

# ParquetDB: A Lightweight Python Parquet-Based Database

Logan L. Lang<sup>1</sup>, Eduardo R. Hernandez<sup>2</sup>, Kamal Choudhary<sup>3</sup>, and Aldo H. Romero<sup>1</sup>

<sup>1</sup> Department of Physics, West Virginia University, Morgantown, United States <sup>2</sup> Instituto de Ciencia de Materiales de Madrid, Madrid, Spain <sup>3</sup> National Institute of Standards and Technology, Gaithersburg, United States

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

## Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

ParquetDB is a Python library serving as a “middleware” solution, bridging the gap between file-based storage and full database systems. A key driver for its development was the need to support iterative research workflows, requiring schema evolvability, the ability to manage complex and evolving nested data structures without predefined rigidity, and the ability to handle-table and field-level metadata. Additionally, its “classically serverless” nature was a crucial design point for deployment in environments such as HPC clusters with limited connectivity. Leveraging Apache Parquet (“Parquet,” n.d.; [Apache Software Foundation, n.d.](#)), it combines file storage portability with advanced querying capabilities, enabling efficient compression and read performance without dedicated server overhead. ParquetDB addresses limitations in both traditional approaches by seamlessly handling complex data types (arrays, nested structures, Python objects), simplifying data interaction compared to direct file manipulation or manual serialization. Performance benchmarks show competitive read/write speeds and effective query performance via predicate pushdown, demonstrating its utility for managing medium-to-large datasets where database complexity is unwarranted but basic file I/O is insufficient.

## Statement of need

The demand for efficient, scalable, and adaptable data storage solutions is critical across research domains. Traditional file formats (e.g., CSV, JSON, TXT) offer simplicity but suffer from inefficiencies, particularly with numerical data due to ASCII/UTF encoding overhead, leading to larger files and slower I/O. While binary formats like HDF5 ([HDF5, n.d.](#)) improve efficiency for large numerical datasets, they function primarily as structured file containers, lacking the rich querying APIs and transactional integrity features common in databases. These file-based approaches often require manual data relationship management and lack built-in indexing, hindering agility as projects scale or require rapid iteration.

Database systems like SQLite ([Allen & Owens, 2010](#)) or MongoDB ([Guo, 2017](#)) provide robust encoding, indexing, and querying. Relational databases ensure integrity via structured schemas but can be rigid when data models evolve ([Pascal, 2000](#)). NoSQL options offer flexibility but may introduce consistency challenges or require complex optimization ([Pivert, 2018](#)). Furthermore, many databases involve server configurations or lack transparent file-based portability, adding overhead unsuitable for lightweight experimentation or simpler deployment scenarios. While SQLite is serverless and ubiquitous, its row-based nature can be less performant for analytical queries scanning wide datasets compared to columnar formats, and managing complex nested data can be cumbersome.

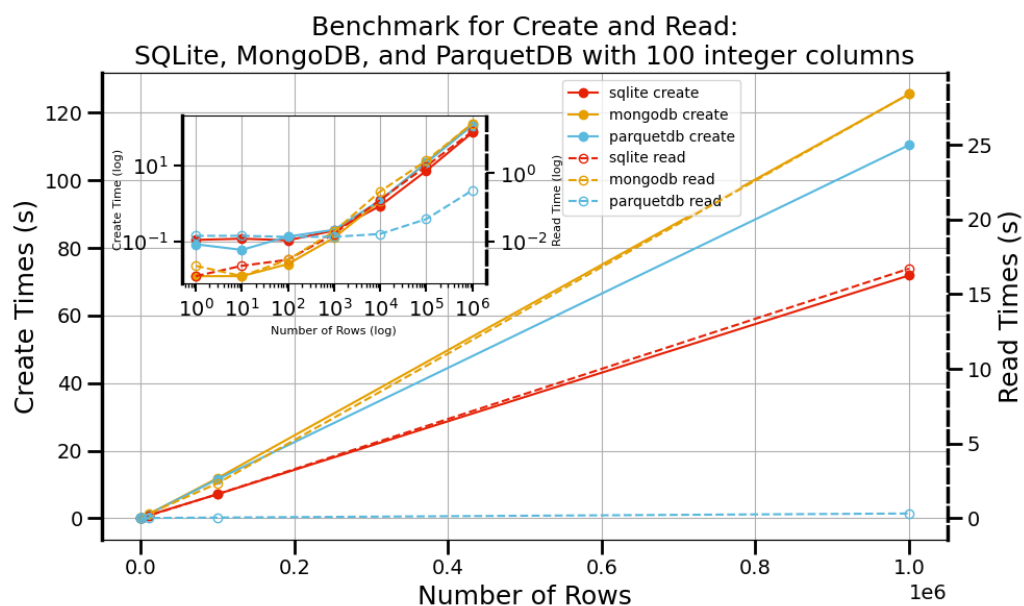
42 Directly using libraries like Apache Arrow (PyArrow) to work with Parquet files offers access  
43 to columnar efficiency and querying primitives like predicate pushdown. However, this still  
44 requires developers to build abstractions for database-like operations (CRUD), manage schema  
45 consistency across multiple files, handle serialization of complex Python objects, and orchestrate  
46 data updates or deletions manually.

47 While powerful dataframe manipulation libraries like Pandas ([Pandas, n.d.](#)), Dask ([Dask, n.d.](#)),  
48 and Polars ([Polars, n.d.](#)), or embedded analytical databases such as DuckDB (["DuckDB,"](#)  
49 [n.d.](#)), are invaluable for many tasks, they may not holistically address the specific needs that  
50 motivated ParquetDB. For researchers dealing with evolving, complexly nested scientific data,  
51 ParquetDB offers a more streamlined approach to schema evolvability and native Python object  
52 persistence directly within a serverless Parquet-based ecosystem. This focus distinguishes it  
53 from tools that might require more manual setup for schema management across multiple files,  
54 or lack the same emphasis on integrated metadata handling and a 'classically serverless' model  
55 for environments like HPC clusters.

56 ParquetDB addresses this gap, providing a "middleware" layer built upon Python and the  
57 Parquet format. It offers a familiar database-like interface (CRUD operations) while leveraging  
58 columnar storage for compression and read performance benefits. Crucially, ParquetDB adds  
59 value beyond direct Parquet file manipulation by automating schema management (including  
60 evolution), simplifying the storage/retrieval of complex Python objects, and providing a unified  
61 API to manage collections of Parquet files as a single logical datastore. It supports predicate  
62 and column pushdown for optimization within a lightweight, serverless architecture, offering a  
63 pragmatic balance for scenarios demanding more than basic files but less than a full database  
64 system, particularly where schema flexibility and ease of use are paramount. For a comprehensive  
65 feature list, visit our documentation (<https://parquetdb.readthedocs.io/en/latest/>).

## 66 Benchmarks

67 We evaluated ParquetDB's performance against SQLite and MongoDB using synthetic datasets  
68 (100 integer columns, varying record counts). Our first experiment compared write and read  
69 performance. ParquetDB's creation times are competitive, performing second best behind  
70 SQLite as dataset size increases. For bulk read operations, ParquetDB initially lags slightly but  
71 significantly outperforms both competitors on larger datasets (beyond several hundred/thousand  
72 rows), benefiting from Parquet's columnar efficiency (see Figure 1).



**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.

73 A “needle-in-a-haystack” benchmark assessed specific record retrieval. While lacking traditional  
 74 B-tree indexes, ParquetDB uses predicate pushdown leveraging Parquet’s field-level statistics  
 75 for efficient filtering without full scans. It is important to note that performance advantages  
 76 depend on the workload; for instance, complex analytical queries involving aggregations or  
 77 returning small, highly filtered results might favor the mature query engine and indexing of  
 78 systems like SQLite. ParquetDB excels when querying or returning substantial portions of wide  
 79 datasets. Detailed benchmarks are in our extended paper (Lang et al., 2025).

## 80 Installation

81 For installation, please use pip:

```
82 pip install parquetdb
```

83 For more details, please visit the GitHub repository: (<https://github.com/lllangWV/ParquetDB>).  
 84 The repository contains additional examples, API documentation, and guidelines for contributing to the project.

## 85 Acknowledgements

86 We thank the Pittsburgh Supercomputer Center (Bridges2) and San Diego Supercomputer  
 87 Center (Expanse) through allocation DMR140031 from the Advanced Cyberinfrastructure  
 88 Coordination Ecosystem: Services & Support (ACCESS) program, which is supported by  
 89 National Science Foundation grants #2138259, #2138286, #2138307, #2137603, and  
 90 #2138296. We gratefully acknowledge the computational resources provided by the WVU  
 91 Research Computing Dolly Sods HPC cluster, partially funded by NSF OAC-2117575.  
 92 Additionally, we recognize the support from the West Virginia Higher Education Policy  
 93 Commission through the Research Challenge Grant Program 2022 (Award RCG 23-007), as  
 94 well as NASA EPSCoR (Award 80NSSC22M0173), for their contributions to this work. The  
 95 work of E.R.H. is supported by MCIN/AEI/ 10.13039/501100011033/FEDER, UE through

96 projects PID2022-139776NB-C66. K.C. thanks funding from the CHIPS Metrology Program,  
97 part of CHIPS for America, National Institute of Standards and Technology, U.S. Department  
98 of Commerce. Certain commercial equipment, instruments, software, or materials are identified  
99 in this paper in order to specify the experimental procedure adequately. Such identifications  
100 are not intended to imply recommendation or endorsement by NIST, nor are they intended  
101 to imply that the materials or equipment identified are necessarily the best available for the  
102 purpose.

## 103 References

- 104 Allen, G., & Owens, M. (2010). *The Definitive Guide to SQLite*. Apress. <https://doi.org/10.1007/978-1-4302-3226-1>  
105
- 106 An in-process SQL OLAP database management system. (n.d.). In *DuckDB*.  
107 <https://duckdb.org/>.
- 108 Dask | Scale the Python tools you love. (n.d.). <https://www.dask.org/>.
- 109 Guo, R. (2017). MongoDB's JavaScript fuzzer. *Commun. ACM*, 60(5), 43–47. <https://doi.org/10.1145/3052937>  
110
- 111 *HDF5 for Python — H5py 3.13.0 documentation*. (n.d.). <https://docs.h5py.org/en/stable/index.html>.
- 112 Lang, L., Hernandez, E., Choudhary, K., & Romero, A. H. (2025). *ParquetDB: A Lightweight*  
113 *Python Parquet-Based Database* (No. arXiv:2502.05311). arXiv. <https://doi.org/10.48550/arXiv.2502.05311>  
114
- 115 *Pandas - Python Data Analysis Library*. (n.d.). <https://pandas.pydata.org/>.
- 116 Parquet. (n.d.). In *Apache Parquet*. <https://parquet.apache.org/>.
- 117 Pascal, F. (2000). *Practical Issues in Database Management: A Reference for the Thinking*  
118 *Practitioner* (1st edition). Addison-Wesley Professional. ISBN: 978-0-201-48555-4
- 119 Pivert, O. (Ed.). (2018). *NoSQL Data Models: Trends and Challenges* (1st edition).  
120 Wiley-ISTE. ISBN: 978-1-78630-364-6
- 121 *Polars*. (n.d.). <https://www.pola.rs/>.
- 122 *Welcome to the Apache Software Foundation*. (n.d.). <https://www.apache.org/>.