

# Datalake File Format: Parquet

[arnaud.nauwynck@gmail.com](mailto:arnaud.nauwynck@gmail.com)

Course Esilv 2023

This document:

[https://github.com/Arnaud-Nauwynck/presentations/tree/main/pres-bigdata/Course-2022-Spark /Datalake-File-Format-Parquet.pdf](https://github.com/Arnaud-Nauwynck/presentations/tree/main/pres-bigdata/Course-2022-Spark/Datalake-File-Format-Parquet.pdf)

# Outline

- Parquet Characteristics
  - Structured, with Schema
  - Data – Metadata (footer)
  - Columnar ... Column Pruning
  - Splittable ... spark Dataset Partitions
  - Compression
  - Encoding
  - Statistics, Bloom ... spark Predicate-Push-Down
  - Optimize write once, read many
    - ... spark dataset.repartition().sortWithinPartition().write

Parquet is a Structured Format  
Strongly Typed (Schema)

# File Formats

## Unstructured

### Text

```
text line 1 \n
text line 2 \n
```

### Csv

```
Col1;Col2;Col3\n
a1;b1;c1\n
a2;b2;c2\n
```

## Semi-Structured

### Json

```
{"a":"a1","b":"b1"}\n
{"a":"a2","b":"b2"}\n
```

### Xml

```
<elt>
  <a>a1</a>
  <b>b1</b>
```

## Structured

### Serialization Structured

#### Avro,Thrift,Protobuf

```
schema:XX,
value:0101010101
```

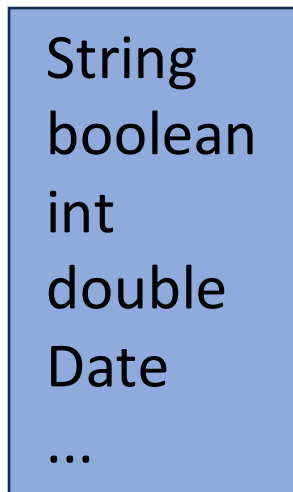
### Columnar Structured

#### Orc, Parquet

# Structured : struct<>, array<>, map<>

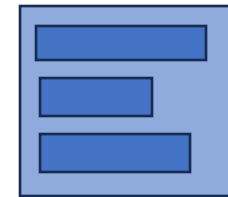
## Scalar Value

(= terminal element in grammar)  
= primitive data-type



## Composite Value

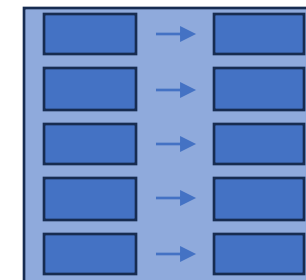
**struct<a:Type1, b:Type2, ...>**



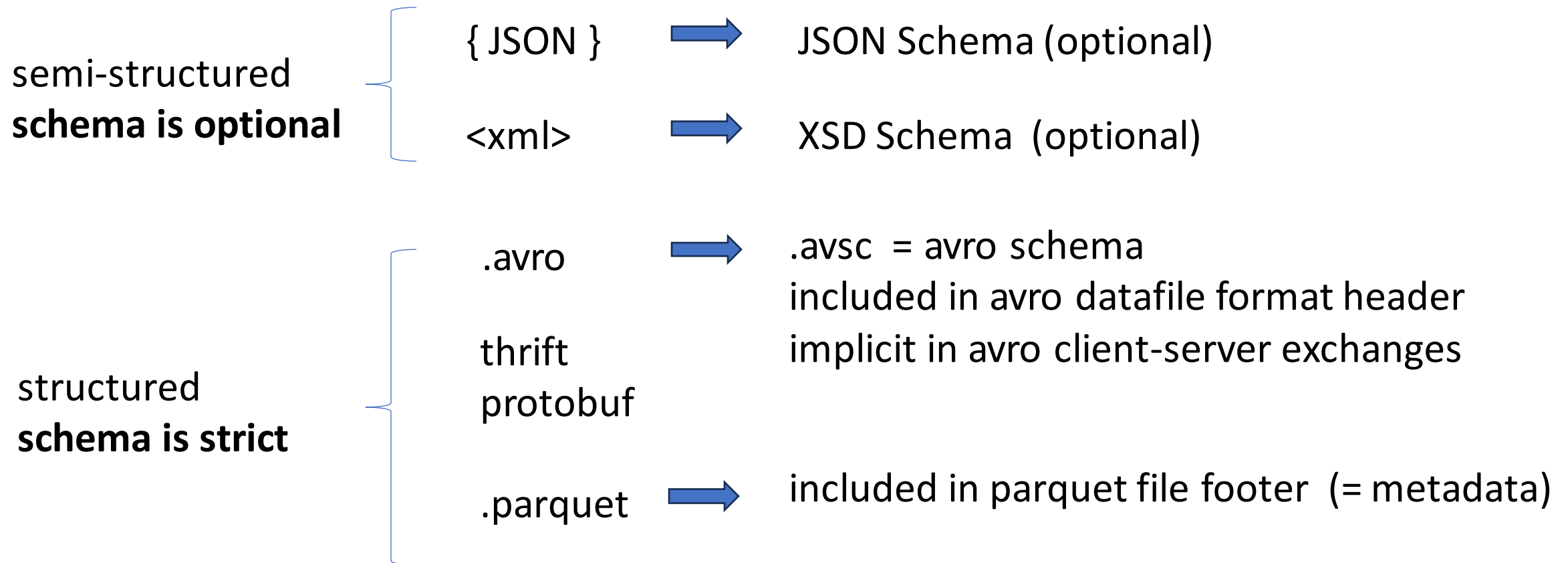
**array<ElementType>**



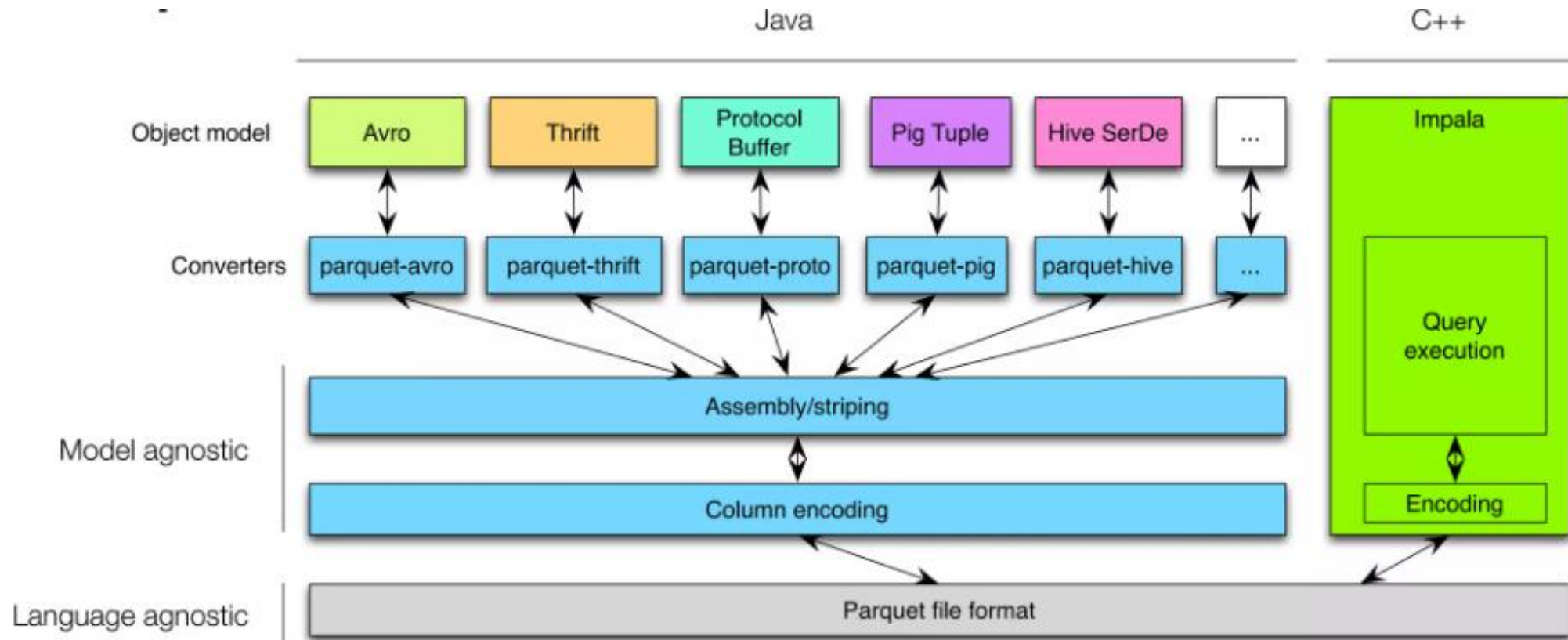
**map<KeyType,ValueType>**



# Type constraint = Schema



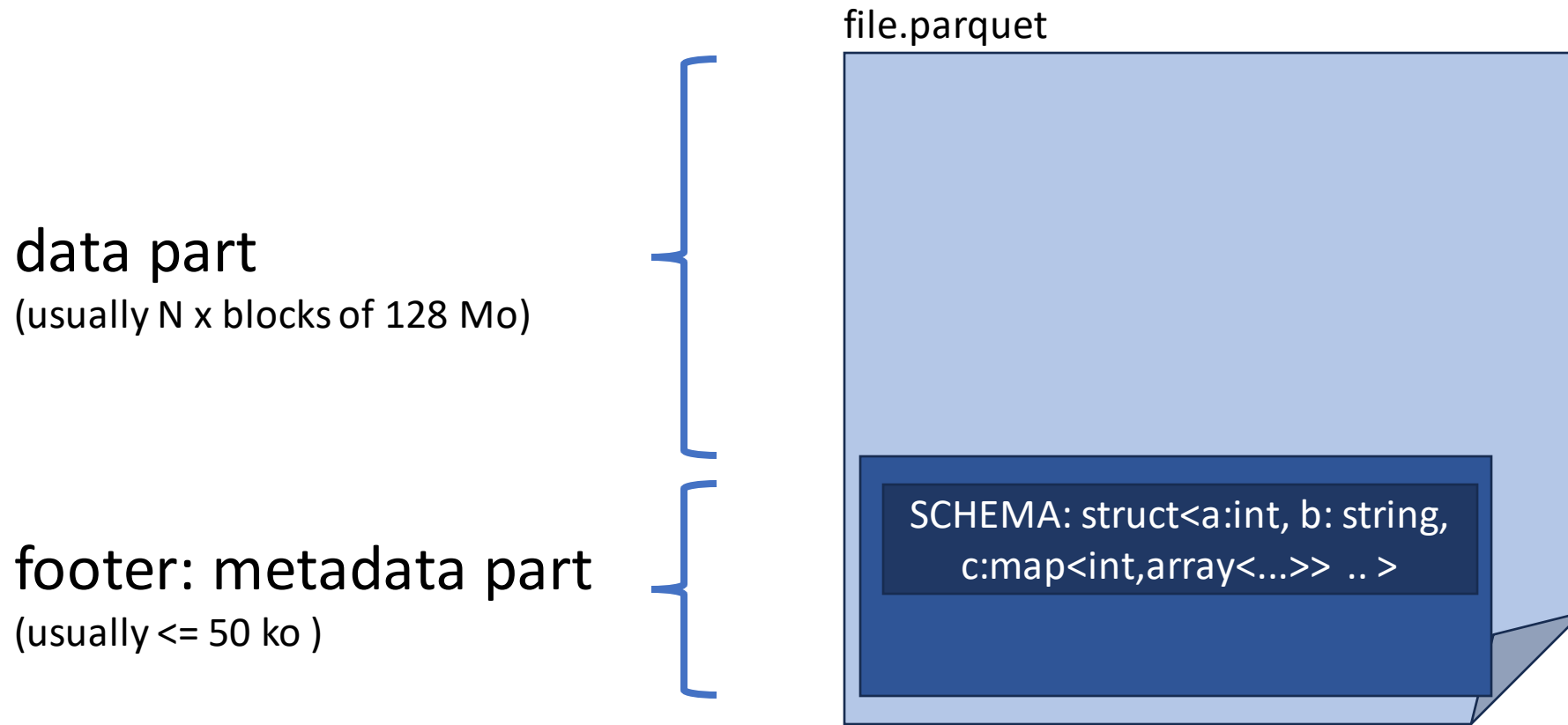
# Parquet SDK / Converter / ObjectModel



Parquet separates Data and MetaData



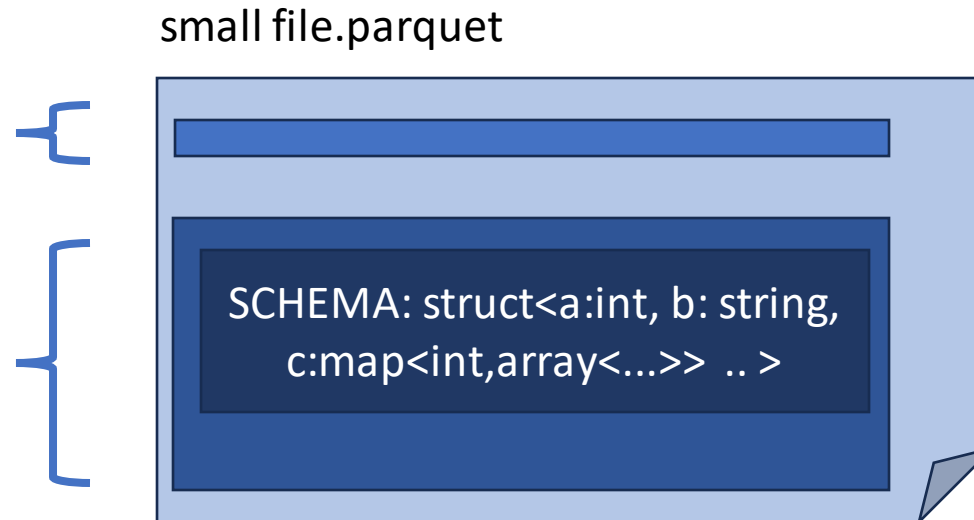
# Parquet Schema : in metadata footer



# Parquet for Small Files ?

data part can be small  
(even empty)

parquet files  
ALWAYS have an  
uncompressable footer,  
containing SCHEMA



Parquet for small data is NOT efficient  
bad ratio of "data / metadata" size

# Read Parquet footer only



1/ query File length => len

3/ read footer [len-F, len]

2/ read last bytes  [len-4, len] => F

file.parquet

SCHEMA: struct<a:int, b: string,  
c:map<int,array<...>> .. >

last bytes: <<footer size>>

# Why Footer instead of Header?

For reader : not a big overhead, ONLY 3 calls to read

For writer : MUCH more practical to "stream" write Nx rows,  
keep in-memory only few metadata,  
to flush write at end

Parquet is Columnar

# Parquet : Columnar File Format

Logical view: rows - columns

a1	b1	c1	d1
a2	b2	c2	d2
a3	b3	c3	d3
a4	b4	c4	d4

On Disk Row Serialization: like CSV, JSON, AVRO,



On Disk: **Columnar**



# Columnar

=> better memory aligned,  
Vectorized CPU pipeline

```
struct A {
```

```
    boolean f1; // <= 1 bit (on 1 byte)
                // in-memory padding 3 bytes
    int      f2; // 4 bytes, aligned on multiple of 4
                // in-memory padding 4 bytes
    long     f3 // 8 bytes
}
```

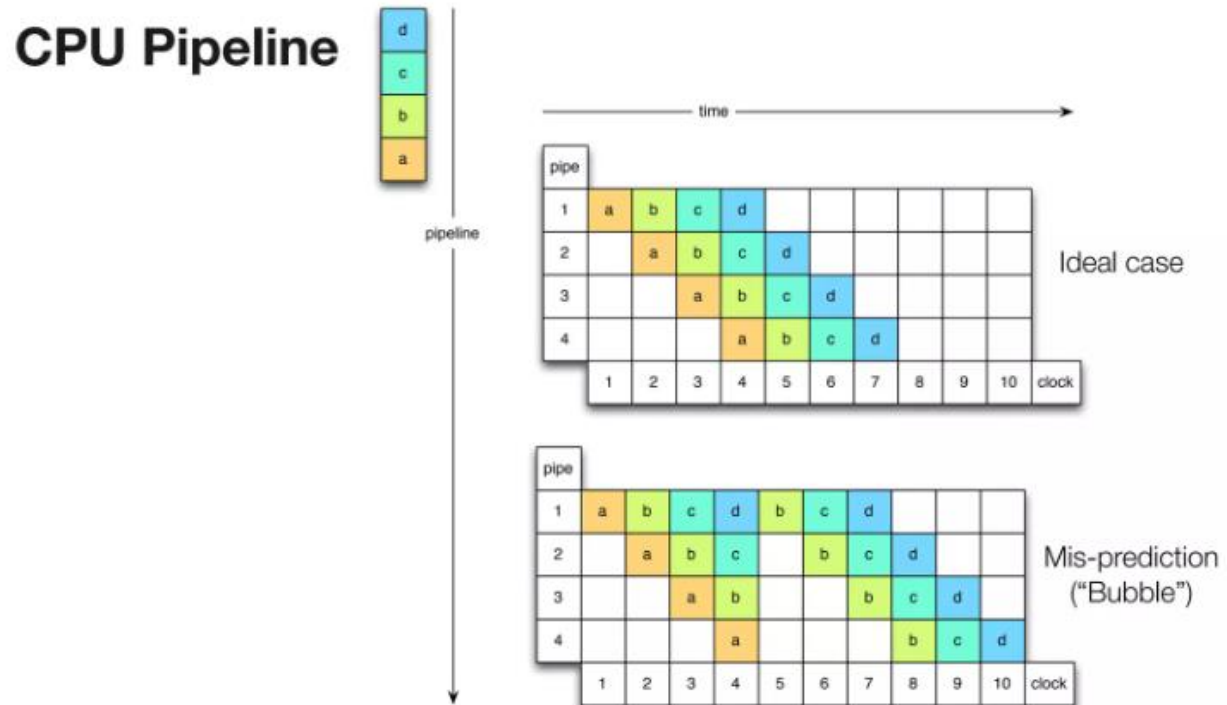
```
array[struct<..>]  <=>  array[boolean]
                    array[int]
                    array[long]
```

## Vectorized Reader ~9x Faster



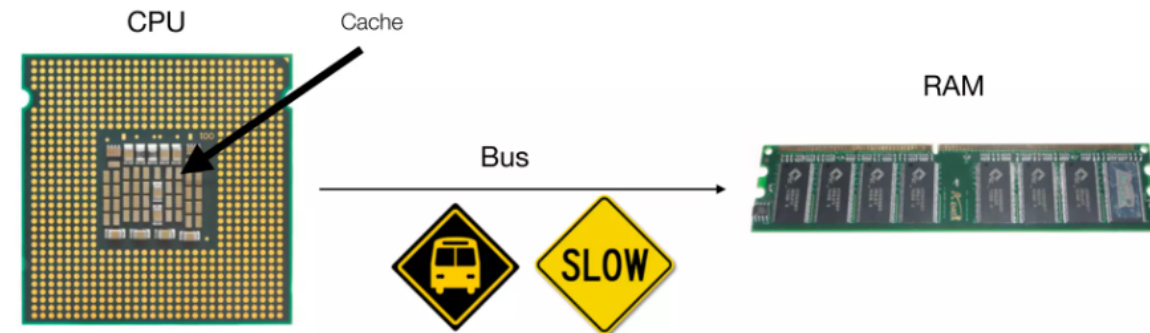
# Vectorized code ... fewer "If", "Loop", "calls"

Better CPU Pipeline



Better Memory Cache

Minimize CPU cache misses



a cache miss costs 10 to 100s cycles depending on the level



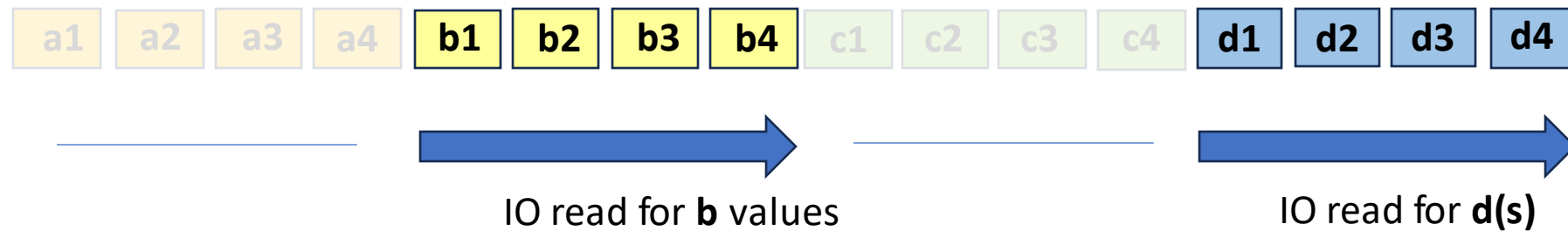
# Column Pruning Optimization

**SELECT** b, d -- ONLY 2 columns  
**FROM** table  
**WHERE** ..

Logical view: rows - columns

a1	<b>b1</b>	c1	<b>d1</b>
a2	<b>b2</b>	c2	<b>d2</b>
a3	<b>b3</b>	c3	<b>d3</b>
a4	<b>b4</b>	c4	<b>d4</b>

On Disk READ: columnar

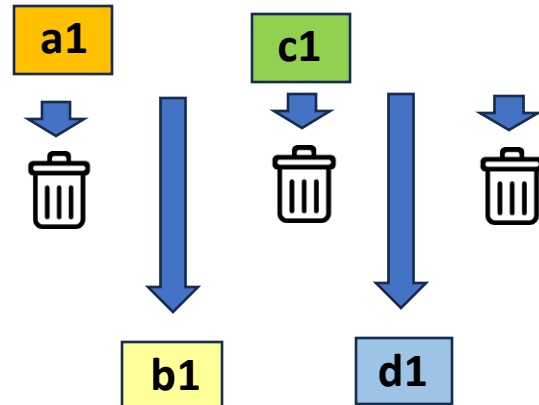


# Compare IO Reads with JSON, CSV, Avro, ..

**SELECT** b, d -- ONLY 2 columns  
**FROM** table  
**WHERE** ..



force sequential READ ALL

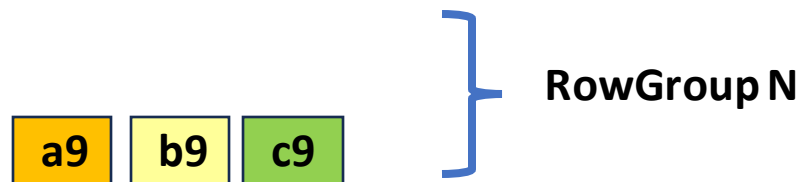
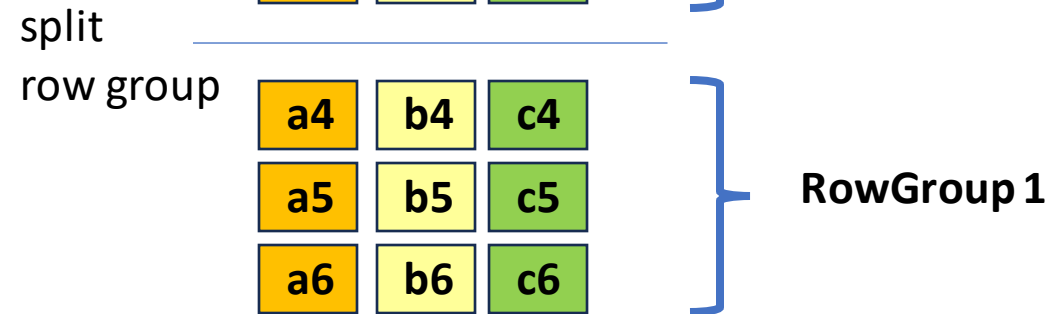
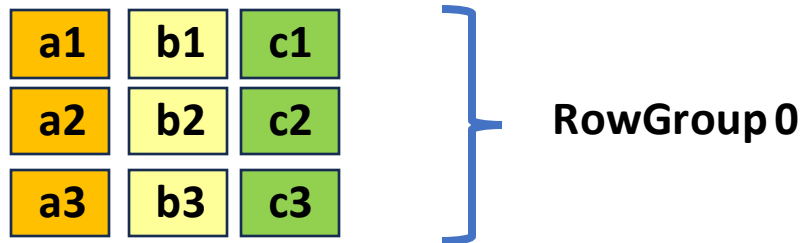


**= IO waste + CPU decoding waste**

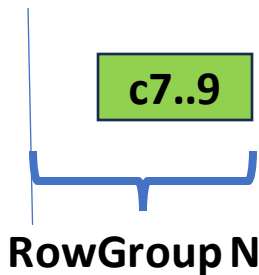
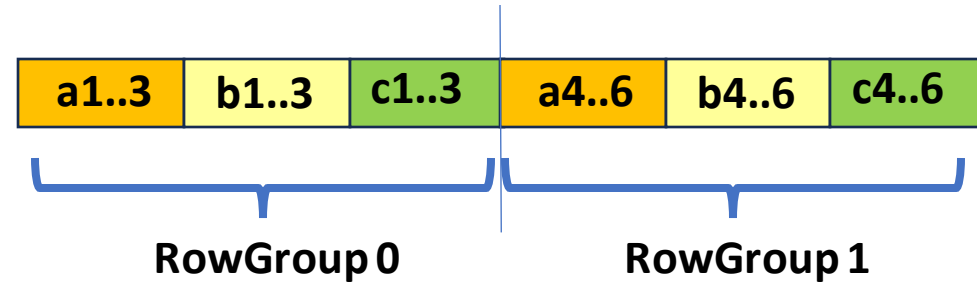
Parquet is Splitteable

# Parquet split rows by RowGroups

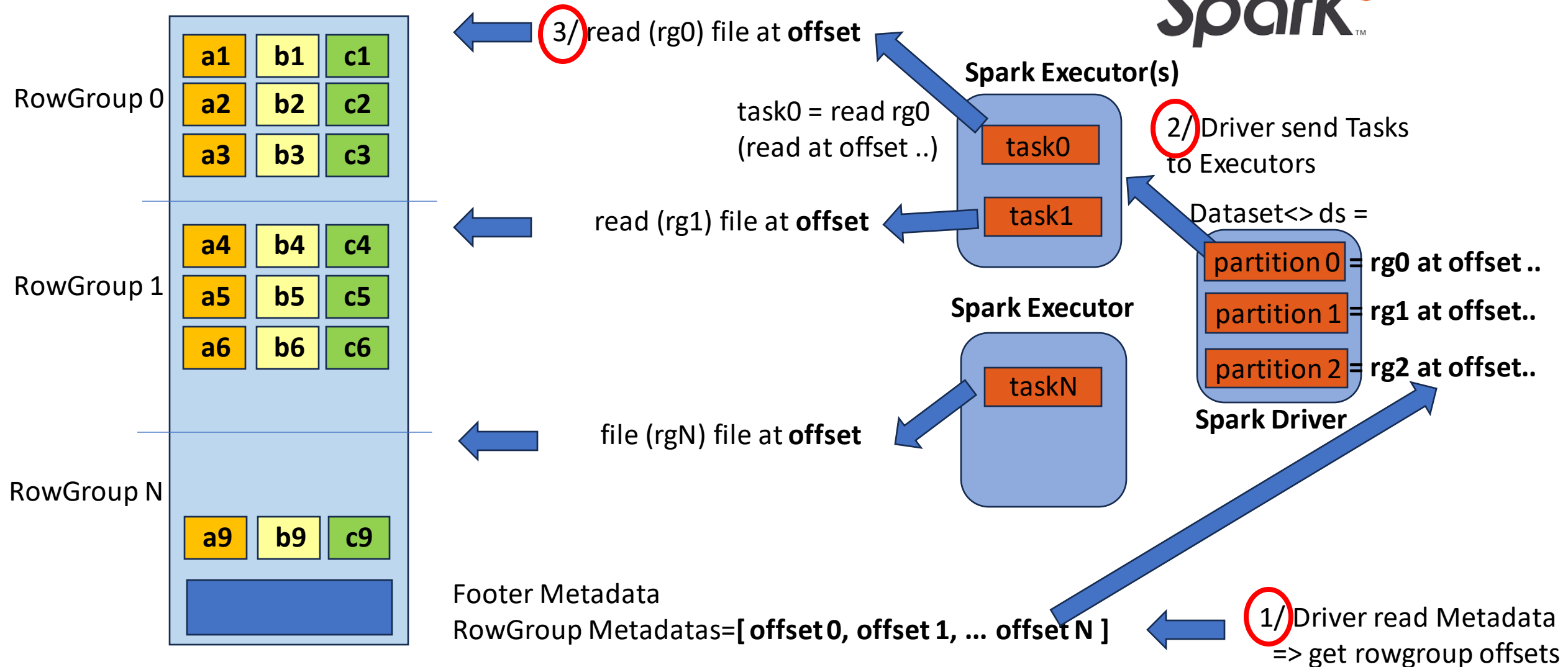
Logical view: rows - columns



On Disk: **columnar**



# RowGroup ~ Spark Partition ~ Executor Thread



# 1 Parquet RowGroup(s)

-> default to 1 Spark DataSet partition

by default,

**parquet.block.size = 128 Mo**

=

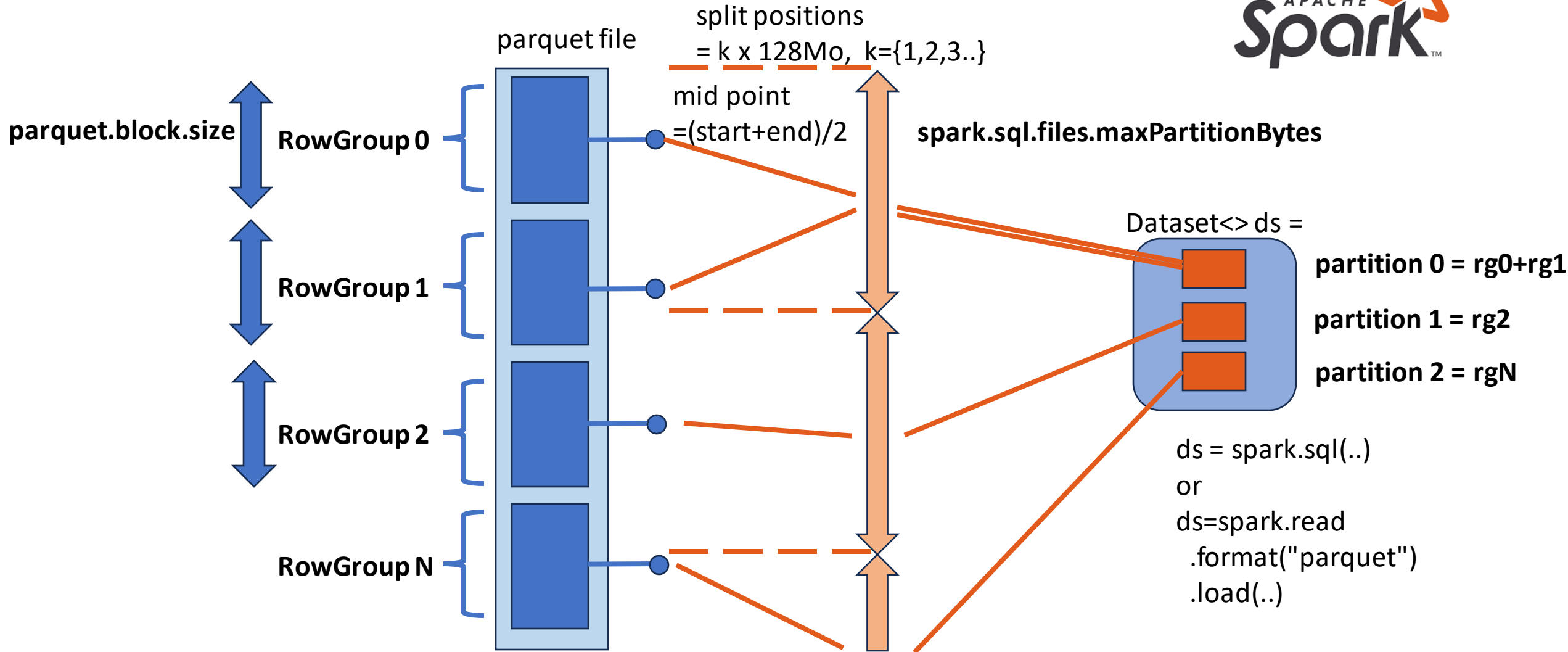
spark.files.maxPartitionBytes = 128 Mo

=

**spark.sql.files.maxPartitionBytes = 128 Mo**

# N Parquet RowGroups

-> FileSplit to P ( $\leq N$ ) Spark Partitions

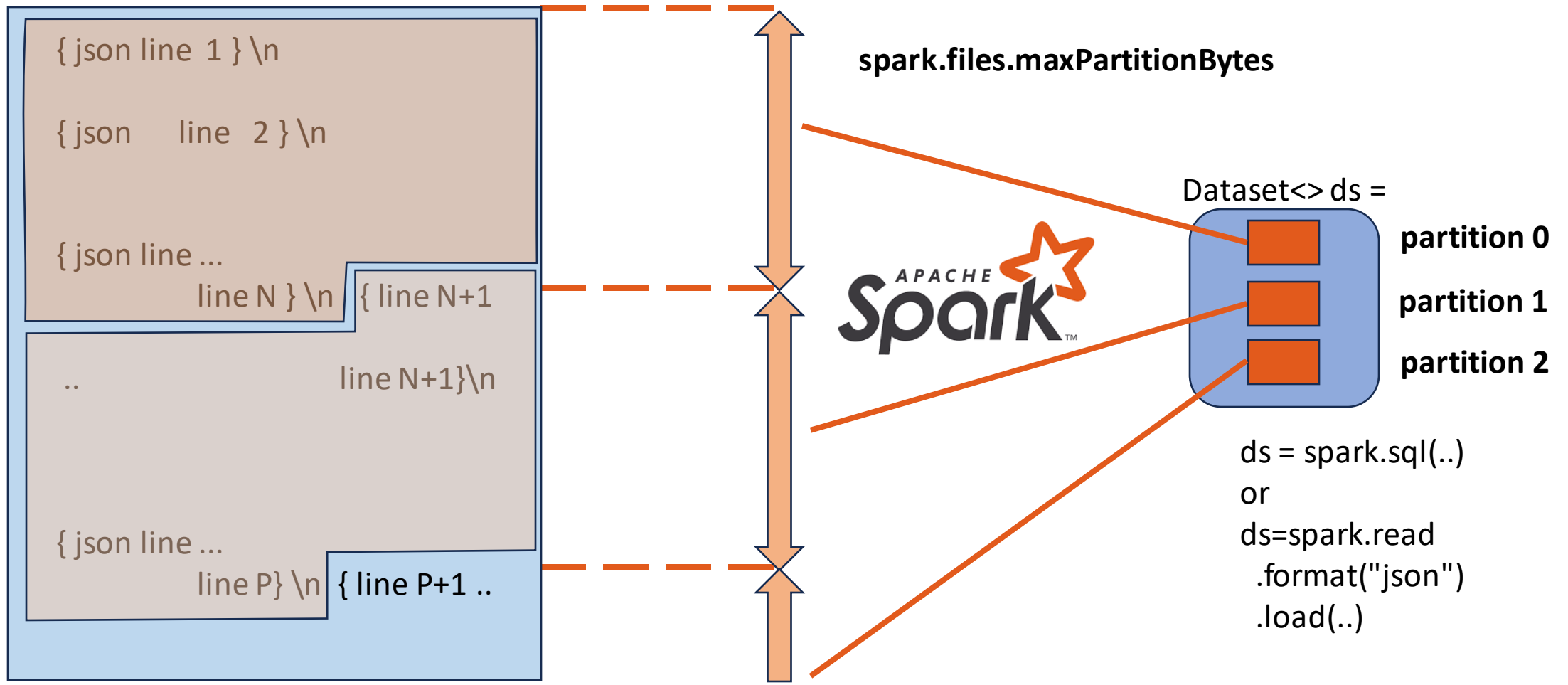


# {Texts,CSV,Json} formats FileSplit

each spark executor reader Thread


at split start => ignore first chars until '\n'

at split end => read extra chars until '\n'





Example : Reading 1 CSV of 3.2 Go  
=> 26 splits = 25 (~128 Mo) + 1 very small

Windows-SSD (C:) > data > OpenData-gouv.fr > bal			
<input type="checkbox"/> Nom	Modifié le	Type	Taille
 adresses-france.csv	08/09/2022 13:49	Fichier CSV Microsoft Excel	3 280 937 Ko


$$3\,280\,937 / (128 * 1024) = 25.03$$

```
scala> val ds = spark.read.format("csv").option("delimiter",";").load("C:/data/OpenData-gouv.fr/bal")
val ds: org.apache.spark.sql.DataFrame = [_c0: string, _c1: string ... 17 more fields]
```

```
scala> ds.toJavaRDD.getNumPartitions
val res0: Int = 26
```

Parquet is Compressable { snappy | gz }

# \*.gz is NOT Splitteable !

Windows-SSD (C:) > data > OpenData-gouv.fr > bal-gz				
<input type="checkbox"/> Nom	Modifié le	Type	Taille	
 adresses-france.csv.gz	08/09/2022 13:49	Dossier d'archive compressé	659 936 Ko	

CSV file was 3.2 Go, now 660 Mo in .gz  
but NOT Splitteable => 1 spark partition, reading (CPU intensive) by 1 Thread only !  
~8 times slower on a 8 cores PC

```
scala> val csvGzDs = spark.read.format("csv").option("delimiter",";").load("C:/data/OpenData-gouv.fr/bal-gz")
23/11/05 20:17:12 WARN ZlibFactory: Failed to load/initialize native-zlib library
val csvGzDs: org.apache.spark.sql.DataFrame = [_c0: string, _c1: string ... 17 more fields]

scala> csvGzDs.count()
[Stage 2:> (0 + 1) / 1]
```

# Compressions Algorithm

data

0101010101001



**.snappy**

010101

fast compression/decompression  
(focus on speed)

**.gz (gzip)**

01101

slower-compression,  
better compression ratio  
(focus on size)

**.lz4**

010101

compromise between  
fast / compression ratio

# PARQUET.{snappy | gz}

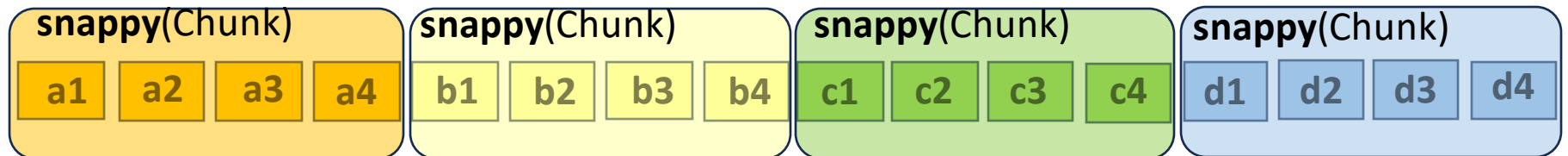
Logical view: rows - columns

a1	b1	c1	d1
a2	b2	c2	d2
a3	b3	c3	d3
a4	b4	c4	d4

On Disk: **Columnar NO compression**



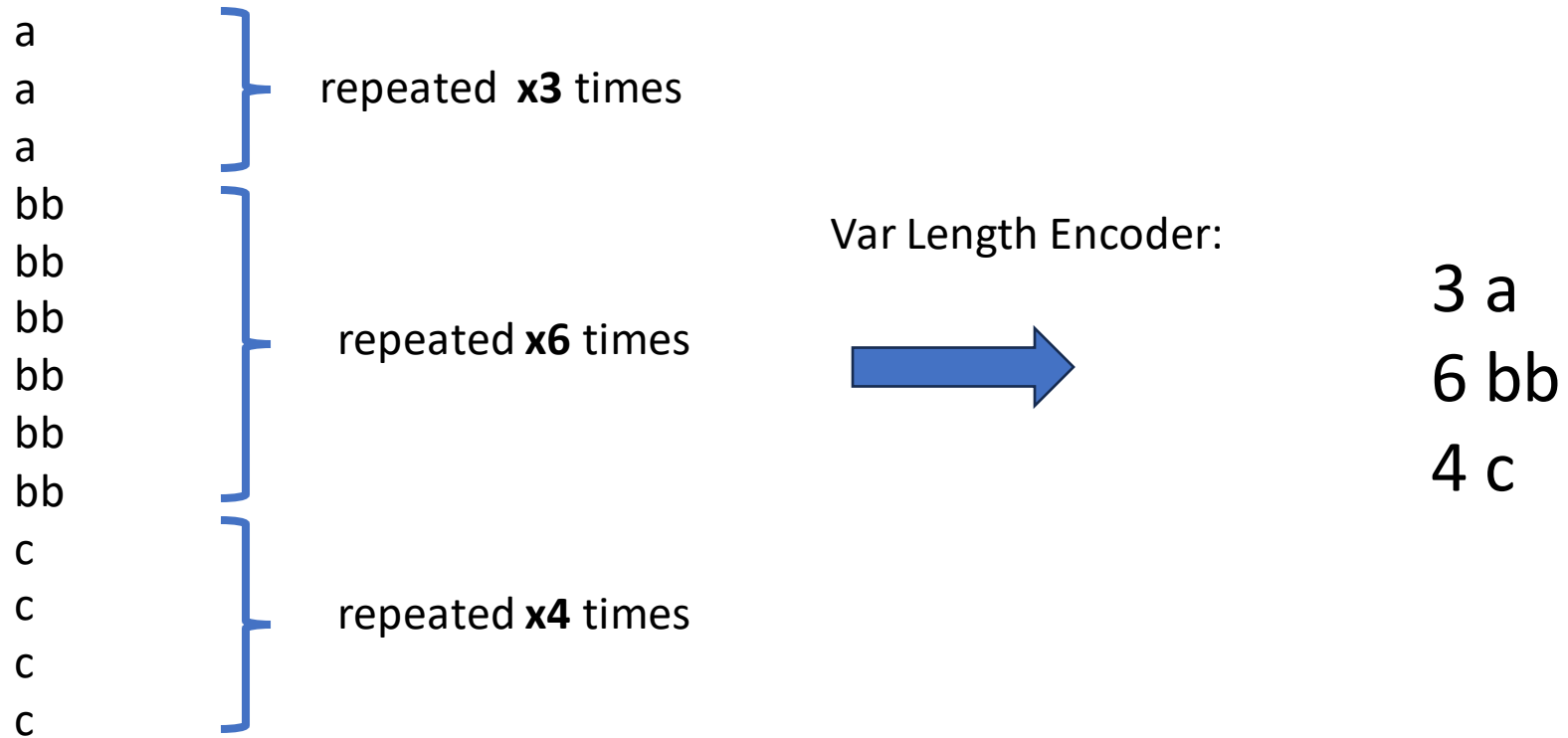
On Disk: **Columnar parquet.{snappy | gz }**



every "Chunk" of column data are compressed INDEPENDENTLY  
=> file is STILL Splitteable

Parquet use Encodings

# Run-length encoding (RLE)



# Size of Adding "Constant" Column

adding a column with only "0" for Billions of rows

=> take only ~100 bytes per RowGroup



# Dictionary Encoding

Manchester City,  
Arsenal,  
Manchester City,  
FC Barcelone,  
Arsenal,  
Newcastle,  
Manchester United,  
Newcastle,  
FC Barcelone,  
Arsenal,  
Manchester United,  
...

Dictionary Encoder:



## **Distinct Dictionary Values:**

1=Manchester City

2=Arsenal

3=FC Barcelone

4=Newcastle

5..

## **Value Indexes:**

1, 2, 1, 3, 4, 5 ..

# Dictionary Size Limit

By default, Dictionary size = max 1 Mega (per RowGroup - Column Chunk )

When using

Huge RowGroup (> 128Mo) => less Dictionaries used

Small RowGroup (32M, 64Mo) => more Dictionaries

Parameters:

`parquet.enable.dictionary=true` (default)

**`parquet.dictionary.page.size=1M`**

# Delta Encoding

**BaseValue=10007**

**MinDelta=-7**

**Delta**

**Delta-MinDelta (>=0)**

10007

-4

+7 →

3

10003

+2

9

10005

-3

4

10002

+7

14

10009

-7

+7 →

0

10002

+6

13

10008

..

..

Delta Encoding



10007, -7, 3,9,4,14,0,13

example size when  
fitting "int"(4 bytes)  
 $7 \times 4 = 28$  bytes

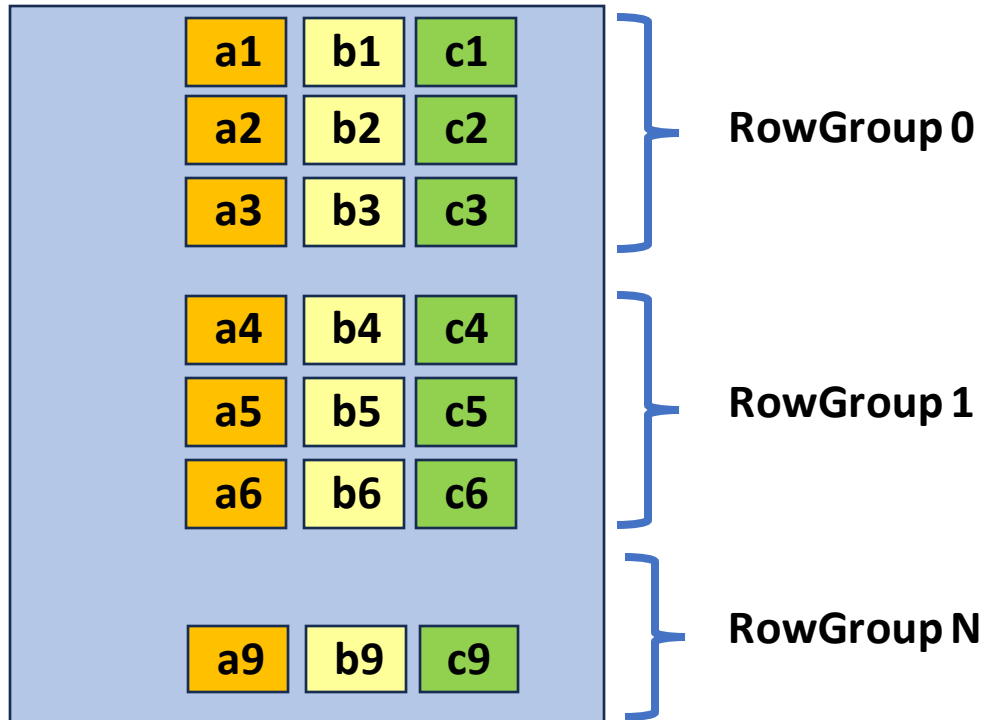
when fitting 2 bytes  
 $4 + 6 \times 2 = 28$  bytes

when fitting 1 bytes  
 $4 + 2 + 6 \times 1 = 12$  bytes

# Parquet use Statistics

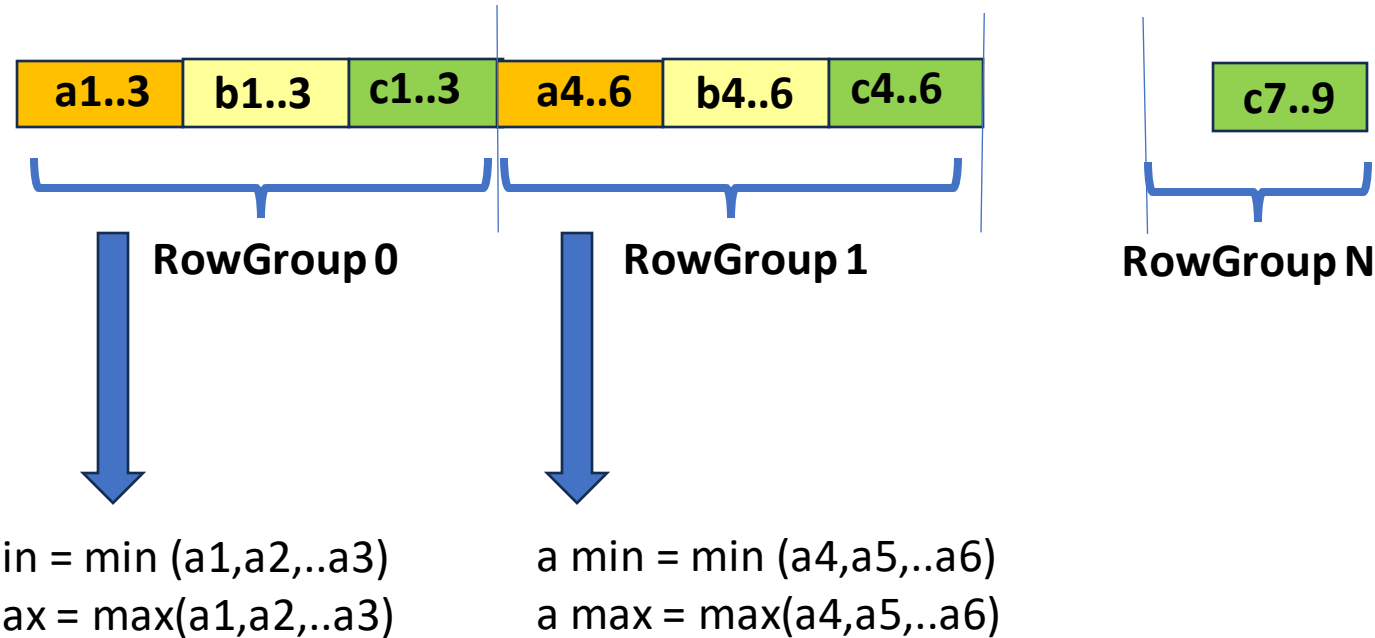
# Column Statistics: min/max Value per RowGroup

Logical view: rows - columns



Schema: {a int, b string, c }  
RowGroups:  
rg0: {offset, a(min,max),  
      b(min,max),c(min,max)}  
rg1: { offset, a(min,max) ..}

On Disk: **columnar**



# Reading Metadata

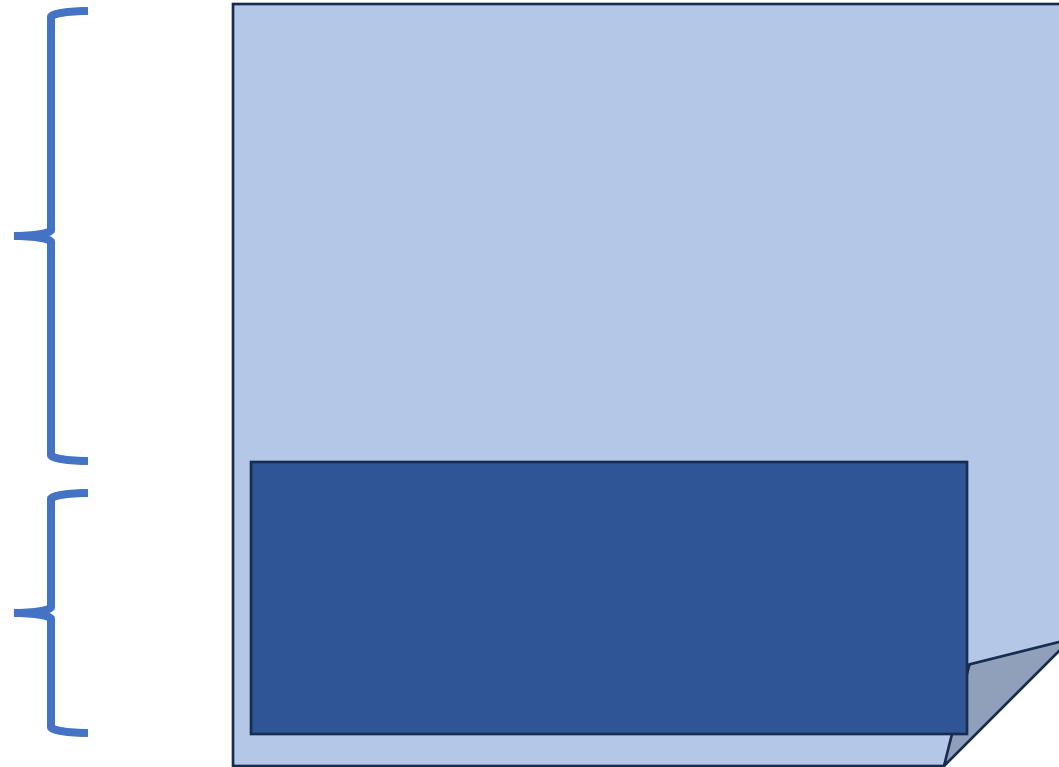
=> read schema + offset + statistics

data part

(usually N x blocks of 128 Mo)

footer: metadata part

(usually  $\leq 50$  ko )

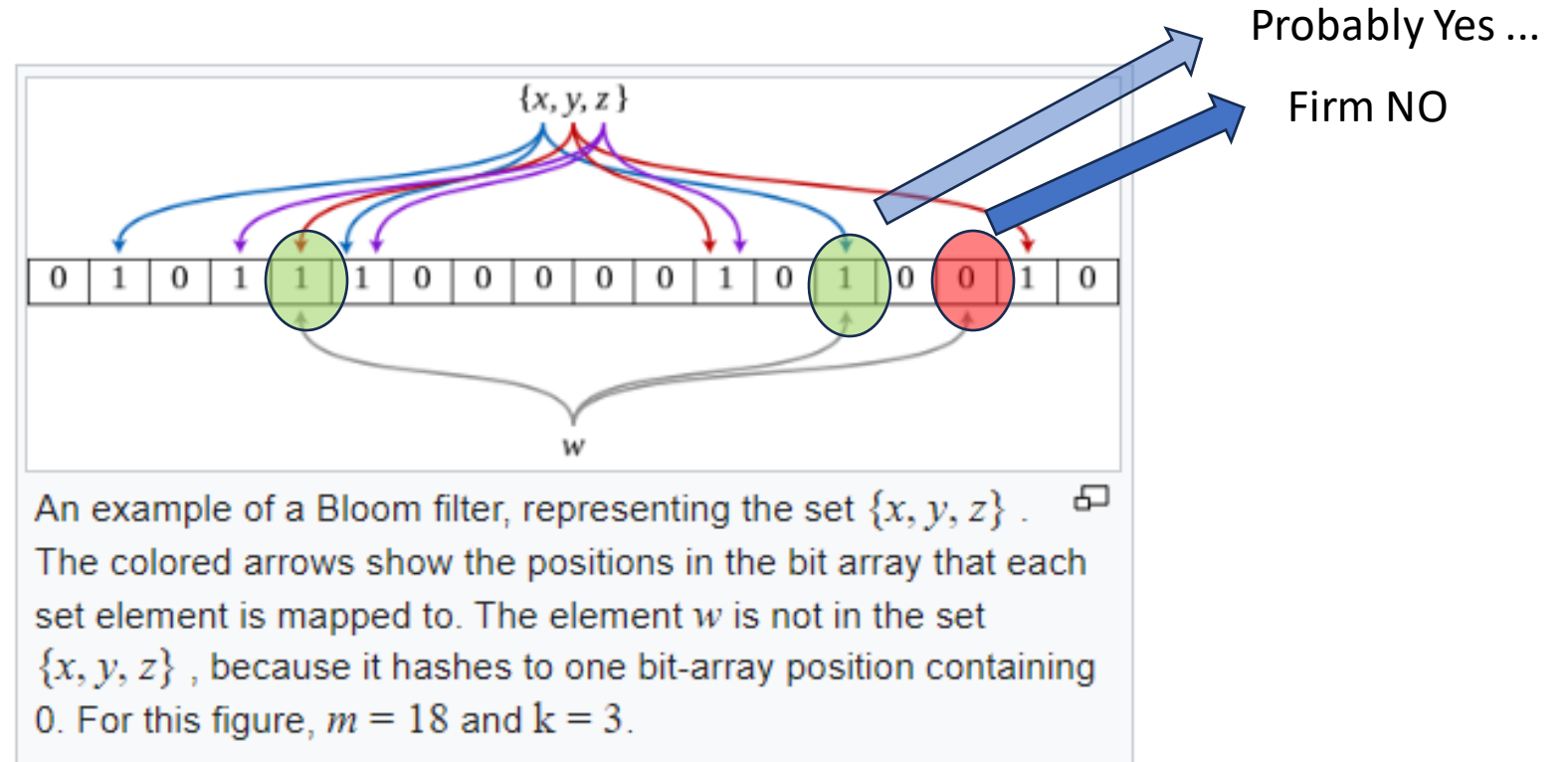


Parquet use Bloom Filter

# Bloom Filter

is  $W$  in the set  $\{x, y, z\}$  ?

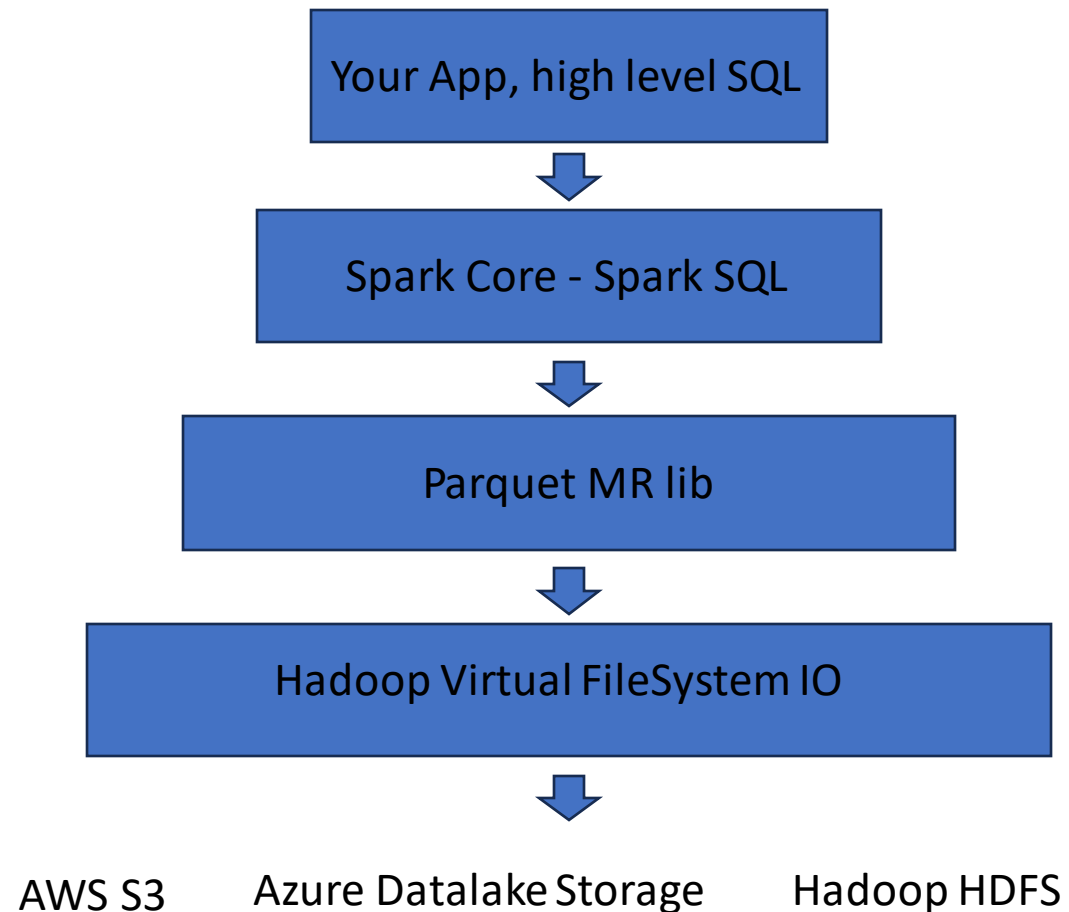
$\text{hash}(W) = 0000100\dots10100$



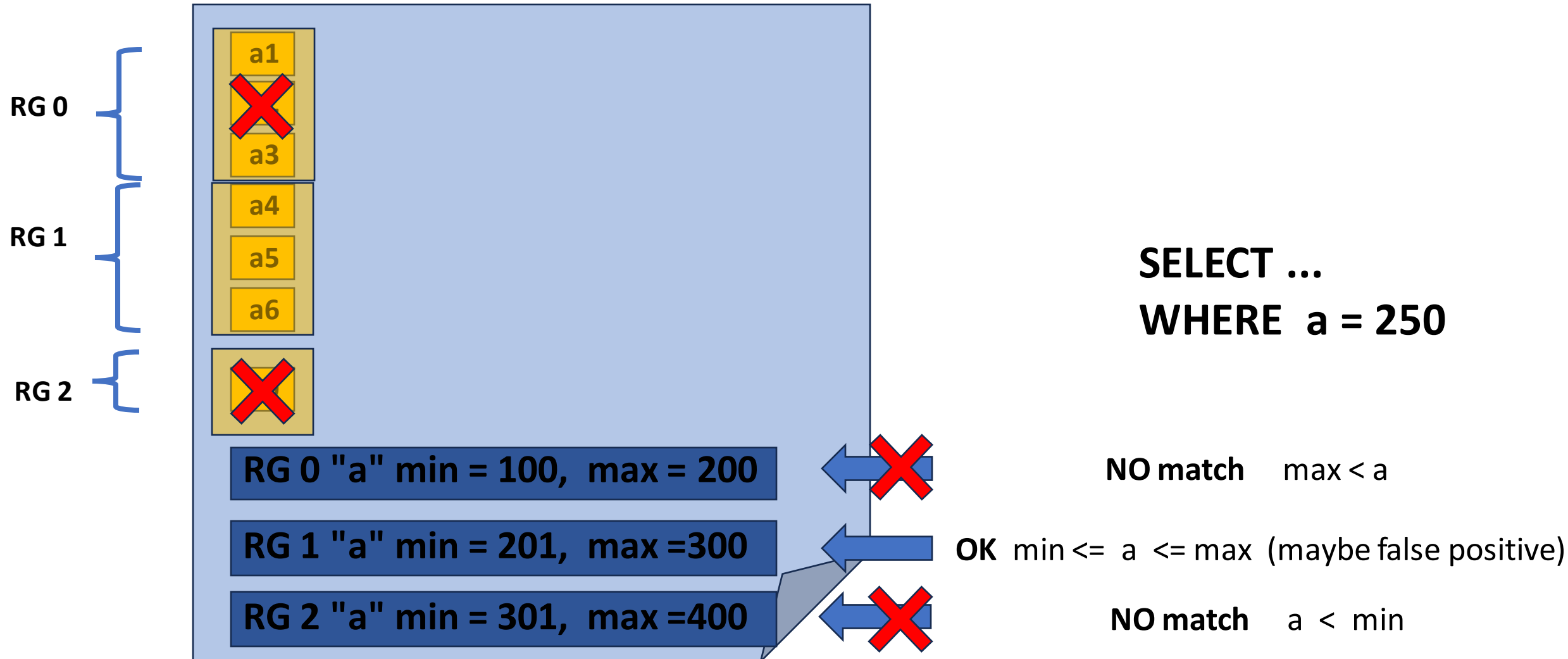


# Parquet Predicate-Push-Down

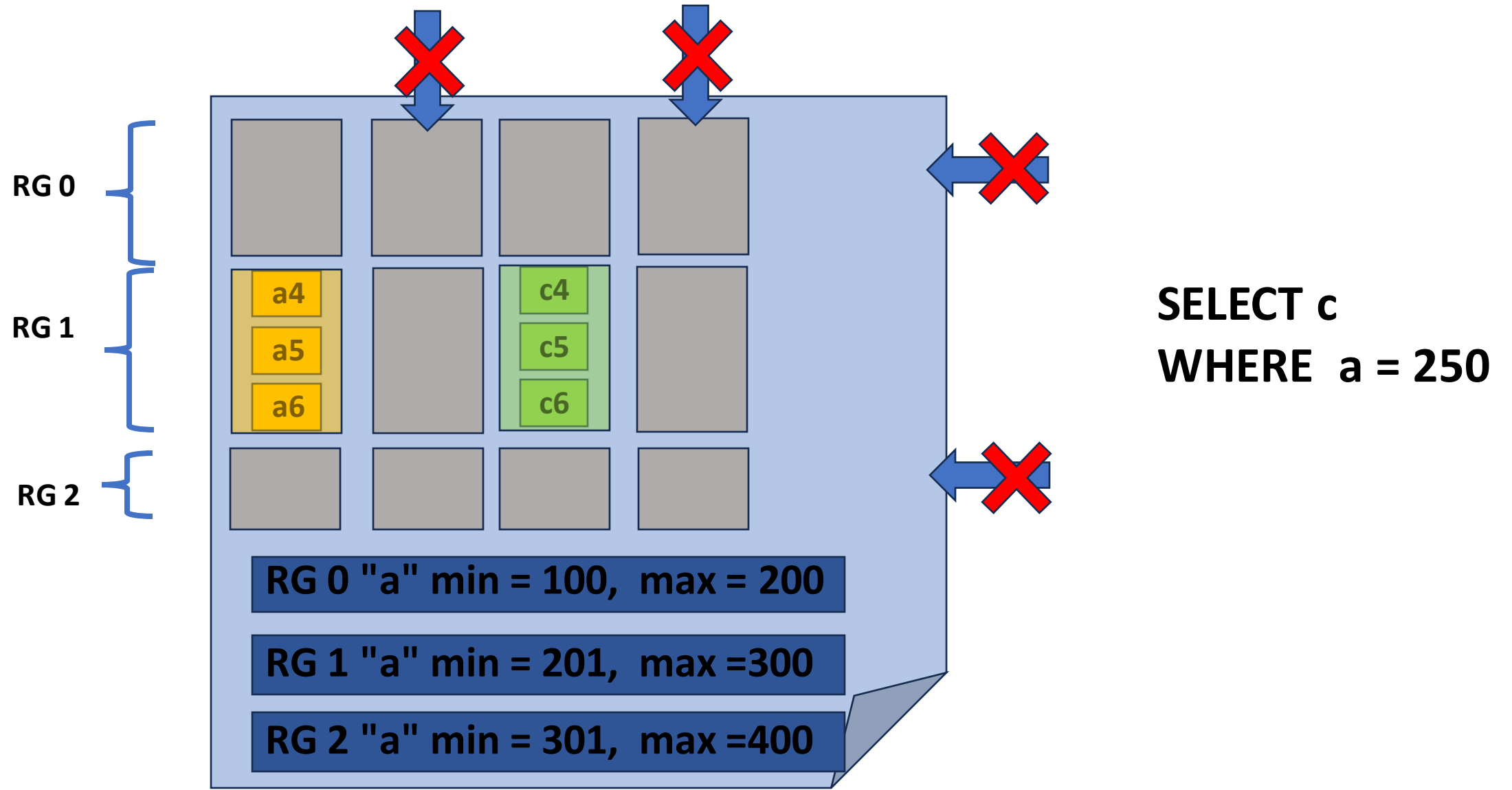
Push-Down :  
from Spark -> Parquet Lib -> IO Storage



# Statistics for skip read RowGroup SELECT .. WHERE column=value

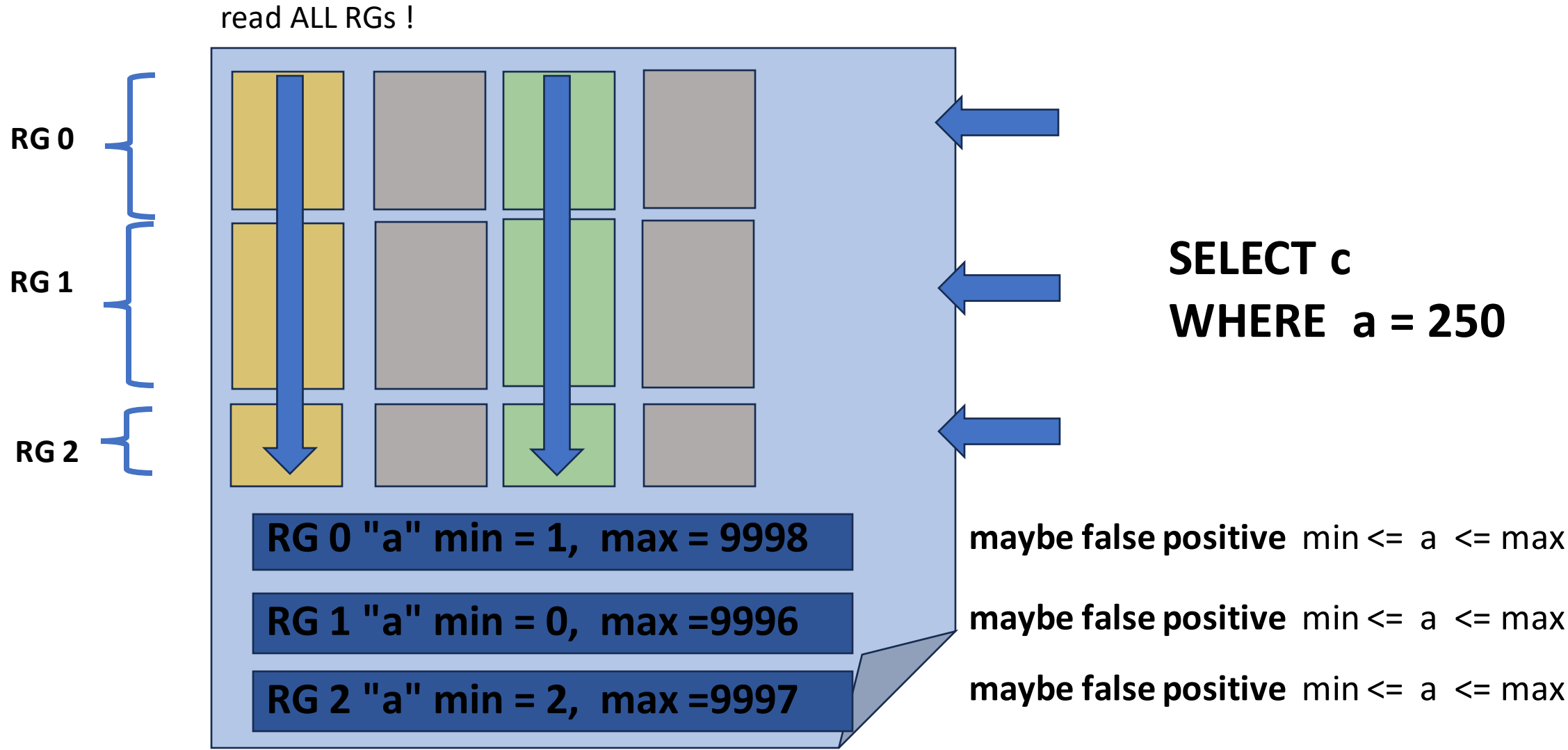


# Column Pruning + RowGroup Pruning(PPD) = minimal Reads



# Badly sorted files

=> Bad min/max statistics => False positives



# Optim Write Once / Read Many

when writing Parquet files ... think how file might be read later !

spend CPU when writing to save CPU / IO later reading

```
dataset
    .repartition(nRepartitionCount)
    // or .repartition("col1", nRepartitionHash)

    .sortWithinPartitions("colA", "colB")

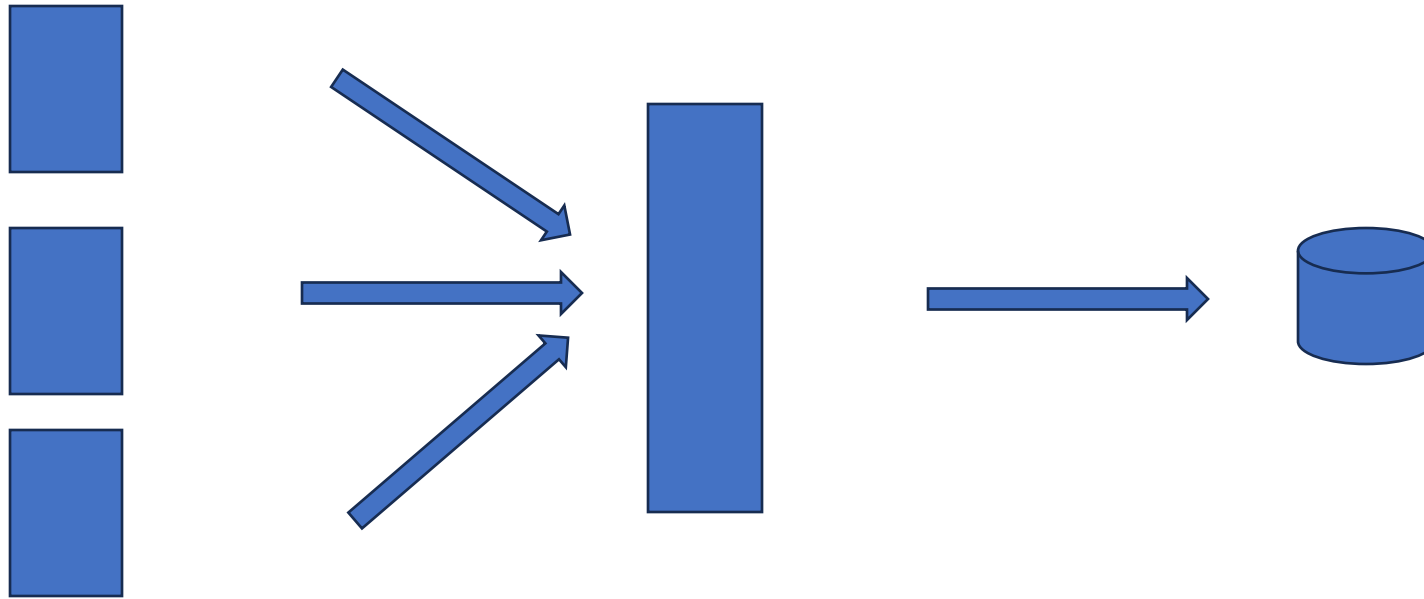
    .write
    .option("parquet.block.size", 32*MEGA)

    .format("parquet")
    .save("file://some-dir")
    // or
    // .format("hive").insertInto("hiveDb.hiveTableName")
```

avoid many small files  
dataset.repartition(1)  
or .coalesce(1)

dataset

**.repartition**(nRepartitionCount) // or .coalesce(nRepartitionCount)



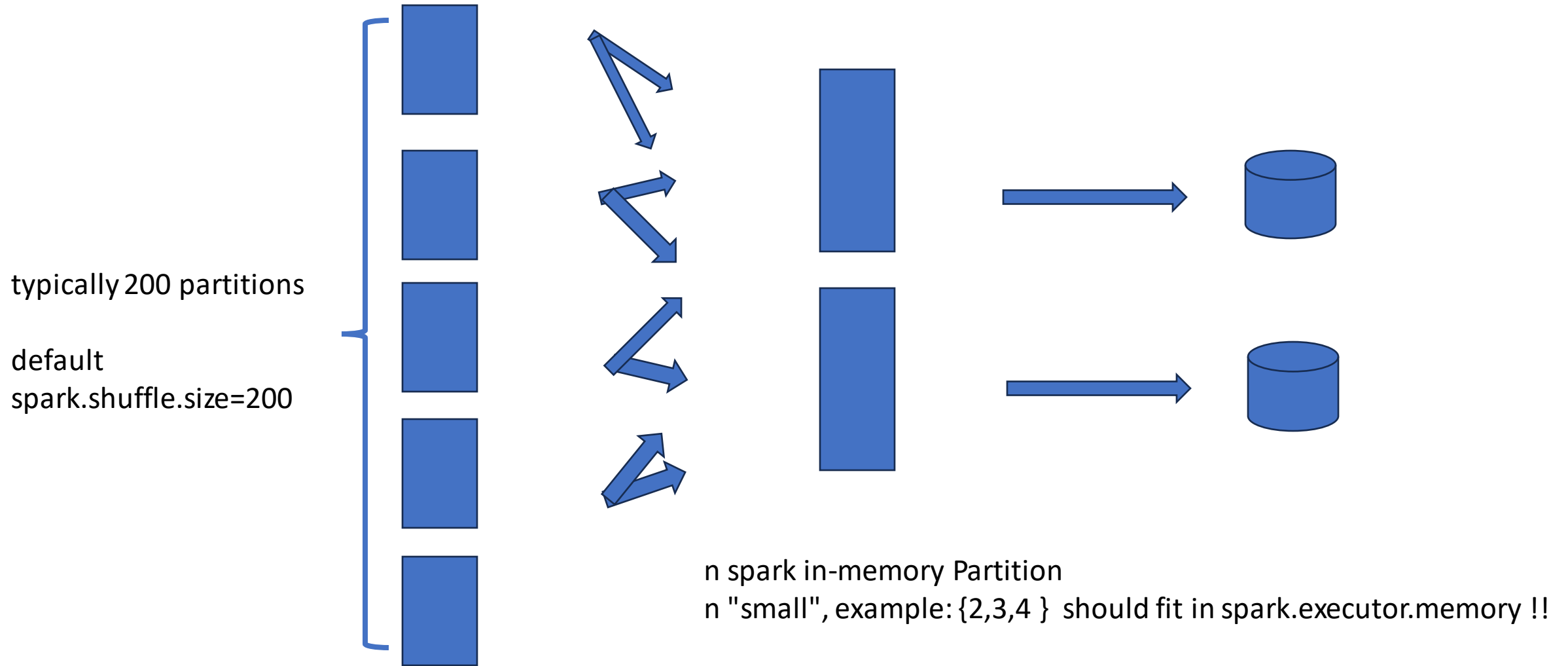
N spark in-memory Partitions  
distributed over N executors

1 spark in-memory Partition  
(should fit in spark.executor.memory !!)

write 1 parquet file



# Does not fit in memory.. compromise to .repartition(smallN)

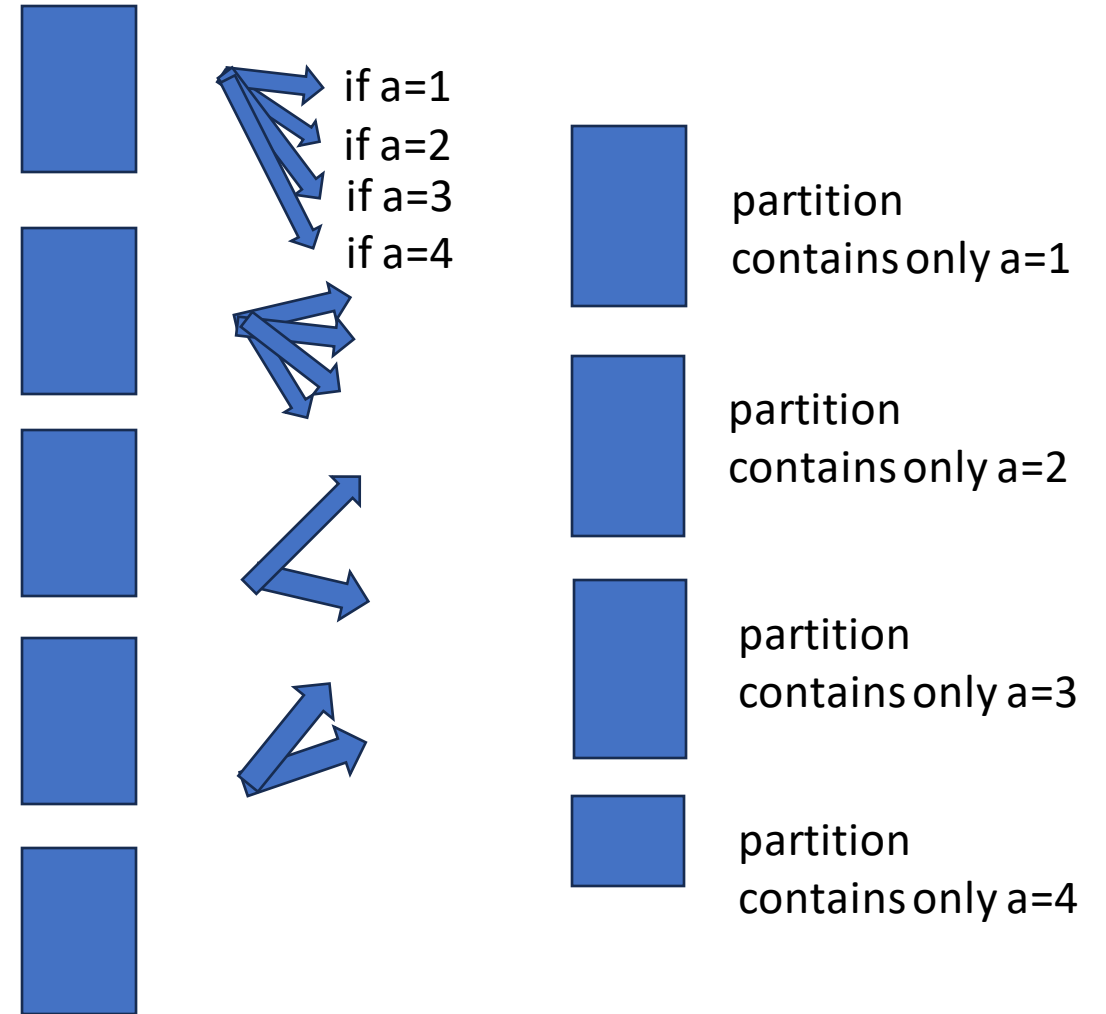


# .repartition("column")

typical usage:

column "a" has FEW distinct values { 1, 2, 3, 4 }

```
ds = dataset.repartition("a")
```



# .repartition("column", nHashCount)

example ... having many distinct values

col "a" values in [ 1, 2, 3, .... 500000]

=>  $h = \text{hash}(a) \% 15$  in [0, 1, .. 14]



partition  
contains only  $a \% 15 == 0$  i.e. {0, 15, 30, 45, ..}



partition  
contains only  $a \% 15 == 1$  i.e. {1, 16, 31, 46, ..}

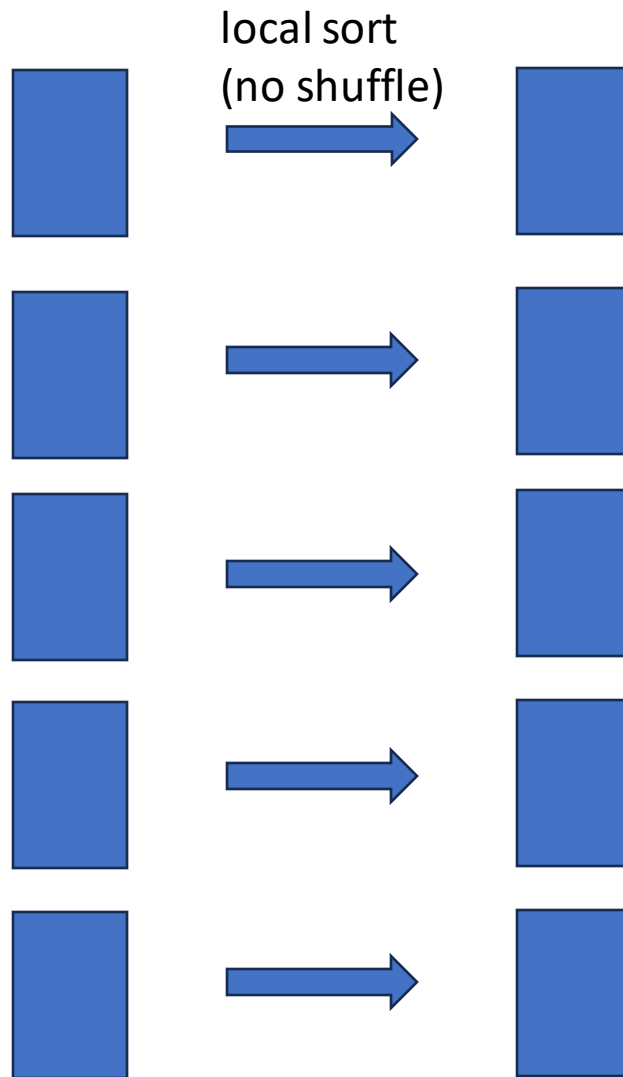


partition  
contains only  $a \% 15 == 2$  i.e. {2, 12, 32, 47, ..}

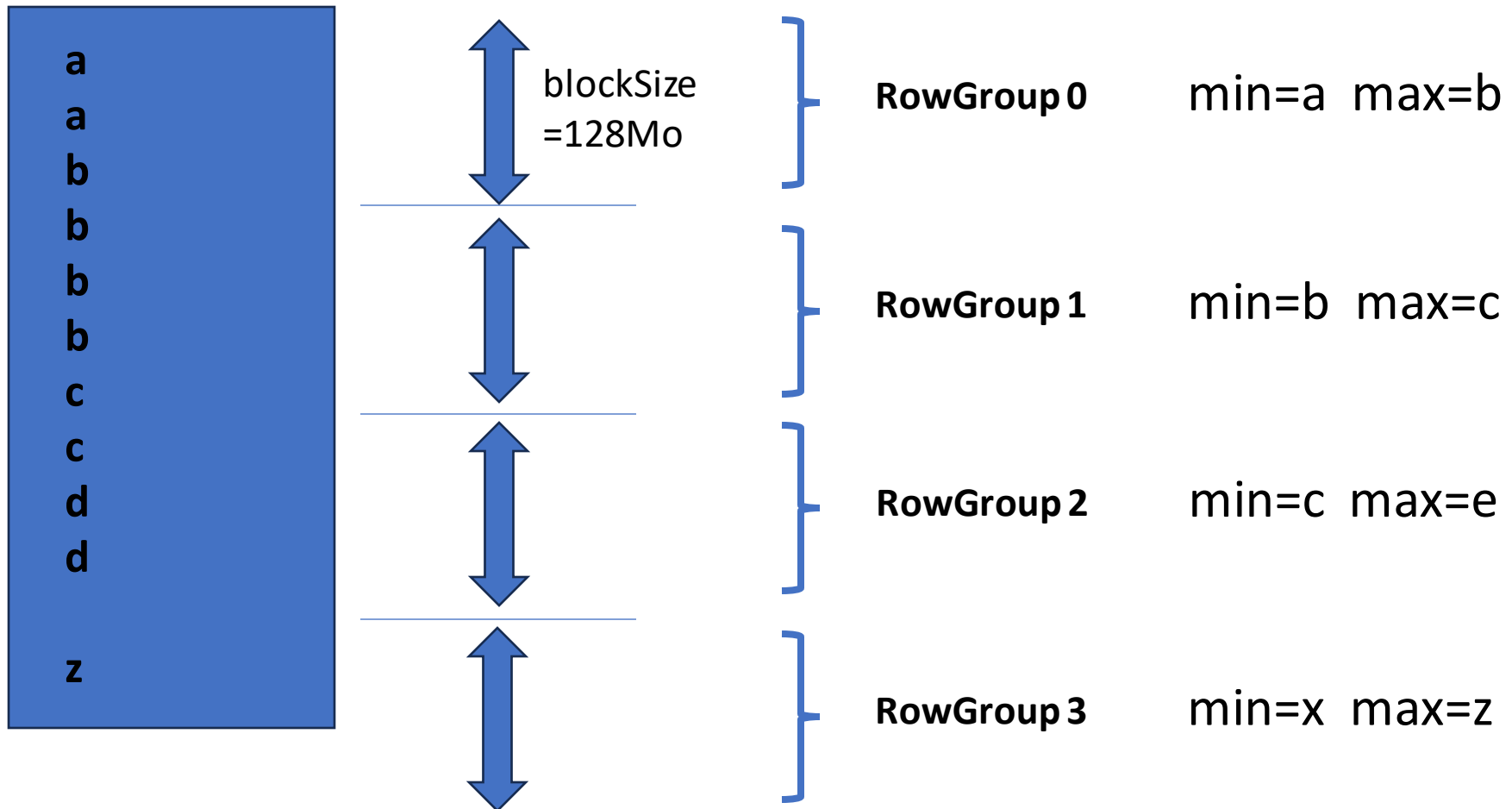


partition  
contains only  $a \% 15 == 14$  i.e. {14, 29, 44, ..}

`.sortWithinPartition("a", "b", ...)`



.sortWithinPartition  
=> RowGroups stats more "compact"



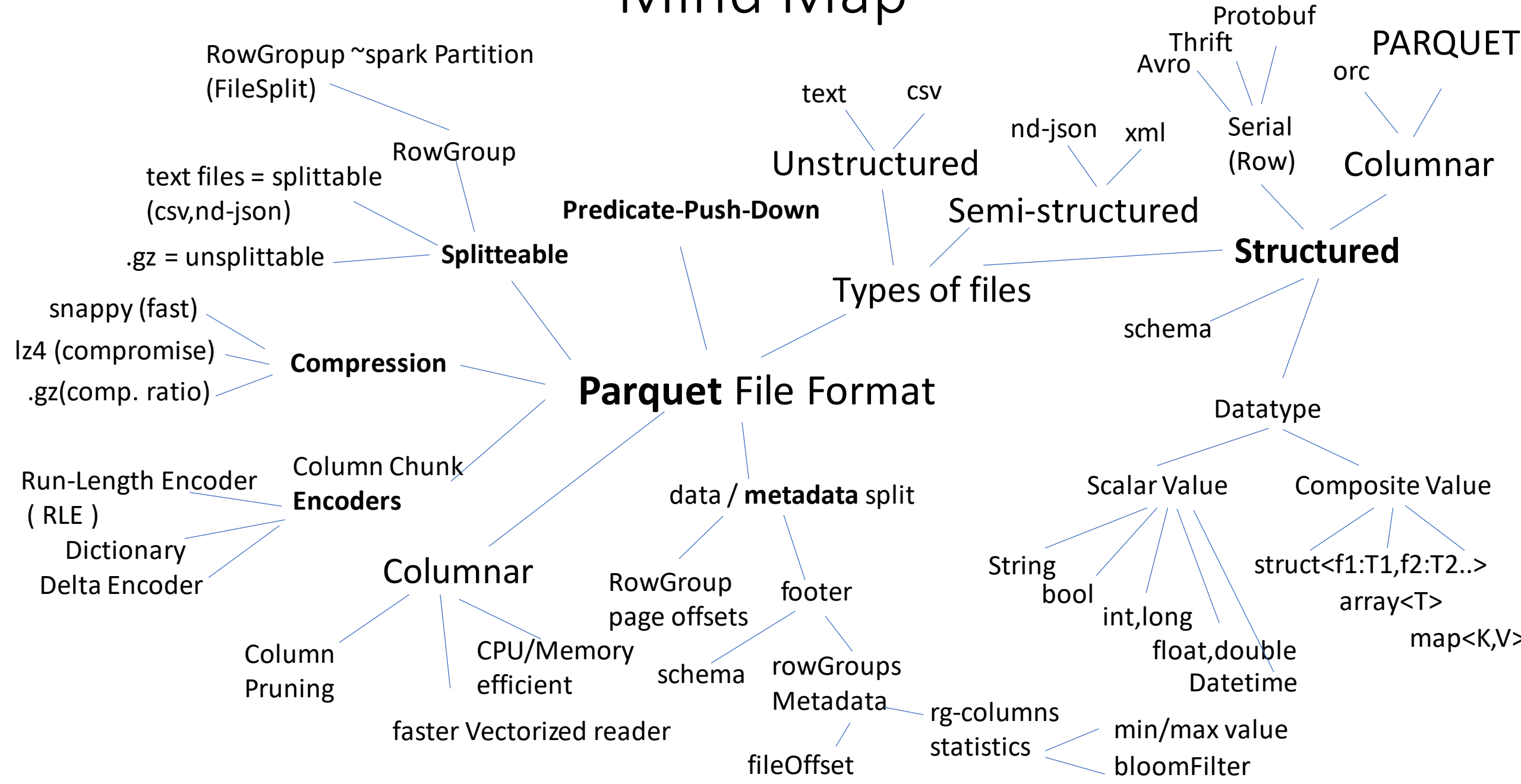
# Conclusion

Parquet File Format is AMAZING

Spark is great using Parquet

Doing BigData processing = doing Spark + Parquet  
with good `.repartition().sortWithinPartition()` ...

# Mind Map



Questions?

[arnaud.nauwynck@gmail.com](mailto:arnaud.nauwynck@gmail.com)