# Applied Comparative Genomics

## Michael Schatz

August 29, 2022

Lecture 1: Course Overview

# Welcome!

*The primary goal of the course is for students to be grounded in theory and leave the course empowered to conduct independent genomic analyses.*

- We will study the leading computational and quantitative approaches for comparing and analyzing genomes starting from raw sequencing data.
- The course will focus on human genomics and human medical applications, but the techniques will be broadly applicable across the tree of life.
- The topics will include genome assembly & comparative genomics, variant identification & analysis, gene expression & regulation, personal genome analysis, and cancer genomics.

*Course Webpage:* https://github.com/schatzlab/appliedgenomics2022
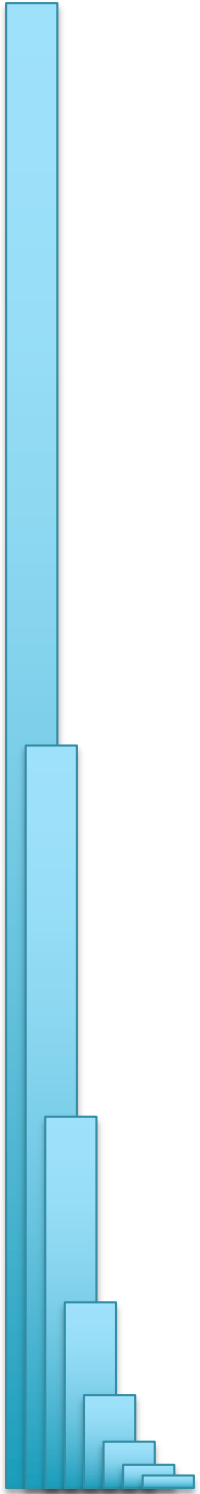
*Course Discussions:* http://piazza.com

*Class Hours:* Mon + Wed @ 1:30p – 2:45p, Gilman 17

*Schatz Office Hours:* TBD and by appointment

*Ni Office Hours:* TBD and by appointment

Please try Piazza first!

# TA: Bohan Ni

# Prerequisites and Resources

*Prerequisites*
- No formal course requirements
- Access to an Apple or Linux Machine, or Install VirtualBox
- Familiarity with the Unix command line for exercises
  - bash, ls, grep, sed, + install published genomics tools
- Familiarity with a major programming language for project
  - C/C++, Java, R, Perl, Python

*Primary Texts*
- None! We will be studying primary research papers

*Other Resources:*
- Google, SEQanswers, Biostars, StackOverflow

- Applied Computational Genomics Course at UU: Spring 2018/2020
- https://github.com/quinlan-lab/applied-computational-genomics

- Ben Langmead's teaching materials:
- http://www.langmead-lab.org/teaching-materials/

# Grading Policies

**Assessments:**

- 5 Assignments:          30%       Due at 11:59pm a week later
  ***Practice using the tools we are discussing***

- 1 Exam:                      30%       Take Home (Tentatively Nov 2)
  ***Assess your performance, focusing on the methods***

- 1 Class Project:          40%       Presented last week of class
  ***Significant project developing a novel analysis/method***

- In-class Participation:   Not graded, but there to help you!

**Policies:**

- Scores assigned relative to the highest points awarded
- Automated testing and grading of assignments
- ***Late Days:***
  - A total of 96 hours (24 x 4) can be used to extend the deadline for assignments, but not the class project, without any penalty; after that time assignments will not be accepted

# Course Webpage



https://github.com/schatzlab/appliedgenomics2022

# Piazza



https://piazza.com/class/l7dg3c82ftw1d/
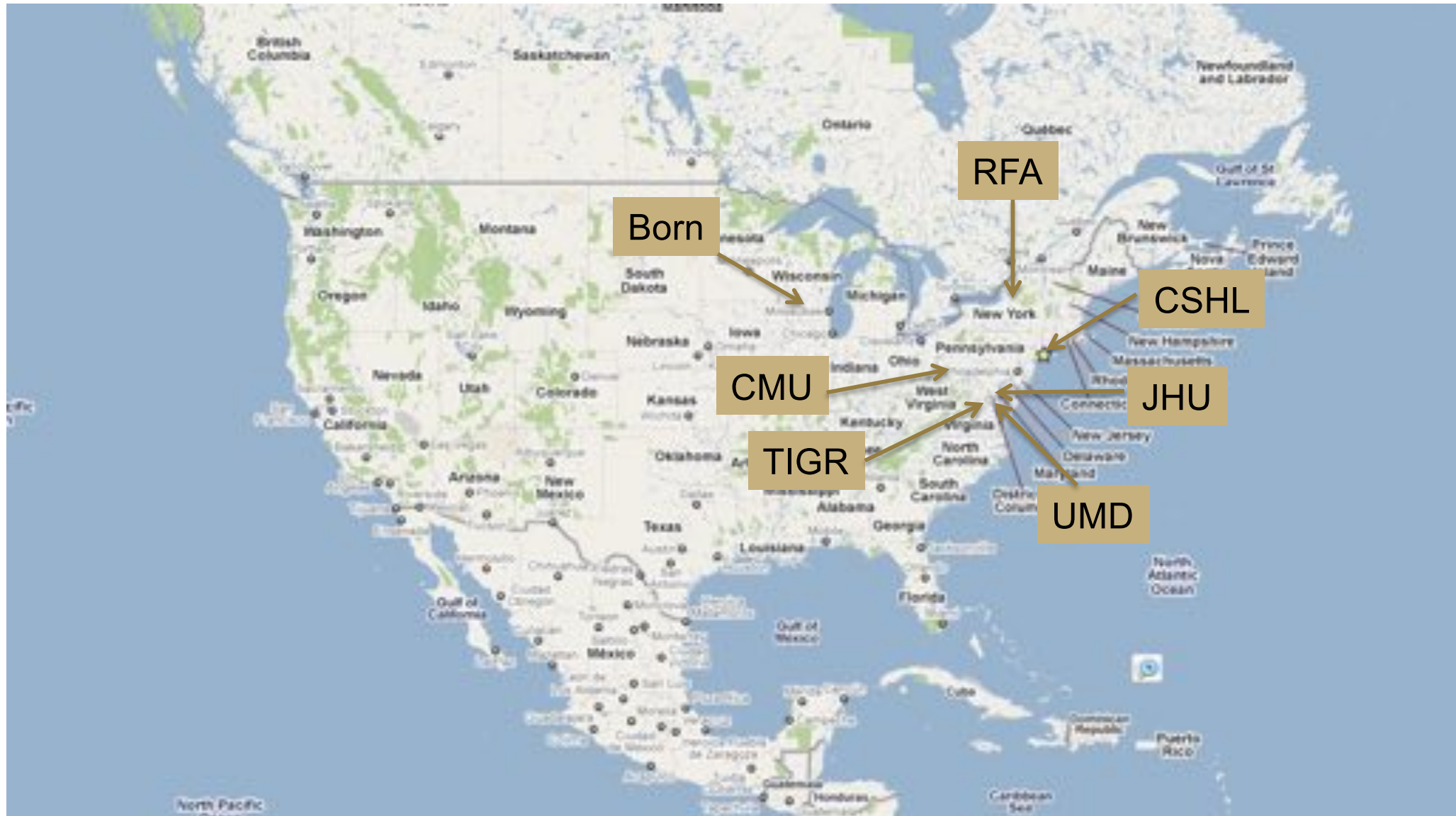
# GradeScope



https://www.gradescope.com/
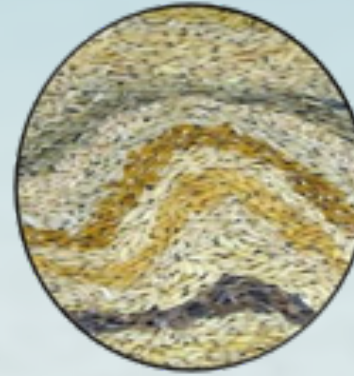Entry Code: **J37JKW**

# A Little About Me

# Schatzlab Overview

**Human Genetics**

Role of mutations in disease

Nurk et al. (2022)
Aganezov *et al.* (2020)

**Agricultural Genomics**
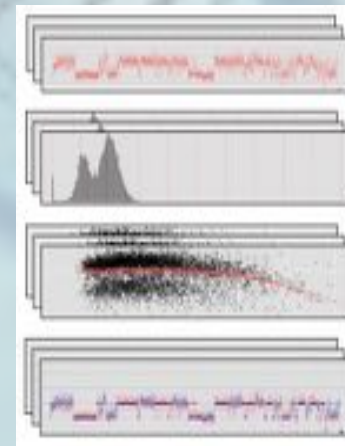
Genomes & Transcriptomes

Alonge *et al.* (2020)
Soyk *et al.* (2019)

**Algorithmics & Systems Research**
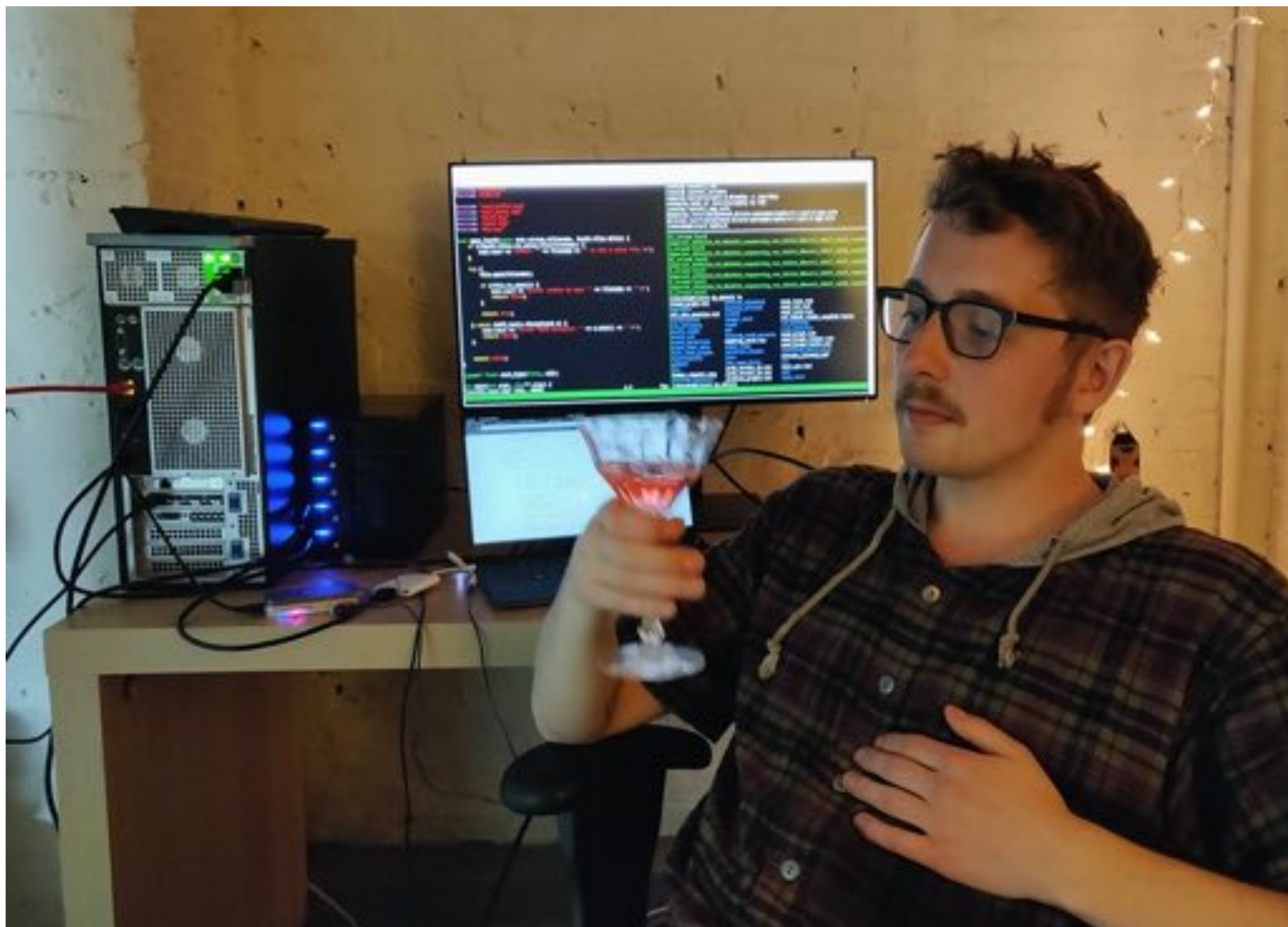
Ultra-large scale biocomputing

Schatz et al. (2022)
Kirsche *et al.* (2020)

**Biotechnology Development**

Single Cell + Single Molecule Sequencing

Kovaka *et al.* (2020)
Sedlazeck *et al.* (2018)

# nature biotechnology

Explore content ⌄     Journal information ⌄     Publish with us ⌄     Subscribe

Sign up for alerts 🔔     RSS feed

nature > nature biotechnology > articles > article

Article | Published: 30 November 2020

# Targeted nanopore sequencing by real-time mapping of raw electrical signal with UNCALLED

Sam Kovaka ✉, Yunfan Fan, Bohan Ni, Winston Timp & Michael C. Schatz

You have full access to this article via **Johns Hopkins Libraries**

Download PDF  ⤓

| Sections | Figures | References |

Abstract

Main

Results

Discussion

Methods

Data availability

Code availability
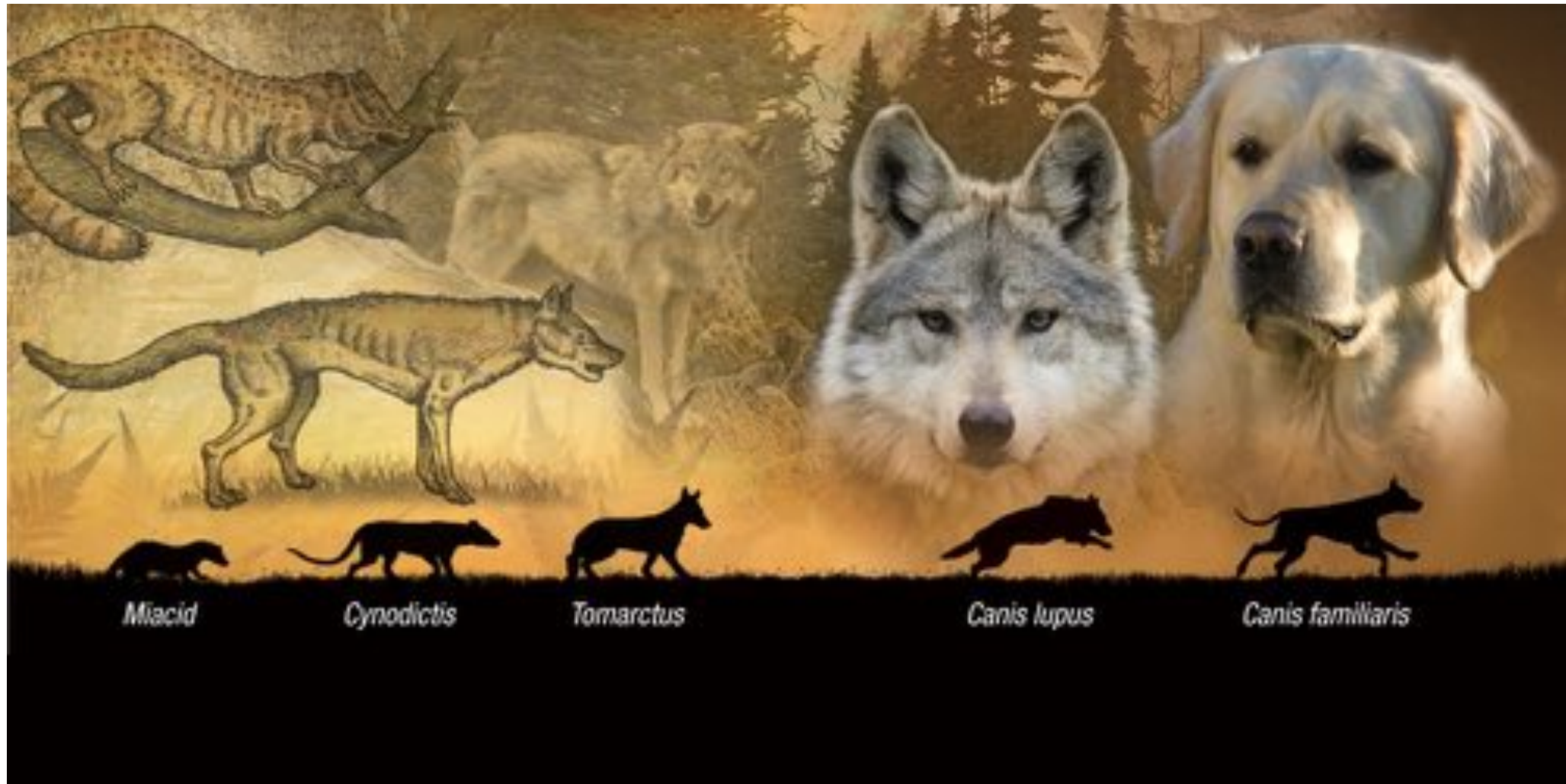
References

Acknowledgements

Author information

Ethics declarations

## Abstract

Conventional targeted sequencing methods eliminate many of the benefits of nanopore sequencing, such as the ability to accurately detect structural variants or epigenetic modifications. The ReadUntil method allows nanopore devices to selectively eject reads from pores in real time, which could enable purely computational targeted sequencing. However, this requires rapid identification of on-target reads while most mapping methods require computationally intensive basecalling. We present UNCALLED (https://github.com/skovaka/UNCALLED), an open source mapper that rapidly matches streaming of nanopore current signals to a reference sequence. UNCALLED probabilistically

# Earliest Genomics

Any Guesses?

# Earliest Genomics



Miacid  Cynodictis  Tomarctus  Canis lupus  Canis familiaris

15,000 to 35,000 YBP

# Earliest Genomics



~1,000 to 10,000 YBP

# Earliest Genomics



Teosinte      ~6,000 to 10,000 YBP      Modern Corn
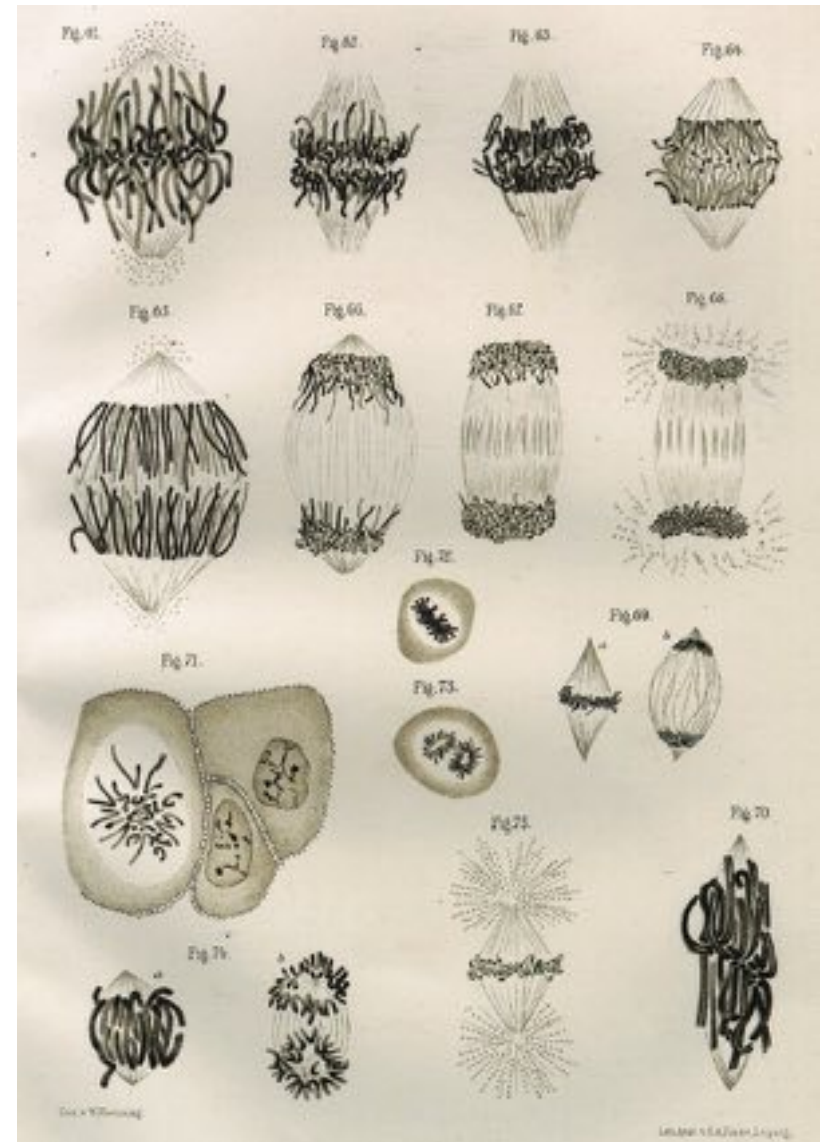
# Angiosperms (Flowering Plants)



~130 Ma

# Discovery of Chromosomes

By the mid-1800s, microscopes were powerful enough to observe the presence of unusual structures called "chromosomes" that seemed to play an important role during cell division.

It was only possible to see the chromosomes unless appropriate stains were used

"Chromosome" comes from the Greek words meaning "color body"

Today, we have much higher resolution microscopes, and a much richer varieties of dies and dying techniques so that we can visualize particular sequence elements.

When you see something unexpected that you think might be interesting, give it a name



*Drawing of mitosis by Walther Flemming.*
Flemming, W. Zellsubstanz, Kern und Zelltheilung (F. C. W. Vogel, Leipzig, 1882).

# The "first" quantitative biologist

Any Guesses?

# Laws of Inheritance





http://en.wikipedia.org/wiki/Experiments_on_Plant_Hybridization

Observations of 29,000 pea plants and 7 traits

| | | | | in Verhältniss gestellt: | | |
|---|---|---|---|---|---|---|
| Generation | A | Aa | a | A : | Aa : | a |
| 1 | 1 | 2 | 1 | 1 : | 2 : | 1 |
| 2 | 6 | 4 | 6 | 3 : | 2 : | 3 |
| 3 | 28 | 8 | 28 | 7 : | 2 : | 7 |
| 4 | 120 | 16 | 120 | 15 : | 2 : | 15 |
| 5 | 496 | 32 | 496 | 31 : | 2 : | 31 |
| n | | | | $2^n-1$ : | 2 : | $2^n-1$ |

*Versuche über Pflanzen-Hybriden. Verh. Naturforsch (Experiments in Plant Hybridization)*
Mendel, G. (1866). Ver. Brünn 4: 3–47 (in English in 1901, J. R. Hortic. Soc. 26: 1–32).

# The first genetic map

Mendel's Second Law (The Law of Independent Assortment) states alleles of one gene sort into gametes independently of the alleles of another gene: **Pr(smooth/wrinkle) is independent of Pr(yellow/green)**

Morgan and Sturtevant noticed that the probability of having one trait given another was **not** always 50/50– those traits are **genetically linked**



http://www.caltech.edu/news/first-genetic-linkage-map-38798

Sturtevant realized the probabilities of co-occurrences could be explained if those alleles were arranged on a linear fashion: traits that are most commonly observed together must be locates closest together



**The Linear Arrangement of Six Sex-Linked Factors in Drosophila as shown by their mode of Association**
Sturtevant, A. H. (1913) *Journal of Experimental Zoology*, 14: 43-59

# Jumping Genes



Previously, genes were considered to be stable entities arranged in an orderly linear pattern on chromosomes, like beads on a string

Careful breeding and cytogenetics revealed that some elements can move (cut-and-paste, DNA transposons) or copy itself (copy-and-paste, retrotransposons)





(Gregory, 2005, Nature Reviews Genetics)

(Much) later analysis revealed that nearly 50% of the human genome is composed of transposable elements, including LINE and SINE elements (long/short interspersed nuclear elements) which can occur in 100k to 1M copies

*"The genome is a graveyard of ancient transposons"*

**The origin and behavior of mutable loci in maize.**
McClintock, B. (1950) *PNAS*. 36(6):344–355.
Nobel Prize in Physiology or Medicine in 1983

# Discovery of the Double Helix



It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material. Full details of the structure, including the con-

*Molecular structure of nucleic acids; a structure for deoxyribose nucleic acid*
Watson JD, Crick FH (1953). *Nature* 171: 737–738.
Nobel Prize in Physiology or Medicine in 1962

# Central Dogma of Molecular Biology

"Once 'information' has passed into protein it cannot get out again. In more detail, the transfer of information **from nucleic acid to nucleic acid, or from nucleic acid to protein may be possible**, but transfer from protein to protein, or from protein to nucleic acid is impossible. Information means here the precise determination of sequence, either of bases in the nucleic acid or of amino acid residues in the protein"



*On Protein Synthesis*
Crick, F.H.C. (1958). Symposia of the Society for Experimental Biology pp. 138–163.

# One Genome, Many Cell Types

Each cell of your body contains an exact copy of your 3 billion base pair genome.

Your body has a few hundred (thousands?) major cell types, largely defined by the gene expression patterns

cell

nucleus

chromosome

gene

DNA

Adapted from National Human Genome Research Institute

# Sequencing Capacity



**Big Data: Astronomical or Genomical?**
*Stephens, Z, et al. (2015) PLOS Biology DOI: 10.1371/journal.pbio.1002195*

# Unsolved Questions in Biology

- What is your genome sequence?
- •
- •
- •
- •
- •
- •
- •
- •
- •
- •
- •
- *Plus thousands and thousands more*

The instruments provide the data, but none of the answers to any of these questions.

*What software and systems will?*

*And who will create them?*

# Who is a Data Scientist?



http://en.wikipedia.org/wiki/Data_science

# Comparative Genomics Technologies



**Results**
Domain Knowledge

**Machine Learning**
classification, modeling, visualization & data Integration

**Scalable Algorithms**
Streaming, Sampling, Indexing, Parallel

**Compute Systems**
CPU, GPU, Distributed, Clouds, Workflows

**IO Systems**
Hardrives, Networking, Databases, Compression, LIMS

**Sensors & Metadata**
Sequencers, Microscopy, Imaging, Mass spec, Metadata & Ontologies

# Next Steps

1. Reflect on the magic and power of DNA ☺

2. Check out the course webpage

3. Register on Piazza

4. Get Ready for assignment 1

   1. Set up Linux, set up Docker
   2. Set up Dropbox for yourself!
   3. Get comfortable on the command line