

# **Scalable and Computationally Reproducible Approaches to Arctic Research**

Matt Jones, Bryce Mecum, Jeanette Clark, Sam Csik

September 19, 2022

# Table of contents

<b>Preface</b>	<b>3</b>
About . . . . .	3
Schedule . . . . .	3
Code of Conduct . . . . .	3
About this book . . . . .	5
<b>1 Welcome and Introductions</b>	<b>6</b>
<b>2 Remote Computing</b>	<b>9</b>
2.1 Learning Objectives . . . . .	9
<b>3 Python Syntax Refresher</b>	<b>10</b>
3.1 Learning Objectives . . . . .	10
<b>4 Pleasingly Parallel Computing</b>	<b>11</b>
4.1 Learning Objectives . . . . .	11
<b>5 Documenting and Publishing Data</b>	<b>12</b>
5.1 Learning Objectives . . . . .	12
<b>References</b>	<b>13</b>

# Preface

## About

This 5-day in-person workshop will provide researchers with an introduction to advanced topics in computationally reproducible research in python and R, 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](#).

	Monday	Tuesday	Wednesday	Thursday	Friday	
08:00-08:30	Coffee (optional)	Coffee (optional)	Coffee (optional)	Coffee (optional)	Coffee (optional)	
08:30-09:00	1. Welcome and Course Overview (Jeanette)	6. Data structures and formats for large data (Bryce)	10. Spatial and Image Data using GeoPandas (Jeanette)	15. Google Earth Engine (Ingmar, Sam)	19. What is cloud computing anyways? (Matt)	
09:00-09:30			11. Data futures: Parquet and Arrow (Jeanette)			
09:30-10:00	2. Remote computing (Sam)					
10:00-10:30						
10:30-11:00	BREAK	BREAK	BREAK	BREAK	BREAK	
11:00-11:30	3. Python programming on clusters (Jeanette)	7. Parallelization with Dask (Bryce)	12. Software Design II (Bryce)	16. Billions of Ice Wedge Polygons (Chandi)	20. Reproducibility redux via containers (Bryce)	
11:30-12:00					Survey Feedback Q & A	
12:00-12:30	Lunch	Lunch	Lunch	Lunch	Adjourn	
12:30-13:00						
13:00-13:30	4. Pleasingly Parallel Programming (Matt)	8. Group project I Data staging and pre-processing (Jeanette)	13. Group project II Parallel data processing (Jeanette)	17. Group project III Visualizing big geospatial data (Jeanette)		
13:30-14:00						
14:00-14:30						
14:30-15:00						
15:00-15:30	Break	Break	Break	Break		
15:30-16:00	5. Documenting and Publishing Data (Daphne)	9. Software design I (Bryce)	14. Data Ethics (Matt)	18. Workflows for data staging and publishing (Jeanette)		
16:00-16:30			Breather Catch-up			
16:30-17:00	Q&A	Q&A	Q&A	Q&A		

## 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](https://creativecommons.org/licenses/by/4.0/).

**Citation:** Matthew B. Jones, Bryce Mecum, S. Jeanette Clark, Samantha Csik. 2022. Scalable and Computationally Reproducible Approaches to Arctic Research.

**Additional contributors:** Amber E. Budden, Natasha Haycock-Chavez, Noor Johnson, Stephanie Hampton, Jim Regetz, Bryce Mecum, Julien Brun, Julie Lowndes, Erin McLean, Andrew Barrett, David LeBauer, Jessica Guo.


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

# 1 Welcome and Introductions



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.

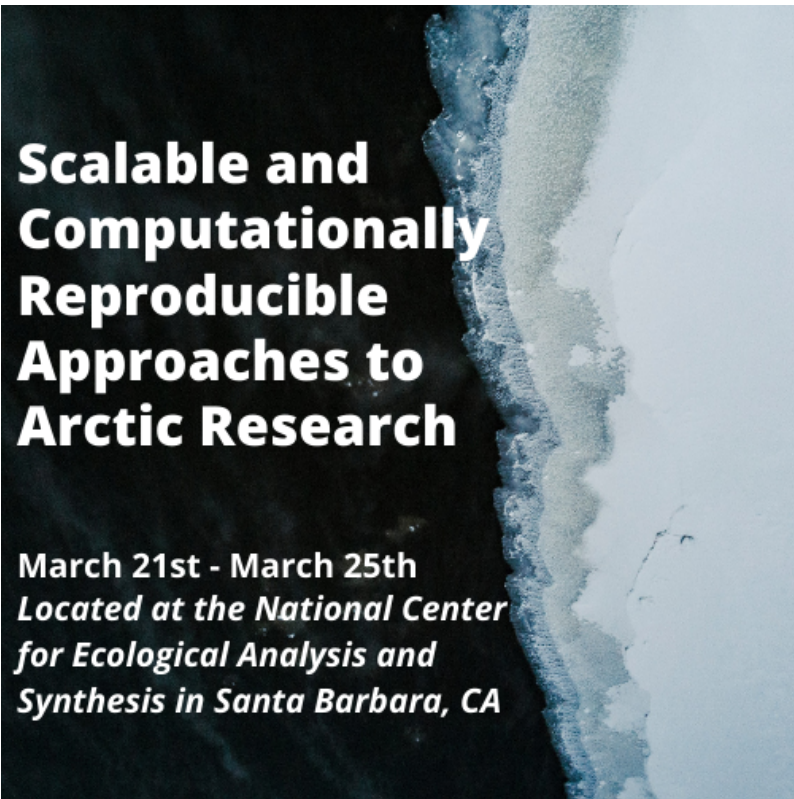




# **Reproducible Practices for Arctic Research Using R**

**February 14th - February  
18th, 2022**

***This course will be taught  
virtually***



# **Scalable and Computationally Reproducible Approaches to Arctic Research**

March 21st - March 25th  
*Located at the National Center  
for Ecological Analysis and  
Synthesis in Santa Barbara, CA*



## **2 Remote Computing**

### **2.1 Learning Objectives**

## **3 Python Syntax Refresher**

### **3.1 Learning Objectives**

## **4 Pleasingly Parallel Computing**

### **4.1 Learning Objectives**

## **5 Documenting and Publishing Data**

### **5.1 Learning Objectives**

## References