

# Parquet Basics

03-06-2025

## Welcome to the fantastic world of Parquet

### What is it?

Apache Parquet is an open source, column-oriented data file format designed for efficient data storage and retrieval. It provides high performance compression and encoding schemes to handle complex data in bulk and is supported in many programming language and analytics tools<sup>1</sup>.

The open-source project to build Apache Parquet began as a joint effort between Twitter[3] and Cloudera.[4] Parquet was designed as an improvement on the Trevni columnar storage format created by Doug Cutting, the creator of Hadoop. The first version, Apache Parquet 1.0, was released in July 2013. Since April 27, 2015, Apache Parquet has been a top-level Apache Software Foundation (ASF)-sponsored project.

### Why should we care?

Apache Parquet is a file format designed to support fast data processing for complex data, with several notable characteristics:

1. **Columnar:** Unlike row-based formats such as CSV, Apache Parquet is column-oriented – meaning the values of each table column are stored next to each other, rather than those of each record:

---

<sup>1</sup><https://parquet.apache.org/docs/overview/>

ROW-BASED STORAGE	COLUMNAR STORAGE
1 MARC, JOHNSON, WASHINGTON, 27	ID: 1 2 3
2 JIM, THOMPSON, DENVER, 33	FIRST NAME: MARC, JIM, JACK
3 JACK, RILEY, SEATTLE, 51	LAST NAME: JOHNSON, THOMPSON, RILEY
	CITY: WASHINGTON, DENVER, SEATTLE
	AGE: 27 33 51

The key difference between a CSV and Parquet file format is how each one is organized. A Parquet file format is structured by row, with every separate column independently accessible from the rest. Since the data in each column is expected to be of the same type, the parquet file format makes encoding, compressing and optimizing data storage possible.

2. **Open-source:** Parquet is free to use and open source under the Apache Hadoop license.

Apache Parquet is a columnar storage format available to any project [...], regardless of the choice of data processing framework, data model or programming language<sup>2</sup>.

3. **Self-describing:** In addition to data, a Parquet file contains metadata including schema and structure. Each file stores both the data and the standards used for accessing each record – making it easier to decouple services that write, store, and read Parquet files.
4. **Binary format:** Parquet file formats store data in binary format, which reduces the overhead of textual representation. It's important to note that Parquet files are not stored in plain text, thus cannot be opened in a text editor.

## Advantages of Parquet Columnar Storage – Why Should You Use It?

The above characteristics of the Apache Parquet file format create several distinct benefits when it comes to storing and analysing large volumes of data.

### Compression

File compression is the act of taking a file and making it smaller. In Parquet, compression is performed column by column and it is built to support flexible compression options and extendable encoding schemas per data type – e.g., different encoding can be used for compressing integer and string data.

---

<sup>2</sup><https://parquet.apache.org/>

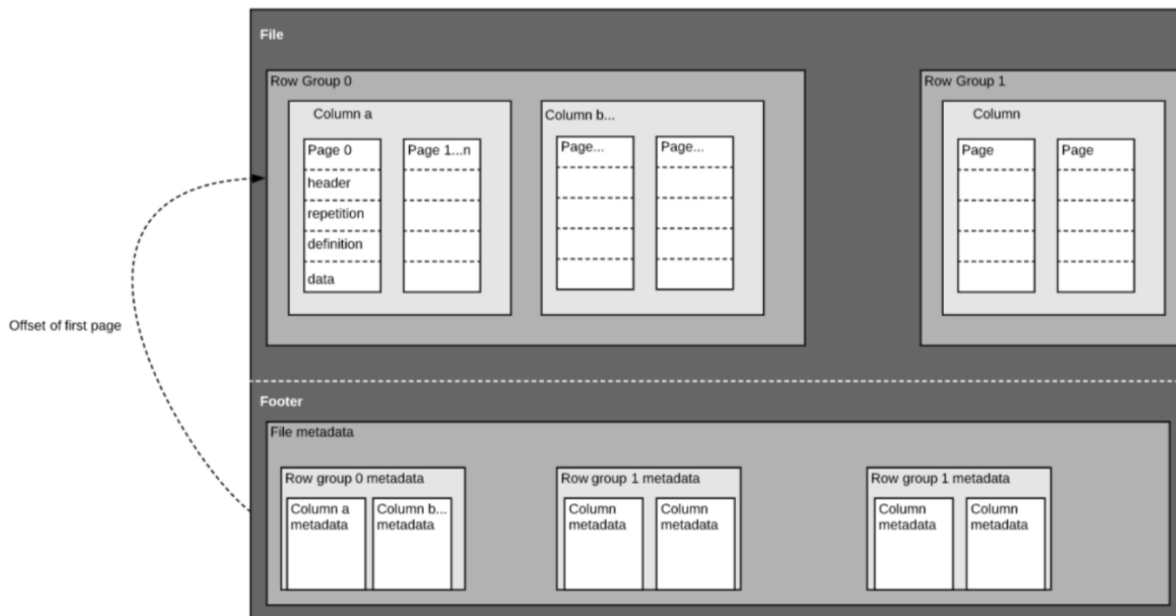
Parquet data can be compressed using these encoding methods:

- **Dictionary encoding:** this is enabled automatically and dynamically for data with a small number of unique values.
- **Bit packing:** Storage of integers is usually done with dedicated 32 or 64 bits per integer. This allows more efficient storage of small integers.
- **Run length encoding (RLE):** when the same value occurs multiple times, a single value is stored once along with the number of occurrences. Parquet implements a combined version of bit packing and RLE, in which the encoding switches based on which produces the best compression results.

## Performance

As opposed to row-based file formats like CSV, Parquet is optimized for performance. When running queries on your Parquet-based file-system, you can focus only on the relevant data very quickly. Moreover, the amount of data scanned will be way smaller and will result in less I/O usage. To understand this, let's look a bit deeper into how Parquet files are structured.

As we mentioned above, Parquet is a self-described format, so each file contains both data and metadata. Parquet files are composed of row groups, header and footer. Each row group contains data from the same columns. The same columns are stored together in each row group:



This structure is well-optimized both for fast query performance, as well as low I/O (minimizing the amount of data scanned). For example, if you have a table with 1000 columns, which you will usually only query using a small subset of columns. Using Parquet files will enable

you to fetch only the required columns and their values, load those in memory and answer the query. If a row-based file format like CSV was used, the entire table would have to have been loaded in memory, resulting in increased I/O and worse performance.

## **Schema evolution**

When using columnar file formats like Parquet, users can start with a simple schema, and gradually add more columns to the schema as needed. In this way, users may end up with multiple Parquet files with different but mutually compatible schemas. In these cases, Parquet supports automatic schema merging among these files.

## **Column-Oriented vs Row-Based Storage for Analytic Querying**

Data is often generated and more easily conceptualized in rows. We are used to thinking in terms of Excel spreadsheets, where we can see all the data relevant to a specific record in one neat and organized row. However, for large-scale analytical querying, columnar storage comes with significant advantages with regards to cost and performance.

Complex data such as logs and event streams would need to be represented as a table with hundreds or thousands of columns, and many millions of rows. Storing this table in a row based format such as CSV would mean:

- Queries will take longer to run since more data needs to be scanned, rather than only querying the subset of columns we need to answer a query (which typically requires aggregating based on dimension or category)
- Storage will be more costly since CSVs are not compressed as efficiently as Parquet

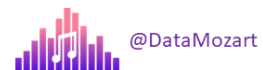
Columnar formats provide better compression and improved performance out-of-the-box, and enable you to query data vertically – column by column.

	Product	Customer	Country	Date	Sales Amount
Row 1	Ball	John Doe	USA	2023-01-01	100
Row 2	T-Shirt	John Doe	USA	2023-01-02	200
Row 3	Socks	Maria Adams	UK	2023-01-01	300
Row 4	Socks	Antonio Grant	USA	2023-01-03	100
Row 5	T-Shirt	Maria Adams	UK	2023-01-02	500
Row 6	Socks	John Doe	USA	2023-01-05	200



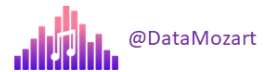
images columns row<sup>3</sup>

	Product	Customer	Country	Date	Sales Amount
Row 1	Ball	John Doe	USA	2023-01-01	100
Row 2	T-Shirt	John Doe	USA	2023-01-02	200
Row 3	Socks	Maria Adams	UK	2023-01-01	300
Row 4	Socks	Antonio Grant	USA	2023-01-03	100
Row 5	T-Shirt	Maria Adams	UK	2023-01-02	500
Row 6	Socks	John Doe	USA	2023-01-05	200



<sup>3</sup><https://data-mozart.com/parquet-file-format-everything-you-need-to-know/>

Column 1	Column 2	Column 3	Column 4	Column 5
Product	Customer	Country	Date	Sales Amount
Ball	John Doe	USA	2023-01-01	100
T-Shirt	John Doe	USA	2023-01-02	200
Socks	Maria Adams	UK	2023-01-01	300
Socks	Antonio Grant	USA	2023-01-03	100
T-Shirt	Maria Adams	UK	2023-01-02	500
Socks	John Doe	USA	2023-01-05	200



	Column 1	Column 2	Column 3	Column 4	Column 5
	Product	Customer	Country	Date	Sales Amount
Row Group 1	Ball	John Doe	USA	2023-01-01	100
	T-Shirt	John Doe	USA	2023-01-02	200
Row Group 2	Socks	Maria Adams	UK	2023-01-01	300
	Socks	Antonio Grant	USA	2023-01-03	100
Row Group 3	T-Shirt	Maria Adams	UK	2023-01-02	500
	Socks	John Doe	USA	2023-01-05	200



	Column 1	Column 2	Column 3	Column 4	Column 5
	Product	Customer	Country	Date	Sales Amount
Row Group 1	Ball	John Doe	USA	2023-01-01	100
	T-Shirt	John Doe	USA	2023-01-02	200
Row Group 2	Socks	Maria Adams	UK	2023-01-01	300
	Socks	John Doe	USA	2023-01-02	100
Row Group 3	T-Shirt	Maria Adams	UK	2023-01-02	500
	Socks	John Doe	USA	2023-01-05	200

The engine will not scan these records



Let's quickly stop here, as I want you to realize the difference between various types of storage in terms of the work that needs to be performed by the engine:

Row store - the engine needs to scan all 5 columns and all 6 rows

Column store - the engine needs to scan 2 columns and all 6 rows

Column store with row groups - the engine needs to scan 2 columns and 4 rows

***"I'm tired of reading non-sense"***

Ok, then let me introduce you to the R packages I've been exploring and how they made my life easier.

## arrow

The R arrow package provides access to many of the features of the Apache Arrow C++ library for R users. The goal of arrow is to provide an Arrow C++ backend to dplyr, and access to the Arrow C++ library through familiar base R and tidyverse functions, or R6 classes. The dedicated R package website is located [here](#).

## What can the arrow package do?

The arrow package provides binding to the C++ functionality for a wide range of data analysis tasks.

- It allows users to read and write data in a variety of formats:
  - Read and write Parquet files, an efficient and widely used columnar format
  - Read and write Arrow (formerly known as Feather) files, a format optimized for speed and interoperability
  - Read and write CSV files with excellent speed and efficiency
  - Read and write multi-file and larger-than-memory datasets
  - Read JSON files
- It provides access to remote file systems and servers:
  - Read and write files in Amazon S3 and Google Cloud Storage buckets
  - Connect to Arrow Flight servers to transport large datasets over networks
- Additional features include:
  - Manipulate and analyse Arrow data with dplyr verbs
  - Zero-copy data sharing between R and Python
  - Fine control over column types to work seamlessly with databases and data warehouses
  - Toolkit for building connectors to other applications and services that use Arrow

## What is Apache Arrow?

Apache Arrow is a cross-language development platform for in-memory and larger-than-memory data. It specifies a standardized language-independent columnar memory format for flat and hierarchical data, organized for efficient analytic operations on modern hardware. It also provides computational libraries and zero-copy streaming, messaging, and interprocess communication.

## Arrow resources

There are a few additional resources that you may find useful for getting started with arrow:

- The official [Arrow R package documentation](#)
- [Arrow for R cheatsheet](#)
- [Apache Arrow R Cookbook](#)
- R for Data Science [Chapter on Arrow](#)
- [Awesome Arrow R](#)



## Installation

The latest release of arrow can be installed from CRAN. In most cases installing the latest release should work without requiring any additional system dependencies, especially if you are using Windows or macOS.

```
install.packages("arrow")
```

If you are having trouble installing from CRAN, then we offer two alternative install options for grabbing the latest arrow release. First, R-universe provides pre-compiled binaries for the most commonly used operating systems.

```
install.packages("arrow", repos = c("https://apache.r-universe.dev", "https://cloud.r-project.org"))
```

## Get started with Arrow

```
library(arrow)
library(dplyr)
```

```
starwars
```

```
# A tibble: 87 x 14
  name      height  mass hair_color skin_color eye_color birth_year sex  gender
  <chr>    <int> <dbl> <chr>      <chr>      <chr>      <dbl> <chr> <chr>
1 Luke Sk~    172    77 blond      fair        blue        19   male masculi~
2 C-3PO      167    75 <NA>        gold        yellow      112  none masculi~
3 R2-D2       96    32 <NA>        white, bl~ red         33  none masculi~
4 Darth V~   202   136 none        white        yellow      41.9 male masculi~
5 Leia Or~   150    49 brown      light        brown        19  fema~ femini~
6 Owen La~   178   120 brown, gr~ light        blue         52  male masculi~
7 Beru Wh~   165    75 brown      light        blue         47  fema~ femini~
8 R5-D4       97    32 <NA>        white, red red         NA  none masculi~
9 Biggs D~   183    84 black      light        brown        24  male masculi~
10 Obi-Wan~   182    77 auburn, w~ fair        blue-gray    57  male masculi~
# i 77 more rows
# i 5 more variables: homeworld <chr>, species <chr>, films <list>,
#   vehicles <list>, starships <list>
```

```
write_parquet(starwars, 'outputs/starwars.parquet')
```

```
sw_frame <- read_parquet('outputs/starwars.parquet')  
sw_table <- read_parquet('outputs/starwars.parquet', as_data_frame = FALSE)  
sw_table
```

```
Table  
87 rows x 14 columns  
$name <string>  
$height <int32>  
$mass <double>  
$hair_color <string>  
$skin_color <string>  
$eye_color <string>  
$birth_year <double>  
$sex <string>  
$gender <string>  
$homeworld <string>  
$species <string>  
$films: list<element <string>>  
$vehicles: list<element <string>>  
$starships: list<element <string>>
```

### Multi-file data sets

```
set.seed(1234)  
nrows <- 100000  
random_data <- data.frame(  
  x = rnorm(nrows),  
  y = rnorm(nrows),  
  subset = sample(10, nrows, replace = TRUE)  
)
```

```
dataset_path <- "outputs/random_data"
```

```
random_data %>%  
  group_by(subset) %>%  
  write_dataset(dataset_path)
```

```
list.files(dataset_path, recursive = TRUE)
```

```
[1] "subset=1/part-0.parquet" "subset=10/part-0.parquet"
[3] "subset=2/part-0.parquet" "subset=3/part-0.parquet"
[5] "subset=4/part-0.parquet" "subset=5/part-0.parquet"
[7] "subset=6/part-0.parquet" "subset=7/part-0.parquet"
[9] "subset=8/part-0.parquet" "subset=9/part-0.parquet"
```

Each of these Parquet files can be opened individually using `read_parquet()` but is often more convenient – especially for very large data sets – to scan the folder and “connect” to the data set without loading it into memory. We can do this using `open_dataset()`:

```
dset <- open_dataset(dataset_path)
dset
```

```
FileSystemDataset with 10 Parquet files
3 columns
x: double
y: double
subset: int32
```

### Analysing Arrow data with dplyr

```
dset %>%
  group_by(subset) %>%
  summarize(mean_x = mean(x), min_y = min(y)) %>%
  filter(mean_x > 0) %>%
  arrange(subset) %>%
  collect()
```

```
# A tibble: 6 x 3
  subset mean_x min_y
  <int>   <dbl> <dbl>
1     2 0.00486 -4.00
2     3 0.00440 -3.86
3     4 0.0125  -3.65
4     6 0.0234  -3.88
5     7 0.00477 -4.65
6     9 0.00557 -3.50
```

### example from real life with ebird pts?

```
ebd_filtered_aus_2025 <- arrow::read_delim_arrow("ebd_filtered_aus_2025.txt",
                                                delim = "\t", quote = "")
checklist <- openxlsx::read.xlsx('https://www.birds.cornell.edu/clementschecklist/wp-content,
janitor::clean_names()

checklist |>
  head() |>
  dplyr::glimpse()
ebd_filtered_aus_2025 |>
  left_join(checklist) |> ## by species_code?
  group_by(family) |>
  write_dataset('ebd_filtered_aus_2025_parquet') ## folder
```

```
ebd_filtered_aus_ds <- open_dataset('ebd_filtered_aus_2025_parquet')
ebd_filtered_aus_ds
```

```
## filter category == species
# unique(ebird_data0$category) ### https://support.ebird.org/en/support/solutions/articles
## filter category == species and select some columns to simplify
# ebird_data <- ebird_data0 |>
#   filter(category == "species")

ebd_filtered_aus_ds |>
  filter() |> ## add filters
  collect()

# ebird_data |>
#   count(exotic_code)
```

### geoparquet - sfarrow

### example from real life with rangemaps?

### tidyverse - parquet

<https://www.tidyverse.org/blog/2025/01/nanoparquet-0-4-0/>