# ParquetDB: A Lightweight Python Parquet-Based Database

**Logan L. Lang** [1], **Eduardo Hernandez** [2], **Kamal Choudhary** [3], **and Aldo H. Romero** [1]

**1** Department of Physics, West Virginia University, Morgantown, WV 26506, United States **2** Instituto de Ciencia de Materiales de Madrid, Campus de Cantoblanco, C. Sor Juana Inés de la Cruz, 3, Fuencarral-El Pardo, Madrid 28049, Spain **3** National Institute of Standards and Technology, 100 Bureau Dr, Gaithersburg, MD 20899, United States

## Summary

ParquetDB is a Python library designed to bridge the gap between traditional file storage and fully fledged databases, all while wrapping the PyArrow library to streamline data input and output. By leveraging the Parquet file format, ParquetDB provides the portability and simplicity of file-based data storage alongside advanced querying features typically found in database systems. Because ParquetDB is built on top of PyArrow, it seamlessly handles data types that need special processing to be compatible with the Parquet format. This reduces manual conversion and boilerplate code, allowing developers to focus on higher-level data operations. In addition, the Parquet format's columnar storage and rich metadata make it possible to efficiently perform predicate and column pushdown, leading to faster queries by reading only the subsets of data you truly need.

## Statement of need

In an era where data is the driving force behind innovation, the demand for highly efficient, scalable, and adaptable storage solutions has never been greater. Traditional file-based storage formats (e.g., CSV, JSON, TXT) and database systems (e.g., SQLite(*About SQLite*, n.d.), MongoDB(*MongoDB*, n.d.)) have historically powered data handling in numerous applications (Habyarimana & Michailidou, 2021) (Jain et al., 2013) (*Well-Known Users Of SQLite*, n.d.) ("Customer Case Studies," n.d.). However, both approaches exhibit unique limitations that can impede rapid experimentation, large-scale research, and data-intensive development.

File-based solutions are popular for their simplicity and portability, often relying on straightforward ASCII/UTF encoding. This design choice, while human-readable, becomes highly inefficient for numerical data. For example, encoding an integer like 127 in ASCII demands three separate bytes (00110001 00110010 00110111), leading to significant overhead as data volumes expand. Such inflation in file size translates into slower input/output (I/O) operations and increased storage requirements, ultimately restricting scalability. Additionally, file-based formats typically lack built-in querying capabilities and indexing features, forcing developers to manage complex data relationships manually. These constraints limit the agility of workflows, especially as projects grow in complexity or require quick iteration cycles.

Conversely, traditional database management systems offer robust encoding, indexing, and querying capabilities out of the box. Relational databases, for instance, enforce structured schemas that ensure data integrity but introduce complexities when the data model evolves over time. Non-relational databases, such as document-oriented or key-value stores, are more flexible but risk data inconsistency and can become cumbersome to optimize for performance(Singh,

). Many of these solutions require dedicated servers or intricate configurations, increasing overhead for lightweight experimentation. Moreover, the underlying architectures—whether row-based or reliant on specialized storage engines—can exhibit performance bottlenecks when handling unstructured or semi-structured data at scale.

ParquetDB is intended to be a "medium ware" solution that sits between these two paradigms. Built in Python and leveraging the Parquet columnar format, ParquetDB combines the efficiency of column-based data storage with the accessibility of file-based approaches. Columnar storage significantly enhances compression and read performance by grouping similar data together, reducing the cost of both serialization and deserialization. Parquet also preserves rich metadata at the table and column levels, enabling advanced features such as schema enforcement and automated indexing without the overhead typically associated with full-fledged database engines. By integrating these capabilities into a lightweight, serverless architecture, ParquetDB mitigates the complexities of evolving data models, which can pose challenges in rigid relational systems or loosely structured NoSQL stores. Researchers and developers benefit from rapid data access, efficient queries, and a streamlined approach to data serialization. ParquetDB's compatibility with Python's scientific ecosystem also facilitates data analysis and machine learning workflows, letting users iterate quickly on large datasets without heavy infrastructure.

Ultimately, ParquetDB delivers a balanced, scalable, and portable platform for data storage and retrieval, addressing the well-known shortcomings of both plain-text file formats and conventional databases. Its columnar design, built-in metadata, and efficient encoding strategies enable fast, cost-effective reads and writes while maintaining a flexible schema structure. By offering this middle ground, ParquetDB helps unify the best of both worlds in data storage, setting a new standard for flexible, high-performance data management that meets the needs of modern analytics and research.

## Features

| Features and Benefits | Description |
| --- | --- |
| Simple Interface | Easy-to-use methods for creating, reading, updating, deleting, and transforming data. |
| High Performance | Utilizes Apache Parquet and PyArrow for efficient data storage and retrieval. |
| Complex Data Types | Handles nested and complex data types (Ex. ndarrays, lists, dictionaries, python functions and classes etc.). |
| Portability | File-based storage, allows for easy transfer. (Any framework that can read a directory of parquet files can read ParquetDB databases) |
| Schema | Contains a schema that describes the data, ensuring consistency. |
| Schema Evolution | Supports adding new fields and updating schemas. |
| Predicate Pushdown | Optimizes queries by reading only relevant data blocks. |
| Column Pushdown | Selects columns to read into memory. |
| Efficient Encoding | Choice of field-level encoding. |
| Efficient Compression | Choice of field-level compression. |
| Metadata Support | Table and field-level metadata support. |
| Batching Support | Files are grouped to facilitate batching. |

## Benchmarks

In this section, we outline the methodology used to evaluate the performance of ParquetDB in comparison to SQLite and MongoDB. Synthetic datasets consisting of 100 integer columns with varying record counts were used to simulate different load levels. The integer values ranged from 0 to one million, with dataset sizes varying from one row to one million rows.

integers were chosen as they are a fundamental data type, allowing us to establish baseline
performance metrics with minimal additional computational complexity. Integers were chosen
as the benchmark data type since they are primitive and other data types can be extrapolated
from their performance. Since integers are more lightweight compared to other data types,
such as strings or nested structures, they provide an initial estimation of database performance
without the variability introduced by more complex data representations. This approach allows
us to evaluate the core efficiency of the database under ideal conditions, providing a best-case
scenario. Subsequently, more complex data types may be expected to degrade performance
due to increased memory usage and processing demands, and these results provide a reference
point for such comparisons



**Figure 1:** Benchmark Create and Read Times for Different Databases. Create time is plotted on the left
y-axis, read time on the right y-axis, and the number of rows on the x-axis. A log plot is shown in the
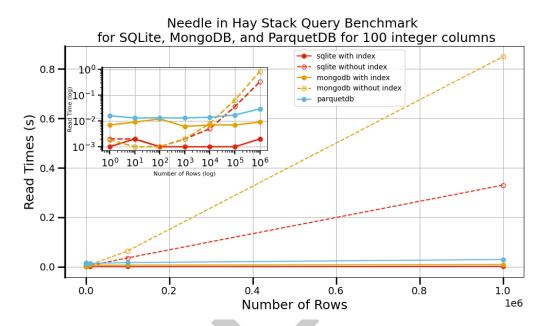inset.

**Figure 2:** Needle-in-a-Haystack Benchmark Results. Time is on the y-axis, number of rows on the x-axis. The log plot is shown in the inset. SQLite and MongoDB are compared with and without indexing.

## Installation

For installation, please use pip:

```
pip install parquetdb
```

For more details, including advanced features and contributions, please visit the [GitHub repository](). The repository contains additional examples, API documentation, and guidelines for contributing to the project.

## Acknowledgements

# References

*About SQLite*. (n.d.). https://www.sqlite.org/about.html.

Customer Case Studies. (n.d.). In *MongoDB*. https://www.mongodb.com/solutions/customer-case-studies.

Habyarimana, E., & Michailidou, S. (2021). Genomics Data. In C. Sodergard, T. Mildorf, E. Habyarimana, A. J. Berre, J. A. Fernandes, & C. Zinke-Wehlmann (Eds.), *Big Data in Bioeconomy: Results from the European DataBio Project* (pp. 69–76). Springer International Publishing. ISBN: 978-3-030-71069-9

Jain, A., Ong, S. P., Hautier, G., Chen, W., Richards, W. D., Dacek, S., Cholia, S., Gunter, D., Skinner, D., Ceder, G., & Persson, K. A. (2013). Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Materials*, *1*(1), 011002. https://doi.org/10.1063/1.4812323

*MongoDB: The Developer Data Platform | MongoDB*. (n.d.). https://www.mongodb.com/.

Singh, A. P. (2024). *15 Types of Databases and When to Use Them*. https://blog.algomaster.io/p/15-types-of-databases.

*Well-Known Users Of SQLite*. (n.d.). https://www.sqlite.org/famous.html.