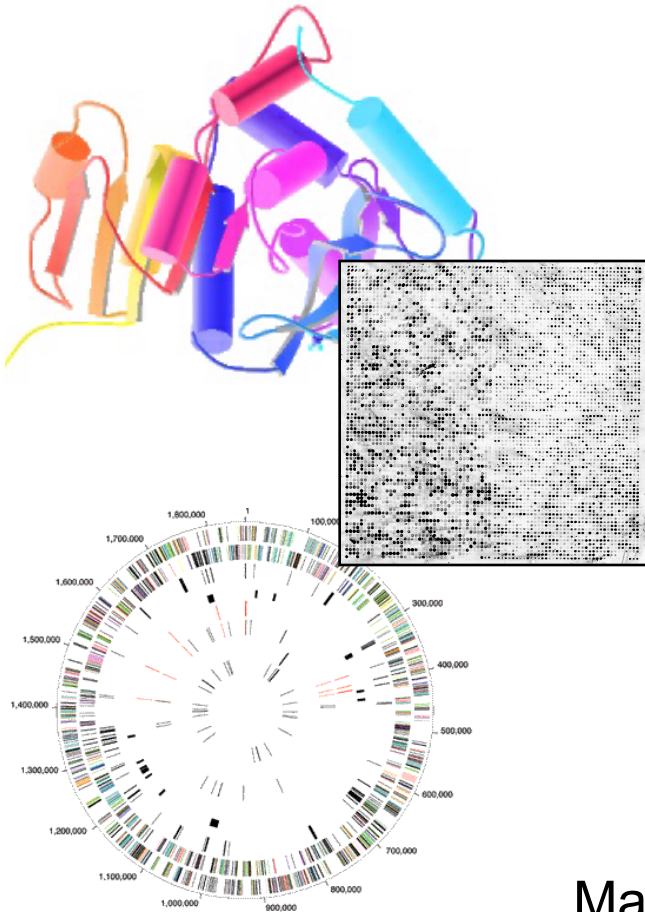


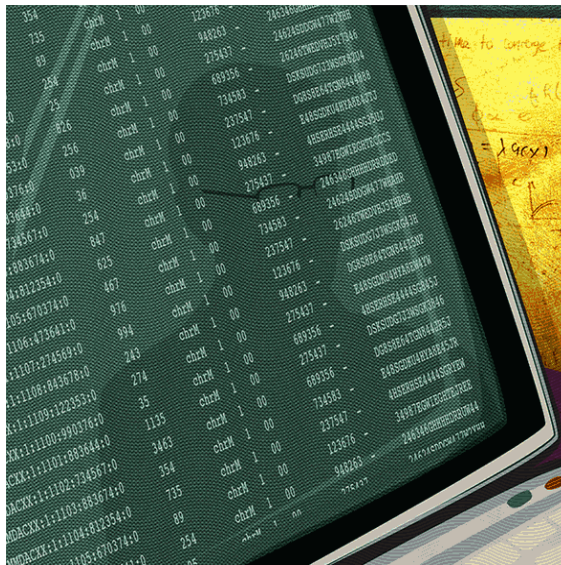
# Biomed. Data Sci. Personal Genomes Intro.



Mark Gerstein, Yale University  
[GersteinLab.org/courses/452](http://GersteinLab.org/courses/452)  
(last edit in spring '20)

# Analyzing Carl Zimmer's genome

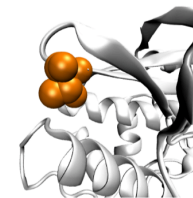
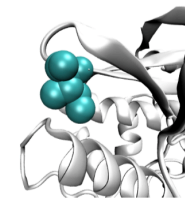
## CARL ZIMMER'S GAME OF GENOMES SEASON 1



SNV

AAGCT → ACGCT

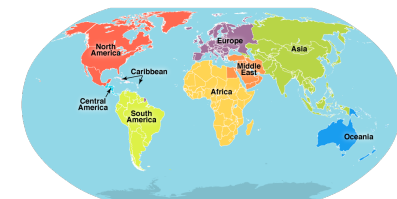
Protein Structure



Wild-type

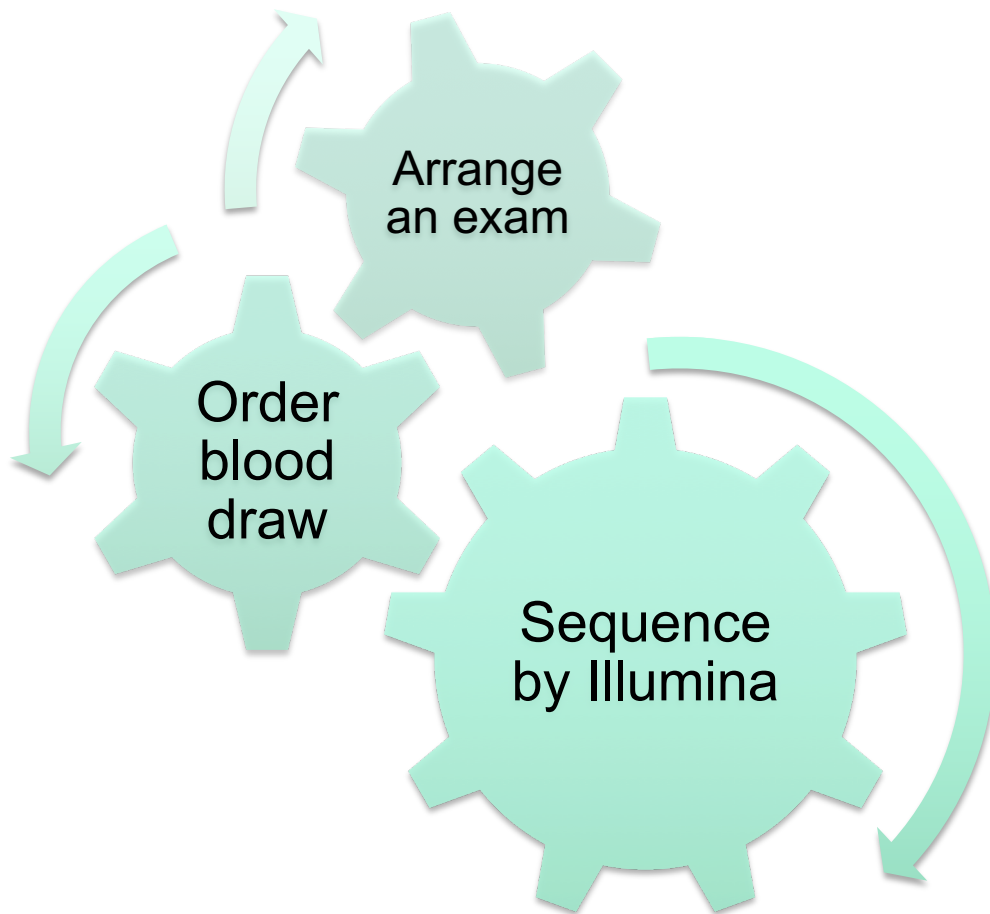
Mutated

Ancestry



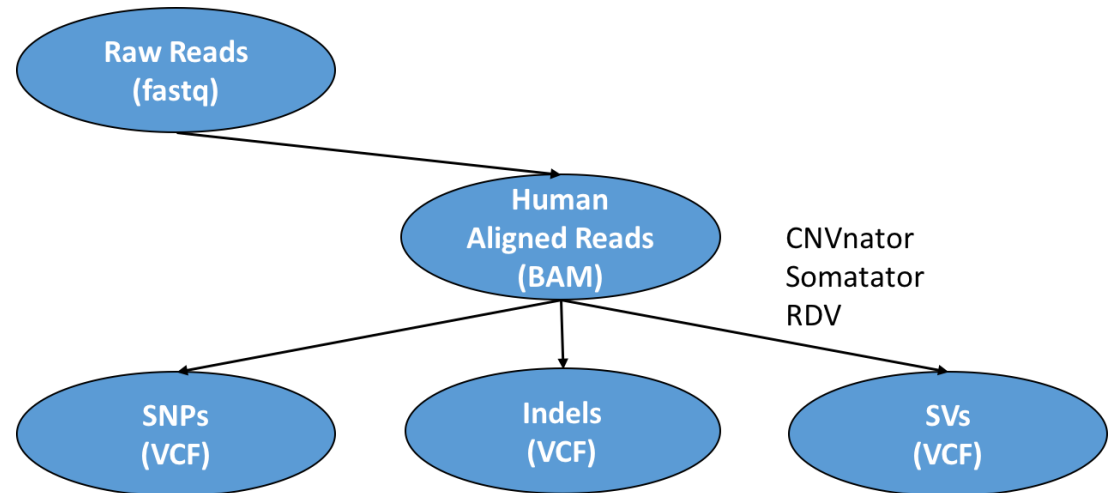


CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



- **Cost: \$3100**
- **Illumina briefly review the sequencing data, evaluating the risk for 1200 disorders, from familiar ones like lung cancer to obscure ones like cherubism**

CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



# Genome Variation

TP53 Sequence:

...GGAGTCTTCCAGTGTGATGATGGTGAGGATGGGCCTCCGGTT...

Single Nucleotide Polymorphism (SNP) – 1nt:

...GGAGTCTTCCAGTGTGATGATGGT**G**AGGATGGGCCTCCGGTT...

T or A or C

Small Insertions and DEletions (INDEL) – 1-10nt:

...GGAGTCTTCCAGTGTGATGATGGT~~GAGGATG~~GGCCTCCGGTT...

Large Structural Variations (SV) -- >100nt:

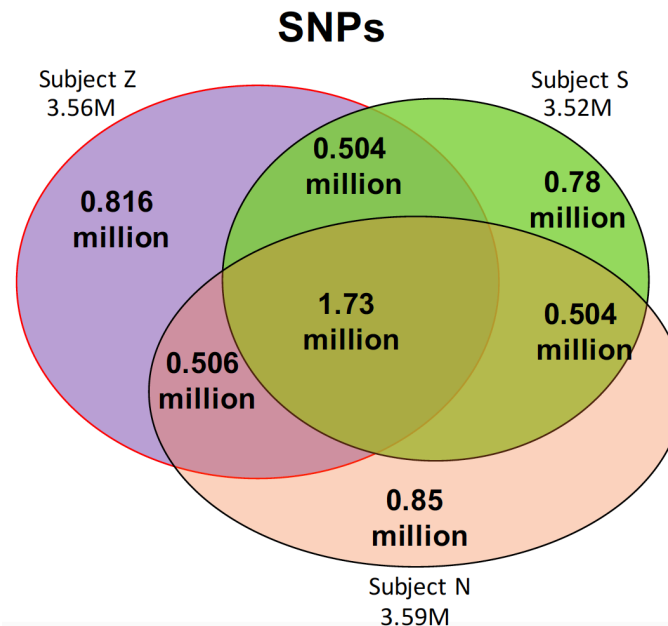
...GGAGTC~~TTCCAGTGTGATGATGGTGAGGATGGGCCTCCGGTT~~...



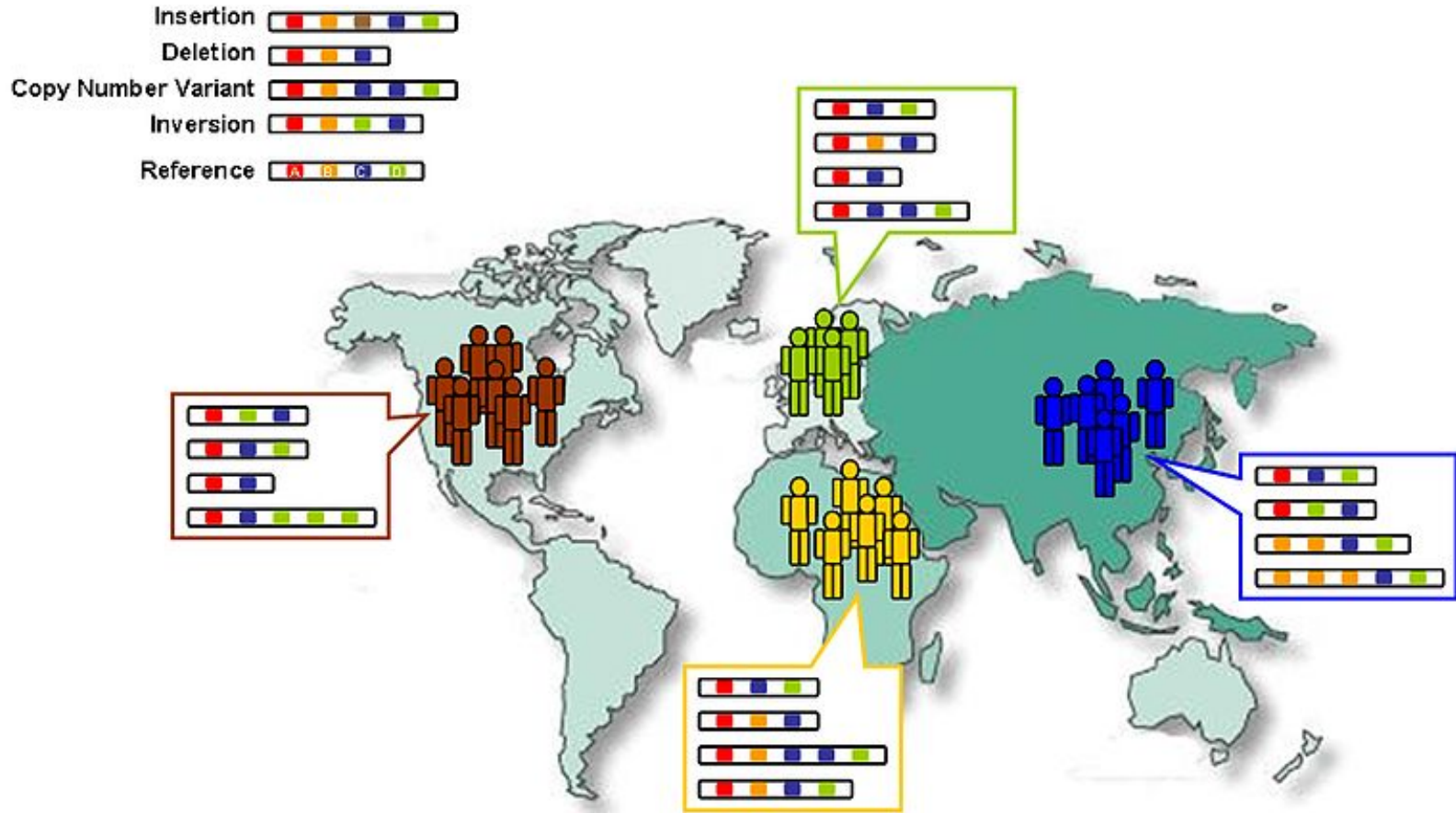
CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



- Normal range of number of SNPs
- Carl's case: more than 3M SNPs
- How do we know if the SNP is harmful?



- Thousand genome project
- Common SNP data base found in the population



# Human Genetic Variation

A Cancer Genome



A Typical Genome

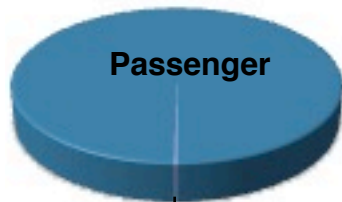


Population of 2,504 peoples



Origin of Variants

	Coding	Non-coding
Germ-line	22K	4.1 – 5M
Somatic	~50	5K



Driver (~0.1%)

Class of Variants

SNP	3.5 – 4.3M
Indel	550 – 625K
SV	2.1 – 2.5K (20Mb)
Total	4.1 – 5M

Prevalence of Variants



Rare\* (1-4%)

SNP	84.7M
Indel	3.6M
SV	60K
Total	88.3M

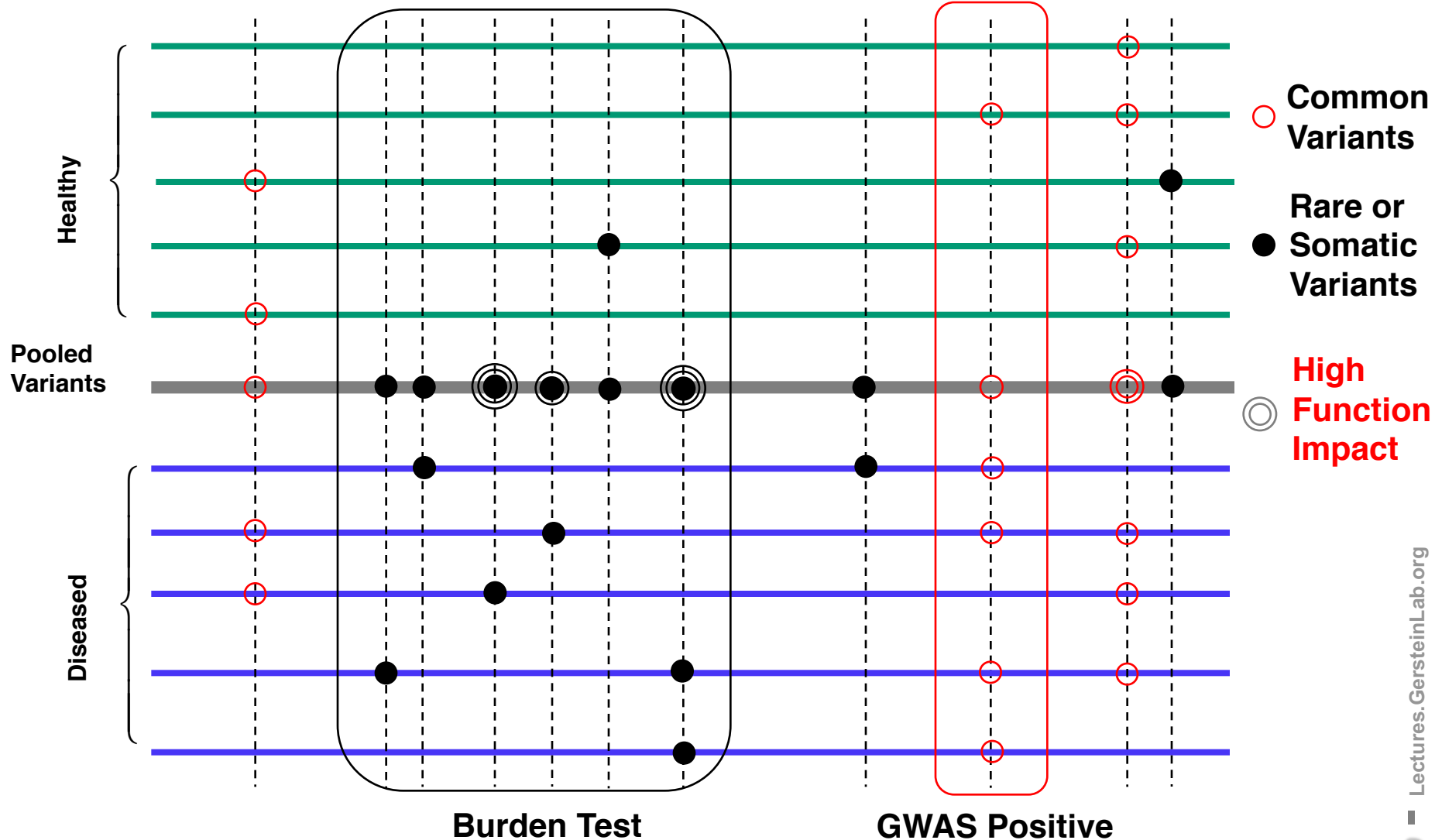


Rare (~75%)

\* Variants with allele frequency < 0.5% are considered as rare variants in 1000 genomes project.



# Association of Variants with Diseases



CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



- **Got a variant in a gene for heart muscles, called DSG2**
- **DSG2 gene encodes a protein in humans called Desmoglein-2**
- **Mutations in desmoglein-2 have been associated with arrhythmogenic right ventricular cardiomyopathy**

**1 in 200**

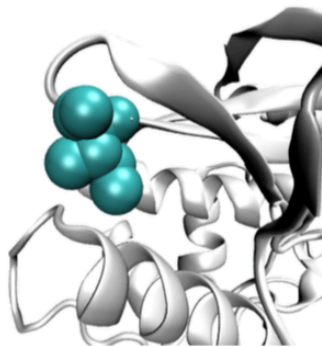
People of European descent carry this variant

**We're all different in our DNA. We're finally starting to understand when those differences matter ---- Carl Zimmer**

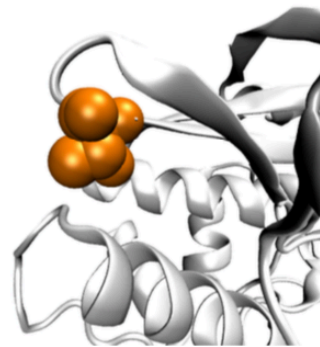
CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



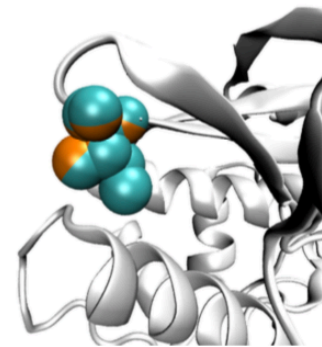
## SNP changing protein structure



*Wild-type*



*Mutated*



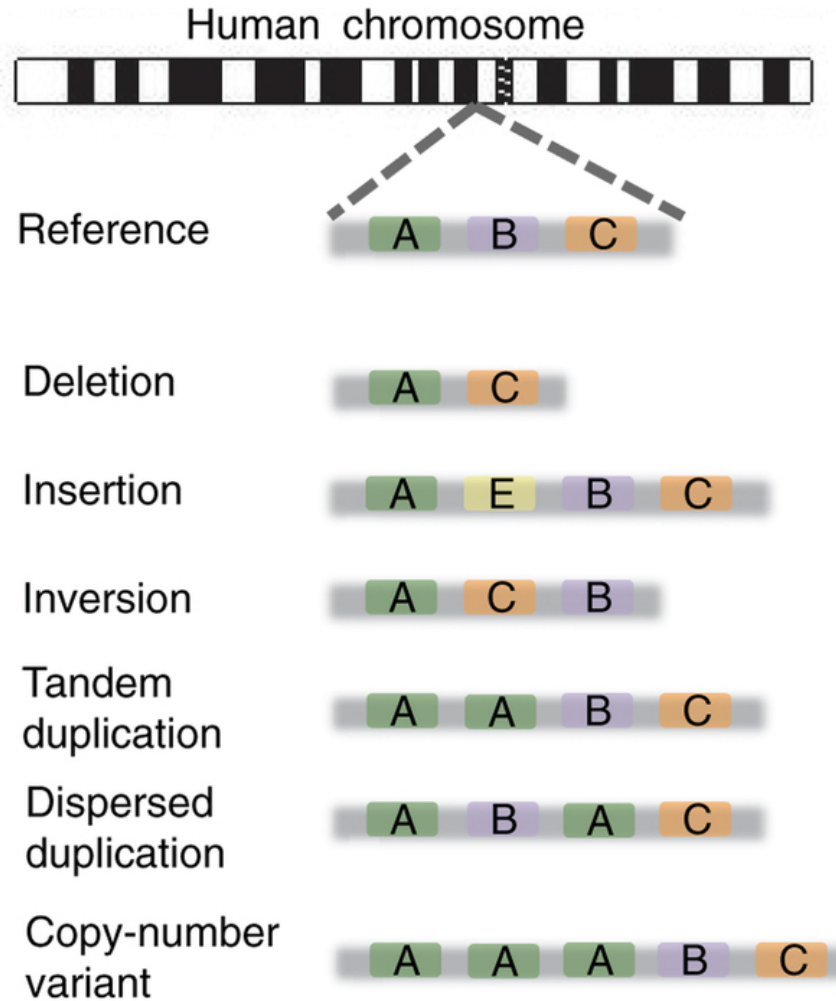
*(superimposed)*

**114: I->T**

- NAT2, an enzyme in the liver that breaks down caffeine and other toxins with a similar molecular structure.
- NAT2 helps break down certain medicines too. The variant puts people at risk of bad side effects from those drugs.



# Structural Variation



CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



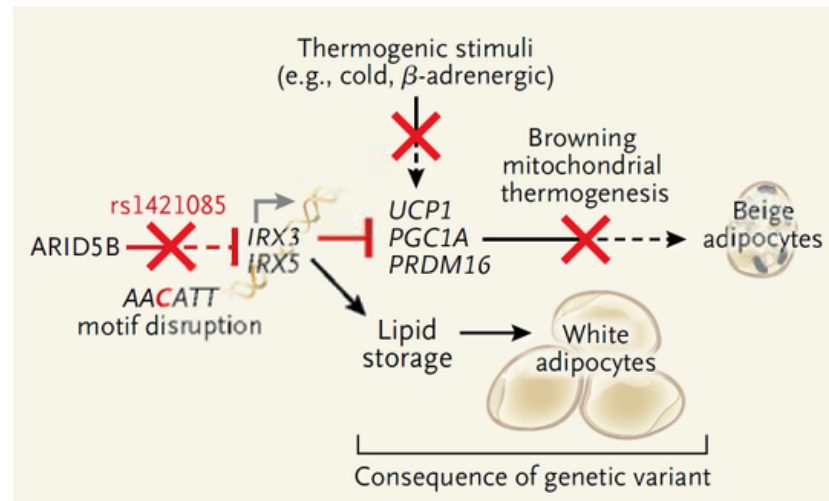
- Structural variation
- Example: HTT
- Certain mutations in HTT cause Huntington's disease.
- Healthy people have a wide range of CAG repeats. It's only when people get 37 or more CAG repeats in HTT that they are at risk of developing Huntington's disease.
- The reference genome has 19 CAG repeats. Carl has 17.

CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 2



## Non-coding variant

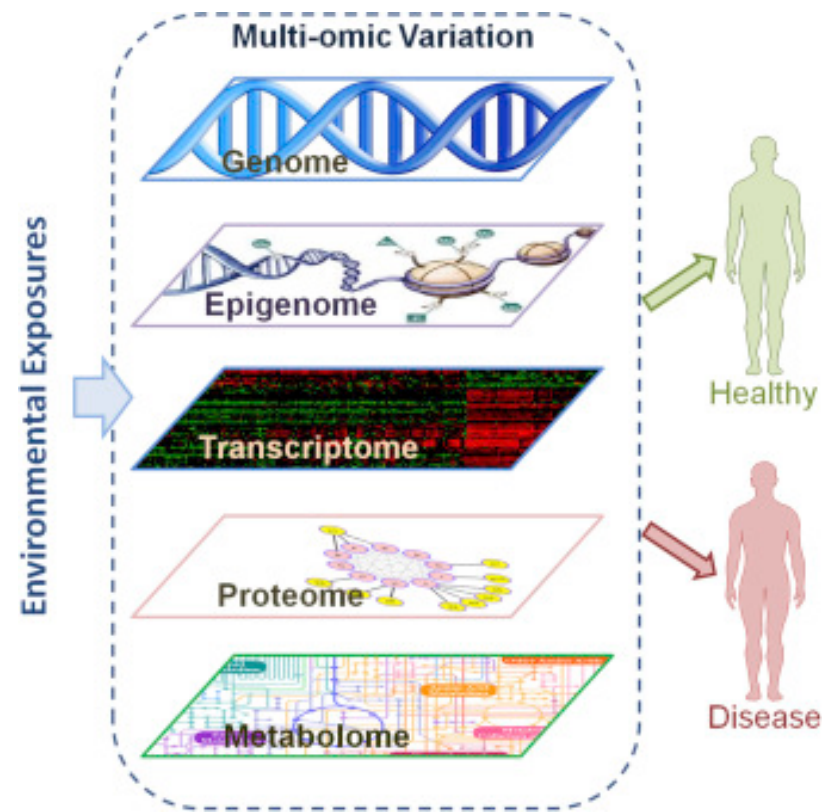
- Variant rs1421085
- Located in a genetic switch that activates several genes in fat cells
- The variant causes people to put on an average of 7 pounds





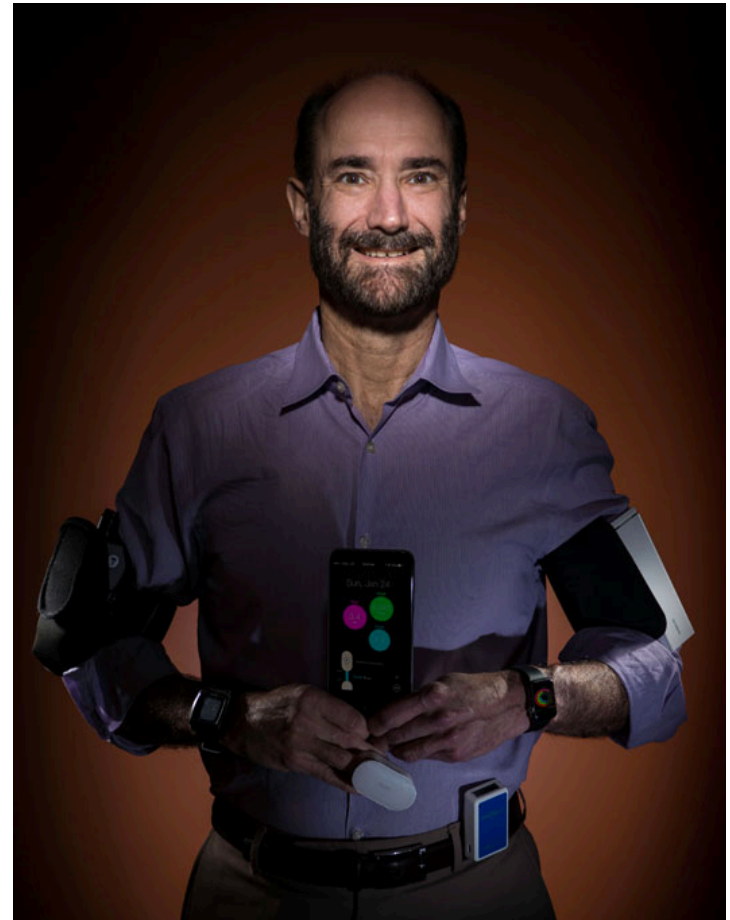
# Integrating environmental factors, genetic background, and large-scale datasets

- Difference between health and disease depends on many factors.
- Environment, genome, cellular contents, etc. all play a role.
- Important to integrate information from multiple large-scale datasets.



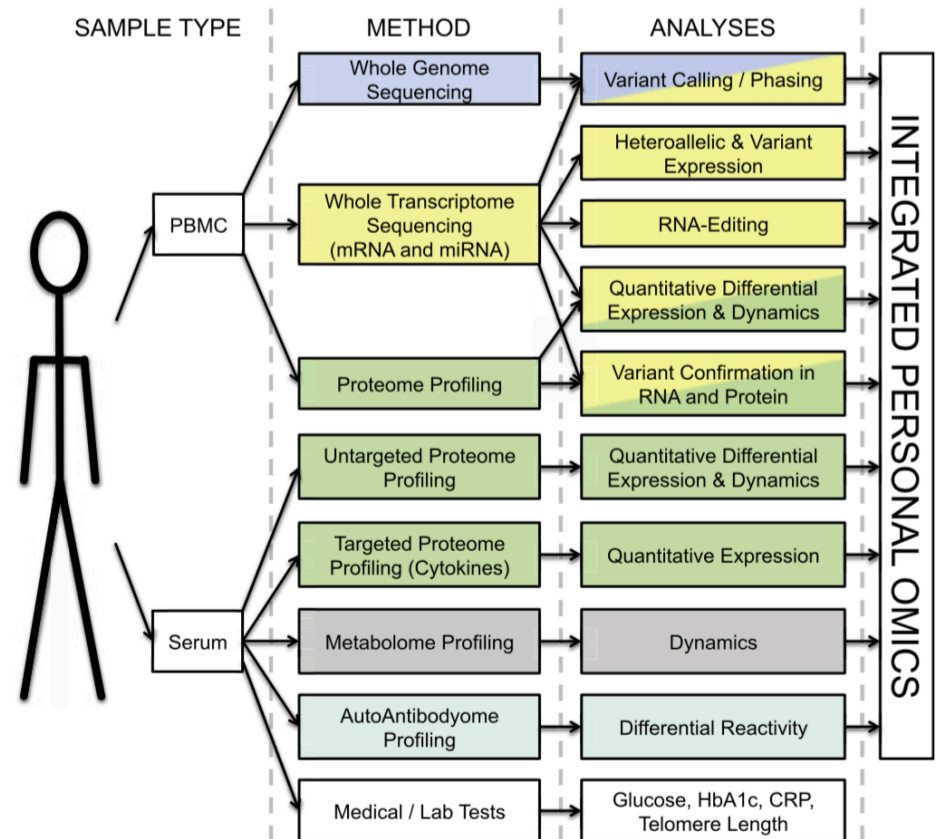
## Expanding personalized medicine beyond the genome.

- 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.



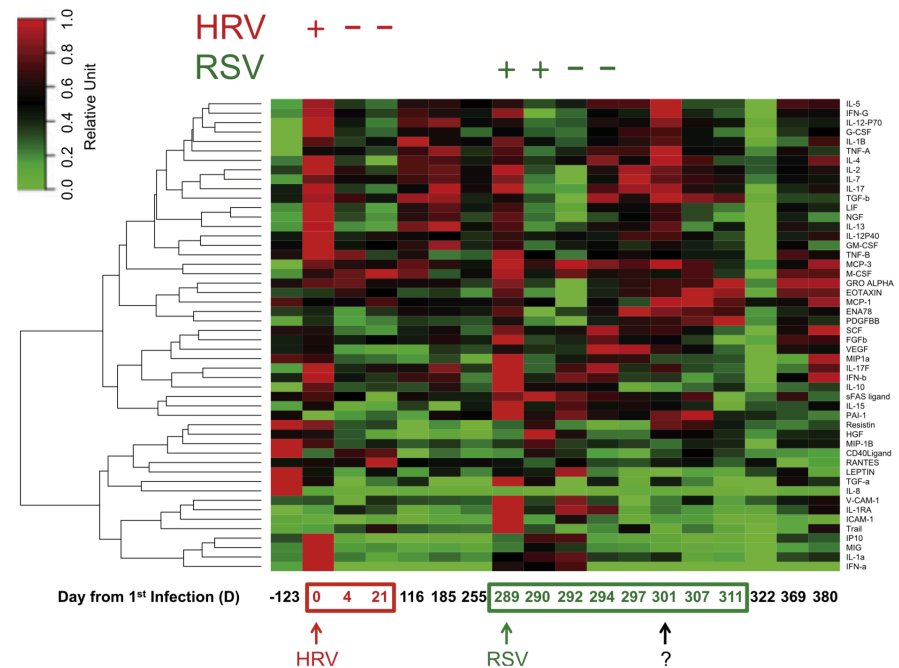
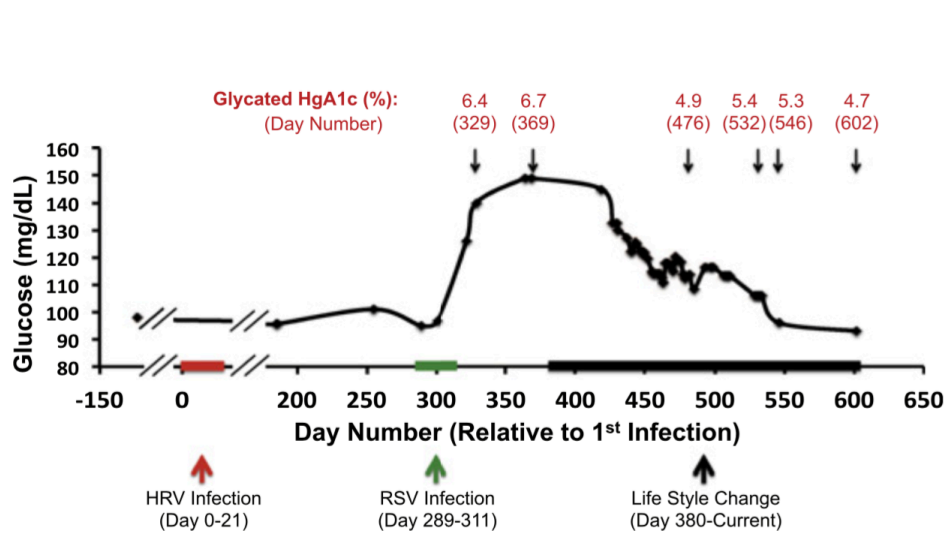
# Integrated personal omics profile (iPOP)

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





# Longitudinal medical data



- Tracking relevant medical (e.g. blood glucose) data over time helps link phenotypic changes with changes at the molecular level.