

Scalable and Computationally Reproducible Approaches to Arctic Research

S. Jeanette Clark, Juliet Cohen, Matthew B. Jones, Anna Liljedahl, Daphne Virlar-Kn

3/27/23

Table of contents

Preface	8
About	8
Schedule	8
Code of Conduct	8
Setting Up	9
Download VS Code and Remote - SSH Extension	9
Log in to the server	10
Install extensions on the server	11
Test your local setup (Optional)	11
Create a (free) Google Earth Engine (GEE) account	13
About this book	14
1 Welcome and Overview	15
1.1 Arctic Data Center Overview	15
1.1.1 Data Discovery Portal	19
1.1.2 Tools and Infrastructure	21
1.1.3 Support Services	23
1.1.4 Training and Outreach	23
1.1.5 Data Rescue	24
1.1.6 Who Must Submit	25
1.1.7 Supporting data reuse	26
1.1.8 Summary	27
1.2 Scalable Computing Topics	28
2 Remote Computing	29
2.1 Introduction	29
2.2 Servers & Networking	29
2.3 IP addressing	30
2.4 Bash Shell Programming	32
2.4.1 Some commonly used (and very helpful) bash commands:	33
2.4.2 General command syntax	34

2.4.3	Some useful keyboard shortcuts	34
2.5	Connecting to a remote computer via a shell	34
2.6	Git via a shell	35
2.7	Let's practice!	36
2.7.1	Exercise 1: Connect to a server using the <code>ssh</code> command (or using VS Code's command palette)	36
2.7.2	Exercise 2: Practice using some common bash commands	37
2.7.3	Exercise 3: Clone a GitHub repository to the server	40
2.7.4	Bonus Exercise: Automate data processing with a Bash script	41
3	Python Programming on Clusters	46
3.1	Introduction	46
3.2	Connect to the server (if you aren't already connected)	47
3.3	Virtual Environments	47
3.4	Brief overview of python syntax	51
3.5	Jupyter notebooks	54
3.5.1	Load libraries	55
3.5.2	Read in a csv	56
3.6	Functions	61
3.7	Summary	63
4	Pleasingly Parallel Programming	64
4.1	Introduction	64
4.2	Why parallelism?	66
4.3	Processors (CPUs), Cores, and Threads	66
4.4	Parallel processing in the shell	70
4.5	Modes of parallelization	70
4.6	Task parallelization in Python	73
4.7	Exercise: Parallel downloads	75
4.7.1	Serial	77
4.7.2	Multi-threaded with <code>concurrent.futures</code>	78
4.7.3	Multi-process with <code>concurrent.futures</code>	79
4.8	Parallel processing with <code>parsl</code>	81
4.9	When to parallelize	84
4.10	Parallel Pitfalls	85
4.11	Summary	86

4.12	Further Reading	86
5	Data Structures and Formats for Large Data	87
5.1	Introduction	87
5.2	NetCDF Data Format	88
5.2.1	Data Model	90
5.2.2	Metadata Standards	91
5.2.3	Exercise	92
5.3	<code>xarray</code>	95
5.3.1	<code>xarray.DataArray</code>	95
5.3.2	<code>xarray.DataSet</code>	103
5.3.3	Exercise	105
5.4	Tabular Data and NetCDF	108
5.4.1	Tabular to NetCDF	108
5.4.2	<code>pandas</code> to <code>xarray</code>	109
6	Parallelization with Dask	111
6.1	Introduction	111
6.2	Dask Cluster	112
6.2.1	Setting up a Local Cluster	112
6.2.2	Dask Dashboard	114
6.3	<code>dask.dataframes</code>	115
6.3.1	Reading a csv	116
7	Encoding?	117
7.0.1	Lazy Computations	118
7.1	<code>dask.arrays</code>	118
7.2	Dask and <code>xarray</code>	120
7.2.1	Open .tif file	120
7.2.2	Calculating NDVI	122
7.3	Best Practices	122
8	Group Project: Staging and Preprocessing	124
8.1	Introduction	124
8.2	Staging and Tiling	126
9	Software Design I: Functions and Concurrency	128
9.1	Why functions?	128
9.2	Functions as objects	129
9.3	Global variables	131
9.4	Pure functions	132
9.5	Task dependencies	133

9.6 Locks	133
9.7 Race conditions	137
9.8 Deadlocks	137
9.9 Further reading	137
10 Spatial and Image Data Using GeoPandas	138
10.1 Introduction	138
10.2 Pre-processing raster data	139
10.3 Pre-processing vector data	143
10.4 Crop data to area of interest	150
10.4.1 Check extents	151
10.5 Calculate total distance per fishing area	152
10.6 Summary	159
11 Billions of Ice Wedge Polygons	160
11.1 Billions of Ice Wedge Polygons	160
11.2 Developing Science Gateways with Clowder for the Permafrost Discover Gateway	161
12 Group Project: Data Processing	162
13 Data Ethics for Scalable Computing	163
13.1 Intro to Data Ethics	163
13.2 Ethics at the Arctic Data Center	166
13.2.1 Ethical Considerations for all Arctic Re- searchers	167
13.2.2 Archaeological and Paleontological Research	168
13.2.3 Human Participation and Sensitive Data .	169
13.2.4 Marines Sciences (e.g. Marine Biology Research)	169
13.2.5 Physical Sciences (e.g. Geology, Glaciol- ogy, and Ice Research)	170
13.2.6 Plant and Soil Research	171
13.2.7 Spatial Data	172
13.2.8 Wildlife Sciences (e.g. Ecology and Biol- ogy Research)	173
13.3 Ethics in Machine Learning	174
13.3.1 ImageNet: A case study of ethics and bias in machine learning	174
13.4 References and Further Reading	175

14 Google Earth Engine	177
14.1 Introduction	177
14.2 Exercise 1: An introductory lesson on using Google Earth Engine	178
14.2.1 Part i. Setup	178
14.2.2 Part ii. Explore the ERA5 Daily Aggregates Data	180
14.2.3 Part iii. Visualize global precipitation using ERA5 Daily Aggregate data	180
14.2.4 Part iv. Extract annual sum of precipitation data from points	184
14.2.5 Part v. Takeaways	185
14.3 Exercise 2: Visualize fire dynamics in the Arctic using GEE	186
14.4 Additional Resources	187
15 Documenting and Publishing Data	189
15.1 Introduction	189
15.2 Metadata	189
15.3 Data Package Structure	191
15.4 Archiving Data: The Large Data Perspective	192
15.4.1 Step 1: The Narrative Metadata Submission	193
15.4.2 Step 2: Adding File & Variable Level Metadata	202
15.4.3 Step 3: Uploading Large Data	206
15.5 Best Practices for Submitting Large Data	209
15.5.1 Data Organization	209
15.6 Data transfer tools	211
15.7 Summary	214
16 Group Project: Visualization	215
17 Software Design II: Modules, Packages, and Testing	216
17.1 Python modules and packages	216
17.2 Anatomy of modules and packages	217
17.3 Managing packages with <i>poetry</i>	219
17.3.1 Get poetry in a new virtual environment .	219
17.4 Setup package workspace	220
17.4.1 Open VS Code window	220
17.4.2 Create a new package	221

17.4.3	Setup the workspace	222
17.5	The <code>adcmodel</code> example package	223
17.5.1	Load software dependencies	223
17.5.2	Adding code to modules	225
17.5.3	Add development packages	226
17.6	Testing	227
17.7	Summary and next steps	229
17.8	Further reading	229
18	What is Cloud Computing Anyways?	230
18.1	What is cloud computing anyways?	230
18.2	Commercial clouds	231
18.3	Academic clouds	231
18.3.1	Cloud deployment options	231
18.4	Service Models	233
18.5	An Analogy: Pizza as a Service	234
18.6	Virtual Machines and Containers	234
19	Reproducibility and Containers	236
19.1	Software collapse	236
19.2	Delaying the inevitable	237
19.3	Hands-off with Containers and Docker	239
19.4	Discussion	241
20	Parquet and Arrow	242
20.1	Introduction	242
20.2	Row major vs column major	243
20.2.1	Row major versus column major files	244
20.3	Parquet	245
20.4	Arrow	245
20.5	Example	246
20.6	Synopsis	249

Preface

About

This 5-day in-person workshop will provide researchers with an introduction to advanced topics in computationally reproducible research in python, including software and techniques for working with very large datasets. This includes working in cloud computing environments, docker containers, and parallel processing using tools like parsl and dask. The workshop will also cover concrete methods for documenting and uploading data to the Arctic Data Center, advanced approaches to tracking data provenance, responsible research and data management practices including data sovereignty and the CARE principles, and ethical concerns with data-intensive modeling and analysis.



Schedule

Code of Conduct

Please note that by participating in this activity you agree to abide by the [NCEAS Code of Conduct](#).

08:00-08:30	Coffee (optional)	Coffee (optional)	Coffee (optional)	Coffee (optional)	Coffee (optional)	08:00-08:30
08:30-09:00	1. Welcome and Course Overview (Matt)	5. Data structures and formats for large data (Jeanette)	9. Spatial and Image Data using GeoPandas (Juliet)	13. Google Earth Engine (Jeanette)	17. What is cloud computing anyways? (Matt)	08:30-09:00
09:00-09:30						09:00-09:30
09:30-10:00	2. Remote computing (Daphne)					09:30-10:00
10:00-10:30					Survey Feedback	10:00-10:30
10:30-11:00	BREAK	BREAK	BREAK	BREAK	BREAK	10:30-11:00
11:00-11:30	3. Python programming on clusters (Jeanette)	6. Parallelization with Dask (Jeanette)	10. Billions of Ice Wedge Polygons (Anna)	14. Documenting and Publishing Data (Daphne)	18. Reproducibility redux via containers (Matt)	11:00-11:30
11:30-12:00					Q & A	11:30-12:00
12:00-12:30	Lunch	Lunch	Lunch	Lunch		12:00-12:30
12:30-13:00					Adjourn	12:30-13:00
13:00-13:30	3. Python programming on clusters (Jeanette)	7. Group project I Data staging and pre-processing (Juliet)	11. Group project II Parallel data processing (Juliet)	15. Group project III Visualizing big geospatial data (Juliet)		13:00-13:30
13:30-14:00						13:30-14:00
14:00-14:30	4. Pleasingly Parallel Programming (Matt)					14:00-14:30
14:30-15:00						14:30-15:00
15:00-15:30	Break	Break	Break	Break		15:00-15:30
15:30-16:00	4. Pleasingly Parallel Programming (Matt)	8. Software design I (Matt)	12. Data Ethics (Daphne)	16. Software design II (Matt)		15:30-16:00
16:00-16:30						16:00-16:30
16:30-17:00	Q&A	Q&A	Q&A	Q&A		16:30-17:00

Setting Up

In this course, we will be using Python (3.9.13) as our primary language, and VS Code as our IDE. Below are instructions on how to get VS Code set up to work for the course. If you are already a regular Python user, you may already have another IDE set up. We strongly encourage you to set up VS Code with us, because we will use your local VS Code instance to write and execute code on one of the NCEAS servers.

Download VS Code and Remote - SSH Extension

First, [download VS Code](#) if you do not already have it installed.

You'll also need to download the [Remote - SSH extension](#).

Log in to the server

To connect to the server using VS Code follow these steps, from the VS Code window:

- open the command palette (Cmd + Shift + P)
- enter “Remote SSH: Connect to Host”

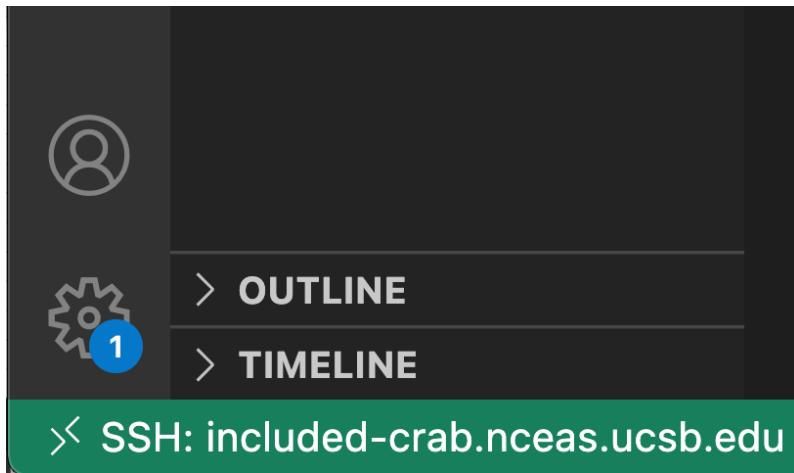


- select “Add New SSH Host”
- enter the ssh command to connect to the host as if in a terminal (`ssh username@included-crab.nceas.ucsb.edu`)
 - Note: you will only need to do this step once



- select the SSH config file to update with the name of the host. You should select the one in your user directory (eg: `/Users/jclark/.ssh/config`)
- click “Connect” in the popup in the lower right hand corner
 - Note: If the dialog box does not appear, reopen the command palette (Cmd + Shift + P), type in “Remote-SH: Connect to Host...”, choose included-crab.nceas.ucsb.edu from the options of configured SSH hosts, then enter your password into the dialog box that appears
- enter your password in the dialog box that pops up

When you are connected, you will see in the lower left hand corner of the window a green bar that says “SSH: included-crab.nceas.ucsb.edu.”



Install extensions on the server

After connecting to the server, in the extensions pane (View > Extensions) search for, and install, the following extensions:

- Python
- Jupyter
- Jupyter Keymap

Note that these extensions will be installed on the server, and not locally.

Test your local setup (Optional)

We are going to be working on the server exclusively, but if you are interested in setting up VS Code to work for you locally with Python, you can follow these instructions. This local setup section summarizes the official VS Code tutorial. For more detailed instructions and screenshots, see the [source material](#). This step is 100% optional, if you already have an IDE set up to work locally that you like, or already have VS code set up to work locally, you are welcome to skip this.

Locally (not connected to the server), check to make sure you have Python installed if you aren't sure you do. File > New Window will open up a new VS Code window locally.

To check your python, from the terminal run:

```
python3 --version
```

If you get an error, it means you need to install Python. Here are instructions for getting installed, depending on your operating system. Note: There are many ways to install and manage your Python installations, and advantages and drawbacks to each. If you are unsure about how to proceed, feel free to reach out to the instructor team for guidance.

- Windows: Download and run an installer from [Python.org](https://www.python.org).
- Mac: Install using [homebrew](#). If you don't have homebrew installed, follow the instructions from their webpage.

```
– brew install python3
```

After you run your install, make sure you check that the install is on your system PATH by running `python3 --version` again.

Next, install the [Python extension for VS Code](#).

Open a terminal window in VS Code from the Terminal drop down in the main window. Run the following commands to initialize a project workspace in a directory called `training`. This example will show you how to do this locally. Later, we will show you how to set it up on the remote server with only one additional step.

```
mkdir training
cd training
code .
```

Next, select the Python interpreter for the project. Open the **Command Palette** using Command + Shift + P (Control + Shift + P for windows). The Command Palette is a handy tool in VS Code that allows you to quickly find commands to VS Code, like editor commands, file edit and open commands, settings, etc. In the Command Palette, type “Python: Select Interpreter.” Push return to select the command, and then

select the interpreter you want to use (your Python 3.X installation).

To make sure you can write and execute code in your project, [create a Hello World test file](#).

- From the File Explorer toolbar, or using the terminal, create a file called `hello.py`
- Add some test code to the file, and save

```
msg = "Hello World"
print(msg)
```

- Execute the script using either the Play button in the upper-right hand side of your window, or by running `python3 hello.py` in the terminal.
 - For more ways to run code in VS Code, see the [tutorial](#)

Finally, to test Jupyter, download the [Jupyter extension](#). You'll also need to install `ipykernel`. From the terminal, run `pip install ipykernel`.

You can create a test Jupyter Notebook document from the command palette by typing “Create: New Jupyter Notebook” and selecting the command. This will open up a code editor pane with a notebook that you can test.

Create a (free) Google Earth Engine (GEE) account

In order to code along during the Google Earth Engine lesson (Ch 15) on Thursday, you'll need to sign up for an account at <https://signup.earthengine.google.com>. Once submitted, you'll receive an email with some helpful links and a message that it may take a few days for your account to be up and running. **Please be sure to do this a few days ahead of needing to use GEE.**

! Important

GEE authentication (more on that in [Lesson 15](#)) uses [Cloud Projects](#). Some organizations control who can create Cloud Projects, which may prevent you from completing the authentication process. To circumvent authentication issues, we recommend creating your GEE account using a non-organizational account (e.g. a personal email account). Check out GEE's [authentication troubleshooting recommendations](#) if you continue to run into issues.

About this book

These written materials reflect the continuous development of learning materials at the Arctic Data Center and NCEAS to support individuals to understand, adopt, and apply ethical open science practices. In bringing these materials together we recognize that many individuals have contributed to their development. The primary authors are listed alphabetically in the citation below, with additional contributors recognized for their role in developing previous iterations of these or similar materials.

This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

Citation: S. Jeanette Clark, Matthew B. Jones, Samantha Csik, Carmen Galaz García, Bryce Mecum, Natasha Haycock-ChavezDaphne Virlar-Knight. 2022. Scalable and Computationally Reproducible Approaches to Arctic Research.

Additional contributors: Amber E. Budden, Noor Johnson

This is a Quarto book. To learn more about Quarto books visit <https://quarto.org/docs/books>.

1 Welcome and Overview



This course is one of three that we are currently offering, covering fundamentals of open data sharing, reproducible research, ethical data use and reuse, and scalable computing for reusing large data sets.



- Welcome and Introductions
- Mission and structure of the Arctic Data Center
- Plan for the Scalable Computing Course

1.1 Arctic Data Center Overview

The Arctic Data Center is the primary data and software repository for the Arctic section of National Science Foundation's Office of Polar Programs (NSF OPP).

We're best known in the research community as a data archive – researchers upload their data to preserve it for the future and make it available for re-use. This isn't the end of that data's life, though. These data can then be downloaded for different

analyses or synthesis projects. In addition to being a data discovery portal, we also offer top-notch tools, support services, and training opportunities. We also provide data rescue services.



Data Archive



Data Discovery Portal



Tools and Infrastructure



Support Services



Training and Outreach

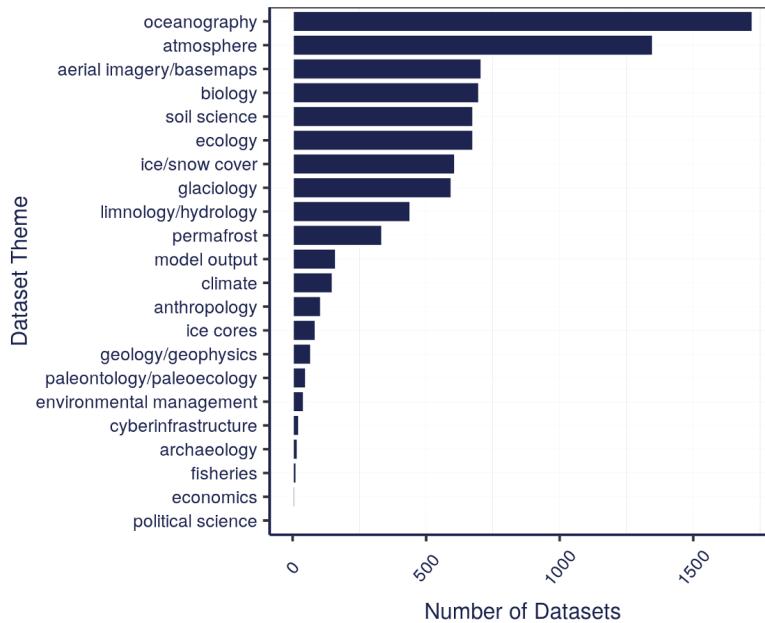


Data Rescue

NSF has long had a commitment to data reuse and sharing. Since our start in 2016, we've grown substantially – from that original 4 TB of data from ACADIS to now over 71 TB. In 2021 alone, we saw 16% growth in dataset count, and about 30% growth in data volume. This increase has come from advances in tools – both ours and of the scientific community, plus active community outreach and a strong culture of data preservation from NSF and from researchers. We plan to add more storage capacity in the coming months, as researchers are coming to us with datasets in the terabytes, and we're excited to preserve these research products in our archive. We're projecting our growth to be around several hundred TB this year, which has a big impact on processing time. Give us a heads up if you're planning on having larger submissions so that we can work with you and be prepared for a large influx of data.



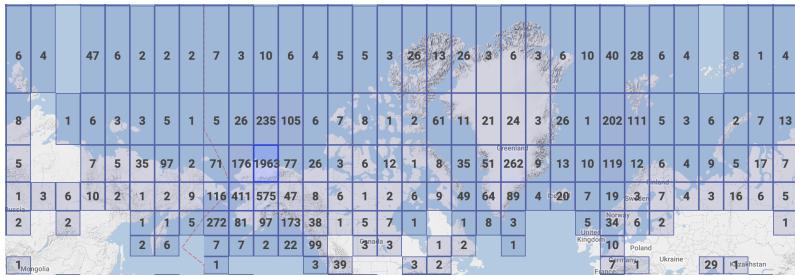
The data that we have in the Arctic Data Center comes from a wide variety of disciplines. These different programs within NSF all have different focuses – the Arctic Observing Network supports scientific and community-based observations of biodiversity, ecosystems, human societies, land, ice, marine and freshwater systems, and the atmosphere as well as their social, natural, and/or physical environments, so that encompasses a lot right there in just that one program. We're also working on a way right now to classify the datasets by discipline, so keep an eye out for that coming soon.



Along with that diversity of disciplines comes a diversity of file types. The most common file type we have are image files in four different file types. Probably less than 200-300 of the datasets have the majority of those images – we have some large datasets that have image and/or audio files from drones. Most of those 6600+ datasets are tabular datasets. There's a large diversity of data files, though, whether you want to look at remote sensing images, listen to passive acoustic audio files, or run applications – or something else entirely. We also cover a long period of time, at least by human standards. The data represented in our repository spans across centuries.



We also have data that spans the entire Arctic, as well as the sub-Arctic, regions.



1.1.1 Data Discovery Portal

To browse the data catalog, navigate to arcticdata.io. Go to the top of the page and under data, go to search. Right now, you're looking at the whole catalog. You can narrow your search down by the map area, a general search, or searching by an attribute.

Clicking on a dataset brings you to this page. You have the option to download all the files by clicking the green “Download All” button, which will zip together all the files in the dataset to your Downloads folder. You can also pick and choose to download just specific files.



All the raw data is in open formats to make it easily accessible and compliant with **FAIR** principles – for example, tabular documents are in .csv (comma separated values) rather than Excel documents.

The metrics at the top give info about the number of citations with this data, the number of downloads, and the number of views. This is what it looks like when you click on the Downloads tab for more information.



Scroll down for more info about the dataset – abstract, keywords. Then you'll see more info about the data itself. This

shows the data with a description, as well as info about the attributes (or variables or parameters) that were measured. The green check mark indicates that those attributes have been annotated, which means the measurements have a precise definition to them. Scrolling further, we also see who collected the data, where they collected it, and when they collected it, as well as any funding information like a grant number. For biological data, there is the option to add taxa.

1.1.2 Tools and Infrastructure

Across all our services and partnership, we are strongly aligned with the community principles of making data FAIR (Findable, Accessible, Interoperable and Reusable).



We have a number of tools available to submitters and researchers who are there to download data. We also partner with other organizations, like [Make Data Count](#) and [DataONE](#), and leverage those partnerships to create a better data experience.

The image displays several screenshots illustrating data management tools and partnerships. At the top left is a screenshot of a 'Metadata Quality Report' showing a progress bar and some text. To its right is a 'Downloads' panel with statistics: 3 views, 852 downloads, 274 downloads, and 209 downloads. Below these are two more screenshots: one for 'Add a Data Set' and another for 'DataONE'. To the right of these screenshots are the logos for 'MAKE DATA COUNT' and 'DataONE'.

One of those tools is provenance tracking. With provenance tracking, users of the Arctic Data Center can see exactly what datasets led to what product, using the particular script that the researcher ran.



Another tool are our Metadata Quality Checks. We know that data quality is important for researchers to find datasets and to have trust in them to use them for another analysis. For every submitted dataset, the metadata is run through a quality check to increase the completeness of submitted metadata records. These checks are seen by the submitter as well as are available to those that view the data, which helps to increase knowledge of how complete their metadata is before submission. That way, the metadata that is uploaded to the Arctic Data Center is as complete as possible, and close to following the guideline of being understandable to any reasonable scientist.



1.1.3 Support Services

Metadata quality checks are the automatic way that we ensure quality of data in the repository, but the real quality and curation support is done by our curation team. The process by which data gets into the Arctic Data Center is iterative, meaning that our team works with the submitter to ensure good quality and completeness of data. When a submitter submits data, our team gets a notification and begins to evaluate the data for upload. They then go in and format it for input into the catalog, communicating back and forth with the researcher if anything is incomplete or not communicated well. This process can take anywhere from a few days or a few weeks, depending on the size of the dataset and how quickly the researcher gets back to us. Once that process has been completed, the dataset is published with a DOI (digital object identifier).



1.1.4 Training and Outreach

In addition to the tools and support services, we also interact with the community via trainings like this one and outreach events. We run workshops at conferences like the American Geophysical Union, Arctic Science Summit Week and others. We also run an intern and fellows program, and webinars with different organizations. We're invested in helping the Arctic science community learn reproducible techniques, since it facilitates a more open culture of data sharing and reuse.



We strive to keep our fingers on the pulse of what researchers like yourselves are looking for in terms of support. We're active on [Twitter](#) to share Arctic updates, data science updates, and specifically Arctic Data Center updates, but we're also happy to feature new papers or successes that you all have had with working with the data. We can also take data science questions if you're running into those in the course of your research, or how to make a quality data management plan. Follow us on Twitter and interact with us – we love to be involved in your research as it's happening as well as after it's completed.



1.1.5 Data Rescue

We also run data rescue operations. We digitized Autin Post's collection of glacier photos that were taken from 1964 to 1997. There were 100,000+ files and almost 5 TB of data to ingest, and we reconstructed flight paths, digitized the images of his notes, and documented image metadata, including the camera specifications.



*Meares Glacier, Prince William Sound, AK
61.187448, -147.457573, taken from 18,000'
December 3, 1995, Roll 3, Frame 110
doi:10.18739/A2FF6Z (NAGAP_95V3_110.jpg)*

1.1.6 Who Must Submit

Projects that have to submit their data include all Arctic Research Opportunities through the NSF Office of Polar Programs. That data has to be uploaded within two years of collection. The Arctic Observing Network has a shorter timeline – their data products must be uploaded within 6 months of collection. Additionally, we have social science data, though that data often has special exceptions due to sensitive human subjects data. At the very least, the metadata has to be deposited with us.

Arctic Research Opportunities (ARC)

- Complete metadata and all appropriate data and derived products
- Within 2 years of collection or before the end of the award, whichever comes first

ARC Arctic Observation Network (AON)

- Complete metadata and all data
- Real-time data made public immediately
- Within 6 months of collection

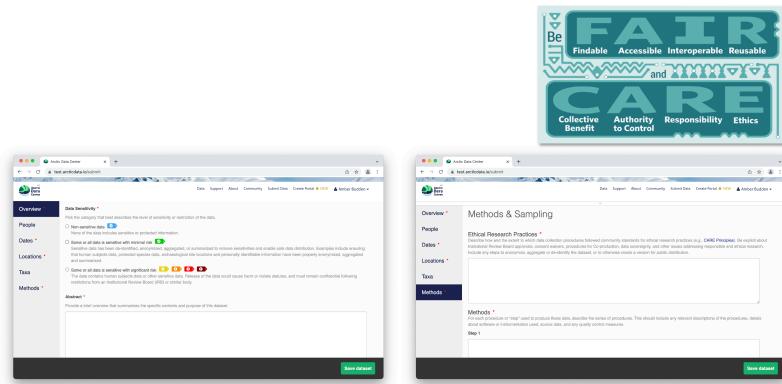
Arctic Social Sciences Program (ASSP)

- NSF policies include special exceptions for ASSP and other awards that contain sensitive data

- Human subjects, governed by an Institutional Review Board, ethically or legally sensitive, at risk of decontextualization
- Metadata record that documents non-sensitive aspects of the project and data
 - Title, Contact information, Abstract, Methods

For more complete information see our “Who Must Submit” [webpage](#)

Recognizing the importance of sensitive data handling and of ethical treatment of all data, the Arctic Data Center submission system provides the opportunity for researchers to document the ethical treatment of data and how collection is aligned with community principles (such as the CARE principles). Submitters may also tag the metadata according to community developed data sensitivity tags. We will go over these features in more detail shortly.



1.1.7 Supporting data reuse

As the Arctic Data Center has grown in size, we have envisioned new capabilities to support efficient reuse of these valuable data. We are developing new tools both within the Arctic Data Center, and through collaborative projects with other researchers in the community. In house, some of the new features we plan to support include:

- Efficient submission and access to multi-Terabyte datasets
- Data Quality assessment services
- Automated workflows for building derived data products

As a concrete example, we are collaborating with the [Permafrost Discovery Gateway](#) to build new capacity for computing and visualization at sub-meter spatial resolution across the global Arctic. The services we are building support workflows for pre-processing massive image data collections to prepare them for modeling, as well as automating machine learning and other computations across large data on distributed compute clusters, and finally post-processing the results to create viable pathways for results mapping and visualization at global scales.



1.1.8 Summary

All the above information can be found on our website or if you need help, ask our support team at support@arcticdata.io or tweet us @arcticdatactr!

6,958 Datasets
71.4 TB

2,584 Creators
16,162 Users

1.2 Scalable Computing Topics

The overall objective of this course is to facilitate effective analysis and modeling of the massive data resources that we curate for the Arctic research community. We assume a baseline proficiency, and plan to build on that to better support time-consuming computations with multi-terabyte datasets.

Activities will include:

- *Technical tutorials* with a lot of hands-on time in Python
- *Lectures and discussions* on core concepts for scalable computing
- *Semi-structured group projects* to practice and cement key skills

Key topics will include:

- Arctic Data Center services and tools
- Overview and review of core computing concepts
- Fundamentals of concurrent programming
 - concurrent.futures
 - Parsl
 - Dask
- Fundamentals of working with big data and imagery
- Good practices in research software design for efficiency, reproducibility, and reuse
- Fundamentals of cloud computing
 - Docker, containers, reproducibility, and more

As this is a lot to cover in a short survey course, we will break it down to fundamentals daily, and provide ample time for hands-on practice.

2 Remote Computing

- Understand the basic architecture of computer networks
- Learn how to connect to a remote computer via a shell
- Become familiarized with Bash Shell programming to navigate your computer's file system, manipulate files and directories, and automate processes

2.1 Introduction

Scientific synthesis and our ability to effectively and efficiently work with big data depends on the use of computers and the internet. Working on a personal computer may be sufficient for many tasks, but as data get larger and analyses more computationally intensive, scientists often find themselves needing more computing resources than they have available locally. Remote computing, or the process of connecting to a computer(s) in another location via a network link is becoming more and more common in overcoming big data challenges.

In this lesson, we'll learn about the architecture of computer networks and explore some of the different remote computing configurations that you may encounter, we'll learn how to securely connect to a remote computer via a shell, and we'll become familiarized with using Bash Shell to efficiently manipulate files and directories. We will begin working in the [VS Code](#) IDE (integrated development environment), which is a versatile code editor that supports many different languages.

2.2 Servers & Networking

Remote computing typically involves communication between two or more “host” computers. Host computers connect via

networking equipment and can send messages to each other over communication protocols (aka an [Internet Protocol](#), or IP). Host computers can take the role of **client** or **server**, where servers share their resources with the client. Importantly, these client and server roles are not inherent properties of a host (i.e. the same machine can play either role).

- **Client:** the host computer *initiating* a request
- **Server:** the host computer *responding* to a request

Fig 1. Examples of different remote computing configurations. (a) A client uses secure shell protocol (SSH) to login/connect to a server over the internet. (b) A client uses SSH to login/connect to a computing cluster (i.e. a set of computers (nodes) that work together so that they can be viewed as a single system) over the internet. In this example, servers A - I are each nodes on this single cluster. The connection is first made through a gateway node (i.e. a computer that routes traffic from one network to another). (c) A client uses SSH to login/connect to a computing cluster where each node is a virtual machine (VM). In this example, the cluster comprises three servers (A, B, and C). VM1 (i.e. node 1) runs on server A while VM4 runs on server B, etc. The connection is first made through a gateway node.

2.3 IP addressing

Hosts are assigned a **unique numerical address** used for all communication and routing called an [Internet Protocol Address \(IP Address\)](#). They look something like this: **128.111.220.7**. Each IP Address can be used to communicate over various “[ports](#)”, which allow multiple applications to communicate with a host without mixing up traffic.

i Port numbers are divided into three ranges:

1. *well-known ports*, range from 0 through 1023 and are reserved for the most commonly used services (see table below for examples of some well-known port

- numbers)
2. *registered ports*, range from 1024 through 49151 and are not assigned or controlled, but can be registered (e.g. by a vendor for use with their own server application) to prevent duplication
 3. *dynamic ports*, range from 49152 through 65535 and are not assigned, controlled, or registered but may instead be used as temporary or private ports

well-known port	assignment
20, 21	File Transfer Protocol (FTP), for transferring files between a client & server
22	secure shell (SSH), to create secure network connections
53	Domain Name System (DNS) service, to match domain names to IP addresses
80	Hypertext Transfer Protocol (HTTP), used in the World Wide Web
443	HTTP Secure (HTTPS), an encrypted version of HTTP

Because IP addresses can be difficult to remember, they are also assigned **hostnames**, which are handled through the global **Domain Name System (DNS)**. Clients first look up a hostname in the DNS to find the IP address, then open a connection to the IP address.

i In order to connect to remote servers, computing clusters, virtual machines, etc., you need to know their IP address (or hostname)

A couple important ones:

1. Throughout this course, we'll be working on a server with the hostname, **included-crab** and IP address, 128.111.85.28 (in just a little bit, we'll learn how to connect to **included-crab** using SSH)

- localhost is a hostname that refers to your local computer and is assigned the IP address 127.0.0.1 – the concept of localhost is important for tasks such as website testing, and is also important to understand when provisioning local execution resources (e.g. we'll practice this during the [section 6 exercise](#) when working with `Parsl`.)

2.4 Bash Shell Programming

What is a shell? From [Wikipedia](#):

“a computer program which exposes an operating system’s services to a human user or other programs. In general, operating system shells use either a command-line interface (CLI) or graphical user interface (GUI), depending on a computer’s role and particular operation.”

What is Bash? Bash, or Bourne-again Shell, is a command line tool (language) commonly used to manipulate files and directories. Accessing and using bash is slightly different depending on what type of machine you work on:

- **Mac:** bash via the [Terminal](#), which comes ready-to-use with all Macs and Linux machines
- **Windows:** running bash depends on which version of Windows you have – newer versions may ship with bash or may require a separate install (e.g. [Windows Subsystem for Linux \(WSL\)](#) or [Git Bash](#)), however there are a number of different (non-bash) shell options as well (they all vary slightly; e.g. [PowerShell](#), [Command Prompt](#)).

i Note

Mac users may have to switch from [Z Shell](#), or zsh, to bash. Use the command `exec bash` to switch your default shell to bash (or `exec zsh` to switch back).

2.4.1 Some commonly used (and very helpful) bash commands:

Below are just a few bash commands that you're likely to use. Some may be extended with options (more on that in the next section) or even piped together (i.e. where the output of one command gets sent to the next command, using the | operator). You can also find some nice bash cheat sheets online, like [this one](#). Alternatively, the [Bash Reference Manual](#) has *all* the content you need, albeit a bit dense.

bash command	what it does
<code>pwd</code>	print your current working directory
<code>cd</code>	change directory
<code>ls</code>	list contents of a directory
<code>tree</code>	display the contents of a directory in the form of a tree structure (not installed by default)
<code>echo</code>	print text that is passed in as an argument
<code>mv</code>	move or rename a file
<code>cp</code>	copy a file(s) or directory(ies)
<code>touch</code>	create a new empty file
<code>mkdir</code>	create a new directory
<code>rm/rmdir</code>	remove a file/ empty directory (be careful – there is no “trash” folder!)
<code>grep</code>	searches a given file(s) for lines containing a match to a given pattern list
<code>awk</code>	a text processing language that can be used in shell scripts or at a shell prompt for actions like pattern matching, printing specified fields, etc.
<code>sed</code>	stands for Stream Editor; a versatile command for editing files
<code>cut</code>	extract a specific portion of text in a file
<code>join</code>	join two files based on a key field present in both
<code>top, htop</code>	view running processes in a Linux system (press Q to quit)

2.4.2 General command syntax

Bash commands are typically written as: `command [options] [arguments]` where the command must be an executable on your PATH and where `options` (settings that change the shell and/or script behavior) take one of two forms: **short form** (e.g. `command -option-abbrev`) or **long form** (e.g. `command --option-name` or `command -o option-name`). An example:

```
# the `ls` command lists the files in a directory
ls file/path/to/directory

# adding on the `--a` or `--all` option lists all files (including hidden files) in a directory
ls -a file/path/to/directory # short form
ls --all file/path/to/directory # long form
ls -o all file/path/to/directory # long form
```

2.4.3 Some useful keyboard shortcuts

It can sometimes feel messy working on the command line. These keyboard shortcuts can make it a little easier:

- `Ctrl + L`: clear your terminal window
- `Ctrl + U`: delete the current line
- `Ctrl + C`: abort a command
- up & down arrow keys: recall previously executed commands in chronological order
- `TAB` key: autocompletion

2.5 Connecting to a remote computer via a shell

In addition to navigating your computer/manipulating your files, you can also use a shell to gain accesss to and remotely control other computers. To do so, you'll need the following:

- a remote computer (e.g. server) which is turned on

- client and server ssh clients installed/enabled
- the IP address or name of the remote computer
- the necessary permissions to access the remote computer

Secure Shell, or SSH, is a network communication protocol that is often used for securely connecting to and running shell commands on a remote host, tremendously simplifying remote computing.

2.6 Git via a shell

[Git](#), a popular version control system and command line tool can be accessed via a shell. While there are lots of graphical user interfaces (GUIs) that facilitate version control with Git, they often only implement a small subset of Git's most-used functionality. By interacting with Git via the command line, you have access to *all* Git commands. While all-things Git is outside the scope of this workshop, we will use some basic Git commands in the shell to clone GitHub (remote) repositories to the server and save/store our changes to files. A few important Git commands:

Git command	what it does
<code>git clone</code>	create a copy (clone) of repository in a new directory in a different location
<code>git add</code>	add a change in the working directory to the staging area
<code>git commit</code>	record a snapshot of a repository; the <code>-m</code> option adds a commit message
<code>git push</code>	send commits from a local repository to a remote repository
<code>git fetch</code>	downloads contents (e.g. files, commits, refs) from a remote repo to a local repo
<code>git pull</code>	fetches contents of a remote repo and merges changes into the local repo

2.7 Let's practice!

We'll now use bash commands to do the following:

- connect to the server (**included-crab**) that we'll be working on for the remainder of this course
- navigate through directories on the server and add/change/manipulate files
- clone a GitHub repository to the server
- automate some of the above processes by writing a bash script

2.7.1 Exercise 1: Connect to a server using the ssh command (or using VS Code's command palette)

Let's connect to a remote computer (**included-crab**) and practice using some of above commands.

1. Launch a terminal in VS Code
 - There are two options to open a terminal window, if a terminal isn't already an open pane at the bottom of VS Code
 - a) Click on Terminal > New Terminal in top menu bar
 - b) Click on the + (dropdown menu) > bash in the bottom right corner

i Note

You don't *need* to use the VS Code terminal to ssh into a remote computer, but it's conveniently located in the same window as your code when working in the VS Code IDE.

2. Connect to a remote server

- You can choose to SSH into the server (included-crab.nceas.ucsb.edu) through (a) the command line by using the `ssh` command, or (b) through VS Code's command palette. If you prefer the latter, please refer back to the [Log in to the server section](#). To do so via the command line, use the `ssh` command followed by `yourusername@included-crab.nceas.ucsb.edu`. You'll be prompted to type/paste your password to complete the login. It should look something like this:

```
yourusername:~$ ssh yourusername@included-crab.nceas.ucsb.edu  
yourusername@included-crab.nceas.ucsb.edu's password:  
yourusername@included-crab:~$
```

! Important

You won't see anything appear as you type or paste your password – this is a security feature! Type or paste your password and press enter/return when done to finish connecting to the server.

i Note

To log out of the server, type `exit` – it should look something like this:

```
yourusername@included-crab.nceas.ucsb.edu:$ exit  
logout  
Connection to included-crab.nceas.ucsb.edu closed.  
(base) .....
```

2.7.2 Exercise 2: Practice using some common bash commands

1. Use the `pwd` command to print your current location, or working directory. You should be in your home directory on the server (e.g. `/home/yourusername`).
2. Use the `ls` command to list the contents (any files or

subdirectories) of your home directory

3. Use the `mkdir` command to create a new directory named `bash_practice`:

```
mkdir bash_practice
```

4. Use the `cd` command to move into your new `bash_practice` directory:

```
# move from /home/yourusername to home/yourusername/bash_practice
cd bash_practice
```

- To move *up* a directory level, use two dots, `..` :

```
# move from /home/yourusername/bash_practice back to /home/yourusername
$ cd ..
```

i Note

To quickly navigate back to your home directory from wherever you may be on your computer, use a tilde, `~` :

```
# e.g. to move from some subdirectory, /home/yourusername/Projects/project1/data, back to home
$ cd ~
```

```
# or use .. to back out three subdirectories
$ cd ../../..
```

5. Add some `.txt` files (`file1.txt`, `file2.txt`, `file3.txt`) to your `bash_practice` subdirectory using the `touch` command (**Note:** be sure to `cd` into `bash_practice` if you're not already there):

```
# add one file at a time
touch file1.txt
touch file2.txt
touch file3.txt
```

```
# or add all files simultanously like this:  
touch file{1..3}.txt
```

```
# or like this:  
touch file1.txt file2.txt file3.txt
```

6. You can also add other file types (e.g. .py, .csv, etc.)

```
touch mypython.py mycsv.csv
```

7. Print out all the .txt files in `bash_practice` using a wildcard, *:

```
ls *.txt
```

8. Count the number of .txt files in `bash_practice` by combining the `ls` and `wc` (word count) funtions using the pipe, |, operator:

```
# `wc` returns a word count (lines, words, chrs)  
# the `--l` option only returns the number of lines  
# use a pipe, `|`, to send the output from `ls *.txt` to `wc -l`  
ls *.txt | wc -l
```

9. Delete `mypython.py` using the `rm` command:

```
rm mypython.py
```

10. Create a new directory inside `bash_practice` called `data` and move `mycsv.csv` into it.

```
mkdir data  
mv mycsv.csv ~/bash_practice/data
```

```
# add the --interactive option (-i for short) to prevent a file from being overwritten by ac  
mv -i mycsv.csv ~/bash_practice/data
```

11. Use `mv` to rename `mycsv.csv` to `mydata.csv`

```
mv mycsv.csv mydata.csv
```

12. Add column headers `col1`, `col2`, `col3` to `mydata.csv`
using `echo` + the `>` operator

```
echo "col1, col2, col3" > mydata.csv
```

 Tip

You can check to see that `mydata.csv` was updated using [GNU nano](#), a text editor for the command line that comes preinstalled on Linux machines (you can edit your file in nano as well). To do so, use the `nano` command followed by the file you want to open/edit:

```
nano mydata.csv
```

To save and quit out of nano, use the `control + X` keyboard shortcut.

You can also create and open a file in nano in just one line of code. For example, running `nano hello_world.sh` is the same as creating the file first using `touch hello_world.sh`, then opening it with nano using `nano hello_world.sh`.

13. Append a row of data to `mydata.csv` using `echo` + the `>>` operator

```
# using `>` will overwrite the contents of an existing file; `>>` appends new information to the end of the file  
echo "1, 2, 3" >> mydata.csv
```

2.7.3 Exercise 3: Clone a GitHub repository to the server

IDEs commonly have helper buttons for cloning (i.e. creating a copy of) remote repositories to your local computer (or in this case, a server), but using git commands in a terminal can be just as easy. We can practice that now, following the steps below:

1. Go to the `scalable-computing-examples` repository on GitHub at <https://github.com/NCEAS/scalable-computing-examples> – this repo contains example files for you to edit and practice in throughout this course. Fork (make your own copy of the repository) this repo by clicking on the **Fork** button (top right corner of the repository's page).
2. Once forked, click on the green **Code** button (from the *forked* version of the GitHub repo) and copy the URL to your clipboard.
3. In the VS Code terminal, use the `git clone` command to create a copy of the `scalable-computing-examples` repository in the top level of your user directory (i.e. your home directory) on the server (**Note:** use `pwd` to check where you are; use `cd ~` to navigate back to your home directory if you find that you're somewhere else).

```
git clone <url-of-forked-repo>
```

4. You should now have a copy of the `scalable-computing-examples` repository to work on on the server. Use the `tree` command to see the structure of the repo (you need to be in the `scalable-computing-examples` directory for this to work) – there should be a subdirectory called `02-bash-babynames` that contains (i) a `README.MD` file, (ii) a `KEY.sh` file (this is a functioning bash script available for reference; we'll be recreating it together in the next exercise) and (iii) a `namesbystate` folder containing 51 `.TXT` files and a `StateReadMe.pdf` file with some metadata.

2.7.4 Bonus Exercise: Automate data processing with a Bash script

As we just demonstrated, we can use bash commands in the terminal to accomplish a variety of tasks like navigating our computer's directories, manipulating/creating/adding files,

and much more. However, writing a bash *script* allows us to gather and save our code for automated execution.

We just cloned the `scalable-computing-examples` GitHub repository to the server in [Exercise 3](#) above. This contains a `02-bash-babynames` folder with 51 .TXT files (one for each of the 50 US states + The District of Columbia), each with the top 1000 most popular baby names in that state. We're going to use some of the bash commands we learned in [Exercise 2](#) to concatenate all rows of data from these 51 files into a single `babynames_allstates.csv` file.

Let's begin by creating a simple bash script that when executed, will print out the message, "Hello, World!" This simple script will help us determine whether or not things are working as expected before writing some more complex (and interesting) code.

1. Open a terminal window and determine where you are by using the `pwd` command – we want to be in `scalable-computing-examples/02-bash-babynames`. If necessary, navigate here using the `cd` command.
2. Next, we'll create a shell script called `mybash.sh` using the `touch` command:

```
$ touch mybash.sh
```

3. There are a number of ways to edit a file or script – we'll use [Nano](#), a terminal-based text editor, as we did earlier. Open your `mybash.sh` with nano by running the following in your terminal:

```
$ nano mybash.sh
```

4. We can now start to write our script. Some important considerations:
 - Anything following a `#` will not be executed as code – these are useful for adding comments to your scripts

- The first line of a Bash script starts with a **shebang**, `#!`, followed by a path to the Bash interpreter – this is used to tell the operating system which interpreter to use to parse the rest of the file. There are two ways to use the shebang to set your interpreter (read up on the pros & cons of both methods on this [Stack Overflow post](#)):

```
# (option a): use the absolute path to the bash binary  
#!/bin/bash
```

```
# (option b): use the env utility to search for the bash executable in the user's $PATH env  
#!/usr/bin/env bash
```

5. We'll first specify our bash interpreter using the shebang, which indicates the start of our script. Then, we'll use the `echo` command, which when executed, will print whatever text is passed as an argument. Type the following into your script (which should be opened with nano), then save (Use the keyboard shortcut `control + X` to exit, then type `Y` when it asks if you'd like to save your work. Press `enter/return` to exit nano).

```
# specify bash as the interpreter  
#!/bin/bash  
  
# print "Hello, World!"  
echo "Hello, World!"
```

6. To execute your script, use the `bash` command followed by the name of your bash script (be sure that you're in the same working directory as your `mybash.sh` file or specify the file path to it). If successful, “Hello, World!” should be printed in your terminal window.

```
bash mybash.sh
```

7. Now let's write our script. Re-open your script in nano by running `nano mybash.sh`. Using what we practiced above and the hints below, write a bash script that does the following:

- prints the number of .TXT files in the `namesbystate` subdirectory
- prints the first 10 rows of data from the `CA.TXT` file (HINT: use the `head` command)
- prints the last 10 rows of data from the `CA.TXT` file (HINT: use the `tail` command)
- creates an empty `babynames_allstates.csv` file in the `namesbystate` subdirectory (this is where the concatenated data will be saved to)
- adds the column `names`, `state`, `gender`, `year`, `firstname`, `count`, in that order, to the `babynames_allstates.csv` file
- concatenates data from all .TXT files in the `namesbystate` subdirectory and appends those data to the `babynames_allstates.csv` file (HINT: use the `cat` command to concatenate files)

Here's a script outline to fill in (**Note:** The `echo` statements below are not necessary but can be included as progress indicators for when the bash script is executed – these also make it easier to diagnose where any errors occur during execution):

```
#!/bin/bash
echo "THIS IS THE START OF MY SCRIPT!"

echo "-----Verify that we have .TXT files for all 50 states + DC-----"
<add your code here>

echo "-----Printing head of CA.TXT-----"
<add your code here>

echo "-----Printing tail of CA.TXT-----"
<add your code here>

echo "-----Creating empty .csv file to concatenate all data-----"
<add your code here>

echo "-----Adding column headers to csv file-----"
<add your code here>

echo "-----Concatenating files-----"
```

```
<add your code here>
```

```
echo "DONE!"
```

💡 Answer

```
#!/bin/bash
echo "THIS IS THE START OF MY SCRIPT!"

echo "-----Verify that we have .TXT files for all 50 states + DC-----"
ls namesbystate/*.TXT | wc -l

echo "-----Printing head of CA.TXT-----"
head namesbystate/CA.TXT

echo "-----Printing tail of CA.TXT-----"
tail namesbystate/CA.TXT

echo "-----Creating empty .csv file to concatenate all data-----"
touch namesbystate/babynames_allstates.csv

echo "-----Adding column headers to csv file-----"
echo "state, gender, year, firstname, count" > namesbystate/babynames_allstates.csv

echo "-----Concatenating files-----"
cat namesbystate/*.TXT >> namesbystate/babynames_allstates.csv

echo "DONE!"
```

3 Python Programming on Clusters

- Basic Python review
- Using virtual environments
- Writing in Jupyter notebooks
- Writing functions in Python

3.1 Introduction

We've chosen to use VS Code in this training, in part, because it has great support for developing on remote machines. Hopefully, your VS Code setup went easily, and you were able to connect to our server `included-crab`. Once connected, the VS Code interface looks just like you were working locally, and connection to the server is seamless.

Other aspects of VS Code that we like: it supports all languages thanks to the extensive free extension library, it has built in version control integration, and it is highly flexible/configurable.

We will also be working quite a bit in Jupyter notebooks in this course. Notebooks are great ways to interleave rich text (mark-down formatted text, equations, images, links) and code in a way that a ‘literate analysis’ is generated. Although Jupyter notebooks are not substitutes for python scripts, they can be great communication tools, and can also be convenient for code development.

3.2 Connect to the server (if you aren't already connected)

To get set up for the course, let's connect to the server again. If you were able to work through the setup for the lesson without difficulty, follow these steps to connect:

- open VS Code
- open the command palette (Cmd + Shift + P)
- enter “Remote SSH: Connect to Host”
- select `included-crab`
- enter your password in the dialog box that pops up

Alternatively, you may see a popup window that says “**Cannot reconnect. Please reload the window**”. Choose the blue **Reload Window** button and enter your password, if prompted.

To open the `scalable-computing-examples` workspace, select File > Open Workspace from File. Then navigate to and select `~/scalable-computing-examples/scalable-computing-examples.code-workspace`. If you cannot find the `scalable-computing-examples` directory, in a terminal you can use the `ls` command to list out the contents of your home directory and ensure that you have a clone of the repository. If not, follow the instructions in [Exercise 3](#) in the [Remote Computing](#) session).

Once connected and in the workspace, use the `pwd` command in the terminal to make sure you're in your project directory (`/home/yourusername/scalable-computing-examples`).

3.3 Virtual Environments

When you install a python library, let's say `pandas`, via `pip`, unless you specify otherwise, `pip` will go out and grab the most recent version of the library, and install it somewhere on your system path (where, exactly, depends highly on how you install python originally, and other factors). That library is going to sit alongside every other python library you have ever installed (probably a lot of them) in your `site-packages` directory. This

might work okay, until you start a new project, have no idea what version of `pandas` you had installed, realize that version conflicts with a new library that you need, and then instead of writing code you are spending days trying to untangle the spaghetti mess of your python install, libraries and their versions. Sound familiar? Virtual environments help to solve this issue without making the all to common situation in the comic above even more complicated.

A virtual environment is a folder structure which creates a symbol-link (pointer) to all of the libraries that you specify into the folder. The three main components will be: the python distribution itself, its configuration, and a site-packages directory (where your libraries like `pandas` live). So the folder is a self contained directory of all the version-specific python software you need for your project.

Virtual environments are very helpful to create reproducible workflows, and we'll talk more about this concept of reproducible environments later in the course. Perhaps most importantly though, virtual environments also help you maintain your sanity when python programming. Because they are just folders, you can create and delete new ones at will, without worrying about bungling your underlying python setup.

In this course, we are going to use `virtualenv` as our tool to create and manage virtual environments. Other virtual environment tools used commonly are `conda` and `pipenv`. One reason we like using `virtualenv` is there is an extension to it called `virtualenvwrapper`, which provides easy to remember wrappers around common `virtualenv` operations that make creating, activating, and deactivating a virtual environment very easy.

First we will create a `.bash_profile` file to create variables that point to the install locations of python and `virtualenvwrapper`. `.bash_profile` is just a text file that contains bash commands that are run every time you start up a new terminal. Although setting up this file is not required to use `virtualenvwrapper`, it is convenient because it allows you to set up some reasonable defaults to the commands (meaning less typing, overall), and it makes sure that the package is available every time you start a new terminal.

3.3.0.1 Setup

- In VS Code, select ‘File > New Text File’
- Paste this text into the file:

```
export VIRTUALENWRAPPER_VIRTUAWORKON_HOME=$HOME/.virtualenvs
source /usr/share/virtualenvwrapper/virtualenvwrapper.sh
```

The first line points `virtualenvwrapper` to the directory where your virtual environments will be stored. We point it to a hidden directory (`.virtualenvs`) in your home directory. The last line sources a bash script that ships with `virtualenvwrapper`, which makes all of `virtualenvwrapper` commands available in your terminal session.

- Save the file in the top of your home directory as `.bash_profile`.
- Restart your terminal (Terminal > New Terminal)
- Check to make sure it was installed and configured correctly by running this in the terminal:

```
mkvirtualenv --version
```

It should return some content that looks like this (with more output, potentially).

```
virtualenv 20.13.0+ds from /usr/lib/python3/dist-packages/virtualenv/__init__.py
```

3.3.0.2 Course environment

Now we can create the virtual environment we will use for the course. In the terminal run:

```
mkvirtualenv -p python3.9 scomp
```

Here, we’ve specified explicitly which python version to use by using the `-p` flag, and the path to the python 3.9 installation on the server. After making a virtual environment, it will automatically be activated. You’ll see the name of the env you are

working in on the left side of your terminal prompt in parentheses. To deactivate your environment (like if you want to work on a different project), just run `deactivate`. To activate it again, run:

```
workon scomp
```

You can get a list of all available environments by just running:

```
workon
```

Now let's install the dependencies for this course into that environment. To install our libraries we'll use `pip`. As of Python 3.4, `pip` is automatically included with your python installation. `pip` is a package manager for python, and you might have used it already to install common python libraries like `pandas` or `numpy`. `pip` goes out to [PyPI](#), the Python Package Index, to download the code and put it in your `site-packages` directory. Note that on this shared server, your user directory will ahve a `site-packages` directory, in addition to one that our systems administrator manages as the root of the system.

```
pip install -r requirements.txt
```

3.3.0.3 Installing locally (optional)

`virtualenvwrapper` was already installed on the server we are working on. To install on your local computer, run:

```
pip3 install virtualenvwrapper
```

And then follow the instructions as described above, making sure that you have the correct paths set when you edit your `.bash_profile`.

3.4 Brief overview of python syntax

We'll very briefly go over some basic python syntax and the base variable types. First, open a python script. From the File menu, select New File, type “python”, then save it as ‘python-intro.py’ in the top level of your directory.

In your file, assign a value to a variable using = and print the result.

```
x = 4  
print(x)
```

4

To run this code in python we can:

- execute `python python-intro.py` in the terminal
- click the Play button in the upper right hand corner of the file editor
- right click any line and select: “Run to line in interactive terminal”

In that interactive window you can then run python code interactively, which is what we'll use for the next bit of exploring data types.

There are 5 standard data types in python

- Number (int, long, float, complex)
- String
- List
- Tuple
- Dictionary

We already saw a number type, here is a string:

```
str = 'Hello World!'  
print(str)
```

Hello World!

Lists in python are very versatile, and are created using square brackets []. Items in a list can be of different data types.

```
list = [100, 50, -20, 'text']
print(list)
```

```
[100, 50, -20, 'text']
```

You can access items in a list by index using the square brackets. Note indexing starts with 0 in python. The slice operator enables you to easily access a portion of the list without needing to specify every index.

```
list[0] # print first element
list[1:3] # print 2nd until 4th elements
list[:2] # print first until the 3rd
list[2:] # print last elements from 3rd
```

```
100
```

```
[50, -20]
```

```
[100, 50]
```

```
[-20, 'text']
```

The + and * operators work on lists by creating a new list using either concatenation (+) or repetition (*).

```
list2 = ['more', 'things']
list + list2
list * 3
```

```
[100, 50, -20, 'text', 'more', 'things']
```

```
[100, 50, -20, 'text', 100, 50, -20, 'text', 100, 50, -20, 'text']
```

Tuples are similar to lists, except the values cannot be changed in place. They are constructed with parentheses.

```
tuple = ('a', 'b', 'c', 'd')
tuple[0]
tuple * 3
tuple + tuple

'a'

('a', 'b', 'c', 'd', 'a', 'b', 'c', 'd', 'a', 'b', 'c', 'd')

('a', 'b', 'c', 'd', 'a', 'b', 'c', 'd')
```

Observe the difference when we try to change the first value. It works for a list:

```
list[0] = 'new value'
list

['new value', 50, -20, 'text']
```

...and errors for a tuple.

```
tuple[0] = 'new value'

TypeError: 'tuple' object does not support item assignment
```

Dictionaries consist of key-value pairs, and are created using the syntax `{key: value}`. Keys are usually numbers or strings, and values can be any data type.

```
dict = {'name': ['Jeanette', 'Matt'],
        'location': ['Tucson', 'Juneau']}

dict['name']
```

```
dict.keys()  
  
['Jeanette', 'Matt']  
  
dict_keys(['name', 'location'])
```

To determine the type of an object, you can use the `type()` method.

```
type(list)  
type(tuple)  
type(dict)
```

list

tuple

dict

3.5 Jupyter notebooks

To create a new notebook, from the file menu select File > New File > Jupyter Notebook. Go ahead and save this notebook at the top level of your `scalable-computing-examples` directory.

At the top of your notebook, add a first level header using a single hash. Practice some markdown text by creating:

- a list
- **bold** text
- a link

Use the [Markdown cheat sheet](#) if needed.

You can click the plus button below any chunk to add a chunk of either markdown or python.

3.5.1 Load libraries

In your first code chunk, lets load in some modules. We'll use `pandas`, `numpy`, `matplotlib.pyplot`, `requests`, `skimpy`, and `exists` from `os.path`.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import urllib
import skimpy
import os
```

A note on style: There are a few ways to construct import statements. The above code uses three of the most common:

```
import module
import module as m
from module import function
```

The first way of importing will make the module a function comes from more explicitly clear, and is the simplest. However for very long module names, or ones that are used very frequently (like `pandas`, `numpy`, and `matplotlib.plot`), the code in the notebook will be more cluttered with constant calls to longer module names. So `module.function()` instead is written as `m.function()`

The second way of importing a module is a good style to use in cases where modules are used frequently, or have extremely long names. If you import every single module with a short name, however, you might have a hard time remembering which modules are named what, and it might be more confusing for others trying to read your code. Many of the most commonly used libraries for python data science have community-driven styling for how they are abbreviated in import statements, and these community norms are generally best followed.

Finally, the last way to import a single object from a module can be helpful if you only need that one piece from a larger module, but again, like the first case, results in less explicit code and

therefore runs the risk of your or someone else misremembering the usage and source.

3.5.2 Read in a csv

Create a new code chunk that will download the csv that we are going to use for this tutorial.

- Navigate to Rohi Muthyala, Åsa Rennermalm, Sasha Leidman, Matthew Cooper, Sarah Cooley, et al. 2022. 62 days of Supraglacial streamflow from June-August, 2016 over southwest Greenland. Arctic Data Center. doi:10.18739/A2XW47X5F.
- Right click the download button for ‘Discharge_timeseries.csv’
- Click ‘copy link address’

Create a variable called URL and assign it the link copied to your clipboard. Then use `urllib.request.urlretrieve` to download the file, and `open` to write it to disk, to a directory called `data/`. We’ll write this bundled in an `if` statement so that we only download the file if it doesn’t yet exist. First, we create the directory if it doesn’t exist:

```
if not os.path.exists ('data/'):
    os.mkdir('data/')

if not os.path.exists('data/discharge_timeseries.csv'):

    url = 'https://arcticdata.io/metacat/d1/mn/v2/object/urn%3Auuid%3Ae248467d-e1f9-4a32

    msg = urllib.request.urlretrieve(url, 'data/discharge_timeseries.csv')
```

Now we can read in the data from the file.

```
df = pd.read_csv('data/discharge_timeseries.csv')
df.head()
```

```
/opt/hostedtoolcache/Python/3.9.16/x64/lib/python3.9/site-packages/IPython/core/formatters.py:
```

```
In future versions `DataFrame.to_latex` is expected to utilise the base implementation of `Sty
```

	Date	Total Pressure [m]	Air Pressure [m]	Stage [m]	Discharge [m ³ /s]	temperature [degree]
0	6/13/2016 0:00	9.816	9.609775	0.206225	0.083531	
1	6/13/2016 0:05	9.810	9.609715	0.200285	0.077785	
2	6/13/2016 0:10	9.804	9.609656	0.194344	0.072278	
3	6/13/2016 0:15	9.800	9.609596	0.190404	0.068756	
4	6/13/2016 0:20	9.793	9.609537	0.183463	0.062804	

The column names are a bit messy so we can use `clean_columns` from `skimpy` to make them cleaner for programming very quickly. We can also use the `skim` function to get a quick summary of the data.

We can see that the `date` column is classed as a string, and not a date, so let's fix that.

If we wanted to calculate the daily mean flow (as opposed to the flow every 5 minutes), we need to:

- create a new column with only the date
- group by that variable
- summarize over it by taking the mean

First we should probably rename our existing date/time column to prevent from getting confused.

```
clean_df = clean_df.rename(columns = {'date': 'datetime'})
```

Now create the new date column

```
clean_df['date'] = clean_df['datetime'].dt.date
```

Finally, we use group by to split the data into groups according to the date. We can then apply the `mean` method to calculate the mean value across all of the columns. Note that there are other methods you can use to calculate different statistics across different columns (eg: `clean_df.groupby('date').agg({'discharge_m_3_s': 'max'})`).

```
daily_flow = clean_df.groupby('date', as_index = False).mean()
```

- create a simple plot

```
clean_df = skimpy.clean_columns(df)
skimpy.skim(clean_df)
```

6 column names have been cleaned

skimpy summary									
Data Summary		Data Types							
dataframe	Values	Column Type	Count						
Number of rows	17856	float64	5						
Number of columns	6	string	1						
number									
column_name	NA	NA %	mean	sd	p0	p25	p75	p100	p00
total_pressure_m	0	0	9.9	0.12	9.6	9.8	10	10	10
air_pressure_m	0	0	9.6	0.06	9.5	9.6	9.7	9.7	9.7
stage_m	0	0	0.28	0.12	0.00056	0.17	0.37	0.37	0.37
discharge_m_3_s	0	0	0.22	0.19	4.7e-08	0.055	0.35	0.35	0.35
temperature_degrees_	8	0.045	-0.034	0.053	-0.1	-0.1	0	0	0
string									
column_name	NA	NA %	words per row				total words		
date	0	0					2		
End									

```
clean_df['date'] = pd.to_datetime(clean_df['date'])
skimpy.skim(clean_df)
```

skimpy summary								
Data Summary		Data Types						
dataframe	Values	Column Type	Count					
Number of rows	17856	float64	5					
Number of columns	6	datetime64	1					
number								
column_name	NA	NA %	mean	sd	p0	p25	p75	p100
total_pressure_m	0	0	9.9	0.12	9.6	9.8	10	10
air_pressure_m	0	0	9.6	0.06	9.5	9.6	9.7	9.7
stage_m	0	0	0.28	0.12	0.00056	0.17	0.37	0.37
discharge_m_3_s	0	0	0.22	0.19	4.7e-08	0.055	0.35	0.35
temperature_degrees_	8	0.045	-0.034	0.053	-0.1	-0.1	0	0
datetime								
column_name	NA	NA %	first	last				
date	0	0	2016-06-13	2016-08-13 23:55:00				
End								

```
var = 'discharge_m_3_s'
var_labs = {'discharge_m_3_s': 'Total Discharge'}

fig, ax = plt.subplots(figsize=(7, 7))
plt.plot(daily_flow['date'], daily_flow[var])
plt.xticks(rotation = 45)
ax.set_ylabel(var_labs.get('discharge_m_3_s'))

(array([16967., 16974., 16983., 16990., 16997., 17004., 17014., 17021.,
       17028.]),
[Text(0, 0, ''),
 Text(0, 0, '')])

Text(0, 0.5, 'Total Discharge')
```



3.6 Functions

The plot we made above is great, but what if we wanted to make it for each variable? We could copy paste it and replace some things, but this violates a core tenet of programming: Don't Repeat Yourself! Instead, we'll create a function called `myplot` that accepts the data frame and variable as arguments.

- create `myplot.py`

```
import matplotlib.pyplot as plt

def myplot(df, var):

    var_labs = {'discharge_m_3_s': 'Total Discharge (m^3/s)',  

                'total_pressure_m': 'Total Pressure (m)',
```

```

    'air_pressure_m': 'Air Pressure (m)',
    'stage_m': 'Stage (m)',
    'temperature_degrees_c': 'Temperature (C)'}

fig, ax = plt.subplots(figsize=(7, 7))
plt.plot(df['date'], df[var])
plt.xticks(rotation = 45)
ax.set_ylabel(var_labs.get(var))

```

- load myplot into jupyter notebook (`from myplot import myplot`)
- add libraries
- replace old plot method with new function

```
myplot(daily_flow, 'temperature_degrees_c')
```



We'll have more on functions in the software design sections.

3.7 Summary

In this lesson we learned all about virtual environments, how to use them, and why. We got our environments set up for the course, did a brief python syntax review, and then jumped into Jupyter notebooks, pandas, and writing functions.

4 Pleasingly Parallel Programming

- Understand what parallel computing is and when it may be useful
- Understand how parallelism can work
- Review sequential loops and map functions
- Build a parallel program using `concurrent.futures`
- Build a parallel program using `parsl`
- Understand Thread Pools and Process pools

4.1 Introduction

Processing large amounts of data with complex models can be time consuming. New types of sensing means the scale of data collection today is massive. And modeled outputs can be large as well. For example, here's a 2 TB (that's Terabyte) set of modeled output data from [Ofir Levy et al. 2016](#) that models 15 environmental variables at hourly time scales for hundreds of years across a regular grid spanning a good chunk of North America:

There are over 400,000 individual netCDF files in the [Levy et al. microclimate data set](#). Processing them would benefit massively from parallelization.

Alternatively, think of remote sensing data. Processing airborne hyperspectral data can involve processing each of hundreds of bands of data for each image in a flight path that is repeated many times over months and years.



Figure 4.1: Levy et al. 2016. doi:10.5063/F1Z899CZ

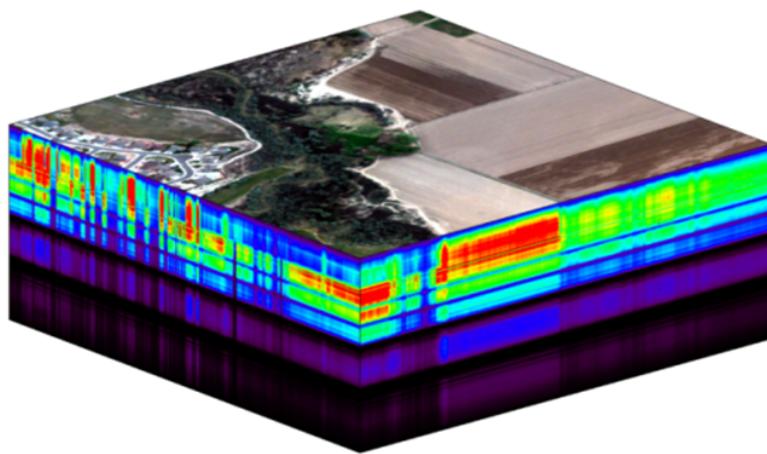


Figure 4.2: NEON Data Cube

4.2 Why parallelism?

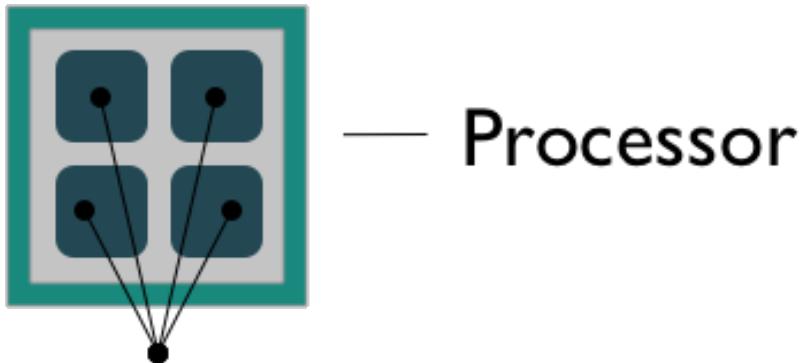
Much R code runs fast and fine on a single processor. But at times, computations can be:

- **cpu-bound**: Take too much cpu time
- **memory-bound**: Take too much memory
- **I/O-bound**: Take too much time to read/write from disk
- **network-bound**: Take too much time to transfer

To help with **cpu-bound** computations, one can take advantage of modern processor architectures that provide multiple cores on a single processor, and thereby enable multiple computations to take place at the same time. In addition, some machines ship with multiple processors, allowing large computations to occur across the entire set of those processors. Plus, these machines also have large amounts of memory to avoid **memory-bound** computing jobs.

4.3 Processors (CPUs), Cores, and Threads

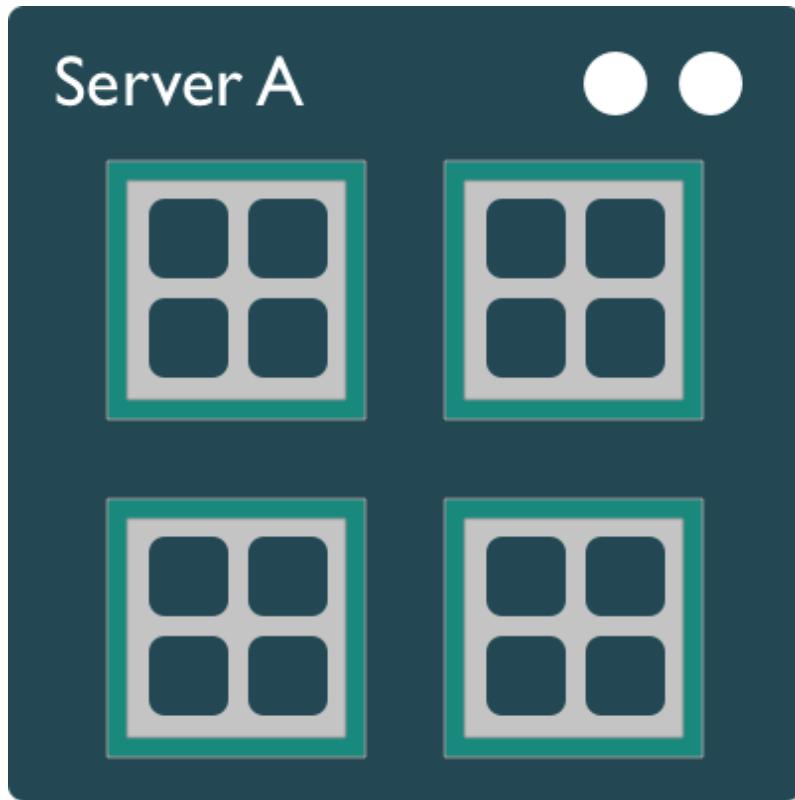
A modern CPU (Central Processing Unit) is at the heart of every computer. While traditional computers had a single CPU, modern computers can ship with multiple processors, each of which in turn can contain multiple cores. These processors and cores are available to perform computations. But, just what's the difference between processors and cores? A computer with one processor may still have 4 cores (quad-core), allowing 4 (or possibly more) computations to be executed at the same time.



Cores

- **Microprocessor:** an integrated circuit that contains the data processing logic and control for a computer.
- **Multi-core processor:** a microprocessor containing multiple processing units (cores) on a single integrated circuit. Each core in a multi-core processor can execute program instructions at the same time.
- **Process:** an instance of a computer program (including instructions, memory, and other resources) that is executed on a microprocessor.
- **Thread:** a thread of execution is the smallest sequence of program instructions that can be executed independently, and is typically a component of a process. The threads in a process can be executed concurrently and typically share the same memory space. They are faster to create than a process.
- **Cluster:** a set of multiple, physically distinct computing systems, each with its own microprocessors, memory, and storage resources, connected together by a (fast) network that allows the nodes to be viewed as a single system.

A typical modern computer has multiple cores, ranging from one or two in laptops to thousands in high performance compute clusters. Here we show four quad-core processors for a total of 16 cores in this machine.



You can think of this as allowing 16 computations to happen at the same time. Theoretically, your computation would take 1/16 of the time (but only theoretically, more on that later).

Historically, many languages only utilized one processor, which makes them single-threaded. Which is a shame, because the 2019 MacBook Pro that I am writing this on is much more powerful than that, and has multiple cores that would support concurrent execution of multiple threads:

```
jones@powder:~$ sysctl hw.ncpu hw.physicalcpu  
hw.ncpu: 12  
hw.physicalcpu: 6
```

To interpret that output, this machine `powder` has 6 physical CPUs, each of which has two processing cores, for a total of 12 cores for computation. I'd sure like my computations to use all of that processing power. Because it's all on one machine, we can

easily use *multicore* processing tools to make use of those cores. Now let's look at the computational server `included-crab` at NCEAS:

```
jones@included-crab:~$ lscpu | egrep 'CPU\(\s\)|per core|per socket'
CPU(s):                      88
On-line CPU(s) list:          0-87
Thread(s) per core:           1
Core(s) per socket:           1
NUMA node0 CPU(s):            0-87
```

Now that's more compute power! `included-crab` has 384 GB of RAM, and ample storage. All still under the control of a single operating system.

Finally, maybe one of these NSF-sponsored high performance computing clusters (HPC) is looking attractive about now:

- [JetStream](#)
 - 640 nodes, 15,360 cores, 80TB RAM
- [Stampede2](#) at TACC
 - 4200 KNL nodes: 285,600 cores
 - 1736 SKX nodes: 83,328 cores
 - 224 ICX nodes: 17,920 cores
 - **TOTAL: 386,848 cores**

Note that these clusters have multiple nodes (hosts), and each host has multiple cores. So this is really multiple computers clustered together to act in a coordinated fashion, but each node runs its own copy of the operating system, and is in many ways independent of the other nodes in the cluster. One way to use such a cluster would be to use just one of the nodes, and use a multi-core approach to parallelization to use all of the cores on that single machine. But to truly make use of the whole cluster, one must use parallelization tools that let us spread out our computations across multiple host nodes in the cluster.

4.4 Parallel processing in the shell

Shell programming helps massively speed up data management tasks, and even more so with simple use of the [GNU parallel](#) utility to execute bash commands in parallel. In its simplest form, this can be used to speed up common file operations, such as renaming, compression, decompression, and file transfer. Let's look at a common example – calculating checksums to verify file integrity for data files. Calculating a hash checksum using the `shasum` command can be time consuming, especially when you have a lot of large files to work on. But it is a classic processor-limited task, and one that can be massively faster using `parallel`.

```
$ for fn in `ls *.gpkg`; do shasum -a 256 ${fn}; done  
  
real    35.081s  
user    32.745s  
system  2.336s  
  
$ ls *.gpkg | parallel "shasum -a 256 {}"  
  
real    2.97s  
user    37.16s  
system  2.70s
```

The first invocation takes *35 seconds* to execute the tasks one at a time serially, while the second version only takes :tada: *3 seconds* :tada: to do the same tasks. Note that the computational time spent in `user` and `system` processing is about the same, with the major difference being that the user-space tasks were conducted on multiple cores in parallel, resulting in more than 10x faster performance. Using `htop`, you can see processor cores spiking in usage when the command is run:

4.5 Modes of parallelization

Several different approaches can be taken to structuring a computer program to take advantage of the hardware capabilities of



Figure 4.3: Processor usage for parallel tasks.

multi-core processors. In the typical, and simplest, case, each task in a computation is executed serially in order of first to last. The total computation time is the sum of the time of all of the subtasks that are executed. In the next figure, a single core of the processor is used to sequentially execute each of the five tasks, with time flowing from left to right.

In comparison, the middle panel shows two approaches to parallelization on a single computer. With **multi-threaded** execution, a separate thread of execution is created for each of the 5 tasks, and these are executed concurrently on 5 of the cores of the processor. All of the threads are in the same process and share the same memory and resources, so one must take care that they do not interfere with each other.

With **multi-process** execution, a separate process is created for each of the 5 tasks, and these are executed concurrently on the cores of the processor. The difference is that each process has its own copy of the program memory, and changes are merged when each child process completes. Because each child process must be created and resources for that process must be marshaled and unmarshaled, there is more overhead in creating a process than a thread.

Finally, **cluster parallel** execution is shown in the last panel, in which a cluster with multiple computers is used to execute multiple processes for each task. Again, there is a setup task associated with creating and marshaling resources for the task, which now includes the overhead of moving data from one machine to the others in the cluster over the network. This further increases the cost of creating and executing multiple processes, but can be highly advantageous when accessing exceedingly large numbers of processing cores on clusters.

Serial



Parallel Threads



Parallel Processes



Cluster Parallel



Figure 4.4: Serial and parallel execution of tasks using threads and processes.

The key to performance gains is to ensure that the overhead associated with creating new threads or processes is small relative to the time it takes to perform a task. Somewhat unintuitively, when the setup overhead time exceeds the task time, parallel execution will likely be slower than serial.

4.6 Task parallelization in Python

Python also provides a number of easy to use packages for concurrent processing. We will review two of these, `concurrent.futures` and `parsl`, to show just how easy it can be to parallelize your programs. `concurrent.futures` is built right into the python3 release, and is a great starting point for learning concurrency.

We're going to start with a task that is a little expensive to compute, and define it in a function. All this `task(x)` function does is to use numpy to create a fairly large range of numbers, and then sum them.

```
def task(x):
    import numpy as np
    result = np.arange(x*10**8).sum()
    return result
```

We can start by executing this task function serially ten times with varying inputs. In this case, we create a function `run_serial` that takes a list of inputs to be run, and it calls the `task` function for each of those inputs. The `@timethis` decorator is a simple way to wrap the function with timing code so that we can see how long it takes to execute.

```
import numpy as np

@timethis
def run_serial(task_list):
    return [task(x) for x in task_list]

run_serial(np.arange(10))
```

```
run_serial: 90626.98006629944 ms
```

```
[0,
4999999950000000,
19999999900000000,
44999999850000000,
79999999800000000,
124999999750000000,
179999999700000000,
244999999650000000,
319999999600000000,
404999999550000000]
```

In this case, it takes around *25 seconds* to execute 10 tasks, depending on what else is happening on the machine and network.

So, can we make this faster using a multi-threaded parallel process? Let's try with `concurrent.futures`. The main concept in this package is one of a `future`, which is a structure which represents the value that will be created in a computation in the future when the function completes execution. With `concurrent.futures`, tasks are scheduled and do not block while they await their turn to be executed. Instead, threads are created and executed *asynchronously*, meaning that the function returns its `future` potentially before the thread has actually been executed. Using this approach, the user schedules a series of tasks to be executed asynchronously, and keeps track of the futures for each task. When the future indicates that the execution has been completed, we can then retrieve the result of the computation.

In practice this is a simple change from our serial implementation. We will use the `ThreadPoolExecutor` to create a pool of workers that are available to process tasks. Each worker is set up in its own thread, so it can execute in parallel with other workers. After setting up the pool of workers, we use `concurrent.futures map()` to schedule each task from our `task_list` (in this case, an input value from 1 to 10) to run on one of the workers. As for all `map()` implementations, we are asking for each value in `task_list` to be executed in the `task` function

we defined above, but in this case it will be executed using one of the workers from the `executor` that we created.

```
from concurrent.futures import ThreadPoolExecutor

@timethis
def run_threaded(task_list):
    with ThreadPoolExecutor(max_workers=20) as executor:
        return executor.map(task, task_list)

results = run_threaded(np.arange(10))
[x for x in results]

run_threaded: 4440.109491348267 ms
[0,
 4999999950000000,
 19999999900000000,
 44999999850000000,
 79999999800000000,
 124999999750000000,
 179999999700000000,
 244999999650000000,
 319999999600000000,
 404999999550000000]
```

This execution took about :tada: 4 seconds :tada:, which is about 6.25x faster than serial. Congratulations, you wrote your a multi-threaded python program!

4.7 Exercise: Parallel downloads

In this exercise, we're going to parallelize a simple task that is often very time consuming – downloading data. And we'll compare performance of simple downloads using first a serial loop, and then using two parallel execution libraries: `concurrent.futures` and `parsl`. We're going to see an example here where parallel execution won't always speed up this task, as this is likely an I/O bound task if you're

downloading a lot of data. But we still should be able to speed things up a lot until we hit the limits of our disk arrays.

The data we are downloading is a pan-Arctic time series of TIF images containing rasterized Arctic surface water indices from:

Elizabeth Webb. 2022. Pan-Arctic surface water (yearly and trend over time) 2000-2022. Arctic Data Center [doi:10.18739/A2NK3665N](https://doi.org/10.18739/A2NK3665N).



Figure 4.5: Webb surface water index data

First, let's download the data serially to set a benchmark. The data files are listed in a table with their filename and identifier, and can be downloaded directly from the Arctic Data Center using their identifier. To make things easier, we've already provided a data frame with the names and identifiers of each file that could be downloaded.

```
/opt/hostedtoolcache/Python/3.9.16/x64/lib/python3.9/site-packages/IPython/core/formatters.py:  
In future versions `DataFrame.to_latex` is expected to utilise the base implementation of `Sty
```

	filename	identifier
0	SWI_2007.tif	urn:uuid:5ee72c9c-789d-4a1c-95d8-cb2b24a20662
1	SWI_2019.tif	urn:uuid:9cd1cdc3-0792-4e61-afff-c11f86d3a9be
2	SWI_2021.tif	urn:uuid:14e1e509-77c0-4646-9cc3-d05f8d84977c
3	SWI_2020.tif	urn:uuid:1ba473ff-8f03-470b-90d1-7be667995ea1
4	SWI_2001.tif	urn:uuid:85150557-05fd-4f52-8bbd-ec5a2c27e23d

4.7.1 Serial

When you have a list of repetitive tasks, you may be able to speed it up by adding more computing power. If each task is completely independent of the others, then it is pleasingly parallel and a prime candidate for executing those tasks in parallel, each on its own core. For example, let's build a simple loop that downloads the data files that we need for an analysis. First, we start with the serial implementation.

```
import urllib

def download_file(row):
    service = "https://arcticdata.io/metacat/d1/mn/v2/object/"
    pid = row[1]['identifier']
    filename = row[1]['filename']
    url = service + pid
    print("Downloading: " + filename)
    msg = urllib.request.urlretrieve(url, filename)
    return filename

@timethis
def download_serial(table):
    results = [download_file(row) for row in table.iterrows()]
    return results

results = download_serial(file_table[0:5])
print(results)
```

Downloading: SWI_2007.tif

Downloading: SWI_2019.tif

```
Downloading: SWI_2021.tif
```

```
Downloading: SWI_2020.tif
```

```
Downloading: SWI_2001.tif
```

```
download_serial: 68629.84848022461 ms
['SWI_2007.tif', 'SWI_2019.tif', 'SWI_2021.tif', 'SWI_2020.tif', 'SWI_2001.tif']
```

In this code, we have one function (`download_file`) that downloads a single data file and saves it to disk. It is called iteratively from the function `download_serial`. The serial execution takes about *20-25 seconds*, but can vary considerably based on network traffic and other factors.

The issue with this loop is that we execute each download task sequentially, which means that only one of our processors on this machine is in use. In order to exploit parallelism, we need to be able to dispatch our tasks and allow each to run at the same time, with one task going to each core. To do that, we can use one of the many parallelization libraries in python to help us out.

4.7.2 Multi-threaded with `concurrent.futures`

In this case, we'll use the same `download_file` function from before, but let's switch up and create a `download_threaded()` function to use `concurrent.futures` with a `ThreadPoolExecutor` just as we did earlier.

```
from concurrent.futures import ThreadPoolExecutor

@timethis
def download_threaded(table):
    with ThreadPoolExecutor(max_workers=15) as executor:
        results = executor.map(download_file, table.iterrows(), timeout=60)
        return results

results = download_threaded(file_table[0:5])
```

```
for result in results:  
    print(result)  
  
Downloading: SWI_2007.tif  
Downloading: SWI_2019.tif  
Downloading: SWI_2021.tif  
Downloading: SWI_2020.tif  
Downloading: SWI_2001.tif  
  
download_threaded: 30993.90411376953 ms  
SWI_2007.tif  
SWI_2019.tif  
SWI_2021.tif  
SWI_2020.tif  
SWI_2001.tif
```

Note how the “Downloading...” messages were printed immediately, but then it still took over 20 seconds to download the 5 files. This could be for several reasons, including that one of the files alone took that long (e.g., due to network congestion), or that there was a bottleneck in writing the files to disk (i.e., we could have been disk I/O limited). Or maybe the multithreaded executor pool didn’t do a good job parallelizing the tasks. The trick is figuring out why you did or didn’t get a speedup when parallelizing. So, let’s try this another way, using a **multi-processing** approach, rather than multi-threading.

4.7.3 Multi-process with concurrent.futures

You’ll remember from earlier that you can execute tasks concurrently by creating multiple threads within one process (multi-threaded), or by creating and executing multiple processes. The latter creates more independence, as each of the executing tasks has their own memory and process space, but it also takes longer to set up. With `concurrent.futures`, we can switch to a multi-process pool by using a `ProcessPoolExecutor`, analogously to how we used `ThreadPoolExecutor` previously. So, some simple changes, and we’re running multiple processes.

```
from concurrent.futures import ProcessPoolExecutor

@timethis
def download_process(table):
    with ProcessPoolExecutor(max_workers=15) as executor:
        results = executor.map(download_file, table.iterrows(), timeout=60)
    return results

results = download_process(file_table[0:5])
for result in results:
    print(result)
```

Downloading: SWI_2020.tif

Downloading: SWI_2007.tif

Downloading: SWI_2001.tif

Downloading: SWI_2021.tif

Downloading: SWI_2019.tif

```
download_process: 15228.946685791016 ms
SWI_2007.tif
SWI_2019.tif
SWI_2021.tif
SWI_2020.tif
SWI_2001.tif
```

Again, the output messages print almost immediately, but then later the processes finish and report that it took between *10 to 15 seconds* to run. Your mileage may vary. When I increase the number of files being downloaded to 10 or even to 20, I notice it is actually about the same, around *10-15 seconds*. So, part of our time now is the overhead of setting up multiple processes. But once we have that infrastructure in place, we can make effective use of the pool of processes to handle many downloads.

4.8 Parallel processing with `parsl`

`concurrent.futures` is great and powerful, but it has its limits. Particularly as you try to scale up into the thousands of concurrent tasks, other libraries like [ParSL \(docs\)](#), [Dask](#), [Ray](#), and others come into play. They all have their strengths, but ParSL makes it particularly easy to build parallel workflows out of existing python code through it's use of decorators on existing python functions.

In addition, ParSL supports a lot of different kinds of `providers`, allowing the same python code to be easily run multi-threaded using a `ThreadPoolExecutor` and via multi-processing on many different cluster computing platforms using the `HighThroughputExecutor`. For example, ParSL includes providers supporting local execution, and on Slurm, Condor, Kubernetes, AWS, and other platforms. And ParSL handles data staging as well across these varied environments, making sure the data is in the right place when it's needed for computations.

Similarly to before, we start by configuring an executor in `parsl`, and loading it. We'll use multiprocessing by configuring the `HighThroughputExecutor` to use our local resources as a cluster, and we'll activate our virtual environment to be sure we're executing in a consistent environment.

```
# Required packages
import parsl
from parsl import python_app
```

```

from parsl.config import Config
from parsl.executors import HighThroughputExecutor
from parsl.providers import LocalProvider

# Configure the parsl executor
activate_env = 'workon scomp'
htex_local = Config(
    executors=[
        HighThroughputExecutor(
            max_workers=15,
            provider=LocalProvider(
                worker_init=activate_env
            )
        )
    ],
)
parsl.clear()
parsl.load(htex_local)

```

<parsl.dataflow.dflow.DataFlowKernel at 0x7f772f79eb20>

We now have a live parsl executor (`htex_local`) that is waiting to execute processes. We tell it to execute processes by annotating functions with decorators that indicate which tasks should be parallelized. Parsl then handles the scheduling and execution of those tasks based on the dependencies between them. In the simplest case, we'll decorate our previous function for downloading a file with the `@python_app` decorator, which tells parsl that any function calls with this function should be run on the default executor (in this case, `htex_local`).

```

# Decorators seem to be ignored as the first line of a cell, so print something first
print("Create decorated function")

@python_app
def download_file_parsl(row):
    import urllib
    service = "https://arcticdata.io/metacat/d1/mn/v2/object/"
    pid = row[1]['identifier']

```

```

filename = row[1]['filename']
url = service + pid
print("Downloading: " + filename)
msg = urllib.request.urlretrieve(url, filename)
return filename

```

Create decorated function

Now we just write regular python code that calls that function, and parsl handles the scheduling. Parsl app executors return an `AppFuture`, and we can call the `AppFuture.done()` function to determine when the future result is ready without blocking. Or, we can just block on `AppFuture.result()` which waits for each of the executions to complete and then returns the result.

```

#! eval: true
print("Define our download_futures function")

@timethis
def download_futures(table):
    results = []
    for row in table.iterrows():
        result = download_file_parsl(row)
        print(result)
        results.append(result)
    return(results)

@timethis
def wait_for_futures(table):
    results = download_futures(table)
    done = [app_future.result() for app_future in results]
    print(done)

wait_for_futures(file_table[0:5])

```

```

Define our download_futures function
<AppFuture at 0x7f777063fe20 state=pending>
<AppFuture at 0x7f772f6b8c70 state=pending>
<AppFuture at 0x7f772f6c61f0 state=pending>

```

```
<AppFuture at 0x7f772f6c6730 state=pending>
<AppFuture at 0x7f772f6c6c70 state=pending>
download_futures: 14.460086822509766 ms

['SWI_2007.tif', 'SWI_2019.tif', 'SWI_2021.tif', 'SWI_2020.tif', 'SWI_2001.tif']
wait_for_futures: 54541.02182388306 ms
```

When we're done, be sure to clean up and shutdown the `htex_local` executor, or it will continue to persist in your environment and utilize resources. Generally, an executor should be created when setting up your environment, and then it can be used repeatedly for many different tasks.

```
htex_local.executors[0].shutdown()
parsl.clear()
```

4.9 When to parallelize

It's not as simple as it may seem. While in theory each added processor would linearly increase the throughput of a computation, there is overhead that reduces that efficiency. For example, the code and, importantly, the data need to be copied to each additional CPU, and this takes time and bandwidth. Plus, new processes and/or threads need to be created by the operating system, which also takes time. This overhead reduces the efficiency enough that realistic performance gains are much less than theoretical, and usually do not scale linearly as a function of processing power. For example, if the time that a computation takes is short, then the overhead of setting up these additional resources may actually overwhelm any advantages of the additional processing power, and the computation could potentially take longer!

In addition, not all of a task can be parallelized. Depending on the proportion, the expected speedup can be significantly reduced. Some propose that this may follow [Amdahl's Law](#), where the speedup of the computation (y-axis) is a function of both the number of cores (x-axis) and the proportion of the computation that can be parallelized (see colored lines):

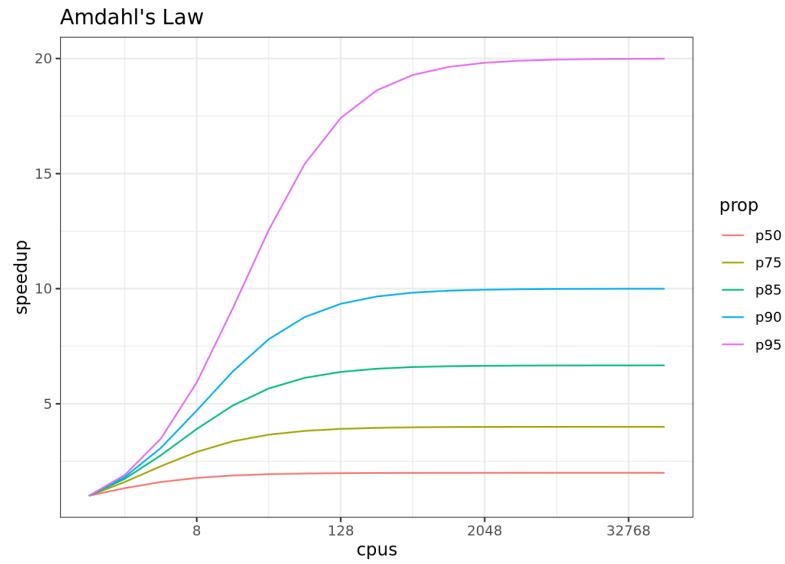


Figure 4.6: Amdahl's Law

So, it is important to evaluate the computational efficiency of requests, and work to ensure that additional compute resources brought to bear will pay off in terms of increased work being done.

4.10 Parallel Pitfalls

A set of tasks is considered ‘pleasingly parallel’ when large portions of the code can be executed independently of the other portions and have few or no dependencies on other parts of the execution. This situation is common, and we can frequently execute parallel tasks on independent subsets of our data. Nevertheless, dependencies among different parts of your computation can definitely create bottlenecks and slow down computations. Some of the challenges you may need to work around include:

- **Task dependencies:** occur when one task in the code depends on the results of another task or computation in the code.

- **Race conditions:** occur when two tasks execute in parallel, but produce different results based on which task finishes first. Ensuring that results are correct under different timing situations requires careful testing.
- **Deadlocks:** occur when two concurrent tasks block on the output of the other. Deadlocks cause parallel programs to lock up indefinitely, and can be difficult to track down.

Even when tasks exhibit strong dependencies, it is frequently possible to still optimize that code by parallelizing explicit code sections, sometimes bringing other concurrency tools into the mix, such as the Message Passing Interface (MPI). But simply improving the efficiency of pleasingly parallel tasks can be liberating.

4.11 Summary

In this lesson, we showed examples of computing tasks that are likely limited by the number of CPU cores that can be applied, and we reviewed the architecture of computers to understand the relationship between CPU processors and cores. Next, we reviewed the way in which traditional sequential loops can be rewritten as functions that are applied to a list of inputs both serially and in parallel to utilize multiple cores to speed up computations. We reviewed the challenges of optimizing code, where one must constantly examine the bottlenecks that arise as we improve cpu-bound, I/O bound, and memory bound computations.

4.12 Further Reading

Ryan Abernathey & Joe Hamman. 2020. [Closed Platforms vs. Open Architectures for Cloud-Native Earth System Analytics](#). Medium.

5 Data Structures and Formats for Large Data

- Learn about the NetCDF data format:
 - Characteristics: self-describing, scalable, portable, appendable, shareable, and archivable
 - Understand the NetCDF data model: what are dimensions, variables, and attributes
 - Advantages and differences between NetCDF and tabular data formats
- Learn how to use the `xarray` Python package to work with NetCDF files:
 - Describe the core `xarray` data structures, the `xarray.DataArray` and the `xarray.Dataset`, and their components, including data variables, dimensions, coordinates, and attributes
 - Create `xarray.DataArrays` and `xarray.DataSets` out of raw `numpy` arrays and save them as netCDF files
 - Load `xarray.DataSets` from netCDF files and understand the attributes view
 - Perform basic indexing, processing, and reduction of `xarray.DataArrays`
 - Convert `pandas.DataFrame`s into `xarray.DataSets`

5.1 Introduction

Efficient and reproducible data analysis begins with choosing a proper format to store our data, particularly when working with large, complex, multi-dimensional datasets. Consider, for example, the following Earth System Data Cube from [Mahecha](#)

[et al. 2020](#), which measures nine environmental variables at high resolution across space and time. We can consider this dataset large (high-resolution means we have a big file), complex (multiple variables), and multi-dimensional (each variable is measured along three dimensions: latitude, longitude, and time). Additionally, necessary metadata must accompany the dataset to make it functional, such as units of measurement for variables, information about the authors, and processing software used.

Keeping complex datasets in a format that facilitates access, processing, sharing, and archiving can be at least as important as how we parallelize the code we use to analyze them. In practice, it is common to convert our data from less efficient formats into more efficient ones before we parallelize any processing. In this lesson, we will

1. familiarize ourselves with the NetCDF data format, which enables us to store large, complex, multi-dimensional data efficiently, and
2. learn to use the `xarray` Python package to read, process, and create NetCDF files.

5.2 NetCDF Data Format

[NetCDF](#) (network Common Data Form) is a set of software libraries and self-describing, machine-independent data formats that support the creation, access, and sharing of array-oriented scientific data. NetCDF was initially developed at the Unidata Program Center and is supported on almost all platforms, and parsers exist for most scientific programming languages.

The [NetCDF documentation](#) outlines that this data format is designed to be:

1. **Self-describing:** Information describing the data contents of the file is embedded within the data file itself. This means that there is a header describing the layout of the rest of the file and arbitrary file metadata.

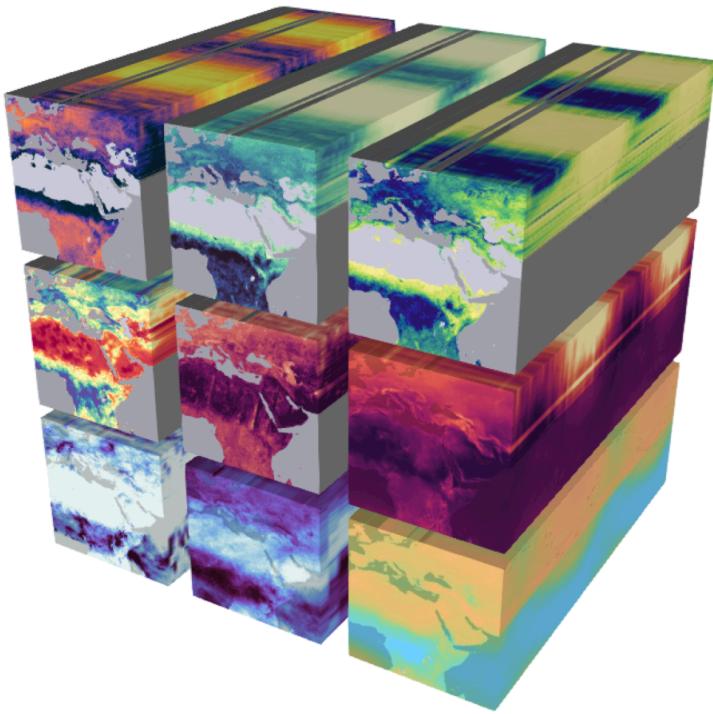


Figure 5.1: Mahecha et al. 2020 . *Visualization of the implemented Earth system data cube. The figure shows from the top left to bottom right the variables sensible heat (H), latent heat (LE), gross primary production (GPP), surface moisture (SM), land surface temperature (LST), air temperature ($Tair$), cloudiness (C), precipitation (P), and water vapour (V). The resolution in space is 0.25° and 8 d in time, and we are inspecting the time from May 2008 to May 2010; the spatial range is from $15^\circ S$ to $60^\circ N$, and $10^\circ E$ to $65^\circ W$.*

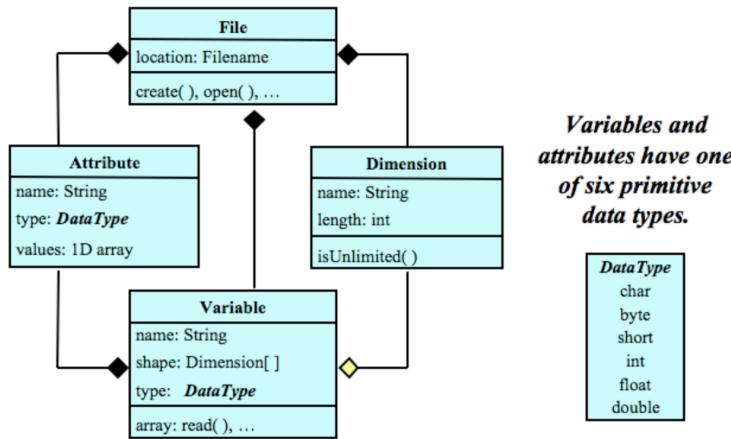
2. **Scalable:** Small subsets of large datasets may be accessed efficiently through netCDF interfaces, even from remote servers.
3. **Portable:** A NetCDF file is machine-independent i.e. it can be accessed by computers with different ways of storing integers, characters, and floating-point numbers.
4. **Appendable:** Data may be appended to a properly structured NetCDF file without copying the dataset or redefining its structure.
5. **Sharable:** One writer and multiple readers may simultaneously access the same NetCDF file.
6. **Archivable:** Access to all earlier forms of NetCDF data will be supported by current and future versions of the software.

5.2.1 Data Model

The NetCDF data model is the way that NetCDF organizes data. This lesson will follow the [Classic NetCDF Data Model](#), which is at the core of all netCDF files.

The model consists of three key components: **variables**, **dimensions**, and **attributes**.

- **Variables** are N-dimensional arrays of data. We can think of these as varying/measured/dependent quantities.
- **Dimensions** describe the axes of the data arrays. A dimension has a name and a length. We can think of these as the constant/fixed/independent quantities at which we measure the variables.
- **Attributes** are small notes or supplementary metadata to annotate a variable or the file as a whole.



A file has named variables, dimensions, and attributes. Variables also have attributes. Variables may share dimensions, indicating a common grid. One dimension may be of unlimited length.

Figure 5.2: Classic NetCDF Data Model ([NetCDF documentation](#))

5.2.2 Metadata Standards

The most commonly used metadata standard for geospatial data is the **Climate and Forecast metadata standard**, also called the [CF conventions](#).

The CF conventions are specifically designed to promote the processing and sharing of files created with the NetCDF API. Principles of CF include self-describing data (no external tables needed for understanding), metadata equally readable by humans and software, minimum redundancy, and maximum simplicity. ([CF conventions FAQ](#))

The CF conventions provide a unique standardized name and precise description of over 1,000 physical variables. To maximize the reusability of our data, it is best to include a variable's standardized name as an attribute called **standard_name**. Variables should also include a **units** attribute. This attribute should be a string that can be recognized by UNIDATA's [UDUNITS package](#). In these links you can find:

- a table with all of the CF convention's standard names, and
- a list of the units found in the UDUNITS database maintained by the North Carolina Institute for Climate Studies.

5.2.3 Exercise

Let's do a short practice now that we have reviewed the classic NetCDF model and know a bit about metadata best practices.

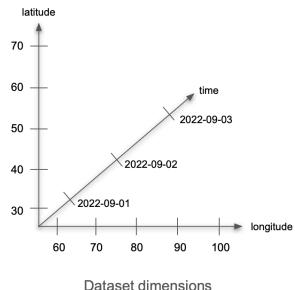
Part 1

Imagine the following scenario: we have a network of 25 weather stations. They are located in a square grid: starting at 30°0 N 60°0 E, there is a station every 10° North and every 10° East. Each station measures the air temperature at a set time for three days, starting on September 1st, 2022. On the first day, all stations record a temperature of 0°C. On the second day, all temperatures are 1°C, and on the third day, all temperatures are 2°C. What are the *variables*, *dimensions* and *attributes* for this data?

Answer

Variables: There is a single variable being measured: temperature. The variable values can be represented as a 5x5x3 array, with constant values for each day.

Dimensions: This dataset has three dimensions: time, latitude, and longitude. Time indicates when the measurement happened, we can encode it as the dates 2022-09-01, 2022-09-02, and 2022-09-03. The pairs of latitude and longitude values indicate the positions of the weather stations. Latitude has values 30, 40, 50, 60, and 70, measured in degrees North. Longitude has values 60, 70, 80, 90, and 100, measured in degrees East.



					2	2	2	2	2
					2	2	2	2	2
			1	1	1	1	1	2	2
		1	1	1	1	1	1	2	2
0	0	0	0	0	1	1	1	2	2
0	0	0	0	0	1	1	1	2	2
0	0	0	0	0	1	1	1	2	2
0	0	0	0	0	1	1	1	2	2
0	0	0	0	0	1	1	1	2	2
0	0	0	0	0	1	1	1	2	2

Temperature variable

Attributes: Let's divide these into attributes for the variable, the dimensions, and the whole dataset:

- Variable attributes:
 - Temperature attributes:
 - * standard_name: air_temperature
 - * units: degree_C
- Dimension attributes:
 - Time attributes:
 - * description: date of measurement
 - Latitude attributes:
 - * standard_name: grid_latitude
 - * units: degrees_N
 - Longitude attributes:
 - * standard_name: grid_longitude
 - * units: degree_E
- Dataset attributes:
 - title: Temperature Measurements at Weather Stations
 - summary: an example of NetCDF data format

Part 2

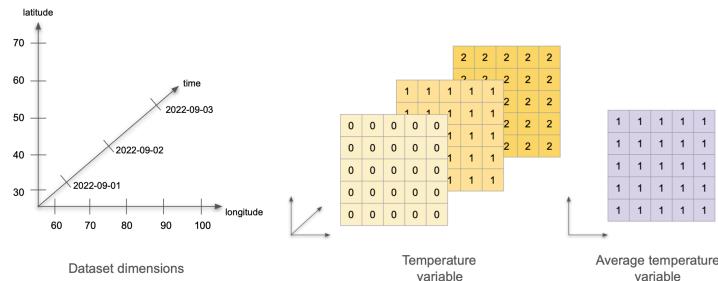
Now imagine we calculate the average temperature over time at each weather station, and we wish to incorporate this data into the same dataset. How will adding the average tempera-

ture data change the dataset's variables, attributes, and dimensions?

Answer

Variables: Now we are measuring two variables: temperature and average temperature. The temperature data stays the same. We can represent the average temperature as a single 5x5 array with value 1 at each cell.

Dimensions: This dataset still has three dimensions: time, latitude, and longitude. The temperature variable uses all three dimensions, and the average temperature variable only uses two (latitude and longitude). This is ok! The dataset's dimensions are the union of the dimensions of all the variables in the dataset. Variables in the same dataset may have all, some, or no dimensions in common.



Attributes: To begin with, we need to keep all the previous attributes. Notice that the dataset's title is general enough that we don't need to update it. The only update we need to do is add the attributes for our new average temperature variable:

- Average temperature attributes:
 - standard_name: average_air_temperature
 - description: average temperature over three days

Our next step is to see how we can translate all this information into something we can store and handle on our computers.

5.3 xarray

xarray is an open source project and Python package that augments NumPy arrays by adding labeled dimensions, coordinates and attributes. xarray is based on the netCDF data model, making it the appropriate tool to open, process, and create datasets in netCDF format.



Figure 5.3: [xarray's development portal](#)

5.3.1 xarray.DataArray

The `xarray.DataArray` is the primary data structure of the xarray package. It is an n-dimensional array with **labeled dimensions**. We can think of it as representing a single variable in the NetCDF data format: it holds the variable's values, dimensions, and attributes.

Apart from variables, dimensions, and attributes, xarray introduces one more piece of information to keep track of a dataset's content: in xarray each dimension has at least one set of **coordinates**. A dimension's coordinates indicate the dimension's values. We can think of the coordinate's values as the tick labels along a dimension. For example, in our previous exercise about temperature measured in weather stations, latitude is a dimension, and the latitude's coordinates are 30, 40, 50, 60, and 70 because those are the latitude values at which we are collecting temperature data. In that same exercise, time is a dimension, and its coordinates are 2022-09-1, 2022-09-02, and 2022-09-03.

Here you can read more about the [xarray terminology](#).

5.3.1.1 Create an `xarray.DataArray`

Let's suppose we want to make an `xarray.DataArray` that includes the information from our previous exercise about measuring temperature across three days. First, we import all the necessary libraries.

```
import os
import urllib
import pandas as pd
import numpy as np

import xarray as xr    # This is the package we'll explore
```

Variable Values

The underlying data in the `xarray.DataArray` is a `numpy.ndarray` that holds the variable values. So we can start by making a `numpy.ndarray` with our mock temperature data:

```
# values of a single variable at each point of the coords
temp_data = np.array([np.zeros((5,5)),
                      np.ones((5,5)),
                      np.ones((5,5))*2]).astype(int)
temp_data

array([[ [0, 0, 0, 0, 0],
         [0, 0, 0, 0, 0],
         [0, 0, 0, 0, 0],
         [0, 0, 0, 0, 0],
         [0, 0, 0, 0, 0]],

        [[1, 1, 1, 1, 1],
         [1, 1, 1, 1, 1],
         [1, 1, 1, 1, 1],
         [1, 1, 1, 1, 1],
         [1, 1, 1, 1, 1]],

        [[2, 2, 2, 2, 2],
         [2, 2, 2, 2, 2],
```

```
[2, 2, 2, 2, 2],  
[2, 2, 2, 2, 2],  
[2, 2, 2, 2]])])
```

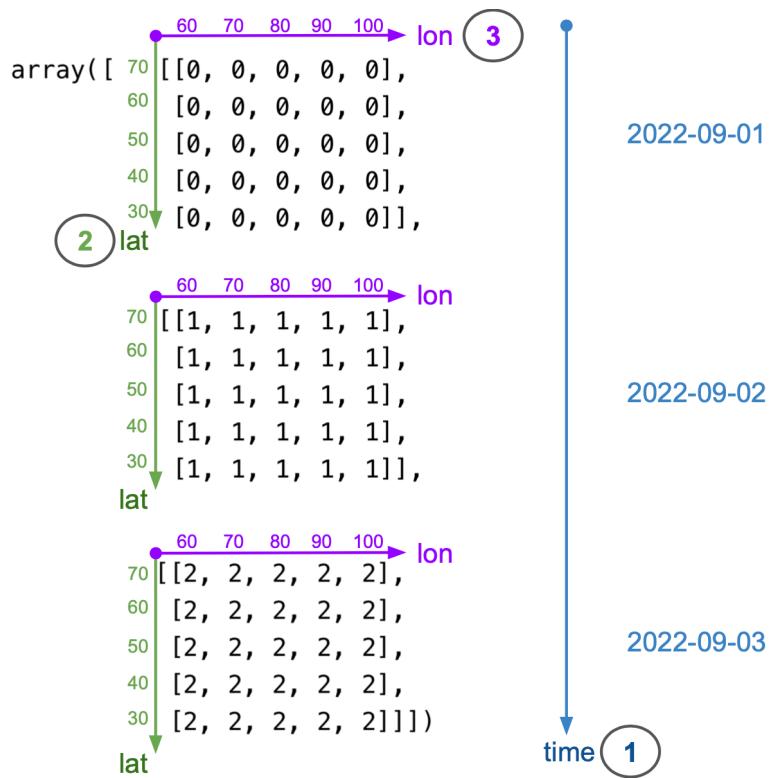
We could think this is “all” we need to represent our data. But if we stopped at this point, we would need to

1. remember that the numbers in this array represent the temperature in degrees Celsius (doesn’t seem too bad),
2. remember that the first dimension of the array represents time, the second latitude and the third longitude (maybe ok), and
3. keep track of the range of values that time, latitude, and longitude take (not so good).

Keeping track of all this information separately could quickly get messy and could make it challenging to share our data and analyses with others. This is what the netCDF data model and `xarray` aim to simplify. We can get data and its descriptors together in an `xarray.DataArray` by adding the dimensions over which the variable is being measured and including attributes that appropriately describe dimensions and variables.

Dimensions and Coordinates

To specify the dimensions of our upcoming `xarray.DataArray`, we must examine how we’ve constructed the `numpy.ndarray` holding the temperature data. The diagram below shows how the dimensions of `temp_data` are ordered: the first dimension is time, the second is latitude, and the third is longitude.



Remember that indexing in 2-dimensional `numpy.ndarrays` starts at the top-left corner of the array, and it is done by rows first and columns second (like matrices). This is why latitude is the second dimension and longitude the third. From the diagram, we can also see that the coordinates (values of each dimension) are as follow:

- date coordinates are 2022-09-01, 2022-09-02, 2022-09-03
- latitude coordinates are 70, 60, 50, 40, 30 (notice decreasing order)
- longitude coordinates are 60, 70, 80, 90, 100 (notice increasing order)

We add the dimensions as a tuple of strings and coordinates as a dictionary:

```
# names of the dimensions in the required order
dims = ('time', 'lat', 'lon')
```

```
# create coordinates to use for indexing along each dimension
coords = {'time' : pd.date_range("2022-09-01", "2022-09-03"),
          'lat' : np.arange(70, 20, -10),
          'lon' : np.arange(60, 110, 10)}
```

Attributes

Next, we add the attributes (metadata) for our temperature data as a dictionary:

```
# attributes (metadata) of the data array
attrs = { 'title' : 'temperature across weather stations',
          'standard_name' : 'air_temperature',
          'units' : 'degree_c'}
```

Putting It All Together

Finally, we put all these pieces together (data, dimensions, coordinates, and attributes) to create an `xarray.DataArray`:

```
# initialize xarray.DataArray
temp = xr.DataArray(data = temp_data,
                     dims = dims,
                     coords = coords,
                     attrs = attrs)
```

```
temp
```

```
<xarray.DataArray (time: 3, lat: 5, lon: 5)>
array([[ [0, 0, 0, 0, 0],
         [0, 0, 0, 0, 0],
         [0, 0, 0, 0, 0],
         [0, 0, 0, 0, 0],
         [0, 0, 0, 0, 0]],

        [[1, 1, 1, 1, 1],
         [1, 1, 1, 1, 1],
         [1, 1, 1, 1, 1],
         [1, 1, 1, 1, 1],
         [1, 1, 1, 1, 1]],

        [[2, 2, 2, 2, 2],
```

```
[2, 2, 2, 2, 2],  
[2, 2, 2, 2, 2],  
[2, 2, 2, 2, 2],  
[2, 2, 2, 2, 2]])
```

Coordinates:

```
* time      (time) datetime64[ns] 2022-09-01 2022-09-02 2022-09-03  
* lat       (lat) int64 70 60 50 40 30  
* lon       (lon) int64 60 70 80 90 100
```

Attributes:

```
title:          temperature across weather stations  
standard_name: air_temperature  
units:          degree_c
```

We can also update the variable's attributes after creating the object. Notice that each of the coordinates is also an `xarray.DataArray`, so we can add attributes to them.

```
# update attributes  
temp.attrs['description'] = 'simple example of an xarray.DataArray'  
  
# add attributes to coordinates  
temp.time.attrs = {'description':'date of measurement'}  
  
temp.lat.attrs['standard_name']= 'grid_latitude'  
temp.lat.attrs['units'] = 'degree_N'  
  
temp.lon.attrs['standard_name']= 'grid_longitude'  
temp.lon.attrs['units'] = 'degree_E'  
temp  
  
<xarray.DataArray (time: 3, lat: 5, lon: 5)>  
array([[ [0, 0, 0, 0, 0],  
        [0, 0, 0, 0, 0],  
        [0, 0, 0, 0, 0],  
        [0, 0, 0, 0, 0],  
        [0, 0, 0, 0, 0]],  
  
       [[1, 1, 1, 1, 1],  
        [1, 1, 1, 1, 1],  
        [1, 1, 1, 1, 1],  
        [1, 1, 1, 1, 1],  
        [1, 1, 1, 1, 1]]])
```

```
[1, 1, 1, 1, 1],  
[1, 1, 1, 1, 1]],  
  
[[2, 2, 2, 2, 2],  
 [2, 2, 2, 2, 2],  
 [2, 2, 2, 2, 2],  
 [2, 2, 2, 2, 2],  
 [2, 2, 2, 2, 2]])
```

Coordinates:

```
* time      (time) datetime64[ns] 2022-09-01 2022-09-02 2022-09-03  
* lat       (lat) int64 70 60 50 40 30  
* lon       (lon) int64 60 70 80 90 100
```

Attributes:

```
title:           temperature across weather stations  
standard_name: air_temperature  
units:          degree_c  
description:    simple example of an xarray.DataArray
```

At this point, since we have a single variable, the dataset attributes and the variable attributes are the same.

5.3.1.2 Indexing

An `xarray.DataArray` allows both positional indexing (like `numpy`) and label-based indexing (like `pandas`). Positional indexing is the most basic, and it's done using Python's `[]` syntax, as in `array[i,j]` with `i` and `j` both integers. **Label-based indexing** takes advantage of dimensions in the array having names and coordinate values that we can use to access data instead of remembering the positional order of each dimension.

As an example, suppose we want to know what was the temperature recorded by the weather station located at 40°N 80°E on September 1st, 2022. By recalling all the information about how the array is setup with respect to the dimensions and coordinates, we can access this data positionally:

```
temp[0,1,2]
```

```
<xarray.DataArray ()>
array(0)
Coordinates:
  time    datetime64[ns] 2022-09-01
  lat      int64 60
  lon      int64 80
Attributes:
  title:          temperature across weather stations
  standard_name: air_temperature
  units:          degree_c
  description:   simple example of an xarray.DataArray
```

Or, we can use the dimensions names and their coordinates to access the same value:

```
temp.sel(time='2022-09-01', lat=40, lon=80)
```

```
<xarray.DataArray ()>
array(0)
Coordinates:
  time    datetime64[ns] 2022-09-01
  lat      int64 40
  lon      int64 80
Attributes:
  title:          temperature across weather stations
  standard_name: air_temperature
  units:          degree_c
  description:   simple example of an xarray.DataArray
```

Notice that the result of this indexing is a 1x1 `xarray.DataArray`. This is because operations on an `xarray.DataArray` (resp. `xarray.DataSet`) always return another `xarray.DataArray` (resp. `xarray.DataSet`). In particular, operations returning scalar values will also produce `xarray` objects, so we need to cast them as numbers manually. See [xarray.DataArray.item](#).

More about [xarray indexing](#).

5.3.1.3 Reduction

xarray has implemented several methods to reduce an xarray.DataArray along any number of dimensions. One of the advantages of xarray.DataArray is that, if we choose to, it can carry over attributes when doing calculations. For example, we can calculate the average temperature at each weather station over time and obtain a new xarray.DataArray.

```
avg_temp = temp.mean(dim = 'time')
# to keep attributes add keep_attrs = True

avg_temp.attrs = {'title':'average temperature over three days'}
avg_temp

<xarray.DataArray (lat: 5, lon: 5)>
array([[1., 1., 1., 1., 1.],
       [1., 1., 1., 1., 1.],
       [1., 1., 1., 1., 1.],
       [1., 1., 1., 1., 1.],
       [1., 1., 1., 1., 1.]])
Coordinates:
 * lat      (lat) int64 70 60 50 40 30
 * lon      (lon) int64 60 70 80 90 100
Attributes:
    title:    average temperature over three days
```

More about [xarray computations](#).

5.3.2 xarray.DataSet

An xarray.DataSet resembles an in-memory representation of a NetCDF file and consists of *multiple* variables (each being an xarray.DataArray), with dimensions, coordinates, and attributes, forming a self-describing dataset. Attributes can be specific to each variable, each dimension, or they can describe the whole dataset. The variables in an xarray.DataSet can have the same dimensions, share some dimensions, or have no dimensions in common. Let's see an example of this.

5.3.2.1 Create an `xarray.Dataset`

Following our previous example, we can create an `xarray.Dataset` by combining the temperature data with the average temperature data. We also add some attributes that now describe the whole dataset, not only each variable.

```
# make dictionaries with variables and attributes
data_vars = {'avg_temp': avg_temp,
             'temp': temp}

attrs = {'title':'temperature data at weather stations: daily and average',
         'description':'simple example of an xarray.Dataset'}

# create xarray.Dataset
temp_dataset = xr.Dataset( data_vars = data_vars,
                           attrs = attrs)
```

Take some time to click through the data viewer and read through the variables and metadata in the dataset. Notice the following:

- `temp_dataset` is a dataset with three dimensions (time, latitude, and longitude),
- `temp` is a variable that uses all three dimensions in the dataset, and
- `aveg_temp` is a variable that only uses two dimensions (latitude and longitude).

```
temp_dataset
```

```
<xarray.Dataset>
Dimensions:  (lat: 5, lon: 5, time: 3)
Coordinates:
 * lat        (lat) int64 70 60 50 40 30
 * lon        (lon) int64 60 70 80 90 100
 * time       (time) datetime64[ns] 2022-09-01 2022-09-02 2022-09-03
Data variables:
 avg_temp   (lat, lon) float64 1.0 1.0 1.0 1.0 1.0 ... 1.0 1.0 1.0 1.0 1.0
```

```

temp      (time, lat, lon) int64 0 0 0 0 0 0 0 0 0 ... 2 2 2 2 2 2 2 2 2
Attributes:
  title:      temperature data at weather stations: daily and average
  description: simple example of an xarray.Dataset

```

5.3.2.2 Save and Reopen

Finally, we want to save our dataset as a NetCDF file. To do this, specify the file path and use the `.nc` extension for the file name. Then save the dataset using the `to_netcdf` method with your file path. Opening NetCDF is similarly straightforward using `xarray.open_dataset()`.

```

# specify file path: don't forget the .nc extension!
fp = os.path.join(os.getcwd(), 'temp_dataset.nc')
# save file
temp_dataset.to_netcdf(fp)

# open to check:
check = xr.open_dataset(fp)
check

<xarray.Dataset>
Dimensions:  (lat: 5, lon: 5, time: 3)
Coordinates:
  * lat      (lat) int64 70 60 50 40 30
  * lon      (lon) int64 60 70 80 90 100
  * time     (time) datetime64[ns] 2022-09-01 2022-09-02 2022-09-03
Data variables:
  avg_temp  (lat, lon) float64 ...
  temp      (time, lat, lon) int64 ...
Attributes:
  title:      temperature data at weather stations: daily and average
  description: simple example of an xarray.Dataset

```

5.3.3 Exercise

For this exercise, we will use a dataset including time series of annual Arctic freshwater fluxes and storage terms. The data

was produced for the publication [Jahn and Laiho, 2020](#) about changes in the Arctic freshwater budget and is archived at the Arctic Data Center [doi:10.18739/A2280504J](#)

```
url = 'https://arcticdata.io/metacat/d1/mn/v2/object/urn%3Auuid%3A792bfc37-416e-409e-80b1-fd

msg = urllib.request.urlretrieve(url, "FW_data_CESM_LW_2006_2100.nc")

fp = os.path.join(os.getcwd(), 'FW_data_CESM_LW_2006_2100.nc')
fw_data = xr.open_dataset(fp)
fw_data

<xarray.Dataset>
Dimensions:                               (time: 95, member: 11)
Coordinates:
  * time                                (time) float64 2.006e+03 ... 2.1e+03
    * member                             (member) float64 1.0 2.0 3.0 ... 10.0 11.0
Data variables: (12/16)
  FW_flux_Fram_annual_net      (time, member) float64 ...
  FW_flux_Barrow_annual_net    (time, member) float64 ...
  FW_flux_Nares_annual_net     (time, member) float64 ...
  FW_flux_Davis_annual_net    (time, member) float64 ...
  FW_flux_BSO_annual_net      (time, member) float64 ...
  FW_flux_Bering_annual_net   (time, member) float64 ...
  ...
  Solid_FW_flux_BSO_annual_net (time, member) float64 ...
  Solid_FW_flux_Bering_annual_net (time, member) float64 ...
  runoff_annual                  (time, member) float64 ...
  netPrec_annual                 (time, member) float64 ...
  Liquid_FW_storage_Arctic_annual (time, member) float64 ...
  Solid_FW_storage_Arctic_annual (time, member) float64 ...
Attributes:
  creation_date: 02-Jun-2020 15:38:31
  author:        Alexandra Jahn, CU Boulder, alexandra.jahn@colorado.edu
  title:         Annual timeseries of freshwater data from the CESM Low W...
  description:   Annual mean Freshwater (FW) fluxes and storage relative ...
  data_structure: The data structure is |Ensemble member | Time (in years)...
```

1. How many dimensions does the `runoff_annual` variable have? What are the coordinates for the second dimension

of this variable?

Answer

We can see in the object viewer that the `runoff_annual` variable has two dimensions: `time` and `member`, in that order. We can also access the dimensions by calling:

```
fw_data.runoff_annual.dims
```

The second dimension is `member`. Near the top of the object viewer, under coordinates, we can see that that member's coordinates is an array from 1 to 11. We can directly see this array by calling:

```
fw_data.member
```

2. Select the values for the second member of the `netPrec_annual` variable.

Answer

```
member2 = fw_data.netPrec_annual.sel(member=2)
```

3. What is the maximum value of the second member of the `netPrec_annual` variable in the time period 2022 to 2100? [Hint](#).

Answer

Based on our previous answer, this maximum is:

```
x_max = member2.loc[2022:2100].max()  
x_max.item()
```

Notice we had to use `item` to transform the array into a number.

5.4 Tabular Data and NetCDF

Undoubtedly, tabular data is one of the most popular data formats. In this last section, we will discuss the relation between tabular data and the NetCDF data format and how to transform a `pandas.DataFrame` into an `xarray.Dataset`.

5.4.1 Tabular to NetCDF

We assume our starting point is tabular data that meets the criteria for **tidy data**, which means:

- Each column holds a different variable.
- Each row holds a different observation.

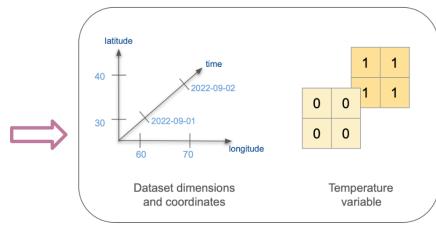
Take, for example, this tidy data subset from our exercise about weather stations measuring temperature:

time	latitude	longitude	temperature
2022-09-01	30	60	0
2022-09-01	30	70	0
2022-09-01	40	60	0
2022-09-01	40	70	0
2022-09-02	30	60	1
2022-09-02	30	70	1
2022-09-02	40	60	1
2022-09-02	40	70	1

To understand how this will transform into NetCDF format, we first need to identify which columns will act as dimensions and which as variables. We can also think of the values of the dimension columns as the coordinates in the `xarray.Dataset`. The diagram below shows how these columns transform into variables, dimensions, and coordinates.

time	latitude	longitude	temperature
2022-09-01	30	60	0
2022-09-01	30	70	0
2022-09-01	40	60	0
2022-09-01	40	70	0
2022-09-02	30	60	1
2022-09-02	30	70	1
2022-09-02	40	60	1
2022-09-02	40	70	1

Tabular data format



NetCDF data format

Tabular formats like csv do not offer an intrinsic way to encode attributes for the dimensions or variables, this is why we don't see any attributes in the resulting NetCDF data. One of the most significant advantages of NetCDF is its self-describing properties.

5.4.2 pandas to xarray

What does the previous example look like when working with `pandas` and `xarray`?

Let's work with a csv file containing the previous temperature measurements. Essentially, we need to read this file as a `pandas.DataFrame` and then use the `pandas.DataFrame.to_xarray()` method, taking into account that the dimensions of the resulting `xarray.DataSet` will be formed using the index column(s) of the `pandas.DataFrame`. In this case, we know the first three columns will be our dimension columns, so we need to group them as a `multindex` for the `pandas.DataFrame`. We can do this by using the `index_col` argument directly when we read in the csv file.

```
fp = os.path.join(os.getcwd(), 'netcdf_temp_data.csv')

# specify columns representing dimensions
dimension_columns = [0,1,2]

# read file
temp = pd.read_csv(fp, index_col=dimension_columns)
temp
```

```
/opt/hostedtoolcache/Python/3.9.16/x64/lib/python3.9/site-packages/IPython/core/formatters.py:1
```

```
In future versions `DataFrame.to_latex` is expected to utilise the base implementation of `Sty
```

time	longitude	latitude	temperature
2022-09-01	30	60	0
		70	0
	40	60	0
		70	0
2022-09-02	30	60	1
		70	1
	40	60	1
		70	1

And this is our resulting `xarray.Dataset`:

```
temp.to_xarray()
```

```
<xarray.Dataset>
Dimensions:      (time: 2, longitude: 2, latitude: 2)
Coordinates:
 * time          (time) object '2022-09-01' '2022-09-02'
 * longitude     (longitude) int64 30 40
 * latitude      (latitude) int64 60 70
Data variables:
    temperature   (time, longitude, latitude) int64 0 0 0 0 1 1 1 1
```

For further reading and examples about switching between `pandas` and `xarray` you can visit the following:

- [xarray's Frequently Asked Questions](#)
- [xarray's documentation about working with pandas](#)
- [pandas.DataFrame.to_xarray documentation](#)

6 Parallelization with Dask

- Become familiar with the Dask processing workflow:
 - What are the client, scheduler, workers, and cluster
 - Understand delayed computations and “lazy” evaluation
 - Obtain information about computations via the Dask dashboard
- Learn to load data and specify partition/chunk sizes of `dask.arrays/dask.dataframes`
- Integrate `xarray` and `rioxarray` with Dask for geospatial computations
- Share best practices and resources for further reading

6.1 Introduction

Dask is a library for parallel computing in Python. It can scale up code to use your personal computer’s full capacity or distribute work in a cloud cluster. By mirroring APIs of other commonly used Python libraries, such as Pandas and NumPy, Dask provides a familiar interface that makes it easier to parallelize your code. In this lesson, we will get acquainted with Dask’s way of distributing and evaluating computations and some of its most commonly used objects.



6.2 Dask Cluster

We can deploy a **Dask cluster** on a single machine or an actual cluster with multiple machines. The cluster has three main components for processing computations in parallel. These are the *client*, the *scheduler* and the *workers*.

- When we code, we communicate directly with the **client**, which is responsible for submitting tasks to be executed to the scheduler.
- After receiving the tasks from the client, the **scheduler** determines how tasks will be distributed among the workers and coordinates them to process tasks in parallel.
- Finally, the **workers** compute tasks and store and return computations results. Workers can be threads, processes, or separate machines in a cluster. Here you can read more about what are [threads and processes](#) and some [best practices for selecting one or the other](#).

To interact with the client and generate tasks that can be processed in parallel we need to use [Dask objects](#) to read and process our data.

6.2.1 Setting up a Local Cluster

First, we need to set up a Dask configuration file. In your `~/.config` directory, create a directory called `dask` (`mkdir -p ~/.config/dask`). In that directory create a file called `dask.yaml` (`touch dask.yaml`). In that file, insert just one line:

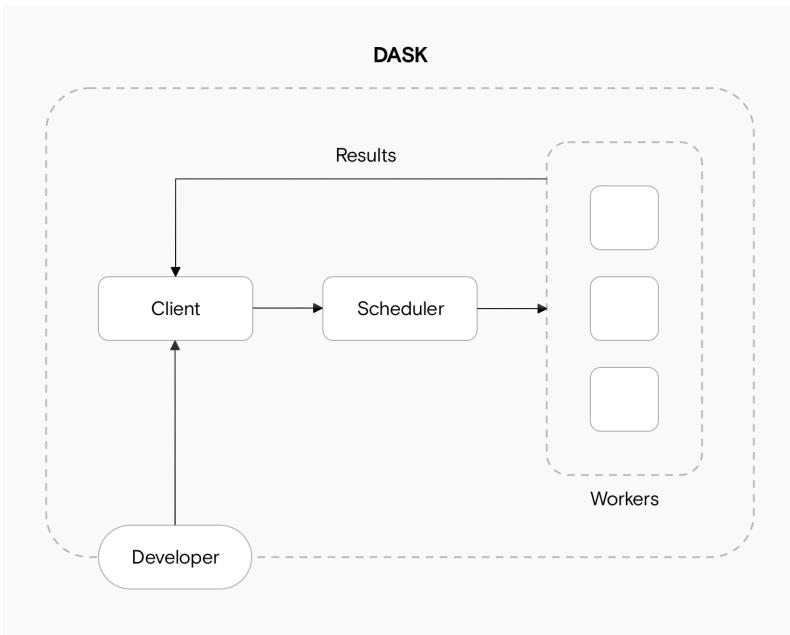


Figure 6.1: M. Schmitt, *Understanding Dask Architecture: The Client, Scheduler and Workers*

```
temporary-directory: ~/tmp
```

Then create that temporary directory (`mkdir ~/tmp`). This will prevent permissions problems while working on the shared cluster.

We can create a **local cluster** as follows:

```
from dask.distributed import LocalCluster, Client

cluster = LocalCluster(n_workers=5, memory_limit=0.1, processes=True)
cluster
```

The screenshot shows a web interface for a LocalCluster. At the top, there are tabs for 'Status' and 'Scaling'. Below the tabs, the cluster name 'LocalCluster' is displayed, along with its ID '26418090'. A 'Dashboard' link is provided, and the number of 'Workers' is listed as '5'. Key performance metrics are shown in two columns: 'Total threads: 90' and 'Total memory: 62.89 GiB'. The 'Status' is noted as 'running', and 'Using processes' is set to 'False'. A 'Scheduler Info' link is also present.

And then we create a client to connect to our cluster:

```
client = Client(cluster)  
client
```

The screenshot shows a 'Client' status page. It displays the client's ID 'Client-13bb684b-306b-11ed-9694-fbcb071c665b' and its 'Connection method: Cluster object'. The 'Cluster type' is identified as 'distributed.LocalCluster'. Below this, a 'Cluster Info' section is expanded, showing details for a 'LocalCluster' with ID '26418090'. These details include the dashboard URL, worker count (5), total threads (90), total memory (62.89 GiB), current status (running), and whether processes are used (False). A 'Scheduler Info' link is also present.

6.2.2 Dask Dashboard

We chose to use the Dask cluster in this lesson instead of the default Dask scheduler to take advantage of the **cluster dashboard**, which offers live monitoring of the performance and progress of our computations. You can learn more about different [Dask clusters here](#).

As seen in the images above, when we set up a cluster we can see the cluster dashboard address by looking at either the client or the cluster. In this example the dashboard address is

<http://128.111.85.28:8787/status>. When we go that address in a web browser we can see the dashboard's main page. This page shows diagnostics about:

- the cluster's and individual worker's memory usage,
- number of tasks being processed by each worker,
- individual tasks being processed across workers, and
- progress towards completion of individual tasks.

There's much to say about interpreting the Dask dashboard's diagnostics. We recommend this documentation to understand the [basics of the dashboard diagnostics](#) and [this video](#) as a deeper dive into the dashboard's functions.

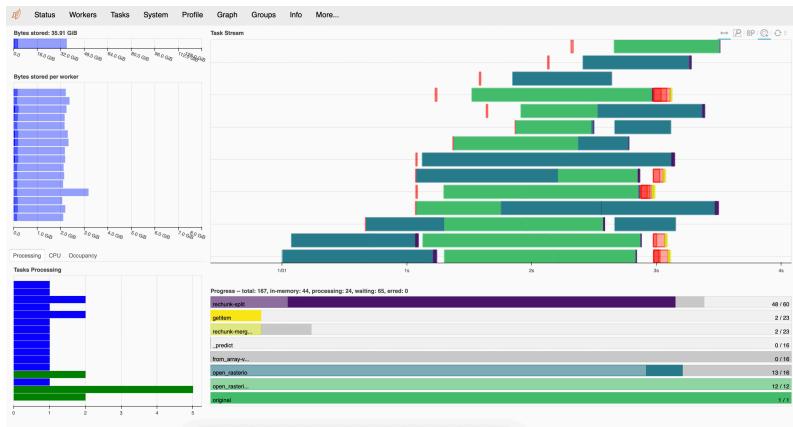


Figure 6.2: A Dask dashboard.

6.3 dask.dataframes

When we analyze tabular data, we usually start our analysis by loading it into memory as a Pandas DataFrame. But what if this data does not fit in memory? Or maybe our analyzes crash because we run out of memory. These scenarios are typical entry points into parallel computing. In such cases, Dask's scalable alternative to a Pandas DataFrame is the `dask.dataframe`.

A `dask.dataframe` comprises many `pd.DataFrame`s, each containing a subset of rows of the original dataset. We call each of these pandas pieces a **partition** of the `dask.dataframe`.

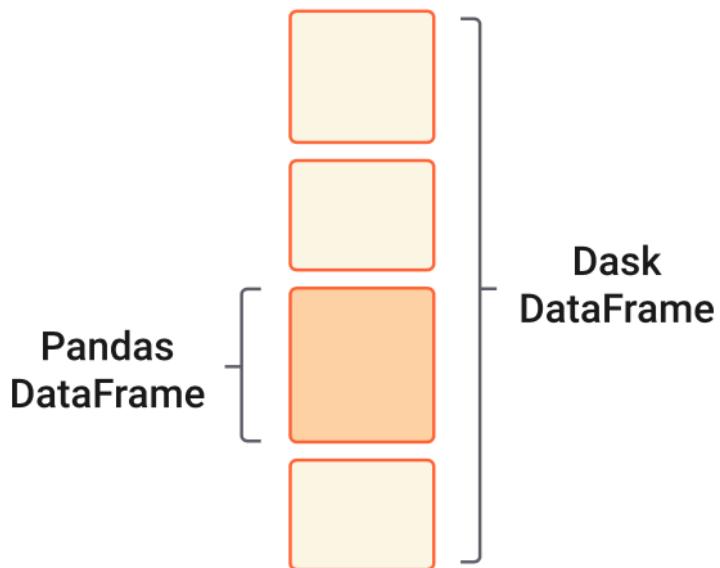


Figure 6.3: Dask Array design ([dask documentation](#))

6.3.1 Reading a csv

To get familiar with `dask.dataframes`, we will use tabular data of soil moisture measurements at six forest stands in north-eastern Siberia. The data has been collected since 2014 and is archived at the Arctic Data Center ([Loranty & Alexander, doi:10.18739/A24B2X59C](#)). Just as we did in the previous lesson, we will download the data using the `requests` package and the data's URL obtained from the Arctic Data Center.

```
import os
import urllib
```

```
import dask.dataframe as dd

url = 'https://arcticdata.io/metacat/d1/mn/v2/object/urn%3Auuid%3A27e4043d-75eb-4c4f-9427-0d

msg = urllib.request.urlretrieve(url, "dg_soil_moisture.csv")
```

In the Arctic Data Center metadata we can see this file is 115 MB. To import this file as a `dask.dataframe` with more than one partition, we need to specify the size of each partition with the `blocksize` parameter. In this example, we will split the data frame into six partitions, meaning a block size of approximately 20 MB.

```
fp = os.path.join(os.getcwd(), 'dg_soil_moisture.csv')
df = dd.read_csv(fp, blocksize = '20MB' , encoding='ISO-8859-1')
df
```

i Note

7 Encoding?

About the encoding parameter: If we try to import the file directly, we will receive an `UnicodeDecodeError`. We can run the following code to find the file's encoding and add the appropriate encoding to `dask.dataframe.read_csv`.

```
import chardet
fp = os.path.join(os.getcwd(), 'dg_soil_moisture.csv')
with open(fp, 'rb') as rawdata:
    result = chardet.detect(rawdata.read(100000))
result
```

Notice that we cannot see any values in the data frame. This is because Dask has not really loaded the data. It will wait until

we explicitly ask it to print or compute something to do so.

However, we can still do `df.head()`. It's not costly for memory to access a few data frame rows.

```
df.head(3)
```

7.0.1 Lazy Computations

The application programming interface (API) of a `dask.dataframe` is a subset of the `pandas.DataFrame` API. So if you are familiar with pandas, many of the core `pandas.DataFrame` methods directly translate to `dask.dataframes`. For example:

```
averages = df.groupby('year').mean()
averages
```

Notice that we cannot see any values in the resulting data frame. A major difference between `pandas.DataFrame`s and `dask.dataframes` is that `dask.dataframes` are “lazy”. This means an object will queue transformations and calculations without executing them until we explicitly ask for the result of that chain of computations using the `compute` method. Once we run `compute`, the scheduler can allocate memory and workers to execute the computations in parallel. This kind of **lazy evaluation** (or **delayed computation**) is how most Dask workloads work. This varies from **eager evaluation** methods and functions, which start computing results right when they are executed.

Before calling `compute` on an object, open the Dask dashboard to see how the parallel computation is happening.

```
averages.compute()
```

7.1 `dask.array`

Another common object we might want to parallelize is a NumPy array. The equivalent Dask object is the `dask.array`,

which coordinates many NumPy arrays that may live on disk or other machines. Each of these NumPy arrays within the `dask.array` is called a **chunk**. Choosing how these chunks are arranged within the `dask.array` and their size can significantly affect the performance of our code. Here you can find more information about [chunks](#).

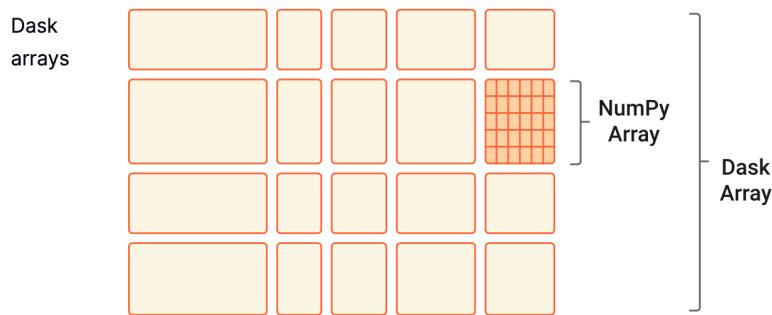


Figure 7.1: Dask Array design ([dask documentation](#))

In this short example we will create a 200x500 `dask.array` by specifying chunk sizes of 100x100.

```
import numpy as np  
  
import dask.array as da  
  
data = np.arange(100_000).reshape(200, 500)  
a = da.from_array(data, chunks=(100, 100))
```

Computations for `dask.arrays` also work lazily. We need to call `compute` to trigger computations and bring the result to memory.

```
a.mean()  
  
a.mean().compute()
```

7.2 Dask and xarray

In the future, it might be more common having to read some big array-like dataset (like a high-resolution multiband raster) than creating one from scratch using NumPy. In this case, it can be useful to use the `xarray` module and its extender `rioxarray` together with Dask. In the previous lesson, *Data Structures and Formats for Large Data*, we explore how to use the `xarray` package to work with labelled arrays. `rioxarray` extends `xarray` with the `rio` accessor, which stands for “raster input and output”.

It is simple to wrap Dask around `xarray` objects. We only need to specify the number of chunks as an argument when we are reading in a dataset (see also [1]).

7.2.1 Open .tif file

As an example, let’s do a Normalized Difference Vegetation Index (NDVI) calculation using remote sensing imagery collected by aerial vehicles over northeastern Siberia ([Lorlanty, Forbath, Talucci, Alexander, DeMarco, et al. 2020. doi:10.18739/A2ZC7RV6H.](#)). The NDVI is an index commonly used to check if an area has live green vegetation or not. It can also show the difference between water, plants, bare soil, and human-made structures, among other things.

The NDVI is calculated using the near-infrared and red bands of the satellite image. The formula is

$$NDVI = \frac{NIR - Red}{NIR + Red}.$$

First, we download the data for the near-infrared (NIR) and red bands from the Arctic Data Center:

```
# download red band
url = 'https://arcticdata.io/metacat/d1/mn/v2/object/urn%3Auuid%3Aac25a399-b174-41c1-b6d3-09
msg = urllib.request.urlretrieve(url, "RU_ANS_TR2_FL005M_red.tif")
```

```
# download nir band
url = 'https://arcticdata.io/metacat/d1/mn/v2/object/urn%3Auuid%3A1762205e-c505-450d-90ed-d4

msg = urllib.request.urlretrieve(url, "RU_ANS_TR2_FL005M_nir.tif")
```

Because these are .tif files and have geospatial metadata, we will use `rioxarray` to read them. You can find more information about `rioxarray` [here](#).

To indicate we will open these .tif files with `dask.array`s as the underlying object to the `xarray.DataArray` (instead of a `numpy.array`), we need to specify either a shape or the size in bytes for each chunk. Both files are 76 MB, so let's have chunks of 15 MB to have roughly six chunks.

```
import rioxarray as rioxr

# read in the file
fp_red = os.path.join(os.getcwd(), "RU_ANS_TR2_FL005M_red.tif")
red = rioxr.open_rasterio(fp_red, chunks = '15MB')
```

We can see a lot of useful information here:

- There are eight chunks in the array. We were aiming for six, but this often happens with how Dask distributes the memory (76MB is not divisible by 6).
- There is geospatial information (transformation, CRS, resolution) and no-data values.
- There is an unnecessary dimension: a constant value for the band. So our next step is to squeeze the array to flatten it.

```
# getting rid of unnecessary dimension
red = red.squeeze()
```

Next, we read in the NIR band and do the same pre-processing:

```
# open data
fp_nir = os.path.join(os.getcwd(), "RU_ANS_TR2_FL005M_nir.tif")
nir = rioxr.open_rasterio(fp_nir, chunks = '15MB')

#squeeze
nir = nir.squeeze()
```

7.2.2 Calculating NDVI

Now we set up the NDVI calculation. This step is easy because we can handle xarrays and Dask arrays as NumPy arrays for arithmetic operations. Also, both bands have values of type float32, so we won't have trouble with the division.

```
ndvi = (nir - red) / (nir + red)
```

When we look at the NDVI we can see the result is another `dask.array`, nothing has been computed yet. Remember, Dask computations are lazy, so we need to call `compute()` to bring the results to memory.

```
ndvi_values = ndvi.compute()
```

And finally, we can see what these look like. Notice that `xarray` uses the value of the dimensions as labels along the x and y axes. We use `robust=True` to ignore the no-data values when plotting.

```
ndvi_values.plot(robust=True)
```

7.3 Best Practices

Dask is an exciting tool for parallel computing, but it may take a while to understand its nuances to make the most of it. There are many best practices and recommendations. These are some of the basic ones to take into consideration:

- For data that fits into RAM, pandas, and NumPy can often be faster and easier to use than Dask workflows. The simplest solution can often be the best.
- While Dask may have similar APIs to pandas and NumPy, there are differences, and not all the methods for the `pandas.DataFrame`s and `numpy.array`s translate in the same way (or with the same efficiency) to Dask objects. When in doubt, always read the documentation.
- Choose appropriate chunk and partition sizes and layouts. This is crucial to best use how the scheduler distributes work. You can read here about [best practices for chunking](#).
- Avoid calling `compute` repeatedly. It is best to group similar computations together and then compute once.

Further reading:

- A friendly article about [common dask mistakes](#)
- [General Dask best practices](#)
- [dask.dataframe best practices](#)
- [dask.array best practices](#)

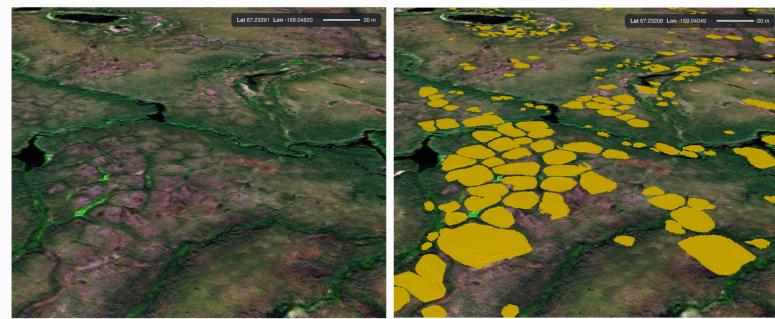
8 Group Project: Staging and Preprocessing

- Get familiarized with the overall group project workflow
- Write a `parsl` app that will stage and tile the IWP example data in parallel

8.1 Introduction

The Permafrost Discovery Gateway is an online platform for archiving, processing, analysis, and visualization of permafrost big imagery products to enable discovery and knowledge-generation. The PDG utilizes and makes available products derived from high resolution satellite imagery from the Polar Geospatial Center, Planet (3 meter resolution), Sentinel (10 meter resolution), Landsat (30 meter resolution), and MODIS (250 meter resolution). One of these products is a dataset showing Ice Wedge Polygons (IWP) that form in melting permafrost.

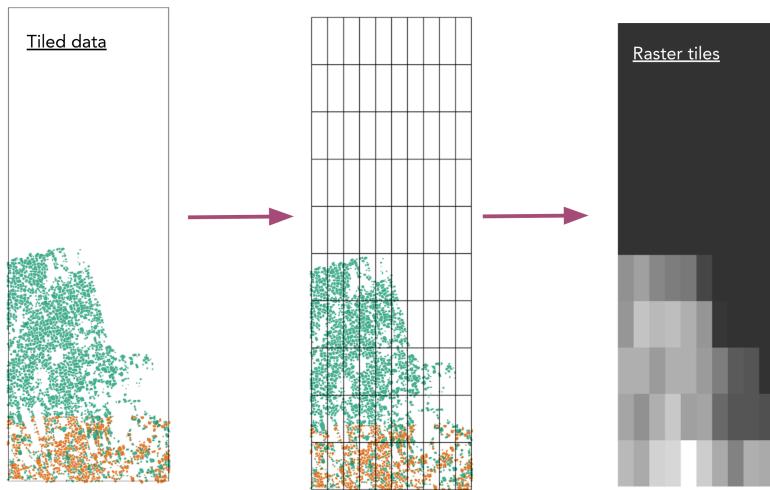
Ice wedges form as a result of cycles of thermal contraction during long-term melting and freezing cycles of permafrost. They can form very distinctive enclosed geometries that are clearly visible in high-resolution satellite images. The PDG team is using advanced analysis and computational tools to process high-resolution satellite imagery and automatically detect where ice wedge polygons form. Below is an example of a satellite image (left) and the detected ice wedge polygons in geospatial vector format (right) of that same image.



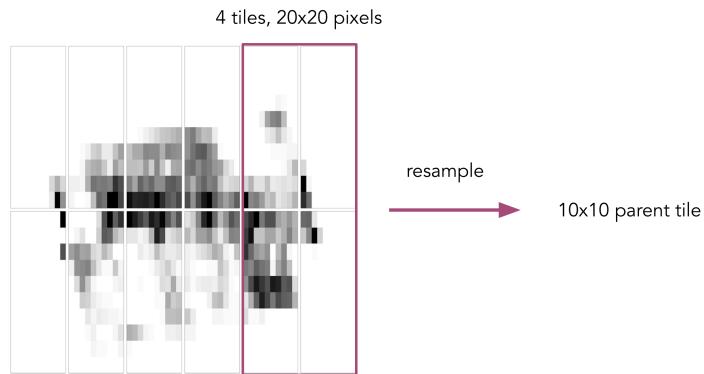
In the group project, we are going to use a subset of the high-resolution dataset of these detected ice wedge polygons in order to learn some of the reproducible, scalable techniques that will allow us to process it. Our workflow will start with a set of large geopackage files that contain the detected ice wedge polygons. These files all have irregular extents due to the variation in satellite coverage, clouds, etc. Our first processing step will take these files and “tile” them into smaller files which have regular extents.



In step two of the workflow, we will take those regularly tiled geopackage files and rasterize them. The files will be regularly gridded, and a summary statistic will be calculated for each grid cell (such as the proportion of pixel area covered by polygons).



In the final step of the workflow, we will take the raster files and resample them to create a set of raster tiles at different resolutions. This last step is what will enable us to visualize our raster data dynamically, such that we look at lower resolutions when very zoomed out (and high resolution data would take too long to load), and higher resolution data when zoomed in and the extent is smaller.



8.2 Staging and Tiling

Today we will undertake the first step of the workflow, staging and tiling the data.



In your fork of the [scalable-computing-examples](#) repository, open the Jupyter notebook in the `group-project` directory called `session-06.ipynb`. This workbook will serve as a skeleton for you to work in. It will load in all the libraries you need, including a few helper functions we wrote for the course, show an example for how to use the method on one file, and then lays out blocks for you and your group to fill in with code that will run that method in parallel.

In your small groups, work together to write the solution, but everyone should aim to have a working solution on their own fork of the repository. In other words, everyone should type out the solution themselves as part of the group effort. Writing the code out yourself (even if others are contributing to the content) is a great way to get “mileage” as you develop these skills.

9 Software Design I: Functions and Concurrency

The goal of this module is to gain some useful approaches to software design and modularity that will help with building scalable, portable, and reusable code. We will cover several key aspects of software design for concurrency:

- Functions as objects
- Global variables
- Pure functions
- Task dependencies
- Locks, deadlocks, and race conditions

9.1 Why functions?

⌚ DRY: Don't Repeat Yourself



By creating small functions that handle only one logical task and do it well, we quickly gain:

- Improved understanding
- Reuse via decomposing tasks into bite-sized chunks
- Improved error testing
- Improved concurrency

When writing functions that are to be used in concurrent programming, it is best to keep them short and focused on a single, well-defined task. This enables you test the function thoroughly, and reuse it in multiple contexts. It also enables you to more easily debug what is happening in the function when trying to understand parallel execution.

9.2 Functions as objects

Functions are first-class objects in Python (and many other languages). This has some real benefits for and implications for parallel programming. Because a function is an object, it means that it can be 1) stored in a variable, and 2) passed as an argument to another function. We saw that in the module on pleasantly parallel codes when we used `ThreadPoolExecutor.map()`, which takes a function and an iterable object as arguments. Let's check out how you can use and manipulate functions as objects. First, let's define a simple function, assign it to another variable, and then use both:

```
def double(x):
    return 2*x

# also assign the function to the `twotimes` variable
twotimes = double
type(twotimes)

function
```

Note that when we print it to screen, we see that `prod` is of type `function`, and when we use the two instances, we get identical results:

```
print(double(7))
print(twotimes(7))
print(double(5) == twotimes(5))
```

```
14  
14  
True
```

This representation of a function as an object comes in handy when we want to invoke a function in multiple different contexts, such as in a parallel execution environment via a `map()` function.

```
list(map(twotimes, [2,3,4]))
```

```
[4, 6, 8]
```

This works because the function `twotimes` can be passed to `map` and executed from within `map`. When you execute a function that is passed in via an argument, it is called **function composition**. We can easily illustrate this by creating some function and passing it to a wrapper function to be executed:

```
def some_function():
    print("Ran some_function")

def wrapper(func_to_run):
    print("Ran wrapper")
    func_to_run()
    print("Finished wrapper")

wrapper(some_function)
```

```
Ran wrapper
Ran some_function
Finished wrapper
```

i Note

Note how we passed the `some_function` as a variable name without the parentheses.

i Decorators

This approach to function composition is exactly what is used by [Python decorator](#) functions.

9.3 Global variables

When executing a function, the variables that are in scope to that function are local, meaning that the presence of another variable with the same name in another scope will not affect a calculation. For example:

```
def do_task():
    x = 10

    x = 5
do_task()
print(x)
```

5

However, if that same variable is declared as global inside the function, then the assignment will have global impact on the value of the parent scope:

```
def do_task():
    global x
    x = 10

    x = 5
do_task()
print(x)
```

10

So, you can see that writing a function that uses global variables can have effects outside of the scope of the function call. This

can have drastic consequences on concurrent code, as the order in which function calls are made when operating concurrently are not deterministic, and so the impact of global variables will also not be deterministic.

A related issue arises when code in a function depends on its enclosing namespace, such as when a function is defined inside of another function. When resolving a variable, python first looks in the Local namespace, and then in the Enclosing namespace, Global namespace, and Built-in namespace. So, even if a variable is not defined locally, it might still be resolved by one of the other namespaces in surprising ways.

```
a = 3
def do_stuff(b):
    return a*b

do_stuff(6)
```

18

9.4 Pure functions

A **pure function** is a function that depends only on its input arguments, and it has no *side effects*. In other words, a pure function returns the same value if called repeatedly with the same arguments. Pure functions are particularly amenable to concurrency. For example, the `double(x)` function above is a pure function, because in all cases calling `double(2)` will always return 4.

In contrast, a non-pure function is a function in which the return value may change if the function is called repeatedly, typically because it depends on some particular state that affects the outcome but is not part of the input arguments. For example, the `time.time()` function returns different values based on the current state of the system clock.

Using a global variable in a function creates a side-effect that makes it an impure function, as would other operations that modify an external state variable.

! Important

Pure functions make writing concurrent code much easier.

9.5 Task dependencies

Task dependencies occur when one task in the code depends on the results of another task or computation in the code.

9.6 Locks

Locks are a mechanism to manage access to a resource so that multiple threads can access the resource. By adding locks to an otherwise parallel process, we introduce a degree of serial execution to the locked portion of the process. Basically, each thread can only access the resource when it has the lock, and only one lock is given out at a time. Take this example of what happens without locking:

```
from concurrent.futures import ProcessPoolExecutor
import time

def hello(i):
    print(i, 'Hello')
    print(i, 'world')

executor = ProcessPoolExecutor()
futures = [executor.submit(hello, i) for i in range(3)]
for future in futures:
    future.result()
```

0

1

Hello

Hello

0

1

world

world

2

Hello

2

```
world
```

You can see that the results come back in a semi-random order, and the call to `sleep` creates a delay between printing the two words, which means that the three messages get jumbled when printed. To fix this, we can introduce a lock from the `multiprocessing` package.

```
from concurrent.futures import ProcessPoolExecutor
import time
import multiprocessing

def hello(i, lock):
    with lock:
        print(i, 'Hello')
        print(i, 'world')

lock = multiprocessing.Manager().Lock()
executor = ProcessPoolExecutor()
futures = [executor.submit(hello, i, lock) for i in range(3)]
for future in futures:
    future.result()
```

```
0
```

```
Hello
```

```
0
```

world

1

Hello

1

world

2

Hello

2

world

The lock is generated and then passed to each invocation of the hello function. Using `with` triggers the use of the context manager for the lock, which allows the manager to synchronize use of the lock. This ensures that only one process can be printing at the same time, ensuring that the outputs are properly ordered.

 **Warning**

Synchronizing with a Lock turns a parallel process back into a serial process, at least while the lock is in use. So use with care lest you lose all benefits of concurrency.

9.7 Race conditions

Race conditions occur when two tasks execute in parallel, but produce different results based on which task finishes first. Ensuring that results are correct under different timing situations requires careful testing.

9.8 Deadlocks

Deadlocks occur when two concurrent tasks block on the output of the other. Deadlocks cause parallel programs to lock up indefinitely, can be difficult to track down, and will often require the program to be killed.

9.9 Further reading

- [With Statement Context Managers](#)

10 Spatial and Image Data Using GeoPandas

- Manipulating raster data with `rasterio`
- Manipulating vector data with `geopandas`
- Working with raster and vector data together

10.1 Introduction

In this lesson, we'll be working with geospatial raster and vector data to do an analysis on vessel traffic in south central Alaska. If you aren't already familiar, geospatial vector data consists of points, lines, and/or polygons, which represent locations on the Earth. Geospatial vector data can have differing geometries, depending on what it is representing (eg: points for cities, lines for rivers, polygons for states.) Raster data uses a set of regularly gridded cells (or pixels) to represent geographic features.

Both geospatial vector and raster data have a coordinate reference system, which describes how the points in the dataset relate to the 3-dimensional spheroid of Earth. A coordinate reference system contains both a datum and a projection. The datum is how you georeference your points (in 3 dimensions!) onto a spheroid. The projection is how these points are mathematically transformed to represent the georeferenced point on a flat piece of paper. All coordinate reference systems require a datum. However, some coordinate reference systems are “unprojected” (also called geographic coordinate systems). Coordinates in latitude/longitude use a geographic (unprojected) coordinate system. One of the most commonly used geographic coordinate systems is WGS 1984.

Coordinate reference systems are often referenced using a short-hand 4 digit code called an EPSG code. We'll be working with two coordinate reference systems in this lesson with the following codes:

- 3338: Alaska Albers
- 4326: WGS84 (World Geodetic System 1984), used in GPS

In this lesson, we are going to take two datasets:

- [Alaskan commercial salmon fishing statistical areas](#)
- [North Pacific and Arctic Marine Vessel Traffic Dataset](#)

and use them to calculate the total distance travelled by ships within each fishing area.

The high level steps will be

- read in the datasets
- reproject them so they are in the same projection
- extract a subset of the raster and vector data using a bounding box
- turn each polygon in the vector data into a raster mask
- use the masks to calculate the total distance travelled (sum of pixels) for each fishing area

10.2 Pre-processing raster data

First we need to load in our libraries. We'll use `geopandas` for vector manipulation, `rasterio` for raster manipulation.

First, we'll use `requests` to download the ship traffic raster from [Kapsar et al.](#). We grab a one month slice from August, 2020 of a coastal subset of data with 1km resolution. To get the URL in the code chunk below, you can right click the download button for the file of interest and select “copy link address.”

```
import urllib

url = 'https://arcticdata.io/metacat/d1/mn/v2/object/urn%3Auuid%3A6b847ab0-9a3d-4534-bf28-3a'
```

```
msg = urllib.request.urlretrieve(url, "Coastal_2020_08.tif")
```

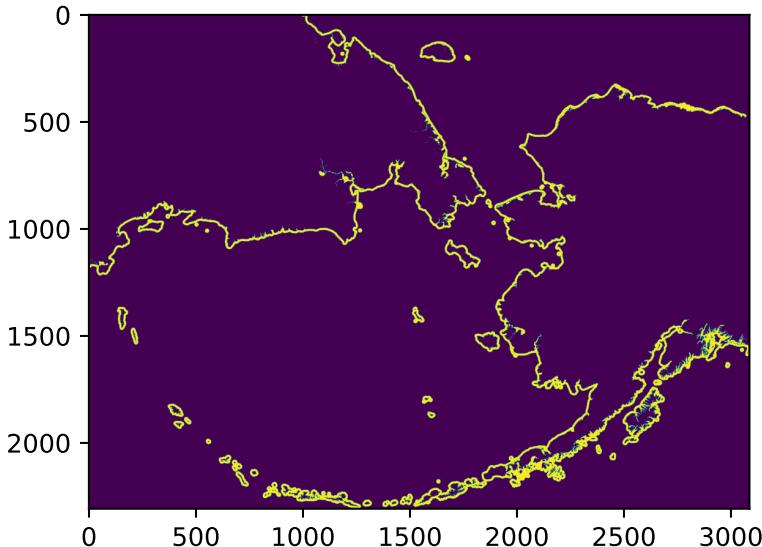
Using `rasterio`, open the raster file, plot it, and look at the metadata. We use the `with` here as a context manager. This ensures that the connection to the raster file is closed and cleaned up when we are done with it.

```
import rasterio
import matplotlib.pyplot as plt

with rasterio.open("Coastal_2020_08.tif") as ship_con:
    # read in raster (1st band)
    ships = ship_con.read(1)
    ships_meta = ship_con.profile

plt.imshow(ships)
print(ships_meta)

{'driver': 'GTiff', 'dtype': 'float32', 'nodata': -3.3999999521443642e+38, 'width': 3087, 'height': 2087, 'affine': Affine(0.000328125, 0.0, -999.9687691991521, 0.0, 2711703.104608573), 'tiled': False, 'compress': 'lzw', 'interlace': 'line'}
```

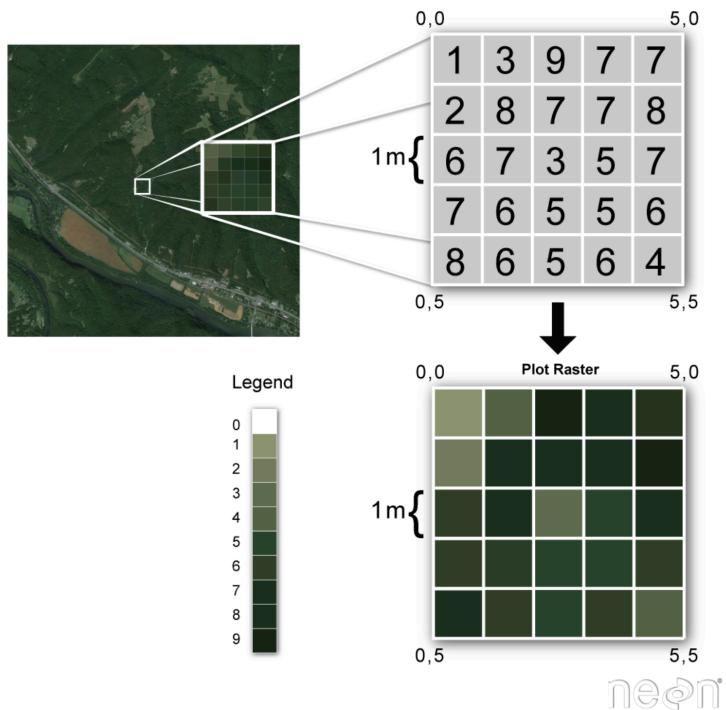


You'll notice that we are saving two objects here, `ships` and

`ships_meta`. Looking at the types of these two objects is useful to understand what `rasterio` is doing.

```
type(ships)  
  
numpy.ndarray  
  
type(ships_meta)  
  
rasterioprofiles.Profile
```

The `ships` object is a `numpy` array, while the `ships_meta` is a special `rasterio` class called `Profile`. To understand why the raster data is represented as an array, and what that profile object is, let's look into what raster data are, exactly.



Source: Leah A. Wasser, Megan A. Jones, Zack Brym, Kristina Riemer, Jason Williams, Jeff Hollister, Mike Smorul. Raster

00: Intro to Raster Data in R. National Evological Observatory Network (NEON).

The upper left panel of the figure above shows some satellite imagery data. These data are in raster format, which when you zoom in, you can see consist of regularly gridded pixels, each of which contains a value. When we plot these data, we can assign a color map to the pixel values, which generates the image we see. The data themselves, though, are just an n-dimensional grid of numbers. Another way we might describe this is...an array! So, this is why raster data is represented in python using a `numpy` array.

This is all great, and the array of values is a lot of information, but there are some key items that are missing. This array isn't imaginary, it represents a physical space on this earth, so where is all of that contextual information? The answer is in the `rasterio` profile object. This object contains all of the metadata needed to interpret the raster array. Here is what our `ships_meta` contains:

```
'driver': 'GTiff',
'dtype': 'float32',
'nodata': -3.3999999521443642e+38,
'width': 3087,
'height': 2308,
'count': 1,
'crs': CRS.from_epsg(3338),
'transform': Affine(999.7994153462766, 0.0, -2550153.29233849, 0.0, -999.9687691991521, 2711700),
'tiled': False,
'compress': 'lzw',
'interleave': 'band'}
```

This object gives us critical information, like the CRS of the data, the no data value, and the transform. The transform is what allows us to move from image pixel (row, column) coordinates to and from geographic/projected (x, y) coordinates. The transform and the CRS are critically important, and related. If the CRS are instructions for how the coordinates can be represented in space and on a flat surface (in the case of projected

coordinate systems), then the transform describes how to locate the raster array positions in the correct coordinates given by the CRS.

Note that since the array and the profile are in separate objects it is easy to lose track of one of them, accidentally overwrite it, etc. Try to adopt a naming convention that works for you because they usually need to work together in geospatial operations.

10.3 Pre-processing vector data

Now download a vector shapefile of commercial fishing districts in Alaska.

```
url = 'https://knb.ecoinformatics.org/knb/d1/mn/v2/object/urn%3Auuid%3A7c942c45-1539-4d47-b4  
msg = urllib.request.urlretrieve(url, "Alaska_Commercial_Salmon_Boundaries.gpkg")
```

Read in the data using `geopandas`.

```
import geopandas as gpd  
  
comm = gpd.read_file("Alaska_Commercial_Salmon_Boundaries.gpkg")
```

Note the “pandas” in the library name “geopandas.” Our `comm` object is really just a special type of pandas data frame called a geodataframe. This means that in addition to any geospatial stuff we need to do, we can also just do regular `pandas` things on this data frame.

For example, we can get a list of column names (there are a lot!)

```
comm.columns.values  
  
array(['OBJECTID', 'GEOMETRY_START_DATE', 'GEOMETRY_END_DATE',  
       'STAT_AREA', 'STAT_AREA_NAME', 'FISHERY_GROUP_CODE',  
       'GIS_SERIES_NAME', 'GIS_SERIES_CODE', 'REGION_CODE',  
       'REGISTRATION_AREA_NAME', 'REGISTRATION_AREA_CODE'],
```

```
'REGISTRATION_AREA_ID', 'REGISTRATION_LOCATION_ABBR',
'MANAGEMENT_AREA_NAME', 'MANAGEMENT_AREA_CODE', 'DISTRICT_NAME',
'DISTRICT_CODE', 'DISTRICT_ID', 'SUBDISTRICT_NAME',
'SUBDISTRICT_CODE', 'SUBDISTRICT_ID', 'SECTION_NAME',
'SECTION_CODE', 'SECTION_ID', 'SUBSECTION_NAME', 'SUBSECTION_CODE',
'SUBSECTION_ID', 'COAR_AREA_CODE', 'CREATOR', 'CREATE_DATE',
'EDITOR', 'EDIT_DATE', 'COMMENTS', 'STAT_AREA_VERSION_ID',
'Shape_Length', 'Shape_Area', 'geometry'], dtype=object)
```

We can also look at the head of the data frame:

```
comm.head
```

		OBJECTID	GEOMETRY_START_DATE	GEOMETRY_END_DATE	STAT_A
0	12	1975-01-01 00:00:00+00:00	NaT	33461	
1	13	1975-01-01 00:00:00+00:00	NaT	33462	
2	18	1978-01-01 00:00:00+00:00	NaT	33431	
3	19	1980-01-01 00:00:00+00:00	NaT	33442	
4	20	1980-01-01 00:00:00+00:00	NaT	33443	
..
860	959	NaT	NaT	19241	
861	960	NaT	NaT	19242	
862	961	1994-01-01 00:00:00+00:00	NaT	19245	
863	962	NaT	NaT	19250	
864	963	1994-01-01 00:00:00+00:00	NaT	18252	
		STAT_AREA_NAME FISHERY_GROUP_CODE \			
0	Tanana River mouth to Kantishna River		B		
1	Kantishna River to Wood River		B		
2	Toklik to Cottonwood Point		B		
3	Right Bank, Bishop Rock to Illinois Creek		B		
4	Left Bank, Cone Point to Illinois Creek		B		
..
860	Kaliakh River		B		
861	Tsiu River		B		
862	Midtimber River		B		
863	Seal River		B		
864	Middle Italio		B		
	GIS_SERIES_NAME GIS_SERIES_CODE REGION_CODE	REGISTRATION_AREA_NAME \			

0	Salmon	B	3	Yukon Area
1	Salmon	B	3	Yukon Area
2	Salmon	B	3	Yukon Area
3	Salmon	B	3	Yukon Area
4	Salmon	B	3	Yukon Area
..

860	Salmon	B	1	Southeastern Alaska Area
861	Salmon	B	1	Southeastern Alaska Area
862	Salmon	B	1	Southeastern Alaska Area
863	Salmon	B	1	Southeastern Alaska Area
864	Salmon	B	1	Southeastern Alaska Area

	COAR_AREA_CODE	CREATOR	CREATE_DATE	\
0	...	YU	Evelyn Russel	2006-03-26 00:00:00+00:00
1	...	YU	Evelyn Russel	2006-03-26 00:00:00+00:00
2	...	YL	Evelyn Russel	2006-03-26 00:00:00+00:00
3	...	YU	Evelyn Russel	2006-03-26 00:00:00+00:00
4	...	YU	Evelyn Russel	2006-03-26 00:00:00+00:00
..
860	...	A2	Evelyn Russel	2006-03-26 00:00:00+00:00
861	...	A2	Evelyn Russel	2006-03-26 00:00:00+00:00
862	...	A2	Evelyn Russel	2006-03-26 00:00:00+00:00
863	...	A2	Evelyn Russel	2006-03-26 00:00:00+00:00
864	...	A2	Sabrina Larsen	2017-05-05 00:00:00+00:00

	EDITOR	EDIT_DATE	\
0	Sabrina Larsen	2017-02-02 00:00:00+00:00	
1	Sabrina Larsen	2017-02-02 00:00:00+00:00	
2	Sabrina Larsen	2017-02-02 00:00:00+00:00	
3	Sabrina Larsen	2017-02-02 00:00:00+00:00	
4	Sabrina Larsen	2017-02-02 00:00:00+00:00	
..
860	Sabrina Larsen		NaT
861	Sabrina Larsen		NaT
862	Sabrina Larsen		NaT
863	Sabrina Larsen		NaT
864	NaN		NaN

	COMMENTS	STAT_AREA_VERSION_ID	\
0	Yukon District, 6 Subdistrict and 6-A Section ...		NaN
1	Yukon District, 6 Subdistrict and 6-B Section ...		NaN

```

2      Yukon District and 3 Subdistrict until 1/1/1980      NaN
3                                         NaN
4                                         NaN
..
860                                         NaN
861                                         NaN
862                                         NaN
863                                         NaN
864                                         NaN

      Shape_Length Shape_Area \
0      4.610183   0.381977
1      3.682421   0.321943
2      2.215641   0.198740
3      9.179852   0.382788
4      9.500826   0.378262
..
860     ...       ...
861     0.223565   0.000408
862     0.030506   0.000006
863     0.019805   0.000012
864     0.096016   0.000238
864     0.016237   0.000006

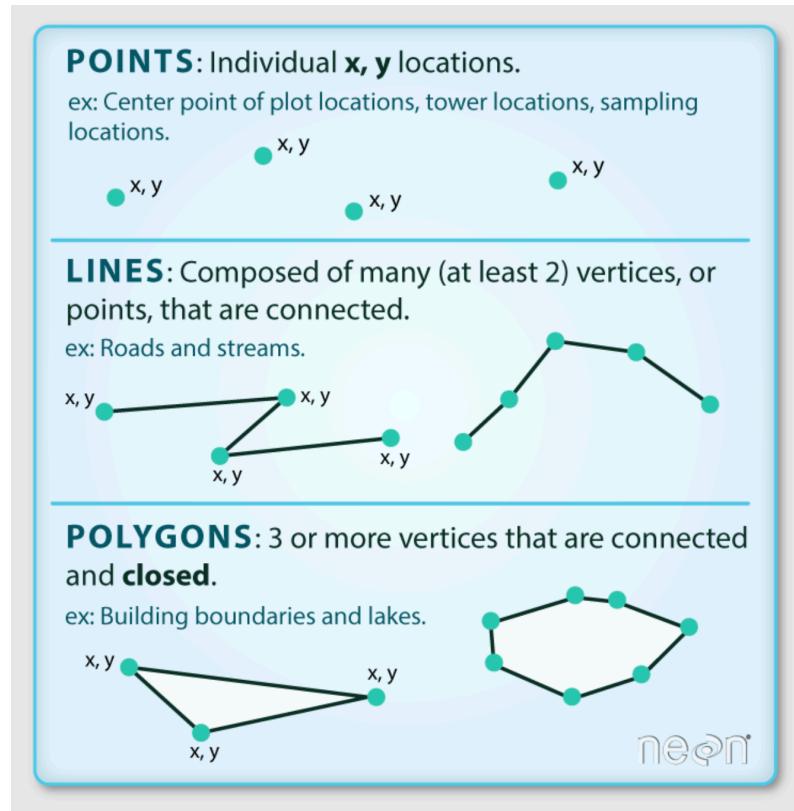
                           geometry
0  MULTIPOLYGON (((-151.32805 64.96913, -151.3150...
1  MULTIPOLYGON (((-149.96255 64.70518, -149.9666...
2  MULTIPOLYGON (((-161.39853 61.55463, -161.4171...
3  MULTIPOLYGON (((-153.15234 65.24944, -153.0761...
4  MULTIPOLYGON (((-152.99905 65.17027, -152.9897...
..
860  MULTIPOLYGON (((-142.90787 60.09177, -142.9051...
861  MULTIPOLYGON (((-143.00416 60.07711, -143.0046...
862  MULTIPOLYGON (((-143.28504 60.05800, -143.2861...
863  MULTIPOLYGON (((-143.49701 60.04832, -143.5083...
864  MULTIPOLYGON (((-139.15063 59.30748, -139.1489...

```

[865 rows x 37 columns]>

Note the existence of the `geometry` column. This is where the actual geospatial points that comprise the vector data are stored, and this brings up the important difference between

raster and vector data - while raster data is regularly gridded at a specific resolution, vector data are just points in space.



Source: Leah A. Wasser, Megan A. Jones, Zack Brym, Kristina Riemer, Jason Williams, Jeff Hollister, Mike Smorul. Raster 00: Intro to Raster Data in R. National Evological Observatory Network (NEON).

The diagram above shows the three different types of geometries that geospatial vector data can take, points, lines or polygons. Whatever the geometry type, the geometry information (the x,y points) is stored in the column named `geometry` in the geopandas data frame. In this example, we have a dataset containing polygons of fishing districts. Each row in the dataset corresponds to a district, with unique attributes (the other columns in the dataset), and it's own set of points defining the boundaries of the district, contained in the `geometry` column.

```
comm['geometry'][:5]
```

```
/opt/hostedtoolcache/Python/3.9.16/x64/lib/python3.9/site-packages/IPython/core/formatters.py:
```

```
In future versions `DataFrame.to_latex` is expected to utilise the base implementation of `Sty
```

geometry

```
0 MULTIPOLYGON (((-151.32805 64.96913, -151.3150...  
1 MULTIPOLYGON (((-149.96255 64.70518, -149.9666...  
2 MULTIPOLYGON (((-161.39853 61.55463, -161.4171...  
3 MULTIPOLYGON (((-153.15234 65.24944, -153.0761...  
4 MULTIPOLYGON (((-152.99905 65.17027, -152.9897...
```

So, now we know where our x,y points are, where is all of the other information like the `crs`? With vector data, all of this information is contained within the geodataframe. We can access the `crs` attribute on the data frame and print it like so:

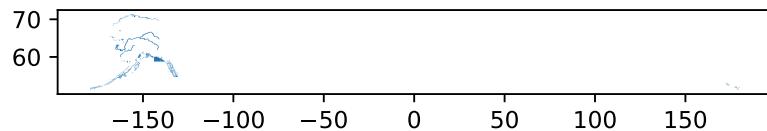
```
comm.crs
```

```
<Geographic 2D CRS: EPSG:4326>  
Name: WGS 84  
Axis Info [ellipsoidal]:  
- Lat[north]: Geodetic latitude (degree)  
- Lon[east]: Geodetic longitude (degree)  
Area of Use:  
- name: World.  
- bounds: (-180.0, -90.0, 180.0, 90.0)  
Datum: World Geodetic System 1984 ensemble  
- Ellipsoid: WGS 84  
- Prime Meridian: Greenwich
```

Now that we know a little about what we are working with, let's get this data ready to analyze. First, we can make a plot of it just to see what we have.

```
comm.plot()
```

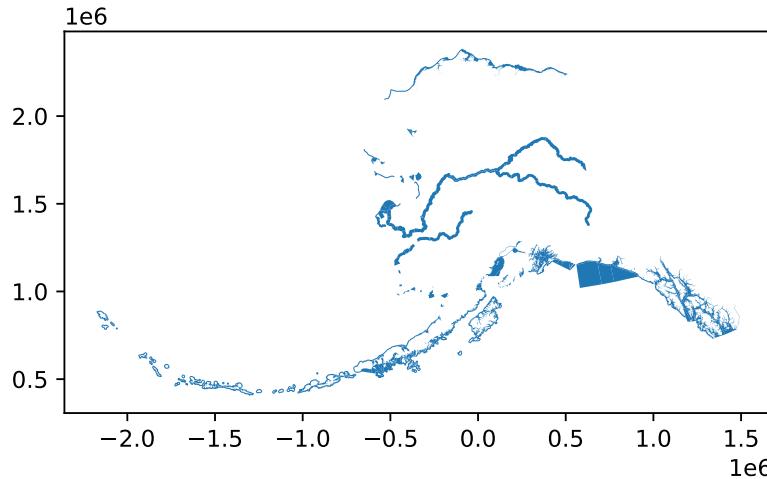
```
<AxesSubplot:>
```



This plot doesn't look so good. Turns out, these data are in WGS 84 (EPSG 4326), as opposed to Alaska Albers (EPSG 3338), which is what our raster data are in. To make pretty plots, and allow our raster data and vector data to be analyzed together, we'll need to reproject the vector data into 3338. To do this, we'll use the `to_crs` method on our `comm` object, and specify as an argument the projection we want to transform to.

```
comm_3338 = comm.to_crs("EPSG:3338")  
comm_3338.plot()
```

```
<AxesSubplot:>
```



Much better!

10.4 Crop data to area of interest

For this example, we are only interested in south central Alaska, encompassing Prince William Sound, Cook Inlet, and Kodiak. Our raster data is significantly larger than that, and the vector data is statewide. So, as a first step we might want to crop our data to the area of interest.

First, we'll need to create a bounding box. We use the `box` function from `shapely` to create the bounding box, then create a `GeoDataFrame` from the points, and finally convert the WGS 84 coordinates to the Alaska Albers projection.

```
from shapely.geometry import box

coord_box = box(-159.5, 55, -144.5, 62)

coord_box_df = gpd.GeoDataFrame(
    crs = 'EPSG:4326',
    geometry = [coord_box]).to_crs("EPSG:3338")
```

Now, we can read in raster again cropped to bounding box. We use the `mask` function from `rasterio.mask`. Note that we apply this to the connection to the raster file (`with rasterio.open(...)`), then update the metadata associated with the raster, because the `mask` function requires as its first `dataset` argument a dataset object opened in `r` mode.

```
import rasterio.mask
import numpy as np

with rasterio.open("Coastal_2020_08.tif") as ship_con:
    shipc_arr, shipc_transform = rasterio.mask.mask(ship_con,
                                                    coord_box_df["geometry"],
                                                    crop=True)

    shipc_meta = ship_con.meta
    # select just the 2-D array (by default a 3-D array is returned even though we only have
    shipc_arr = shipc_arr[0,:,:]
    # turn the no-data values into NaNs.
    shipc_arr[shipc_arr == ship_con.nodata] = np.nan
```

```
shipc_meta.update({"driver": "GTiff",
                   "height": shipc_arr.shape[0],
                   "width": shipc_arr.shape[1],
                   "transform": shipc_transform,
                   "compress": "lzw"})
```

Now we'll do a similar task with the vector data. In this case, we use a spatial join. The join will be an inner join, and select only rows from the left side (our fishing districts) that are **within** the right side (our bounding box). I chose this method as opposed to a clipping type operation because it ensures that we don't end up with any open polygons at the boundaries of our box, which could cause problems for us down the road.

```
comm_clip = gpd.sjoin(comm_3338,
                      coord_box_df,
                      how='inner',
                      predicate='within')
```

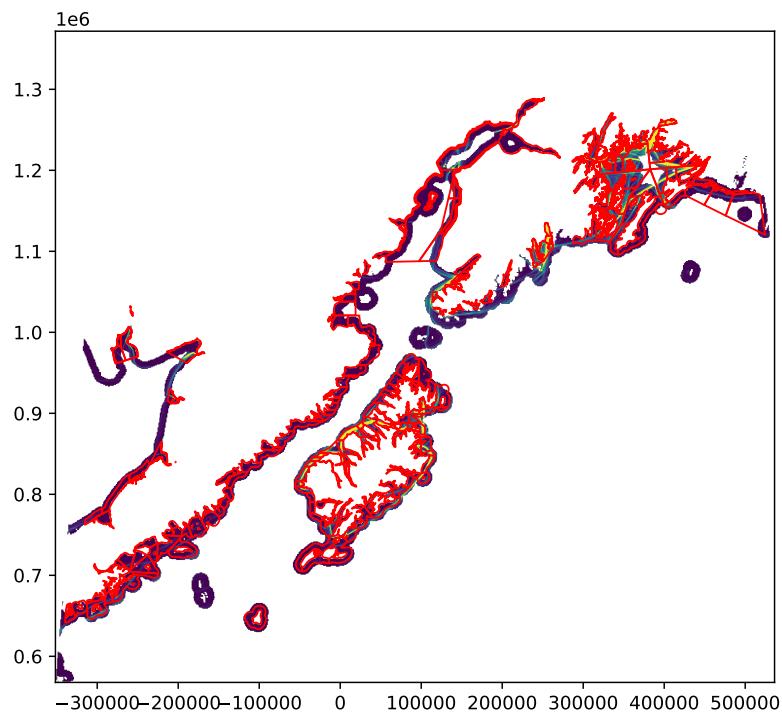
10.4.1 Check extents

Now let's look at the two cropped datasets overlayed on each other to ensure that the extents look right.

```
import rasterio.plot

# set up plot
fig, ax = plt.subplots(figsize=(7, 7))
# plot the raster
rasterio.plot.show(shipc_arr,
                    ax=ax,
                    vmin = 0,
                    vmax = 50000,
                    transform = shipc_transform)
# plot the vector
comm_clip.plot(ax=ax, facecolor='none', edgecolor='red')
```

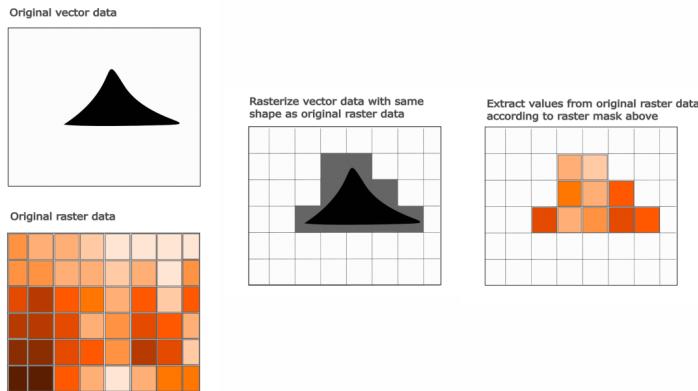
```
<AxesSubplot:>
```



10.5 Calculate total distance per fishing area

In this step, we rasterize each polygon in the shapefile, such that pixels in or touching the polygon get a value of 1, and pixels not touching it get a value of 0. Then, for each polygon, we extract the indices of the raster array that are equal to 1. We then extract the values of these indicies from the original ship traffic raster data, and calculate the sum of the values over all of those pixels.

Here is a simplified diagram of the process:



10.5.0.1 Zonal statistics over one polygon

Let's look at how this works over just one fishing area first. We use the `rasterize` method from the `features` module in `rasterio`. This takes as arguments the data to rasterize (in this case the 40th row of our dataset), and the shape and transform the output raster will take. We also set the `all_touched` argument to true, which means any pixel that touches a boundary of our vector will be burned into the mask.

```
from rasterio import features

r40 = features.rasterize(comm_clip['geometry'][40].geoms,
                        out_shape=shipc_arr.shape,
                        transform=shipc_meta['transform'],
                        all_touched=True)
```

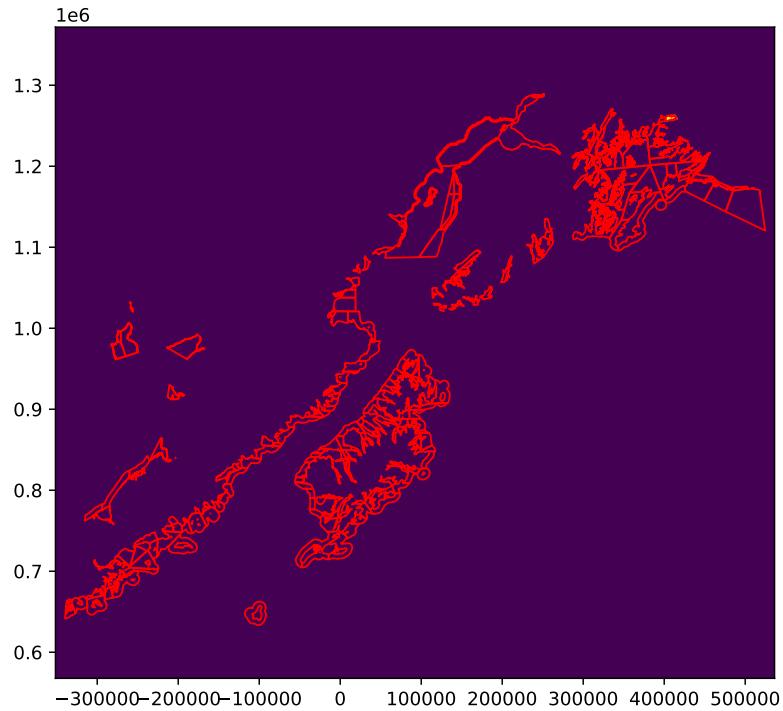
If we have a look at a plot of our rasterized version of the single fishing district, we can see that instead of a vector, we now have a raster showing the rasterized district (with pixel values of 1) and any area not in the district has a pixel value of 0.

```
# set up plot
fig, ax = plt.subplots(figsize=(7, 7))
```

```
# plot the raster
rasterio.plot.show(r40,
                    ax=ax,
                    vmin = 0,
                    vmax = 1,
                    transform = shipc_meta['transform'])

# plot the vector
comm_clip.plot(ax=ax, facecolor='none', edgecolor='red')
```

<AxesSubplot:>



A quick call to `np.unique` shows our unique values are 0 or 1, which is what we expect.

```
np.unique(r40)

array([0, 1], dtype=uint8)
```

Finally, we need to know is the indices of the original raster where the fishing district is. We can use `np.where` to extract this information

```
r40_index = np.where(r40 == 1)
print(r40_index)
```

```
(array([108, 108, 108, 108, 108, 109, 109, 109, 109, 109, 109,
       109, 109, 110, 110, 110, 110, 110, 110, 110, 110, 110, 110,
       110, 110, 110, 111, 111, 111, 111, 111, 111, 111, 111, 111,
       111, 111, 111, 111, 111, 112, 112, 112, 112, 112, 112, 112,
       112, 112, 112, 112, 112, 113, 113, 113, 113, 113, 113, 113,
       113, 113, 113, 113, 113, 113, 113, 114, 114, 114, 114,
       114, 114, 114, 114, 114, 115, 115, 115, 115, 115, 115, 115,
       115, 115, 115, 115, 115, 116]), array([759, 760, 762,
       764, 765, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762,
       763, 764, 765, 766, 753, 754, 755, 756, 757, 758, 759, 760,
       761, 762, 763, 764, 765, 766, 767, 753, 754, 755, 756, 757,
       758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 753, 754,
       755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766,
       767, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763,
       764, 765, 766, 767, 753, 754, 755, 756, 757, 758, 759, 760,
       761, 762, 763, 764, 765, 766, 767, 753, 754, 755, 756, 757,
       758, 759, 760, 761, 762, 763, 753, 754, 755, 756, 757, 758,
       759]))
```

In the last step, we'll use these indices to extract the values of the data from the fishing raster, and sum them to get a total distance travelled.

```
np.nansum(shipc_arr[r40_index])
```

```
14369028.0
```

Now that we know the individual steps, let's run this over all of the districts. First we'll create an `id` column in the vector data frame. This will help us track unique fishing districts later.

10.5.0.2 Zonal statistics over all polygons

```
comm_clip['id'] = range(0,len(comm_clip))
```

For each district (with `geometry` and `id`), we run the `features.rasterize` function. Then we calculate the sum of the values of the shipping raster `r_array` based on the indicies in the raster where the district is located.

```
distance_dict = {}
for geom, idx in zip(comm_clip['geometry'], comm_clip['id']):
    rasterized = features.rasterize(geom.geoms,
                                    out_shape=shipc_arr.shape,
                                    transform=shipc_meta['transform'],
                                    all_touched=True)

    r_index = np.where(rasterized == 1)
    distance_dict[idx] = np.nansum(shipc_arr[r_index])
```

Now we just create a data frame from that dictionary, and join it to the vector data using `pandas` operations.

```
import pandas as pd

# create a data frame from the result
distance_df = pd.DataFrame.from_dict(distance_dict,
                                      orient='index',
                                      columns=['distance'])

# extract the index of the data frame as a column to use in a join and convert distance to km
distance_df['id'] = distance_df.index
distance_df['distance'] = distance_df['distance']/1000
```

Now we join the result to the original geodataframe.

```
# join the sums to the original data frame
res_full = comm_clip.merge(distance_df,
                           on = "id",
                           how = 'inner')
```

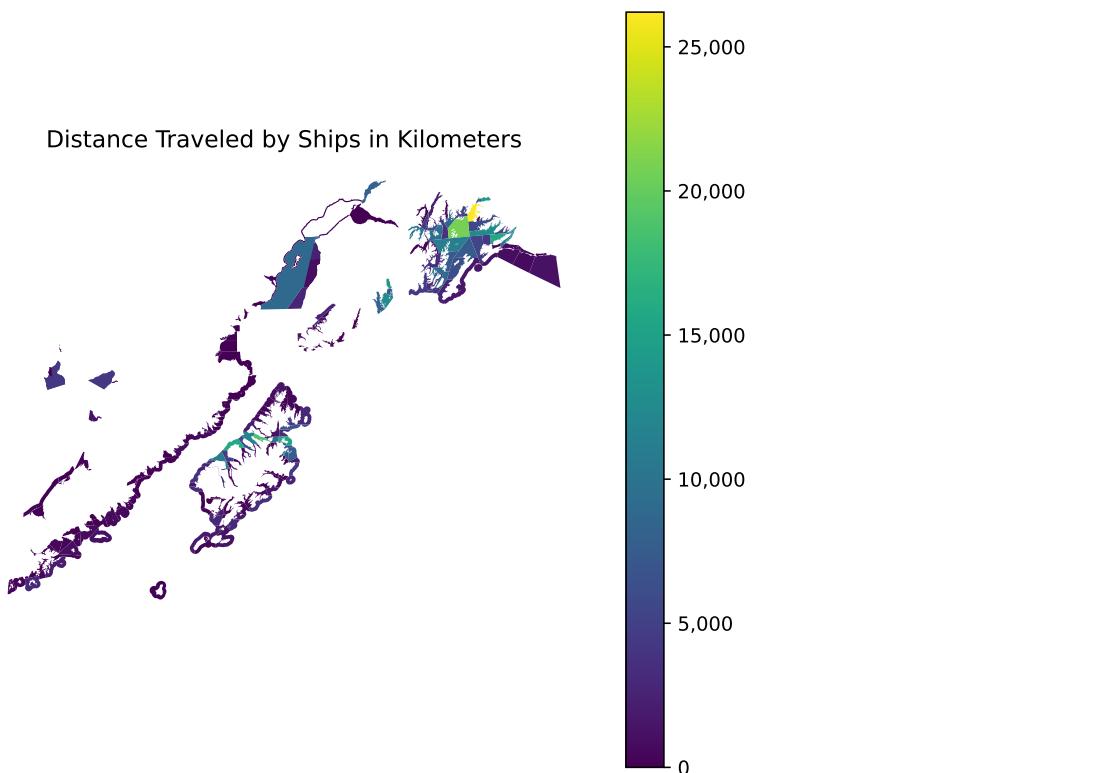
Finally, we can plot our result!

```
import matplotlib.ticker
fig, ax = plt.subplots(figsize=(7, 7))
```

```

ax = res_full.plot(column = "distance", legend = True, ax = ax)
fig = ax.figure
label_format = '{:.0f}'
cb_ax = fig.axes[1]
ticks_loc = cb_ax.get_yticks().tolist()
cb_ax.yaxis.set_major_locator(matplotlib.ticker.FixedLocator(ticks_loc))
cb_ax.set_yticklabels([label_format.format(x) for x in ticks_loc])
ax.set_axis_off()
ax.set_title("Distance Traveled by Ships in Kilometers")
plt.show()

```



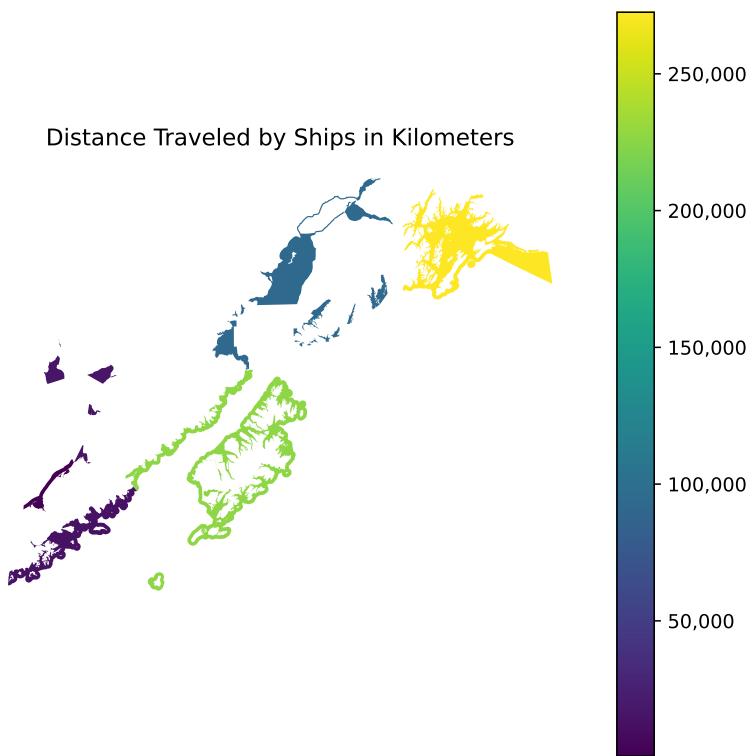
From here we can do any additional `geopandas` operations we might be interested in. For example, what if we want to calculate the total distance by registration area (a superset of fishing district). We can do that using `dissolve` from `geopandas`.

```
reg_area = res_full.dissolve(by = "REGISTRATION_AREA_NAME",
                             aggfunc = 'sum')
```

Let's have a look at the same plot as before, but this time over our aggregated data.

```
fig, ax = plt.subplots(figsize=(7, 7))

ax = reg_area.plot(column = "distance", legend = True, ax = ax)
fig = ax.figure
label_format = '{:,.0f}'
cb_ax = fig.axes[1]
ticks_loc = cb_ax.get_yticks().tolist()
cb_ax.yaxis.set_major_locator(matplotlib.ticker.FixedLocator(ticks_loc))
cb_ax.set_yticklabels([label_format.format(x) for x in ticks_loc])
ax.set_axis_off()
ax.set_title("Distance Traveled by Ships in Kilometers")
plt.show()
```



10.6 Summary

We covered a lot of ground here, so let's recap some of the high level points:

- Raster data consist of regularly gridded values, and can be represented in python as an array
- Vector data consist of any number of points, that might be connected, and is represented in python as a geodataframe
- We can do geospatial operations like changing the projection or cropping the data to a particular extent on both raster and vector data
- You can use vector data to help analyze raster data (and vice versa!) by rasterizing the vector data and using numpy operations on the resulting array.

11 Billions of Ice Wedge Polygons

The Permafrost Discovery Gateway (PDG) is an online platform for creation, archiving, processing, analysis, and visualization of permafrost big imagery products to enable discovery and knowledge-generation. The new online scientific gateway will make information of changing permafrost conditions available throughout the Arctic by providing access to very high resolution satellite data products and new visualization tools that will allow exploration and discovery for researchers, educators, and the public at large.

Knowledge discovery through AI, big imagery, high-performance computing (HPC) is just starting to be realized in Arctic permafrost science. Chandi Witharana and the PDG team have developed a novel high- performance image analysis framework – Mapping application for Arctic Permafrost Land Environment (MAPLE) that enables the integration of operational-scale GeoAI capabilities into Arctic permafrost modeling. Interoperability across heterogeneous HPC systems, optimal usage of computational resource, and extensibility of mapping workflows are some of the key design goals of MAPLE. MAPLE was recently deployed across the Arctic tundra to map ice-wedge polygons from thousands of Maxar imagery. This mapping exercise has produced the first pan-Arctic ice-wedge polygon map, which consists of more than one billion individual ice-wedge polygons.

11.1 Billions of Ice Wedge Polygons

Chandi's slides are available to download [here](#)

11.2 Developing Science Gateways with Clowder for the Permafrost Discover Gateway

Luigi's slides are available to download [here](#)

12 Group Project: Data Processing

In your fork of the [scalable-computing-examples](#) repository, open the Jupyter notebook in the `group-project` directory called `session-13.ipynb`. This workbook will serve as a skeleton for you to work in. It will load in all the libraries you need, including a few helper functions we wrote for the course, show an example for how to use the method on one file, and then lays out blocks for you and your group to fill in with code that will run that method in parallel.

In your small groups, work together to write the solution, but everyone should aim to have a working solution on their own fork of the repository. In other words, everyone should type out the solution themselves as part of the group effort. Writing the code out yourself (even if others are contributing to the content) is a great way to get “mileage” as you develop these skills.

13 Data Ethics for Scalable Computing

- Review FAIR and CARE Principles, and their relevance to data ethics
- Examine how ethical considerations are shared and considered at the Arctic Data Center
- Discuss ethical considerations in machine learning

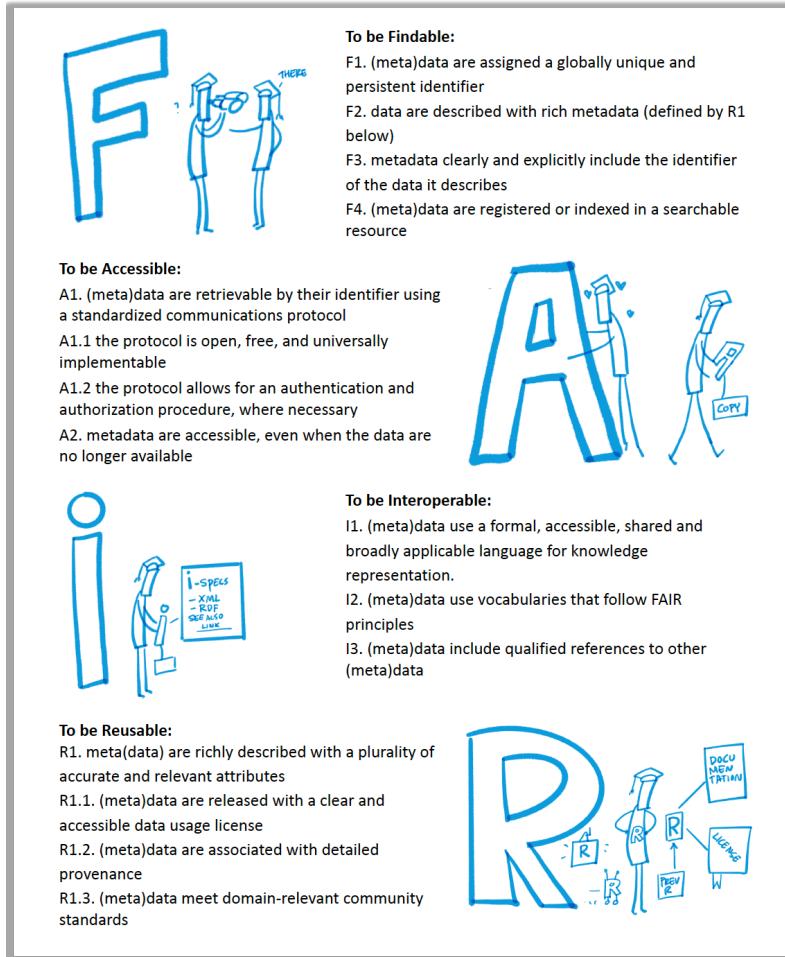
13.1 Intro to Data Ethics



To recap, the Arctic Data Center is an openly-accessible data repository and the data published through the repository is open for anyone to reuse, subject to conditions of the license (at the Arctic Data Center, data is released under one of two licenses: CC-0 Public Domain and CC-By Attribution 4.0). In facilitating use of data resources, the data stewardship community has converged on principles surrounding best practices for open data management.

Two principles that the Arctic Data Center explicitly adopts are FAIR Principles (Findable, Accessible, Interoperable, and Reproducible) and CARE Principles for Indigenous Governance (Collective Benefit, Authority to Control, Responsibility, Ethics).

FAIR and CARE principles are relevant in the context of data ethics for multiple reasons. FAIR speaks to how metadata is managed, stored, and shared.



FAIR principles and open science are overlapping concepts, but are distinctive from one another. Open science supports a culture of sharing research outputs and data, and FAIR focuses on how to prepare the data. **The FAIR principles place emphasis on machine readability, “distinct from peer initiatives that focus on the human scholar”** (Wilkinson et al 2016) and as such, do not fully engage with sensitive data considerations and with Indigenous rights and interests (Research Data Alliance International Indigenous Data Sovereignty Interest Group, 2019). Metadata

can be FAIR but not open. For example, sensitive data (data that contains personal information) may not be appropriate to share, however sharing the anonymized metadata that is easily understandable will reduce research redundancy.



Research has historically perpetuated colonialism and represented extractive practices, meaning that the research results were not mutually beneficial. These issues also related to how data was owned, shared, and used. **To address issues like these, the Global Indigenous Data Alliance (GIDA) introduced CARE Principles for Indigenous Data Governance to support Indigenous data sovereignty.** CARE Principles speak directly to how the data is stored and shared in the context of Indigenous data sovereignty. CARE Principles stand for:

- *Collective Benefit* - Data ecosystems shall be designed and function in ways that enable Indigenous Peoples to derive benefit from the data
- *Authority to Control* - Indigenous Peoples' rights and interests in Indigenous data must be recognized and their authority to control such data be empowered. Indigenous data governance enables Indigenous Peoples and governing bodies to determine how Indigenous Peoples, as well as Indigenous lands, territories, resources, knowledges and geographical indicators, are represented and identified within data.

- *Responsibility* - Those working with Indigenous data have a responsibility to share how those data are used to support Indigenous Peoples' self-determination and collective benefit. Accountability requires meaningful and openly available evidence of these efforts and the benefits accruing to Indigenous Peoples.
- *Ethics* - Indigenous Peoples' rights and wellbeing should be the primary concern at all stages of the data life cycle and across the data ecosystem. To many, the FAIR and CARE principles are viewed by many as complementary: CARE aligns with FAIR by outlining guidelines for publishing data that contributes to open-science and at the same time, accounts for Indigenous Peoples' rights and interests.

13.2 Ethics at the Arctic Data Center

Transparency in data ethics is a vital part of open science. Regardless of discipline, various ethical concerns are always present, including professional ethics such as plagiarism, false authorship, or falsification of data, to ethics regarding the handling of animals, to concerns relevant to human subjects research. As the primary repository for the Arctic program of the National Science Foundation, the Arctic Data Center accepts Arctic data from all disciplines. Recently, a new submission feature was released which asks researchers to describe the ethical considerations that are apparent in their research. This question is asked to all researchers, regardless of disciplines.

Sharing ethical practices openly, similar in the way that data is shared, enables deeper discussion about data management practices, data reuse, sensitivity, sovereignty and other considerations. Further, such transparency promotes awareness and adoption of ethical practices.

Inspired by CARE Principles for Indigenous Data Governance (Collective Benefit, Authority to Control, Responsibility, Ethics) and FAIR Principles (Findable, Accessible, Interoperable, Reproducible), we include a space in the data submission process for researchers to describe their ethical research

practices. These statements are published with each dataset, and the purpose of these statements is to promote greater transparency in data collection and to guide other researchers. For more information about the ethical research practices statement, check out this blog.

To help guide researchers as they write their ethical research statements, we have listed the following ethical considerations that are available on our website. The concerns are organized first by concerns that should be addressed by all researchers, and then by discipline.

Consider the following ethical considerations that are relevant for your field of research.

13.2.1 Ethical Considerations for all Arctic Researchers



Research Planning

1. Were any permits required for your research?
2. Was there a code of conduct for the research team decided upon prior to beginning data collection?
3. Was institutional or local permission required for sampling?
4. What impact will your research have on local communities or nearby communities (meaning the nearest community within a 100 mile radius)?

Data Collection

5. Were any local community members involved at any point of the research process, including study site identification, sampling, camp setup, consultation or synthesis?
6. Were the sample sites near or on Indigenous land or communities?

Data Sharing and Publication

7. How were the following concerns accounted for: misrepresentation of results, misrepresentation of experience, plagiarism, improper authorship, or the falsification or data?
8. If this data is intended for publication, are authorship expectations clear for everyone involved? Other professional ethics can be found here

13.2.2 Archaeological and Paleontological Research



Research Planning

1. Were there any cultural practices relevant to the study site? If yes, how were these practices accounted for by the research methodologies.

Data Collection

2. Did your research include the removal or artifacts?
3. Were there any contingencies made for the excavation and return of samples after cleaning, processing, and analysis?

Data Sharing and Publication 4. Were the samples deposited to a physical repository? 5. Were there any steps taken to account for looting threats? Please explain why or why not?

13.2.3 Human Participation and Sensitive Data



Research Planning

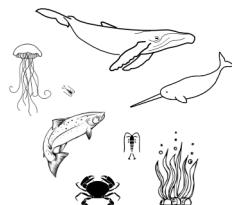
1. Please describe the institutional IRB approval that was required for this research.
2. Was any knowledge provided by community members?

Data Collection 3. Did participants receive compensation for their participation? 4. Were decolonization methods used?

Data Sharing and Publication

5. Have you shared this data with the community or participants involved?

13.2.4 Marines Sciences (e.g. Marine Biology Research)



Research Planning

1. Were any of the study sites or species under federal or local protection?

Data Collection

2. Were endangered, threatened, or otherwise special-status species collected?
3. How were samples collected? Please describe any handling practices used to collect data.
4. What safety measures were in place to keep researchers, research assistants, technicians, etc., out of harms way during research?
5. How were animal care procedures evaluated, and do they follow community norms for organismal care?

Data Sharing and Publication 6. Did the species or study area represent any cultural importance to local communities, or include culturally sensitive information? Please explain how you came to this conclusion and how any cultural sensitivity was accounted for.

13.2.5 Physical Sciences (e.g. Geology, Glaciology, and Ice Research)



Research Planning 1. Was any knowledge provided by community members, including information regarding the study site?

Data Collection

2. What safety measures were in place to keep researchers, research assistants, technicians, etc., out of harms way during research?
3. Were there any impacts to the environment/habitat before, during or after data collection?

Data Sharing and Publication

4. Is there any sensitive information including information on sensitive sites, valuable samples, or culturally sensitive information?

13.2.6 Plant and Soil Research



Research Planning

1. Were any of the study sites protected under local or federal regulation?
2. Was any knowledge provided by nearby community members, including information regarding the study site?

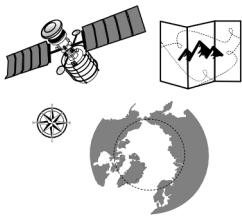
Data Collection

3. Did sample collection result in erosion of soil or other physical damage? If so, how was this addressed?
4. What safety measures were in place to keep researchers, research assistants, technicians, etc., out of harms way during research?

Data Sharing and Publication

5. Do any of the study sites or specimens represent culturally sensitive areas or species? Please explain how you came to this conclusion, and if yes, how was this accounted for?

13.2.7 Spatial Data



Research Planning

1. Were any land permits required for this research?

Data Collection

2. Were any data collected using citizen science or community participation?
3. If yes, were community members compensated for their time and made aware of their data being used and for what purpose?

Data Sharing and Publication

4. If data were ground-truthed, was institutional or local permissions required and/or obtained for land/property access?
5. Have you shared this data with the community or participants involved?
6. If location sensitive data was obtained (endangered/threatened flora & fauna location, archaeological and historical sites, identifiable ships, sensitive spatial information), how were the data desensitized?

13.2.8 Wildlife Sciences (e.g. Ecology and Biology Research)



Research Planning

1. Were any permits required for data sampling?

Data Collection

2. Were endangered, threatened, or otherwise special-status plants or animal species collected?
3. How were samples collected? Please describe any handling practices used to collect data.
4. What safety measures were in place to keep researchers, research assistants, technicians, etc., out of harms way during research?
5. How were animal care procedures evaluated, and do they follow community norms for organism care?

Data Sharing and Publication

6. Do any of the study sites or specimens represent culturally sensitive areas or species? Please explain how you came to this conclusion, and if yes, how was this accounted for?

Menti question:

1. Have you thought about any of the ethical considerations listed above before?
2. Were any of the considerations new or surprising?
3. Are there any for your relevant discipline that are missing?

13.3 Ethics in Machine Learning

Menti poll

1. What is your level of familiarity with machine learning
2. Have you thought about ethics in machine learning prior to this lesson?
3. Can anyone list potential ethical considerations in machine learning?

What comes to mind when considering ethics in machine learning?

The stories that hit the news are often of privacy breaches or biases seeping into the training data. Bias can enter at any point of the research project, from preparing the training data, designing the algorithms, to collecting and interpreting the data. When working with sensitive data, a question to also consider is how to deanonymize, anonymized data. A unique aspect to machine learning is how personal bias can influence the analysis and outcomes. A great example of this is the case of ImageNet.

13.3.1 ImageNet: A case study of ethics and bias in machine learning

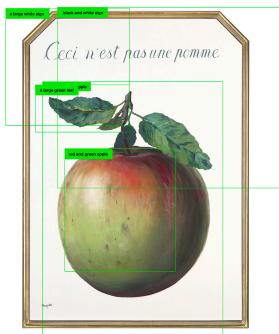


Image source: Kate Crawford and Trevor Paglen, “Excavating AI: The Politics of Training Sets for Machine Learning” (September 19, 2019).

ImageNet is a great example of how personal bias can enter machine learning through the training data. ImageNet was a

training data set of photos that was used to train image classifiers. The data set was initially created as a large collection of pictures, which were mainly used to identify objects, but some included images of people. The creators of the data set created labels to categorize the images, and through crowdsourcing, people from the internet labeled these images. (This example is from Kate Crawford and Trevor Paglen, “Excavating AI: The Politics of Training Sets for Machine Learning”, September 19, 2019).

13.3.1.1 Discussion:

1. Where are the two areas bias could enter this scenario?
2. Are there any ways that this bias could be avoided?
3. While this example is specific to images, can you think of any room for bias in your research?

13.4 References and Further Reading

Carroll, S.R., Herczog, E., Hudson, M. et al. (2021) Operationalizing the CARE and FAIR Principles for Indigenous data futures. *Sci Data* 8, 108 <https://doi.org/10.1038/s41597-021-00892-0>

Chen, W., & Quan-Haase, A. (2020) Big Data Ethics and Politics: Towards New Understandings. *Social Science Computer Review*. <https://journals.sagepub.com/doi/10.1177/0894439318810734>

Crawford, K., & Paglen, T. (2019) Excavating AI: The Politics of Training Sets for Machine Learning. <https://excavating.ai/>

Gray, J., & Witt, A. (2021) A feminist data ethics of care framework for machine learning: The what, why, who and how. *First Monday*, 26(12), Article number: 11833

Puebla, I., & Lowenberg, D. (2021) Recommendations for the Handling for Ethical Concerns Relating to the Publication of Research Data. *FORCE 11*. <https://force11.org/post/recommendations-for-the-handling-of-ethical-concerns-relating-to-the-publication-of-research-data/>

Research Data Alliance International Indigenous Data Sovereignty Interest Group. (2019). “CARE Principles for Indigenous Data Governance.” The Global Indigenous Data Alliance. GIDA-global.org

Wilkinson, M., Dumontier, M., Aalbersberg, I. et al. (2016) The FAIR Guiding Principles for scientific data management and stewardship. Sci Data 3, 160018. <https://doi.org/10.1038/sdata.2016.18>

Zwitter, A., Big Data ethics. (2014) Big Data and Society. DOI: 10.1177/2053951714559253

14 Google Earth Engine

- Understand what Google Earth Engine provides and its applications
- Learn how to search for, import, manipulate, and visualize Google Earth Engine Data
- Learn about some real-world applications of Google Earth Engine in the geosciences

14.1 Introduction

[Google Earth Engine](#) (GEE) is a geospatial processing platform powered by Google Cloud Platform. It contains over 30 years (and multiple petabytes) of satellite imagery and geospatial datasets that are continually updated and available instantly. Users can process data using Google Cloud Platform and built-in algorithms or by using the Earth Engine API, which is available in Python (and JavaScript) for anyone with an account (Earth Engine is free to use for research, education, and non-profit use).

Image Source: [Earth Engine Data Catalog](#)

So what's so exciting about platforms like GEE? Ryan Abernathey frames this nicely in his blogpost [Closed Platform vs. Open Architectures for Cloud-Native Earth System Analytics...](#)

- as Earth System data have gotten larger, the typical download-data-work-locally workflow is no longer always feasible
- those data are also produced and distributed by lots of different organizations (e.g. [NASA](#), [NOAA](#), [Copernicus](#))

GEE is just one of a number of cloud platform solutions developed for climate and geoscience research. Others include [Microsoft Planetary Computer](#), [Pangeo](#), & [Amazon Sustainability Data Initiative \(ADSI\)](#)

- researchers often need to apply a wide range of analytical methods to those data, ranging from simple stats to machine learning approaches

GEE offers web access (i.e. no need to download data to your computer) to an extensive catalog of analysis-ready geospatial data (from many different organizations) and scalable computing power via their cloud service, making global-scale analyses and visualizations possible for anyone with an account ([sign up here!](#)). Explore the public [Earth Engine Data Catalog](#) which includes a variety of standard Earth science raster datasets. Browse by [dataset tags](#) or by satellite ([Landsat](#), [MODIS](#), [Sentinel](#)).

In this lesson, we'll first get some hands-on practice connecting to and using Google Earth Engine to visualize global precipitation data. We'll then walk through a demonstration using GEE to visualize and analyze fire dynamics in the Arctic.

14.2 Exercise 1: An introductory lesson on using Google Earth Engine

14.2.1 Part i. Setup

1. Create a Google Earth Engine account (if you haven't already done so)
 - Please refer back to the [Preface](#) to find instructions on creating a GEE account.
2. Load libraries

```
import ee
import geemap
import pandas as pd
```

3. Authenticate your GEE account

- In order to begin using GEE, you'll need to connect your environment (`scomp`) to the authentication credentials associated with your Google account. This will need to be done each time you connect to GEE, (but only be done once per session).

```
ee Authenticate() # triggers the authentication process
```

- This should launch a browser window where you can login with your Google account to the Google Earth Engine Authenticator. Following the prompts will generate a code, which you'll then need to copy and paste into the VS Code command palette (at the top of the IDE). This will be saved as an authentication token so you won't need to go through this process again until the next time you start a new session. The browser-based authentication steps will look something like this:
 - a. **Notebook Authenticator:** choose an active Google account and Cloud Project (you may have to create one if this is your first time authenticating) and click “Generate Token”
 - b. **Choose an account:** if prompted, select the same Google account as above
 - c. **Google hasn't verified this app:** You may be tempted to click the blue “Back to safety” button, but don't! Click “Continue”
 - d. **Select what Earth Engine Notebook Client can access:** click both check boxes, then “Continue”
 - e. **Copy your authorization code** to your clipboard to paste into the VS Code command palette
 4. Lastly, initialize. This verifies that valid credentials have been created and populates the Python client library with methods that the backend server supports.

```
ee Initialize()
```

If successful, you're now ready to begin working with Earth Engine data!

14.2.2 Part ii. Explore the ERA5 Daily Aggregates Data

We'll be using the ERA5 daily aggregates reanalysis dataset, produced by the [European Centre for Medium-Range Weather Forecasts](#) (ECMWF), found [here](#), which models atmospheric weather observations.

[ERA5 Daily Aggregates](#) dataset, available via the Earth Engine Data Catalog

Take a few moments to explore the metadata record for this dataset. You'll notice that it includes a bunch of important information, including:

- **Dataset Availability:** the date range
- **Dataset Provider:** where the data come from
- **Earth Engine Snippet:** a code snippet used for loading the dataset
- **Description (tab):** get to know a bit about the data
- **Bands (tab):** the variables present in the dataset; each band has its own name, data type, scale, mask and projection
- **Image Properties:** metadata available for each image band
- **Example Code:** a script to load and visualize ERA5 climate reanalysis parameters in Google Earth Engine (JavaScript)

Reanalysis combines observation data with model data to provide the most complete picture of past weather and climate. To read more about reanalyses, check out the [EWCMWF website](#).

14.2.3 Part iii. Visualize global precipitation using ERA5 Daily Aggregate data

Content for this section was adapted from Dr. Sam Stevenson's [Visualizing global precipitation using Google Earth Engine lesson](#), given in her [EDS 220 course](#) in Fall 2021.

1. Create an interactive basemap

- The default basemap is (you guessed it) Google Maps. The following code displays an empty Google Map that you can manipulate just like you would in the typical Google Maps interface. Do this using the `Map` method from the `geemap` library. We'll also center the map at a specified latitude and longitude (here, 50N, 151E), set a zoom level, and save our map as an object called `myMap`.

```
myMap = geemap.Map(center = [60, -151], zoom = 4)  
myMap
```

2. Load the ERA5 Image Collection from GEE

- Next, we need to tell GEE what data we want to layer on top of our basemap. The `ImageCollection` method extracts a set of individual images that satisfies some criterion that you pass to GEE through the `ee` package. This is stored as an `ImageCollection` object which can be filtered and processed in various ways. We can pass the `ImageCollection` method arguments to tell GEE which data we want to retrieve. Below, we retrieve all daily ERA5 data.

💡 Earth Engine Snippets make importing `ImageCollection`s easy!

To import an `ImageCollection`, copy and paste the Earth Engine Snippet for your dataset of interest. For example, the Earth Engine Snippet to import the ERA5 daily aggregates data can be found on the [dataset page](#).

```
weatherData = ee.ImageCollection('ECMWF/ERA5/DAILY')
```

3. Select an image to plot

- To plot a map over our Google Maps basemap, we need an *Image* rather than an *ImageCollection*. ERA5 contains many different climate variables – explore which variables the dataset contains under the `Bands` tab. We'll use the

`select` method to choose the parameter(s) we're interested in from our `weatherData` object. Let's select the `total_precipitation` band.

```
# select desired bands (total_precipitation)
precip = weatherData.select("total_precipitation")
```

- We can look at our `precip` object metadata using the `print` method to see that we've isolated the `total_precipitation` band, but it's still an `ImageCollection`.

```
print(precip)
```

i Note

You may see a message in that says, “**Output exceeds the size limit. Open the full output data in a text editor**” when printing your image object metadata. Click here to see the entire output, which includes date range information.

- Let's say that we want to look at data for a particular time of interest – e.g. January 1, 2019 - December 31, 2019. We can apply the `filterDate` method to our selected `total_precipitation` parameter to filter for data from our chosen date range. We can also apply the `mean` method, which takes whatever precedes it and calculates the average – this step reduces our `ImageCollection` to a single `Image`.

```
# initial date of interest (inclusive)
i_date = '2019-01-01'

# final date of interest (exclusive)
f_date = '2020-01-01'

# select desired bands (total_precipitation), dates, and calculate total precipitation across
```

Plotting an `Image` will produce a static visualization (e.g. total precipitation at a particular point in time, or the average precipitation over a specified date range), while an `ImageCollection` can be visualized as either an animation or as a series of thumbnails (aka a “filmstrip”), such as this animation showing a three-day progression of Atlantic hurricanes in September, 2017 (source: [Google Earth Engine](#)).

```
precip = weatherData.select("total_precipitation").filterDate(i_date, f_date).sum()
```

- Use the `print` method again to check out your new `precip` object – notice that it's now an `ee.Image` (rather than `ee.ImageCollection`) and the start and end date values over which the average is taken are as we specified.

```
print(precip)
```

4. Add the precipitation Image to the basemap

- First, set a color palette to use when plotting the data layer. The following is a palette specified for ERA5 precipitation data (scroll down to the example code, available on the [landing page for the ERA5 metadata](#) in the Earth Engine Data Catalog). Here, we adjusted the max value to change the range of pixel values to which the palette should be applied – this will make our colors stand out a bit more when we layer our precipitation data on our basemap, below.

```
precip_palette = {  
    'min':0,  
    'max':5,  
    'palette': ['#FFFFFF', '#00FFFF', '#0080FF', '#DA00FF', '#FFA400', '#FF0000']  
}
```

Learn more about GEE color palettes and Image visualization [here](#).

i Note

GEE has lots of pre-defined color palettes to choose from based on the type of data you want to visualize. Also check out [Cramer et al. 2020](#) for recommended best practices when choosing color gradients that accurately represent data and are readable by those with color vision deficiencies.

- Finally, plot our filtered data, `precip`, on top of our basemap using the `addLayer` method. We'll also pass it our visualization parameters (colors and ranges stored

in `precip_palette`, the name of the data field, `total precipitation`, and opacity (so that we can see the basemap underneath).

```
myMap.addLayer(precip, precip_palette, 'total precipitation', opacity = 0.7)  
myMap
```

14.2.4 Part iv. Extract annual sum of precipitation data from points

We probably don't want to just plot the data, extracting it at points of interest to use in analysis would be helpful too. I'll show a simple example here of how to extract raster values from a set of point geometries.

First, we'll create a `pandas` data frame with some locations in Alaska.

```
# initialize data frame of points  
data = [['anc', 61.2, -150.1], ['fai', 64.8, -147.6], ['jun', 58.3, -134.4]]  
  
# Create the pandas DataFrame  
df = pd.DataFrame(data, columns=['loc', 'lat', 'lon'])
```

Next we need to convert the rows in the data frame to a GEE `FeatureCollection`. To do this, we will run a loop over every row in the data frame which:

1. creates a point geometry from the lat/lon coordinates
2. creates an attribute dictionary from the row
3. uses the attribute dictionary and geom in the `Feature` function

```
# convert data frame to list of GEE features  
features=[]  
  
for index, row in df.iterrows():  
    p_geom = ee.Geometry.Point([row['lon'], row['lat']])  
    # construct attributes for each row  
    p_props = dict(row)
```

```
p_feature = ee.Feature(p_geom, p_props)
features.append(p_feature)
```

We then pass that features list into the `FeatureCollection` function

```
ee_fc = ee.FeatureCollection(features)
```

Now, we can use the `sampleRegions` method on our `precip` image, passing the feature collection we just created to the `collection` argument.

```
res = precip.sampleRegions(collection = ee_fc)
```

This returns another feature collection, which we can convert back to a `pandas` data frame with just one function call.

```
res_df = geemap.ee_to_pandas(res)
```

Of note, by default GEE data is in maps mercator (EPSG:3857). [The documentation](#) has many more details on how to handle projections, but in this case it is okay for us to use the default value since our input data were also in EPSG:3857.

This was a simple example but you can do much more with these methods, including extracing timeseries data, extracting data from regions with polygon geometries (as opposed to point geometries), buffering features, and more. The GEE documentation is excellent, and examples abound.

14.2.5 Part v. Takeaways

In just about five lines of code (and mere seconds of execution time), we've applied and visualized a global precipitation model (that's pretty cool, right??). We can zoom in/out across our interactive map while Google Earth Engine recalculates and revisualizes our model in near real-time.

GEE does have some limitations. To name a few:

- it's a closed platform service, so GEE software can only be run on Google's infrastructure, as opposed to other open source options, like Pangeo
- it's only free for non-commercial use
- there are memory and storage limits

Despite these limitations, GEE's data catalog and cloud computing resources can vastly streamline and expedite analyzing and visualizing large geospatial datasets.

14.3 Exercise 2: Visualize fire dynamics in the Arctic using GEE

In this exercise, use the [Google Earth Engine datasets listing](#) to find a dataset showing burn areas, that includes a day burned layer.

Create a map of your dataset using methods similar to above. Look to the dataset landing page on GEE for palette help, or create your own!

Below is one solution. Don't peek unless you are stuck!

```

burn = ee.ImageCollection("ESA/CCI/FireCCI/5_1")
burn_2020 = burn.select('BurnDate').filterDate('2019-01-01', '2019-12-31').max()
fire_pal = {
  'min': 1,
  'max': 366,
  'palette': [
    'ff0000', 'fd4100', 'fb8200', 'f9c400', 'f2ff00', 'b6ff05',
    '7aff0a', '3eff0f', '02ff15', '00ff55', '00ff99', '00ffd9',
    '00ddff', '0098ff', '0052ff', '0210ff', '3a0dfb', '7209f6',
    'a905f1', 'e102ed', 'ff00cc', 'ff0089', 'ff0047', 'ff0004'
  ]
}

myMap2 = geemap.Map(center = [66,-145], zoom = 7)

myMap2.addLayer(burn_2020, fire_pal, "Burned area with date burned", opacity = 1)

```

```

legend_dict = {
    "January 1" : '#ff0000',
    "March 1" : '#7aff0a',
    "June 1" : '#00ddff',
    "Sept 1" : '#a905f1'
}

myMap2.add_legend(legend_title="Burn Day", legend_dict=legend_dict)

```

14.4 Additional Resources

Tools:

- [GEE Code Editor](#), a web-based IDE for using GEE (JavaScript)

Data:

- [Earth Engine Data Catalog](#), the main resource for “official” GEE Datasets
- [awesome-gee-community-datasets](#), the main resource for “community” GEE Datasets

Documentation, Tutorials, & Help:

- [earthengine-api](#) installation instructions
- [Creating and managing Google Cloud projects](#)
- [Troubleshooting authentication issues](#)
- [An Intro to the Earth Engine Python API](#)
- [geemap documentation](#)
- [Qiusheng Wu’s YouTube channel for GEE & geemap Python tutorials](#)
- [GEE on StackExchange](#)

Other:

- [Climate Engine App](#), a no-code user interface to Google Earth Engine for quickly and easily visualizing various Earth observation processes and variables

15 Documenting and Publishing Data

- Become familiar with the submission process
- Understand what constitutes as “large data”
 - Know when to reach out to support team for help
- Learn how data & code can be documented and published in open data archives

15.1 Introduction

A data repository is a database infrastructure that collects, manages, and stores data. In addition to the [Arctic Data Center](#), there are many other repositories dedicated to archiving data, code, and creating rich metadata. The Knowledge Network for Biocomplexity (KNB), the Digital Archaeological Record (tDAR), Environmental Data Initiative (EDI), and Zenodo are all examples of dedicated data repositories.



15.2 Metadata

Metadata are documentation describing the content, context, and structure of data to enable future interpretation and reuse of the data. Generally, metadata describe *who* collected the data, *what* data were collected, *when* and *where* they were collected, and *why* they were collected.

For consistency, metadata are typically structured following metadata content standards such as the [Ecological Metadata Language \(EML\)](#). For example, here's an excerpt of the machine-readable version of the metadata for a [sockeye salmon dataset](#):

```
<?xml version="1.0" encoding="UTF-8"?>
<eml:eml packageId="df35d.442.6" system="knb"
  xmlns:eml="eml://ecoinformatics.org/eml-2.1.1">
  <dataset>
    <title>Improving Preseason Forecasts of Sockeye Salmon Runs through
      Salmon Smolt Monitoring in Kenai River, Alaska: 2005 - 2007</title>
    <creator id="1385594069457">
      <individualName>
        <givenName>Mark</givenName>
        <surName>Willette</surName>
      </individualName>
      <organizationName>Alaska Department of Fish and Game</organizationName>
      <positionName>Fishery Biologist</positionName>
      <address>
        <city>Soldotna</city>
        <administrativeArea>Alaska</administrativeArea>
        <country>USA</country>
      </address>
      <phone phonetype="voice">(907) 260-2911</phone>
      <electronicMailAddress>mark.willette@alaska.gov</electronicMailAddress>
    </creator>
    ...
  </dataset>
</eml:eml>
```

Alternatively, the same metadata document can be converted to HTML format and displayed in a more readable form on the [web](#):

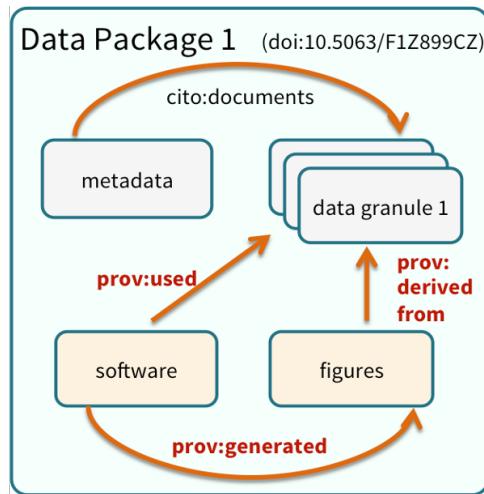
The screenshot shows a dataset page on the KNB (Knowledge Network for Biocomplexity) website. At the top, there are links for ABOUT, DATA, SHARE, TOOLS, and SIGN IN. Below the header, a navigation bar includes 'Back to search', 'Home', 'Search / Metadata', 'Copy Citation', and 'Quality report'. The main content area displays a table of files in the dataset, with columns for Name, File type, Size, Downloads, and a 'Download All' button. The table lists several files: 'Metadata: Improving Preseason Forecasts of Sockeye Salmon Runs through Salmon Smolt Monitoring in Kenai River, Alaska: 2005 - 2007.' (EML v2.1.1, 70 KB, 869 views), 'smoltSample05_07.csv' (text/csv, 21 KB, 781 downloads), 'smoltSample06_07.csv' (text/csv, 43 KB, 9 downloads), and 'smoltCount05_07.csv' (text/csv, 78 KB, 7 downloads). A link 'Show 6 more items in this data set' is visible. Below the table, a 'General' section contains an 'Identifier' field (df35d.442.6) and an 'Abstract' field containing a detailed scientific text about sockeye salmon smolt abundance estimates from 2005-2007.

As you can see from the picture above, users can download either the whole dataset or its individual components. This makes the dataset and its associated data reusable.

Additionally, the repository tracks how many times each file has been downloaded, which gives great feedback to researchers on the activity for their published data.

15.3 Data Package Structure

Note that the dataset above lists a collection of files that are contained within the dataset. We define a data package as a scientifically useful collection of data and metadata that a researcher wants to preserve. Sometimes a data package represents all of the data from a particular experiment, while at other times it might be all of the data from a grant, or on a topic, or associated with a paper. Whatever the extent, we define a data package as having one or more data files, software files, and other scientific products such as graphs and images, all tied together with a descriptive metadata document.



These data repositories all assign a unique identifier to *every version* of every data file, similarly to how it works with source code commits in GitHub. Those identifiers usually take one of two forms. A DOI (digital object identifier) identifier is often assigned to the metadata and becomes the publicly citable identifier for the package. Each of the other files gets a global identifier, often a UUID (universally unique identifier) that is globally unique. In the example above, the package can be cited with the DOI [doi:10.5063/F1F18WN4](https://doi.org/10.5063/F1F18WN4), and each of the individual files have their own identifiers as well.

15.4 Archiving Data: The Large Data Perspective

There are two components to any data package archived with the Arctic Data Center: the metadata & the data themselves. Data can be images, plain text documents, tabular data, spatial data, scripts used to analyze the data, a readme file, and more. To the best of your ability, please make sure that the data uploaded are in an open format, rather than proprietary format. We strongly recommend using open, self-documenting binary formats for large data archival. NetCDF, HDF, .mat (v7.3) and Parquet files are all examples of “self-documenting” files. In the case of a NetCDF file, users can input the attribute name, attribute description, missing value codes, units, and more into

the file itself. When these data are well-documented within themselves, it can save the time when users submit their data to us, since the documentation for variable level information is already mostly complete. We'll discuss NetCDF and metadata more in Session 8. For geospatial data, we recommend using geotiff for raster files, and geopackage files for vector files.

This section provides an overview of some highlights within the data submission process, and will specifically address issues related to datasets with large amounts of data, whether that be in number of files or cumulative file size.

First we'll go over the metadata submission; then learn how to upload the data using a secure File Transfer Protocol; and finally how to add attribute information to the data.

15.4.1 Step 1: The Narrative Metadata Submission

15.4.1.1 ORCiDs

In order to archive data with the Arctic Data Center, you must log in with your ORCID account. If you do not have one, you can create at <https://orcid.org/>. ORCID is a non-profit organization made up of research institutions, funders, publishers and other stakeholders in the research space. ORCID stands for Open Researcher and Contributor ID. The purpose of ORCID is to give researchers a unique identifier which then helps highlight and give credit to researchers for their work. If you click on someone's ORCID, their work and research contributions will show up (as long as the researcher used ORCID to publish or post their work).

Once you're logged into the Arctic Data Center with your ORCID, you can access the data submission form by clicking "Submit Data" in the navigation bar. For most dataset submissions, you would submit your data and metadata at the same using the "Add Files" buttons seen in the image below. However, when you know you have a large quantity of files or large cumulative file size, you should focus only on submitting metadata through the web form. We'll discuss how to submit large quantities of data in the next section.

The screenshot shows the Arctic Data Center's dataset submission interface. At the top, there are navigation links: Data, Support, About, Community, Submit Data, Create Portal (NEW), and a user profile for Natasha Haycock-Chavez. Below the header, it says "Natasha Haycock-Chavez. Untitled dataset. um:uuid:7e9b23a1-1ae8-4cc4-8f27-d54f9ce6d13a." A table lists one file: "Untitled dataset" (Size: 0, Type: Text, Share: off, Status: Pending). Below the table is a message: "Add files to start your dataset" with a "+ Add Files" button. On the left, a sidebar menu includes "Overview" (selected), People, Dates, Locations, Taxa, and Methods. The "Overview" tab displays fields for "Title" (with placeholder "Example: Greater Yellowstone Rivers from 1:126,700 U.S. Forest Service Visitor Maps (1981-1983)"), "Data Sensitivity" (with options: "Non-sensitive data" (selected) and "Some or all data is sensitive with minimal risk" (with a note about de-identification)), and other descriptive text.

15.4.1.2 Overview Section

In the overview section, you will include a descriptive title of your data set, select the appropriate data sensitivity tag, an abstract of the data set, keywords, funding information, and a license.

In general, if your data has been anonymized or de-identified in any way, your submission is no longer considered to have “Non-sensitive data”. If you have not had to de-identify your data or through an Institutional Review Board process, you should select the “Non-sensitive data” tag. You can find a more in-depth review of the data sensitivity tag in Chapter 12 of our [Fundamentals in Data Management](#) coursebook.



Overview Overview

People

Dates *

A title for this dataset. Include the topic, geographic location, dates, and if applicable, the scale of the data. Write out all abbreviations.

Islands of Four Mountains Artifact Analysis, Aleutian Islands, Alaska, 2013

Locations *

Taxa

Methods *

Title *

A title for this dataset. Include the topic, geographic location, dates, and if applicable, the scale of the data. Write out all abbreviations.

Islands of Four Mountains Artifact Analysis, Aleutian Islands, Alaska, 2013

Data Sensitivity *

Pick the category that best describes the level of sensitivity or restriction of the data.

Non-sensitive data

None of the data includes sensitive or protected information.

Proceed with uploading data.

Some or all data is sensitive with minimal risk

Sensitive data has been de-identified, anonymized, aggregated, or summarized to remove sensitivities and enable safe data distribution. Examples include ensuring that human subjects data, protected species data, archaeological site locations and personally identifiable information have been properly anonymized, aggregated and summarized.

Some or all data is sensitive with significant risk

The data contains human subjects data or other sensitive data. Release of the data could cause harm or violate statutes, and must remain confidential following restrictions from an Institutional Review Board (IRB) or similar body.

Abstract *

Provide a brief overview that summarizes the specific contents and purpose of this dataset.

volcanoes, the Four Mountains provided a superlative opportunity to assess the development of prehistoric human risk management of, and adaptations to, environmental instability (climatic change, sea level fluctuations) through deep time. These sites are highly unusual in terms of the paleo-data indicating human activity during human migration and societal development in the Aleutians. Extensive new radiocarbon, paleoenvironmental, and cultural data extracted from these sites yielded novel insights into North Pacific Rim regional interactions. Understanding coping mechanisms, changing subsistence, and adaptations during the prehistoric and European contact periods. The primary research goals in the Four Mountains region was understanding: (1) how human cultures and behavior have been shaped by Holocene climatic, biotic and geologic change; and (2) how human cultures have used and impacted biotic environments. Researchers tested and documented: a) relationships of and interactions among human groups who peopled the Aleutians; b) long-term change in Holocene environments and consequent change in terrestrial and marine animal populations and diversity.

You also must enter a funding award number and choose a license. The funding field will search for an NSF award identifier based on words in its title or the number itself. When including funding information not from an NSF award, please make sure to add an award number, title, and organization if possible.

The licensing options are CC-0 and CC-BY, both of which allow your data to be downloaded and re-used by other researchers.

- CC-0 Public Domain Dedication: "...can copy, modify, distribute and perform the work, even for commercial purposes, all without asking permission."
- CC-BY: Attribution 4.0 International License: "...free to...copy,...redistribute,...remix, transform, and build upon the material for any purpose, even commercially,...[but] must give appropriate credit, provide a link to the license, and indicate if changes were made."

The screenshot shows the 'Overview' tab of a dataset submission form. The 'Funding' section displays a search result for 'NSF Award 1301929 (Collaborative Research: Geological Hazards, Climate Change, and Human/Ecosystems Resilience in the Islands of the Four Mountains)'. The 'People' section is populated with 'Natasha Haycock-Chavez'. The 'Locations' section is empty. The 'Publication Date' field contains 'YYYY-MM-DD or YYYY'. The 'Methods' section is empty. The 'Usage Rights' section has 'Creative Commons Public Domain' selected. The 'Alternate Identifiers' section is empty. At the bottom right is a green 'Save dataset' button.

15.4.1.3 People Information

Information about the people associated with the dataset is essential to provide credit through citation and to help people understand who made contributions to the product. Enter information for the following people:

- Creators - **all the people who should be in the citation for the dataset**
- Contacts - one is required, but defaults to the dataset submitter if omitted
- Principal Investigators
- Any others that are relevant

For each, please provide their [ORCID](#) identifier, which helps link this dataset to their other scholarly works.

The screenshot shows a web-based form for submitting data. On the left, there's a sidebar with categories: Overview, People (which is selected), Dates, Locations, Taxa, and Methods. The main area is titled 'Co-Principal Investigators'. It contains two sets of input fields for people and organizations. The first set for a person includes fields for First name (Virginia), Last name (Hatfield), Email (virginia@hatfielddarh.com), Street address (314 Salmon Way), and Address line 2. The second set for an organization includes fields for Position name (Executive Director), Website (www.aleutians.org), and Address line 2. Below these are additional fields for a person: First name, Last name, Email, Street address, Position name, Website, Organization name, Phone, Fax, City, State or province, Postal code, and Country. A green button labeled 'ID' with the URL https://orcid.org/0000-0003-1426-8369 is present. At the bottom, there's a dropdown menu 'Choose new person or organization role ...', another set of person-related fields, and a 'Save dataset' button.

15.4.1.4 Temporal Information

Add the temporal coverage of the data, which represents the time period when data was collected. If your sampling over time was discontinuous and you require multiple date ranges to represent your data, please email us the date ranges and we will add that to the submission on the back-end.

The screenshot shows the Arctic Data Center's dataset editor. On the left is a sidebar with categories: Overview, People, Dates (which is selected and highlighted in dark blue), Locations, Taxa, and Methods. The main content area is titled 'Dates' and contains instructions: 'Specify the exact date and time of the collection or creation of this data set. Specify a year only (YYYY) or a year, month, and date (YYYY-MM-DD) and optionally, an exact time (HH:MM:SS).'. It has two sections: 'Begin' and 'End'. Under 'Begin', there are fields for 'Date (YYYY-MM-DD or YYYY)' containing '2014-07-01' and 'Time (HH:MM:SS)'. Under 'End', there are fields for 'Date (YYYY-MM-DD or YYYY)' containing '2015-08-31' and 'Time (HH:MM:SS)'. At the bottom right of the content area is a green 'Save dataset' button.

15.4.1.5 Location Information

The geospatial location that the data were collected is critical for discovery and interpretation of the data. Coordinates are entered in decimal degrees, and be sure to use **negative values** for West longitudes. The editor allows you to enter multiple locations, which you should do if you had noncontiguous sampling locations. This is particularly important if your sites are separated by large distances, so that spatial search will be more precise.

Fix the errors flagged below before submitting:

- Provide a title.
- Provide the date(s) for this data set.
- At least one method step is required.
- Pick the category that best describes the level of sensitivity or restriction of the data.
- At least one location is required.

Need help? Email us at support@arcticdata.io

The screenshot shows a form with several required fields highlighted in red. At the top, a pink banner lists five errors: 'Provide a title.', 'Provide the date(s) for this data set.', 'At least one method step is required.', 'Pick the category that best describes the level of sensitivity or restriction of the data.', and 'At least one location is required.' Below the banner, there are tabs for 'Dates', 'Locations' (which is selected and highlighted in red), 'Taxa', and 'Methods'. The 'Locations' tab has a sub-section for 'Description' and 'Bounding Box Coordinates'. The 'Description' section contains a text input field and a note: 'Provide a short and comprehensive geographic description of the sampling site(s) or location(s) where these data were collected. For example, "Teshikpuk Lake, Alaska".' The 'Bounding Box Coordinates' section contains two sets of input fields: 'Northwest coordinates' (Lat: 52.8756, Long: -169.794) and 'Southeast coordinates' (Lat: [empty], Long: [empty]).

Note that, if you miss fields that are required, they will be highlighted in red to draw your attention. In this case, for the description, provide a comma-separated place name, ordered from the local to global:

- Mission Canyon, Santa Barbara, California, USA

The screenshot shows the same dataset submission form as the previous one, but now all fields are filled correctly. The 'Description' field contains 'Aleutian Islands, Islands of the Four Mountains, Carlisle, Herbert, and Chugnatak Islands'. The 'Northwest coordinates' are now filled with '52.8756' and 'Long: -169.794'. The 'Save dataset' button at the bottom is visible.

15.4.1.6 Methods

Methods are critical to the accurate interpretation and reuse of your data. The editor allows you to add multiple different methods sections, so that you can include details of sampling methods, experimental design, quality assurance procedures, and/or computational techniques and software. Please be complete with your methods sections, as they are fundamentally important to reuse of the data. Ideally, enough detail should be provided such that a reasonable scientist could interpret the study and data for reuse without needing to consult the researchers or any other resources.

Included in the methods section is a question that asks users about the ethical research practices that may or may not have been considered throughout the research process. You will learn more about this in Chapter 14, and can find more in-depth information on our website's [data ethics](#) page.

The ethical research practices response box must be filled out in order to save your dataset. If users feel as though this question is not applicable to their research, we encourage them to discuss why that is rather than simply stating "Not Applicable," or some variation thereof. For example, say a researcher has compiled satellite imagery of weather patterns over the open ocean. Rather than respond with "N/A", a user should instead include something to the effect of "Dataset contains satellite imagery of weather patterns over the Atlantic Ocean. Imagery was downloaded from NASA and NOAA, and does not contain any individual's identifiable information." When in doubt, you can always email the support team at support@arcticdata.io.

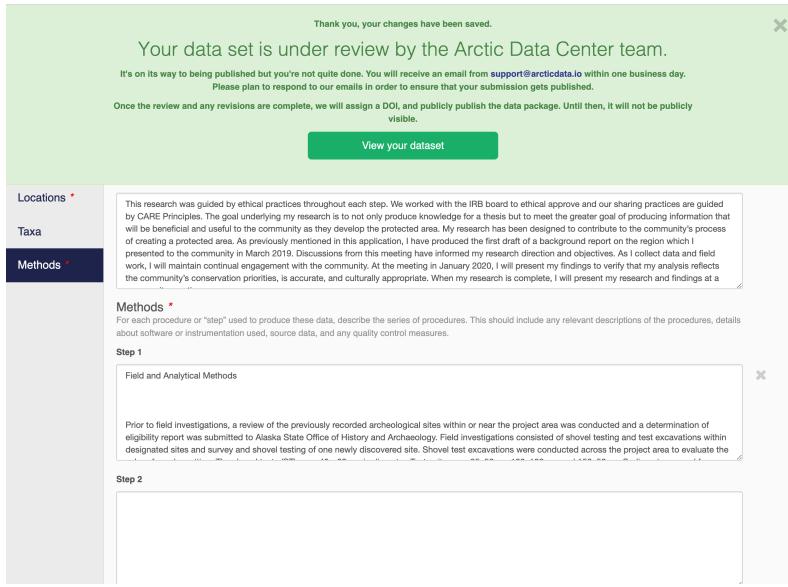
The screenshot shows a web-based dataset submission form for the Arctic Data Center. The top navigation bar includes links for Data, Support, About, Community, Submit Data, Create Portal (NEW), and a user profile for Natasha Haycock-Chavez. On the left, a sidebar lists categories: Overview, People, Dates, Locations, Taxa, and Methods (which is currently selected). The main content area is titled "Methods & Sampling". A sub-section titled "Ethical Research Practices" contains a detailed text box describing research practices, mentioning CARE Principles and community involvement. Below this is a "Methods" section with a "Step 1" sub-section containing a text box about artifact processing. At the bottom right of the main content area is a green "Save dataset" button.

15.4.1.7 Save Metadata Submission

When you're finished editing the narrative metadata, click the *Save Dataset* button at the bottom right of your screen.

If there are errors or missing fields, they will be highlighted with a red banner as seen earlier. Correct those, and then try submitting again. If the save button disappears after making corrections, add a space in the abstract and the save button should reappear. If not, please reach out to the support team for assistance.

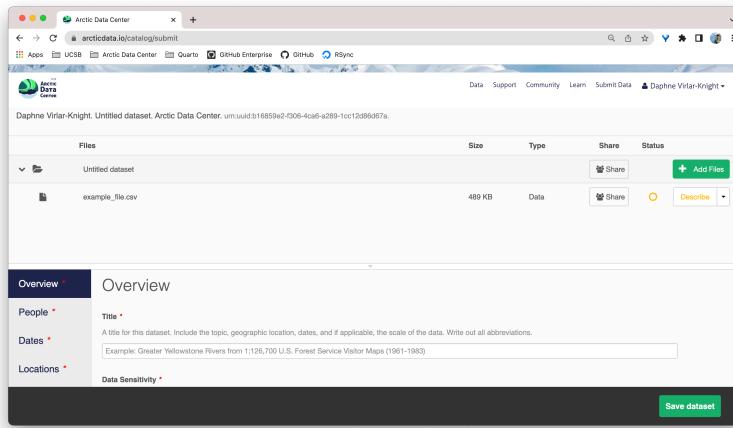
When you are successful, you should see a large green banner with a link to the current dataset view. Click the X to close that banner if you want to continue editing metadata.



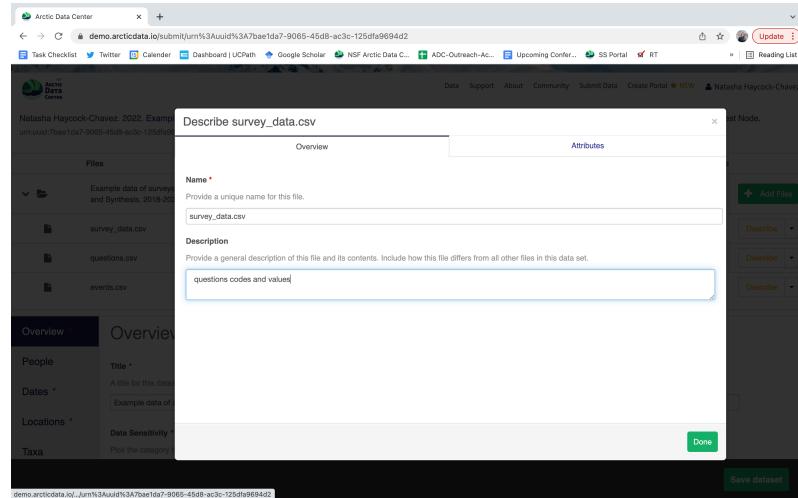
15.4.2 Step 2: Adding File & Variable Level Metadata

The final major section of metadata concerns the structure and content of your data files. Assuming there are many files (and not a few very large ones), it would be unreasonable for users to input file and variable level metadata for each file. When this situation occurs, we encourage users to fill out as much information as possible for each *unique* type of file. Once that is completed, usually with some assistance from the Data Team, we will then programmatically carry over the information to other relevant files.

When your data are associated with your metadata submission, they will appear in the data section at the top of the page when you go to edit your dataset. Choose which file you would like to begin editing by selecting the “Describe” button to the right of the file name.

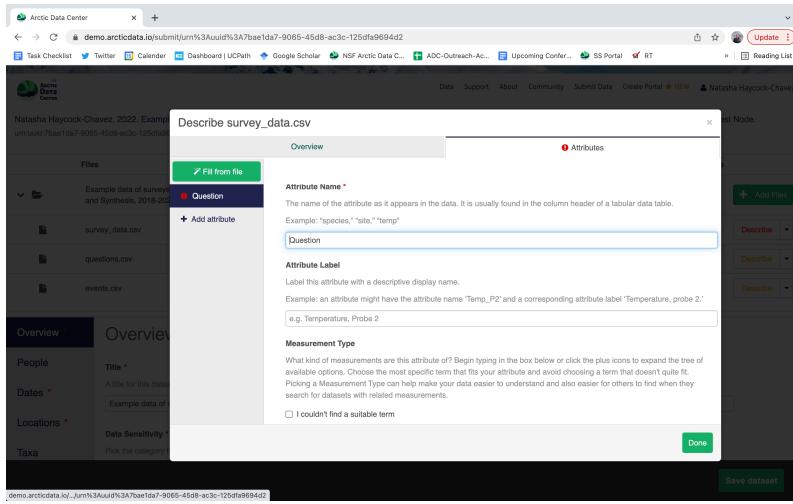


Once there, you will see the following screen. In the Overview section, we recommend ***not*** editing the file name, and instead add a descriptive overview of the file. Once done, click the Attributes tab.

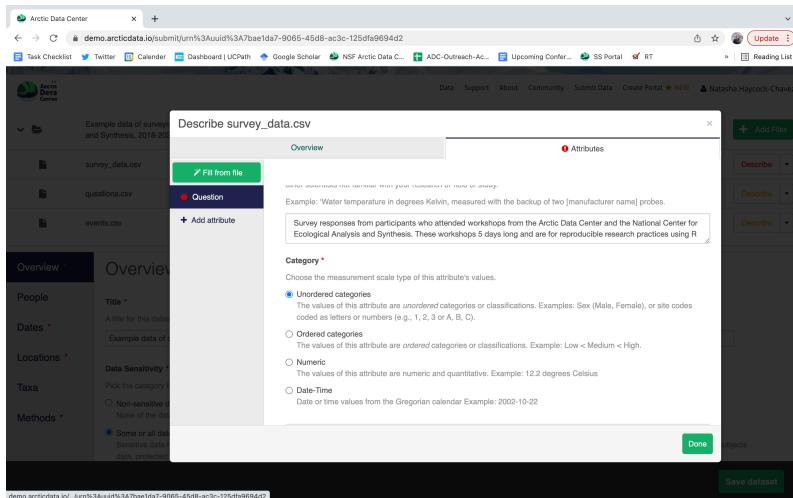


The **Attributes** tab is where you enter variable (aka attribute) information, including:

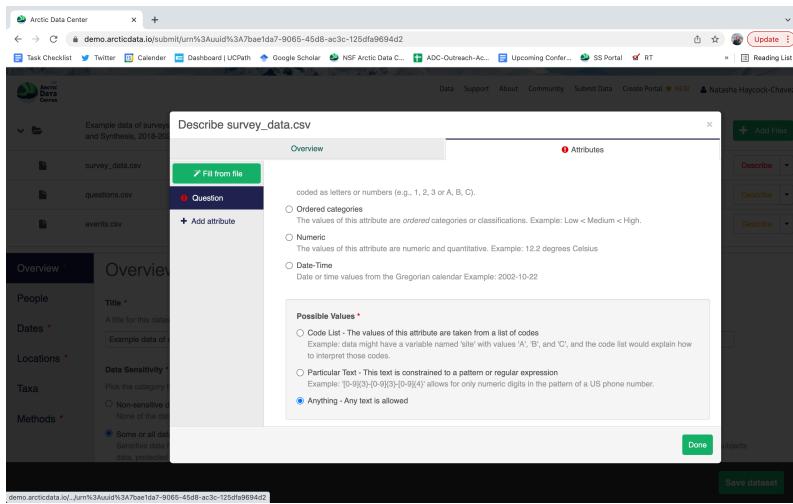
- attribute name (for programs)
- attribute label (for display)



- variable definition (be specific)
- type of measurement



- units & code definitions



Users will need to add these definitions for every variable (column) in the file. When done, click **Done**. Now, the list of data files will show a green checkbox indicating that you have fully described that file's internal structure. Proceed with documenting other unique file types, and then click **Save Dataset** to save all of these changes.

Tip

When dealing with large datasets, the Data Team can help users fill out metadata for similar file types with identical internal structures. An example of this will be discussed in the Best Practices section.

The screenshot shows a web browser window for the Arctic Data Center. The URL is <https://demo.arcticedata.io/submit/urn%3Auuid%3A7bae1da7-9065-45db-ac3c-125dfa9694d2>. The page title is "Arctic Data Center". The top navigation bar includes links for Task Checklist, Twitter, Calendar, Dashboard | UCPath, Google Scholar, NSF Arctic Data C..., ADC-Outreach-A..., Upcoming Confer..., SS Portal, RT, and Reading List. The user is logged in as Natasha Haycock-Chavez.

The main content area displays a dataset overview. The "Files" section lists four files: "Example data of surveys of remote and in-person workshop participants from the National Center of Ecological Analysis and Synthesis, 2018-2022." (5 KB), "survey_data.csv" (33 KB, Data), "questions.csv" (294 B, Data), and "events.csv" (224 B, Data). Each file has a "Share" button and a "Describe" dropdown menu.

The "Overview" tab is selected, showing fields for People, Dates, Locations, and Taxa. The "Title" field contains the dataset title. The "Data Sensitivity" field is set to "Public". A "Save dataset" button is at the bottom right.

After you get the big green success message, you can visit your dataset and review all of the information that you provided. If you find any errors, simply click **Edit** again to make changes.

15.4.3 Step 3: Uploading Large Data

In order to submit your large data files to the Arctic Data Center repository, we encourage users to directly upload their data to the Data Team's servers using a secure file transfer protocol (SFTP). There are a number of GUI driven and command line programs out there that all work well. For a GUI program, our team uses and recommends the free program [Cyberduck](#). We will discuss command line programs in [Session 18](#) in more detail, including `rsync` and Globus.

Before we begin, let's answer the following question: Why would a user want to upload their data through a separate process, rather than the web form when they submit their metadata?

Depending on your internet connection, the number of files you have, and the cumulative size of the data, users may experience difficulty uploading their data through the submission form. These difficulties are most often significantly reduced upload speeds and submission malfunctions. As such, it is best for all parties for large quantities of data to be uploaded directly to our server through an FTP.

The second question is, under what circumstances should I consider uploading data directly to the server, versus the webform? Although server uploads are more efficient for large datasets, for smaller datasets it is much more efficient to go through the webform. So, what constitutes a large dataset and what tools should you consider? Here is a helpful diagram:

	Small File Size	Large File Size
Few Files	Web Editor	SFTP (GUI or CLI)
Many Files	SFTP (GUI or CLI)	GridFTP (Globus)

- Web Editor: The web editor is most often used by those with less than a hundred files and a small cumulative file size (0-5 GBs). Overall, this option is best for those who have less than 250 files with a small cumulative file size.
- SFTP: For users that expect to upload more than 250 files and have a medium cumulative file size (10-100 GBs), uploading data to our servers via SFTP is the recommended method. This can be done through the command line, or a program like Cyberduck. If you find yourself considering uploading a zip file through the web editor, you should instead upload your files using this method.
- GridFTP: For users that expect to upload hundreds or thousands of files, with a cumulative file size of hundreds of GB to TBs, you will likely want to make use of GridFTP through Globus. Jeanette will be talking about data transfers in more depth on Thursday.

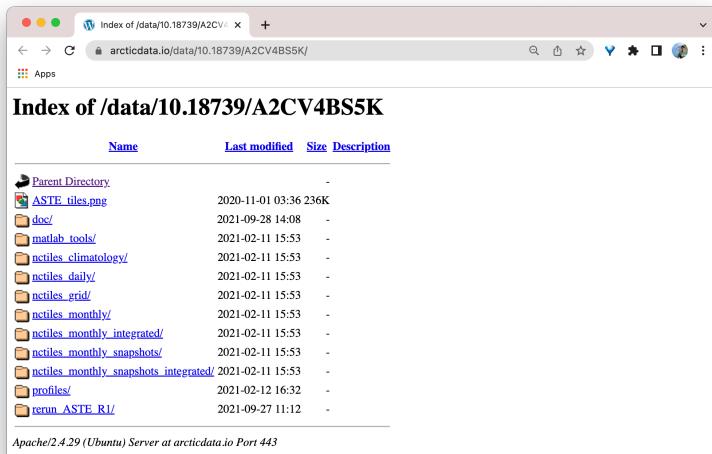
Before you can upload your data to the Data Team's server, make sure to email us at support@arcticdata.io to retrieve the login password. Once you have that, you can proceed through the following steps.

💡 Tip

If you know that you will need to use this process for more than one dataset, we suggest creating folders with the **same name** as the associated dataset's title. This way, it will be clear as to which submission the data should be associated with.

Once you have finished uploading your data to our servers, please let the Data Team know via email so that we can continue associate your uploaded data with your metadata submission.

As mentioned in [Step 1: The Narrative Metadata Submission](#) section above, when the data package is finalized and made public, there will be a sentence in the abstract that directs users to a separate page where your data will live. The following image is an example of where the data from [this](#) dataset live.



15.5 Best Practices for Submitting Large Data

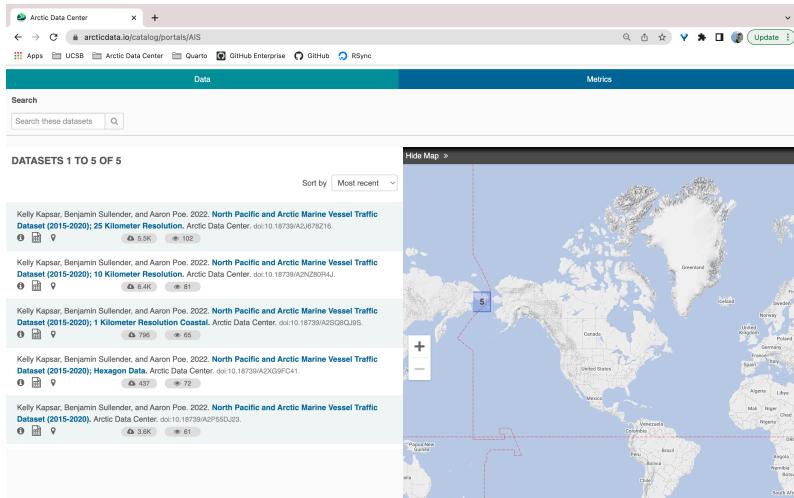
15.5.1 Data Organization

The data archival process is much smoother when users already have a well structured and organized file management system. We strongly recommend starting to work with the Arctic Data Center towards the beginning of your project, instead of the very end. Our curation team can help you publish your data faster, and more smoothly, the earlier you talk to us. Regardless of where you are in your project, however, considering your organization strategy is best done *before* data upload.

Take a look at this recently made [data portal](#). Originally, this user had only submitted one dataset with all of her files (over 700, at varying resolutions). The data contain rasters and shapefiles depicting vessel paths in the North Pacific at many different resolutions. Upon review, the submitter and curation team together determined that it would be best if other users could download the data at specific resolution scales. As a solution, we created separate datasets according to resolution and data type. With the addition of the data portal, the user was able to house their related datasets at a single URL, which allows for search and discovery within only the datasets for that project.

Note

To find out more about portals, visit <https://arcticdata.io/data-portals/>. You can view our catalog of portals [here](#). We recommend reviewing the Distributed Biological Observatory's [portal](#) to see what a more fleshed out portal can look like.



Tip

This example also illustrates that users can create multiple datasets for a single NSF award. It's understandable that you may think you have to submit all your data for a given award under one dataset. However, we think breaking up data into separate datasets (in logical, structured ways) increases the discoverability and usability of the data. This is especially true when users also create a data portal for their project.

When uploading your data files using SFTP, either through the command line or a GUI, we recommend creating folders with the same name as your dataset's title so that it is easier to associate the data with their associated metadata submission.

For more complex folder structures, it is wise to include a README file that explicitly walks users through the structure and provides a basic understanding of what is available. Generally speaking, we recommend structuring your data in an easy to understand way such that a README isn't completely necessary.

15.6 Data transfer tools

Now that we've talked about what types of large datasets you might have that need to get published on the Arctic Data Center, let's discuss how to actually get the data there. If you have even on the order of only 50GB, or more than 500 files, it will likely be more expedient for you to transfer your files via a command line tool than uploading them via our webform. So you know that you need to move a lot of data, how are you going to do it? More importantly, how can you do it in an efficient way?

There are three key elements to data transfer efficiency:

- endpoints
- network
- transfer tool

Endpoints

The from and to locations of the transfer, an endpoint is a remote computing device that can communicate back and forth with the network to which it is connected. The speed with which an endpoint can communicate with the network varies depending on how it is configured. Performance depends on the CPU, RAM, OS, and disk configuration. One key factor that affects data transfer speed is how quickly that machine can write data to disk. Slow write speeds will throttle a data transfer on even the fastest internet connection with the most streamlined transfer tool. Examples of endpoints could be:

- NCEAS included-crab server
- Your standard laptop
- A cloud service like AWS

Network speed

Network speed determines how quickly information can be sent between endpoints. It is largely, but not entirely, dependent

on what you pay for. Importantly, not all networks are created equal, even if they nominally have the same speed capability. Wired networks get significantly more speed than wireless. Networks with lots of “stuff” along the pipe (like switches or firewalls) can perform worse than those that don’t. Even the length and type of network cabling used can matter.

Transfer tools

Poll: what data transfer tools do you use regularly?

Finally, the tool or software that you use to transfer data can also significantly affect your transfer speed. There are a lot of tools out there that can move data around, both GUI driven and command line. We’ll discuss a few here, and their pros and cons.

15.6.0.0.1 * scp

`scp` or secure copy uses `ssh` for authentication and transfer, and it is included with both unix and linux. It requires no setup (unless you are on a Windows machine and need to install), and if you can `ssh` to a server, you can probably use `scp` to move files without any other setup. `scp` copies all files linearly and simply. If a transfer fails in the middle, it is difficult to know exactly what files didn’t make it, so you might have to start the whole thing over and re-transfer all the files. This, obviously, would not be ideal for large data transfers. For a file or two, `scp` is a fine tool to use.

15.6.0.0.2 * rsync

`rsync` is similar to `scp`, but syncs files/directories as opposed to copying. This means that `rsync` checks the destination to see if that file (with the same size and modified date) already exists. If it does, `rsync` will skip the file. This means that if an `rsync` transfer fails, it can be restarted again and will pick up where it left off, essentially. Neat!

15.6.0.3 * Globus

Globus is a software that uses multiple network sockets simultaneously on endpoints, such that data transfers can run in parallel. As you can imagine, that parallelization can dramatically speed up data transfers. Globus, like `rsync` can also fail gracefully, and even restart itself. Globus does require that each endpoint be configured as a Globus node, which is more setup than is required of either `scp` or `rsync`. Many institutions computing resources may have endpoints already configured as Globus endpoints, so it is always worth checking in with any existing resources that might already be set up before setting up your own. Although Globus is a free software, there are [paid options](#) which provide support for configuring your local workstation as a Globus node. Globus is a fantastic tool, but remember the other two factors controlling data transfer, it can only help so much in overcoming slow network or write speeds.

15.6.0.1 AWS sync

Amazon Web Services (AWS) has a Command Line Interface (CLI) that includes a `sync` utility. This works much like `rsync` does in that it only copies new or updated files to the destination. The difference, of course, is that AWS `sync` is specifically built to work with interacting with the AWS cloud, and is compatible with S3 buckets.

15.6.0.2 nc

`nc` (or netcat) is a low level file transfer utility that is extremely efficient when moving files around on nodes in a cluster. It is not the easiest of these tools to use, however, in certain situations it might be the best option because it has the least overhead, and therefore can run extremely efficiently.

15.7 Summary

In this lesson we learned about metadata and the structure of a data package; how to submit narrative metadata; how to upload data and when to use different services; how to document file and attribute-level metadata; and best practices for data organization and management. As mentioned earlier, more in-depth information on uploading data to our servers using Globus and RSync can be found in [Session 18](#).

16 Group Project: Visualization

In your fork of the [scalable-computing-examples](#) repository, open the Jupyter notebook in the `group-project` directory called `session-17.ipynb`. This workbook will serve as a skeleton for you to work in. It will load in all the libraries you need, including a few helper functions we wrote for the course, show an example for how to use the method on one file, and then lays out blocks for you and your group to fill in with code that will run that method in parallel.

In your small groups, work together to write the solution, but everyone should aim to have a working solution on their own fork of the repository. In other words, everyone should type out the solution themselves as part of the group effort. Writing the code out yourself (even if others are contributing to the content) is a great way to get “mileage” as you develop these skills.

17 Software Design II: Modules, Packages, and Testing

- Python modules and packages
- Managing packages with `poetry`
- Testing: why, how, and when

17.1 Python modules and packages

The [Python Package Index](#) currently houses over 400,000 software projects that you can install, reuse, and extend to accelerate your work. If you've worked in python for long at all, you have certainly used the `import` statement to load **modules** that you need in your code.

```
import numpy as np
np.pi*np.power(3, 2)
```

28.274333882308138

This lesson is a first primer on building your own modules and packages in python, and provides some core tools to get started building your own packages. For a deeper dive, check out the [Python Packages](#) book, which is freely available online. First, a few definitions:

Module A [python module](#) is a file containing commands that create functions, variables, and other definitions that can be reused in other code that you write.

Package A [python package](#) represents a way to bundle python modules, documentation, and related resources as reusable and sharable code, both for your own use and to collaborate with others.

Building reusable code though modules and packages is key to the scientific reproducibility theme that weaves throughout this course. Creating your own modules and bundling them into packages will help you (and future you) in so many ways:

- to keep you organized, making it easier for you and others to understand and use your code
- to provide consistent approaches for documentation
- to easily import your code into your own projects in a portable way that promotes code reuse
- to formally document and manage software dependencies
- to share your software with colleagues and the community

17.2 Anatomy of modules and packages

Modules in python are loaded from the definitions stored in files of the same name with a .py file extension. Within each file are definitions of objects which can be accessed when the module is imported. For example, imagine you have a file named *hello.py* in your current directory that contains a function called **helloworld** with the following contents:

```
def helloworld():
    """
    Print a hello message
    """
    print("Hello world!")
```

With that in place, you can then import and use that function in other code you write:

```
import helloworld
helloworld()
```

Packages help organize these modules into a directory structure so that they can be predictably named and utilized. Packages can contain subpackages to keep modules organized into logical groups. Packages also provide structured metadata so that we can document the package and its dependencies so that it can be reliably installed and correctly used. The structure of a package is a simple directory tree containing the package modules and subpackages. Two common structures are supported. First, with modules in the root directory:

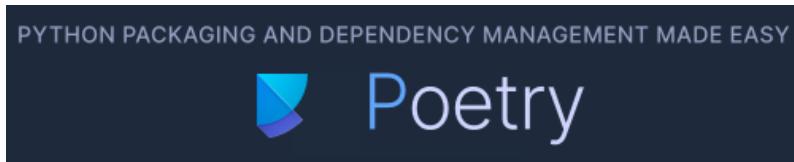
```
adcmodel
  README.md
  adcmodel
    __init__.py
    affine.py
    data.py
    reports
      __init__.py
      pdf.py
  tests
```

Alternatively, the module code is often put in it's own `src` directory, like:

```
adcmodel
  README.md
  src
    adcmodel
      __init__.py
      affine.py
      data.py
      reports
        __init__.py
        pdf.py
  tests
```

The presence of the file `__init__.py` in each module directory indicates that that directory is a module. We won't get into the details of this special file, but it is often used to load code to initialize a package.

17.3 Managing packages with poetry



While this structure and organization can be created manually, tools like [Python Poetry](#) provide a simple way to create and manage packages and their dependencies. Poetry is but one of many tools for managing packages in the python ecosystem, but we think it is useful and straightforward to use, and it tames some of the complexity of the python packaging ecosystem.

In the rest of this tutorial we will:

- Create a new virtual environment for our package
- Set up a new VS Code workspace to develop the package
- Get Poetry installed and create a skeleton package
- Add some code to the package and run it
- Talk through testing strategies for the package

17.3.1 Get poetry in a new virtual environment

To get started, we're using `virtualenvwrapper` to manage virtual environments in the course, so let's be sure it is installed in your environment. We're starting from the `scalable-computing-examples` workspace, so you should be in the `scomp` virtual environment in your terminal. So, make sure `virtualenvwrapper` is installed, and install it if it is not in `scomp`, using:

```
$ pip install virtualenvwrapper
```

The example package we will create is called `adcmodel`, and we will create a virtual environment of the same name to manage its dependencies. We'll also install poetry into that `adcmodel` virtual environment.

```
$ cd ~  
$ mkvirtualenv -p python3.9 adcmodel  
$ pip install poetry
```

Once that is done, we're going to return our current workspace to the `scomp` virtual environment (so that future lessons continue to work), and create a new workspace for our project work. To reset the venv, run the following from your terminal:

```
$ workon scomp  
$ cd ~/scalable-computing-examples
```

Now we have a virtual environment. Next we'll use it to create our new package.

17.4 Setup package workspace

17.4.1 Open VS Code window

Now that we have a virtual environment set up, let's create a new VS Code window, connect to 'included-crab', and get our workspace set up.

- Start by opening a new VS Code window by connecting to included-crab using `Command-Shift-P` and then select "Connect to Host..." and choose `included-crab.nceas.ucsb.edu`
- Once the window opens, open a new Terminal if one doesn't open, and switch to our new `adcmodel` venv using:

```
workon adcmodel
```

Finally, choose the Python version to use with `Command-Shift-P`, then choose `Python: Select Interpreter`, and choose the python version that lists your `adcmodel` virtual environment.

When that is complete, your VS Code window should look like:

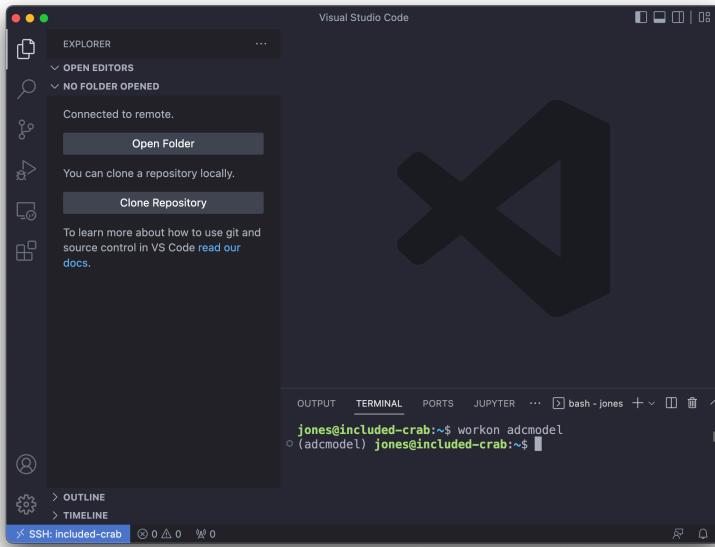


Figure 17.1: VS Code remote with a virtual environment enabled.

17.4.2 Create a new package

Now we will create a skeleton for our package with `poetry new`, and take a look at what it produces:

```
(adcmodel) jones@included-crab:~$ poetry new adcmodel
Created package adcmodel in adcmodel
(adcmodel) jones@included-crab:~$ cd adcmodel/
(adcmodel) jones@included-crab:~/adcmodel$ tree
.
  README.md
  adcmodel
    __init__.py
    pyproject.toml
    tests
      __init__.py

2 directories, 4 files
```

You'll see it created a README, an `adcmodel` directory for our source code (with the accompanying `__init__.py` to indicate this is a module), a metadata file defining our project metadata and dependencies, and a `tests` directory that we'll come back to later when we explore testing.

17.4.3 Setup the workspace

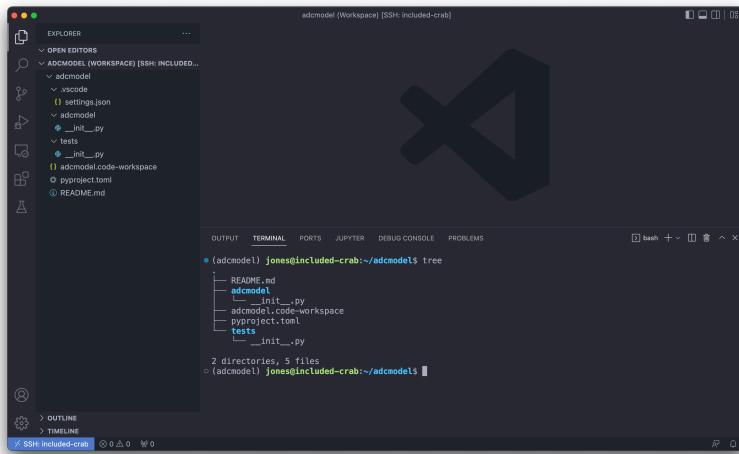
To make it easy to open up the project in VS Code, it is helpful to setup a workspace file, and configure VS Code to automatically activate our `adcmodel` virtual environment. First, let's set up the workspace by opening the project folder and saving a workspace file:

- Click “Open Folder” to open the `adcmodel` folder and display it in the VS Code file explorer
- Select “File | Save Workspace As...” to save the workspace configuration file
- Next, open the Terminal window and type `workon adcmodel` if it hasn't already been done
- Finally, save a new settings file for our virtual environment:
 - create a directory called `.vscode` using `mkdir .vscode`
 - Save a new settings file as `.vscode/settings.json` with the following content

```
{  
    "python.terminal.activateEnvInCurrentTerminal": true,  
    "python.defaultInterpreterPath": "~/.virtualenvs/adcmodel/bin/python"  
}
```

- Close the window, and open “New Window” under VS Code, then using `Ctrl-R` shortcut, open the recent `adcmodel` workspace on included-crab. If this all goes well, your session should be set to open up this `adcmodel` virtualenv each time that you open the workspace, and you should be able to see the contents of the skeleton package that you created.

Let's do some coding!



17.5 The adcmode1 example package

So, imagine this `adcmode1` package is intended to hold the code for one of our research projects in which we need to load data from a variety of sources, handle quality analysis and data cleaning, and then integrate, analyze and visualize those data. While there are many useful and general functions you might want to build into a package such as this, we're going to keep it simple and show how to build a single function that consistently loads and caches CSV-formatted data from the Arctic Data Center based on the identifier of the data object.

For this, we'll create a module called `adcmode1.data` that contains our data handling functions, one of which will be a `load_adc_csv()` function that takes an Arctic Data Center Identifier as input.

17.5.1 Load software dependencies

To create this functionality, we'll make use of existing data loading functions and data retrieval functions from the `pandas` package, and operating system utilities from the `os` package. Almost all packages and functions you write will make use of

the existing ecosystem of python packages, and we need to keep track of those dependencies in our package.

Luckily, Poetry provides simple command line utilities for adding, removing, and installing dependent packages into your virtual environment. For each package that you need, use `poetry add` to register that package in your dependencies list and install it in your `adcmodel` virtual environment. We'll do this for `pandas`, but don't need to register `urllib` or `os` because these base packages ship with python.

```
• (adcmodel) jones@included-crab:~/adcmodel$ poetry add pandas
  Using version ^1.5.0 for pandas

  Updating dependencies
  Resolving dependencies... (0.6s)

  Writing lock file

  Package operations: 4 installs, 0 updates, 0 removals

  • Installing numpy (1.23.3)
  • Installing python-dateutil (2.8.2)
  • Installing pytz (2022.2.1)
  • Installing pandas (1.5.0)
```

i The `pyproject.toml` configuration file

Dependency metadata that is added to the project and other information about the project is found in the `pyproject.toml` configuration file that is managed by poetry. In this file you'll find key metadata about the project like it's name and description and authors, as well as the list of specific dependencies that the package needs installed to work correctly. You should open this file and edit the key metadata, but `poetry add` and `poetry remove` are probably better for adding and removing dependencies.

```

[tool.poetry]
name = "adcmodel"
version = "0.1.0"
description = "Arctic Data Center example package for scalable computing course"
authors = ["Matt Jones <gitcode@magisa.org>"]
readme = "README.md"

[tool.poetry.dependencies]
python = "^3.9"
pandas = "^1.5.0"

[build-system]
requires = ["poetry-core"]
build-backend = "poetry.core.masonry.api"

```

17.5.2 Adding code to modules

Let's write some module code! In this case, we want to implement the `adcmodel.data.load_adc_csv(identifier)` function that we'll use for data loading and caching. We do this by creating a new file called `data.py` in the `adcmodel` directory, and then implementing the function. Create a new python file called `data.py` and save it to the `adcmodel` module directory, with the following function implementation:

```

import pandas as pd
import urllib.request
import hashlib
import os

def load_adc_csv(identifier=None):
    """
    Load and cache a CSV data file from the Arctic Data Center as a pandas.DataFrame
    """
    hashval = hashlib.sha256(bytes(identifier, 'utf-8')).hexdigest()

```

```
path = "data/" + hashval
if not os.path.exists("data/"):
    os.mkdir("data/")
if not os.path.exists(path):
    service = "https://arcticdata.io/metacat/d1/mn/v2/object/"
    url = service + identifier
    msg = urllib.request.urlretrieve(url, path)
df = pd.read_csv(path)
return df
```

Once you save the file, we can use `poetry install` to install the package and try running the function in an interactive session. Check it out by opening a new Jupyter notebook and running the following in the first cell (after you install the needed jupyter packages):

```
from adcmodel.data import load_adc_csv

df = load_adc_csv("urn:uuid:e248467d-e1f9-4a32-9e38-a9b4fb17cefb")
df.head()
```

:tada: Congratulations, you've created your first package!
:tada:

17.5.3 Add development packages

Utilities like python linters and code formatters can make coding faster, more consistent, and more reliable. Poetry supports the use of development packages that are available in the virtual environment during the development cycle but are not required for runtime use of the package. Like regular dependencies, development packages can be added with `poetry add`, but we also pass in a qualifier to add it to `--group dev`, which marks it as only for use during development.

```
poetry add --group dev pytest black
```

You can now use `black` to format your code (via Command-Shift-P Format Document) and `pytest` to run tests. Next, let's get into testing.

17.6 Testing

A major advantage of writing packages is that you can thoroughly test each of the functions that you write, thereby ensuring that they operate correctly both with good data and bad data. While it's always helpful to confirm that a function returns the correct value within the bounds that it is normally meant to operate, it is also important to pass it data that it wouldn't normally see. For example, does your linear model produce the correct results when it is passed a normal dataset? How about if it is passed a dataset with only one observation, or with all the same value, or with 0 or negative numbers? Does it produce an understandable error if some of the values are character values when they should be numeric? Does it continue operating in the face of bad data? Should it, or should it not produce an exception? These and others are the questions you need to ask when testing your packages.

Python provides a number of robust testing frameworks, but we will work with `pytest`, one of the common ones used for unit testing. The `poetry new` command already set up our package to be ready for testing, so now all we need to do is to add some test files in the `tests` directory.

In `pytest` and other frameworks, the common pattern is to use the python `assert` function to determine whether a value produced from your code matches expectations. For example, given a function that produces a product of two numbers, we could assert that `product(5, 6) == 30` and that `product(-3, -2) == 6`.

Generally, it is good practice to produce at least one test file for each of your package source files, and their names are prefixed with `test_`. So, for our `data.py` source file, we might expect to build tests in a file called `test_data.py`. Let's create that file in the `tests` directory and walk through populating it with some simple tests. The functions in the test file are also named with `test_` as the prefix on the function name, and generally each test function atomically tests one aspect of the code. That way, when a test function fails, it should be quickly clear what type of an error would produce that kind of failure, and how to locate it in the code.

To get started, create `tests/test_data.py` and populate it with our first test function:

```
from adcmodel.data import load_adc_csv

def test_load_adc_csv():
    df = load_adc_csv("urn:uuid:e248467d-e1f9-4a32-9e38-a9b4fb17cefb")
    assert df is not None
```

You can run the tests after saving the file by running `pytest` on the commandline, which will produce output indicating which tests passed and failed.

We can add other tests that test multiple conditions, and we can use any of the features of python to creatively test that our code will produce correct results. Here's another test to check that the data file we loaded has the right number of rows and columns:

```
def test_load_adc_csv_shape():
    df = load_adc_csv("urn:uuid:e248467d-e1f9-4a32-9e38-a9b4fb17cefb")
    assert df is not None
    assert (df.shape[0] == 17856) & (df.shape[1] == 6)
```

```
• (adcmodel) jones@includer-crab:~/adcmodel$ pytest
=====
test session starts =====
platform linux -- Python 3.9.13, pytest-7.1.3, pluggy-1.0.0
rootdir: /home/jones/adcmodel
collected 2 items

tests/test_data.py .. [100%]

===== 2 passed in 0.38s =====
○ (adcmodel) jones@includer-crab:~/adcmodel$
```

Finally, what happens when tests are run that fail? Try adding the following test, and determine if the failure is due to the code, the test, or the data? What would you change to make this test reliably pass, or fail more gracefully?

```
def test_load_adc_csv_inputs():
    df = load_adc_csv()
    assert df is not None
```

The answer to this and other more detailed questions on testing are covered in the excellent overview of testing provided by the

“Chapter 5: Testing” of “Python Packages”.

17.7 Summary and next steps

Python packages are an excellent mechanism to organize your modules for reuse both within your local code projects and with your colleagues and the broader science community. Adopting a structured approach to organizing, documenting, and testing your code can ultimately save you a lot of time and also increase the quality of your work. Scientific software is also a first-class product of the research enterprise, and is critical for fully understanding the lineage of derived data and research results. The emerging consensus is that publishing open source software packages to GitHub or other source code control systems, and then archiving them in repositories like [Software Heritage](#), the [Arctic Data Center](#), [Zenodo](#), or the [Python Package Index](#) increases the robustness, transparency, and reproducibility of research results. And it also provides a robust starting point for research teams to build upon and amplify each other’s work. Using tools like Poetry and pytest to organize and build packages makes the process quick and efficient.

17.8 Further reading

- [Python Packages](#) by Tomas Beuzen & Tiffany Timbers
- [Python Poetry](#)

18 What is Cloud Computing Anyways?

- Learn different uses of the term *cloud computing*
- Understand the services of cloud providers
- Gain familiarity with containerized computing
- Work hands-on running a docker container, and building your own and running it

18.1 What is cloud computing anyways?

The buzzword we all hear, but maybe don't quite understand.

..

Cloud computing is a lot of things... but generally speaking:

 Cloud computing is the delivery of on-demand computer resources over the Internet. Or just “using someone else’s computer”.

 Cloud computing services are typically offered using a “pay-as-you-go” pricing model, which in some scenarios may reduce capital expenses.

 “The Cloud” is powered by a global network of data centers which house the hardware (servers), power, and backup systems, etc. These data centers and infrastructure are managed by cloud providers

 Cloud computing is a technology approach to using lightweight virtualization services to share large physical computing clusters across many users

Check out this article by Ingrid Burrington in The Atlantic, [Why Amazon’s Data Centers are Hidden in Spy Country](#), for

some interesting insight into one of the largest cloud provider's data centers.

18.2 Commercial clouds

There are a lots of different cloud computing platforms, but the big ones are:

- [Amazon Web Services \(AWS\)](#)
- [Google Cloud Platform \(GCP\)](#)
- [Microsoft Azure](#)



There are *many* other cloud service providers that offer varying degrees of infrastructure, ease of use, and cost. Check out [DigitalOcean](#), [Kamatera](#), and [Vultr](#) to start.

18.3 Academic clouds

And the [Pangeo](#) project is creating an open community focused on maintaining, supporting, and deploying open infrastructure for cloud computing. They support key [scientific software packages](#) used throughout the cloud community, including `xarray` and `dask`, and generally are broadening capacity for large-scale, impactful research.

18.3.1 Cloud deployment options

Different cloud service and deployment models offer a suite of options to fit client needs

- ☰ **Service Models:** When you work in “the cloud” you’re using resources – including servers, storage, networks, applications, services, (and more!) – from a very large resource pool



Advancing
Innovation



Federal agencies in the US and other institutions also support massive computing facilities supporting cloud computing. While there are too many to fully list, programs such as the National Science Foundation's [ACCESS](#) program, the Department of Energy's [National Energy Research Scientific Computing Center \(NERSC\)](#), and the [CyVerse](#) platform provide massive cloud computing resources to academic and agency researchers. The huge advantage is that these resources are generally free-for-use for affiliated researchers. When you need access to massive CPU and GPU hours, an application for access to these facilities can be extremely effective.

that is managed by you or the cloud service provider. Three cloud service models describe to what extent your resources are managed by yourself or by your cloud service providers.

- Infrastructure as a Service (IaaS)
- Platform as a Service (PaaS)
- Software as a Service (SaaS)

☞ **Deployment Models:** Cloud deployment models describe the type of cloud environment based on ownership, scale, and access.

- Private Cloud
- Public Cloud
- Hybrid Cloud

18.4 Service Models

➤ **Infrastructure as a Service (IaaS)** provides users with computing resources like processing power, data storage capacity, and networking. IaaS platforms offer an alternative to on-premise physical infrastructure, which can be costly and labor-intensive. In comparison, IaaS platforms are more cost-effective (pay-as-you-go), flexible, and scalable.

One example of IaaS is [Amazon EC2](#), which allows users to rent virtual computers on which to run their own computer applications (e.g. R/RStudio).

➤ **Platform as a Service (PaaS)** provides developers with a framework and tools for creating unique applications and software. A benefit of SaaS is that developers don't need to worry about managing servers and underlying infrastructure (e.g. managing software updates or security patches). Rather, they can focus on the development, testing, and deploying of their application/software.

One example of SaaS is [AWS Lambda](#), a serverless, event-driven compute service that lets you run code for virtually any type of application or backend service without provisioning or managing servers.

➤ **Software as a Service (SaaS)** makes software available to users via the internet. With SaaS, users don't need to install and run software on their computers. Rather, they can access everything they need over the internet by logging into their account(s). The software/application owner does not have any control over the backend except for application-related management tasks.

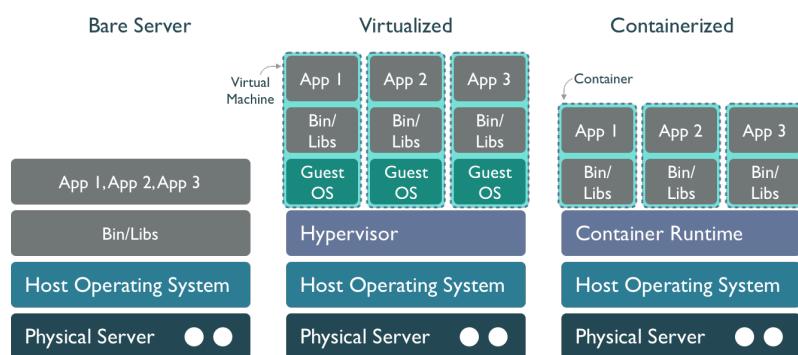
Some examples of SaaS applications include [Dropbox](#), [Slack](#), and [DocuSign](#).

18.5 🍕 An Analogy: Pizza as a Service

Image Source: David Ng, Oursky

18.6 Virtual Machines and Containers

As servers grow in size, we have increasing amounts of power and resources, but also a larger space to manage. Traditional operating systems use a common memory and process management model that is shared by all users and applications, which can cause some issues if one of the users consumes all of the memory, fills the disk, or causes a kernel panic. When running on a bare server, all of the processes from all users are mixed together and are not isolated, so the actions of one process can have large consequences for all of the others.



Virtual machines Virtualization is an approach to isolate the environments of the various users and services of a system so that we can make better use of the resource, and protect processes. In a virtualized environment, the host server still runs a host operating system, which includes a hypervisor process that can mediate between guest hosts on the machine and the underlying host operating system and hardware. This is effective at creating multiple **virtual machines** (VMs) running side-by-side on the same hardware. From the outside, and to most users, virtual machines appear to be a regular host on the network, when in fact they are virtual hosts that share the same underlying physical resources. But it also results in a fair amount of redundancy, in that each virtual machine must have its own operating system, libraries, and other resources. And calls pass through the guest operating system through the hypervisor to the physical layer, which can impose a performance penalty.

Containers A further step down the isolation road is to use a Container Runtime such as [containerd](#) or [Docker Engine](#). Like virtual machines, containers provide mechanisms to create images that can be executed by a container runtime, and which provide stronger isolation among deployments. But they are also more lightweight, as the container only contains the libraries and executables needed to execute a target application, and not an entire guest operating system. They also are built using a layered file system, which allows multiple images to be layered together to create a composite that provides rich services without as much duplication. This means that applications run with fewer resources, start up and shut down more quickly, and can be migrated easily to other hosts in a network.



19 Reproducibility and Containers

- Think about dependency management, reproducibility, and software
- Become familiar with containers as a tool to improve computational reproducibility
- Discuss how the techniques from this class can improve reproducibility
- Slides: Accelerating synthesis science through reproducible science

19.1 Software collapse

There are essentially two tasks that need to be organized and financed: preventing collapse due to changes in layers 1 and 2, and implementing new models and methods as the scientific state of the art advances. These two tasks go on in parallel and are often executed by the same people, but in principle they are separate and one could concentrate on just one or the other.

The common problem of both communities is collapse, and the common enemy is changes in the foundations that scientists and developers build on.

While most reproducibility work has focused at the project-specific software level, Hinsen argues that there is tremendous risk deep down in the software ecosystem. There are many examples of compatibility decisions in scientific software that create problems for long-term stability, even if they are the

In his blog post on [software collapse](#), Konrad Hinsen outlines the fragile state of scientific software today.

While we like to think that our software is robust, in reality it is being constantly revised, expanded, broken, and abandoned, such that our use of any software today is almost certainly doomed to dysfunction shortly. Hinsen describes the science software ecosystem as having 4 layers, from totally general operating systems at the bottom, to

scientific libraries, discipline-specific libraries, and finally project-specific software at the top. Hinsen states:

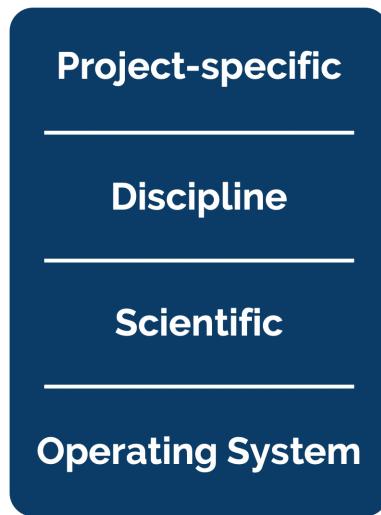


Figure 19.1: Adapted from Hinsen's software stack

right decision at the time. For example, see the compatibility discussions from [numpy](#) and [pandas](#) to get a sense of the challenges faced when maintaining science software. This issue with abandoned software and lonely maintainers has become a prevalent meme in the open source software world.

19.2 Delaying the inevitable

While software is destined to continue to rot, projects like [Wholetale](#), [repo2docker](#), and [Binder](#) are working towards tools that can help capture more of the software stack needed to reproduce computations, thereby enabling us to freeze older versions of software and replay them in compatible combinations. Or, barring full re-execution, at least to understand them. The idea is to take advantage of containerized computing environments, which can be described as ‘software as code’ because containers can be declaratively described in a configuration file. While project-specific software and possibly disciplinary software de-

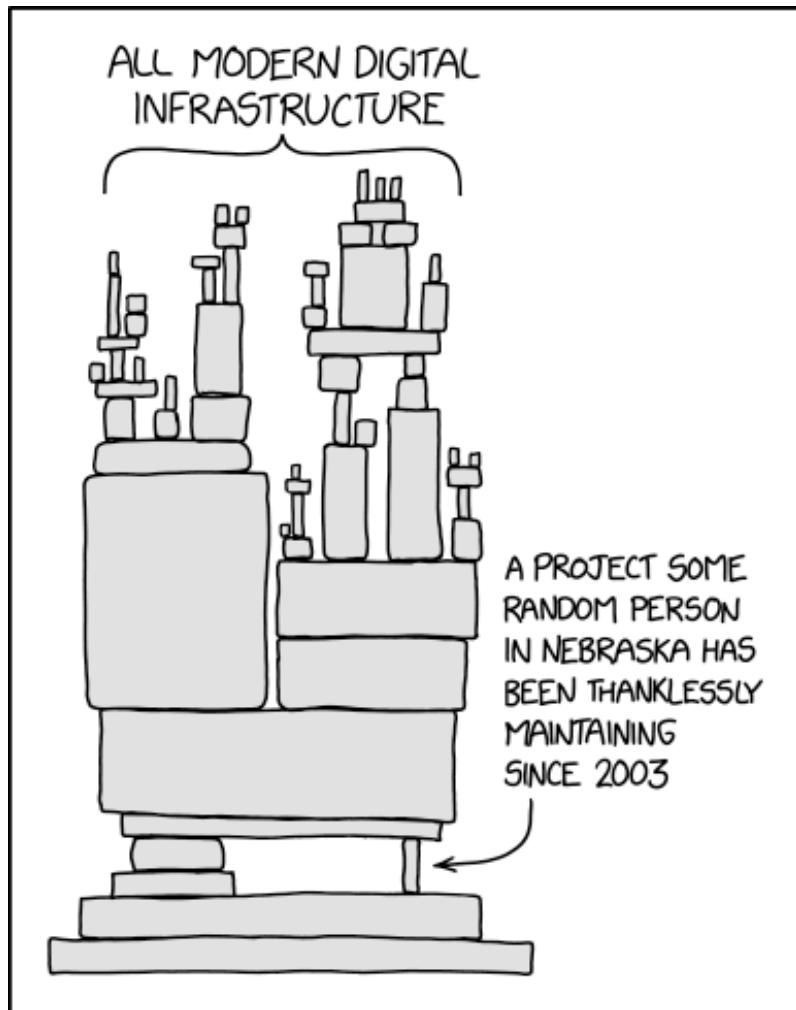
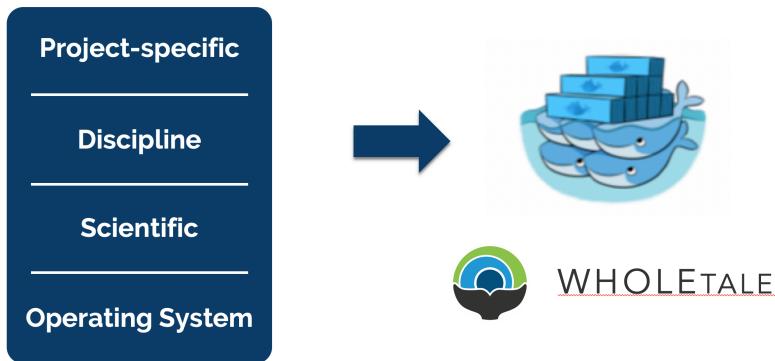


Figure 19.2: [Dependency](#)

dependencies might be captured in a `requirements.txt` file for python, we also need to know all of the details of lower layers, including system libraries and kernel versions. Plus data dependencies, processing workflows, and provenance. And it turns out that these very details can be captured quite effectively in a virtualized container. For example, WholeTale uses container images to record not just project-specific software dependencies, but also operating system details, data dependencies, and runtime provenance.



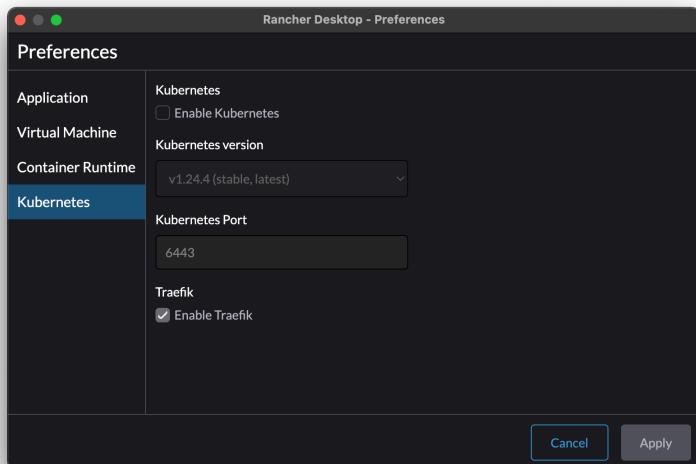
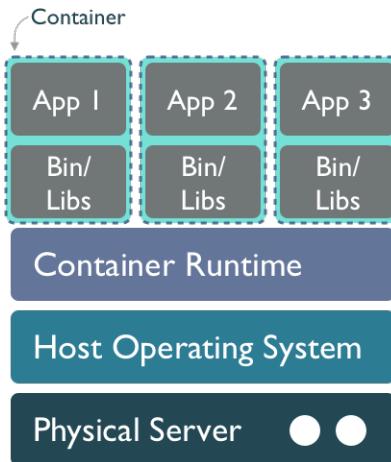
This is accomplished by first configuring a base image that contains core software from layers 1 and 2, and then researchers work within common tools like Jupyter and RStudio to create their scientific products at layers 3 and 4. These can then be captured in virtual environments, which get added to the container configuration, along with key details such as external data dependencies and process metadata. All of this can be serialized as a “Tale”, which is basically an archival package bundled around a container image definition.

19.3 Hands-off with Containers and Docker

Working with docker or containers requires a container runtime. One of the nicest lately is [Rancher Desktop](#). Install the binary for your platform, and then after it starts, enter the Configuration Preferences, and then deselect “Kubernetes” to disable the kubernetes distribution, which takes up a lot of resources.

This approach combines the evolving approach to using a **Dockerfile** to precisely configure a container environment with the use of structured metadata and provenance information to create an archival-quality research artifact. These WholeTale **tales** can be downloaded to any machine with a container runtime and executed under the same conditions that were used when the original tale was archived.

Containerized



There are many different tools you can use with docker, including the Docker engine system from Docker.com using the **docker** client tool, and the **containerd** ecosystem using **nerdctl** as a client tool. Both the **docker** command and the **nerdctl** command share the same commands. The basics we will practice will be:

- `nerdctl pull python:3.9`: to grab an existing python image from the DockerHub repository
- `nerdctl run -it python:3.9 -- python`: to start a standard python interpreter
- `nerdctl build`: to build a new image from a Dockerfile configuration file

19.4 Discussion

To wrap up this week, let's kick off a discussion with a couple of key questions.

- As we learn new tools for scalable and reproducible computing, what can we as software creators do to improve robustness and ease maintenance of our packages?
- Given the fragility of software ecosystems, are you worried about investing a lot of time in learning and building code for proprietary cloud systems? Can we compel vendors to keep their systems open?
- How much can technological solutions such as containers truly address the issues around dependencies, maintainability, and sustainability of scientific software?
- What relative fraction of the research budget should funders invest in software infrastructure, data infrastructure, and research outcomes?

20 Parquet and Arrow

- The difference between column major and row major data
- Speed advantages to columnar data storage
- How `arrow` enables faster processing

20.1 Introduction

Parallelization is great, and can greatly help you in working with large data. However, it might not help you with every processing problem. Like we talked about with Dask, sometimes your data are too large to be read into memory, or you have I/O limitations. Parquet and `pyarrow` are newer, powerful tools that are designed to help overcome some of these problems. `pyarrow` and Parquet are newer technologies, and are a bit on the ‘bleeding edge’, but there is a lot of excitement about the possibility these tools provide.

Before jumping into those tools, however, first let’s discuss system calls. These are calls that are run by the operating system within their own process. There are several that are relevant to reading and writing data: open, read, write, seek, and close. Open establishes a connection with a file for reading, writing, or both. On open, a file offset points to the beginning of the file. After reading or writing n bytes, the offset will move n bytes forward to prepare for the next operation. Read will read data from the file into a memory buffer, and write will write data from a memory buffer to a file. Seek is used to change the location of the offset pointer, for either reading or writing purposes. Finally, close closes the connection to the file.

If you’ve worked with even moderately sized datasets, you may have encountered an “out of memory” error. Memory is where a computer stores the information needed immediately for processes. This is in contrast to storage, which is typically slower

to access than memory, but has a much larger capacity. When you `open` a file, you are establishing a connection between your processor and the information in storage. On `read`, the data is read into memory that is then available to your python process, for example.

So what happens if the data you need to read in are larger than your memory? My brand new M1 MacBook Pro has 16 GB of memory, but this would be considered a modestly sized dataset by this course's standards. There are a number of solutions to this problem, which don't involve just buying a computer with more memory. In this lesson we'll discuss the difference between row major and column major file formats, and how leveraging column major formats can increase memory efficiency. We'll also learn about another python library called `pyarrow`, which has a memory format that allows for "zero copy" read.

20.2 Row major vs column major

The difference between row major and column major is in the ordering of items in the array when they are read into memory.

Take the array:

```
a11 a12 a13
```

```
a21 a22 a23
```

This array in a row-major order would be read in as:

```
a11, a12, a13, a21, a22, a23
```

You could also read it in column-major order as:

```
a11, a21, a12, a22, a13, a33
```

By default, C and SAS use row major order for arrays, and column major is used by Fortran, MATLAB, R, and Julia.

Python uses neither, instead representing arrays as lists of lists, though `numpy` uses row-major order.

20.2.1 Row major versus column major files

The same concept can be applied to file formats as the example with arrays above. In row-major file formats, the values (bytes) of each record are read sequentially.

Name	Location	Age
John	Washington	40
Mariah	Texas	21
Allison	Oregon	57

In the above row major example, data are read in the order:
John, Washington, 40 \n Mariah, Texas, 21.

This means that getting a subset of rows with all the columns would be easy; you can specify to read in only the first X rows (utilizing the seek system call). However, if we are only interested in Name and Location, we would still have to read in all of the rows before discarding the Age column.

If these data were organized in a column major format, they might look like this:

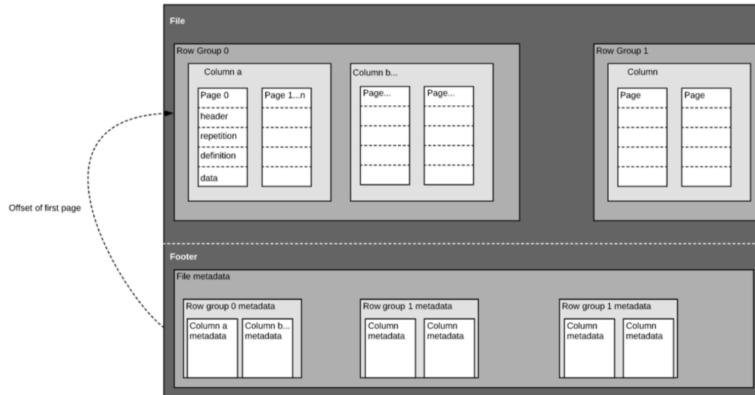
```
Name: John, Mariah, Allison  
Location: Washington, Texas, Oregon  
Age: 40, 21, 57
```

And the read order would first be the names, then the locations, then the age. This means that selecting all values from a set of columns is quite easy (all of the Names and Ages, or all Names and Locations), but reading in only the first few records from each column would require reading in the entire dataset. Another advantage to column major formats is that compression is more efficient since compression can be done across each column, where the data type is uniform, as opposed to across rows with many data types.

20.3 Parquet

Parquet is an open-source binary file format that stores data in a column-major format. The format contains several key components:

- row group
- column
- page
- footer



Row groups are blocks of data over a set number of rows that contain data from the same columns. Within each row group, data are organized in column-major format, and within each column are pages that are typically 1MB. The footer of the file contains metadata like the schema, encodings, unique values in each column, etc.

The parquet format has many tricks to increase storage efficiency, and is increasingly being used to handle large datasets.

20.4 Arrow

So far, we have discussed the difference between organizing information in row-major and column-major format, how that

applies to arrays, and how it applies to data storage on disk using Parquet.

Arrow is a language-agnostic specification that enables representation of column-major information in memory without having to serialize data from disk. The Arrow project provides implementation of this specification in a number of languages, including Python.

Let's say that you have utilized the Parquet data format for more efficient storage of your data on disk. At some point, you'll need to read that data into memory in order to do analysis on it. Arrow enables data transfer between the on disk Parquet files and in-memory Python computations, via the `pyarrow` library.

`pyarrow` is great, but relatively low level. It supports basic group by and aggregate functions, as well as table and dataset joins, but it does not support the full operations that `pandas` does.

20.5 Example

In this example, we'll read in a dataset of fish abundance in the San Francisco Estuary, which is published in csv format on the [Environmental Data Initiative](#). This dataset isn't huge, but it is big enough (3 GB) that working with it locally can be fairly taxing on memory. Motivated by user difficulties in actually working with the data, the [deltafish R](#) package was written using the R implementation of `arrow`. It works by downloading the EDI repository data, writing it to a local cache in parquet format, and using `arrow` to query it. In this example, I've put the Parquet files in a sharable location so we can explore them using `pyarrow`.

First, we'll load the modules we need.

```
import pyarrow.dataset as ds
import numpy as np
import pandas as pd
```

Next we can read in the data using `ds.dataset()`, passing it the path to the parquet directory and how the data are partitioned.

```
deltafish = ds.dataset("/home/shares/deltafish/fish",
                      format="parquet",
                      partitioning='hive')
```

You can check out a file listing using the `files` method. Another great feature of parquet files is that they allow you to partition the data across variables of the dataset. These partitions mean that, in this case, data from each species of fish is written to its own file. This allows for even faster operations down the road, since we know that users will commonly need to filter on the species variable. Even though the data are partitioned into different files, `pyarrow` knows that this is a single dataset, and you still work with it by referencing just the directory in which all of the partitioned files live.

```
deltafish.files
```

```
['/home/shares/deltafish/fish/Taxa=Acanthogobius flavimanus/part-0.parquet',
 '/home/shares/deltafish/fish/Taxa=Acipenser medirostris/part-0.parquet',
 '/home/shares/deltafish/fish/Taxa=Acipenser transmontanus/part-0.parquet',
 '/home/shares/deltafish/fish/Taxa=Acipenser/part-0.parquet'...]
```

You can view the columns of a dataset using `schema.to_string()`

```
deltafish.schema.to_string()
```

```
SampleID: string
Length: double
Count: double
Notes_catch: string
Species: string
```

If we are only interested in a few species, we can do a filter:

```
expr = ((ds.field("Taxa")=="Dorosoma petenense") |  
        (ds.field("Taxa")=="Morone saxatilis") |  
        (ds.field("Taxa")=="Spirinchus thaleichthys"))  
  
fishf = deltafish.to_table(filter = expr,  
                           columns =['SampleID', 'Length', 'Count', 'Taxa'])
```

There is another dataset included, the survey information. To do a join, we can just use the `join` method on the `arrow` dataset.

First read in the survey dataset.

```
survey = ds.dataset("/home/jclark/deltafish/survey",  
                    format="parquet",  
                    partitioning='hive')
```

Take a look at the columns again:

```
survey.schema.to_string()
```

Let's pick out only the ones we are interested in.

```
survey_s = survey.to_table(columns=['SampleID', 'Datetime', 'Station', 'Longitude', 'Latitude'])
```

Then do the join, and convert to a pandas `data.frame`.

```
fish_j = fishf.join(survey_s, "SampleID").to_pandas()
```

Note that when we did our first manipulation of this dataset, we went from working with a `FileSystemDataset`, which is a representation of a dataset on disk without reading it into memory, to a `Table`, which is read into memory. `pyarrow` has a [number of functions](#) that do computations on datasets without reading them into memory. However these are evaluated “eagerly,” as opposed to “lazily.” These are useful in some cases, like above, where we want to take a larger than memory dataset and generate a smaller dataset (via filter, or group by/summarize), but are not as useful if we need to do a join before our summarization/filter.

More functionality for lazy evaluation is on the horizon for `pyarrow` though, by leveraging [Ibis](#).

20.6 Synopsis

In this lesson we learned:

- the difference between row major and column major formats
- under what circumstances a column major data format can improve memory efficiency
- how `pyarrow` can interact with Parquet files to analyze data