# Data Management

## Introduction to High Performance Computing

# Instructors

– – – –

**Julie Goldman**

Research Data Services Librarian

**Meghan Kerr**

Archivist and Records Manager

**Radhika Khetani**

Training Director, Research Scientist
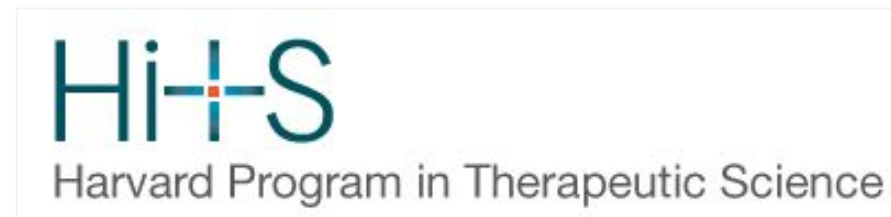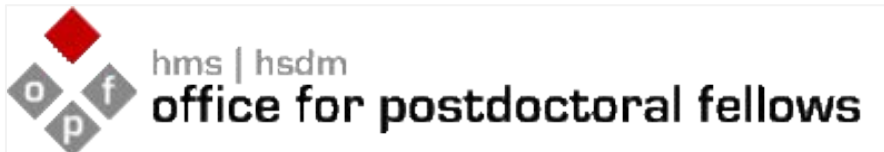
**Meeta Mistry**

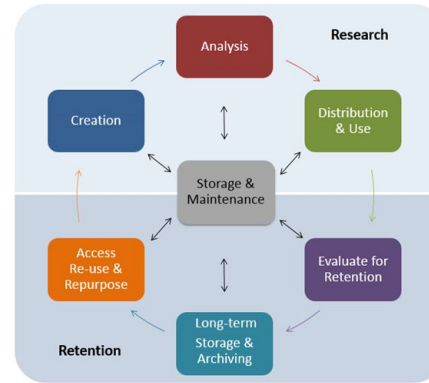Research Scientist

**Kathleen Keating**

Research Computing Consultant

Slides: https://datamanagement.hms.harvard.edu/class-materials

# Harvard Medical School | Data Management Working Group

## Countway Library of Medicine
*An Alliance of the Harvard Medical School and Boston Medical Library*

## Center *for the* History *of* Medicine

# Harvard Chan Bioinformatics Core

## hms | hsdm
### office for postdoctoral fellows

## HARVARD MEDICAL SCHOOL | OFFICE FOR Academic and Research Integrity

## Department of Systems Biology

## HARVARD MEDICAL SCHOOL
### Research Information Technology Solutions - RITS

## HMS Information Technology

## ICCB-Longwood Screening Facility

## DRSC/TRiP Functional Genomics

## Hi+S
### Harvard Program in Therapeutic Science

# Harvard Biomedical Data Management Website
https://datamanagement.hms.harvard.edu

# Why Manage Data?

----

- Running the same workflow can be labour intensive

- Manual manipulation of data files:

    - is often not captured in documentation

    - is hard to reproduce

    - is hard to troubleshoot, review, or improve

- Hard to find poorly organized, documented data

- Hard to analyze poorly recorded workflows

https://datamanagement.hms.harvard.edu/overview

# Why HPC?

----

- High Performance Computing makes workflows more efficient

    - If you work with a lot of data or you have really complex computations, scheduling scripts reduces computation time

    - Automated workflows makes you more productive and also improves the reproducibility of your work by allowing you to save and repeat them

- Using a command line interface to work with files

    - Every step can be captured in the shell script and allow reproducibility and easy troubleshooting

- Offers storage space for active data files and shared drives for sharing data between labs

# Training Materials
https://tinyurl.com/hpc-july11

# Workshop Outline

| Lessons | Estimated Duration |
|---|---|
| Intro to High-Performance Computing | 25 min |
| Intro to O2 | 55 min |

# Tying it Together

----

**Why Data Management:**

Not a prerequisite of HPC, but data should be organized in a clear and predictable manner.

Taking the time to structure your research data and filenaming conventions in a consistent and predictable manner is certainly a significant step towards getting the most out of data analysis.

**Why HPC:**

Allows you to reduce computation time and help make analyses more efficient. Using a cluster offers advantages such as: speed, volume, efficiency, cost, and convenience.

Automate repetitive tasks and capture small data manipulation steps that are normally not recorded to make research reproducible.

# Questions?



https://datamanagement.hms.harvard.edu

# Upcoming Workshops / Seminars

————

**Creating Meaningful Data: Metadata Essentials**

Thursday, August 8
12:30 - 1:30 pm
Countway Library 403 Classroom

bit.ly/RDM-Summer19

**Version Control for Scripts, Data & Text documents**

Wednesday, August 21
1:30 - 3:00 pm
TMEC 227 Mini amphitheater

bit.ly/RDM-Summer19

bit.ly/rdm-survey