

## Parquet

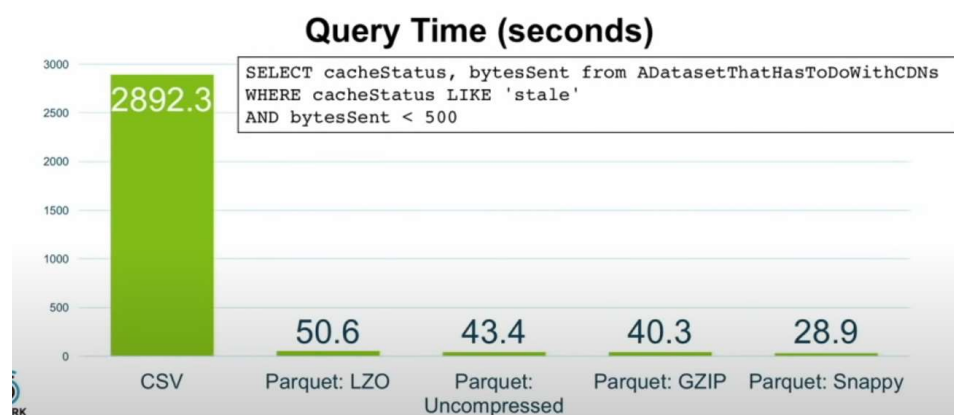
Columnar storage format available to any project in the Hadoop ecosystem, regardless of the choice of data processing framework, data model or programming language.

- Binary format
- API for JVM/Hadoop, C++
- Columnar
- Encoded
- Compressed
- Machine-friendly

## Options For Multi-PB Data Lake Storage

|                     | Files                               | Compressed Files     | Databases                                 |
|---------------------|-------------------------------------|----------------------|---|
| Usability           | Great!                              | Great!               | OK to <b>BAD</b> (not as easy as a file!) |
| Administration      | None!                               | None!                | <b>LOTS</b>                               |
| Spark Integration   | Great!                              | Great!               | Varies                                    |
| Resource Efficiency | <b>BAD</b> (Big storage, heavy I/O) | OK... (Less storage) | <b>BAD</b> (Requires storage AND CPU)     |
| Scalability         | Good-ish                            | Good-ish             | <b>BAD</b> (For multi-petabyte!)          |
| CO\$\$\$\$T         | OK...                               | OK...                | <b>TERRIBLE</b>                           |
| QUERY TIME          | <b>TERRIBLE</b>                     | <b>BAD</b>           | Good!                                     |

## CSV vs. Parquet Column Selection Query



# CSV vs. Parquet Table Scan Query

