BEYOND FAIR What data infrastructure does open science need?

Ryan Abernathey

PANGEO

OOIFB Data Systems Committee Meeting, 2022

THE OPEN SCIENCE VISION



https://earthdata.nasa.gov/esds/open-science



for in everyone_in_the_world:

- for in all_scientific_knowledge:
 - 👰 .verify())

discovery = @.extend()

This would transform the S by allowing all of humanity to participate in the scientific process.

What are the barriers to realizing this vision?

12

THE OPEN SCIENCE VISION



https://earthdata.nasa.gov/esds/open-science



for in everyone in the world:

for 🔄 in everyone else in the world:

This would transform the So by allowing all of humanity to participate in the scientific process.

What are the barriers to realizing this vision?





FAIR = Findable, Accessible, Interoperable, Reusable



- FAIR is great. Nobody disagrees with FAIR.
- But making data-intensive scientific workflows FAIR is easier said than done. FAIR does not specify the protocols, technologies, or infrastructure.
 - FAIR is not a platform.

Simulation

$$ho rac{\mathrm{D} \mathbf{u}}{\mathrm{D} t} = -
abla p +
abla \cdot oldsymbol{ au} +
ho \, \mathbf{g}$$



• optimized, scalable algorithm







Big Data

https://figshare.com/articles/figure/Earth_Data_Cube/4822930/2

Data-Intensive Science

open-ended problem

VS.

- exploratory analysis
- "human in the loop"
- ML model development & training
- visualization needed
- highly varied computational patterns / algorithms



Big Data





Insights Predictions







CLAIM: OPEN SCIENCE NEEDS A PLATFORM

- The word "platform" is terribly overloaded.
- A platform is something you can build on—specifically, new scientific discoveries and new translational applications. Let's call these *projects*.
- For open science to take off at a global scale, everyone in the world needs access to the platform (like Facebook)
- This is why we are excited about cloud, but cloud as-is (e.g. AWS) is not itself an open-science platform.
- Does the open science platform need to be open?









OUTLINE

10 mins	The status quo of da
10 mins	Cloud computing and
10 mins	From Software to Sa
10 mins	Where are things hea



- ata-intensive scientific infrastructure
- d Pangeo
- aS: Pangeo Forge and Earthmover
- aded?



PART I: THE STATUS QUO

DATA-INTENSIVE SCIENCE INFRASTRUCTURE: THE STATUS QUO*





Personal Laptop

Group Server

more storage, more CPU, more security, more constraints







Department Cluster

Agency Supercomputer





STATUS QUO: WHAT INFRASTRUCTURE CAN WE RELY ON?



Files / POSIX filesystems

V Programming languages: C, FORTRAN, Python, R, Julia





Batch queuing system HPC only

1 The internet Not on HPC nodes!

I Globus for file transfer 🔁 Not supported everywhere

K High level data services, APIs, etc. Virtually unknown in my world





THE "DOWNLOAD" MODEL









THE "DOWNLOAD" MODEL







THE "DOWNLOAD" MODEL





THE "DOWNLOAD" MODEL





THE "DOWNLOAD" MODEL

PRIVILEGED INSTITUTIONS CREATE "DATA FORTRESSES*"

Image credit: Moahim, CC BY-SA 4.0, via Wikimedia Commons



*Coined by Chelle Gentemann





PRIVILEGED INSTITUTIONS CREATE "DATA FORTRESSES*"



Data

Image credit: Moahim, CC BY-SA 4.0, via Wikimedia Commons





open("/some/random/files/on/my/cluster")

*Coined by Chelle Gentemann





PROBLEMS WITH THE STATUS QUO

- even simulation output.
- can't be shared.
- \$ a big agency supercomputer. This really limits participation.
- feature, not a bug. Restricts collaboration and reproducibility!



Each fortress is a special snowflake. Code developed in one will not run inside another.



Emphasis on *files* as a medium of data exchange creates lots of work for individual scientists (downloading, organizing, cleaning). Most file-based datasets are a mess-

left the hard work of data wrangling is rarely collaborative. Outputs are not reusable and

Doing data-intensive science requires either expensive local infrastructure or access to

2 Data intensive science is locked inside *data fortresses*. Limiting access to outsiders is a











PART II: CLOUD COMPUTING AND PANGE0







WHAT ABOUT CLOUD?

Can we create a "data watering hole"* instead of a fortress?

Research









Education & Outreach



Coastal Ocean Environment Summer School in Ghana

Industry Partners

OPTION A: VERTICALLY INTEGRATED PLATFORM



All the data All the compute

Θ PANGEO

Contact Us

Solutions

Demos v

Observed data is.

cleaned, calibrated,

Company v





OPTION B: INTEROPERABLE CLOUD-NATIVE DATA, SOFTWARE, AND SERVICES

Community-Maintained ARCO Data Lake[s]





Data Provider's Resources







22



Data Consumer's Resources



THE PANGEO COMMUNITY PROCESS

Scientific users / use cases







- Define science questions
- Use software / infrastructure
- Identify bugs / bottlenecks
- Provide feedback to developers



Agile development

Open-source software libraries



- Contribute widely the the open source scientific python ecosystem
- Maintain / extend existing libraries, start new ones reluctantly

• Solve integration challenges

HPC and cloud infrastructure

- Deploy interactive analysis environments
- Curate analysis-ready datasets
- Platform agnostic

Azure Microso

Google Cloud Platform





ACCIDENTAL PIVOT TO CLOUD

Uz, Baris M <bmuz@nsf.gov> to me, kpaul@ucar.edu -

Dear Ryan and Kevin,

I am delighted to give you good news about your recent EarthCube proposal. Based on supportive reviews and our excitement about the proposal, we are prepared to make a funding recommendation for it. The total budget in your proposal was for \$1.46M. Unfortunately, it looks like we are only going to be able to allocate \$1.2M to this project. This would be about an 18% reduction in the total budget. I would like to know how much of an impact this would have in the scope and/or chance of success of your project.

Ryan Abernathey <rpa@ldeo.columbia.edu> to Baris, kpaul@ucar.edu 👻

Hi Mete,

Thanks for this excellent news. We would define way to accomplish this.

In the meantime, I have a few questions: - What is the time frame for us to submit the revised budget? - Do you have any specific suggestions about how to rebduget, perhaps informed by feedback from the reviews? - Is there any chance of us getting access to commercial cloud computing resources via NSF (as with the BIGDATA program https://www.nsf.gov/funding/pgm_summ.jsp?pims_id=504767)? That could help reduce our computing budget.

Thanks,

Ryan Abernathey



Mon, Jun 26, 2017, 7:01 PM

Mon, Jun 26, 2017, 8:13 PM

5

Thanks for this excellent news. We would definitely like to rebudget and accept the \$1.2M allocation. I will confer with my co-I's about the best



THE PANGEO CLOUD-NATIVE STACK







Rich interactive computing environment in the web browser.











Domain specific packages







Etc.

High-level API for analysis of multidimensional labelled arrays.

Cloud Infrastructure

Kubernetes

Flexible, general-purpose parallel computing framework.

Cloud-optimized storage for multidimensional arrays.









ANALYSIS-READY, CLOUD OPTIMIZED: ARCO DATA

THEME ARTICLE: JUPYTER IN COMPUTATIONAL SCIENCE

Cloud-Native Repositories for Big Scientific Data

Ryan P. Abernathey[®], Charles C. Blackmon-Luca, Timothy J. Crone, Naomi Henderson, Chiara Lepore[®], Lamont–Doherty Earth Observatory of Columbia University, Palisades, NY, 10964, USA

Tom Augspurger ⁽⁰⁾, Anaconda, Des Moines, IA, 50313, USA

Anderson Banihirwe [©] and Joseph J. Hamman [©], National Center for Atmospheric Research, Boulder, CO, 80305, USA

Chelle L. Gentemann, Farallon Institute, Petaluma, CA, 94952, USA

Theo A. McCaie ¹⁰ and Niall H. Robinson, *Met Office, University of Exeter, Exeter EX4 4PY, U.K.*

Richard P. Signell ¹⁰, US Geological Survey, Woods Hole, MA, 02543, USA

https://doi.org/10.1109/MCSE.2021.3059437

This also demonstrates the potential of the "hybrid cloud" model with OSN.

B PANGEO

A array A array C arr Z arr Z arr Object Object Object Object Object Object Object Object





Andrew Pauling Dandrewp109

#cmip6hack is just wrapping up, and has changed the way I will think about, and hopefully do, climate model analysis in the future. The <a>opangeo_data infrastructure makes it all so easy.

5:29 PM - Oct 18, 2019 - Twitter Web App

6 Retweets 24 Likes



Erin Dougherty @edougherty_

My #dayofscience: paper revisions •• and using @pangeo_data to analyze massive amounts of high-res climate data to understand floods in a current and future climate over the U.S. When not doing this, I 💙 observing #wx directly via field work and watching storms.



2:38 PM · Oct 15, 2019 · Twitter for iPhone

5 Retweets 51 Likes

October 20, 2021 (v1) Presentation Open Access

D Augspurger, Tom; Morris, Dan; Emanuele, Rob; McFarland, Matt; Sanchez-Andrade Nuno, Bruno;

We'll see how Microsoft's Planetary Computer team has embraced the ideas and techniques pioneered by pangeo, putting them to work for sustainability. We'll see the use of STAC for cataloging and querying data, analysis-ready data in blob storage, and JupyterHub-based environments for sca

Uploaded on October 25, 2021

November 3, 2021 (v1) Presentation Open Access

xcube - Python package for Earth Observation data cubes

Brandt, Gunnar; Norman Fomferra;

xcube is a software package and toolkit for generation and exploitation of analysis-ready data cubes from Earth Observation and other geographical data. It is based on Python's popular data science technology stack, particularly on xarray, dask, and zarr, and is freely available as open-source s

Uploaded on November 11, 2021



PANGEO HAS BROAD ADOPTION

Scalable Sustainability with the Planetary Computer





European Space Agency





LIMITATIONS OF THE PANGEO APPROACH

- Pangeo software can be deployed as a platform: sources
- Sharing projects between these hubs is still very hard
- contractor services). There is no "Pangeo as a Service"
- Getting data into the cloud in ARCO format is hard and full of toil



JupyterHub in the cloud with Xarray, Dask, etc., connected to ARCO data

• But there are many distinct deployments of this platform - dozens of similar yet distinct JupyterHubs with different configurations, environments, capabilities, etc

• Deploying hubs generally requires DevOps work (billed a developer time or

CHALLENGES WITH CLOUD IN GENERAI



- Non-agency scientists have many barriers to adopting cloud: Overhead policies, purchasing challenges, lack of IT support, etc.
- Cloud is too complicated! The services offered are not useful to scientists: An extra layer of science-oriented services must be developed
- Europe basically forbids scientists from using US-based cloud providers
- Not much has changed for university scientists since 2017



Enabling Computer and Information Science and Engineering Research and Education in the Cloud • May 2018

Publisher: National Science Foundation, Arlington, Va, United States

🗏 Purchase this Technical Report

in Recommend ACM DL

ALREADY A SUBSCRIBER? SIGN IN



PART III: FROM SOFTWARE TO SAAS

PANGEO FORGE & EARTHMOVER



TOOLS FOR COLLABORATION



These are all proprietary SaaS (Software as a Service) applications. They may use open standards, but they are not open source.

We (or our institutions) have no problem paying for them.



Some of the most impactful services used in open science...

THF MODFRN DATA STACK

🔀 Continual



https://continual.ai/post/the-modern-data-stack-ecosystem-fall-2021-edition

32



- In the past 5 years, a platform has emerged for *enterprise data science* called the Modern Data Stack
- The MDS is centered around a "data lake" or "data warehouse"
- Different platform elements are provided by different SaaS companies; integration through standards and APIs
- No one in science uses any of this stuff

- Embrace commercial SaaS: a Modern Data Stack for Science
- Cultivate community-operated SaaS: e.g. Wikipedia, Conda Forge, Binder, 2i2c Hubs, Pangeo Forge
- We probably need a mix of both



mmunity-operated SaaS for ETL (Extract / Transform / Load) of ARCO Data



HOW CAN WE DELIVER AN OPEN SCIENCE PLATFORM IN A SCALABLE, SUSTAINABLE WAY?



Our new startup. Building a commercial **cloud data** lake platform for scientific data.

ARCO DATA Analysis Ready, Cloud Optimzed

What is "Analysis Ready"?

- Think in "datasets" not "data files"
- No need for tedious homogenizing / cleaning steps
- Curated and cataloged





- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

How do data scientists spend their time? Crowdflower Data Science Report (2016)

What is "Cloud Optimized"?

- Compatible with object storage (access via HTTP)
- Supports lazy access and intelligent subsetting
- Integrates with high-level analysis libraries and distributed frameworks





THEME ARTICLE: JUPYTER IN COMPUTATIONAL SCIENCE

Cloud-Native Repositories for Big Scientific Data

Ryan P. Abernathey [©], Charles C. Blackmon-Luca, Timothy J. Crone, Naomi Henderson, Chiara Lepore [©], Lamont–Doherty Earth Observatory of Columbia University, Palisades, NY, 10964, USA

Tom Augspurger ⁽⁰⁾, Anaconda, Des Moines, IA, 50313, USA

Anderson Banihirwe ⁽⁾ and Joseph J. Hamman ⁽⁾, National Center for Atmospheric Research, Boulder, CO, 80305, USA

Chelle L. Gentemann, Farallon Institute, Petaluma, CA, 94952, USA

Theo A. McCaie ¹⁰ and Niall H. Robinson, *Met Office, University of Exeter, Exeter EX4 4PY, U.K.*

Richard P. Signell ¹⁰, US Geological Survey, Woods Hole, MA, 02543, USA

https://doi.org/10.1109/MCSE.2021.3059437

This also demonstrates the potential of the "hybrid cloud" model with OSN.



ARCO DATA IS FAST!





PROBLEM: Making ARCO Data is Hard!

To produce useful ARCO data, you must have:

Domain Expertise: How to find, clean, and homogenize data

Tech Knowledge: How to efficiently produce cloud-optimized formats



Compute Resources: A place where to stage and upload the ARCO data

Analysis Skills: To validate and make use of the ARCO data.



Data Scientist

WHOSE JOB IS IT TO MAKE ARCO DATA?



EILU Diff





Data providers are concerned with preservation and archival quality.

Scientists users know what they need to make the data analysis-ready.

Let's democratize the production of ARCO data!

Domain Expertise: How to find, clean, and homogenize data



39



PANGEO FORGE





CONDA-FORGE

A community-led collection of recipes, build infrastructure and distributions for the conda package manager.

40



INSPIRATION: CONDA FORGE





README.md Public

Hi there 👏

conda-forge is a community led collection of recipes, build infrastructure and distributions for the conda package manager

Important git repositories

- staged-recipes place to submit new recipes
- miniforge An installer with conda-forge as the default channel
- conda-forge.github.io website and tracker for general conda-forge problems and enhancements
- feedstocks A monorepo containing all feedstocks as submodules
- conda-smithy The tool for managing conda-forge feedstocks.
- admin-requests Github repo to ask conda-forge/core to mark a package as broken
- · conda-forge-repodata-patches-feedstock Git repo that track hotfixing of package metadata
- conda-forge-pinning-feedstock Global pinnings in conda-forge and migration information





Pangeo Forge Recipes

https://github.com/pangeo-forge/pangeo-forge-recipes



Open source python package for describing and running data pipelines ("recipes")



Pangeo Forge Cloud

https://pangeo-forge.org/



Cloud platform for automatically executing recipes stored in GitHub repos.



PANGEO FORGE RECIPES https://pangeo-forge.readthedocs.io/

Describes where to find the source files which are the inputs to the recipe



Describes where to store the outputs of our recipe



StorageConfig



A complete, selfcontained representation of the pipeline



PANGEO FORGE CLOUD https://pangeo-forge.org/

Contains the code and metadata for one or more Recipes



Feedstock



terraclimate-feedstock

A pangeo-smithy repository for the terraclimate dataset.

● Python Apache-2.0 😚 3 🏠 2 🕛 1 🎝 3 Updated on Jan 1

noaa-oisst-avhrr-feedstock

● Python Apache-2.0 % 2 ☆ 1 ① 0 \$ 4 Updated on Jan 1

Runs the recipes in the cloud using elastic scaling clusters













GCS





VISION: COLLABORATIVE DATA CURATION







Feedstock

GitHub

A pangeo-smithy repository for the terraclimate dataset.

● Python ▲ Apache-2.0 % 3 ☆ 2 ① 1 \$3 Updated on Jan 1

...Oh the metadata need an update.

● Python 小 Apache-2.0 % 2 ☆ 1 ① 0 \$ 4 Updated on Jan 1





Data Producer











Oceanographers building a full-stack cloud SaaS automation platform.

https://twitter.com/Colinoscopy/status/1255890780641689601

Pangeo Forge Cloud is live and open for business! pangeo-forge.org

• Our recipes often contain 10000+ tasks. We are hitting the limits on Prefect as a workflow engine. Currently refactoring to move to Apache Beam.

Data has lots of edge cases! This is really hard. X

But we remain very excited about the potential of Pangeo Forge to transform how scientists interact with data.



CURRENT STATUS

EARTHMOVER A Public Benefit Corporation

Founders:





Ryan Abernathey

Joe Hamman



Mission: To empower people to use scientific data to solve humanity's greatest challenges

Product: ArrayLake

- High performance for analytics (based) on Zarr data model)
- Ingest and index data from archival formats (NetCDF, HDF, GRIB, etc.)
- Automatic background optimizations
- Versioning / snapshots / time travel
- Data Governance
- Compare to Databricks, Snowflake, Dremio



PART IV: WHERE ARE WE HEADING?



PILLARS OF CLOUD NATIVE SCIENTIFIC DATA ANALYTICS

1. Analysis-Ready, Cloud-Optimized Data

49

2. Data-Proximate Computing

xarray.Dataset					
Dimensions:	(latitude:	720, longitude:	1440, nv : 2, t	t ime : 8901)	
 Coordinates: 					
crs	0		int32	m	
lat_bnds	(time, latitude, nv) (latitude) (longitude, nv)		float32	dask.array <chunksize=(-89.875 -89.625 89, dask.array<chunksize=(< td=""><td rowspan="3"></td></chunksize=(<></chunksize=(
latitude			float32		
lon_bnds			float32		
longitude	(longitude)	float32	0.125/0.375 359.625	65
nv	(nv)		int32	Ø 1	8
time	(time) dateti		datetime64[ns]	1993-01-01 2017-05	85
axis : long_name : standard_na	T Time time				
Data variables:					
adt	(time, latit	ude, longitude)	float64	dask.array <chunksize=(< td=""><td></td></chunksize=(<>	
	Bytes Shape Count Type	Array 73.83 GB (8901, 720, 14 1782 Tasks float64	Chunk 41.47 M 440) (5, 720 1781 Ch numpy.r	1B b, 1440) hunks hdarray 1440	728
err	(time, latitude, longitude)		float64	dask.array <chunksize=(,< td=""><td></td></chunksize=(,<>	
sla	(time, latitude, longitude)		float64	dask.array <chunksize=(< td=""><td></td></chunksize=(<>	
ugos	(time, latitude, longitude)		float64	dask.array <chunksize=(< td=""><td>85</td></chunksize=(<>	85
ugosa	(time, latitude, longitude)		float64	dask.array <chunksize=(< td=""><td></td></chunksize=(<>	
vgos	(time, latitude, longitude)		float64	dask.array <chunksize=(< td=""><td>BS</td></chunksize=(<>	BS
	(time, latitude, longitude)				



3. Elastic Distributed Processing





end user



compute node

compute node



SEPARATION OF STORAGE AND COMPUTF



Storage costs are steady.

Data provider pays for storage costs.

May be subsidized by cloud provider. (Thanks AWS, GCP, Azure!

Or can live outside the cloud (e.g. Wasabi, OSN)



This is completely different from the status quo infrastructure!

OPEN SCIENCE PLATFORM





Google Cloud Platform



Runs on any modern cloud-like platform or on premises data center





FEDERATED, EXTENSIBLE MODEL





DATA GRAVITY

"Data gravity is the ability of a body of data to attract applications, services and other data." - Dave McCrory







DATA GRAVITY

What is the stable steady-state solution?

?



Google Cloud



DOE







HECC

XSEDE

WE NEED A GLOBAL SCIENTIFIC DATA COMMONS

Need to be exploring: edge storage, decentralized web, web3





SHOUT OUTS



















SGCI

Science Gateways Community Institute

the international interactive computing collaboration











unidata





DISCUSSION QUESTIONS

- What's the right model to deliver data and computing services to the research community? Commercial vendors? Co-ops?
- How can we avoid recreating existing silos in the cloud?
- Who should we pay for cloud infrastructure for the science community? University? Agency? PI?
- How can we make cloud interoperate more with HPC and on-premises computing resources?



LEARN MORE

















http://pangeo.io

https://discourse.pangeo.io/

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

https://medium.com/pangeo

@pangeo_data