

# DaskDB: Revolutionizing Cloud-Based Data Analytics with Distributed Serverless Spatial and SQL Query Processing

Ronnit Peter, Nikhil Chhadikar, Prasanna Kumar Batlanki

## Introduction

The digital era's exponential data growth poses both opportunities and challenges for analytics, especially in cloud environments. "DaskDB"[1][2] emerges as a groundbreaking solution, utilizing a distributed, serverless architecture for large-scale data analytics. This project proposal outlines DaskDB's development and its potential to enhance cloud-based analytics through seamless integration with AWS NoSQL databases and other platforms.

## Objectives

- **Making DaskDB Serverless:** Integration DaskDB with AWS services (S3, CloudWatch) to ensure compatibility within the AWS ecosystem. To evaluate the performance improvements in data processing and retrieval times. Thus, ensuring efficient scaling and minimal data movement and reduction in significant query processing times (subjected to resources available on AWS) due to the nature of distributed data workload possible through DaskDB.

## Methodology

The project employs a detailed experimental approach to validate DaskDB's performance, scalability, and integration with cloud-based data storage:

1. **Scalability Tests:** Evaluate DaskDB's AWS EC2 performance under increasing workloads.
2. **Performance and Integration Tests:** Measure efficiency and seamless AWS ecosystem operation
3. **Front-End Implementation:** Develop a web-based user interface for querying the dataset, utilizing AWS S3 for storage. Access to the front-end page will be restricted as needed via IAM policies. Additional AWS resources such as Lambda S3 Buckets and IAM policies will be deployed to facilitate data retrieval and interaction with the DaskDB backend.

## Literature Review

The recent launch of spatial data processing capabilities in cloud platforms, notably Snowflake's introduction for paid members in January 2024, underscores a significant advancement in cloud-based data analytics. Unlike earlier technologies, which also support scalable spatial querying[5], Snowflake's move illustrates the novel and proprietary nature of these enhancements, protected by a paywall.

DaskDB distinguishes itself by proposing a serverless, distributed architecture that integrates both spatial[3] and SQL query processing[4]. This innovation seeks to make advanced data analytics more accessible, bypassing the constraints of proprietary systems. By doing so, DaskDB positions itself as a groundbreaking solution in the landscape of cloud information management systems, offering a scalable, open alternative to the paywall-protected advancements in spatial data processing[6].

## Conclusion

DaskDB is a pioneering cloud analytics solution, addressing scalability, efficiency, and AWS integration challenges. Its innovative architecture promises to transform analytics, providing a unified, scalable solution for evolving disciplines.

## References

- [1] Suprio Ray , Suvam Kumar Das, Ronnit Peter Scalable Spatial Analytics and In Situ Query Processing in DaskDB
- [2] Md. Mahbub Alam, Suprio Ray, and Virendra C. Bhavsar. 2018. A Performance Study of Big Spatial Data Systems. In Proceedings of ACM SIGSPATIAL International Workshop on Analytics for Big Geospatial Data, BigSpatial@SIGSPATIAL. 1–9.
- [3] Dask-GeoPandas [n.d.]. Dask-GeoPandas. <https://daskgeopandas.readthedocs.io/en/stable/>.
- [4] GeoNB [n.d.]. NB Dataset. <http://www.snb.ca/geonb1/e/DC/catalogue-E.asp>
- [5] GeoPandas [n.d.]. GeoPandas. <https://geopandas.org/en/stable/>.
- [6] Suprio Ray, Bogdan Simion, and Angela Demke Brown. 2011. Jackpine: A benchmark to evaluate spatial database performance. In 2011 IEEE 27th International Conference on Data Engineering. 1139–1150.