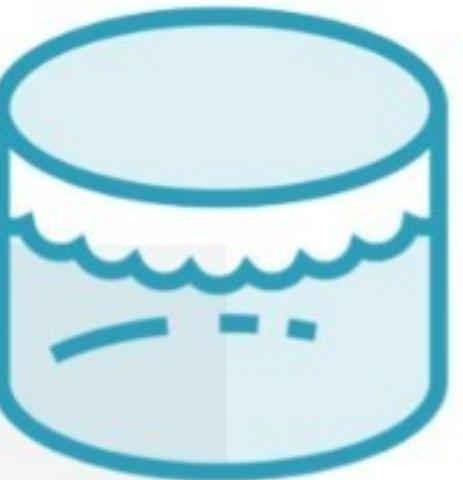
`.format("delta"):`



```
(  
    dataframe  
        .write  
            .format("delta")      # like other formats - csv, parquet etc.  
            .save(filepath)  
)
```

Save DataFrame in Delta Format

# Customer (folder)



**Write Operation 1**  
(write 2 part files)

part-000.parquet

part-001.parquet

\_delta\_log (subfolder)

000.json

**Write Operation 2**  
(append 1 part file)

part-002.parquet

001.json

**Write Operation 3**  
(append 1 part file)

part-003.parquet

002.json

**For each write operation:**

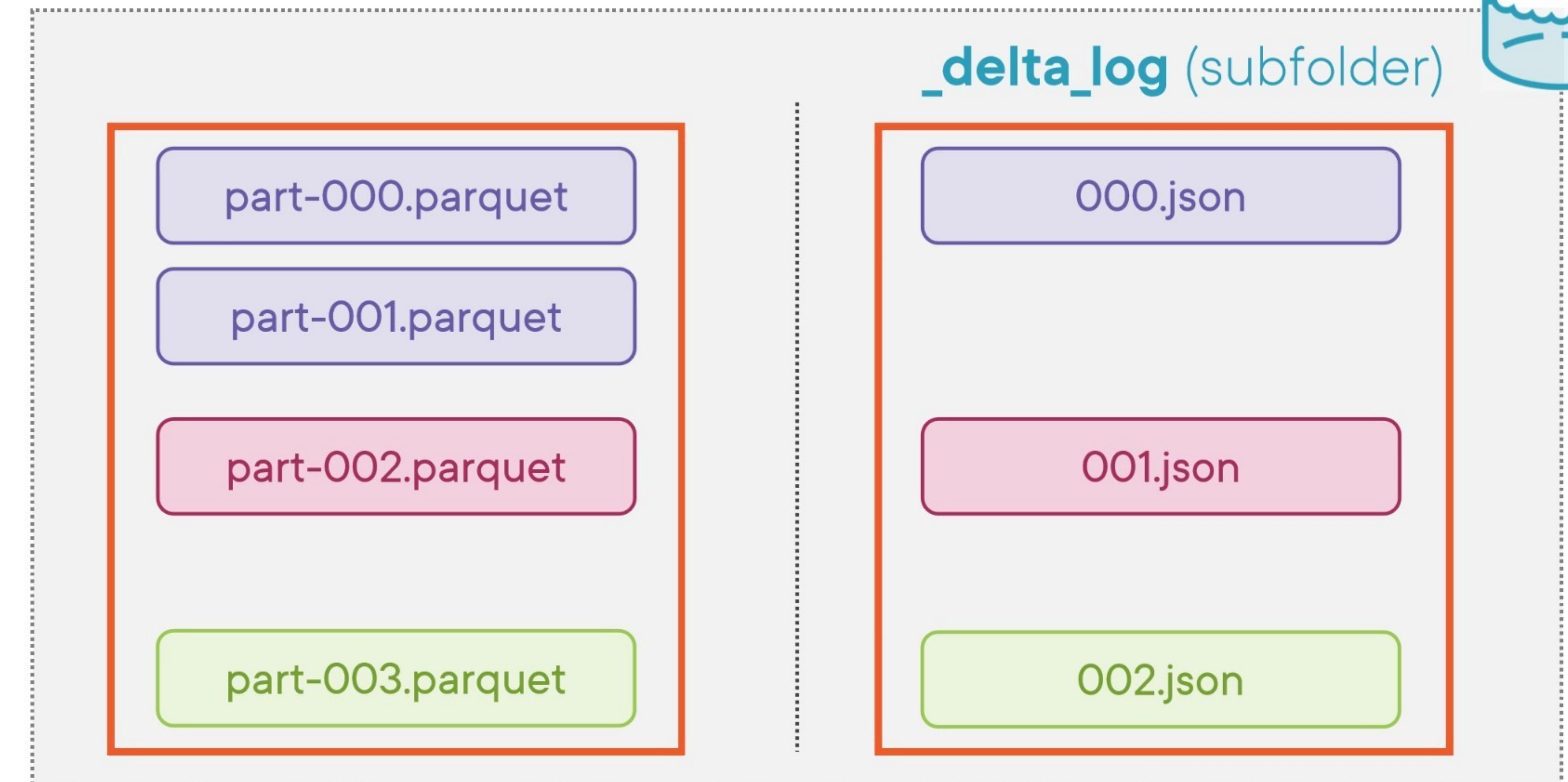
- Part files are written first
- A transaction log file is added to \_delta\_log folder in JSON format

## Customer (folder)

### Read Operation

Reads log files -  
000.json, 001.json,  
002.json

Reads 4 part files,  
based on log  
information



### For each read operation:

- Transaction log files are read first
- Part files are then read based on log files

## Customer (folder)

**Create Operation**  
(with 2 part files)

Id	Name
1	A
2	B

Part 1

Id	Name
3	C
4	D

Part 2

## \_delta\_log (subfolder)

Operation	File Name
Add <del>C</del>	Part 1 <del>C</del>
Add	Part 2

000.json

**Insert Operation**  
(append 1 part file)

Id	Name
5	E
6	F

Part 3

Operation	File Name
Add	Part 3

001.json

**Update Operation**  
(Change name from  
A to AA where Id = 1)

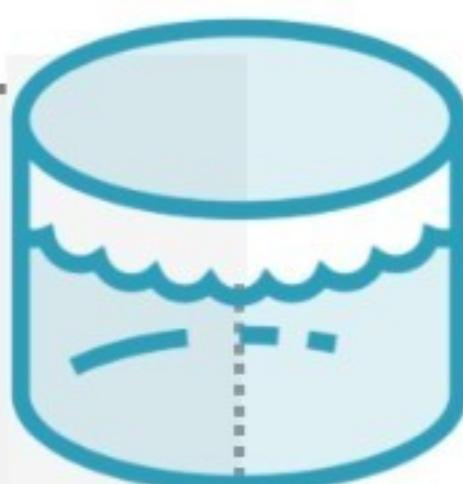
Id	Name
1	AA
2	B

Part 4

Add modified row  
Copy  
unchanged rows

Operation	File Name
Remove	Part 1
Add	Part 4

002.json



# Customer (folder)

**Read all records**

~~Part 1~~

**Part 2**

**Part 3**

**Part 4**

Id	Name	Id	Name
1	A	3	C
2	B	4	D

Part 1      Part 2

Id	Name
5	E
6	F

Part 3

Id	Name
1	AA
2	B

Part 4

# \_delta\_log (subfolder)

Operation	File Name
Add	Part 1
Add	Part 2

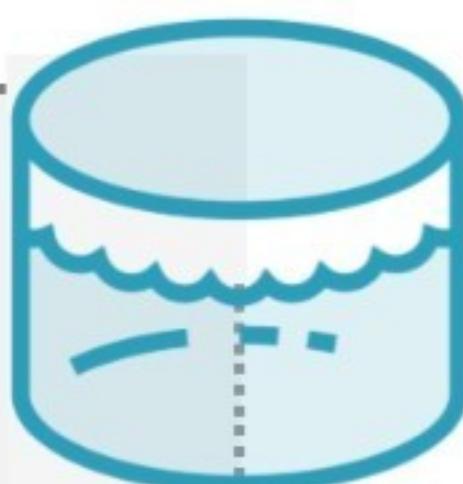
000.json

Operation	File Name
Add	Part 3

001.json

Operation	File Name
Remove	Part 1
Add	Part 4

002.json



## Customer (folder)

**Read all records**

Id	Name
3	C
4	D
5	E
6	F
1	AA
2	B

Id	Name
1	A
2	B
3	C
4	D

Part 1

Id	Name
3	C
4	D

Part 2

Id	Name
5	E
6	F

Part 3

Id	Name
1	AA
2	B

Part 4

## \_delta\_log (subfolder)

Operation	File Name
Add	Part 1
Add	Part 2

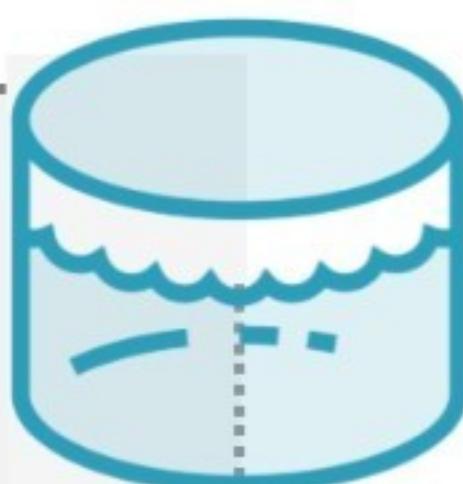
000.json

Operation	File Name
Add	Part 3

001.json

Operation	File Name
Remove	Part 1
Add	Part 4

002.json



# Delta Lake Features

## **Provides ACID Guarantees**

- No data corruption because of failures
- Data consistency while reading data

## **Perform inserts / updates / deletes**

## **Schema enforcement and evolution**

## **Protect data using Time Travel**

## **Handle batch & streaming data together**

## **Apply data quality checks**

## **Performance improvements using statistics**

***...and much more***

# ACID Guarantees on Delta Lake

---

# 1 – Job Failure in Appending Data

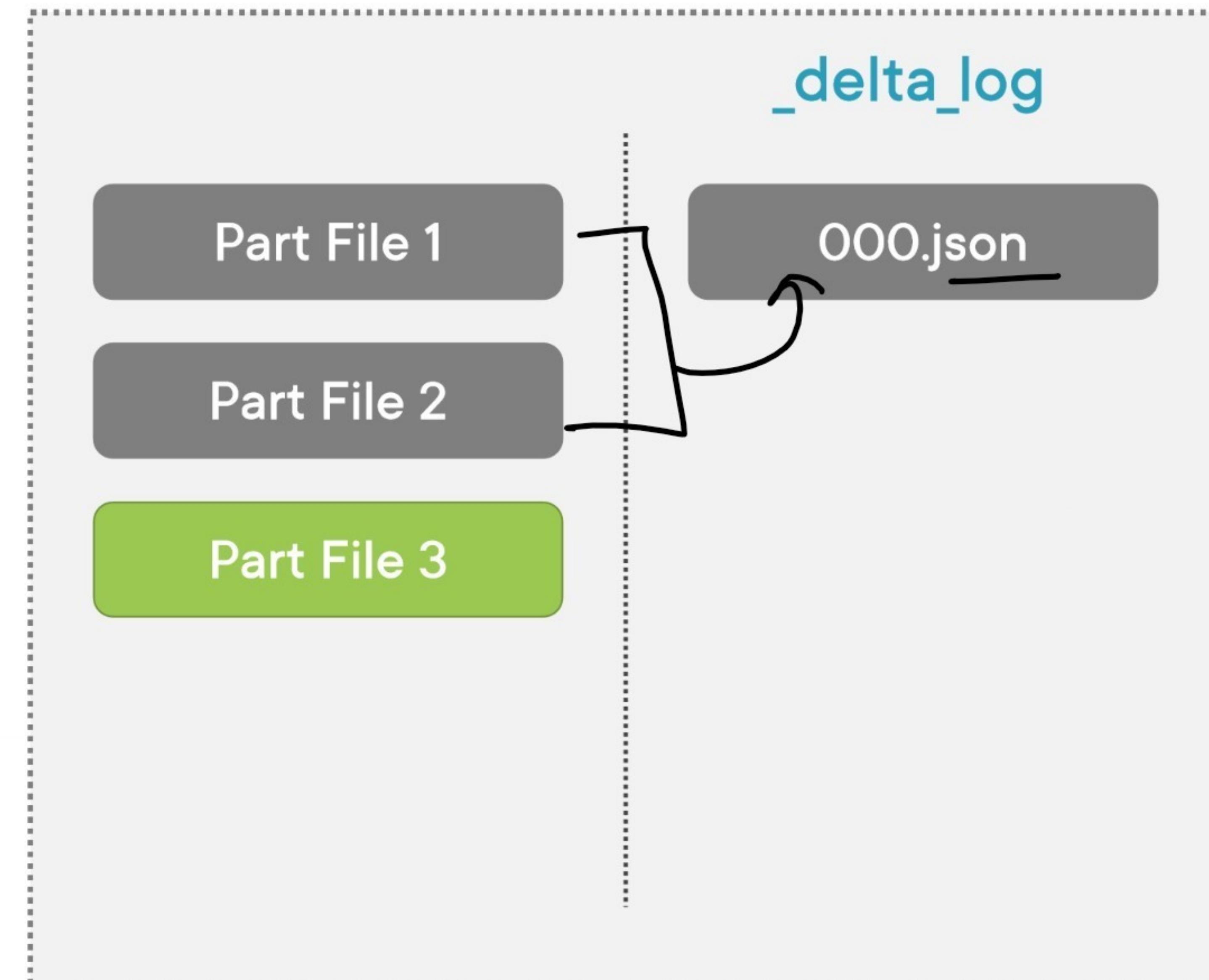


**Writer**

Append 2 more part files  
using 2 tasks

Job failed with runtime  
error – Part file 4 could  
not be written

Part File 4



**Reader**

Reads log file  
(000.json)

Only reads Part files 1 & 2  
Reads consistent data

## 2 – Job Failure in Overwriting Data



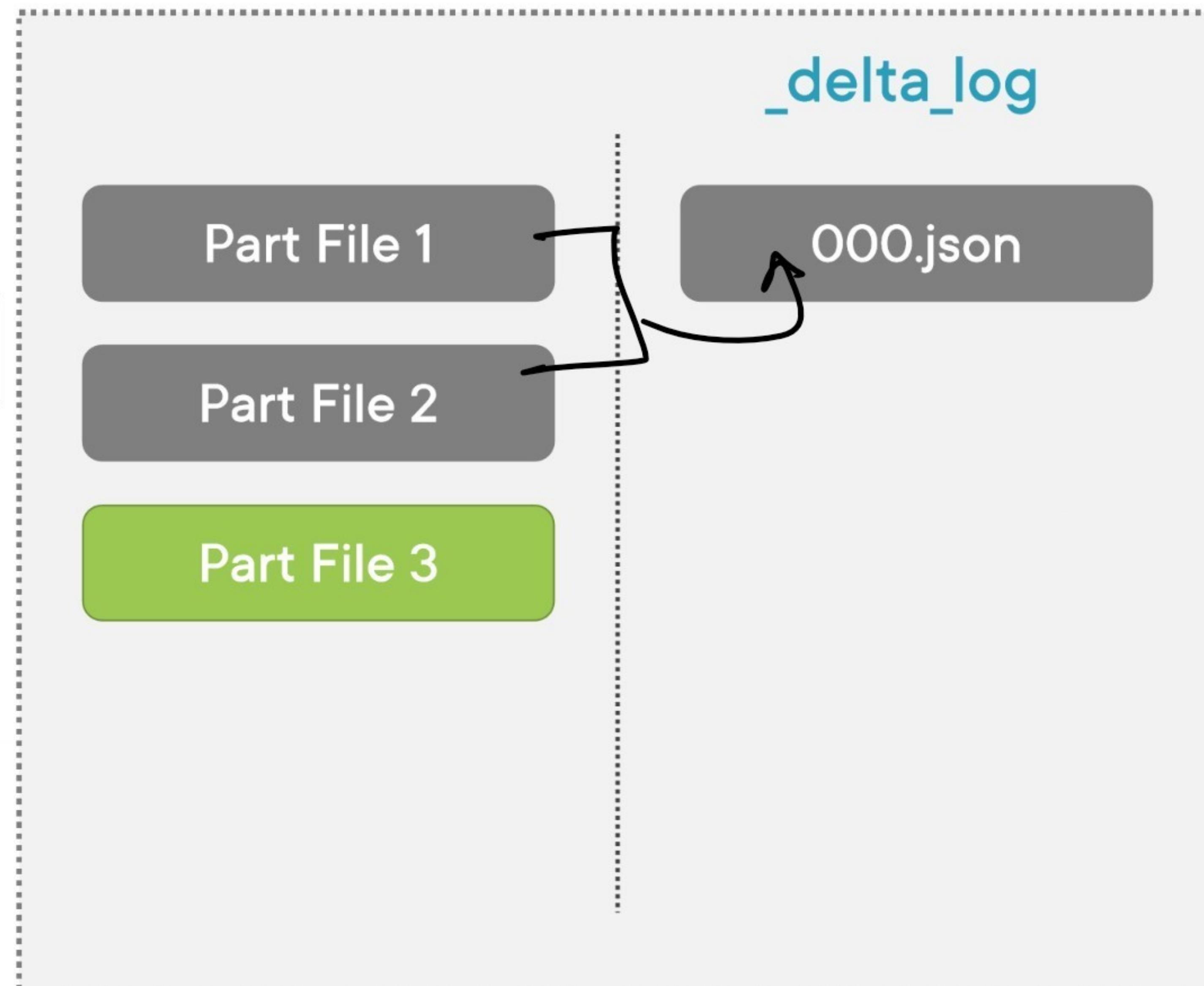
### Writer

Overwrite existing data by writing 2 new part files

**Existing files are not deleted!**

Job failed with runtime error – New part file 4 could not be written

Part File 4



### Reader

Reads log file (000.json)

Only reads Part files 1 & 2  
Reads consistent data  
+ Previous data is not lost

# 3 – Simultaneous Read / Write

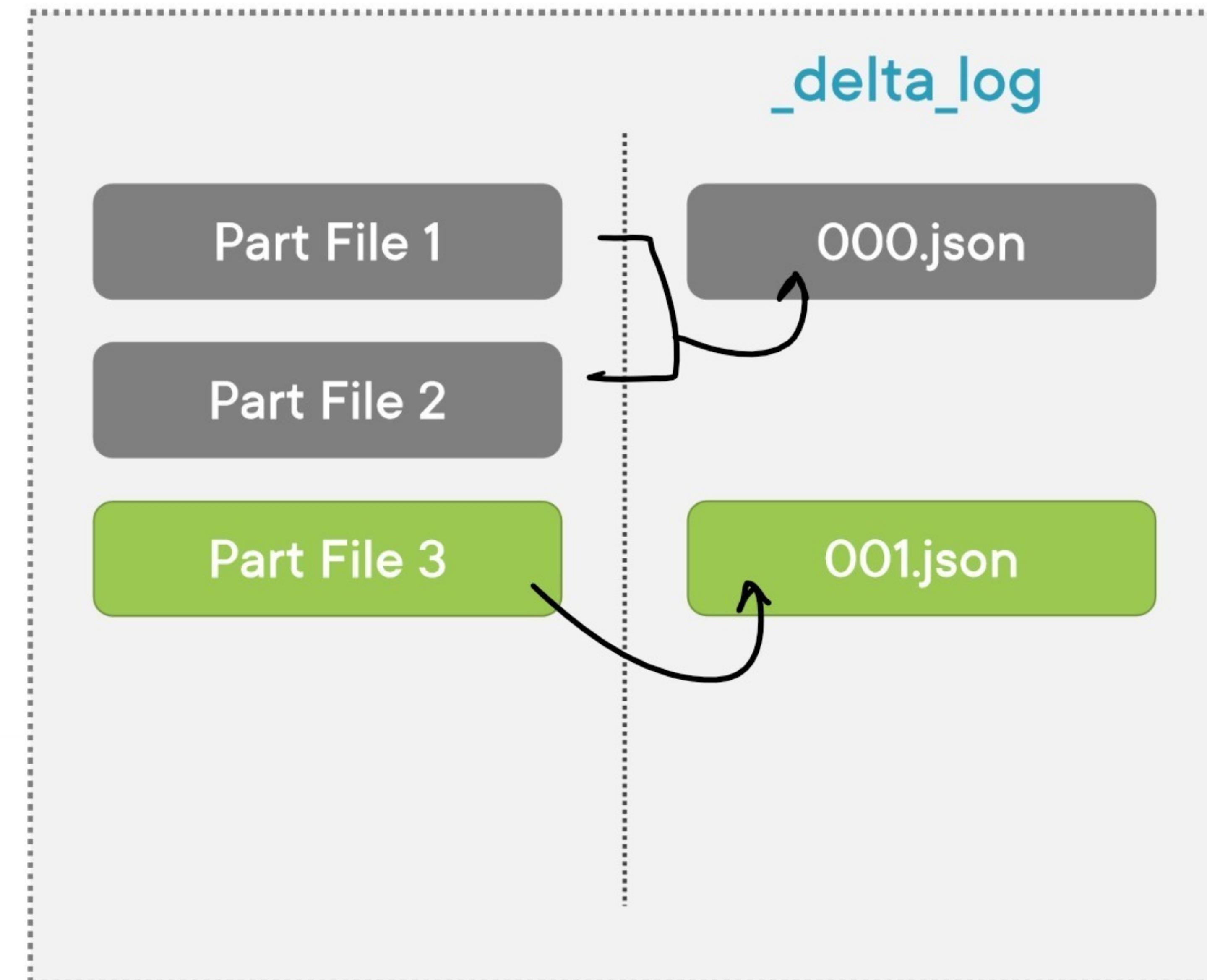


**Writer**

Write 2 more part files  
using 2 tasks

Only Part file 3 is written.  
Part 4 is still getting  
processed

Part File 4



**Reader**

Reads log file while writing  
is in progress  
(reads 000.json)

Only reads Part files 1 & 2  
No dirty reads while  
writing is in progress

# Creating Delta Tables

---

# Sample Transaction Log Entry

```
{"commitInfo": {"timestamp": 1640067393982, "userId": "1472626815582251",  
  "operation": "WRITE", "operationParameters": {"mode": "Overwrite"},  
  "operationMetrics": {"numFiles": "2", "numOutputRows": "200", "numOutputBytes": "14629"} }  
  
{ "metaData": { "schemaString": {  
    "fields": [  
      { "name": "Day", "type": "integer", "nullable": true, "metadata": {} },  
      { "name": "RideId", "type": "integer", "nullable": true, "metadata": {} },  
      { "name": "Amount", "type": "double", "nullable": true, "metadata": {} }  
    ]  
  } }  
  
{ "add": { "path": "part-000.parquet", "size": 7284 } }  
{ "add": { "path": "part-001.parquet", "size": 7345 } }
```

# Sample Transaction Log Entry

```
{"commitInfo": {"timestamp": 1640067393982, "userId": "1472626815582251",  
  "operation": "WRITE", "operationParameters": {"mode": "Overwrite"},  
  "operationMetrics": {"numFiles": "2", "numOutputRows": "200", "numOutputBytes": "14629"}}}
```

```
{"metaData": {"schemaString": {  
    "fields": [  
      {"name": "Day", "type": "integer", "nullable": true, "metadata": {}},  
      {"name": "RideId", "type": "integer", "nullable": true, "metadata": {}},  
      {"name": "Amount", "type": "double", "nullable": true, "metadata": {}}  
    ]  
  }}}
```

```
{"add": {"path": "part-000.parquet", "size": 7284 }},  
 {"add": {"path": "part-001.parquet", "size": 7345 }},
```

# Sample Transaction Log Entry

```
{"commitInfo": {"timestamp": 1640067393982, "userId": "1472626815582251",  
    "operation": "WRITE", "operationParameters": {"mode": "Overwrite"},  
    "operationMetrics": {"numFiles": "2", "numOutputRows": "200", "numOutputBytes": "14629"}}}
```

```
{"metaData": {"schemaString": {  
    "fields": [  
        {"name": "Day", "type": "integer", "nullable": true, "metadata": {}},  
        {"name": "RideId", "type": "integer", "nullable": true, "metadata": {}},  
        {"name": "Amount", "type": "double", "nullable": true, "metadata": {}}  
    ]  
}}}
```

```
{"add": {"path": "part-000.parquet", "size": 7284 }},  
{"add": {"path": "part-001.parquet", "size": 7345 }}
```

# Sample Transaction Log Entry

```
{"commitInfo": {"timestamp": 1640067393982, "userId": "1472626815582251",  
    "operation": "WRITE", "operationParameters": {"mode": "Append"},  
    "operationMetrics": {"numFiles": "2", "numOutputBytes": "14629", "numOutputRows": "200"}}}
```

```
{"metaData": {"schemaString": {  
    "fields": [  
        {"name": "Day", "type": "integer", "nullable": true, "metadata": {}},  
        {"name": "RideId", "type": "integer", "nullable": true, "metadata": {}},  
        {"name": "Amount", "type": "double", "nullable": true, "metadata": {}}  
    ]  
}}}
```

```
{"add": {"path": "part-000.parquet", "size": 7284 }},  
{"add": {"path": "part-001.parquet", "size": 7345 }}
```

# Inserting Data to Delta Table

---

# Options to Insert Data

**INSERT Command**  
(SQL)

**Append DataFrame**  
(PySpark / Scala)

# Performing DML Operations

---

# Applying Table Constraints

---

# Table Constraints

## **NOT NULL**

To prevent column from having  
NULL values

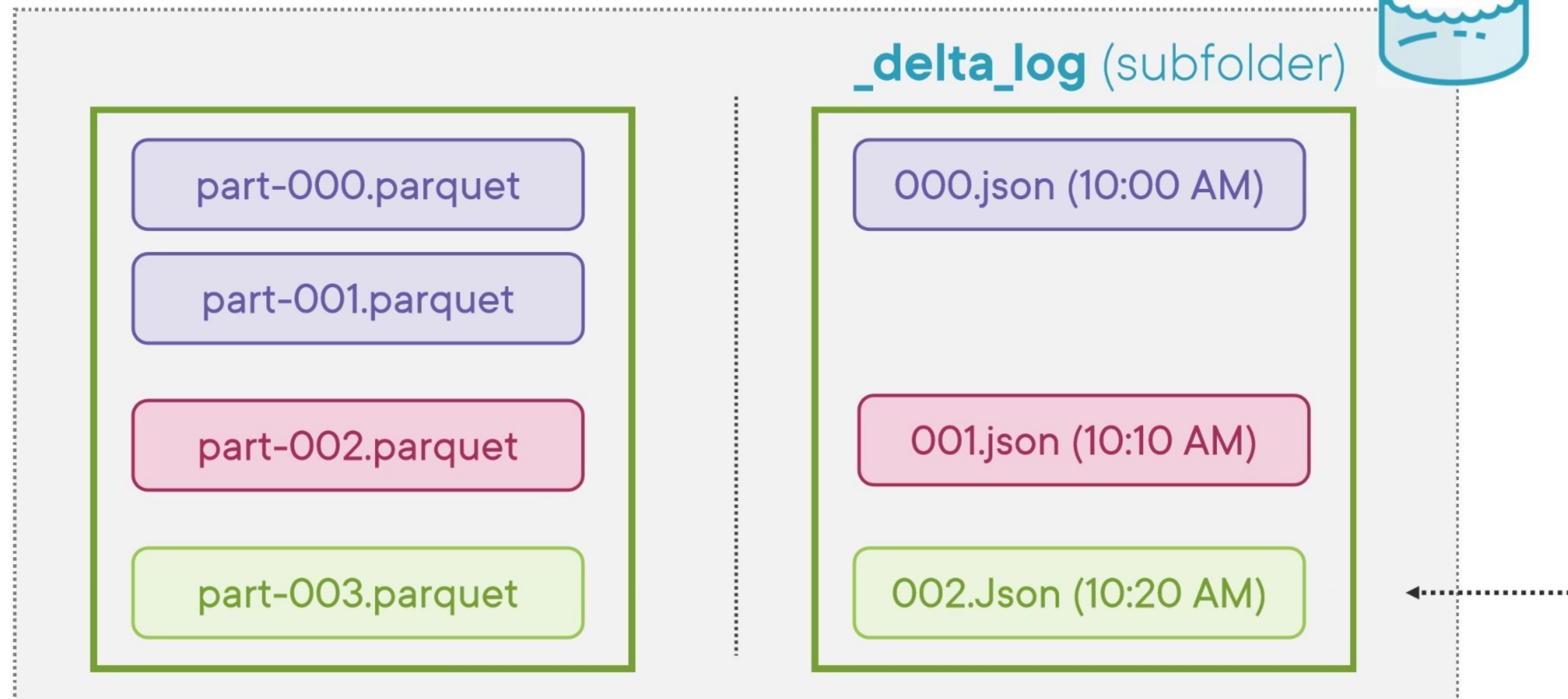
## **CHECK**

Define conditions to enforce on  
data in the table

# Accessing Data with Time Travel

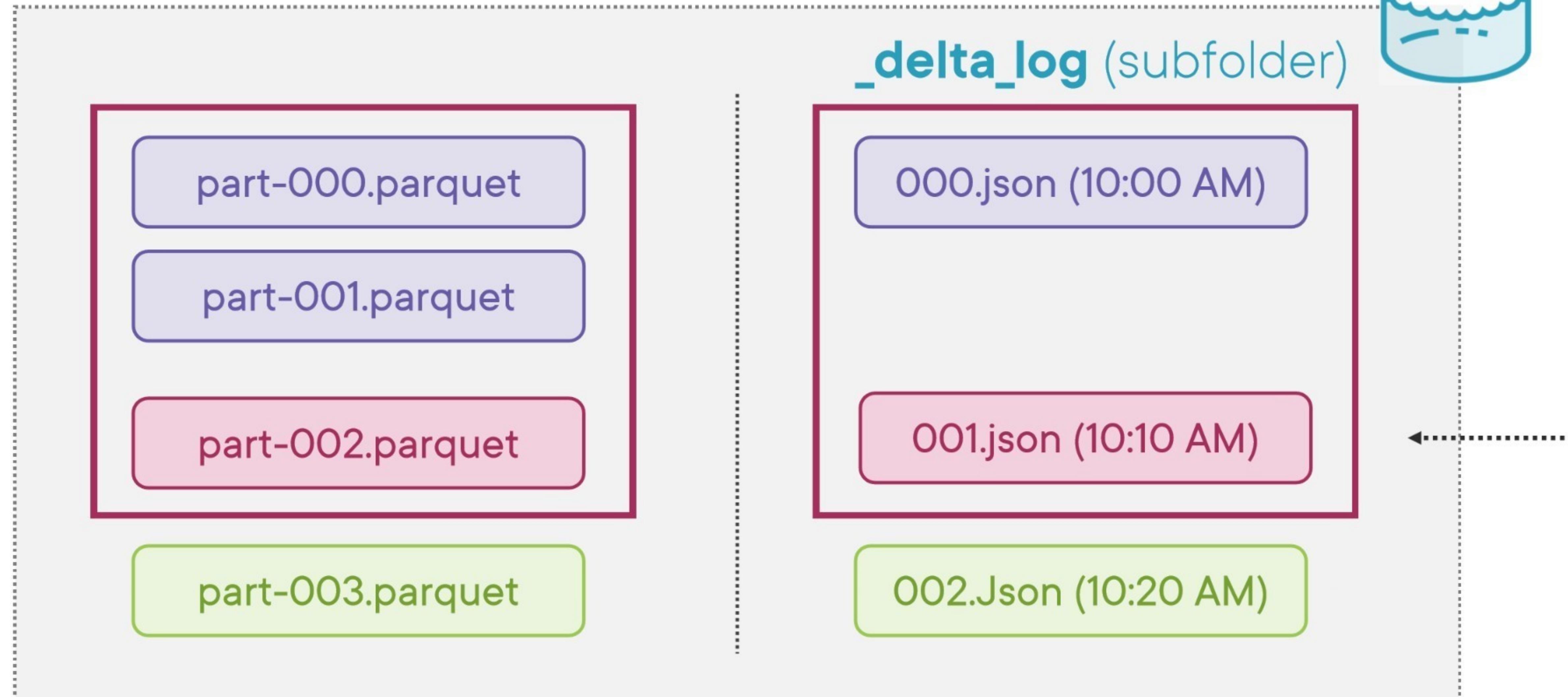
---

## Customers (folder)



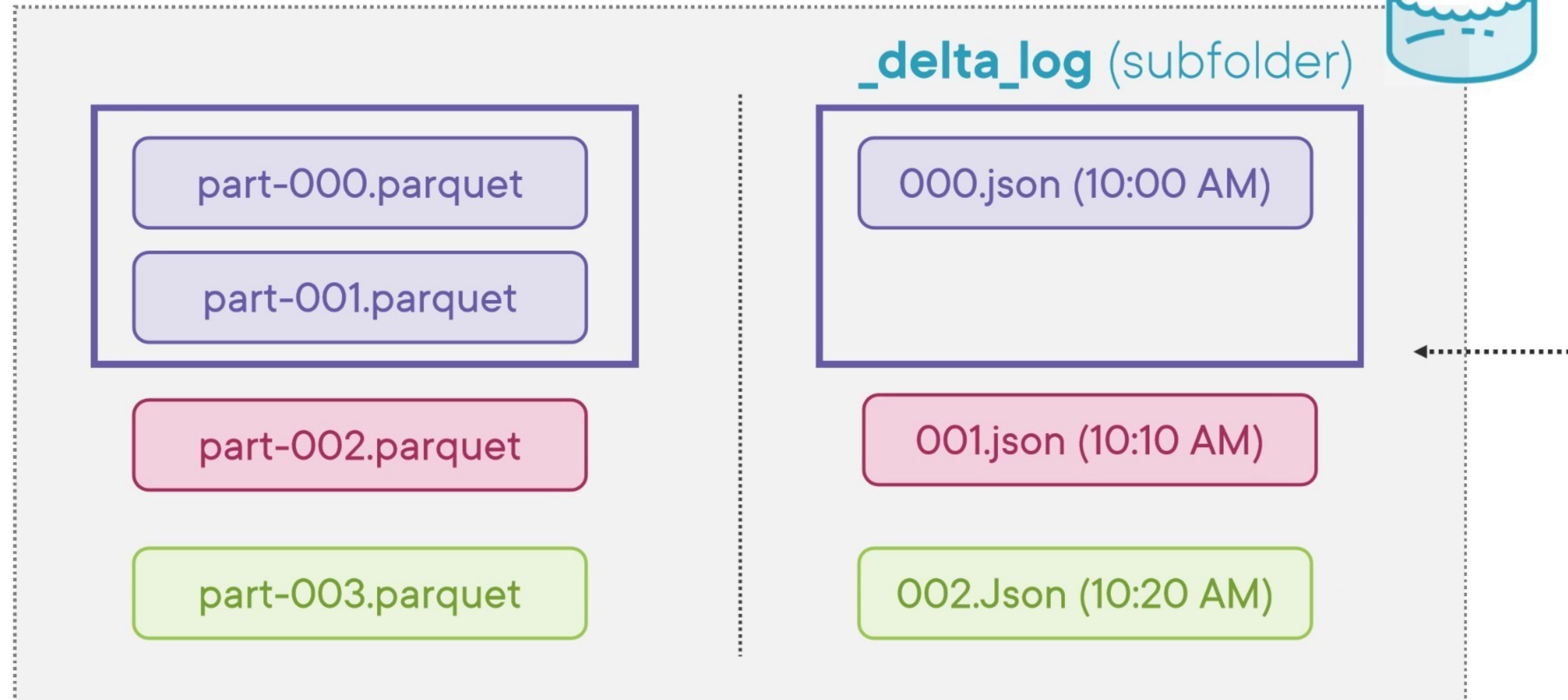
```
SELECT *\n\nFROM Customers
```

## Customers (folder)



```
SELECT *  
FROM Customers VERSION AS OF 1
```

## Customers (folder)



```
SELECT *
FROM Customers TIMESTAMP AS OF '2022-03-01 10:05'
```

**Time Travel** allows to access/restore previous snapshot of data, even if data has been modified or deleted

## Summary



**Delta Lake is storage layer bringing reliability to Lakes**

- Stores data in parquet format + Transaction log
- Provides ACID guarantees

**Operations**

- Write operation: First write the files, then the log
- Read operation: First read the log, then the files

**Create Delta Table by storing data in delta format**

**Various options to add data to Delta Table**

- Append DataFrame (Python) & Insert (SQL)

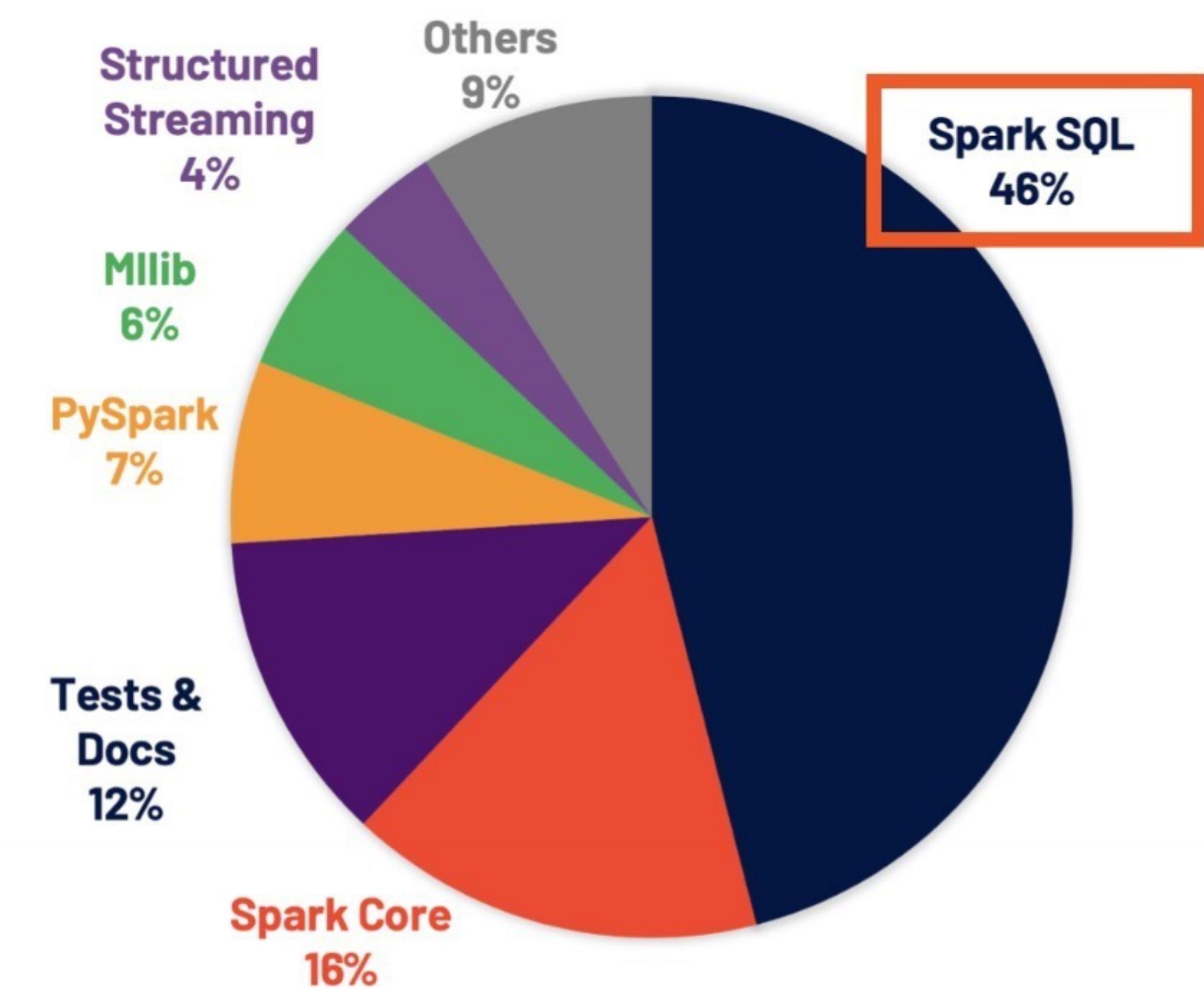
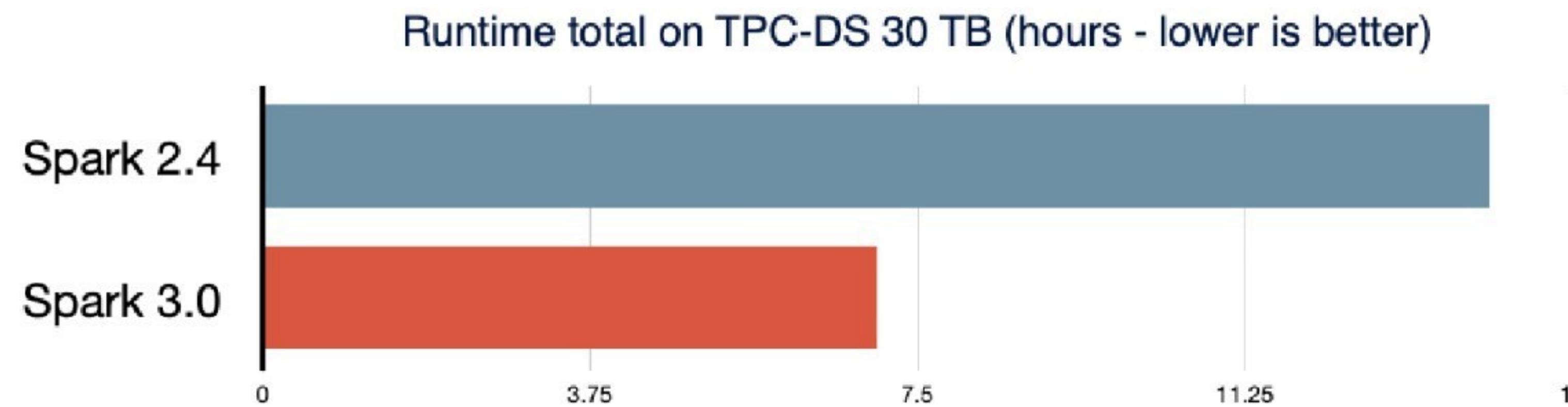
**Perform DML operations - Update, Delete & Merge**

**Apply table constraints – Not Null & Check**

**Access data using time travel**

- Query using version or timestamp, or restore table

# Apache Spark 3

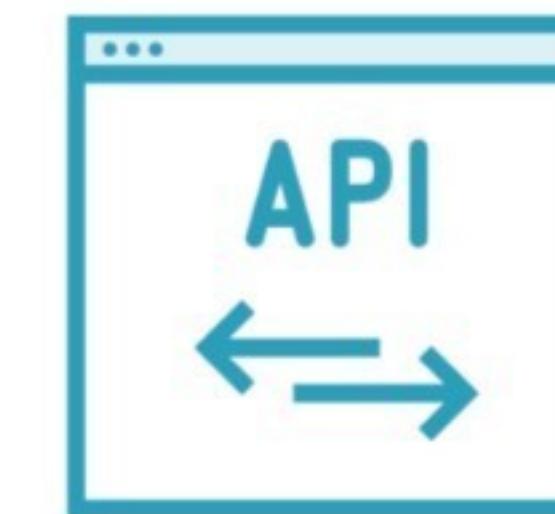


\*Images Source: [databricks.com/blog](https://databricks.com/blog)

# Apache Spark 3 Features



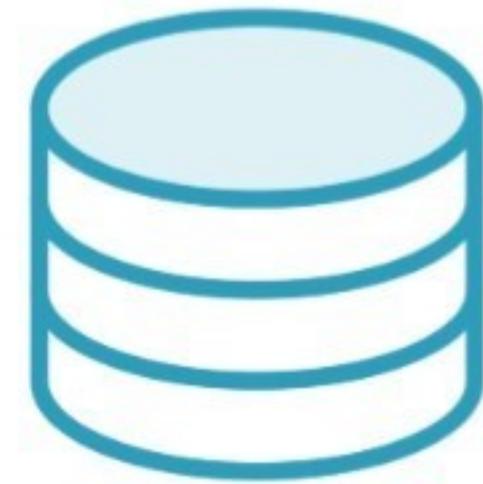
**Performance**



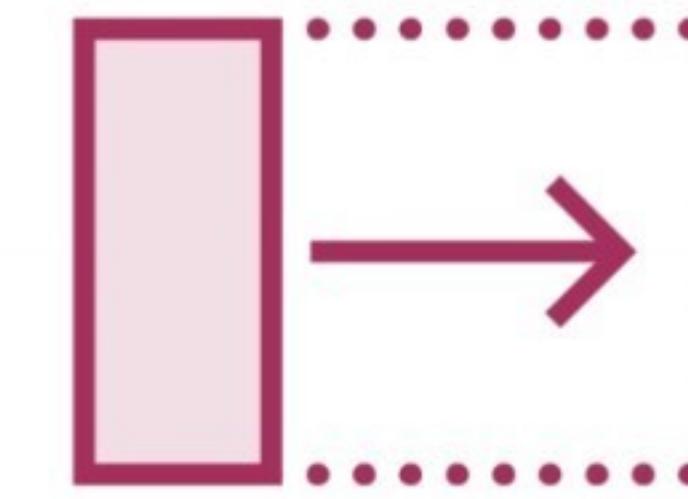
**API Enhancements**



**SQL Compliance**



**Data Sources**



**Extensibility**



**Monitoring**

# Performance Improvements

1

## Adaptive Query Execution (AQE) framework

- Reoptimizes query plan at runtime based on stats
- a) Dynamically coalescing shuffle partitions
- b) Dynamically switching join strategy
- c) Handling data skew in joins

## Dynamic Partition Pruning (DPP)

- Improves on Partition Pruning technique

## SQL join hints

- Spark has multiple join strategies
- Join Hint allows to enforce a particular join strategy
- Join Hints available for each join strategy

## Faster query compilation

# APIs, SQL and Monitoring

2

## 30+ new built-in functions

### PySpark enhancements

- Pandas API enhancement & support
- Better error handling

### Deep Learning improvements

- Project Hydrogen aims to unify data processing & deep learning

### ANSI SQL compliance

- `spark.sql.ansi.enabled` (*disabled by default*)

### Monitoring

- Better UI for Spark Structured Streaming
- Observable Metrics

# Data Sources, Extensibility & Ecosystem

3

## Built-in Data Sources

- New data sources like Apache Iceberg
- Performance improvements to existing sources like Parquet, Kafka, Delta Lake etc.

## Extensibility

- Catalog API to use external catalog for managing tables (instead of Hive)

## Ecosystem

- Support for Hadoop 3+, Hive 3+, JDK 11+

# Overview



## Adaptive Query Execution

- Dynamically coalescing shuffle partitions
- Dynamically switching join strategy
- Handling data skew in joins

## Dynamic Partition Pruning

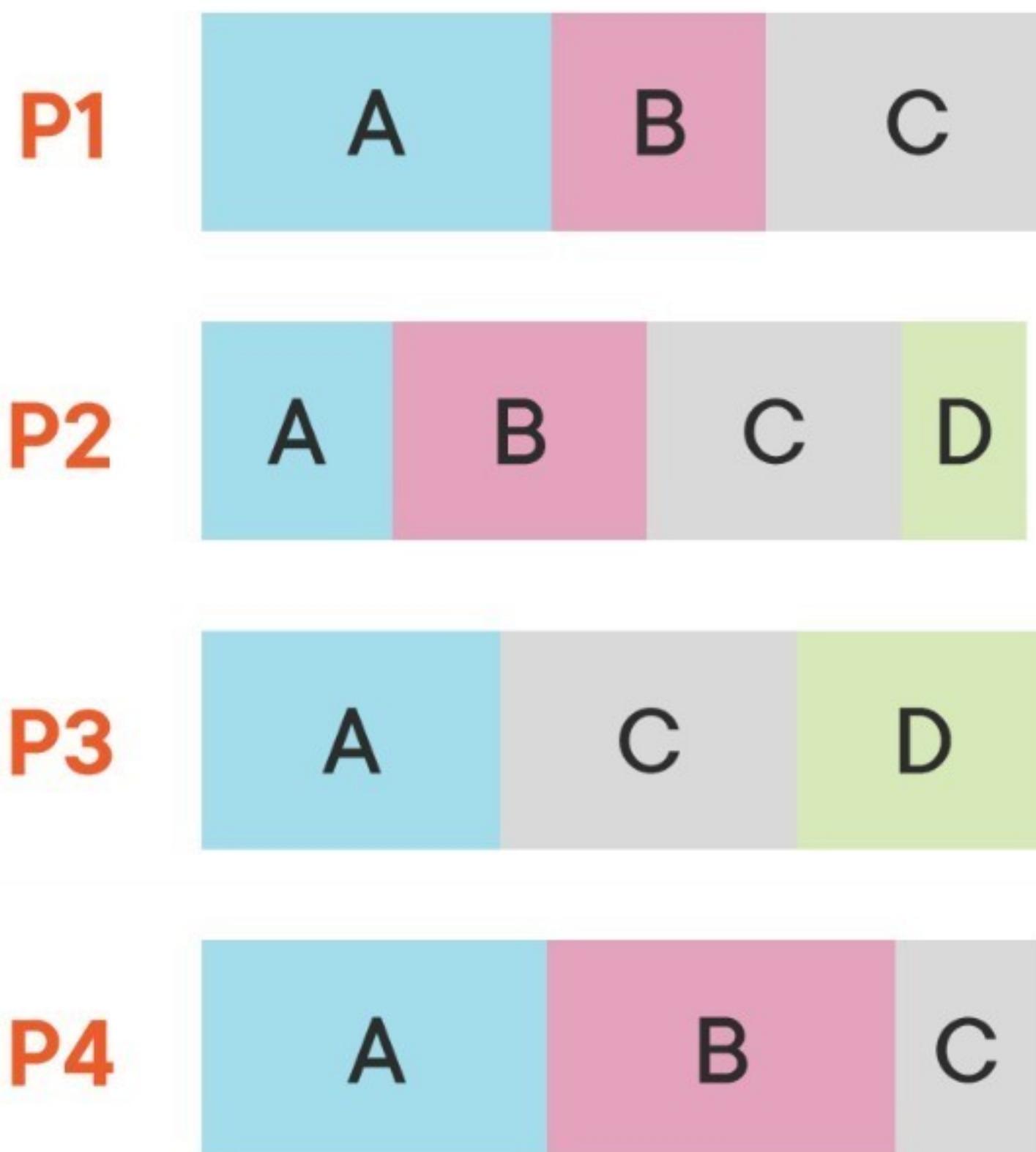
# Adaptive Query Execution: Dynamic Coalescing

---

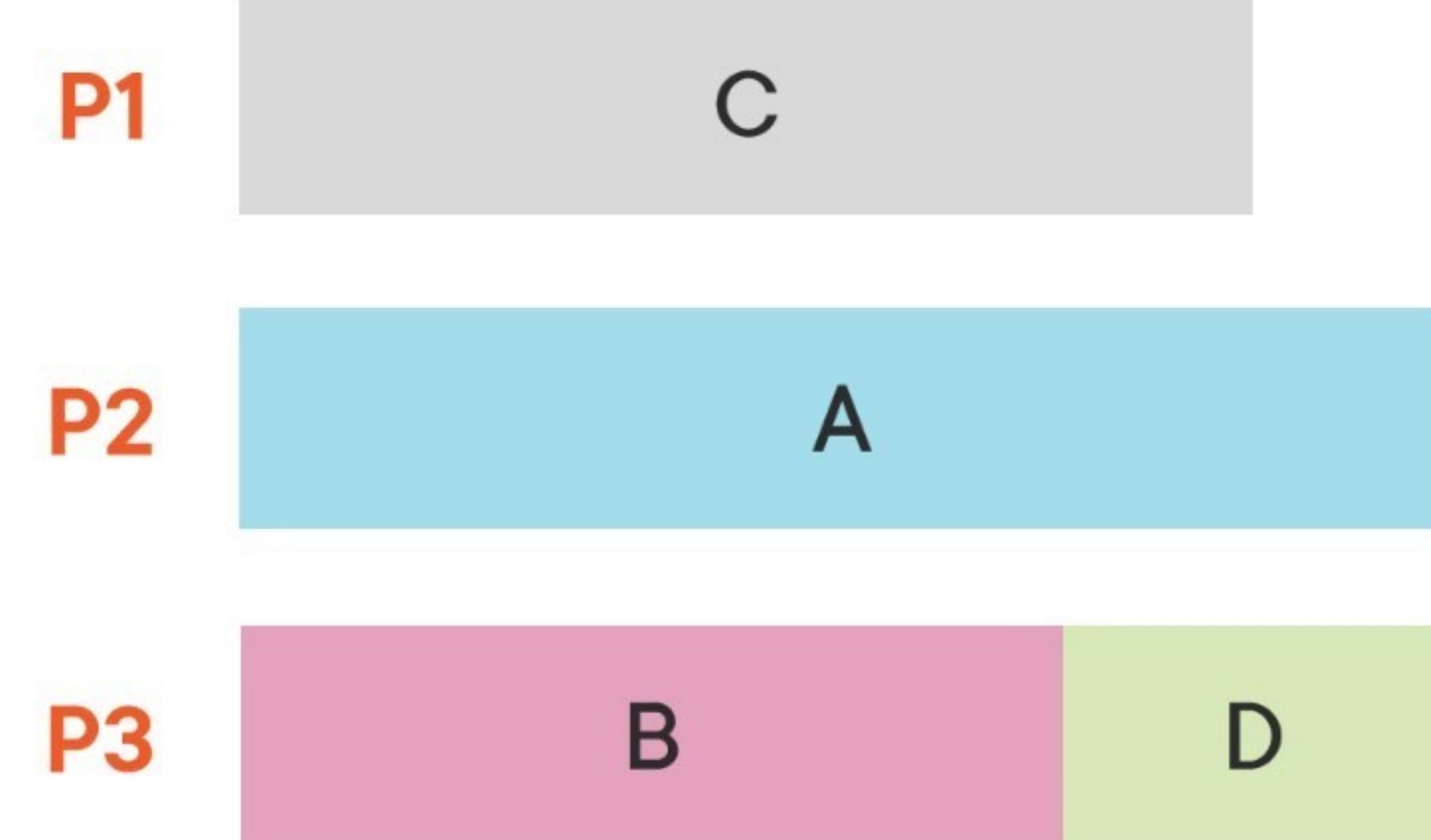
# Shuffling Data

`spark.sql.shuffle.partitions = 3`

[default = 200]



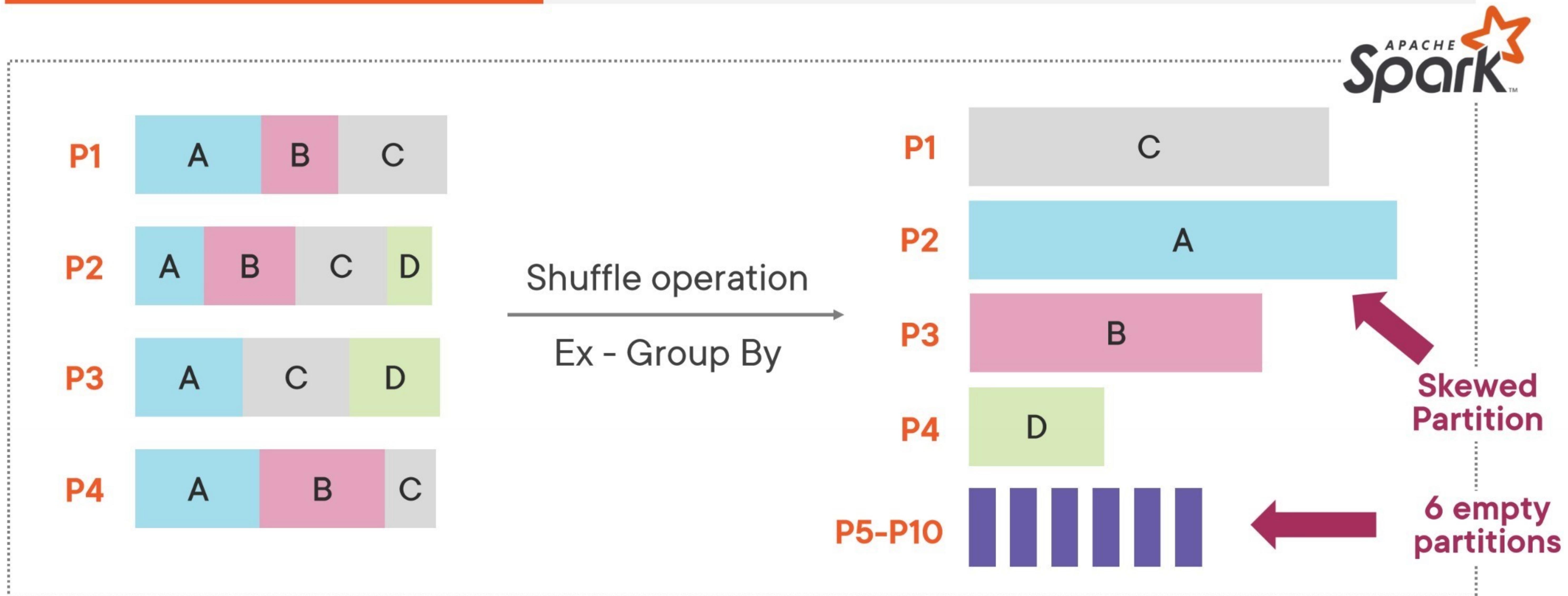
Shuffle operation  
Ex - Group By



# Shuffling Data

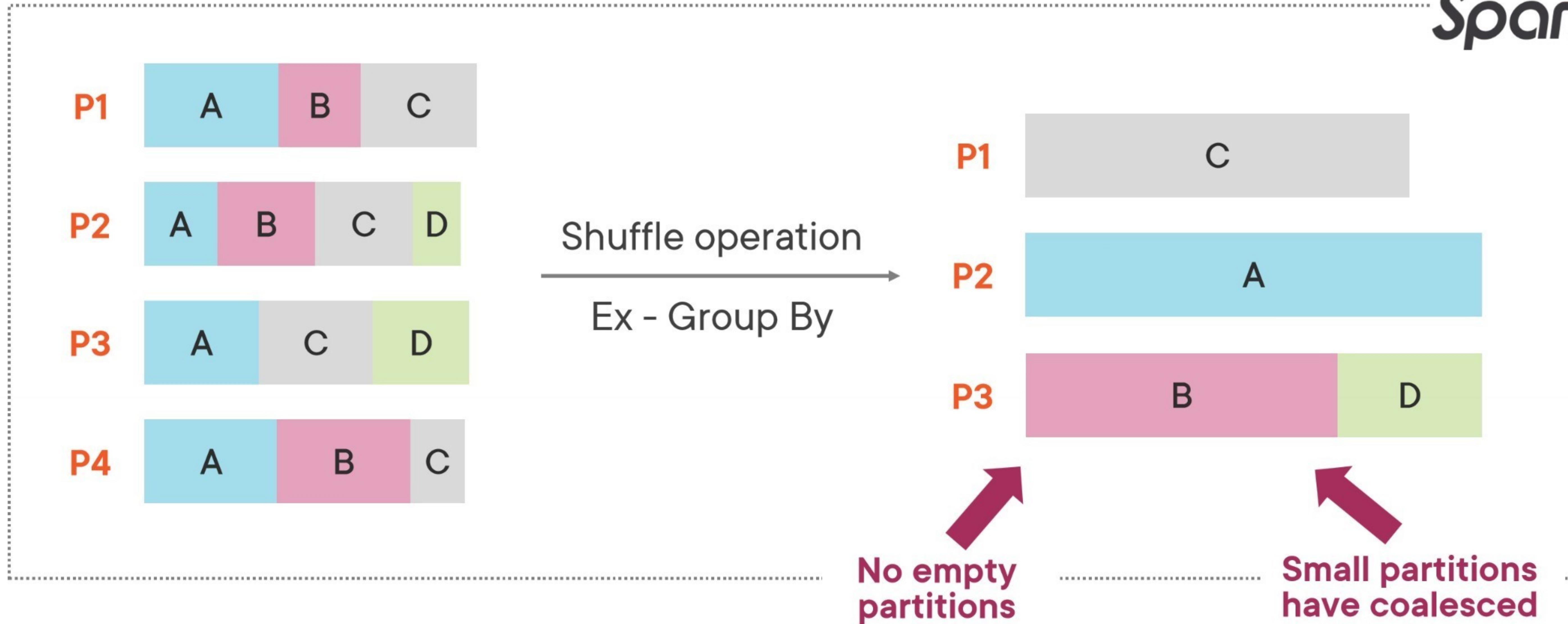
spark.sql.shuffle.partitions = 10

[default = 200]

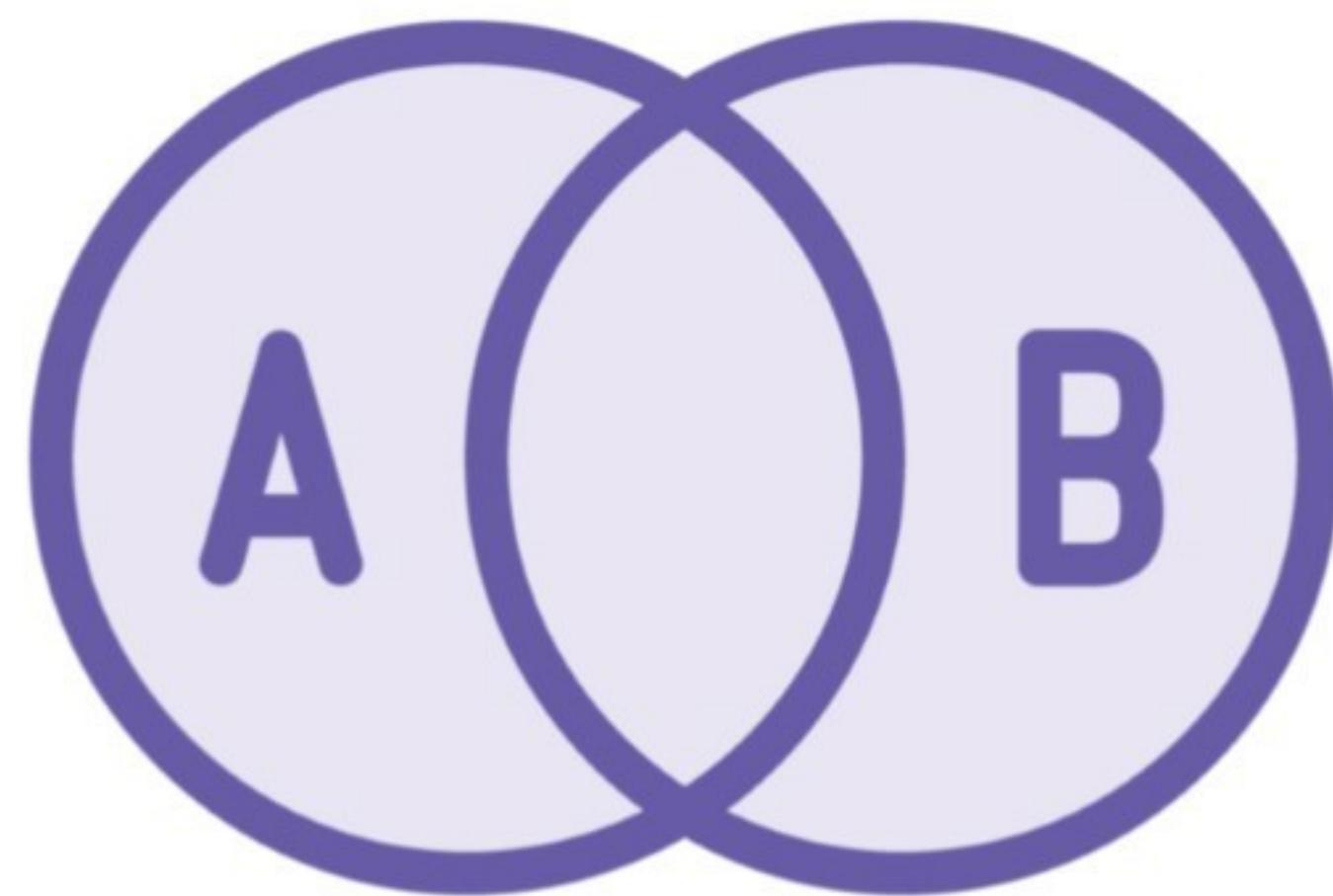


# Shuffling With AQE

spark.sql.shuffle.partitions	= 10	[default = 200]
spark.sql.adaptive.enabled	= true	[default = true]
spark.sql.adaptive.coalescePartitions.enabled	= true	[default = true]



# Dynamic Coalescing Shuffle Partitions



**Having empty or lot of small partitions:**

- Too many tasks are created
- Reduces parallelism and consumes time/resources

**AQE dynamically coalesces shuffle partitions**

- Removes empty partitions
- Combines small partitions to produce optimal sized shuffle partitions

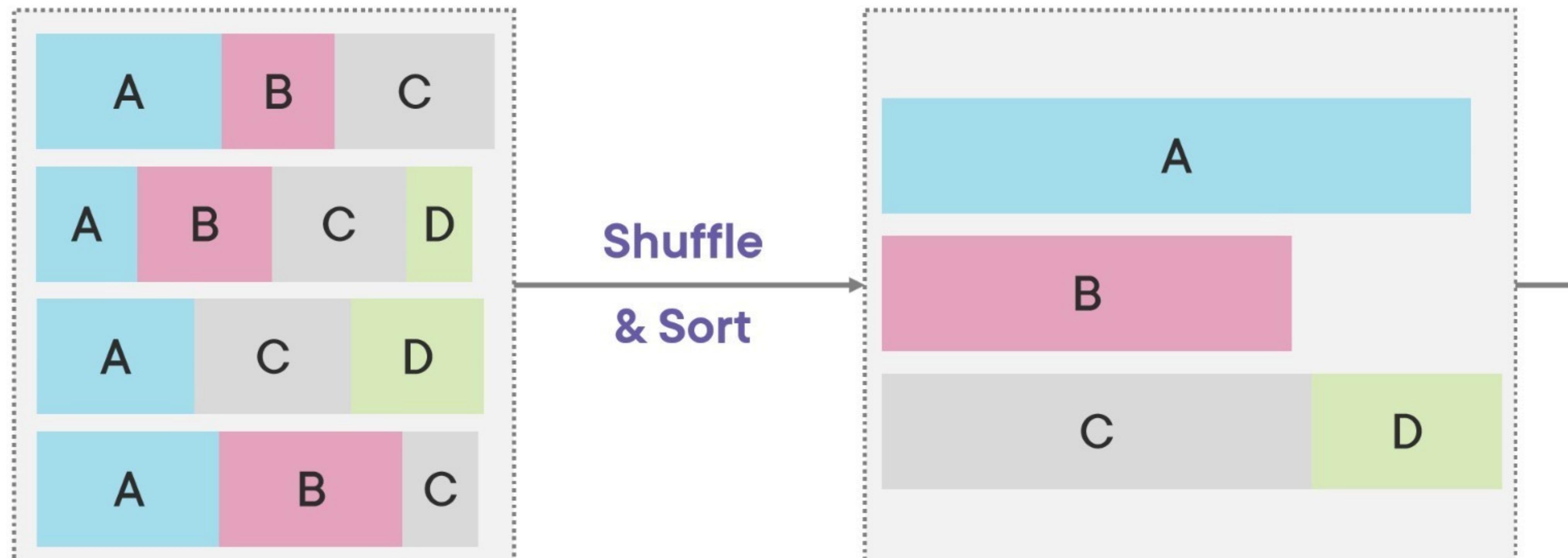
Datasets are big → Sort merge

One dataset  
is small and → Broadcast  
one is big

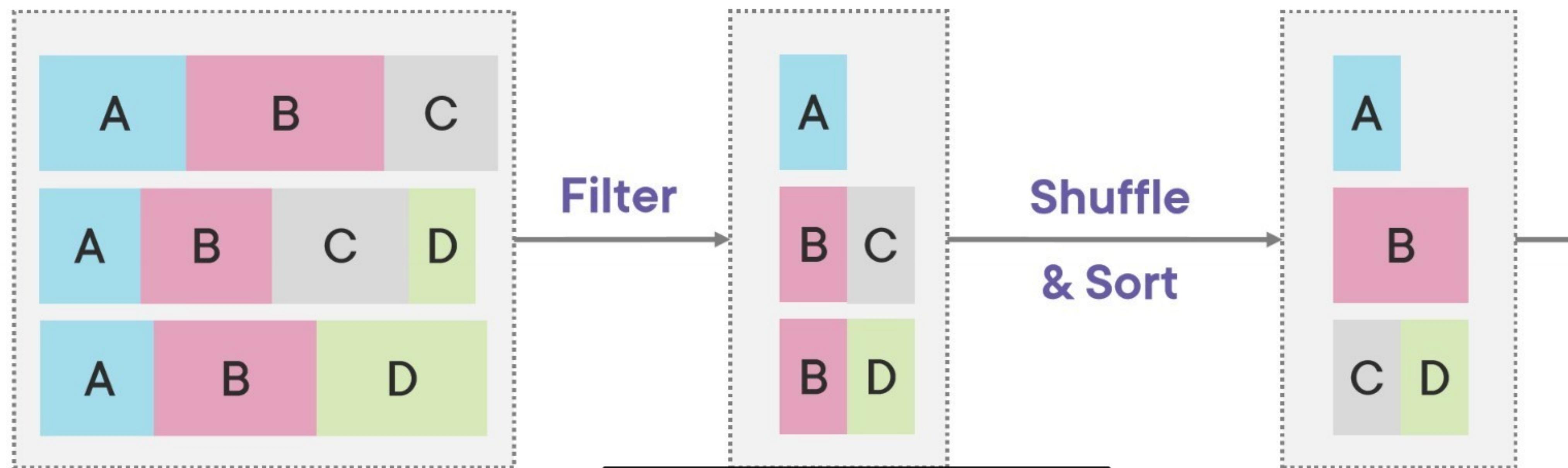
## Adaptive Query Execution: Dynamic Join

---

## Without AQE

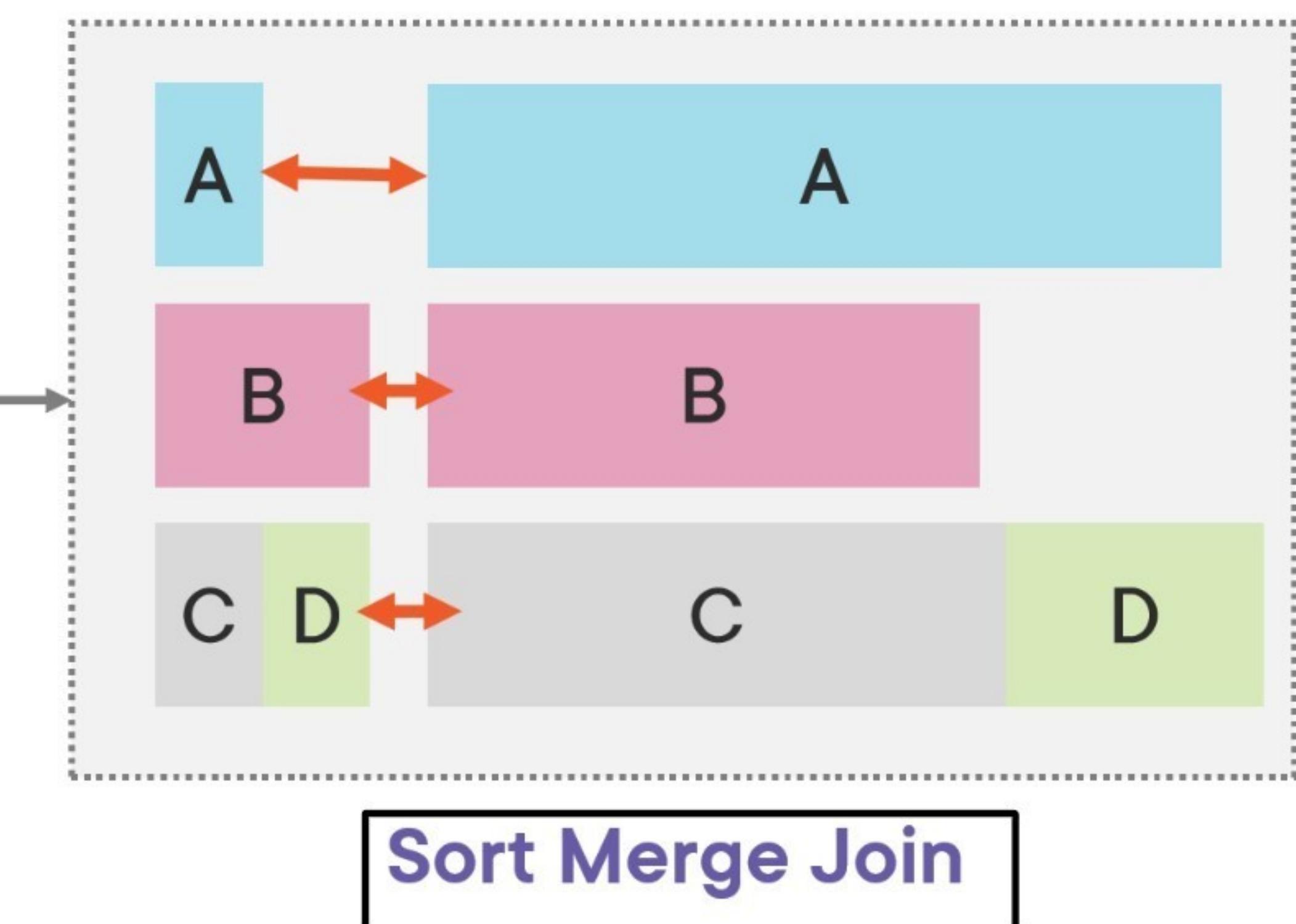


Sales DF  
100 mn records

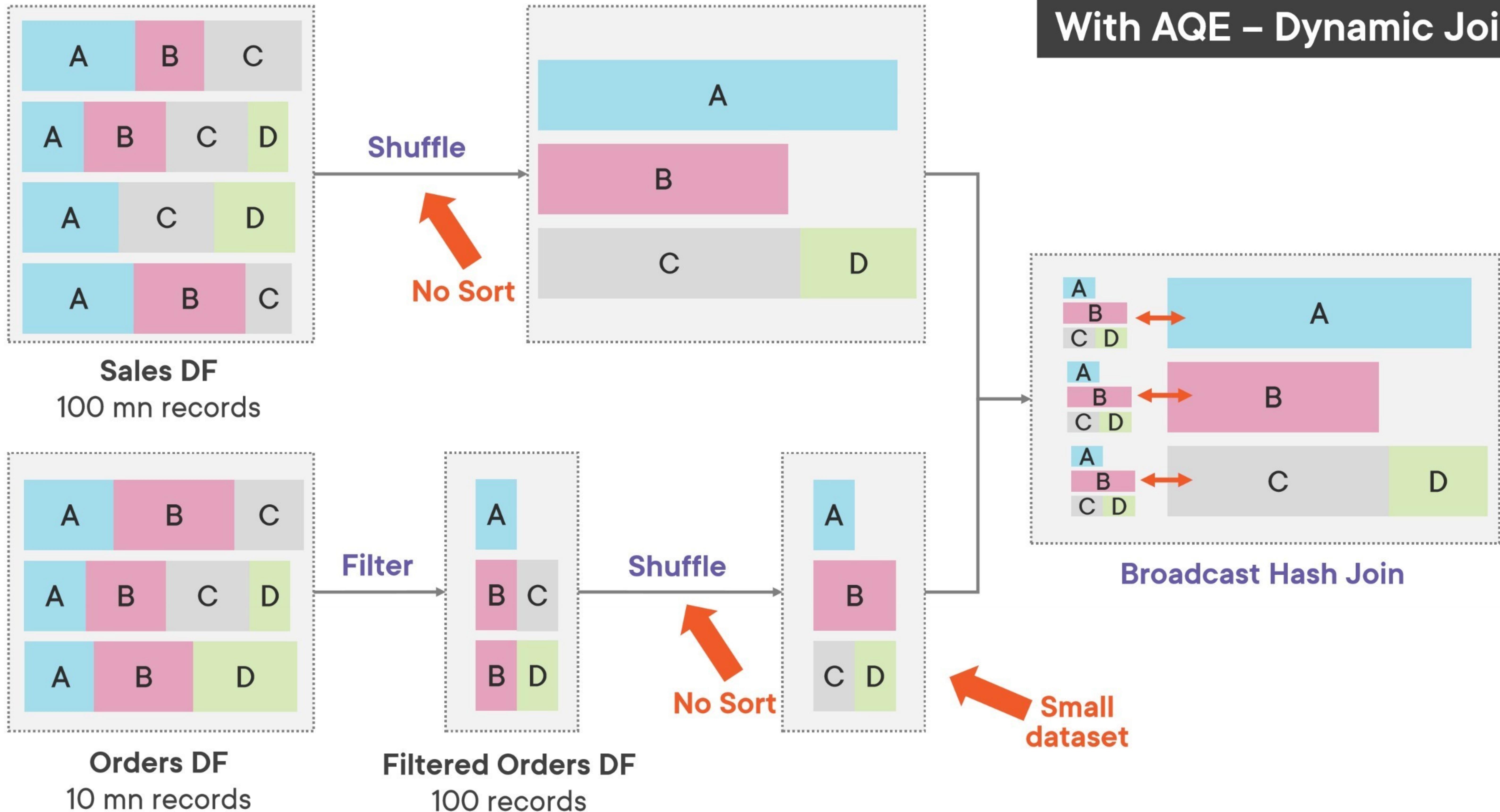


Orders DF  
10 mn records

Filtered Orders DF  
100 records



## With AQE – Dynamic Join



# AQE – Dynamic Switching Join Strategy OR Manual Broadcast

Which is better?

# Dynamically Switching Join Strategies

**For large datasets, Shuffle Sort Merge Join is performed**

**For Broadcast Hash Join, one dataset must be small**

- Small dataset should be less than setting `spark.sql.autoBroadcastJoinThreshold`
- Default threshold is 10 MB

**If highly selective filter is applied on large dataset**

- It may become smaller than broadcast threshold
- Since execution plan is ready, still performs Shuffle Sort Merge Join (if joined with large dataset)

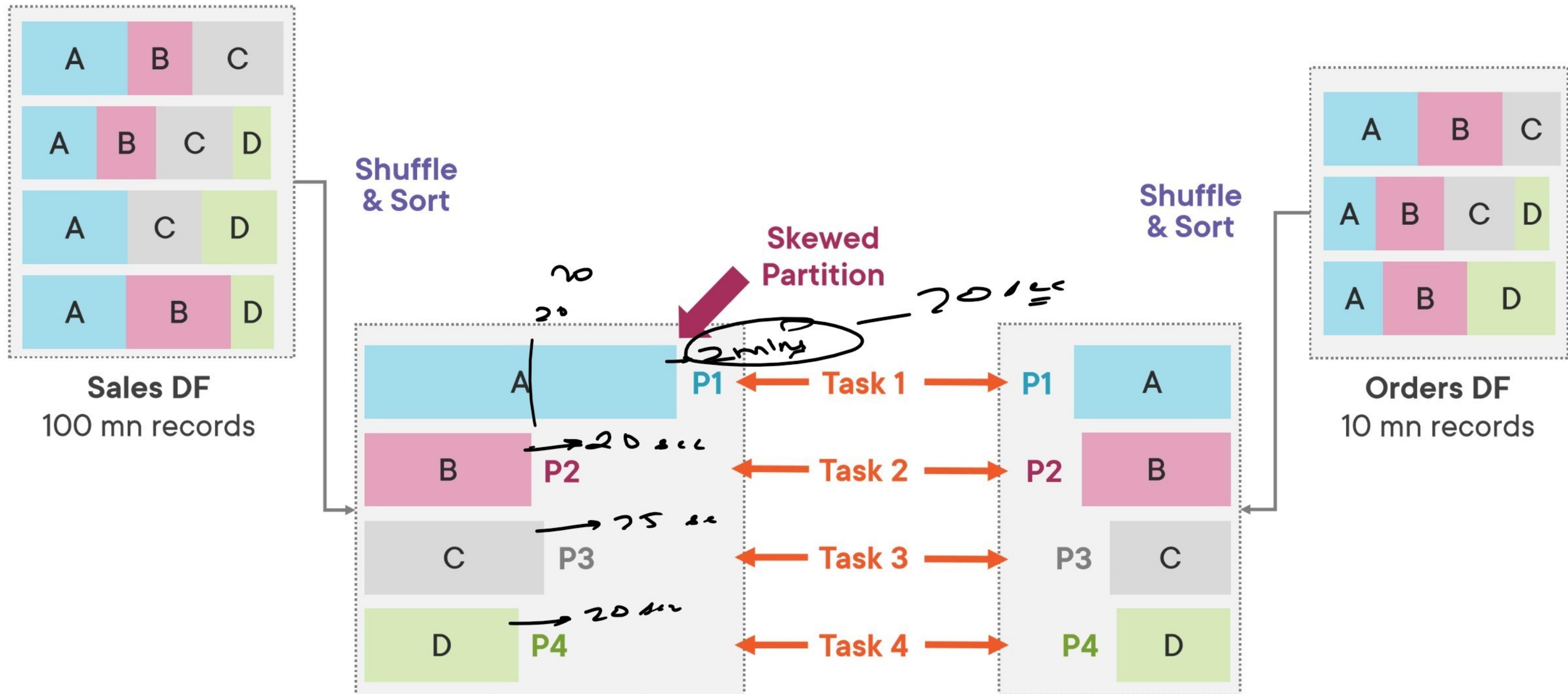
**If enabled, AQE dynamically switches from Sort Merge Join to Broadcast Hash Join at runtime**

- AQE checks for dataset size after shuffle
- If size of a dataset is now less than broadcast threshold, switches to Broadcast Hash Join

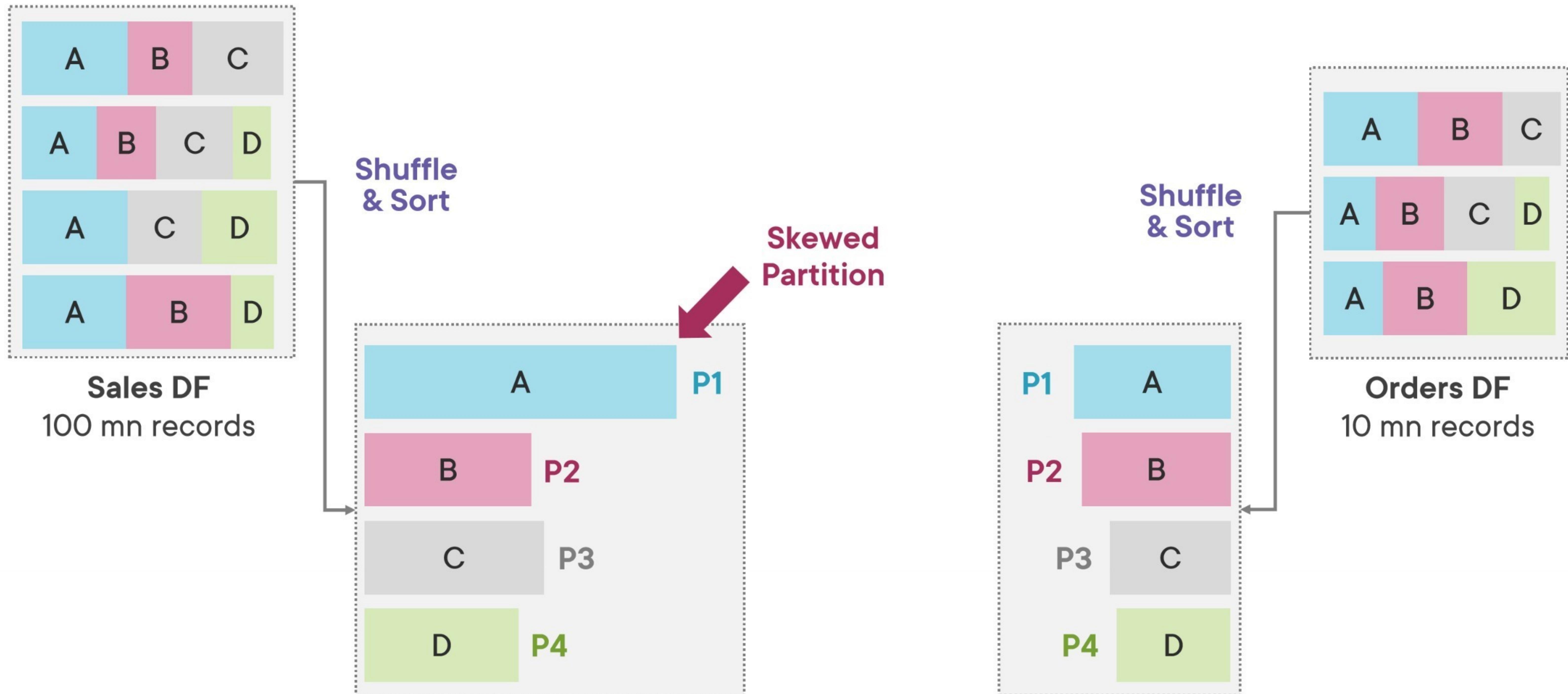
# Adaptive Query Execution: Handling Skew

---

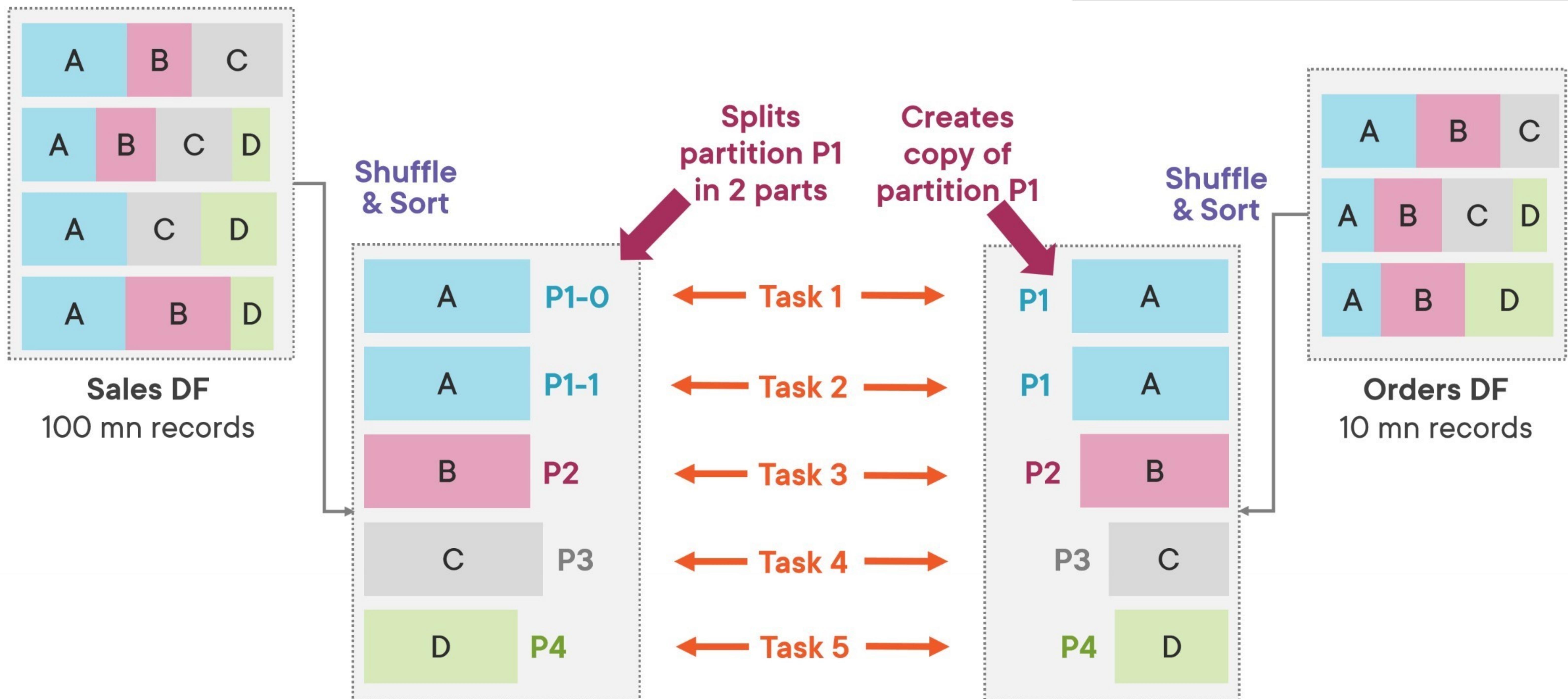
# Without AQE



# With AQE: Handling Skew



## With AQE: Handling Skew



# Dynamically Optimizing Skew Joins

## Data Skew

- One partition has much more data than others

## In join operations:

- After shuffle, data may be unevenly distributed among partitions
- Data skew can impact query performance
- Tasks processing larger partitions will take more time than ones handling smaller partitions

## If enabled, AQE dynamically optimizes data skew in joins

- AQE checks for partition sizes after shuffle
- Splits skewed partitions into smaller sub-partitions
- Creates copy of corresponding partition on other side
- Number of tasks increase, but each one will almost take same time to finish

# Dynamic Partition Pruning

---

SALES.CSV	
	ProductId=1
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=3
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=4
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=5
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=6
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=7
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv

## Disk Partitions



```
(  
  salesDF  
    .write  
    .partitionBy("ProductId")  
    .saveAsTable("Sales")  
)
```

SALES.CSV	
	ProductId=1
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=3
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=4
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=5
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=6
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	ProductId=7
	part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
	part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv

## Partition Pruning

```
SELECT *
FROM Sales
WHERE ProductId = 3
```

## Sales

100 mn records - Partitioned

SALES.CSV	
▼ ProductId=1	
part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
▼ ProductId=3	
part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
▼ ProductId=4	
part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
▼ ProductId=5	
part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
▼ ProductId=6	
part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	
part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv	

PRODUCTS.CSV	
	part-00000-34eb7004-74e7-40e0-885c-96f3820ad4df-c000.csv

## Products

1000 records – Non-partitioned

Partition Pruning  
Works!

SELECT /\*+ BROADCASTJOIN(p) \*/ \*

FROM Sales s  
JOIN Products p  
ON s.ProductId = p.ProductId

WHERE s.ProductId = 3

Small Table

## Sales

100 mn records - Partitioned

### SALES.CSV

- ▼ ProductId=1
  - part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
  - part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
- ▼ ProductId=3
  - part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
  - part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
- ▼ ProductId=4
  - part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
  - part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
- ▼ ProductId=5
  - part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
  - part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
- ▼ ProductId=6
  - part-00000-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv
  - part-00001-21020a36-b4ab-4c7f-9f17-a9c79a3f47d8.c000.csv

### PRODUCTS.CSV

- part-00000-34eb7004-74e7-40e0-885c-96f3820ad4df-c000.csv

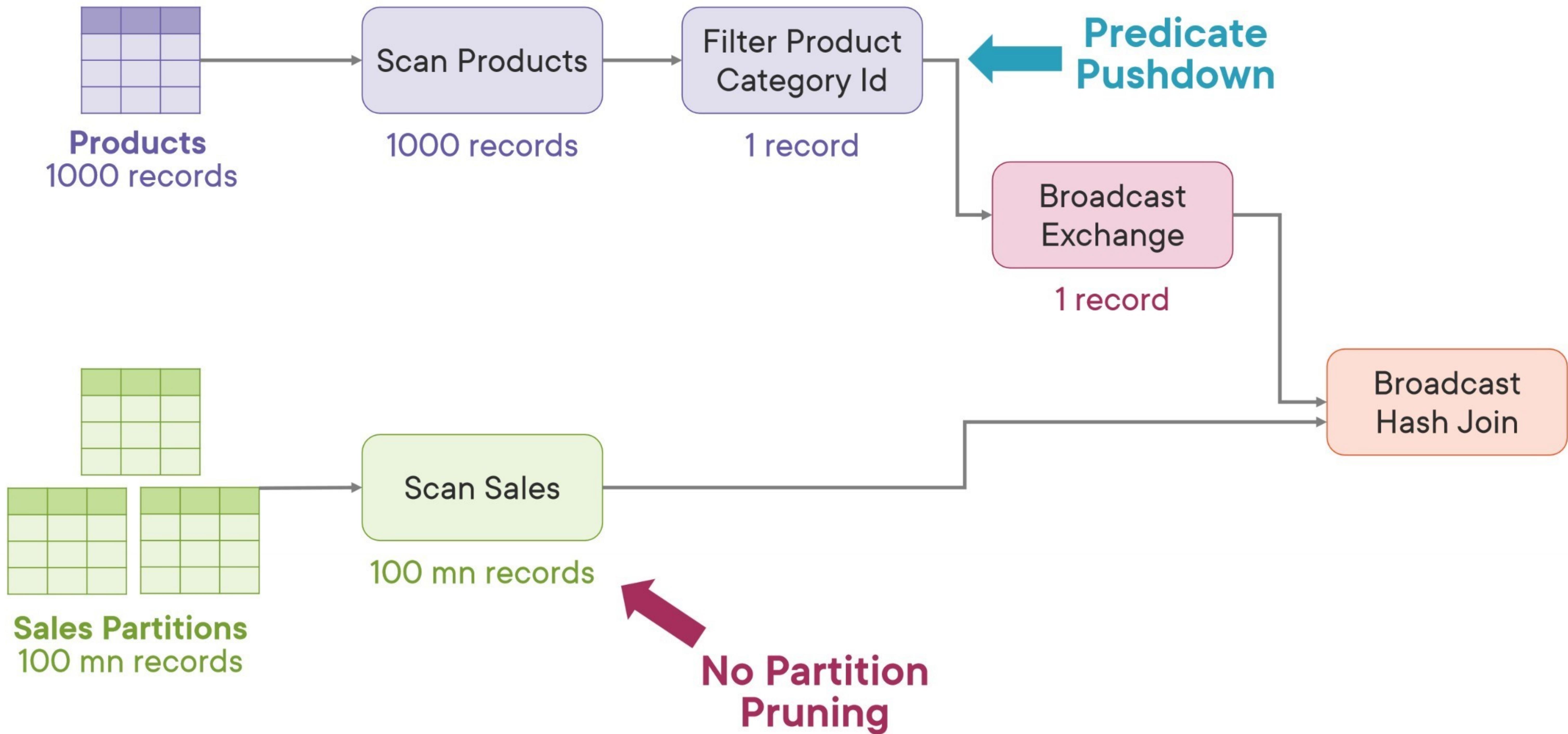
## Products

1000 records – Non-partitioned

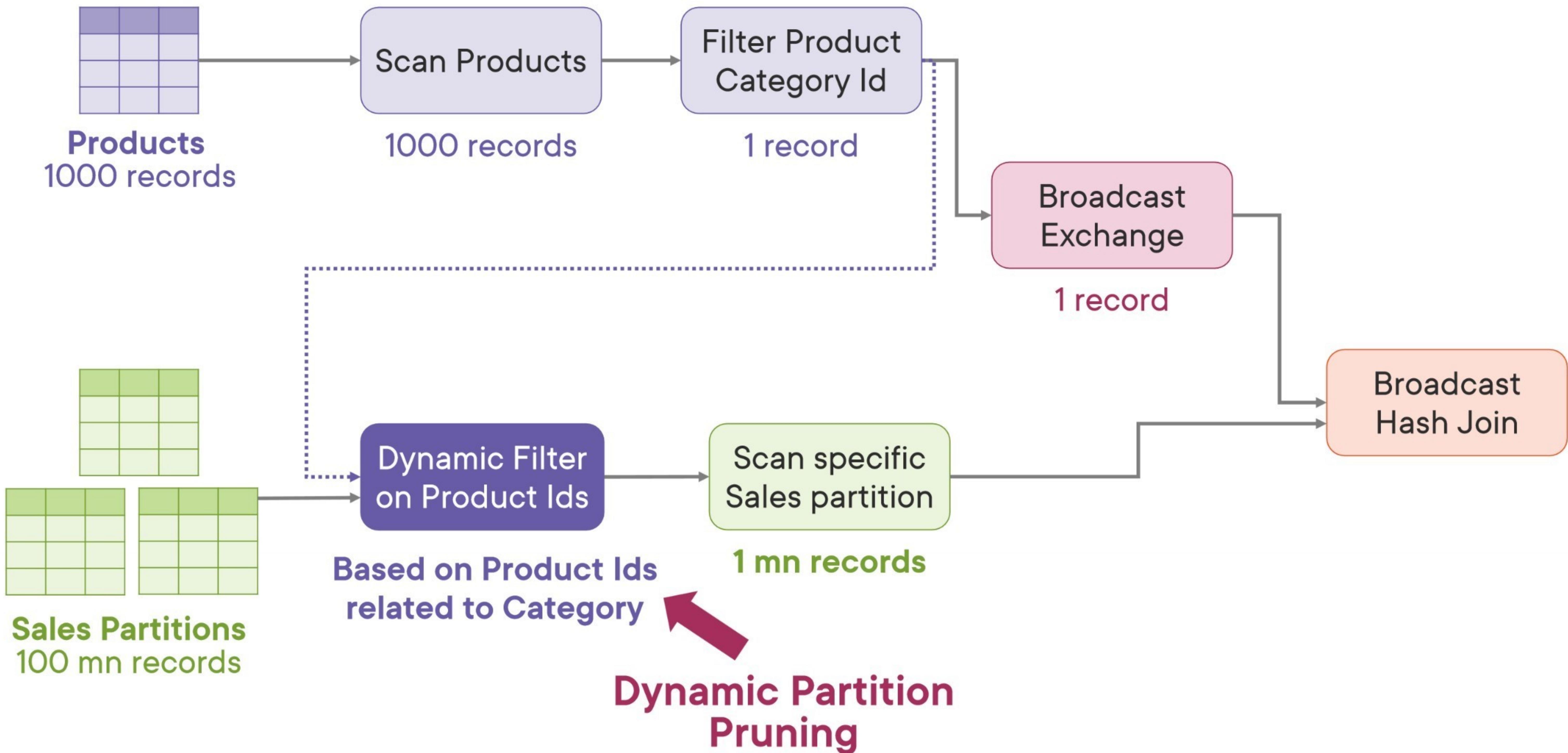
## Partition Pruning Does Not Work!

```
SELECT /*+ BROADCASTJOIN(p) */ *
  FROM Sales s
  JOIN Products p
    ON s.ProductId = p.ProductId
   WHERE p.ProductCategoryId = 3
```

## Without Dynamic Partition Pruning



## With Dynamic Partition Pruning



## Conditions

**Dynamic Partition Pruning must be enabled**

- `spark.sql.optimizer.dynamicPartitionPruning.enabled`

**Large table must have disk partitions**

**During join, small table should be broadcasted**

## Summary



## Spark 3 has several optimizations

### Adaptive Query Execution

- Reoptimizes query plan at runtime based on stats

#### 1. AQE: Dynamically coalescing shuffle partitions

- After shuffle, empty shuffle partitions are removed & small ones are merged

#### 2. AQE: Dynamically switching join strategy

- At runtime, if one dataset becomes small, AQE can switch from Shuffle Sort to Broadcast Hash join

#### 3. AQE: Handling data skew in joins

- After shuffle, if any partition is skewed, AQE can split that partition to multiple smaller partitions

### Dynamic Partition Pruning

- During join, Spark can dynamically skip disk partitions at runtime that are not required by the query