# Introductions!

Shannan Ho Sui
*Director*

Meeta Mistry
*Associate Director*

Lorena Pantano
*Director of Bioinformatics Platform*

John Quackenbush
*Faculty Advisor*

Upen Bhattarai

Heather Wick

Will Gammerdinger

Noor Sohail

Elizabeth Partan

Alex Bartlett

Emma Berdan

James Billingsley

Zhu Zhuo

Maria Simoneau

Shannan Ho Sui
*Director*

Meeta Mistry
*Associate Director*

Lorena Pantano
*Director of Bioinformatics Platform*

John Quackenbush
*Faculty Advisor*

Upen Bhattarai

Heather Wick

Will Gammerdinger

Noor Sohail

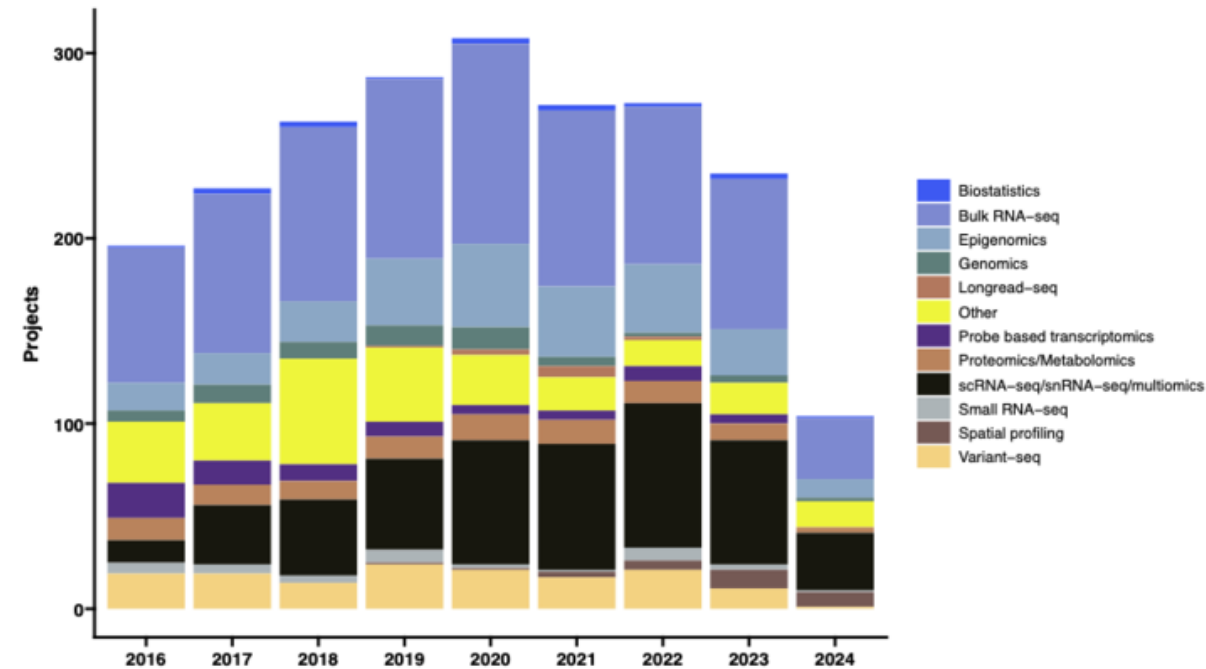Elizabeth Partan

Alex Bartlett

Emma Berdan

James Billingsley

Zhu Zhuo

Maria Simoneau

# Consulting

❖ Transcriptomics: Bulk, single cell, small RNA

❖ Epigenomics: ChIP-seq, CUT&RUN, ATAC-seq, DNA methylation

❖ Variant discovery: WGS, resequencing, exome-seq and CNV

❖ Multiomics integration

❖ Spatial biology

❖ Experimental design and grant support

# Consulting

❖ Transcriptomics: Bulk, single cell, small RNA

❖ Epigenomics: ChIP-seq, CUT&RUN, ATAC-seq, DNA methylation

❖ Variant discovery: WGS, resequencing, exome-seq and CNV

❖ Multiomics integration

❖ Spatial biology

❖ Experimental design and grant support

NIEHS

# Training

❖ Hands-on workshops design to reflect best practices,

reproducibility and an emphasis on experimental design

❖ Basic Data Skills

- ❖ Shell

- ❖ R

❖ Advanced Topics: Analysis of high-throughput sequencing data

- ❖ Chromatin Biology

- ❖ Bulk RNA-seq

- ❖ Differential Gene Expression

- ❖ scRNA-seq

- ❖ Variant Calling

❖ Current Topics in Bioinformatics

**https://bioinformatics.sph.harvard.edu/training**

# Training

❖ Hands-on workshops design to reflect best practices, reproducibility and an emphasis on experimental design

  ❖ Basic Data Skills

    ❖ Shell

    ❖ R

  ❖ Advanced Topics: Analysis of high-throughput sequencing data

    ❖ Chromatin Biology

    ❖ Bulk RNA-seq

    ❖ Differential Gene Expression

    ❖ scRNA-seq

    ❖ Variant Calling

  ❖ Current Topics in Bioinformatics

**https://bioinformatics.sph.harvard.edu/training**

Setting up to perform Bioinformatics analysis

# Setting up…

❖ Introduction to the command-line interface (shell, Unix, Linux)

   ❖ Dealing with large data files

   ❖ Performing bioinformatics analysis

   ❖ Using tools

   ❖ Accessing and using compute clusters

❖ Introduction to R

   ❖ Parsing and working with smaller BED files

   ❖ Statistical analysis, e.g. differential binding analysis

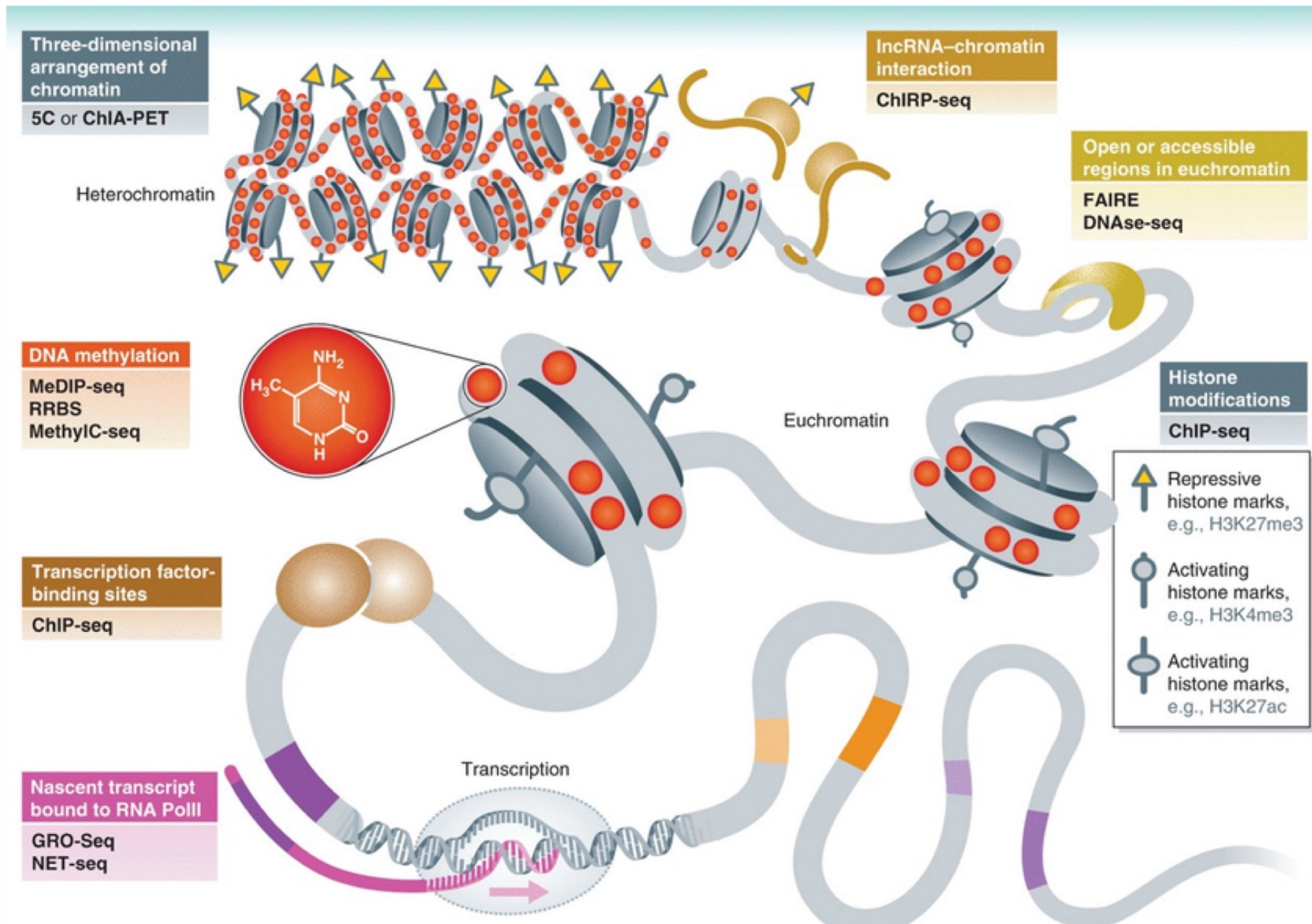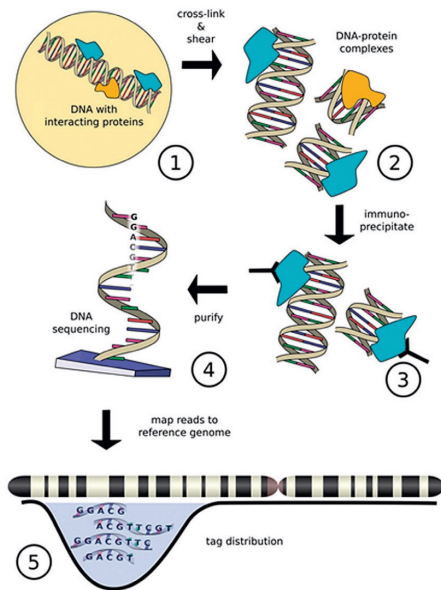   ❖ Generating figures from complex data

# Workshop scope

Bioinformatic Data Analysis

**Three-dimensional arrangement of chromatin**
5C or ChIA-PET

Heterochromatin

**lncRNA–chromatin interaction**
ChIRP-seq

**Open or accessible regions in euchromatin**
FAIRE
DNAse-seq

**DNA methylation**
MeDIP-seq
RRBS
MethylC-seq

Euchromatin

**Histone modifications**
ChIP-seq

Repressive histone marks, e.g., H3K27me3

Activating histone marks, e.g., H3K4me3

Activating histone marks, e.g., H3K27ac

**Transcription factor-binding sites**
ChIP-seq

**Nascent transcript bound to RNA PolII**
GRO-Seq
NET-seq

Transcription

Figure adapted from Soon WW, Hariharan M, Snyder MP, "High throughput sequencing for biology and medicine". Molecular Systems Biology 9:640 2013

# Genomic Methods for Profiling Chromatin

## ChIP-seq

## ATAC-seq

## CUT&RUN

# Learning Objectives

❖ Describe important considerations for setting up a successful ChIP-seq, CUT&RUN or ATAC-seq experiment

❖ Describe the steps in an ChIP-seq analysis workflow (from sequence data to peak calls) and contrast any differences for CUT&RUN and ATAC-seq analyses

❖ Learn how to handle various file formats encountered when analyzing ChIP-seq and related data

❖ Implement shell scripts on a high-performance compute cluster to perform the above steps

The Workflow

Sequence Reads

FASTQ → Sequence Data QC

Alignment to Genome

SAM/BAM

Filter out Duplicates & Multi-mappers*

BAM

* optionally filter out blacklisted regions

Peak Calling

Peak Call QC
Replicate Concordance
Qualitative Assessment
Quantitative Assessment

BED

Differential Enrichment (*between* groups)

Combine Replicates (*within* group)

Characterize Binding Profile

Annotation    Functional Enrichment    Motif Analysis

# The Workflow

*Boxes in green represent parts of the workflow that will not be covered in this workshop*

The Workflow

*Boxes in green represent parts of the workflow that will not be covered in this workshop*

# Logistics

# Course schedule

Pre-reading:

- Please **study the contents** and **work through all the exercises** within the following lessons:
    - Shell basics review
    - Best Practices in Research Data Management (RDM)
    - Working in an HPC environment
    - A review of high-throughput sequencing methods for understanding chromatin biology

## Day 1

| Time | Topic | Instructor |
|------|-------|------------|
| 09:30 - 09:45 | Workshop Introduction | Meeta |
| 09:45 - 11:00 | Understanding chromatin biology using high-throughput sequencing | Dr. Shannan Ho Sui |
| 11:00- 11:05 | Break | |
| 11:05 - 11:20 | HPC review Q&A | Will |
| 11:20 - 11:50 | Dataset overview and project organization | Will |

**https://tinyurl.com/hbc-chipseq**

# Course materials

❖ We continuously update our materials to reflect changes in the field/software

**Peak calling with MACS2**

View on GitHub

Contributors: Meeta Mistry, Jihe Liu, Radhika Khetani, Mary Piper, Will Gammerdinger

Approximate time: 60 minutes

## Learning Objectives

- Describe the different components of the MACS2 peak calling algorithm
- Describe the parameters involved in running MACS2
- List and describe the output files from MACS2

**https://tinyurl.com/hbc-chipseq**

# Single Screen & 3 Windows

# Single Screen & 3 Windows



Zoom

Our
Recommendation

# Single Screen & 3 Windows



*Our Recommendation*

# Single Screen & 3 Windows



*Our Recommendation*

**Terminal**

# Single Screen & 3 Windows

# Course participation

❖ Mandatory review of self-learning lessons and assignments

❖ Attendance required for all classes

❖ Your questions and active participation drive learning

❖ **We look forward to all of your questions!**

# Course participation

❖ At-home lessons and exercises after each session

❖ Cover material not previously discussed

❖ Provides us feedback to help pace the course appropriately

❖ 3-5 hours to complete

❖ Homework load is heavier in the beginning of this workshop series and tapers off

# Using AI for Assignments

❖ Do
  ❖ Try to resolve error messages with it
  ❖ Test code written by AI on a dataset where you have expected results
  ❖ Take the time to review the generated code line-by-line

❖ Don't
  ❖ Implement it in replacement to learning
  ❖ Write code that you don't understand
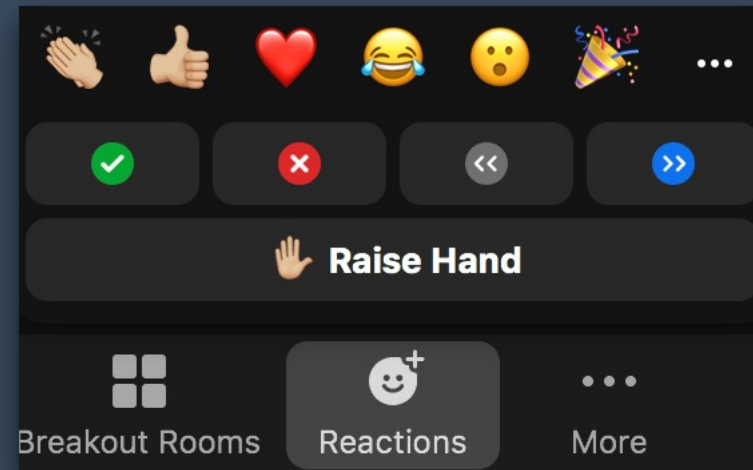  ❖ Assume the output from an AI process is correct

# Odds & Ends

❖ Quit/minimize all applications that are not required for class
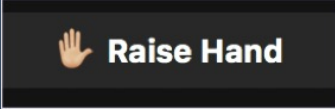
❖ Are you all set?

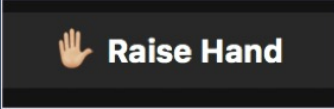    ❖  ✅  = "agree", "I'm all set"

    ❖  ❌  = "disagree", "I need help"

# Odds & Ends

❖ Questions for the presenter?

   ❖ Post the question in the Chat window OR

   ❖ ✋ Raise Hand when the presenter asks for questions

   ❖ Let the Moderator know

# Odds & Ends

❖ Questions for the presenter?

  ❖ Post the question in the Chat window OR

  ❖ ✋ Raise Hand  when the presenter asks for questions

  ❖ Let the Moderator know

❖ Technical difficulties with software?

  ❖ Start a private chat with the Troubleshooter with a description of the problem

# Thanks!

❖ Kathleen Chappell and Andy Bergman from HMS-RC

❖ Data Carpentry

# Contact Us

HBC

Harvard Chan Bioinformatics Core

- ❖ *HBC training team:* hbctraining@hsph.harvard.edu

- ❖ *HBC consulting:* bioinformatics@hsph.harvard.edu

- ❖ *O2 (HMS-RC):* rchelp@hms.harvard.edu