# Basic Computing Skills

## MLGEO Lecture 2

Marine Denolle @ UW mdenolle@uw.edu

# Goals of class

1. Literature review Guidelines & Goals
2. Textbook and other resources
   a. Our class working textbook
   b. Borrowed materials from EarthDataScience.
3. CyberInfrastructure
4. Get to know your cohort and brainstorm on cool data sets. Check MLGEO-data set and suggest a contribution.


5. Login to the Hub: https://jupyter.rttl.uw.edu/2024-autumn-ess-469-a
6. Basic Bash: example in jupyterhub
7. Basic Jupyter Environment: getting_started_guide by the UW IT
8. Python Environments: exercise, make your own environment and test in python.

# Literature Review

1. Full citation (including doi):
2. Scientific Motivation:
3. Data Source, Type, and modality:
4. Method [*]: if they use ML, what is the baseline model? What model architecture did they use? How did they train their model?
5. Key points of the research findings:
6. [If applicable] Data: does the study use open-access data (check for public repositories)? What is the guidance on how to access the data?
7. [If applicable] Is the workflow described, and could it be reproduced? (describe data and computational workflow
8. [If applicable]  Does the manuscript provide a link (or zip file) to a code or notebook that can reproduce the work?
9. Optional: What type of open access does this journal offer?

[*] for review papers, did they, and how many papers did they use to review? Did they have any systematic meta-analysis process?

# Ethics in publishing

## CRediT – Contributor Roles Taxonomy

CRediT (Contributor Roles Taxonomy) is high-level taxonomy, including 14 roles, that can be used to represent the roles typically played by contributors to scientific scholarly output. The roles describe each contributor's specific contribution to the scholarly output.

### 14 Contributor Roles

| | |
|---|---|
| Conceptualization | Resources |
| Data curation | Software |
| Formal Analysis | Supervision |
| Funding acquisition | Validation |
| Investigation | Visualization |
| Methodology | Writing – original draft |
| Project administration | Writing – review & editing |

Marine Denolle @ UW mdenolle@uw.edu

**Conceptualization** – Ideas; formulation or evolution of overarching research goals and aims.
**Data curation** – Management activities to annotate (produce metadata), scrub data and maintain research data (including software code, where it is necessary for interpreting the data itself) for initial use and later re-use.
**Formal analysis** – Application of statistical, mathematical, computational, or other formal techniques to analyze or synthesize study data.
**Funding acquisition** - Acquisition of the financial support for the project leading to this publication.
**Investigation** – Conducting a research and investigation process, specifically performing the experiments, or data/evidence collection.
**Methodology** – Development or design of methodology; creation of models.
**Project administration** – Management and coordination responsibility for the research activity planning and execution.
**Resources** – Provision of study materials, reagents, materials, patients, laboratory samples, animals, instrumentation, computing resources, or other analysis tools.
**Software** – Programming, software development; designing computer programs; implementation of the computer code and supporting algorithms; testing of existing code components.
**Supervision** – Oversight and leadership responsibility for the research activity planning and execution, including mentorship external to the core team.
**Validation** – Verification, whether as a part of the activity or separate, of the overall replication/reproducibility of results/experiments and other research outputs.
**Visualization** – Preparation, creation and/or presentation of the published work, specifically visualization/data presentation.
**Writing** – original draft – Preparation, creation and/or presentation of the published work, specifically writing the initial draft (including substantive translation).
**Writing** – review & editing – Preparation, creation and/or presentation of the published work by those from the original research group, specifically critical review, commentary or revision – including pre- or post-publication stages.

# CyberInfrastructure - CI

1. Hardware
   a. Local hardware, on-prem (laptop, workstations, towers, …)
   b. Clusters: large computers dedicated to large scale simulation. Hyak.
   c. Cloud Computing (next)
2. Software
   a. Programming languages (python)
   b. Operating System: OSX, Linux, Windows
   c. Scripted code (e.g., *.py)
   d. Containers are lightweight, portable, self-contained environment that include applications and all dependencies necessary to run across different computing environment. Docker is a tool to manage containers.

Marine Denolle @ UW mdenolle@uw.edu

# What is a cloud?



*Huge buildings filled with servers owned by Microsoft, Amazon, Google, etc.*

Marine Denolle
mdenolle@uw.edu
Zoe Krauss
zkrauss@uw.edu
Naomi Alterman
naomila@uw.edu

**Compute**

"Virtual machines" (VMs)



**Storage**

File stores, "Object storage"

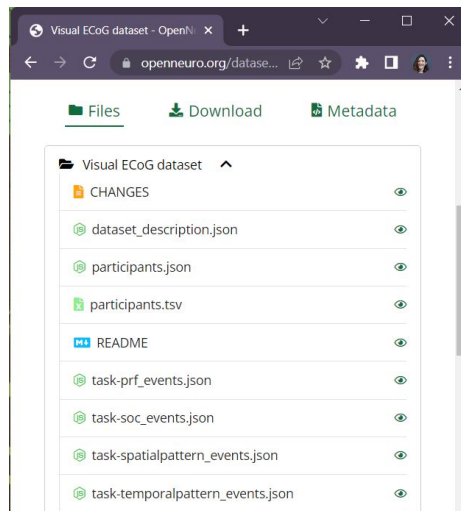# Virtual machines: emulators of computer system

Instance type

CPU   Memory   GPU   Storage

Image

Operating system
Python
Libraries
Home directory

Marine Denolle mdenolle@uw.edu
Naomi Alterman naomila@uw.edu
Rob Fatland rob5@uw.edu

Marine Denolle mdenolle@uw.edu
Naomi Alterman naomila@uw.edu
Rob Fatland rob5@uw.edu

# Storage



## Block Storage
Fastest, usually one per-VM,
Usually 4 GB - 512 GB
$$$ ~$0.16 / GB

## File Storage
Most flexible
< 100 TB
$$ ~$0.06 / GB

## Object Storage
Most cost-effective
Unlimited
$ ~$0.01 / GB

# How to access data?

**Scedc** data ~ 100TB
**Poro tomo DAS** data ~ 50TB
(soon) **Earthscope/IRIS** data

**NASA**
**SAR**

Blob Storage

**Sentinel**
**Landsat**
**NOAA**

Planetary Computer

Earth Engine

70 PBs of geospatial data
**Landast**
**Sentinel**
**DEMs**
**Surface water**

10

Marine Denolle mdenolle@uw.edu

# How to access compute?

## Single VM
### Exploration & ML training



- Ssh to Linux OS - install all packages or preload Machine Image (AMI)



- Use institution-hosted jupyterHub (at cost to institutions)

## Many VMs
### Deployment @ scale



- Containerize software using Docker



- DIY cluster with kubernetes



- Batch services to ease cluster management

11

Marine Denolle mdenolle@uw.edu

# Cloud Concepts with AWS

## EC2: Elastic Compute.



Storage optimized instances are designed for workloads that require high, sequential read and write access to very large data sets on local storage. They are optimized to deliver tens of thousands of low-latency, random I/O operations per second (IOPS) to applications.

**PAGE CONTENT**

General Purpose

Compute Optimized

Memory Optimized

Accelerated Computing

**Storage Optimized**

HPC Optimized

Instance Features

Measuring Instance Performance

| I4g | Im4gn | Is4gen | I4i | I3 | I3en | D3 | D3en | D2 | H1 |

Amazon EC2 Im4gn instances are powered by AWS Graviton2 processors and provide the best price performance for storage-intensive workloads in Amazon EC2. They provide up to 40% better price performance, up to 44% lower cost per TB of storage over I3 instances.

**Features:**

- Powered by AWS Graviton2 processors
- Featuring up to 30 TB of NVMe SSD instance storage with AWS Nitro SSDs that provide up to 60% lower I/O latency and up to 75% reduced latency variability compared to I3 and I3en instances and feature always-on encryption
- Optimized for workloads that map to 4 GB of memory per vCPU
- 2x NVMe SSD storage density per vCPU compared to I3 instances
- Up to 100 Gbps of network bandwidth using Elastic Network Adapter (ENA)-based Enhanced Networking
- Support for Elastic Fabric Adapter on im4gn.16xlarge
- Up to 38 Gbps of bandwidth to the Amazon Elastic Block Store
- Built on the AWS Nitro System, a combination of dedicated hardware and lightweight hypervisor
- Support for Torn Write Prevention (TWP) to enable additional performance and reduce latencies with database workloads such as MySQL and MariaDB.

| Instance Size | vCPU | Memory (GiB) | Instance Storage (GB) | Network Bandwidth (Gbps)*** | EBS Bandwidth (Gbps) |
|---|---|---|---|---|---|
| Im4gn.large | 2 | 8 | 1 x 937 AWS Nitro SSD | Up to 25 | Up to 9.5 |
| Im4gn.xlarge | 4 | 16 | 1 x 1875 AWS Nitro SSD | Up to 25 | Up to 9.5 |
| Im4gn.2xlarge | 8 | 32 | 1 x 3750 AWS Nitro SSD | Up to 25 | Up to 9.5 |
| Im4gn.4xlarge | 16 | 64 | 1 x 7500 AWS Nitro SSD | 25 | 9.5 |
| Im4gn.8xlarge | 32 | 128 | 2 x 7500 AWS Nitro SSD | 50 | 19 |
| Im4gn.16xlarge | 64 | 256 | 4 x 7500 AWS Nitro SSD | 100 | 38 |

**$ 0.18 / hour**

**$ 5.8 / hour**

**Compute**

"Virtual machines" (VMs)

Marine Denolle mdenolle@uw.edu

# General purpose EC2

# Small - testing - microservice

$0.0058 / hour

$0.3712 / hour

Marine Denolle mdenolle@uw.edu

Amazon EC2 T2 instances are Burstable Performance Instances that provide a baseline level of CPU performance with the ability to burst above the baseline.

T2 Unlimited instances can sustain high CPU performance for as long as a workload needs it. For most general-purpose workloads, T2 Unlimited instances will provide ample performance without any additional charges. If the instance needs to run at higher CPU utilization for a prolonged period, it can also do so at a flat additional charge of 5 cents per vCPU-hour.

The baseline performance and ability to burst are governed by CPU Credits. T2 instances receive CPU Credits continuously at a set rate depending on the instance size, accumulating CPU Credits when they are idle, and consuming CPU credits when they are active. T2 instances are a good choice for a variety of general-purpose workloads including micro-services, low-latency interactive applications, small and medium databases, virtual desktops, development, build and stage environments, code repositories, and product prototypes. For more information see Burstable Performance Instances.

**Features:**

- Up to 3.3 GHz Intel Xeon Scalable processor (Haswell E5-2676 v3 or Broadwell E5-2686 v4)

- High frequency Intel Xeon processors

- Burstable CPU, governed by CPU Credits, and consistent baseline performance

- Low-cost general purpose instance type, and Free Tier eligible*

- Balance of compute, memory, and network resources

* t2.micro only. If configured as T2 Unlimited, charges may apply if average CPU utilization exceeds the baseline of the instance. See documentation for more details.

| Instance | vCPU* | CPU Credits / hour | Mem (GiB) | Storage | Network Performance |
|----------|-------|--------------------|-----------|---------|---------------------|
| t2.nano | 1 | 3 | 0.5 | EBS-Only | Low |
| t2.micro | 1 | 6 | 1 | EBS-Only | Low to Moderate |
| t2.small | 1 | 12 | 2 | EBS-Only | Low to Moderate |
| t2.medium | 2 | 24 | 4 | EBS-Only | Low to Moderate |
| t2.large | 2 | 36 | 8 | EBS-Only | Low to Moderate |
| t2.xlarge | 4 | 54 | 16 | EBS-Only | Moderate |
| t2.2xlarge | 8 | 81 | 32 | EBS-Only | Moderate |

13

# Heterogeneous computing EC2 instances.

# Machine Learning training or big Xcorr-ready workflows

**$98.32 / hour**

Marine Denolle mdenolle@uw.edu

# Cloud Concepts with AWS

**EC2: Elastic Compute At Scale**

**Kubernetes**

Kubernetes is a container orchestration platform that automates the deployment, scaling, and management of containerized applications.

**AWS Batch**

AWS Batch is a fully managed service that enables users to run batch computing jobs on the AWS cloud efficiently and at any scale.

*Kubernetes* is a versatile platform for containerized applications with complex orchestration needs, while *AWS Batch* is tailored for efficiently running and managing batch processing jobs with minimal overhead.

### *This course will show how to use AWS Batch*

Marine Denolle mdenolle@uw.edu

# Cloud Concepts with AWS

**Pricing**

| Amazon S3 | Overview | Features ▾ | Storage classes ▾ | Pricing | Security | Resources ▾ | FAQs |
|-----------|----------|-----------|-------------------|---------|----------|-------------|------|

Please note that we list Storage Requests and Data Retrievals Pricing below the Storage Pricing table.

Region: [ US West (Oregon) ◆ ]

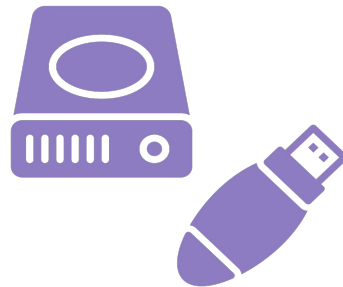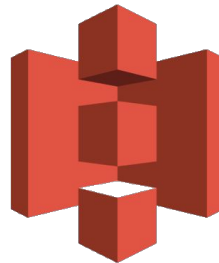| | Storage pricing |
|---|---|
| **S3 Standard** – General purpose storage for any type of data, typically used for frequently accessed data | |
| First 50 TB / Month | $0.023 per GB |
| Next 450 TB / Month | $0.022 per GB |
| Over 500 TB / Month | $0.021 per GB |
| **S3 Intelligent - Tiering *** - Automatic cost savings for data with unknown or changing access patterns | |
| Monitoring and Automation, All Storage / Month (Objects > 128 KB) | $0.0025 per 1,000 objects |
| Frequent Access Tier, First 50 TB / Month | $0.023 per GB |
| Frequent Access Tier, Next 450 TB / Month | $0.022 per GB |
| Frequent Access Tier, Over 500 TB / Month | $0.021 per GB |
| Infrequent Access Tier, All Storage / Month | $0.0125 per GB |
| Archive Instant Access Tier, All Storage / Month | $0.004 per GB |
| **S3 Intelligent - Tiering *** - Optional asynchronous Archive Access tiers | |
| Archive Access Tier, All Storage / Month | $0.0036 per GB |
| Deep Archive Access Tier, All Storage / Month | $0.00099 per GB |

**50 TB = 1,150$/mo**

**50 TB = 625$/mo**

**Storage**

File stores, "Object storage"

Amazon S3

Marine Denolle mdenolle@uw.edu

# Who Pays for Compute?

**You**. You own the account, you pay with a credit card.

**CloudBank**. An [NSF project](#) to support the adoption of cloud computing in NSF supported research and education. You may apply to cloudbank credits when submitting a proposal to specific RFPs.

**Your institution.** They may support cloud credit applications through partnerships between cloud providers and the institutions.

**The national archives.** It is unlikely that archives such as Earthscope Data Services will pay for large-scale computing, but they will offer basic jupyter Hubs for data explorations.

Marine Denolle mdenolle@uw.edu

# Abstracting the Cloud

Infrastructure-as-a-service (Iaas) is a way to abstract the compute and make the experience of running a job as *local* as possible.

Coiled is an example for Dask, Python users.

- It borrows from Dask.
- It maps your local environment and run it there (but the local environment has to be simple)

No additional fee for small compute but the Cloud costs.

Marine Denolle mdenolle@uw.edu

Marine Denolle mdenolle@uw.edu

**What resources are available to me on the cloud?**

aws    Amazon Web Services

*Ranges from:*
Fully "built" systems with operating systems, programs like Jupyter, access to data, already installed

Microsoft Azure

Google Cloud

Empty "from scratch" instances- you choose CPU, GPU, memory, etc. Infinitely scalable

Marine Denolle mdenolle@uw.edu

# Software

Programming language of the class: Python

Everybody opens up: Terminal, VsCode, of the JupyterHub.

We will create a conda environment and add the dependencies in the python environment.

Sign up for Github Education account if you can, or sign up for CoPilot

Marine Denolle mdenolle@uw.edu

# Get up to speed

Some useful resources

- [https://www.earthdatascience.org/](https://www.earthdatascience.org/)

- [https://software-carpentry.org/lessons/](https://software-carpentry.org/lessons/)

The tutorials and course materials will be done in Python.

Marine Denolle [mdenolle@uw.edu](mailto:mdenolle@uw.edu)

# Build Cohort - Chat about Data & cool problems

## Get up to speed with computing

# 3. Basic shell

GUI: Graphical User Interface
     Instructions from human to computer are done by clicking buttons on a mouse
     Fun, easy, but not good if you want to repeat your task 1000s (watch out for carpal
     tunnel syndrome!)
CLI: Command Line Interface
     Use terminal, rudimentary commands

Tutorials: EarthDataScience + software carpentry

pwd
ls
cd
mkdir
rmdir
cp
rm

Marine Denolle mdenolle@uw.edu

# 4. Git, GitHub, GitLab

**Git** is an *open source version control tool*, **GitHub** and **GitLab** are companies that hosts Git repositories in the web and provides a web interface to interact with repos they host.

More in depth tutorials:

Software Carpentry

EarthDataScience

Marine Denolle mdenolle@uw.edu

# 4. Version Control:

**Iterations over docs**
**Collaborative comments and edits**
**How to keep track of changes?**

**Code Version Control is the scripting analog of Tracked changes**



Marine Denolle mdenolle@uw.edu