# Biomedical Data Science (GersteinLab.org/courses/452)
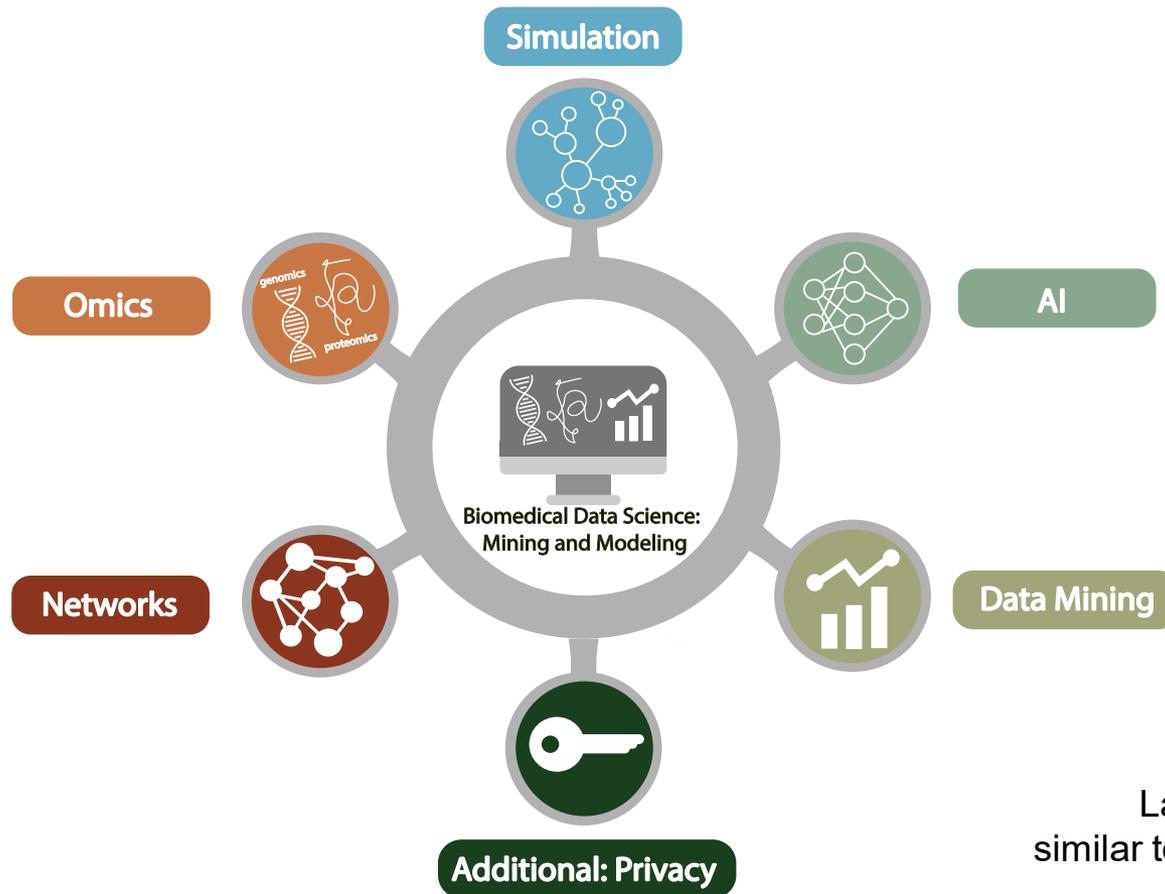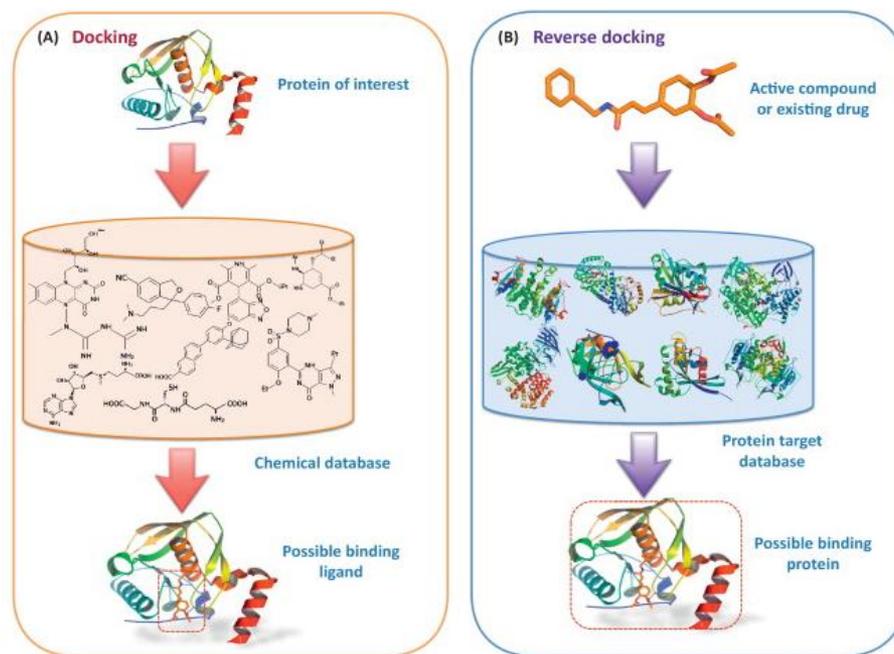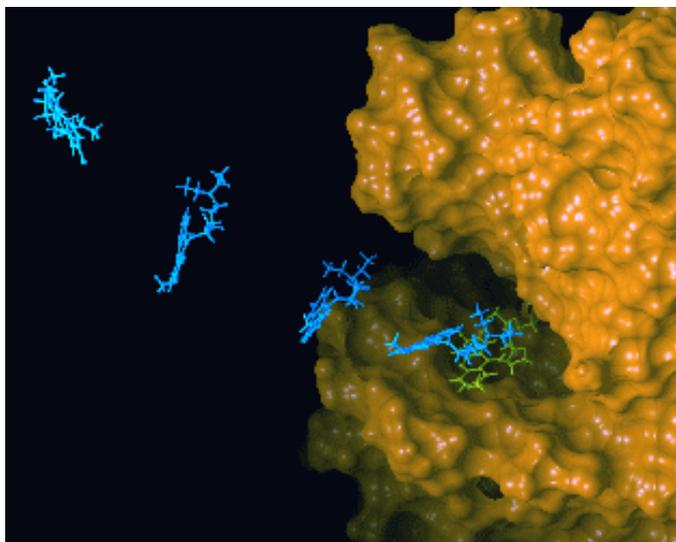# Introduction to Personal Genomes (23i2a)
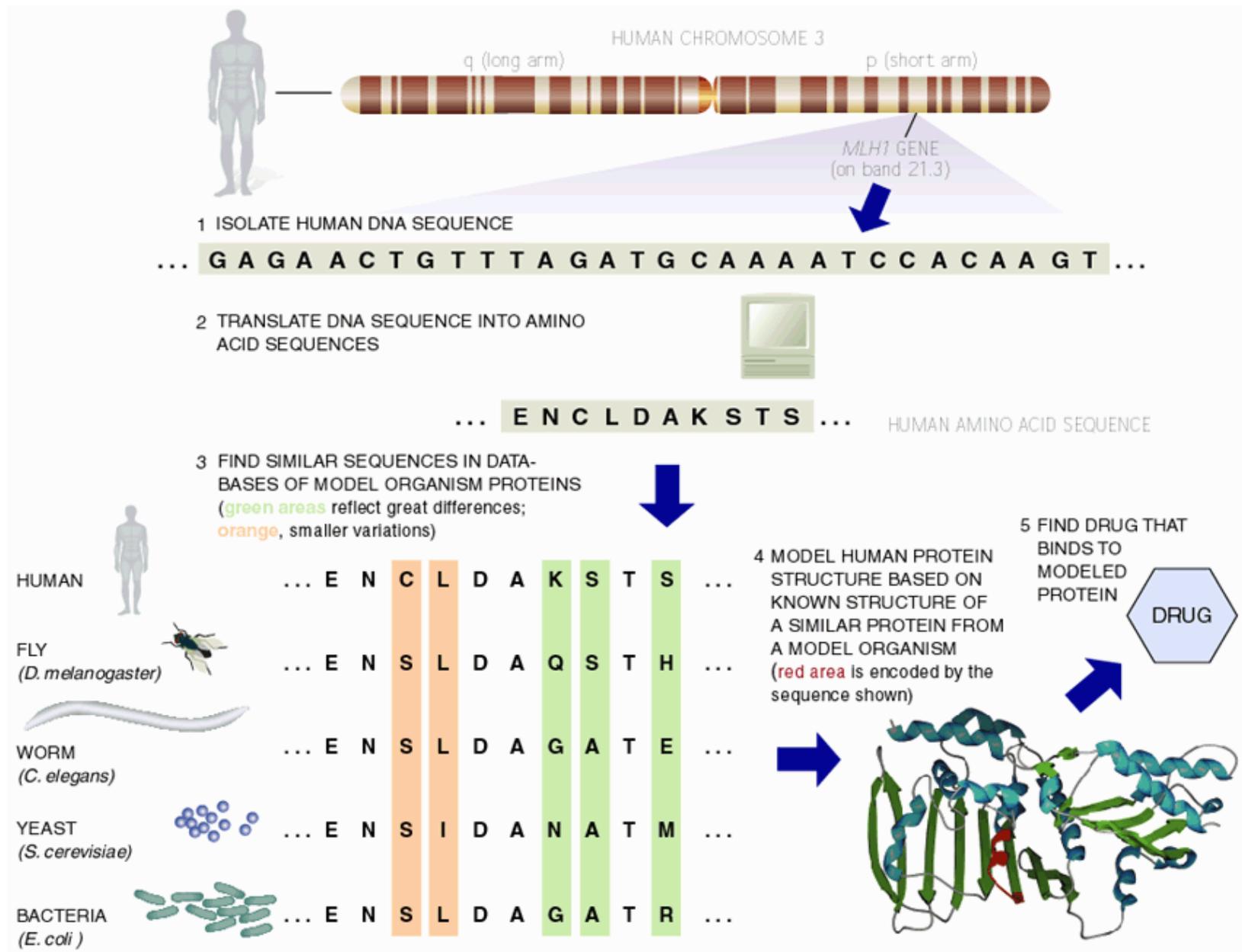


Mark Gerstein
Yale U.

Last edit in spring '23, similar to 22i2a & 2021's I2a [which has a video], with additions beginning at slide 10, describing Zimmerome history & assignment. Slight modification at slide 5 too.

# Major Application I:
## Designing Drugs from Structural Targets

- Understanding how structures bind other molecules
- Designing inhibitors using docking, structure modeling
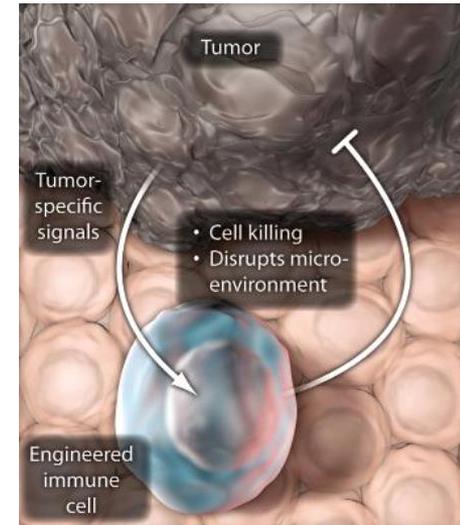- *In silico* screens of chemical and protein databases
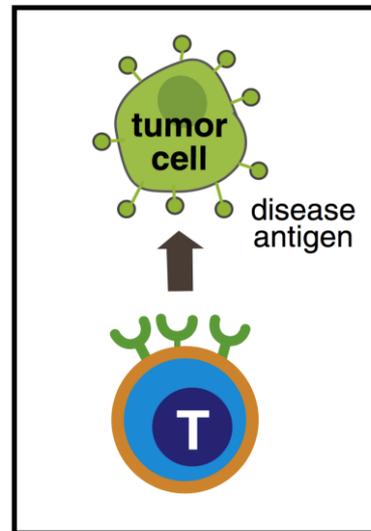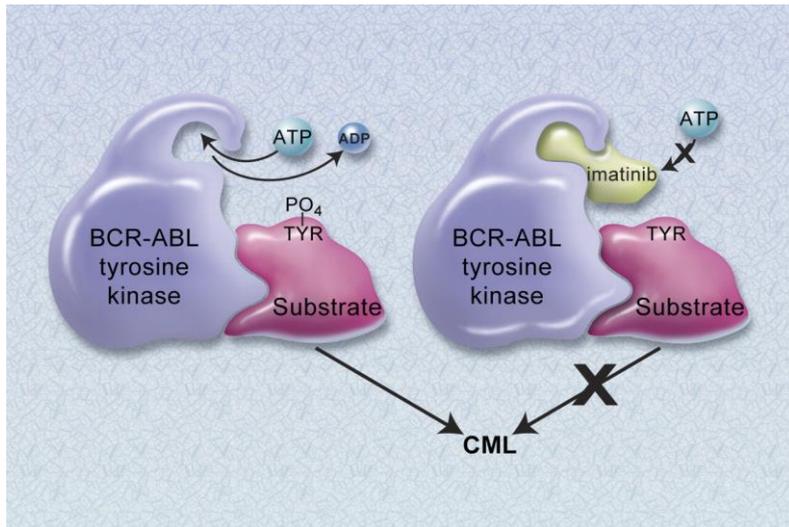
(From left to right, figures adapted from Olsen Group Docking Page at Scripps, Zheng et al. Trends in Pharmacological Sciences 2013)

# Major Application II: Finding Homologs, to Find Experimentally Tractable Gene Targets

# Major Application III:
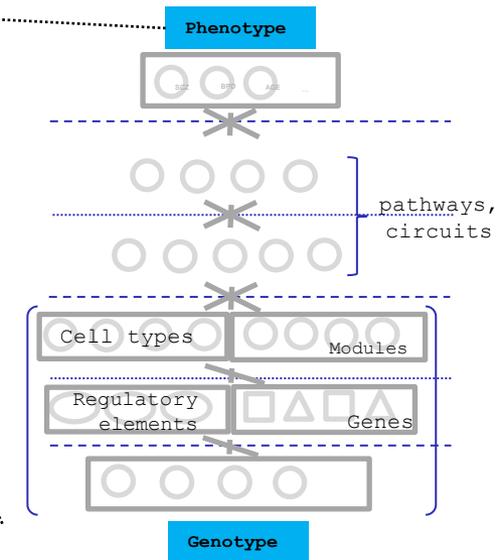# Customizing treatment in oncology

- Identifying disease causing mutations in individual patients

- Designing targeted therapeutics
  - e.g. BCR-abl and Gleevec
  - Cancer immunotherapies targeting neoantigens



**(From left to right, figures adapted from Druker BJ. Blood 2008 and the Lim Lab at UCSF)**

# Major Application IV:
# Finding molecular mechanisms & drug targets
# for diseases we know little about (Neuro-psychiatic Diseases)

| Disease | Heritability* | Molecular **Mechanisms** |
|---|---|---|
| **Schizophrenia** | **81%** | **-** |
| **Bipolar disorder** | 70% | **-** |
| **Alzheimer's disease** | 58 - 79% | Apolipoprotein E (APOE), Tau |
| **Hypertension** | 30% | Renin–angiotensin–aldosterone |
| **Heart disease** | 34-53% | Atherosclerosis, VCAM-1 |
| **Stroke** | 32% | Reactive oxygen species (ROS), Ischemia |
| **Type-2 diabetes** | 26% | Insulin resistance |
| **Breast Cancer** | 25-56% | BRCA, PTEN |



Many psychiatric conditions are highly heritable
      Schizophrenia: up to 80%
 But we don't understand basic molecular mechanisms underpinning this association
      (in contrast to many other diseases such as cancer & heart disease)
Moreover, current models substantially underestimate heritability using genetic data
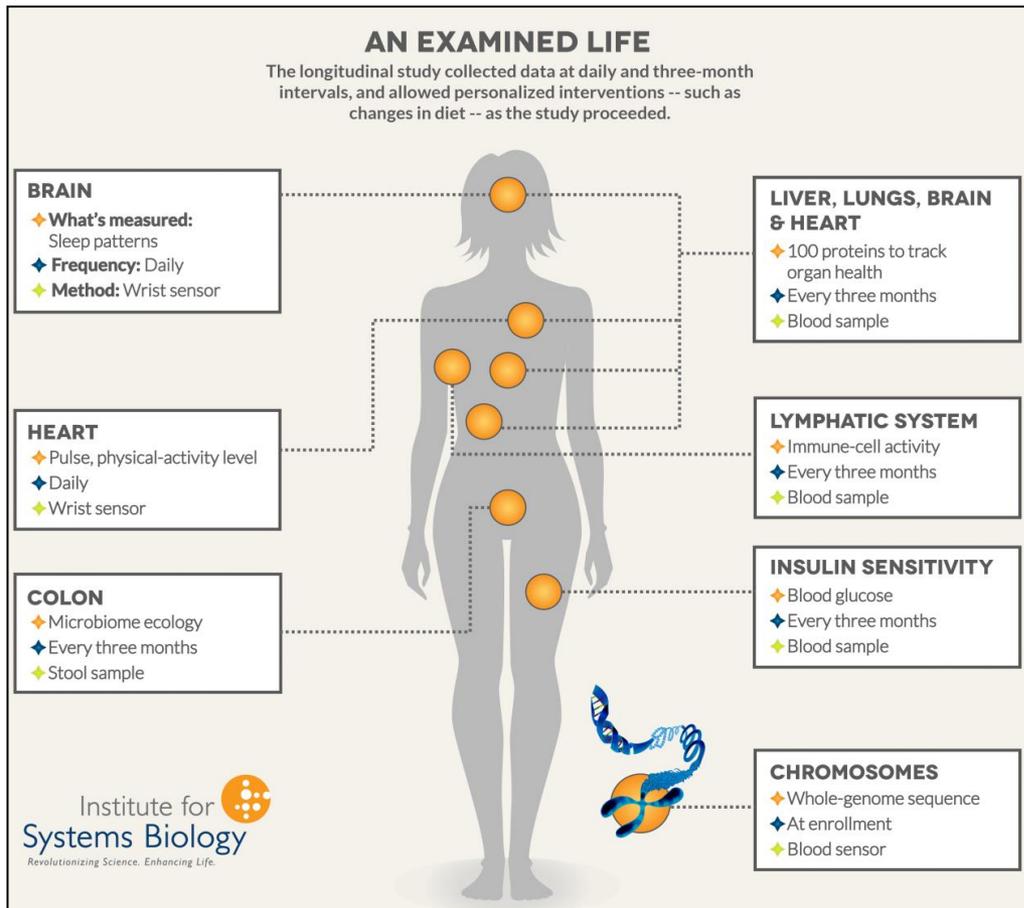      Schizophrenia : ~25%
Thus, interested in developing predictive models of psychiatric traits which:
      Use observations at intermediate (molecular levels) levels to inform latent
      structure.
      Use the predictive features of these "molecular endo phenotypes" to begin to
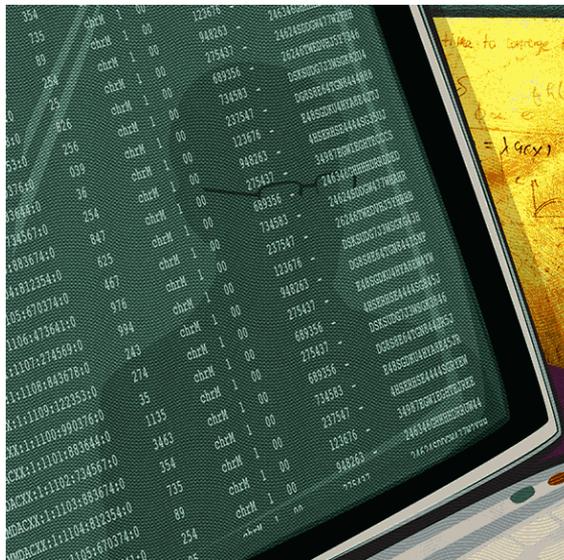      suggest actors involved in mechanism

# Major Application V: Holistic Personal Genome Characterization, in Normal Individuals



**AN EXAMINED LIFE**
The longitudinal study collected data at daily and three-month intervals, and allowed personalized interventions -- such as changes in diet -- as the study proceeded.

**BRAIN**
- What's measured: Sleep patterns
- Frequency: Daily
- Method: Wrist sensor

**HEART**
- Pulse, physical-activity level
- Daily
- Wrist sensor

**COLON**
- Microbiome ecology
- Every three months
- Stool sample

**LIVER, LUNGS, BRAIN & HEART**
- 100 proteins to track organ health
- Every three months
- Blood sample

**LYMPHATIC SYSTEM**
- Immune-cell activity
- Every three months
- Blood sample

**INSULIN SENSITIVITY**
- Blood glucose
- Every three months
- Blood sample

**CHROMOSOMES**
- Whole-genome sequence
- At enrollment
- Blood sensor

Institute for Systems Biology
*Revolutionizing Science. Enhancing Life.*

**(Figure from Institute for Systems Biology)**

- Mental disease & cancer are two extremes with respect to genomics (CEN, 92: 26)
  - Many other conditions in between, often involving interaction with the environment
- Pers. Genome Characterization
  - Identify mutations in personal genomes (SNPs, SVs, &c)
  - Estimate phenotypic (deleterious or protective) impact of variants.
  - Compare one person to wider population.
- Track changes over time & consider interaction w/ environment
  - Transcriptome studies
  - Longitudinal health studies (e.g. 100K wellness project, Framingham Heart Study)
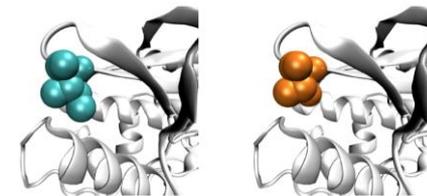
# Analyzing Carl Zimmer's genome



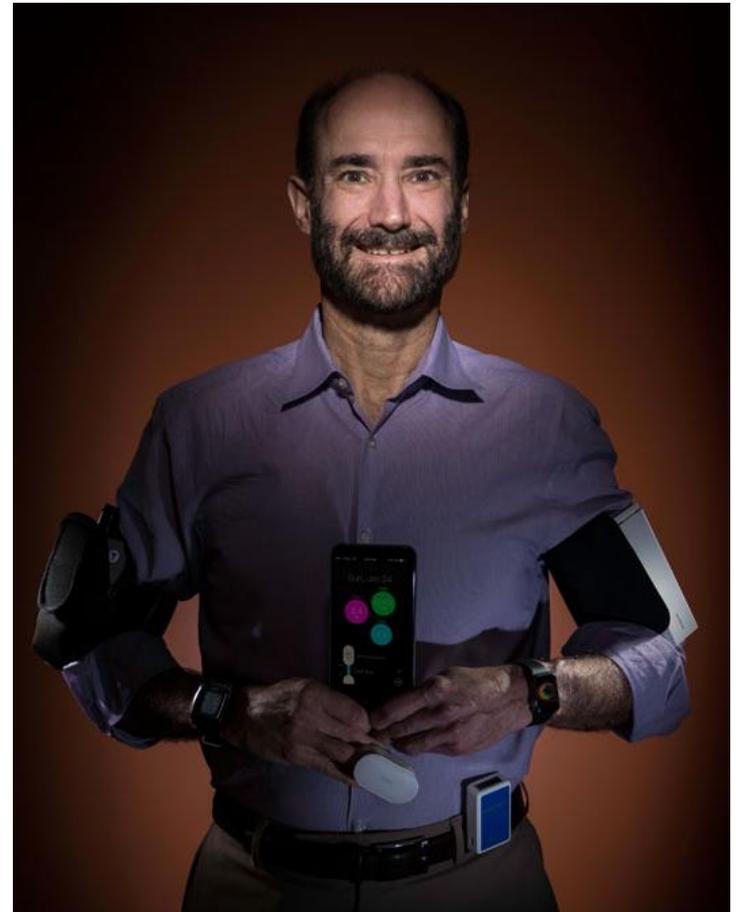SNV  AAGCT → ACGCT

Protein
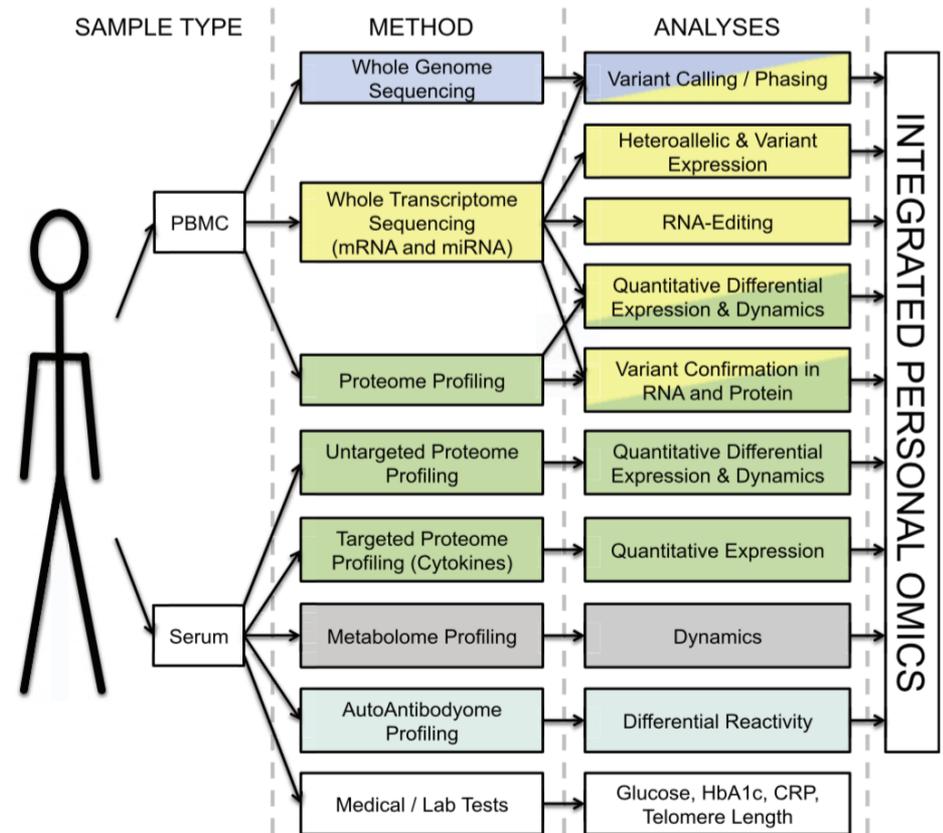Structure

Wild-type    Mutated

Ancestry

# Expanding personalized medicine beyond the genome.

- An integrated personal omics profile (iPOP) is an example of a more comprehensive version of personalized medicine.

- Michael Snyder had his genome sequenced and collected many other large scale datasets over an extended period of time.

# Integrated personal omics profile (iPOP)

- Numerous types of data were collected, primarily from blood samples. The datasets include:
  - Transcriptomic
  - Proteomic
  - Metabolomic
  - Cytokine profiling
  - Autoantibody profiling
  - Medical exams

Chen et al. Cell 2012

# History of the Analysis of the "Zimmerome" in the Class

## 2017

- Each group created a GitHub page detailing the work of each team
  - Additionally, each group has a power point presentation:

- Topics of projects include:
  - Comparative analysis of personal genomes
  - Personal genomes and personalized medicine (CRISPR)
  - Network analysis of personal genomes
  - Structure analysis

## 2018

- Each group had a power point presentation and a writeup
- Topics of projects include:
  - Finding how much of your genetic material comes from the Neanderthals
  - Using carl's genome to predict differences in gene expression from the average human and infer possible changes in physiology from these differences (GTEx analysis)
  - Predicting gene expression values from Carl's SNP information
  - Finding a common variant associated with inflammatory response in Carl
  - Calculating Zimmer's risk for Alzheimer's disease
  - Identifying significant protein-coding mutations in Carl's genome
  - polygenic risk score prediction in coronary artery disease, type II diabetes, and schizophrenia for Carl

# History of the Analysis of the "Zimmerome" in the Class

## 2019 - 2022

- **Each group had a power point presentation and a write-up**
- **Started analyzing Carl's gene chromosome by chromosome**
  - **Part 1: Prioritization of 10 genes**
  - **Part 2: In-depth Analysis of prioritized genes:**
    - **Gene expression analysis**
    - **Network analysis**
    - **Protein structure analysis**
    - **Text mining analysis**

## Genes prioritized

# History of the Analysis of the "Zimmerome" in the Class

26 groups (2019-2022)

↓

10 groups identified SNPs

↓

Total of 36 SNPs are disease associated

## Disease Associated SNPs



| Disease Type |
|---|

Schwartz Jampel Syndrome: 6
Eczema: 1
Metabolic Syndrome: 1
Atrial Filbrillation and Ataxia Telangiectasia: 1
Obesity: 1
Cleft Palate: 1
Abetalipoproteinemia: 1
Lower height and weight: 3
Asthenozoospermia: 1

# This year's Zimmerome Assignment: Investigate and Analyze a Personal Genome Using Bioinformatic/Biomedical Tools





**Team based approach**
- **Assigned Teams (4-5 people in your section, assigned by TFs)**
- **Each team focuses on a single chromosome**
- **Cross-disciplinary**

**1. Computational**
- **Leveraging tools to prioritize genes or variants**
- **Pipeline Development**

**2. Biological/Biomedical**
- **Interpretation of prioritized genes or loci**

**3. Written and Oral**
- **Communication of project and results through written report**

# 1. Computational Pipeline Development

**1**

**VCF to BED**

Converted Zimmer SNV VCF file for ease of use; filtered for Ch17 (*BEDOPS*)

**2**

**GENCODE**

Took GTF file for Gencode (GRCh37) and converted to BED (*BEDOPS*)

**3**

**Filtering**

Extracted CDS regions only; eliminated repeat entries; kept position/category/gene info

**7**

**Future Direction**

Weight variants with other variant prioritization tools or databases

Noncoding analysis

**4**

**Intersect Files**

Intersected annotation file with variant file (*BEDTools*), created gene-SNV barcode

**5**

**Removing Duplicates**

Eliminated repeat position entries from gene isoforms using barcodes

**6**

**Compile Data**

Sum mutations by gene, sort high to low, extract top 10; convert file to VCF

→ **GTEx**

**Computational Pipeline**
- **Full code/software/script package**
- **GitHub**
- **Data files**
- **readme**

# 2. Biological/Biomedical Interpretation

Tissue Specific
Expression
Extracted from
GTEx



**Interpretation of Results**

- Biological interpretation of prioritized genes or loci

- Leveraging public omics or biomedical data

- Further discussion of results

# 3. Oral Presentation and Written Report

**Oral Presentation**
- **April 26**
    - **2 minute mp4 recording per group (nominate 1 person to make the recording)**
    - **We will play these recordings in class on 4/26**
- **April 27 or 28 in your Section**
    - **10 minute presentation by other members of the group**

**Written Report**
- **Due: May 10, 2023**
- **At least 1000 words**

**Summary Slide**
- **1 summary slide giving an overview of your project**
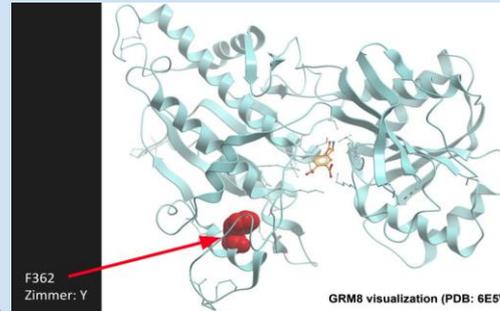
**Summary Metadata File**
- **A single text file containing relevant information**
- **More description in assignment file**
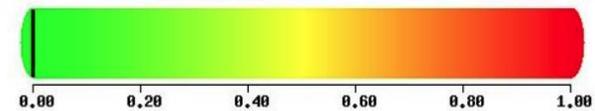
## 2020 group 2 (chr 7)

Top 10 Prioritized Genes

1. CNTNAP2
2. MAGI2
3. PTPRN2
4. DPP6
5. SDK1
6. DGKB
7. AUTS2
8. HDAC9
9. GRM8
10. PDE1C

Summary Figure:



F362
Zimmer: Y

GRM8 visualization (PDB: 6E5V)

### PolyPhen-2:

This mutation is predicted to be **BENIGN** with a score of **0.000** (sensitivity: **1.00**; specificity: **0.00**)

0.00   0.20   0.40   0.60   0.80   1.00

Curr Protoc Hum Genet. 2013 Jan; 0 7: Unit7.20.

Summary:

1. Prioritization approach: mutational burden
2. Downstream analysis: PDB structural analysis
3. Findings:
   a. Among the top 10 most mutated genes on chromosome 7, there are 6 missense variants within 4 genes
   b. Only one variant, conferring a protein, is characterized: GRM8
   c. PolyPhen analysis shows that substitution at pos 362 from F to Y is predicted to be tolerated

## 2023 group N (chr Z)

**Top 10 Prioritized Genes**

1.
2.
3.
4.
5.
6.
7.
8.
9.
10.

Summary Figure:

Summary:
1. Prioritization approach:
2. Downstream analysis::
3. Findings:
   a.
   b.