

# Accessing Big Vector Data on the Cloud using Arrow Parquet



Johannes Heisig

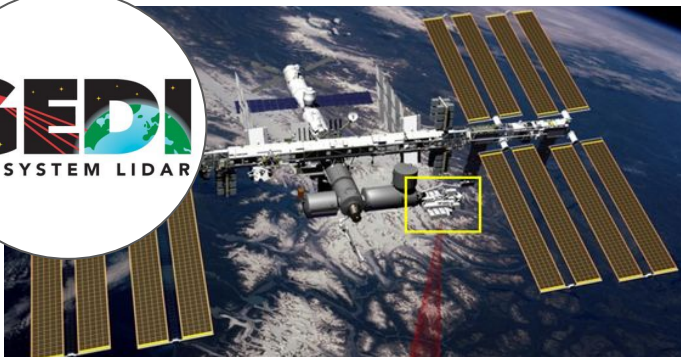


Yu-Feng Ho

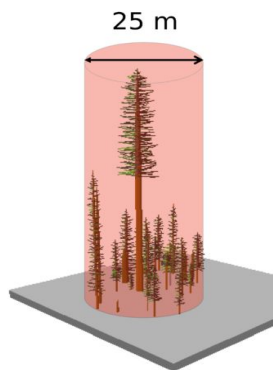
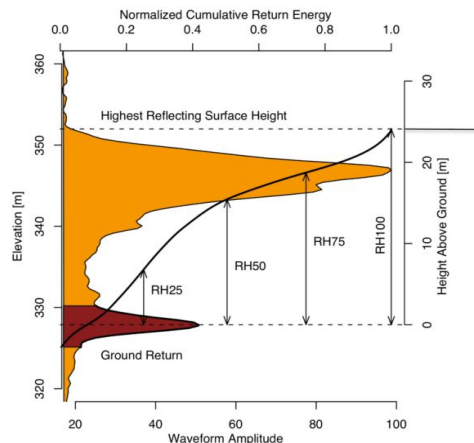




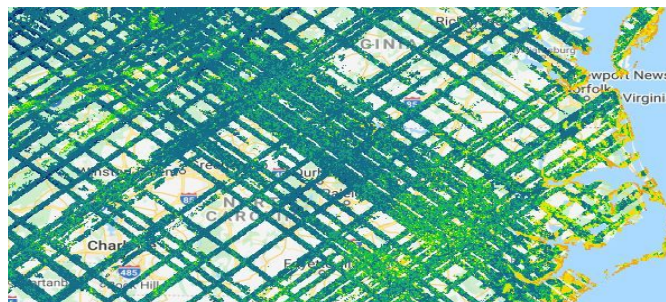
1. The data
2. Apache Arrow & Parquet
3. Cloud Setup
4. Notebook examples



- Full-waveform LiDAR system
- Data since 2019
- $N \sim 5.4$  billion points
- Level 2A+B data, 94 variables
- Point geometries (lon/lat)



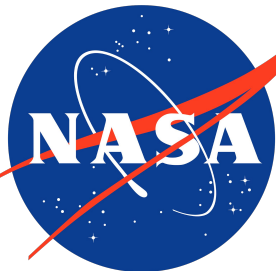
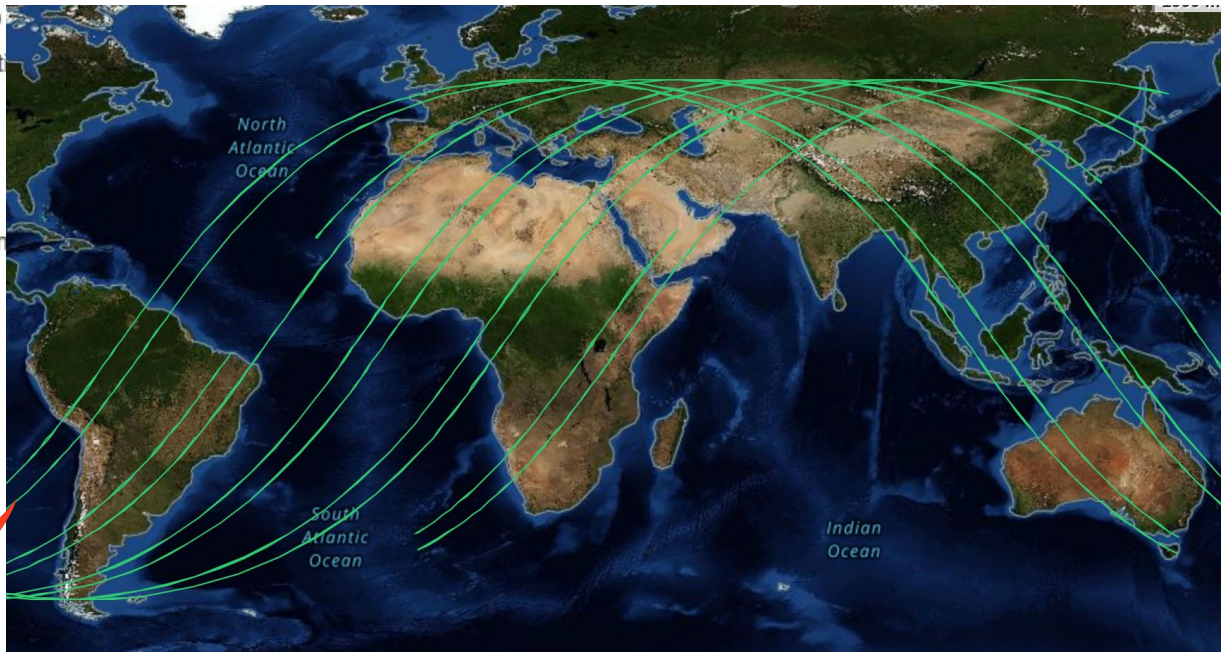
Sampling pattern





In this example of a Version 2 Level 01B product, the filename **GEDI01\_B\_2019110110221\_001997\_03\_T03335\_02\_005\_01\_V002.h5** indicates:

- **GEDI01\_B** = Product Short Name
- **2019110** = Julian Date of Acquisition in YYYYDDD
- **110221** = Hours, Minutes, and Seconds of Acquisition
- **001997** = O = Orbit, 01997 = Orbit Number
- **03** = Sub-Orbit Granule Number
- **T03335** = T = Track, 033335 = Track Number
- **02** = Positioning and Pointing Determination System
- **005** = PGEVersion = SDPS Release Number
- **01** = Granule Production Version
- **V002** = Version Number
- **.h5** = Data Format





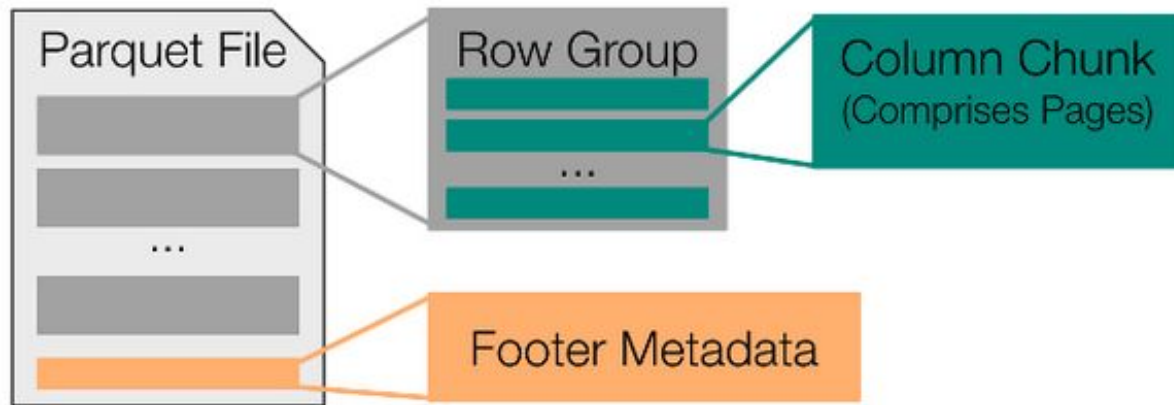


Figure 1. High-level visual representation of the Parquet file format

Dataset

|--- Partition  
|--- File  
|--- Row group  
|--- Column chunk  
|--- Page

Source:

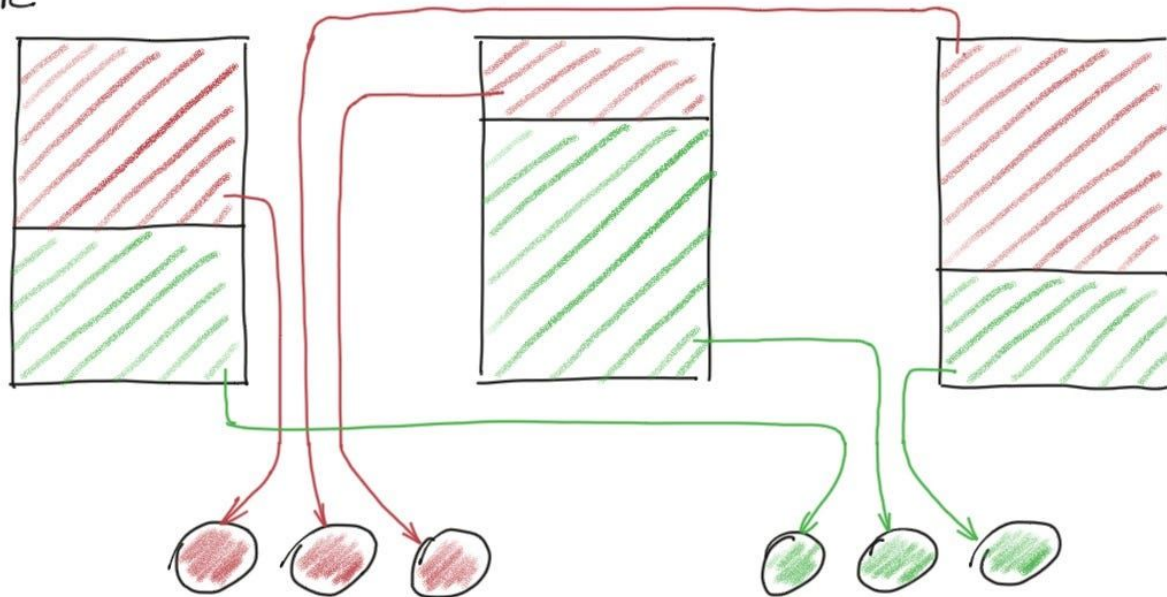
(1) [Optimizing Access to Parquet Data with fsspec](#)

```
https = S3FileSystem(
    endpoint_url='https://s3.eu-central-1.wasabisys.com',
    anon=True
)
pq_dataset = pq.ParquetDataset(
    path_or_paths='gedi-ard/level2/l2v002.gedi_20190418_20230316_go_epsg.4326_v20240622',
    filesystem=https)
```

CPU times: user 634 ms, sys: 0 ns, total: 634 ms  
Wall time: 1.35 s



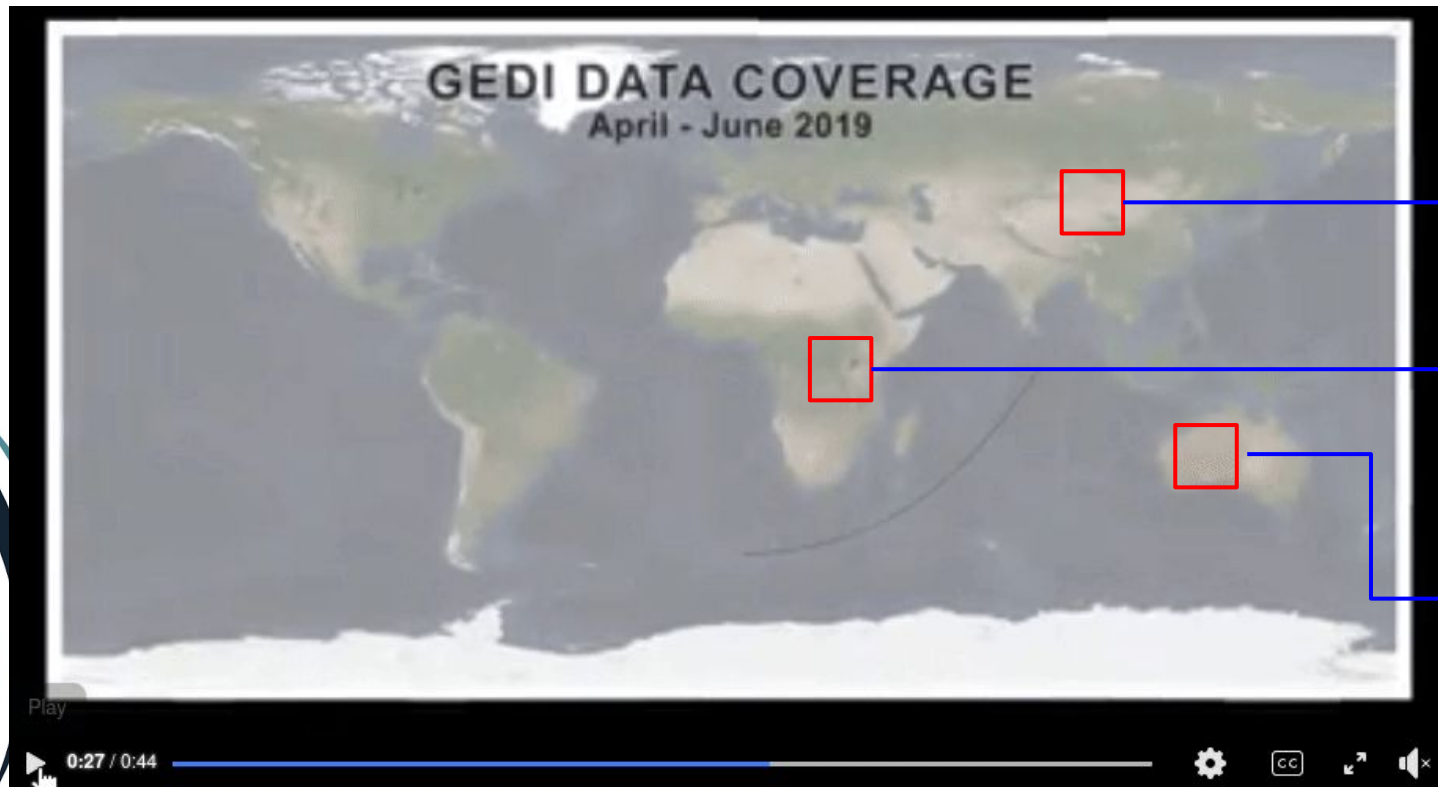
Spark



S3

Source:

[Spark partitioning: the fine print, Vladimir Prus, 2021](#)



***partition-1***

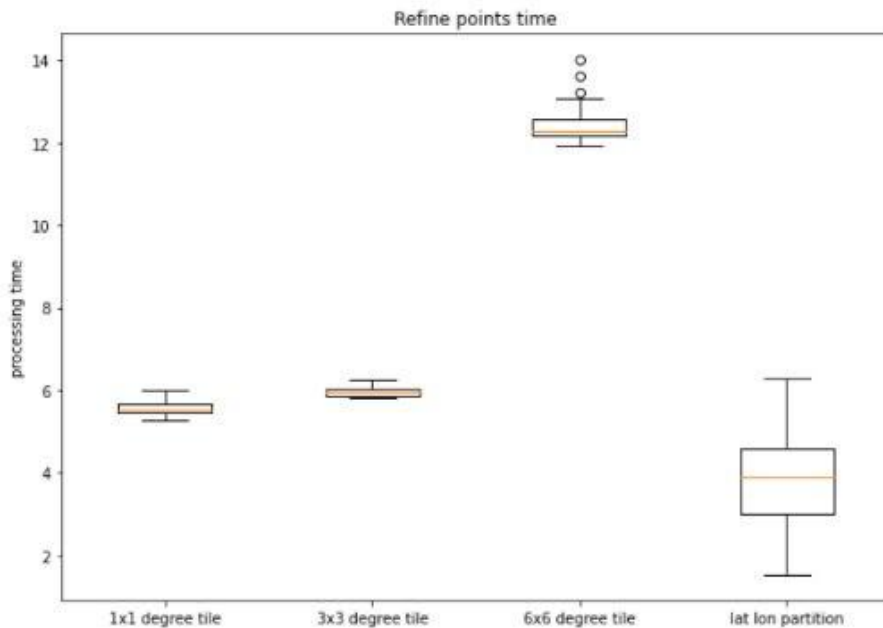
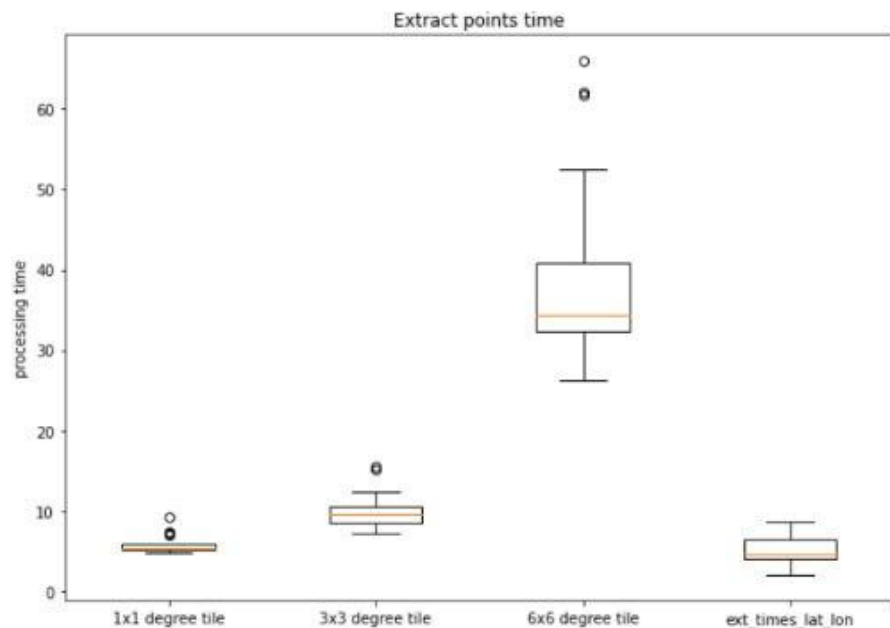
- gedi-o1a.parquet
- gedi-o2a.parquet

***partition-2***

- gedi-o1b.parquet
- gedi-o2b.parquet
- gedi-o3a.parquet
- gedi-o4a.parquet

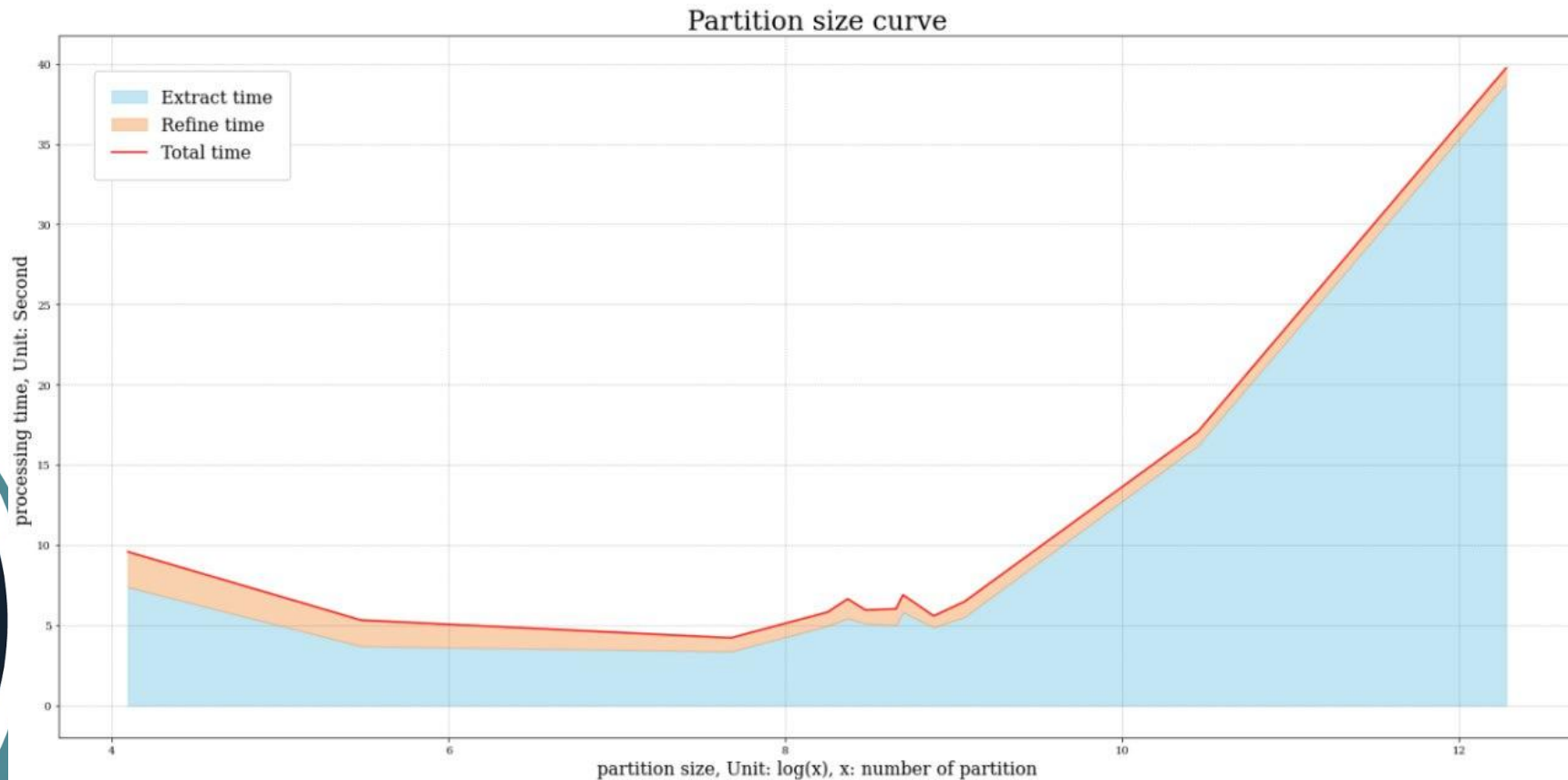
***partition-3***

- gedi-o3b.parquet
- gedi-o4b.parquet
- gedi-o5a.parquet



The multi depth structure(lat in 1 degree / lon in 1 degree)  
outperforms flattened structure (1x1 degree grid)







## Comparison of Python Parquet Engine in Parallelisation and Cloud Object Storage

Stats for all stories

6 min read · Dec 1, 2023



### Lifetime

Dec 1, 2023 – Today · Updated every 24 hours

536

Views

337

Reads

63%

Read ratio ⓘ

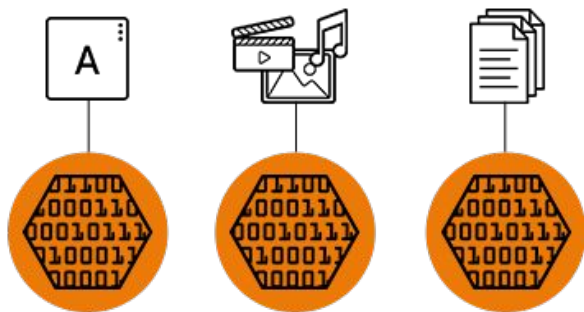
### pyarrow.dataset.write\_dataset

```
pyarrow.dataset.write_dataset(data, base_dir, *, basename_template=None,  
format=None, partitioning=None, partitioning_flavor=None, schema=None,  
filesystem=None, file_options=None, use_threads=True, max_partitions=None,  
max_open_files=None, max_rows_per_file=None, min_rows_per_group=None,  
max_rows_per_group=None, file_visitor=None, existing_data_behavior='error',  
create_dir=True)
```

[\[source\]](#)

Write a dataset to a given format and partitioning.





## Cloud Object Storage

```
url = 'http://s3.eu-central-1.wasabisys.com/gedi-ard/level2/l2v002.gedi_20190418_20230316_go_epsg.4326_v20240614/lon=15/lat=5/
year=2021/gedi_l2_0.parquet'
q = pl.scan_parquet(url)
if cols is not None:
    q = q.select(cols)
q = q.filter((pl.col('longitude') >= bbox[0]) & (pl.col('longitude') <= bbox[2]) &
             (pl.col('latitude') >= bbox[1]) & (pl.col('latitude') <= bbox[3]))
```

Source:

- (1) [What is object storage?](#) and [File storage, block storage, or object storage?](#)
- (2) [Polars — DataFrames for the new era](#)



This initiative is inclined to address global vector (point) data accessing and hosting in a cloud-native environment.

Readme

Apache-2.0 license

Activity

Custom properties

1 star

1 watching

0 forks

Report repository



yu-feng-ho Update README.md

f0112ca · 12 minutes ago

14 Commits

Python	python tutorial update	50 minutes ago
R	automatic package install	1 hour ago
.gitignore	rearrange repo structure	18 hours ago
LICENSE	Initial commit	5 months ago
Readme.md	Update README.md	12 minutes ago

README

Apache-2.0 license



## GlobalEarthPoint

GlobalEarthPoint is a **Open Source Data Service Library** that specializes in managing large datasets of geographical points, providing a portal to access cloud infrastructure, subset and retrieve data efficiently.

### Key Features

- Cloud Optimization:** Data is stored using [ArrowParquet](#) in the partitioning structure, and retrieved with Lazy evaluation.
- High Efficiency:** Designed to handle massive datasets by [Parquet](#) format, minimizing data size, latency and maximizing throughput.
- Easy Integration:** Compatible with popular data processing frameworks and geospatial tools, facilitating easy

Access the repository here:  
<https://github.com/Open-Earth-Monitor/GlobalEarthPoint>

[Create a new release](#)

### Packages

No packages published

[Publish your first package](#)

### Contributors 2



yu-feng-ho



johesig Johannes Heisig



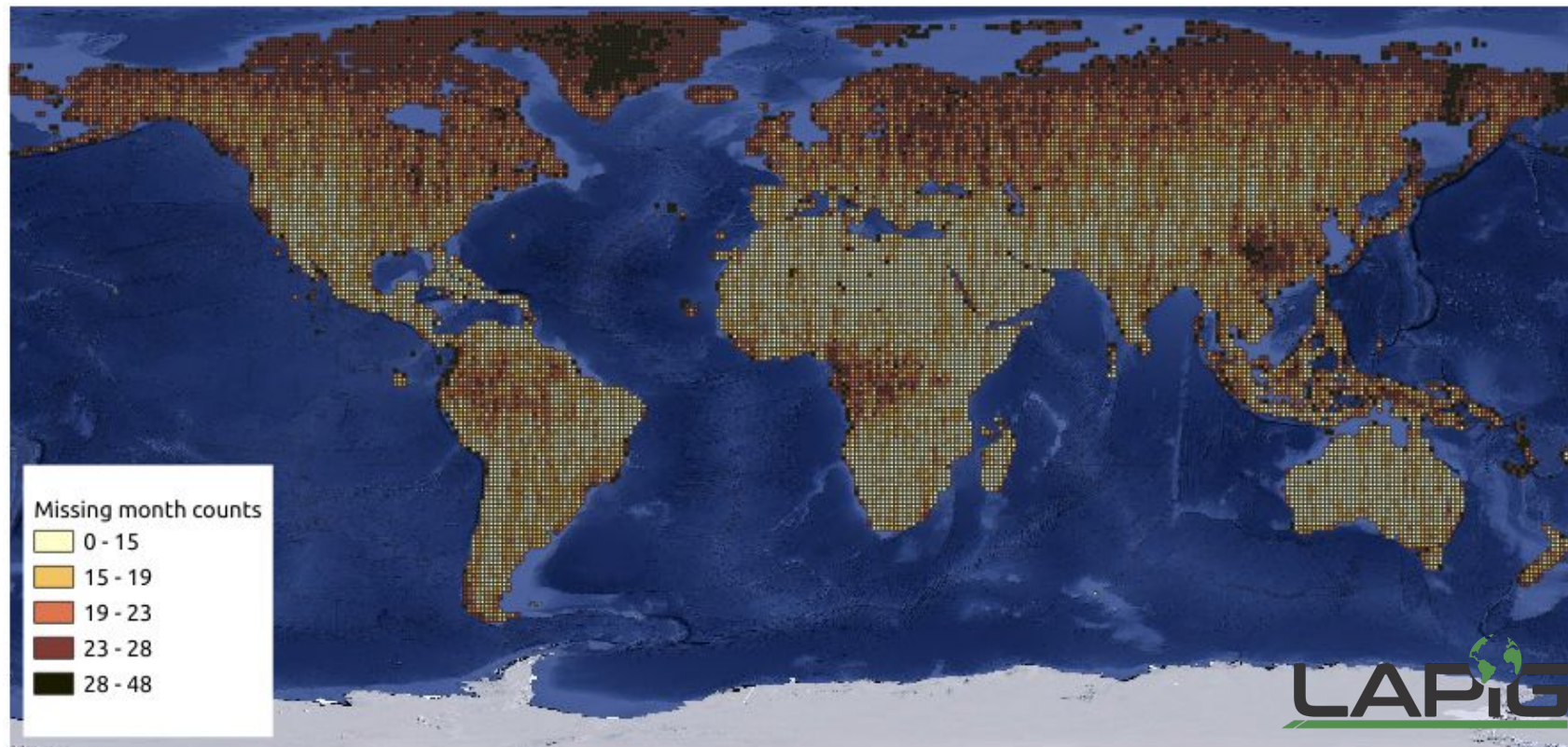
Institut für Geoinformatik  
Universität Münster



OpenGeoHUB  
Connect · Create · Share · Repeat











Icesat-2 ATL08v006 in 20m segment

24.8 billions

Generating extra attributes

~ 4TB

- **n\_ph\_20m**: Number of photons in 20m
- **med\_ht**: Median of vegetation photons relative height

Filter

By:

- **n\_ph\_20m**  $\geq 3$ : number of photons in 20m greater than or equal to 3  
\*separate segments with valid data
- **med\_ht**  $\leq 5$ : median less than or equal to 5m  
\*focus sample on areas dominated by low vegetation
- **night\_flag**  $= 1$ : data acquired at night  
\*reduction in RMSE of more than 1m

**Processing in 2 hours with a HPC 96 cores**

Icesat-2 ATL08v006 in 20m segment  
for Short Vegetation Mapping

4.9 billions

Production Process





# Questions & Feedback

