# SCALING ANALYTICS WITH PANGEO: XARRAY + DASK + ZARR

RYAN ABERNATHEY

CEOS WGISS-49

# WHO AM I?

Physical Oceanographer

Core developer of Zarr

Ph.D. From MIT, 2012

Core developer of Xarray

Associate Prof. at Columbia / LDEO

Co-founder of Pangeo

https://ocean-transport.github.io/

Open Source Advocate

# WHAT IS PANGEO?

*"A community platform for Big Data geoscience"*

- Open Community

- Open Source Software
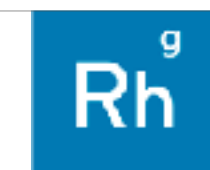
- Open Source Infrastructure

# PANGEO COMMUNITY



HTTP://PANGEO.IO

# GRASS-ROOTS ADOPTION

# PANGEO SOFTWARE ECOSYSTEM



**Inspiration: Stephan Hoyer, Jake Vanderplas (SciPy 2015)**

# PANGEO SOFTWARE ECOSYSTEM



**Inspiration: Stephan Hoyer, Jake Vanderplas (SciPy 2015)**

# JUPYTER



*"Project Jupyter exists to develop open-source software, open-standards, and services for interactive computing across dozens of programming languages."*

XARRAY

https://github.com/pydata/xarray

*"netCDF meets pandas.DataFrame"*

# XARRAY: EXPRESSIVE & HIGH-LEVEL

```python
sst_clim = sst.groupby('time.month').mean(dim='time')
sst_anom = sst.groupby('time.month') - sst_clim
nino34_index = (sst_anom.sel(lat=slice(-5, 5), lon=slice(190, 240))
                        .mean(dim=('lon', 'lat'))
                        .rolling(time=3).mean(dim='time'))

nino34_index.plot()
```



SST Anomaly in Nino 3.4 Region (5N-5S,120-170W)

https://github.com/dask/dask/

Flexible, general-purpose parallel computing framework.



ND-Arrays are split into chunks that comfortably fit in memory

Complex computations represented as a graph of individual tasks.

Scheduler optimizes execution of graph.

https://github.com/dask/dask/

Flexible, general-purpose parallel computing framework.



ND-Arrays are split into chunks that comfortably fit in memory

Complex computations represented as a graph of individual tasks.

Scheduler optimizes execution of graph.

# ZARR

**https://zarr.readthedocs.io/**

**Zarr Group:** *group_name*

.zgroup

.zattrs

**Zarr Array:** *array_name*

.zarray

.zattrs

| 0.0 | 0.1 |
| 1.0 | 1.1 |
| 2.0 | 2.1 |

- Open source library for storage of chunked, compressed ND-arrays

- Created by Alistair Miles (Imperial) for genomics research (*@alimanfoo*); now community supported standard

- Arrays are split into user-defined chunks; each chunk is optional compressed (zlib, zstd, etc.)

- Can store arrays in memory, directories, zip files, or any python mutable mapping interface (dictionary)

- External libraries (s3fs, gcsf) provide a way to store directly into cloud object storage

- Implementations in Python, C++, Java (N5), Julia, Javascript

# PANGEO ARCHITECTURE



Cloud / HPC

Distributed storage

dask

xarray

jupyter

web browser

end user

Jupyter for interactive access remote systems

"Analysis Ready Data" stored on globally-available distributed storage.

Zarr.

Parallel computing system allows users deploy clusters of compute nodes for data processing.

Dask tells the nodes what to do.

Xarray provides data structures and intuitive interface for interacting with datasets

# PANGEO DEPLOYMENTS



NASA Pleiades

NCAR Cheyenne

OCEAN.PANGEO.IO

Google Cloud Platform

Microsoft Azure

aws

HTTP://PANGEO.IO/DEPLOYMENTS.HTML

# Live example at
# [http://gallery.pangeo.io/repos/pangeo-gallery/physical-oceanography/01_sea-surface-height.html](http://gallery.pangeo.io/repos/pangeo-gallery/physical-oceanography/01_sea-surface-height.html)

```
[4]: from intake import open_catalog
     cat = open_catalog("https://raw.githubusercontent.com/pangeo-data/pangeo-datastore/master/intake-catalogs/ocean.yaml")
     ds  = cat["sea_surface_height"].to_dask()
     ds
```

[4]: xarray.Dataset

▶ Dimensions:          (**latitude**: 720, **longitude**: 1440, **nv**: 2, **time**: 8901)

▼ Coordinates:

| | | | | |
|---|---|---|---|---|
| crs | () | int32 | ... | |
| lat_bnds | (time, latitude, nv) | float32 | dask.array<chunksize=(5, 720, ... | |
| **latitude** | (latitude) | float32 | -89.875 -89.625 ... 89.625 89... | |
| lon_bnds | (longitude, nv) | float32 | dask.array<chunksize=(1440, 2... | |
| **longitude** | (longitude) | float32 | 0.125 0.375 ... 359.625 359.875 | |
| **nv** | (nv) | int32 | 0 1 | |
| **time** | (time) | datetime64[ns] | 1993-01-01 ... 2017-05-15 | |

▼ Data variables:

| | | | | |
|---|---|---|---|---|
| adt | (time, latitude, longitude) | float64 | dask.array<chunksize=(5, 720, ... | |

| | Array | Chunk |
|---|---|---|
| **Bytes** | 73.83 GB | 41.47 MB |
| **Shape** | (8901, 720, 1440) | (5, 720, 1440) |
| **Count** | 1782 Tasks | 1781 Chunks |
| **Type** | float64 | numpy.ndarray |

8901  1440  720

```
[6]: ds.sla.hvplot.image('longitude', 'latitude',
                          rasterize=True, dynamic=True, width=800, height=450,
                          widget_type='scrubber', widget_location='bottom', cmap='RdBu_r')
```

[6]:

## Create and Connect to Dask Distributed Cluster

```
[4]:  from dask_gateway import Gateway
      from dask.distributed import Client

      gateway = Gateway()
      cluster = gateway.new_cluster()
      cluster.adapt(minimum=1, maximum=20)
      cluster
```

### GatewayCluster

| | | |
|---|---|---|
| **Workers** | 0 | ▸ **Manual Scaling** |
| **Cores** | 0 | |
| **Memory** | 0 B | ▸ **Adaptive Scaling** |

**Name:** prod.f855cbed758f4a628bfc7190df860ad3

**Dashboard:** https://hub.binder.pangeo.io/services/dask-gateway/clusters/prod.f855cbed758f4a628bfc7190df860ad3/status

```
[7]: # the computationally intensive step
     sla_timeseries = ds.sla.mean(dim=('latitude', 'longitude')).load()
```

```
[9]: sla_timeseries.plot(label='full data')
     sla_timeseries.rolling(time=365, center=True).mean().plot(label='rolling annual mean')
     plt.ylabel('Sea Level Anomaly [m]')
     plt.title('Global Mean Sea Level')
     plt.legend()
     plt.grid()
```

# SHARING DATA IN THE CLOUD ERA

## Pangeo Approach: Direct Access to Cloud Object Storage

# Do we need an API for data access?

# OPENDAP

# XARRAY + ZARR

☑ Access remote netCDF-style datasets over HTTP

☑ Subset based on coordinates / variables

☑ Load data lazily

✖ Requires a server

☑ Access remote netCDF-style datasets over HTTP

☑ Subset based on coordinates / variables

☑ Load data lazily

☑ Serverless (only uses S3)

# OPENDAP



ESGF UCAR OPeNDAP

❌ Requires a server
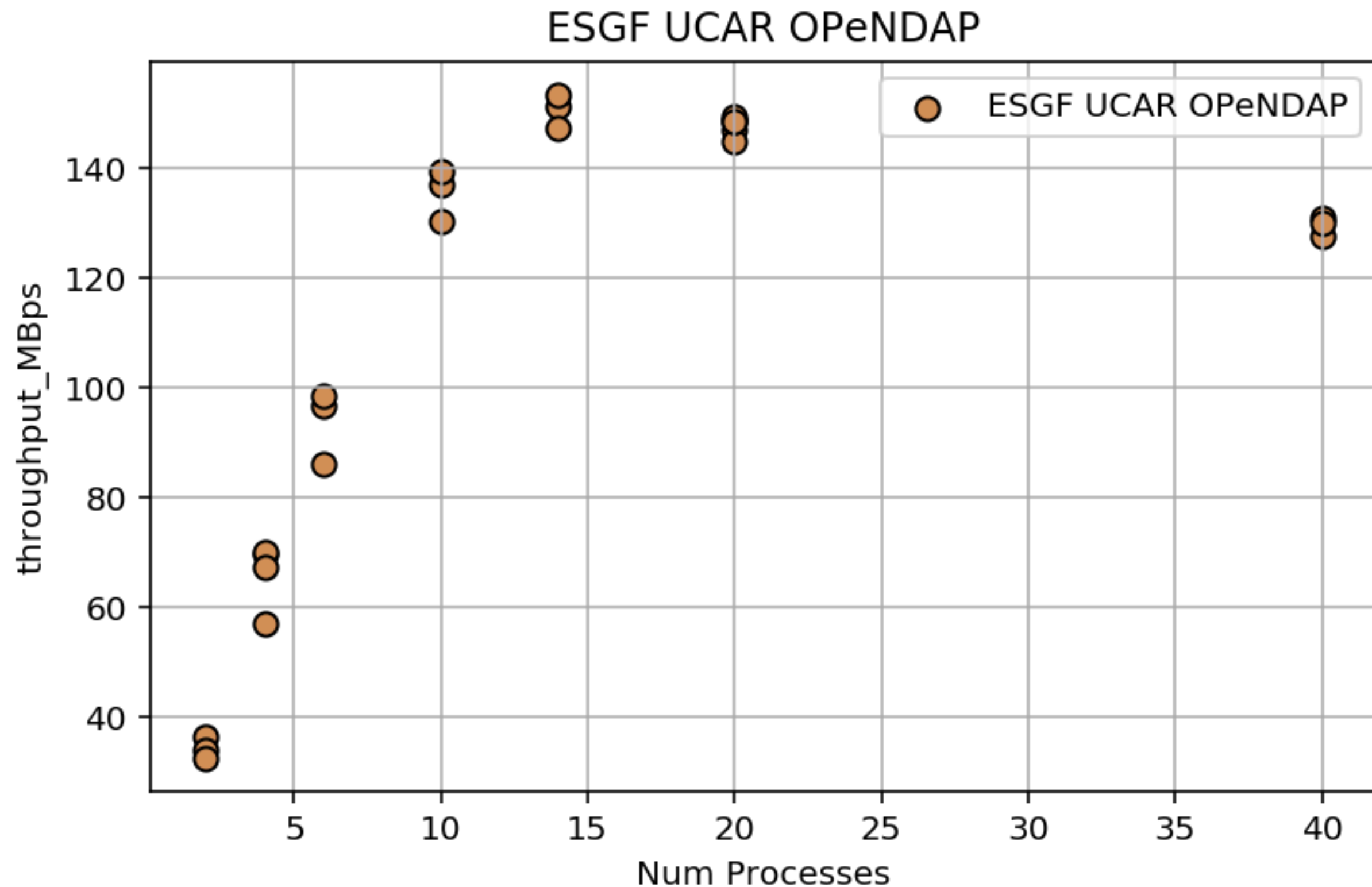
# XARRAY + ZARR

✅ Access remote netCDF-style datasets over HTTP

✅ Subset based on coordinates / variables

✅ Load data lazily
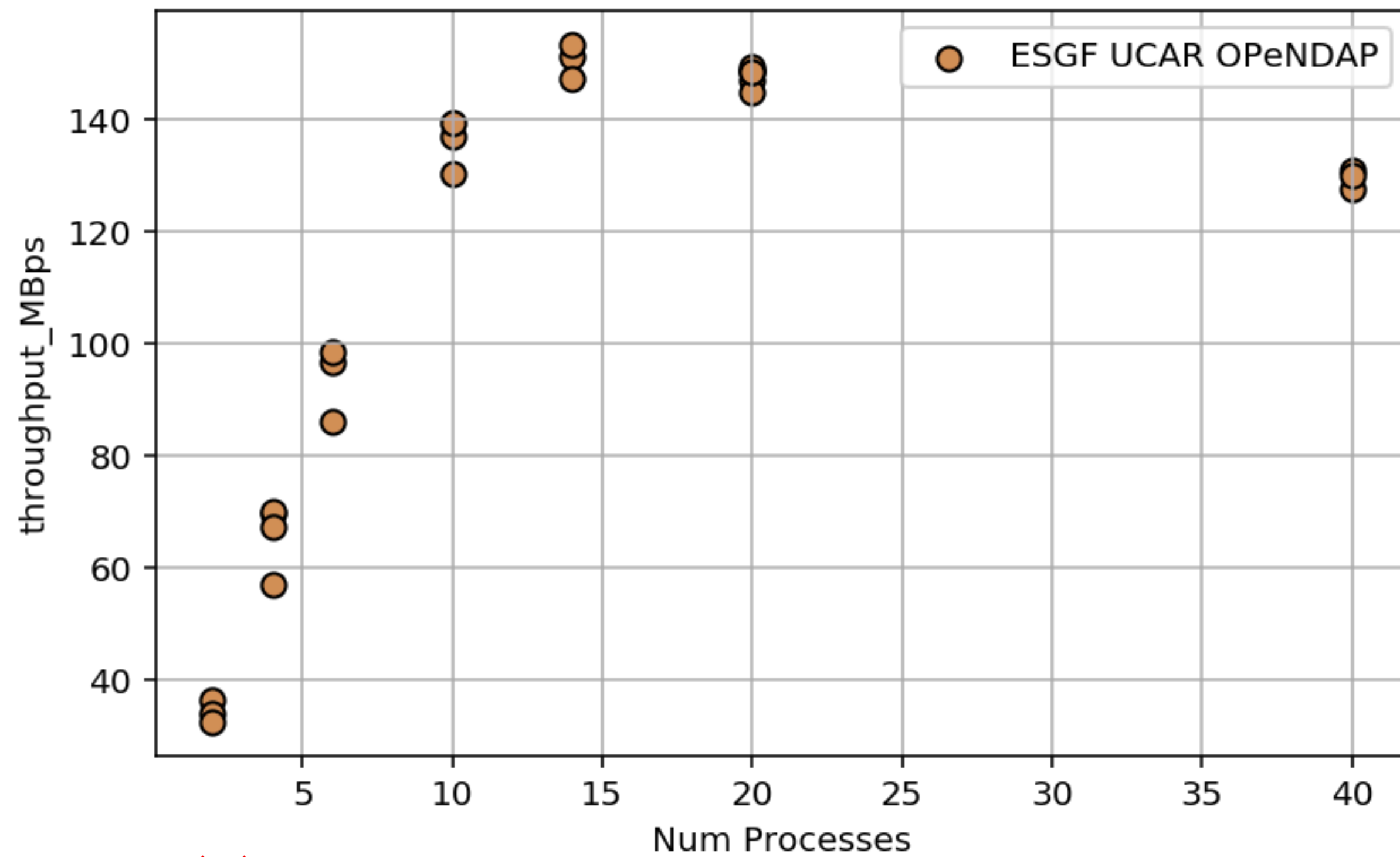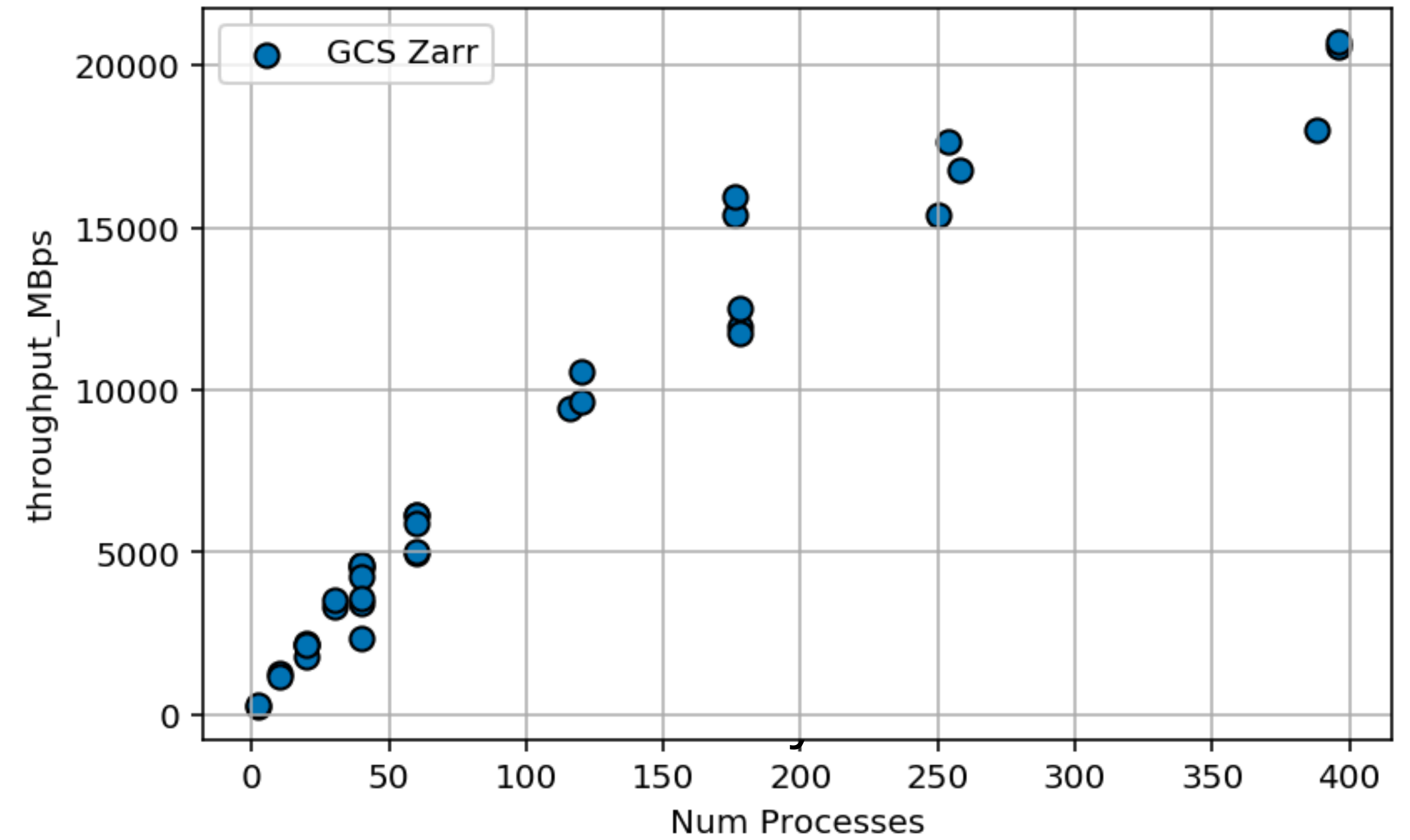
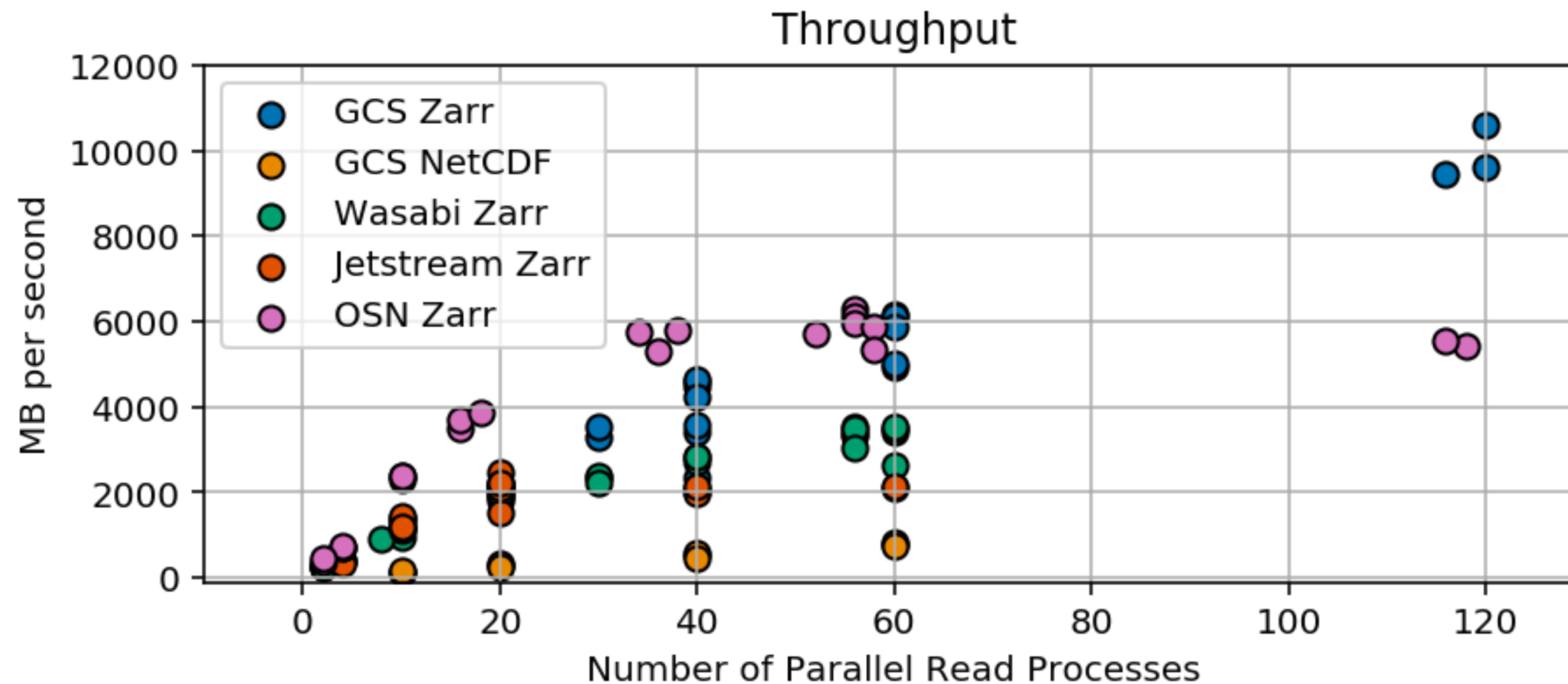✅ Serverless (only uses S3)

# OPENDAP

# XARRAY + ZARR



❌ Requires a server

✅ Serverless (only uses S3)

# CLOUD STORAGE THROUGHPUT BENCHMARKS



**To Google Cloud Region US-CENTRAL**

# PANGEO CLOUD DATA CATALOG

## CATALOG.PANGEO.IO

# TAKE-HOME MESSAGE TO WGISS

- Open source scientific python provides a great foundation for scalable earth system analytics (especially in the cloud). **Use it, don't reinvent it!**

- How do we support / sustain open-source foundational software tools? (No agency or lab "owns" these, but they are critical infrastructure.)

  **idea**: pay your staff to contribute to Xarray, Dask, etc., don't just build on top of them

- The best way to take advantage of cloud is to **give users direct access to analysis-ready data in object storage**. Don't hide it behind an API.

# LEARN MORE

http://pangeo.io

https://github.com/pangeo-data/

https://medium.com/pangeo

@pangeo_data

**Extra Slides**

# PANGEO CLOUD STACK



High-level API for analysis of multidimensional labelled arrays.

Flexible, general-purpose parallel computing framework.

Cloud-optimized storage for multidimensional arrays.

**Kubernetes**

**Object Storage**

# ZARR

**Zarr Group:** *group_name*

.zgroup     .zattrs

**Zarr Array:** *array_name*

.zarray     .zattrs

| 0.0 | 0.1 |
| 1.0 | 1.1 |
| 2.0 | 2.1 |

## Example .zarray file (json)

```json
{
    "chunks": [                  "dtype": "<f8",
        5,                       "fill_value": "NaN",
        720,                     "filters": null,
        1440                     "order": "C",
    ],                           "shape": [
    "compressor": {                  8901,
        "blocksize": 0,              720,
        "clevel": 3,                 1440
        "cname": "zstd",         ],
        "id": "blosc",           "zarr_format": 2
        "shuffle": 2         }
    },
}
```

# ZARR

**Zarr Group:** *group_name*

.zgroup  .zattrs

**Zarr Array:** *array_name*

.zarray  .zattrs

0.0  0.1

1.0  1.1

2.0  2.1

## Example .attrs file (json)

```json
{
    "_ARRAY_DIMENSIONS": [
        "time",
        "latitude",
        "longitude"
    ],
    "comment": "The sea level anomaly
is the sea surface height above mean
sea surface; it is referenced to the
[1993, 2012] period; see the product
user manual for details",
    "coordinates": "crs",
    "grid_mapping": "crs",
    "long_name": "Sea level anomaly",
    "standard_name":
"sea_surface_height_above_sea_level",
    "units": "m"
}
```