

# Fundamentals of population genomics

Reid Brennan & Jenny Nascimento-Schulze

[rbrennan@geomar.de](mailto:rbrennan@geomar.de)

# Goals of course

1. Basics of population and quantitative genetic theory
2. Gain understanding of the technologies and data underlying genomics
3. Experience using the command line and R to analyze data
4. Ability to analyze and interpret population genetics data

# Course format

- Lectures introducing fundamentals
  - Mostly short lectures
- Hands on tutorial
  - → analyzing population genomic dataset together
- Independent project
  - In groups, analyze population genomic data and present results via poster

# Independent project

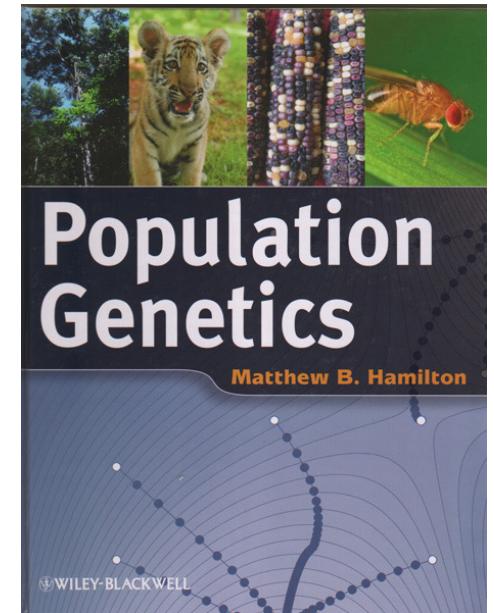
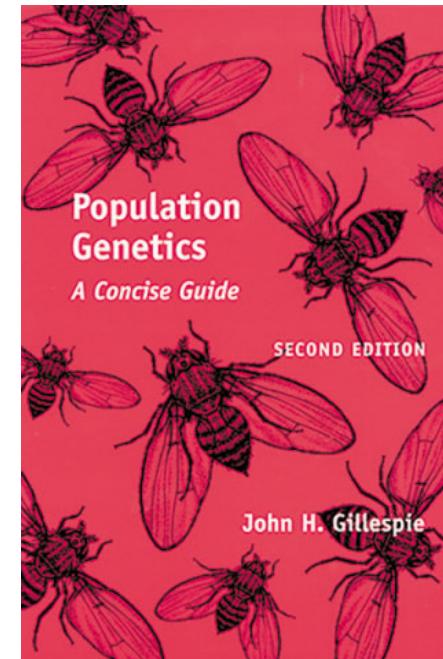
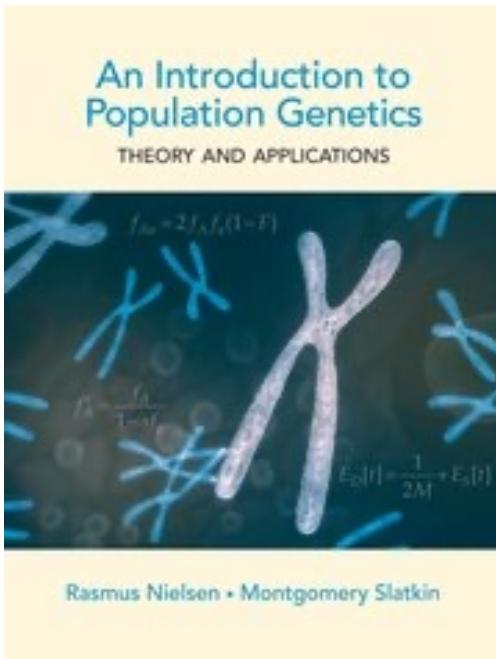
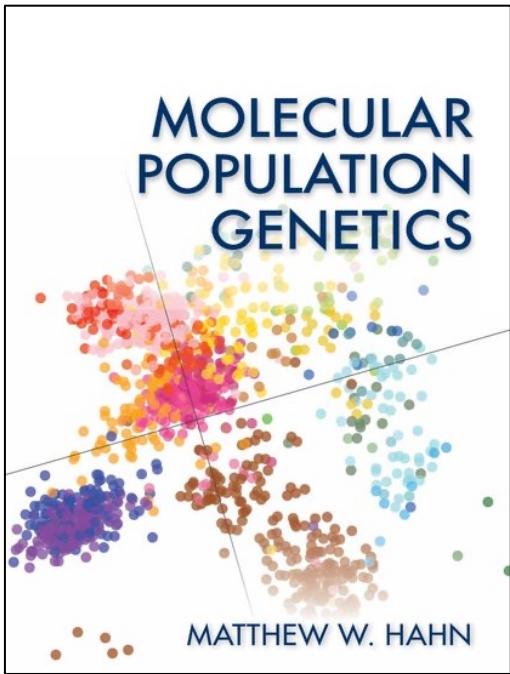
- In groups of 2-3
- From VCF, analyze data to make population genetic inferences
- Poster on July 4
- Grading:
  - Use of tools covered in the course
  - Appropriate conclusions from your results
  - Overall quality of analysis, interpretation, and presentation
- You must hand in:
  - 1. Poster
  - 2. mark down/html document of your analysis.

# Schedule

- Day 1-5: Tutorial
- Day 6-10: Project
- For the most up to date information, slides, tutorials, see the course website:
  - [https://rsbrennan.github.io/EvolutionaryGenomics\\_2024/](https://rsbrennan.github.io/EvolutionaryGenomics_2024/)

# Resources

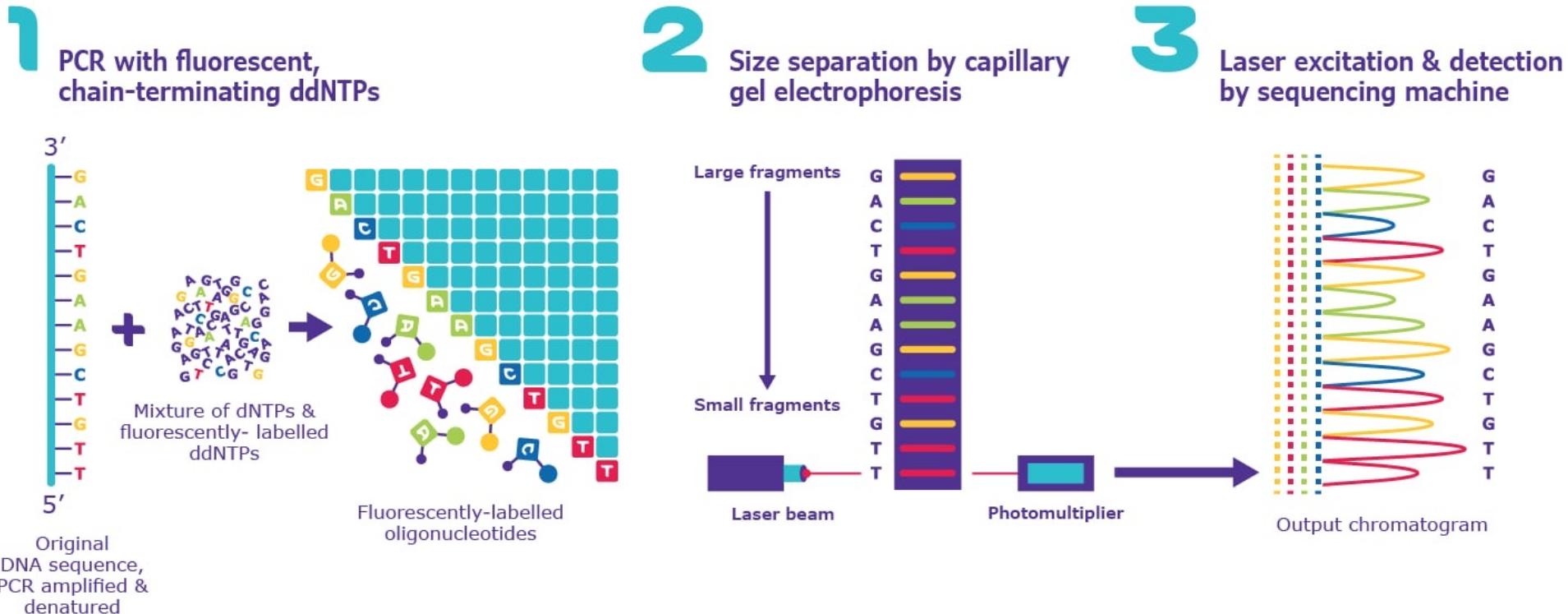
- Graham Coop's Population and Quantitative Genetics notes
  - <https://github.com/cooplabs/popgen-notes/releases>



And many others....

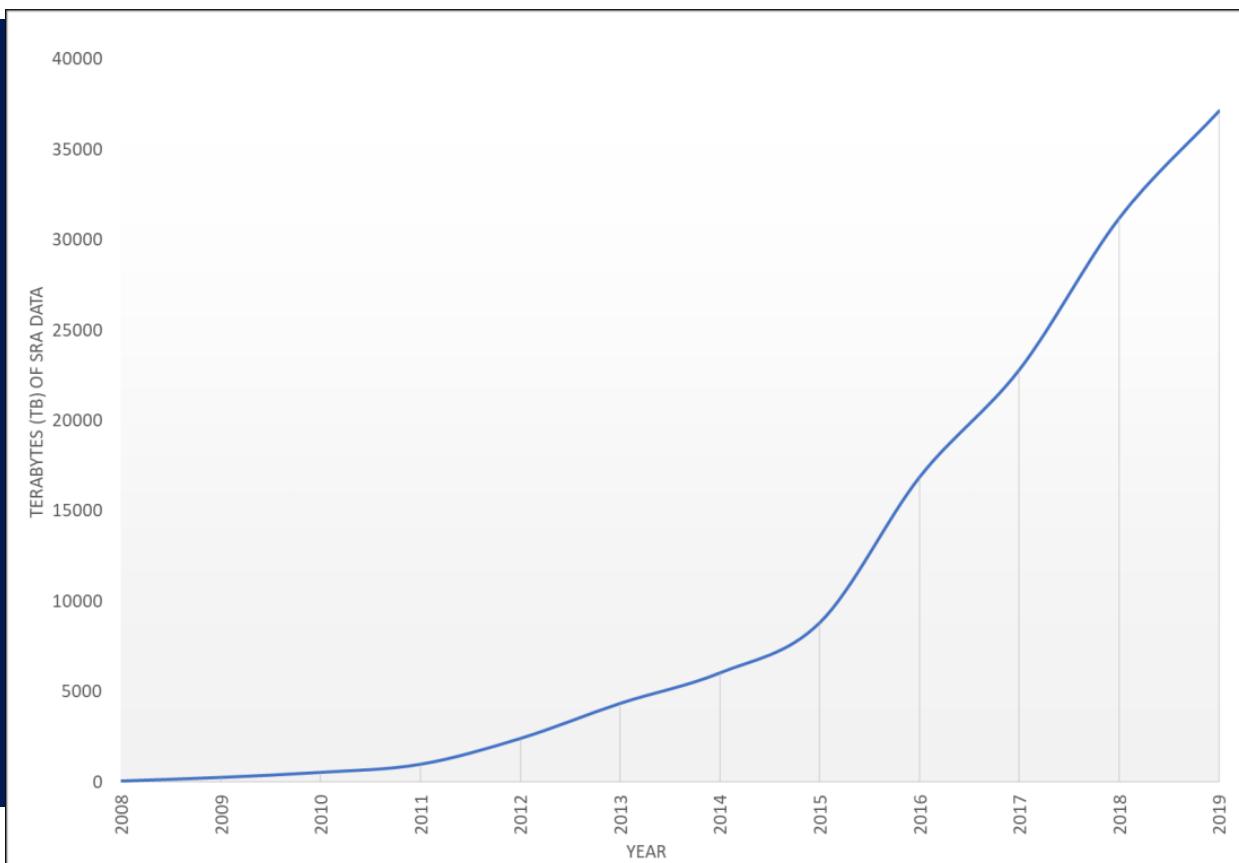
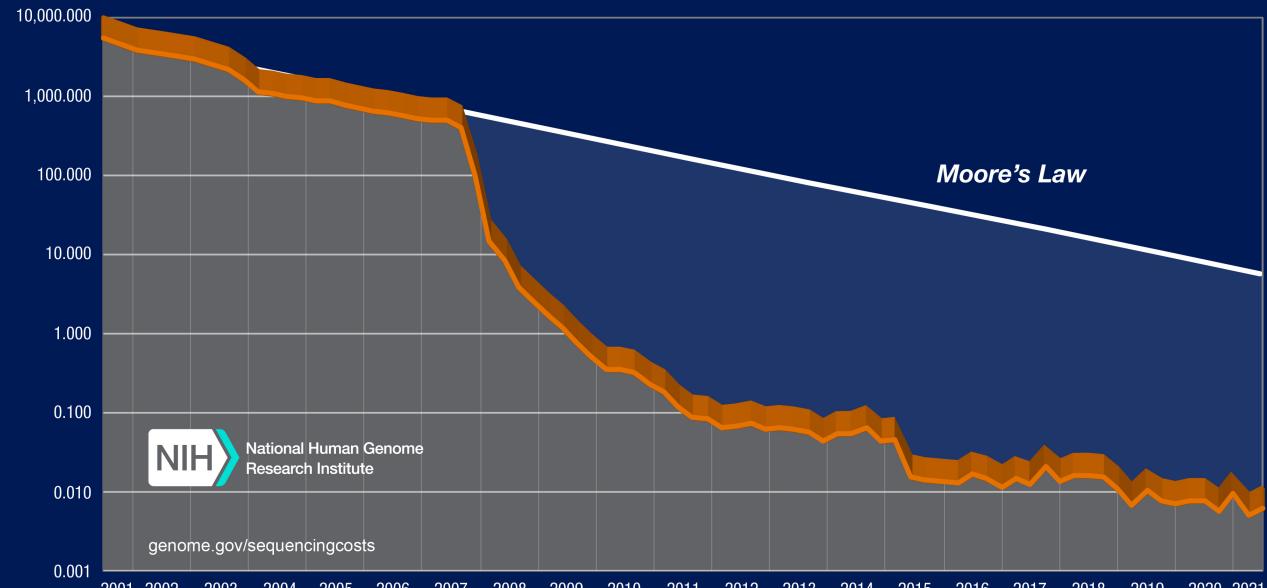
What is genomics anyway?

# Sanger sequencing



- A few hundred to ~1000 bp
- Low throughput
- Accurate

### **Cost per Raw Megabase of DNA Sequence**



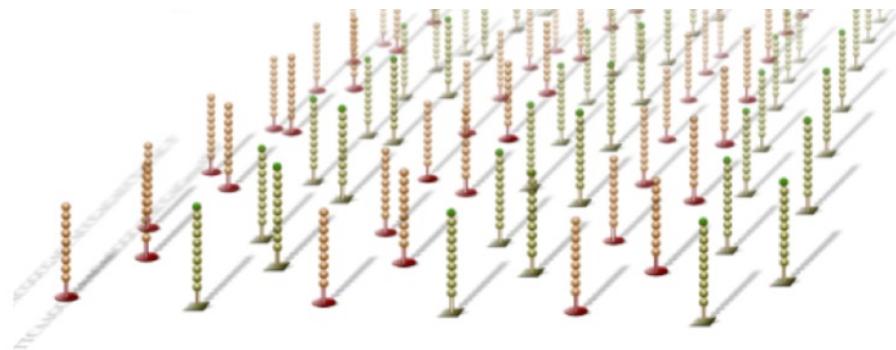
**Data in NCBI**

# Massively parallel sequencing

Illumina flow cell



- Aka
  - Next generation sequencing
  - Second generation sequencing
- Dominated by Illumina currently
- Short reads
  - 50-150 bp
- Very high throughput



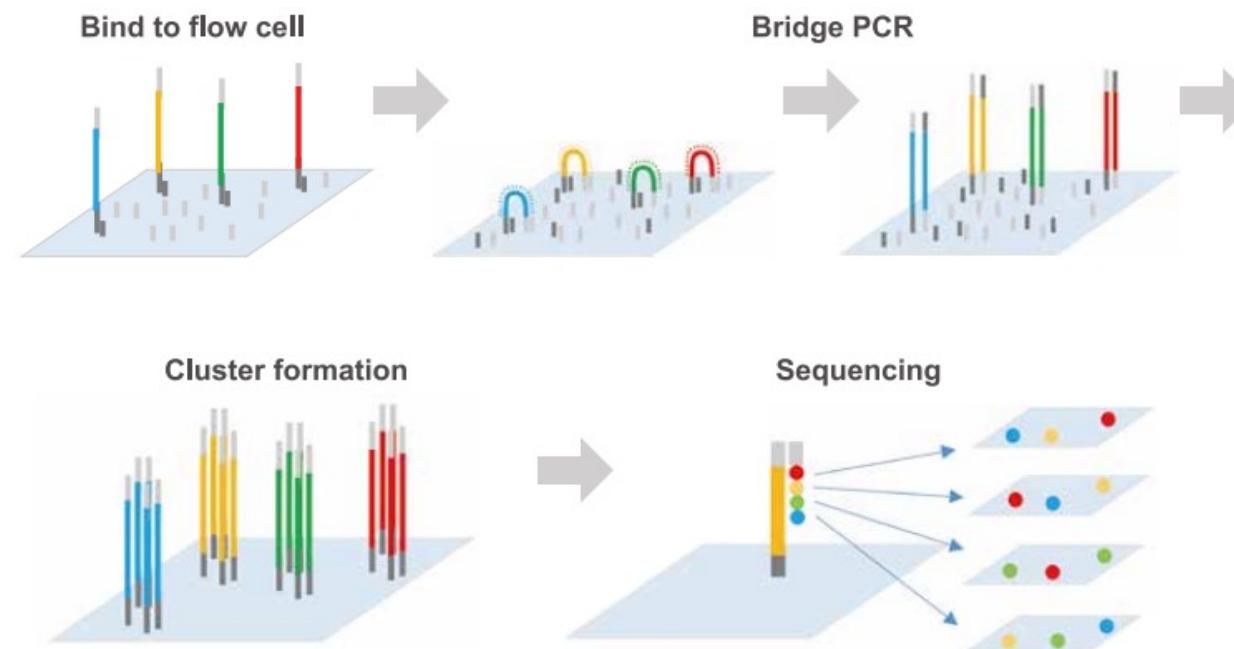
Fragment your DNA



Add adapters and amplify

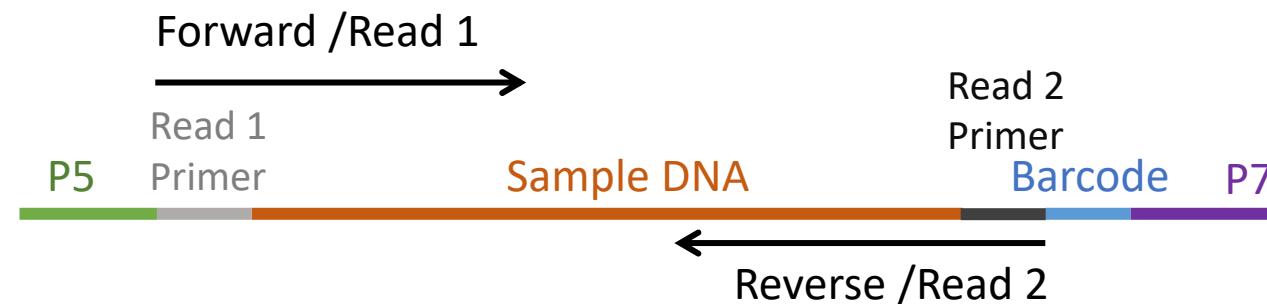


Sequencing by synthesis



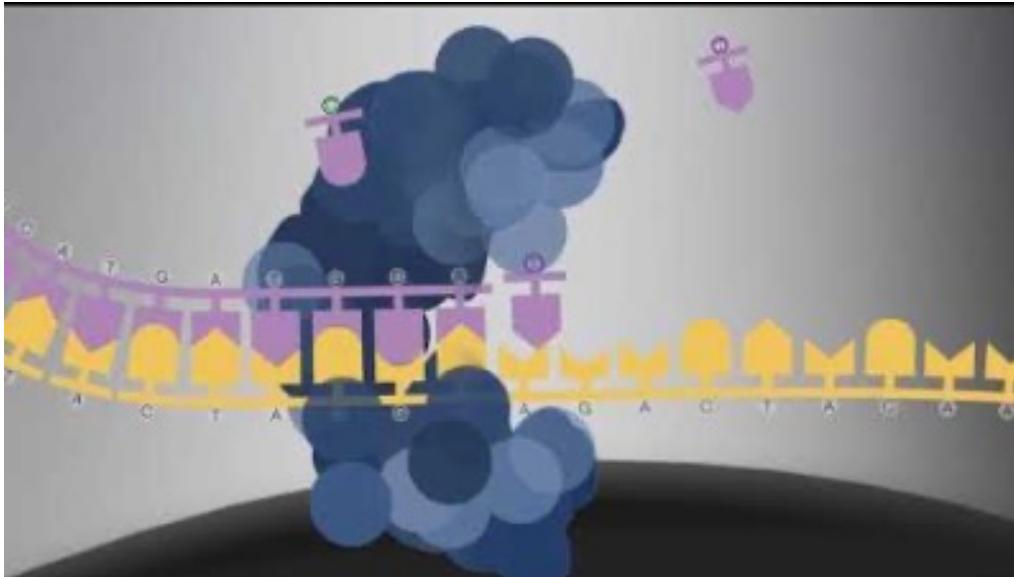
# Some important concepts to understand

- Library
  - The prepared sequences ready for sequencing that represent your sample
  - We modify the DNA to make a sequencing library
- Read
  - One sequence from a library



# Long read technologies

- PacBio
  - 10-50kb read length
  - Medium output
  - Read length limited by the longevity of the polymerase.
- Oxford Nanopore
  - 100kb or more
  - Low output
  - Read length limited by the length of the input DNA and decreasing yield with increasing DNA length



[https://www.youtube.com/watch?v=\\_ID8JyAbwEo](https://www.youtube.com/watch?v=_ID8JyAbwEo)



<https://www.youtube.com/watch?v=E9-Rm5AoZGw>

### Illumina

150-300bp read length  
High output

Population genetics,  
SNPs, RNAseq

### PacBio

10-25kb read length  
Medium output

Structural variation,  
genome assembly

### Nanopore

100kb or more  
Low output

Structural variation,  
genome assembly,  
in field analysis

What would we use each technology for?

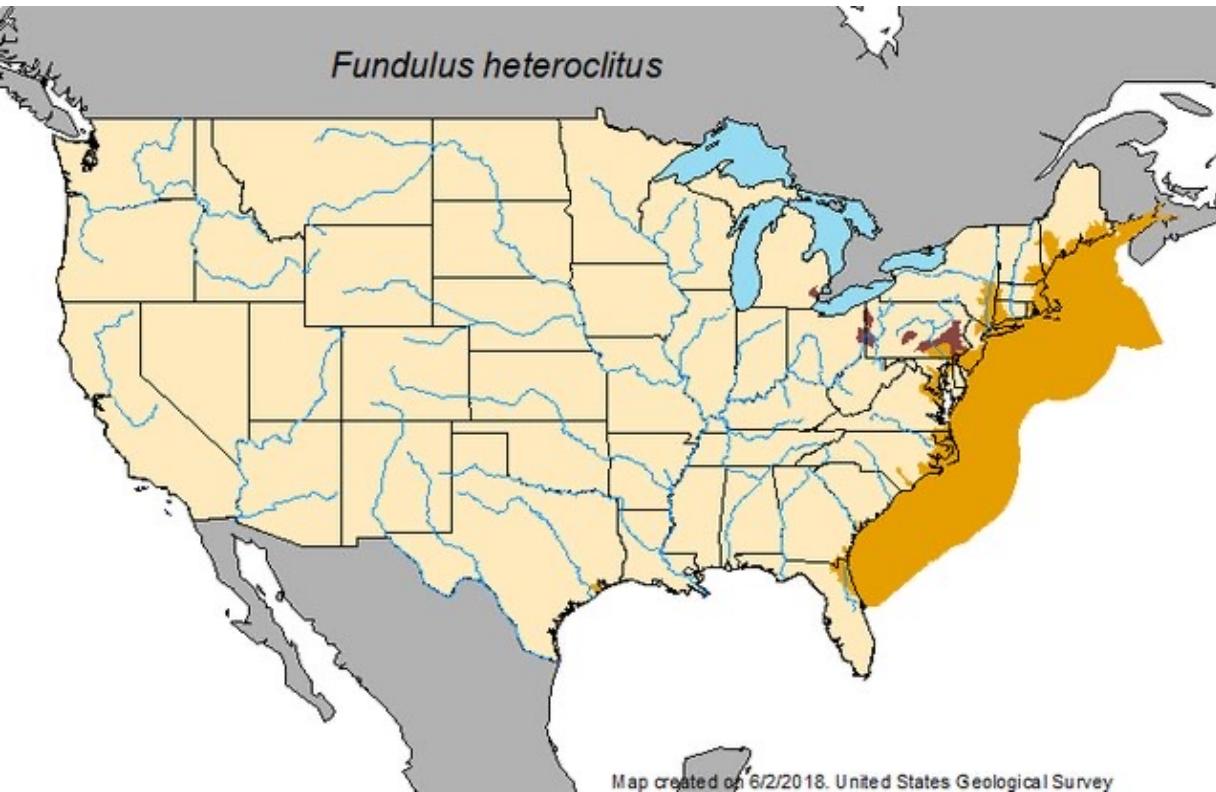
To study population genetics, we need to know allele frequencies → genomics!

Study system for the  
practical

# *Fundulus heteroclitus* study system



# *Fundulus heteroclitus* study system



Alcaraz-Hernandez and Garcia-Berthou

- Large population sizes
- Low movement
- Huge range of environments
  - Temperature
  - Salinity
  - Hypoxia
- Given the high population sizes and low migration:
  - what might we predict about populations of killifish from different regions of their distribution?

Received: 3 May 2018 | Accepted: 6 June 2018

DOI: 10.1111/gcb.14386

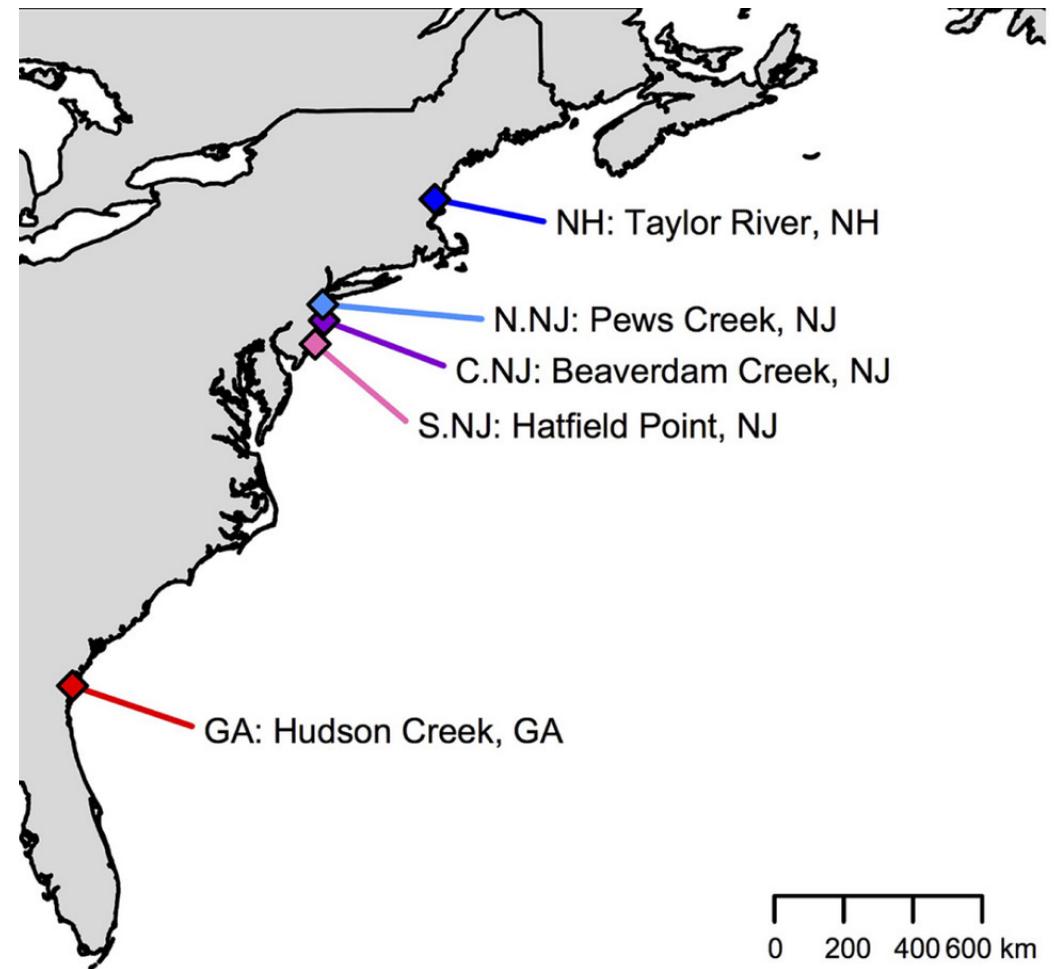
**PRIMARY RESEARCH ARTICLE**

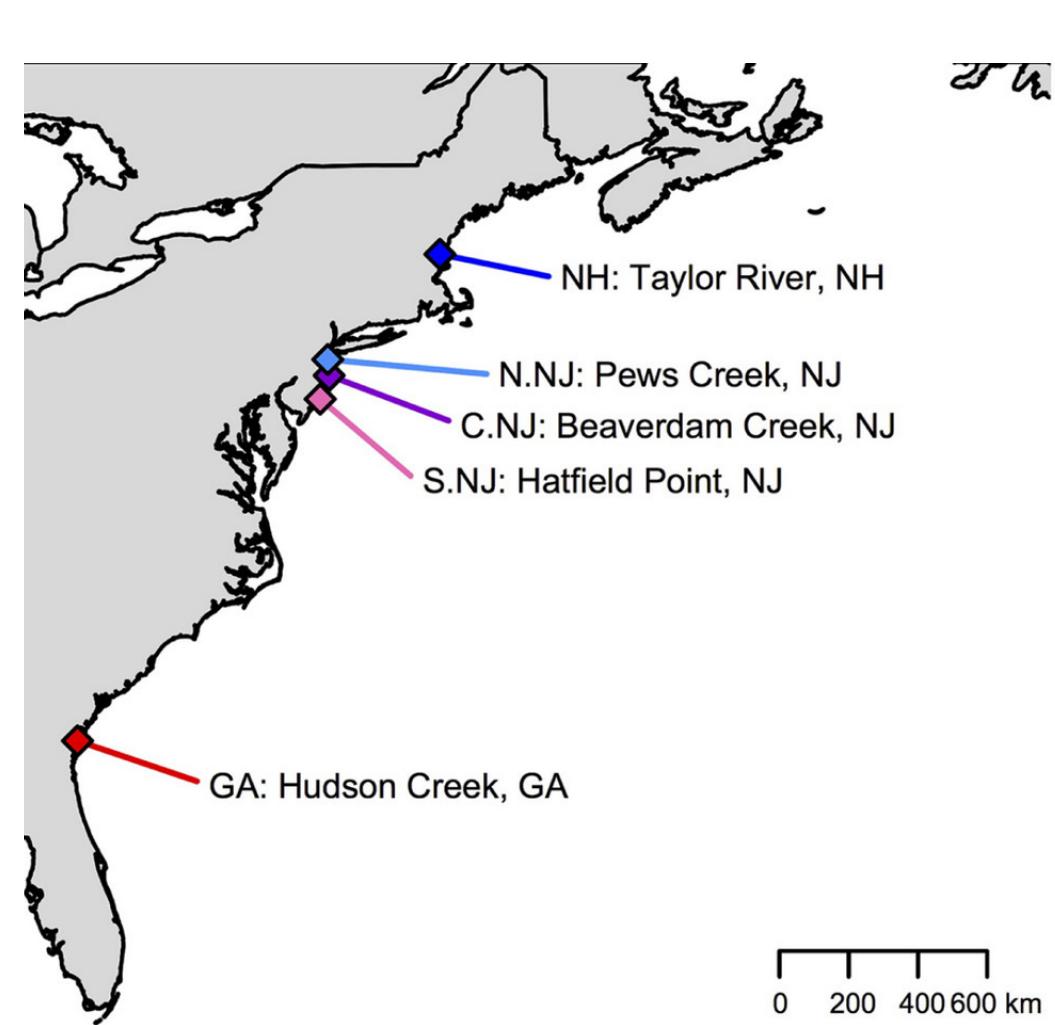
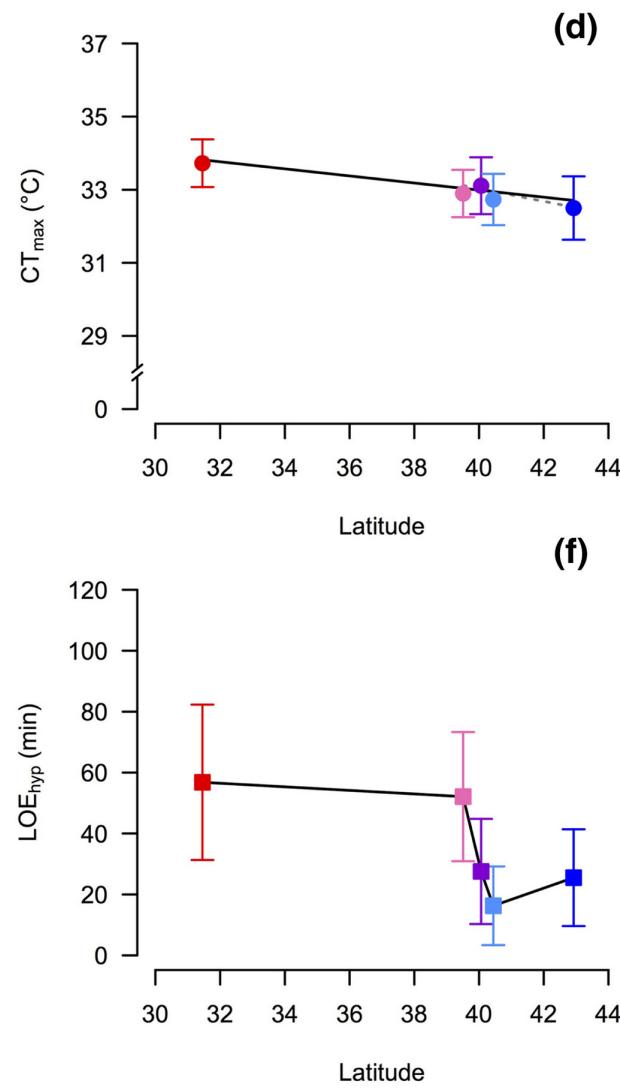
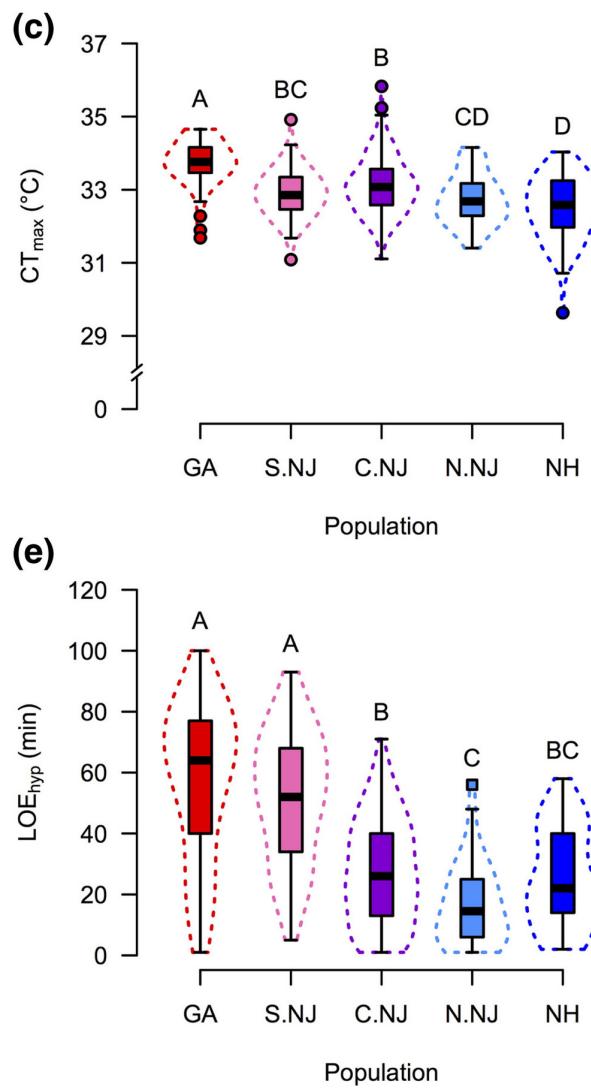
WILEY  Global Change Biology

## Tolerance traits related to climate change resilience are independent and polygenic

Timothy M. Healy<sup>1</sup>  | Reid S. Brennan<sup>2</sup>  | Andrew Whitehead<sup>2</sup>  |

Patricia M. Schulte<sup>1</sup> 

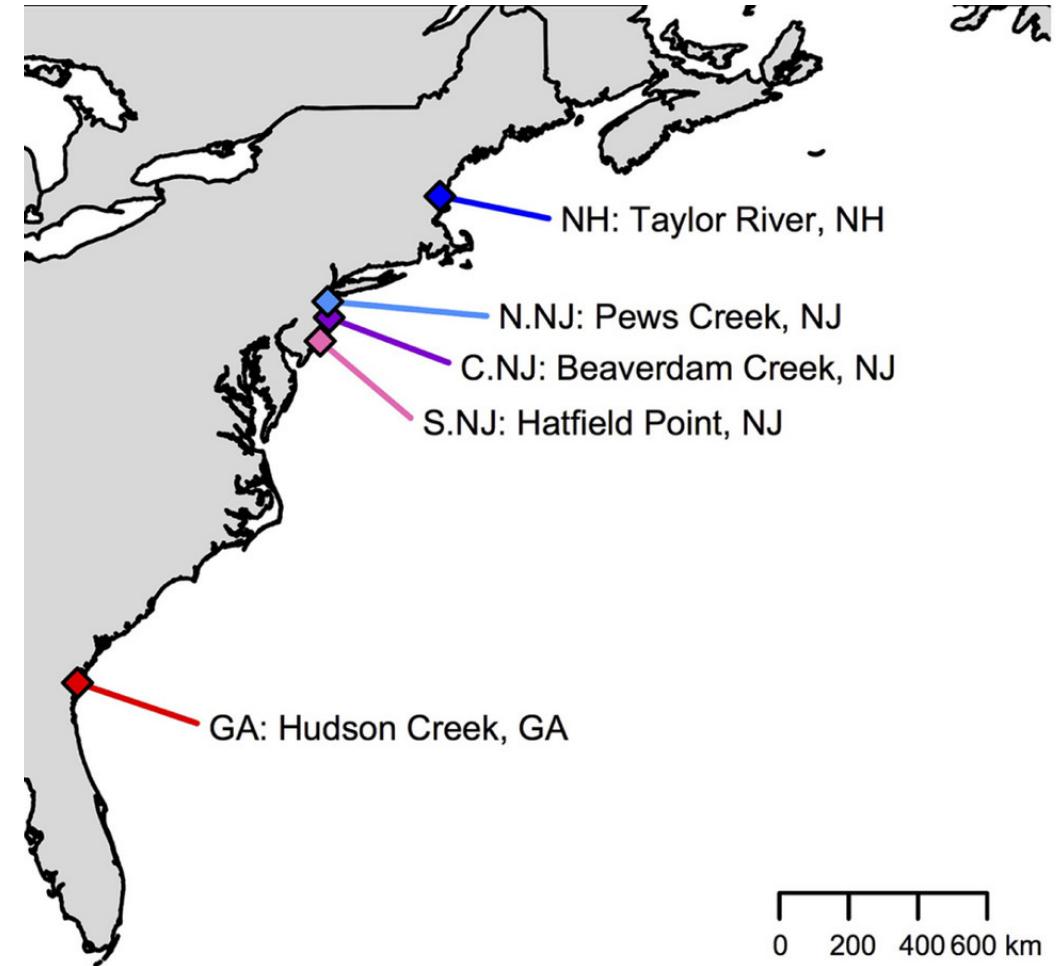




# Dataset

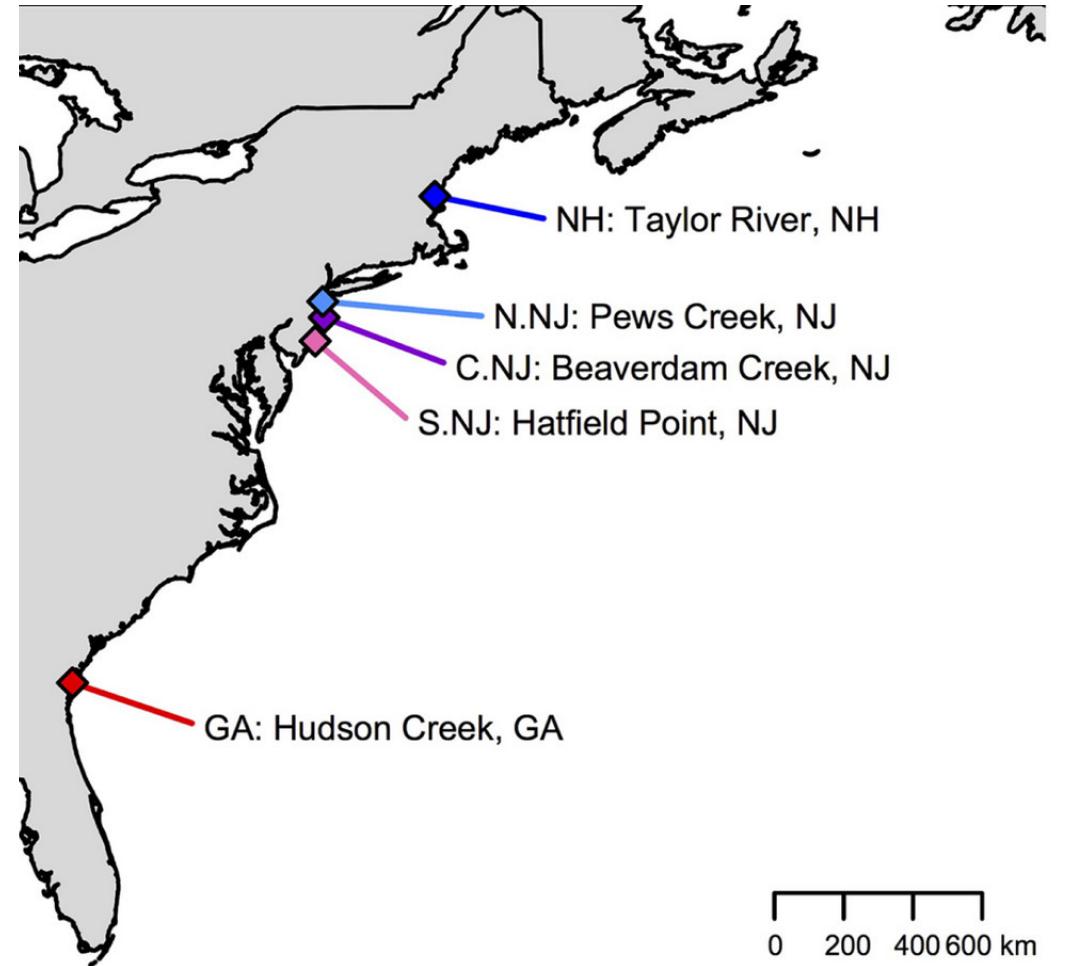
Population Location	Population ID	N
GA	GA	25
S.NJ	HP	46
C.NJ	BC	50
N.NJ	PC	44
NH	TR	45

RAD-seq



# Dataset

- Population structure
- Genetic diversity
- Selection



# Reduced representation

- In an ideal world...
  - Sequence full genomes of all individuals from many populations
  - Expensive, sometimes unnecessary, limited genomic resources
- RADseq, GBS, sequence capture, ddRAD, 2b-RAD...
- Useful for population structure and similar
- But... misses lots of the genome
- Cut anywhere you see the restriction site:
  - CCTGCAGG

---

## Methods

### Rapid and cost-effective polymorphism identification and genotyping using restriction site associated DNA (RAD) markers

Michael R. Miller,<sup>1</sup> Joseph P. Dunham,<sup>2</sup> Angel Amores,<sup>3</sup> William A. Cresko,<sup>2</sup> and Eric A. Johnson<sup>1,4</sup>

<sup>1</sup>Institute for Molecular Biology, University of Oregon, Eugene, Oregon 97403, USA; <sup>2</sup>Center for Ecology & Evolutionary Biology, University of Oregon, Eugene, Oregon 97403, USA; <sup>3</sup>Institute of Neuroscience, University of Oregon, Eugene, Oregon 97403, USA

2007

# Reduced representation

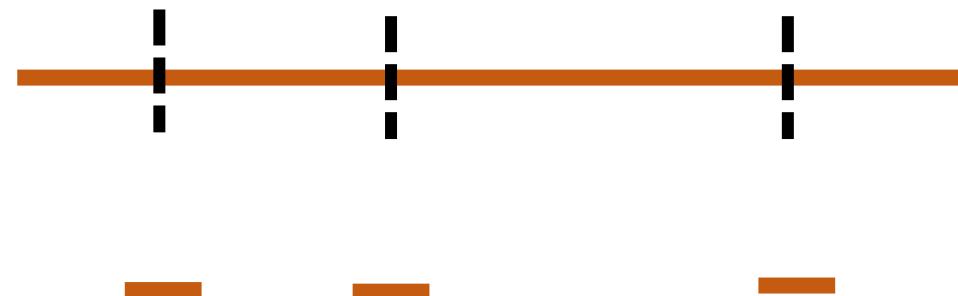
- In an ideal world...
  - Sequence full genomes of all individuals from many populations
  - Expensive, sometimes unnecessary, limited genomic resources
- RADseq, GBS, sequence capture, ddRAD, 2b-RAD...
- Useful for population structure and similar
- But... misses lots of the genome
- Cut anywhere you see the restriction site:
  - CCTGCAGG

5' CCTGCAGG 3'  
3' GGACGTCC 3'

Ligate  
adapters  
for  
sequencing

5' CCTGCA|GG 3'  
3' GG|ACGTCC 3'

GG.....  
ACGTCC.....



# How much of the genome does RADseq cover?

- How big is the genome?
  - 1.5 billion bases
- How often do we find a cut site?
  - Assume bases are randomly distributed
  - $0.25^8$
- $1.5e+09 * (0.25^8) = 22,888$

# Unix, cloudb, and R

# We have lots of data

- Normal to get back 100's of GB of data.
  - You can't have this on your computer!
- Computer clusters are needed
  - These require command line, unix, etc.

# Why unix/command line?

- Most software is written for unix
- Computing clusters use unix
- Good at working with large data

# Why R?

- Very popular in biology (good community)
- Good for statistics
- Good for plotting

# Why Both?

- Reproducible!!

# How to learn?

- By trying
- Workshops/courses
- Google
  - Stackoverflow; seqanswers; others
- ChatGPT

# “the cloud”

- Aka, a computer somewhere else



- We will use Cloudlab → hosted by CAU

# R Markdown

- You want to have a record of your analyses → markdown
- plain text formatting syntax designed to be converted to HTML and others
- R markdown is an extended markdown, that can embed and run R code, generate plots, etc.
- Why?
  - Reproducibility

You will turn in a R markdown file with your final project

metadata

text

code chunks

Source Visual

Outline

```
1 ---  
2 title: "Example file"  
3 output: html_document  
4 ---  
5  
6 This is an R Markdown document. This is just plain text.  
7  
8  
9 Below, this is a code block that will be executed  
10 ````{r}  
11  
12 plot(seq(1,10,1), seq(10,1,-1))  
13  
14 ````
```

## Example file

This is an R Markdown document. This is just plain text.

Below, this is a code block that will be executed

```
plot(seq(1,10,1), seq(10,1,-1))
```

A scatter plot with the x-axis labeled 'seq(1, 10, 1)' and the y-axis labeled 'seq(10, 1, -1)'. The plot shows a series of open circles forming a straight line with a negative slope, starting at approximately (1, 10) and ending at approximately (10, 1).

x	y
1	10
2	9
3	8
4	7
5	6
6	5
7	4
8	3
9	2
10	1

# Today's tutorial

- Go to cloudlab together.
- Start exploring Rmarkdown
- Unix for biologists