

Cloud Giovanni: Reining in Costs and Improving Performance with Analytical Data Stores using Scalable Serverless Architecture

AGU100 ADVANCING EARTH AND SPACE SCIENCE

IN13B-0708

NASA/Goddard EARTH SCIENCES DATA and INFORMATION SERVICES CENTER (GES DISC)

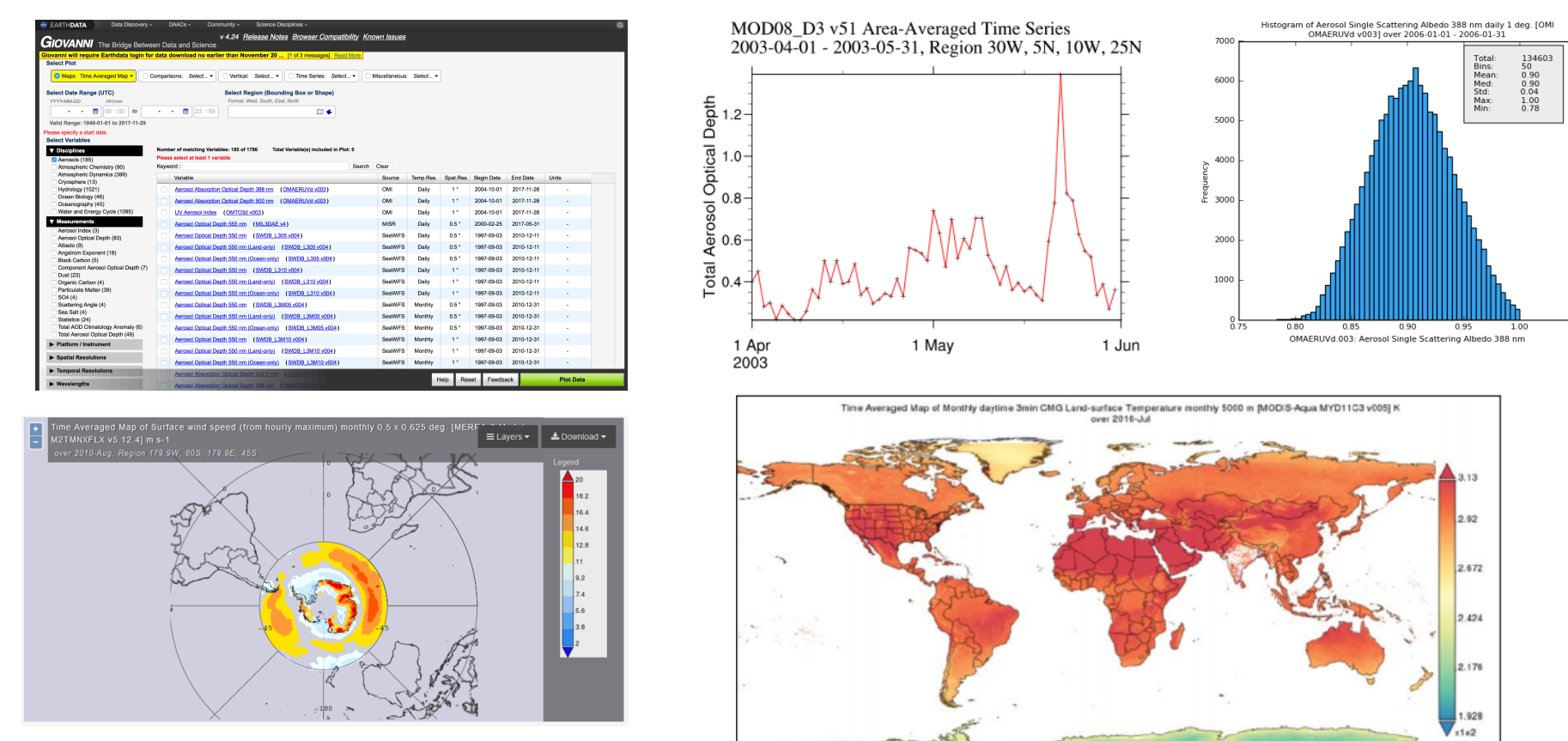
Hailiang Zhang^{1,2}, Mahabal Hegde^{1,2}, Christine Smit^{1,3}, Maksym Petrenko^{1,2} and Long Pham¹

¹ NASA Goddard Space Flight Center, ² ADNET Systems Inc., ³ Telophas Corp

Background

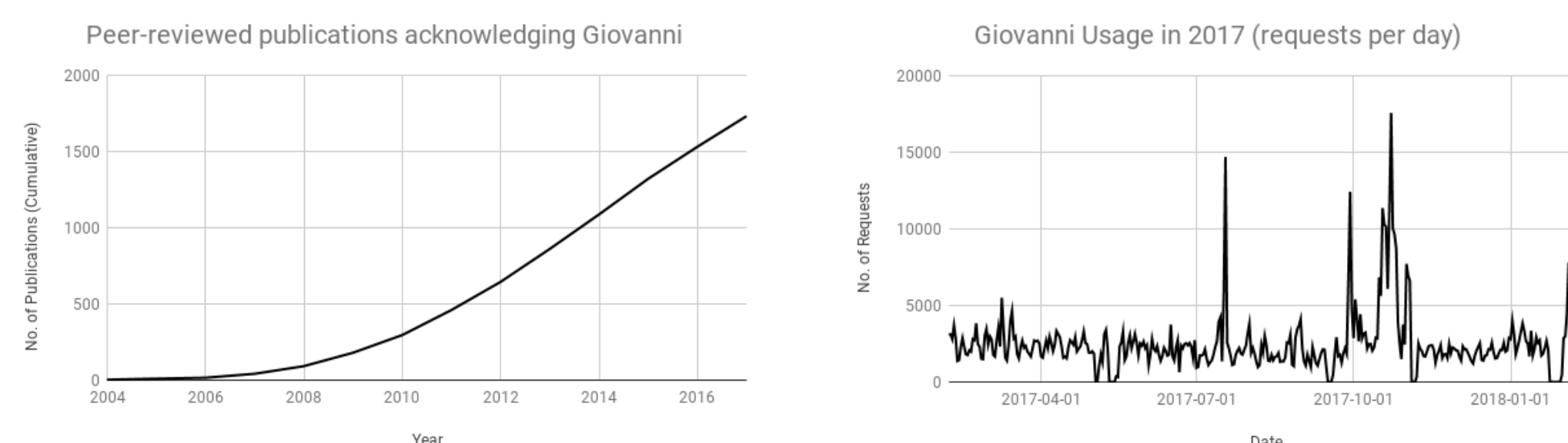
Giovanni (<https://giovanni.gsfc.nasa.gov>) is the Geospatial Interactive Online Visualization ANd aNalysis Infrastructure developed at NASA's GES DISC which provides a simple and intuitive way to visualize, analyze, and access vast amounts of Earth science data:

- Twenty-two (22) analysis and visualization services at the click of a button
- Access to over 1477 data variables
- Persistent URLs for sharing data and visualizations



Challenges

- Emphasis on feature set rather than reliability and performance, the two key pillars of a well architected framework
- Unable to meet spikes in demand during training, and "seasonal" events such as conferences and end of academic terms
- Increased demand on resources due to higher resolution data and user demand for data statistics



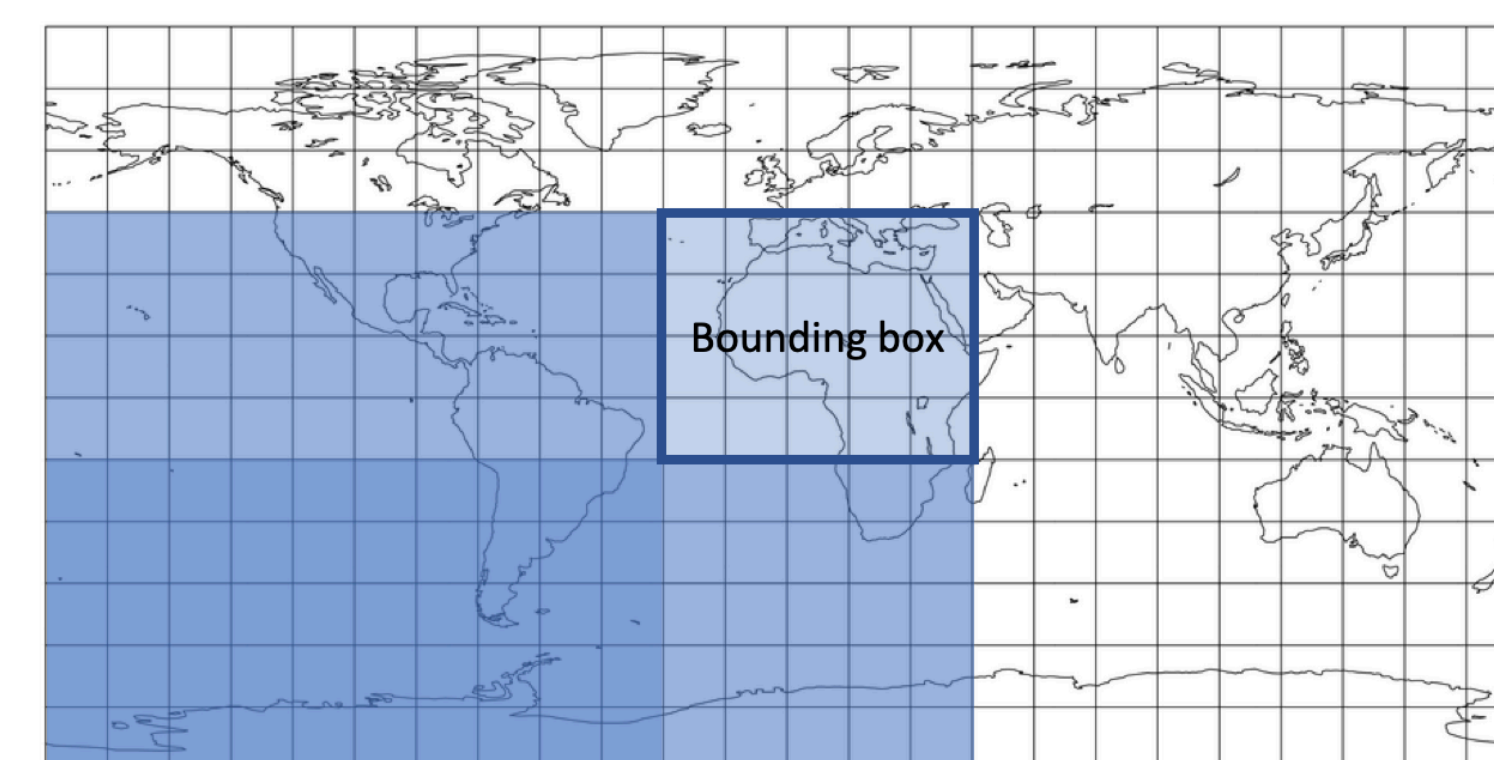
Algorithm

Multi-dimensional global averaging $\iiint \text{variable}(\text{time}, \text{lat}, \text{lon} \dots) \cdot \text{weight}(\text{time}, \text{lat}, \text{lon} \dots) dV$

2-dimensional averaging (time series)

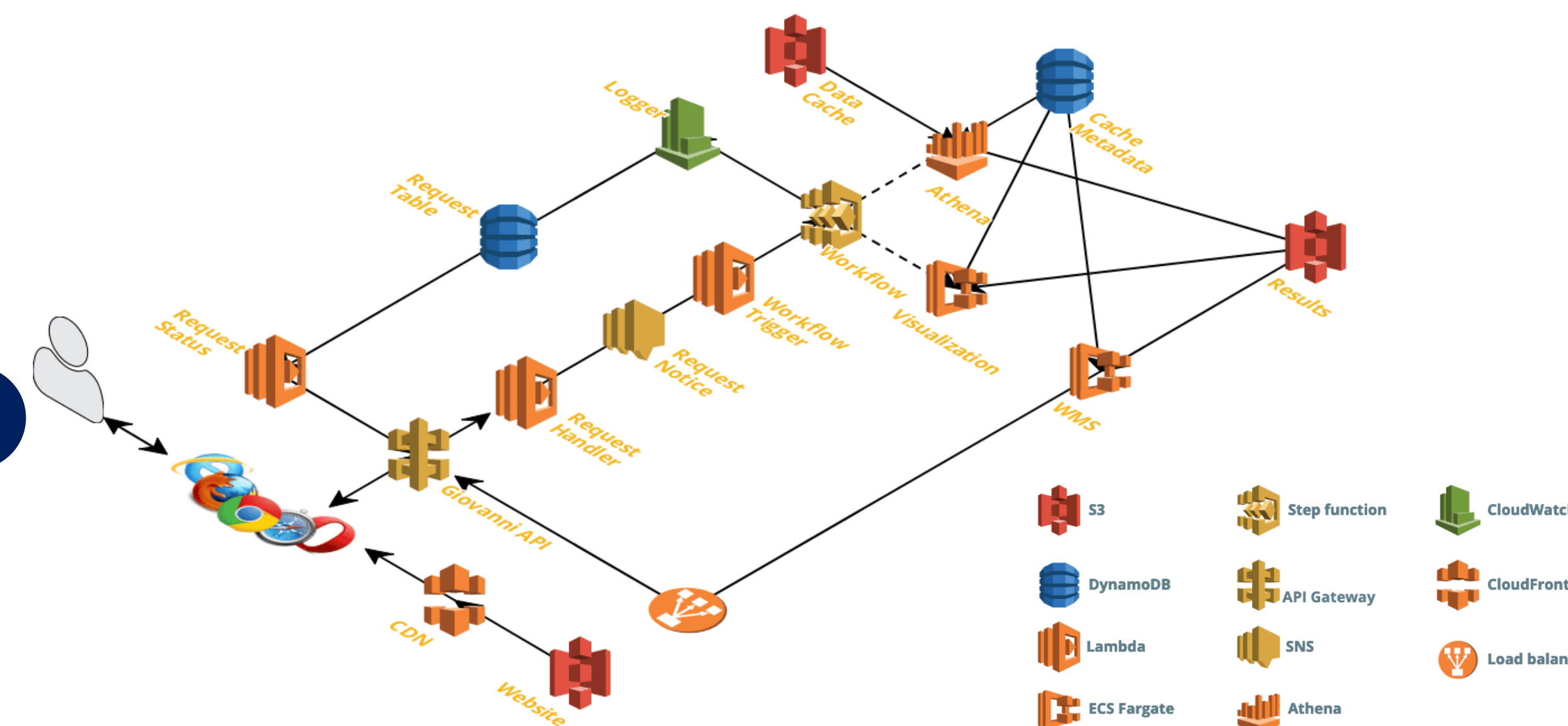
$$\iint \text{weight}(\text{time}, \text{lat}, \text{lon} \dots) dV$$

$$\frac{\sum_{t=0}^{\text{lat}} \sum_{l=0}^{\text{lon}} y_{t,l} \cdot \text{cos}(\text{lat}) - \sum_{t=0}^{\text{lat}} \sum_{l=0}^{\text{lon}} y_{t,l} \cdot \text{cos}(\text{lat}) + \sum_{t=0}^{\text{lat}} \sum_{l=0}^{\text{lon}} y_{t,l} \cdot \text{cos}(\text{lat})}{\sum_{t=0}^{\text{lat}} \sum_{l=0}^{\text{lon}} \text{cos}(\text{lat}) - \sum_{t=0}^{\text{lat}} \sum_{l=0}^{\text{lon}} \text{cos}(\text{lat}) + \sum_{t=0}^{\text{lat}} \sum_{l=0}^{\text{lon}} \text{cos}(\text{lat})}$$



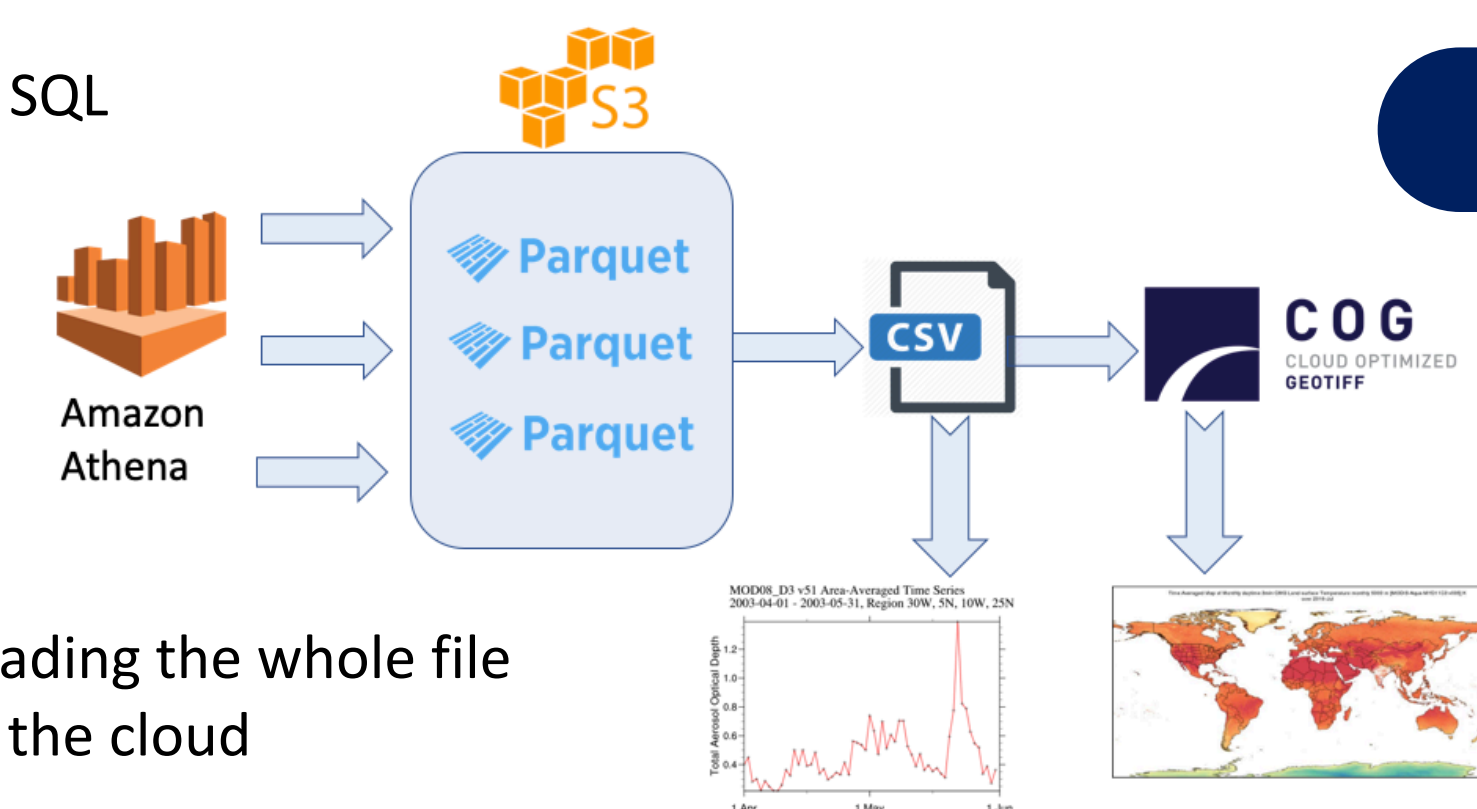
Architecture

AWS stack



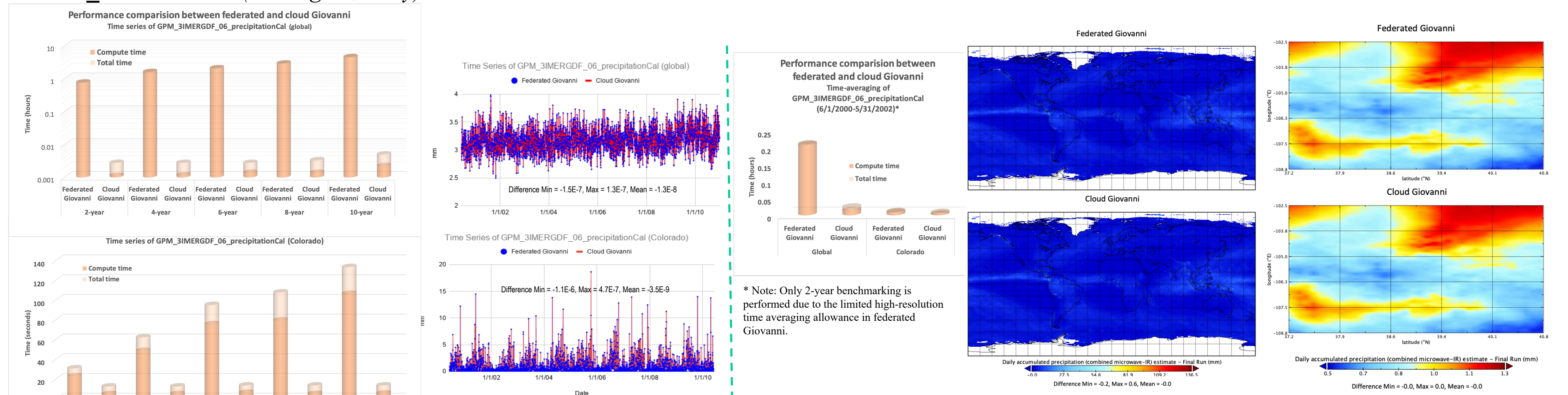
Data analysis frameworks

- AWS Athena
 - Efficient query service using standard SQL
 - High performance
 - Cost efficient: \$5/TB
- Parquet-based cloud object storage
 - Columnar data store
 - Highly compressible
- Cloud optimized Geotiff
 - Streaming GeoTIFF instead of downloading the whole file
 - Enabling more efficient workflows on the cloud

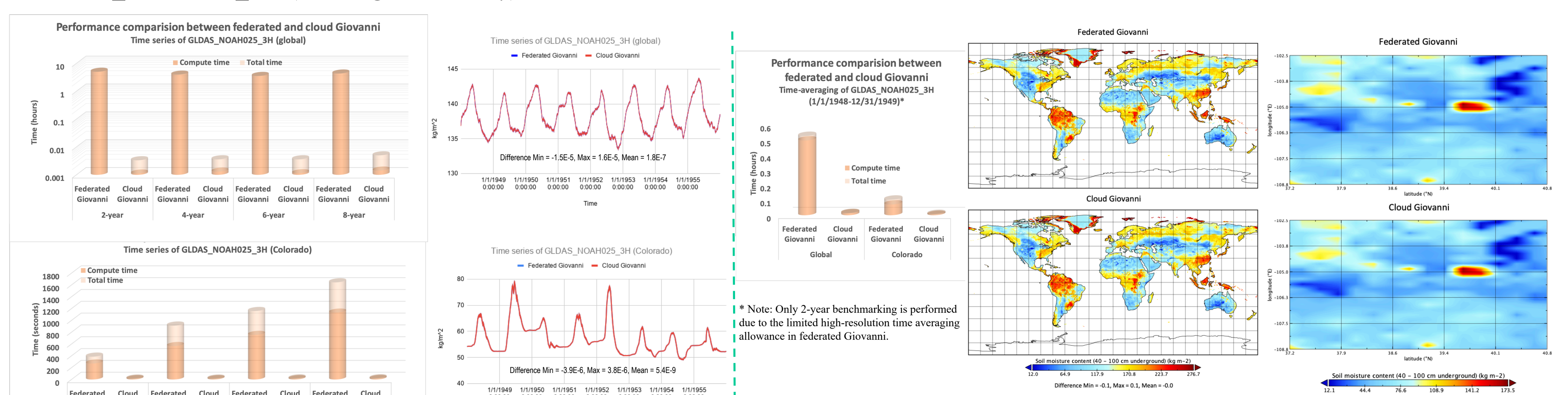


Results and Performance

GPM_3IMERGDF (0.1 degree, daily)



GLDAS_NOAH025_3H (0.25 degree, 3-hourly)



Cost Analysis

GPM_3IMERGDF (0.1 degree, daily)

Service type	Duration	Total cost (\$)/1000 requests	Global	Colorado
Area-averaged time series	2-year	0.4	0.58	0.58
	4-year	0.58	0.94	0.94
	6-year	0.75	1.28	1.28
Time-averaged map	2-year	12.12	1.63	1.63

AWS total cost for multi-dimensional accumulation-based cloud Giovanni

GLDAS_NOAH025_3H (0.25 degree, 3-hourly)

Service type	Duration	Total cost (\$)/1000 requests	Global	Colorado
Area-averaged time series	2-year	0.79	1.22	1.22
	4-year	1.33	2.18	2.18
	6-year	1.85	3.11	3.11
Time-averaged map	2-year	10.07	0.57	0.57

Conclusion

- Multi-dimensional accumulation-based cloud Giovanni significantly improved the performance of Giovanni services including area-averaged time series and time-averaged map. It can improve the performance by >1000 times for time series of variables with high spatial and temporal resolution.
- Multi-dimensional accumulation-based cloud Giovanni provides nearly identical results compared to federated Giovanni.
- Cloud Giovanni incurs reasonable AWS cost for highly intensive averaging services.
- Work is currently underway to further improve cloud Giovanni performance and reduce AWS costs.