

# Introduction to R

## Harvard Chan Bioinformatics Core

<https://tinyurl.com/hbc-r-online>

Sponsored by DF/HCC, CFAR, and HMS Foundry



Shannan Ho Sui  
*Director*



Meeta Mistry  
*Associate Director*



John Quackenbush  
*Faculty Advisor*



Emma Berdan



Heather Wick



Will Gammerdinger



Noor Sohail



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

<http://bioinformatics.sph.harvard.edu/>



NIEHS

---



# Training

A key component of the HBC's mission is its training initiative. Our dedicated training team holds workshop to help researchers at Harvard better understand analytical methods for NGS data.

HBC's training team is made up of four PhD-level scientists who devote substantial time to material development, training and community building/outreach. All members of the training team also participate in consultations on research projects to ensure they remain up-to-date on current best practices in NGS analysis.

Our hands-on workshops focus on **basic data skills** and **analysis of high-throughput sequencing data**, with an emphasis on **experimental design**, current **best practices** and **reproducibility**. Our workshops are designed for **wet-lab biologists** aiming to independently design sequencing-based experiments and analysing the resulting data.

We offer three types of workshops:

1. Short, 3-hour monthly workshops (*Current topics in bioinformatics*)
2. Basic Data Skills\*\*
3. Advanced Topics: Analysis of high-throughput sequencing (NGS) data\*\*

*\*\*The basic data skills workshops serve as the foundation for the advanced workshops.*

<http://bioinformatics.sph.harvard.edu/training/>

<https://hbctraining.github.io/main/>

# Training

A key component of the HBC's mission is to provide training for researchers at Harvard and beyond.

HBC's training team is made up of experts in training and community building who work on research projects to ensure our training is effective.

Our hands-on workshops focus on **bioinformatics**, with an emphasis on **experimental design** and **data analysis**. We also offer **wet-lab biologists** and **clinicians** training in **bioinformatics** and **data analysis**.

We offer three types of workshops:

1. Short, 3-hour monthly workshops
2. Basic Data Skills\*\*
3. Advanced Topics: Analysis of high-throughput sequencing data

\*\*The basic data skills workshop is designed for researchers with no prior experience in bioinformatics.



**HARVARD  
T.H. CHAN  
SCHOOL OF PUBLIC HEALTH**

**DF/HCC**  
DANA-FARBER / HARVARD CANCER CENTER



THE HARVARD CLINICAL  
AND TRANSLATIONAL  
SCIENCE CENTER



Our dedicated training team holds workshops to help researchers learn how to analyze **bioinformatics** or **NGS** data.

The training team also devote substantial time to material development, consulting, and teaching. Our training team also participate in consultations on best practices in NGS analysis.

Our workshops focus on **bioinformatics**, with an emphasis on **experimental design** and **reproducibility**. Our workshops are designed for **wet-lab biologists** and **clinicians** to learn about **bioinformatics** and **data analysis** in the context of **high-throughput sequencing** experiments and analysing the resulting **bioinformatics** and **NGS** data.

**bioinformatics**)

**bioinformatics** and **NGS** data)\*\*

and **bioinformatics** for the advanced workshops.

<http://bioinformatics.sph.harvard.edu/training/>

<https://hbctraining.github.io/main/>

# Training

A key component of the HBC's mission is to support researchers at Harvard by providing training.

HBC's training team is made up of experts in training and community based research projects to ensure that our trainees are well prepared for their future careers.

Our hands-on workshops focus on practical skills, with an emphasis on **experimental design** and **bioinformatics**, for **wet-lab biologists** and **bioinformaticians** alike.

We offer three types of workshops:

1. Short, 3-hour monthly workshops
2. Basic Data Skills\*\*
3. Advanced Topics: Analysis of high-throughput sequencing data

\*\*The basic data skills workshop is designed for researchers who have no prior experience with bioinformatics or NGS data analysis.



**HARVARD  
T.H. CHAN  
SCHOOL OF PUBLIC HEALTH**

**DF/HCC**  
DANA-FARBER / HARVARD CANCER CENTER



THE HARVARD CLINICAL  
AND TRANSLATIONAL  
SCIENCE CENTER



Our dedicated training team holds workshops to help researchers learn how to analyze high-throughput sequencing (NGS) data.

In addition to devote substantial time to material development, the training team also participate in consultations on best practices in NGS analysis.

The workshops focus on the analysis of high-throughput sequencing data, with an emphasis on **experimental design**, **bioinformatics**, and **reproducibility**. Our workshops are designed to help researchers design experiments and analyse the resulting data.

**bioinformatics**)

**bioinformatics (NGS) data**\*\*

or the advanced workshops.

<http://bioinformatics.sph.harvard.edu/training/>

<https://hbctraining.github.io/main/>

# Introductions!



Shannan Ho Sui  
*Director*



Meeta Mistry  
*Associate Director*



John Quackenbush  
*Faculty Advisor*



Emma Berdan



Heather Wick



Will Gammerdinger



Noor Sohail



James Billingsley

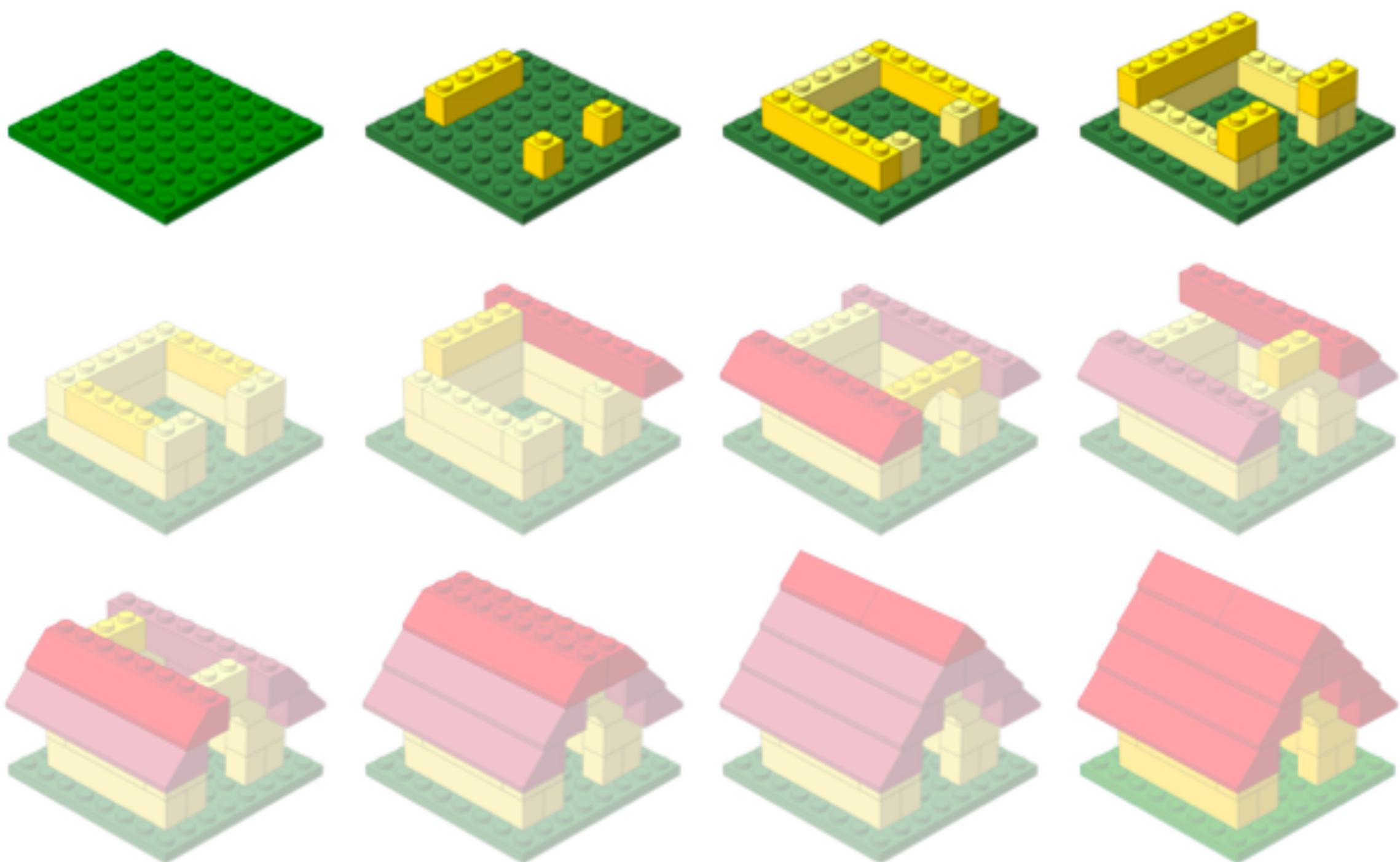


Zhu Zhuo



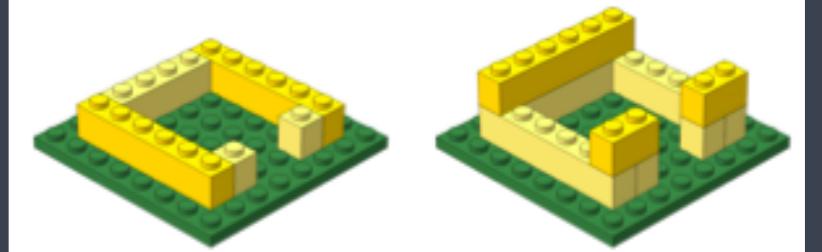
Maria Simoneau

# Workshop Scope...



# Learning R

# Workshop Scope



- ✓ Comfortably use RStudio (a graphical interface for R)
- ✓ Fluently interact with R using RStudio
- ✓ Become familiar with R syntax
- ✓ Understand data structures in R
- ✓ Inspect and manipulate data structures
- ✓ Install packages and use functions in R

# CRAN

## (Comprehensive R Archive Network)



**Available CRAN Packages By Name**

[A](#) [B](#) [C](#) [D](#) [E](#) [F](#) [G](#) [H](#) [I](#) [J](#) [K](#) [L](#) [M](#) [N](#) [O](#) [P](#) [Q](#) [R](#) [S](#) [T](#) [U](#) [V](#) [W](#) [X](#) [Y](#) [Z](#)

<a href="#">A3</a>	Accurate, Adaptable, and Accessible Error Metrics for Predictive Models
<a href="#">abbyyR</a>	Access to Abbyy Optical Character Recognition (OCR) API
<a href="#">abc</a>	Tools for Approximate Bayesian Computation (ABC)
<a href="#">ABCanalysis</a>	Computed ABC Analysis
<a href="#">abc.data</a>	Data Only: Tools for Approximate Bayesian Computation (ABC)
<a href="#">abcdeFBA</a>	ABCDE_FBA: A-Biologist-Can-Do-Everything of Flux Balance Analysis with this package
<a href="#">ABCOptim</a>	Implementation of Artificial Bee Colony (ABC) Optimization
<a href="#">ABCp2</a>	Approximate Bayesian Computational Model for Estimating P2
<a href="#">abcrf</a>	Approximate Bayesian Computation via Random Forests

*CRAN  
Mirrors  
What's new?  
Task Views  
Search  
  
About R  
R Homepage  
The R Journal*

- The main repository for R packages
- Easy to install

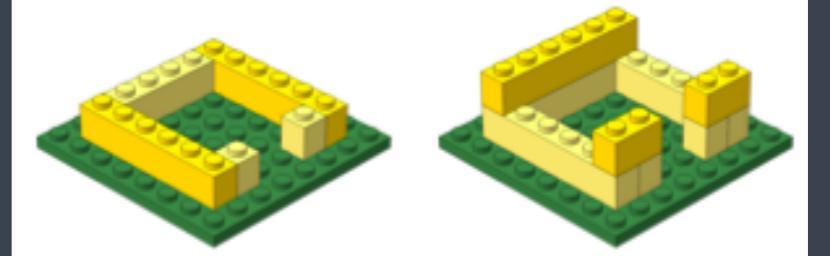
<https://cran.r-project.org/>



- An alternative package repository; “..provides tools for the analysis and comprehension of *high-throughput genomic data*.”
- Includes (but is not limited to) tools for:
  - performing statistical analysis
  - accessing public datasets
- Open source and open development
- Free

[www.bioconductor.org](http://www.bioconductor.org)

# Workshop Scope



- Comfortably use RStudio (a graphical interface for R)
  - Fluently interact with R using RStudio
  - Become familiar with R syntax
  - Understand data structures in R
  - Inspect and manipulate data structures
  - Install packages and use functions in R
- ✓ Visualize data using *ggplot2*
- ✓ Utilize pipes, tibbles and functions from the Tidyverse package suite

# Logistics

# Course webpage

<https://tinyurl.com/hbc-r-online>

# Course webpage

## Introduction to DGE

[View on GitHub](#)

Approximate time: 60 minutes

### Learning Objectives

- Explore different types of normalization methods
- Become familiar with the `DESeqDataSet` object
- Understand how to normalize counts using DESeq2

### Normalization

The first step in the DE analysis workflow is count normalization, which is necessary to make accurate comparisons of gene expression between samples.

```
graph TD; A["Pseudocounts with  
Kallisto, Sailfish, Salmon"] --> B["Read counts  
associated with genes"]; B --> C["Normalization"]; C --> D["Unsupervised clustering analyses"]; C -.-> E["Quality control"]
```

The diagram illustrates the DE analysis workflow. It starts with 'Pseudocounts with Kallisto, Sailfish, Salmon', followed by 'Read counts associated with genes'. This leads to 'Normalization', which then leads to 'Unsupervised clustering analyses'. A bracket on the right side groups 'Normalization' and 'Unsupervised clustering analyses' under the heading 'Quality control'.

# Course schedule online

## Workshop Schedule

### Day 1

Time	Topic	Instructor
10:00 - 10:30	Workshop Introduction	Jihe
10:30 - 11:45	Introduction to R and RStudio	Radhika
11:45 - 12:00	Overview of self-learning materials and homework submission	Mary

### Before the next class:

1. Please **study the contents** and **work through all the code** within the following lessons:
  - o [R Syntax and Data Structure](#)
  - o [Functions and Arguments](#)
  - o [Reading in and inspecting data](#)
2. **Complete the exercises:**
  - o Each lesson above contain exercises; please go through each of them.
  - o **Copy over** your code from the exercises into a text file.
  - o **Upload the saved text file** to [Dropbox](#) the **day before the next class**.

# Course participation

- ▶ Mandatory review of self-learning lessons and assignments
- ▶ Attendance required for all classes
- ▶ Your questions and active participation drive learning and discussion
- ▶ Have fun dabbling with R!



# Homework and Expectations

- ❖ 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

# Contact us!

*HBC training team:* [hbctraining@hsph.harvard.edu](mailto:hbctraining@hsph.harvard.edu)

*HBC consulting:* [bioinformatics@hsph.harvard.edu](mailto:bioinformatics@hsph.harvard.edu)

Twitter

[@bioinfocore](https://twitter.com/bioinfocore)