# Pangeo -TileDB

8/22/2019
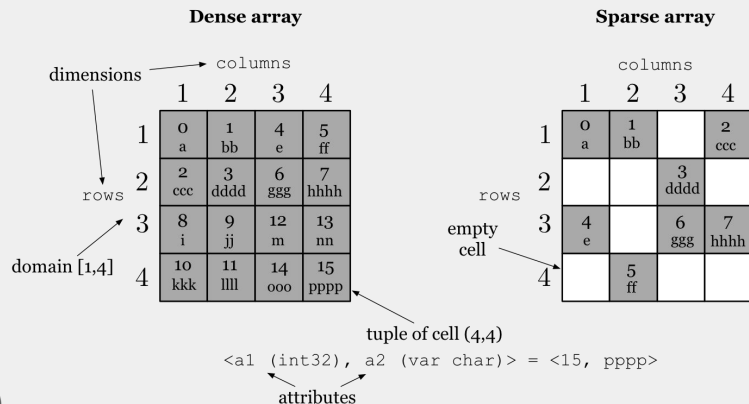
# About TileDB

- TileDB introduces a novel format and a powerful storage library for storing and accessing massive dense and sparse multi-dimensional array data
- The TileDB storage engine is optimized for the cloud
- TileDB, Inc. was founded in 2017 to further develop and maintain the TileDB project, which was a research project at Intel Labs and MIT
- TileDB has raised $4 million to date in Seed funding by Nexus Venture Partners and Intel Capital

[tile]DB

# TileDB Features



Dense array / Sparse array diagram showing dimensions, columns, rows, domain [1,4], empty cell, tuple of cell (4,4), `<a1 (int32), a2 (var char)> = <15, pppp>`, attributes

- Native support for sparse and dense arrays

- Rapid updates (and appends on dimensions)

- Multiple attributes for a single cell

- Parallel reads/writes for network storage, e.g. S3

- Integration with Python, Dask, R, Java, Spark, Presto, Go, PDAL and GDAL

- Well documented

- MIT license

[tile]DB

# TileDB for Geospatial Data Processing

- Can model complex raster data, with multiple attributes / pixel
- GDAL integration makes common geospatial ETL operations (re-projection, registration, translation, etc.) easy and efficient.
  - Full support for complex data types (I / Q) in the same array
- Supports efficient updates / appends / labeling.
- Full support for multi-variable datasets, e.g. HDF 4/5 , netCDF and NITF
- Attribute labelling of values in the array
- Time series in the same array (appending datasets)

[tile]DB

# Dask/TileDB

- GDAL code modified to support parallel writes to an existing TileDB array
- Dask
  - Parallel computations that scale up to thousands of nodes
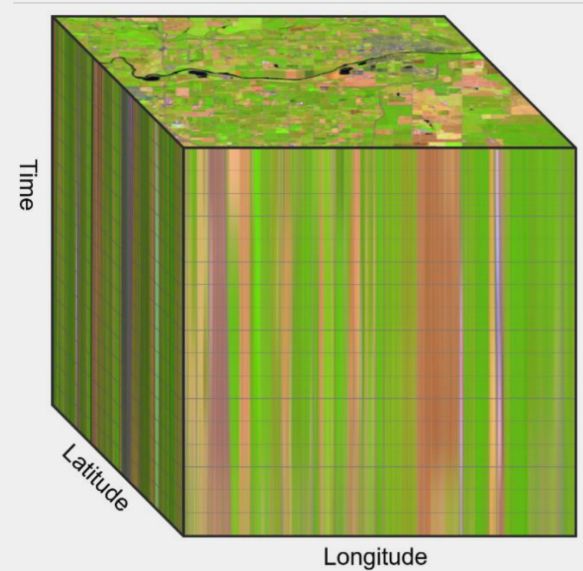  - https://docs.dask.org/en/latest/

TileDB enables fast parallel reads/writes with Dask

TileDB and the GDAL driver support multiple attributes per pixel value e.g. latitude, longitude and complex value

[tile]DB

# Applicability to ARDs

Analysis Ready Datasets

- TileDB Differentiators
  - Multi-dimensional support
    - Multiple attributes per pixel
    - Dense arrays
    - Sparse arrays
  - Append
  - Update
  - Parallel processing/tooling around these areas
  - No additional metadata indexing required



Achieving the Full Vision of Earth
Observation Data Cubes
https://www.mdpi.com/2306-5729/4/3/94/htm

[tile]DB

# COGs / TileDB

- COGs can be a data management headache
- COGs are limited by the options available in the TIFF spec
  - COGDumper
  - Libtiff
  - GDAL

TileDB provides an open source SDK in multiple languages, Python, Java, R, Spark in addition to a novel format that matches the needs of an ARD

[tile]DB

# TileDB / ARDs

Current state of play

- GDAL
  - Multi-dimensional support and translation from HDF / netCDF
  - Parallel appends / updates
  - Support for complex data types
- PDAL
  - Sparse array support with parallel appends / updates suitable near real time collects

[tile]DB