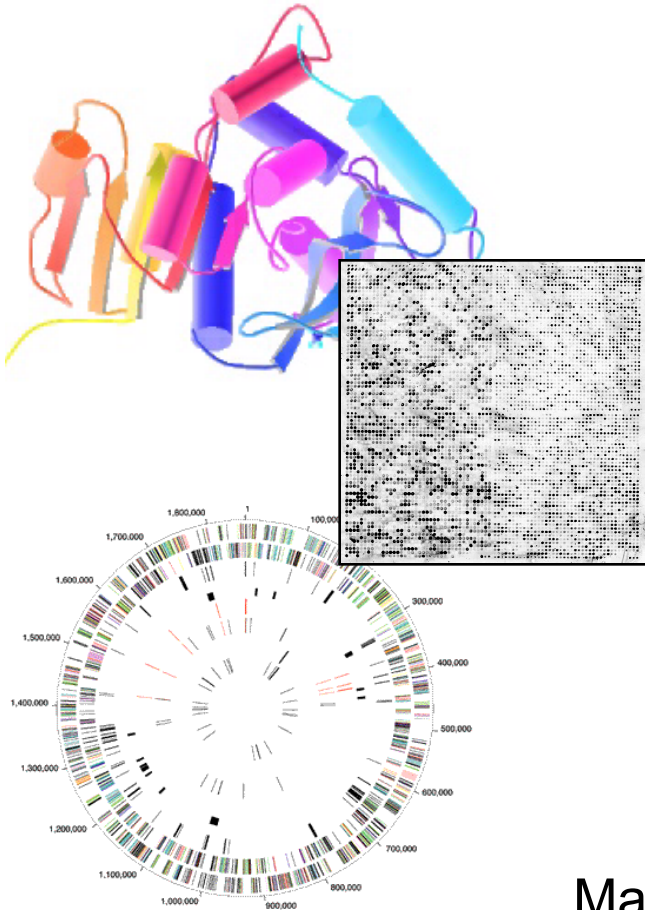


# BIOINFORMATICS

## Personal Genomes Intro.



Mark Gerstein, Yale University  
[GersteinLab.org/courses/452](http://GersteinLab.org/courses/452)  
(last edit in spring '18, MG lecture #2, final edit)

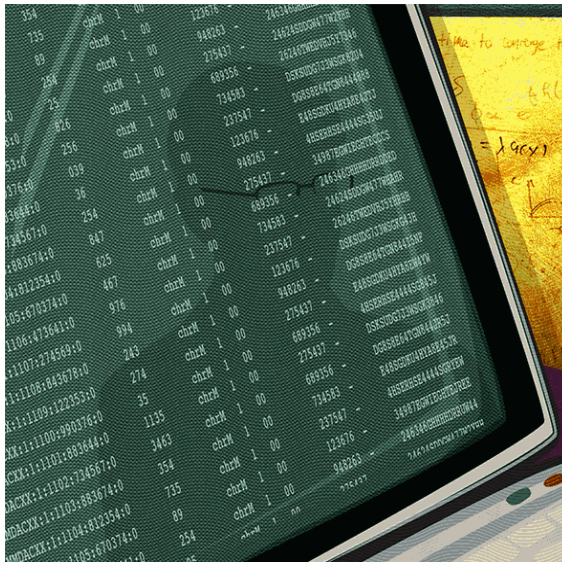
# Personal Genomics

## as an an organizing theme for this class

- A personal genome can reveal a lot about an individual.
  - Disease risks, ancestry, personal traits, etc.
- Personal genome annotation combined with multi-omic and longitudinal health data can inform new links between genotype and phenotype relevant to an individual and the larger population.
- Genomic privacy will become increasingly important as precision medicine becomes more common.
- In this class, we will look at how to identify key genomic variants with the most impact.
- We will also use analysis techniques including systems and network modeling as well as structural modeling to contextualize and interpret the mechanisms through which these variants impact health.

# Analyzing Carl Zimmer's genome

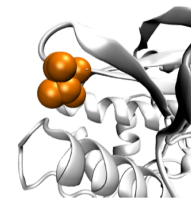
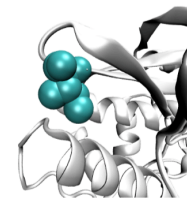
## CARL ZIMMER'S GAME OF GENOMES SEASON 1



SNV

AAGCT → ACGCT

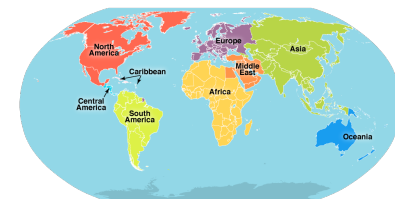
Protein  
Structure



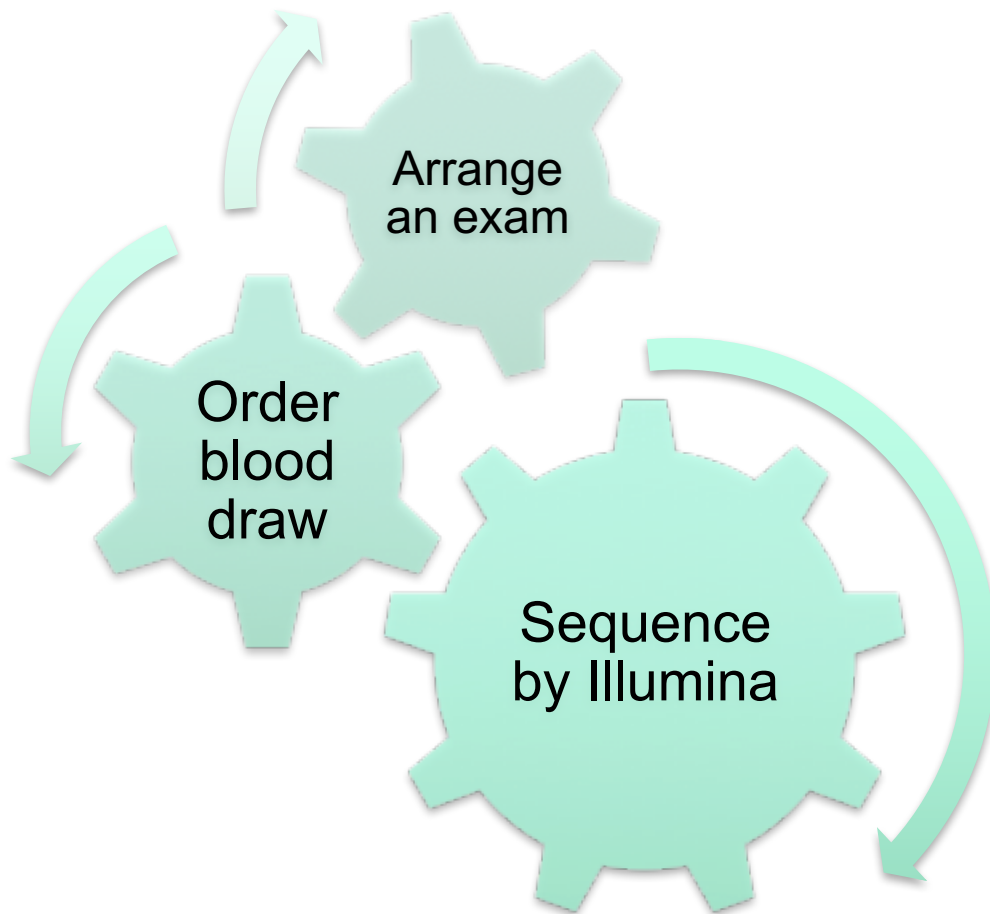
Wild-type

Mutated

Ancestry

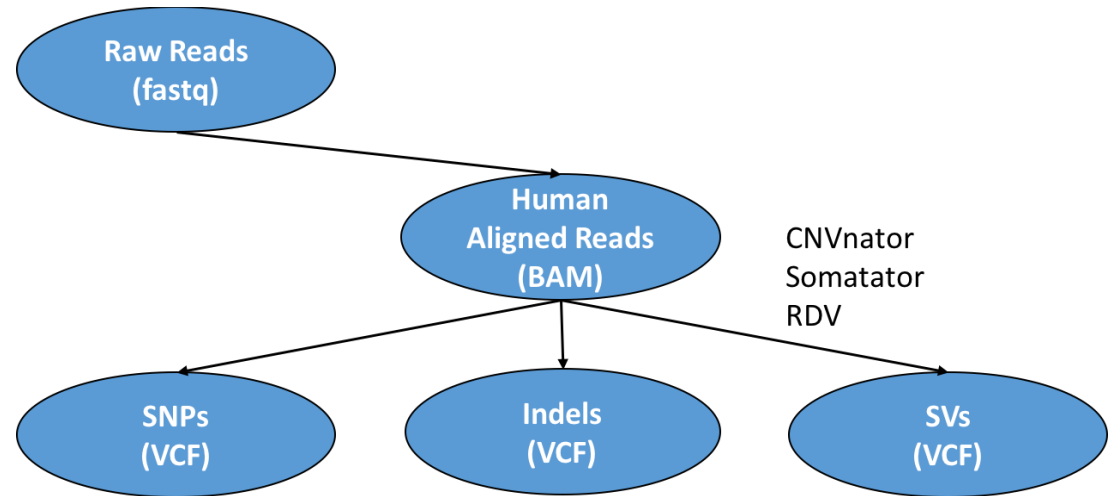


CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



- **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

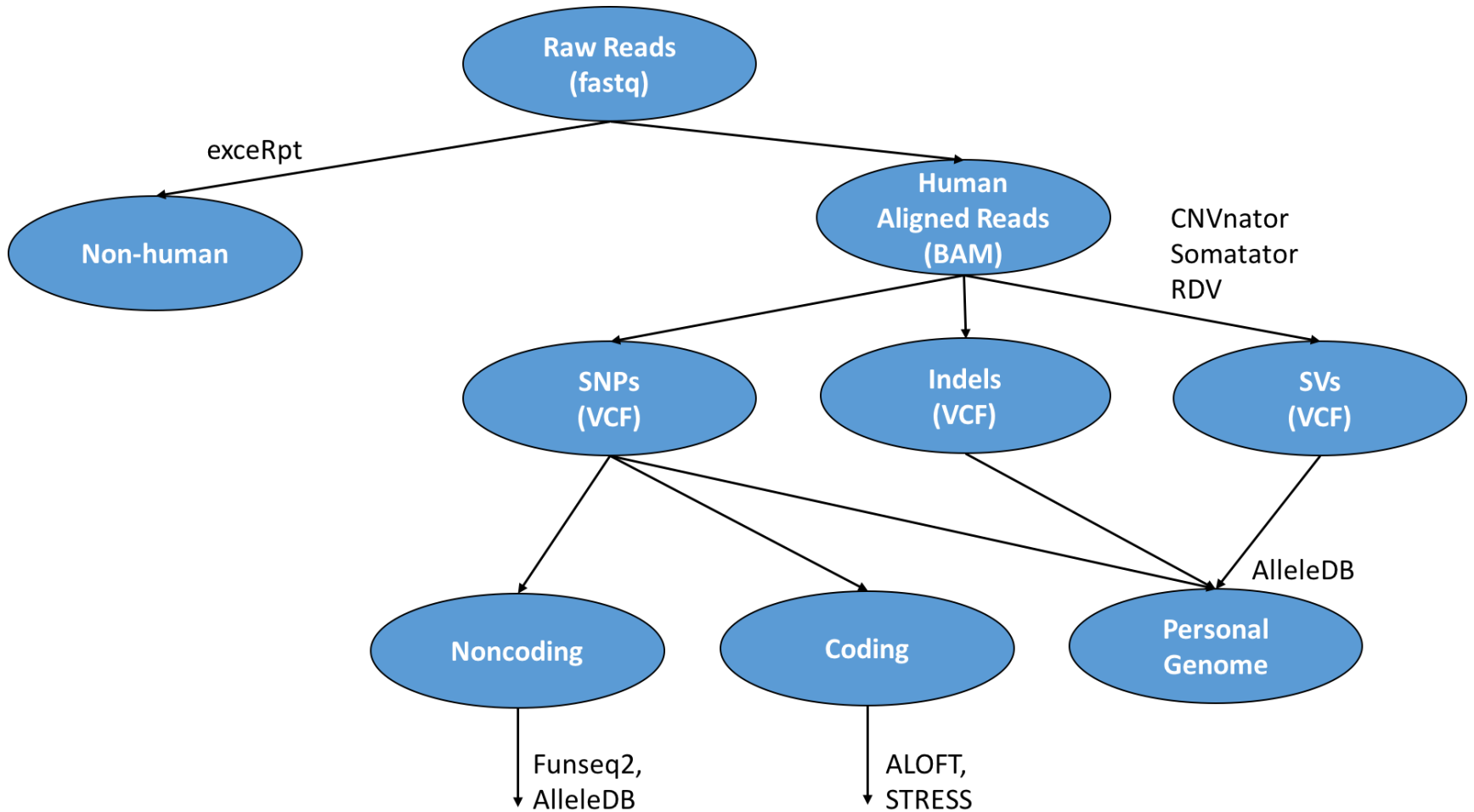
Small Insertions and DEletions (INDEL) – 1-10nt:

...GGAGTCTTCCAGTGTGATGATGGT ~~GAGGATG~~GGGCCTCCGGTT...

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

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

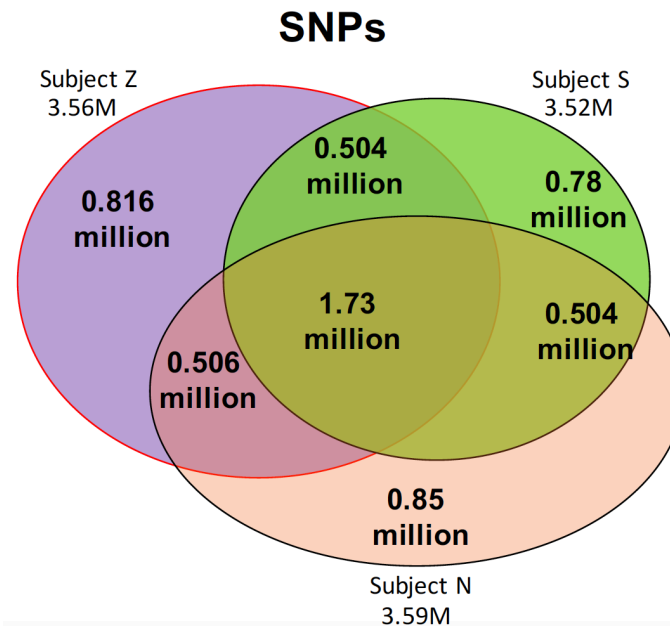
CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



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

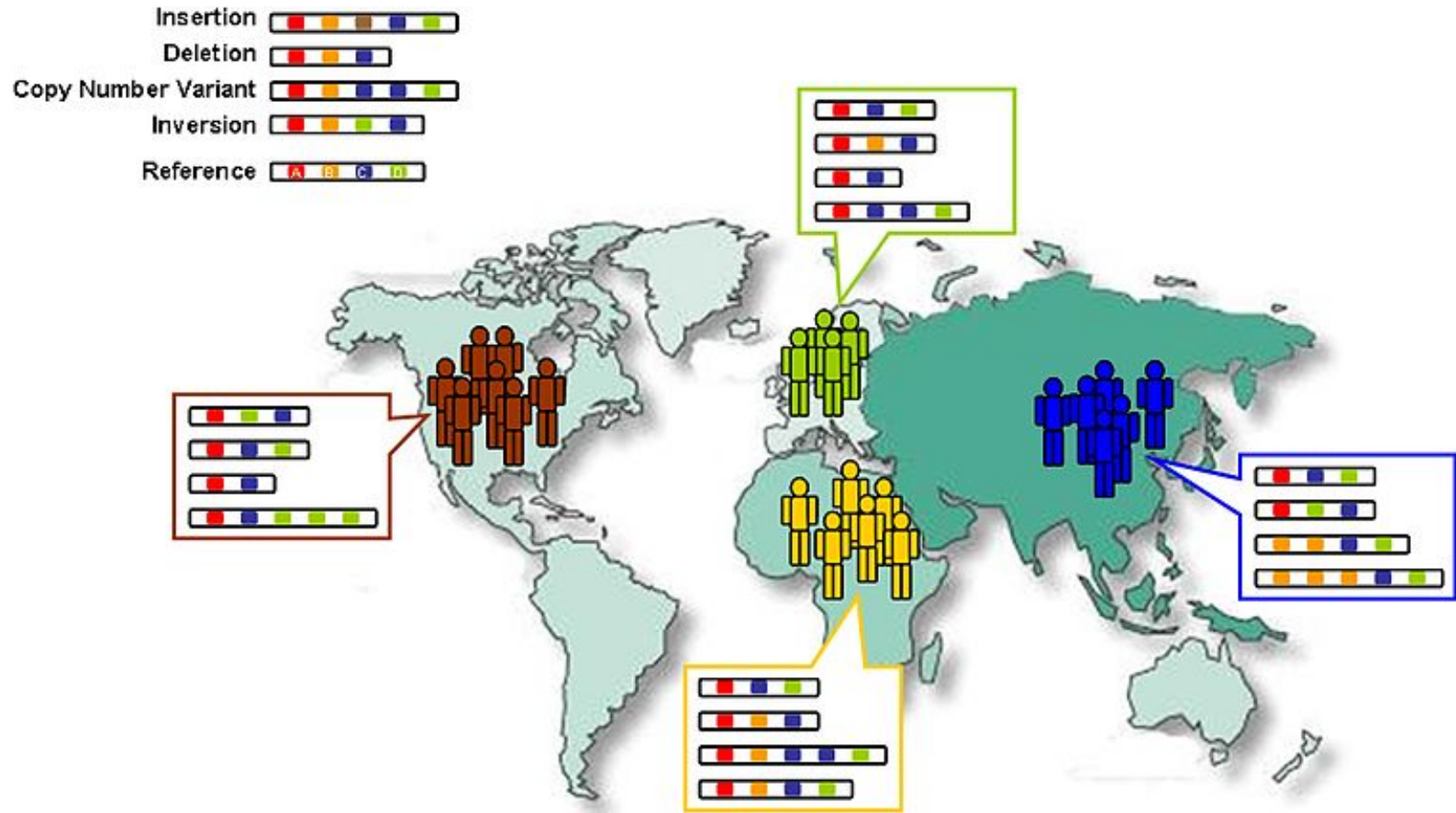


- Normal range of number of SNPs
- Carl's case: more than 3M SNPs
- How do we know if the SNP is harmful?





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



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**

# 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%)

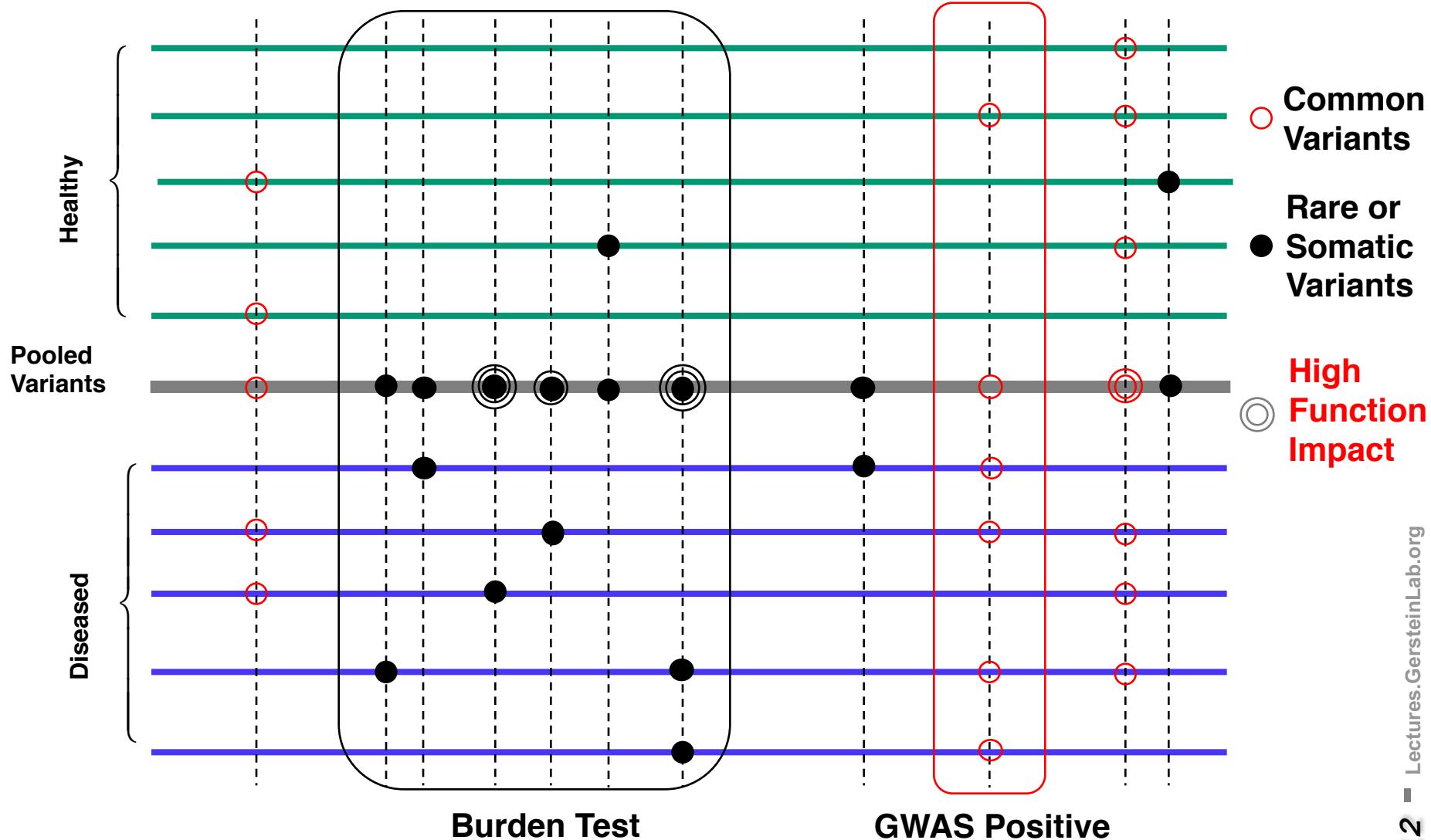
SNP	84.7M
Indel	3.6M
SV	60K
Total	88.3M



Rare (~75%)

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

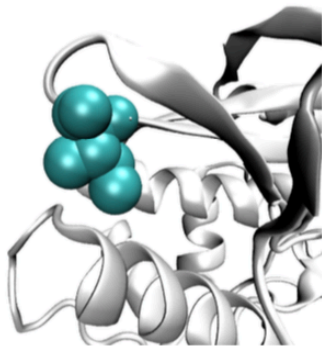
# Association of Variants with Diseases



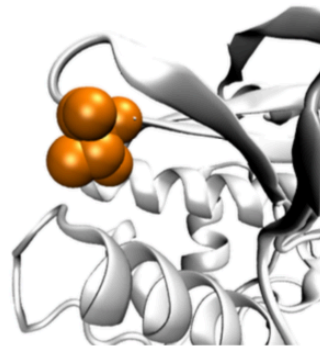
CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



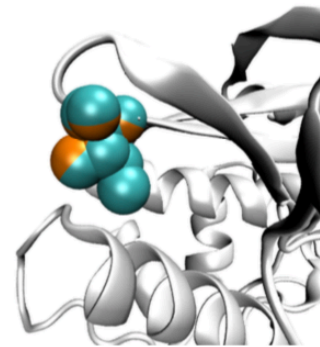
## SNP changing protein structure



*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.

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



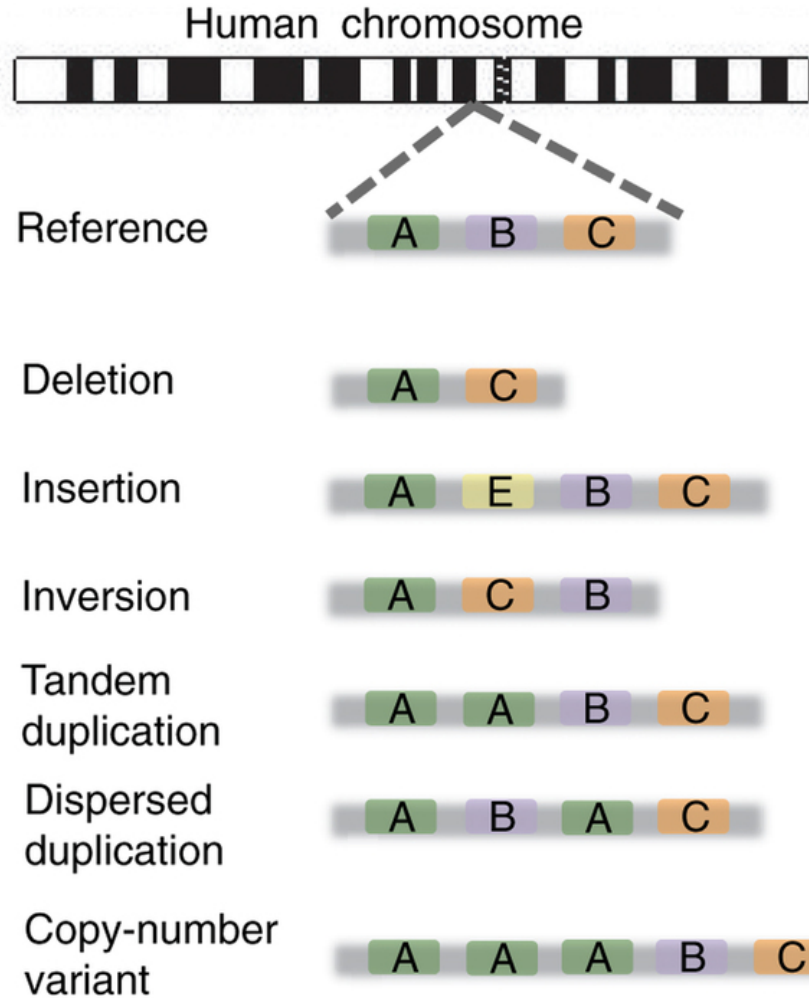
## Indels (Insertions/deletions)

- In coding regions, unless the length of an indel is a multiple of 3, it will produce frameshift mutation
- Likely to disrupt genes (loss-of-function variant)

Example: Non-functional F8 gene

- Can't make essential clotting protein
- Get hemophilia and can bleed to death from a little cut

# Structural Variation



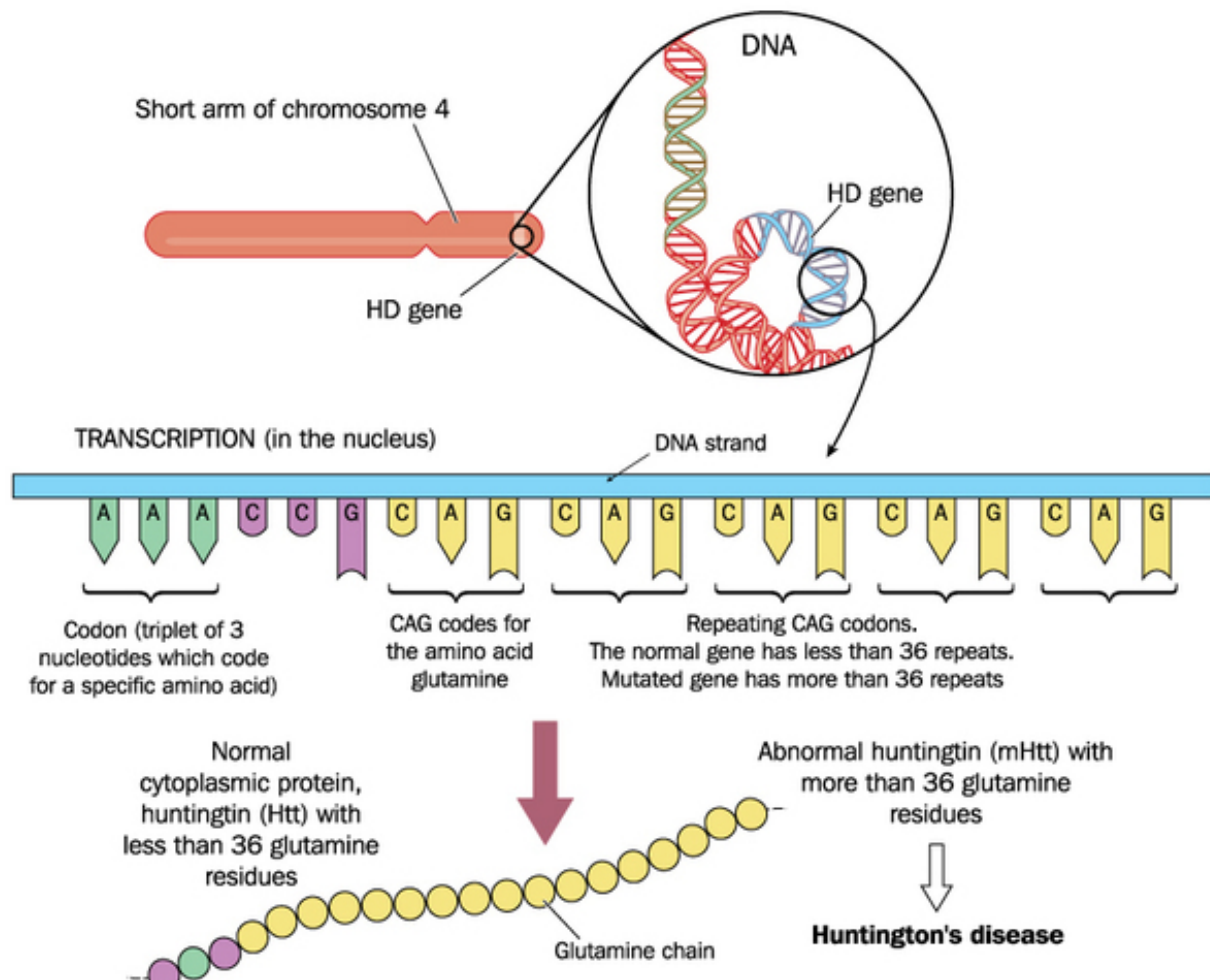
CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



- 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 1



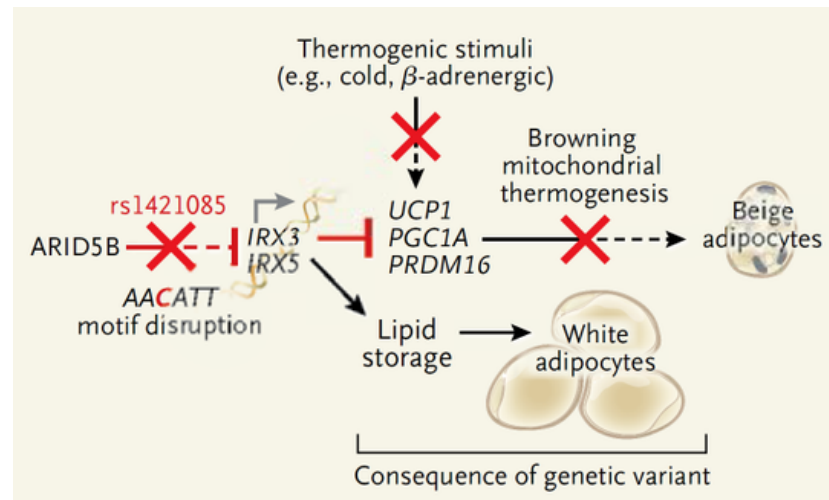
<https://ghr.nlm.nih.gov/condition/huntington-disease>

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

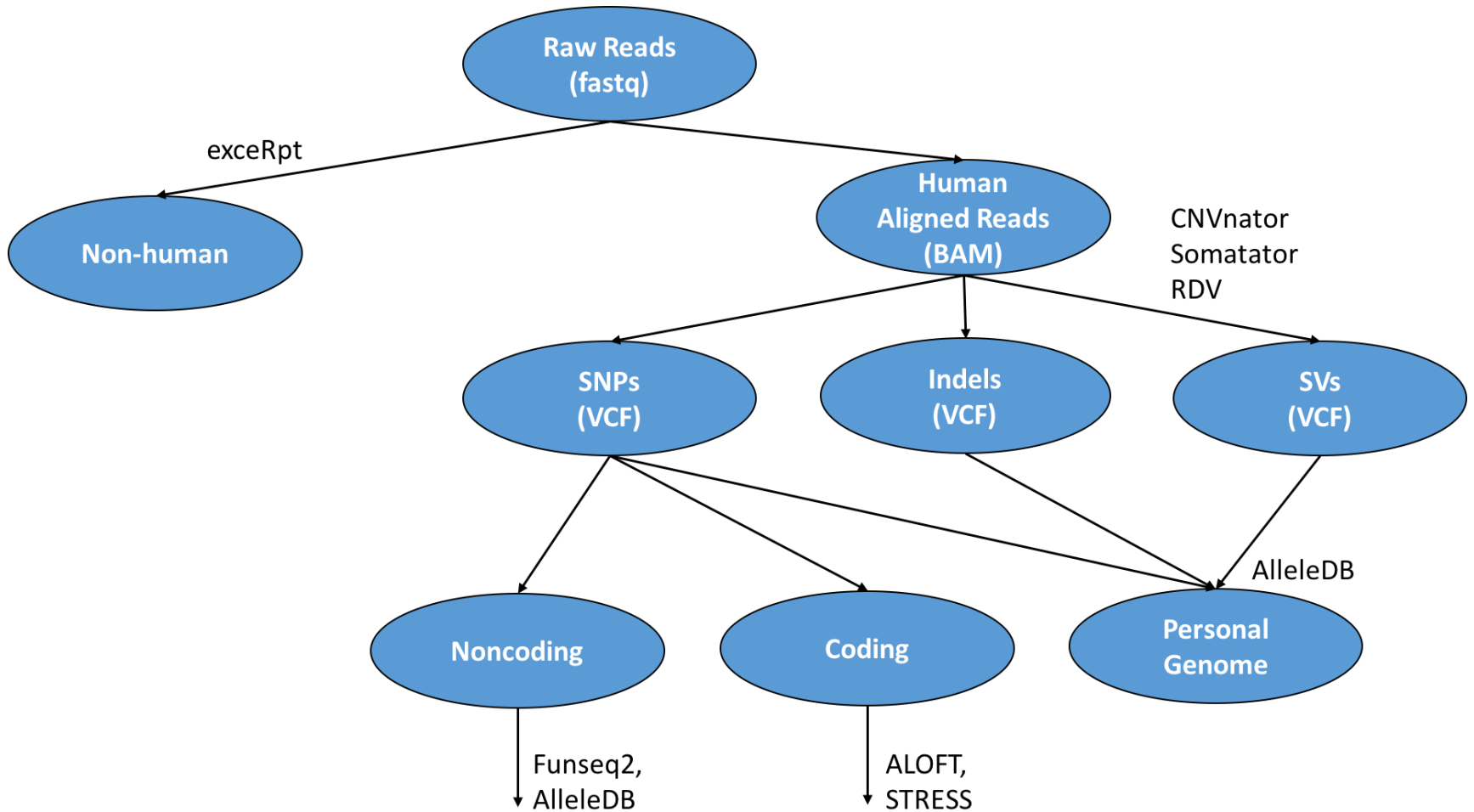


## Non-coding variant

- 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



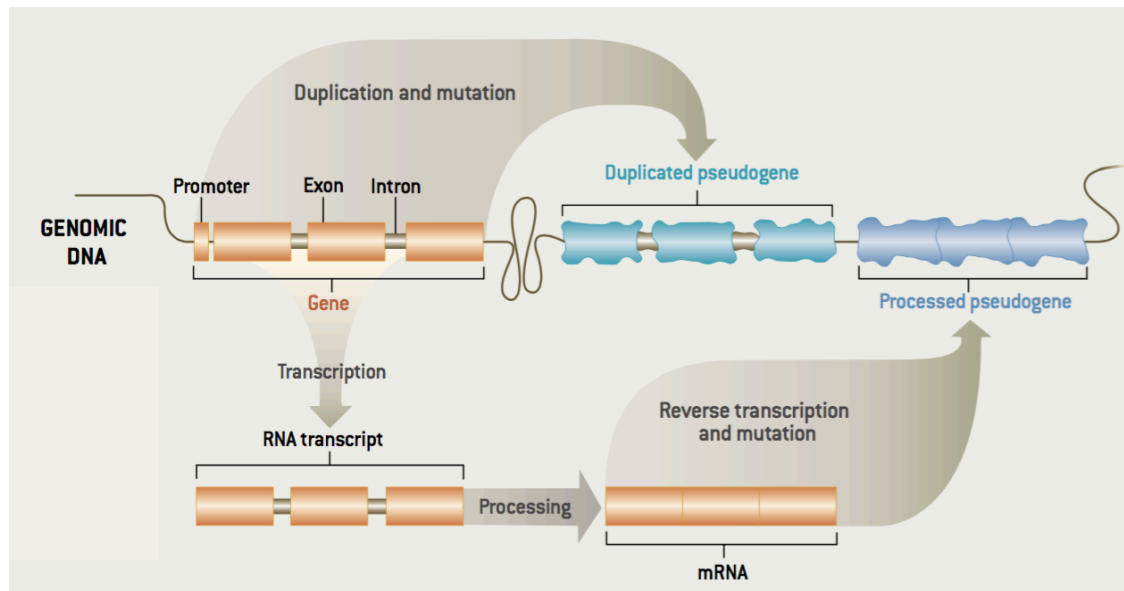
CARL ZIMMER'S  
**GAME OF GENOMES**  
SEASON 1



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

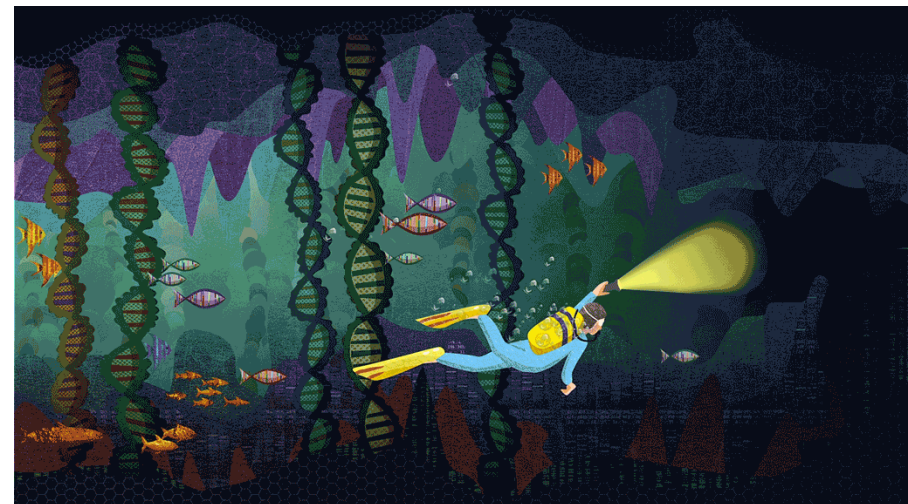


- What else are in the genome?
- Pseudogenes
- About 14000 pseudogenes carried in our genome



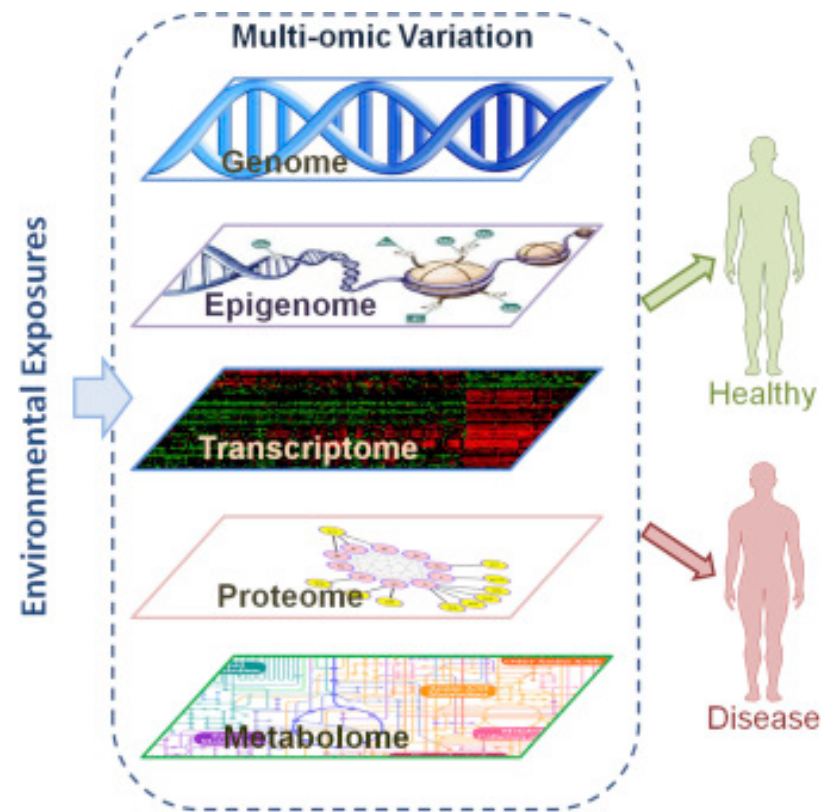
## Looking beyond the genome

- In the Game of Genomes Carl Zimmer explored his genomic sequence.
- The genome provides a mostly static view.
- Misses the active regulation, transcription, and translation

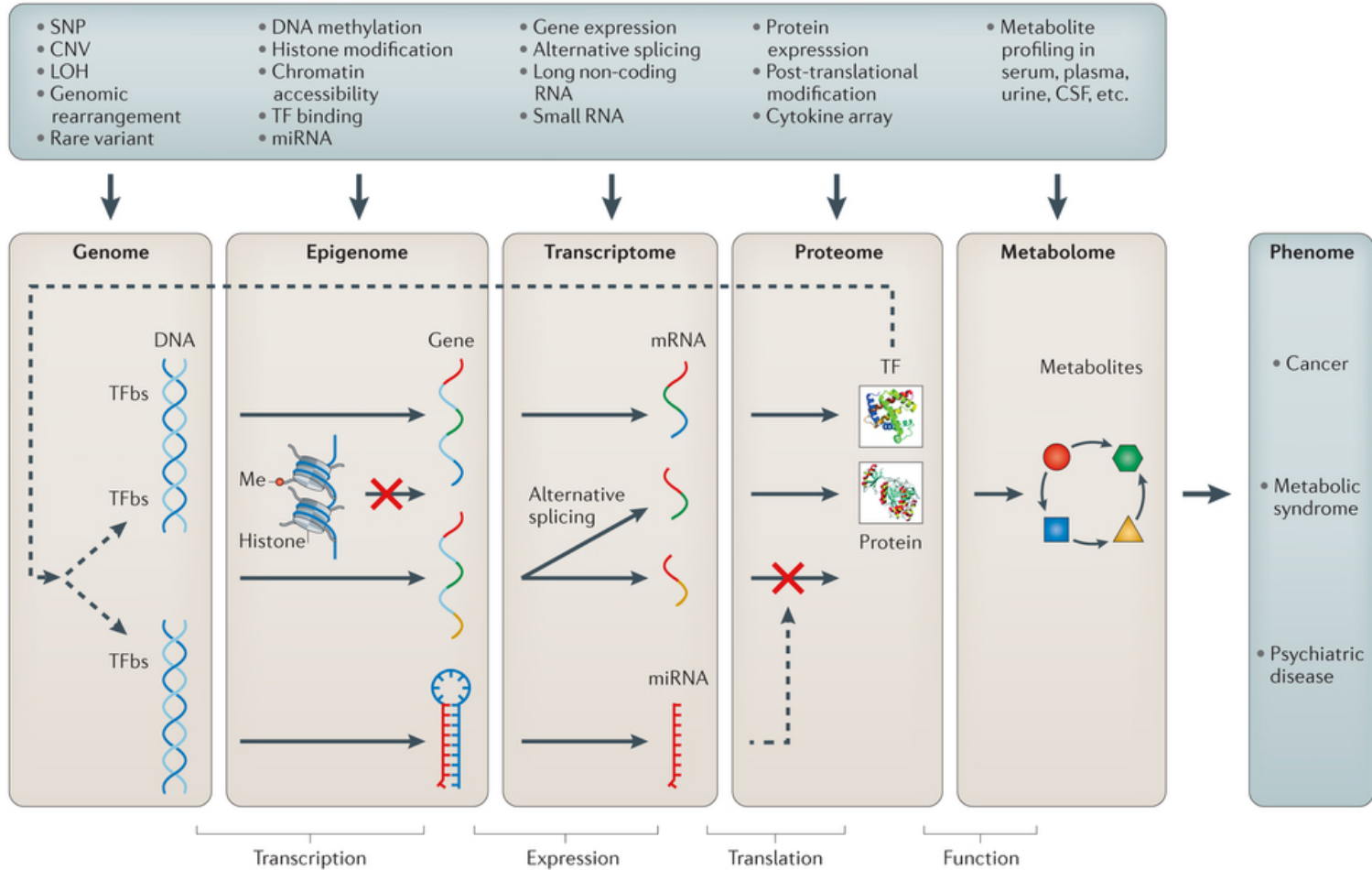


# Integrating environmental factors, genetic background, and large scale datasets

- 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.



# Different large-scale assays provide information on many types of biological regulation



## Expanding personalized medicine beyond the genome.

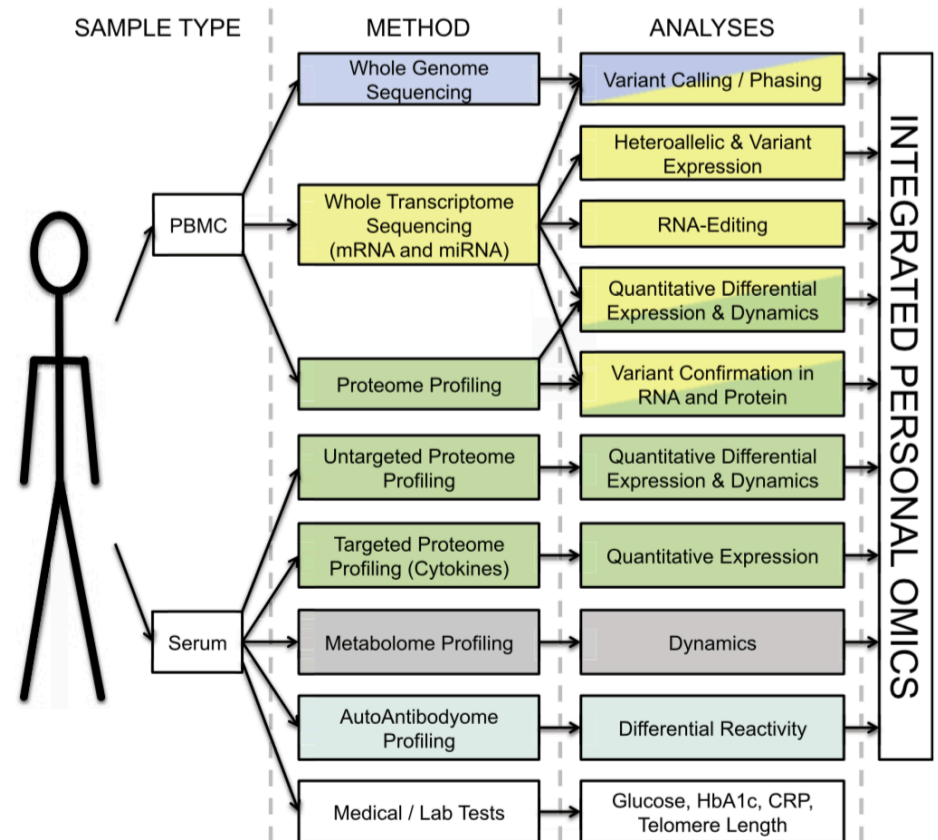
- An integrated personal omics profile (iPOP) is an example of a more comprehensive version of personalized medicine.
- Michael Snyder had his genome sequenced and collected many other large scale datasets over an extended period of time.





# 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



# Michael Snyder's personal genome: a starting point

**Table 1. Summary and Breakdown of DNA Variants**

Type	Total Variants	Total High Confidence	Heterozygous High Confidence	Homozygous High Confidence
Total SNVs	3,739,701	3,301,521	1,971,629	1,329,892
Total gene-associated SNVs	1,312,780	1,183,847	717,485	466,362
Total coding/UTR	49,017	44,542	27,383	17,159
Missense	10,592	9,683	5,944	3,739
Nonsense	83	73	49	24
Synonymous	11,459	10,864	6,747	4,117
5'UTR	4,085	2,978	1,802	1,176
3'UTR	22,798	20,944	12,841	8,103
Intron	1,263,763	1,139,305	690,102	449,203
Ts/Tv	—	2.14	—	—
dbSNP	3,493,748	3,167,180	—	—
Candidate private SNV	245,953	134,341	—	—
Indels (−107~ +36 bp)	1,022,901	216,776	—	—
Coding	3,263	302	—	—
Structural variants (>50 bp)	44,781	2,566	—	—
In 1000G project <sup>a</sup>	4,434	1,967	—	—

# Prioritizing variants by leveraging mutation databases

- Using existing databases of population level genetic variation, rare and disease associated variants could be identified.
- Helped prioritize medical conditions for monitoring (e.g. glucose for diabetes)

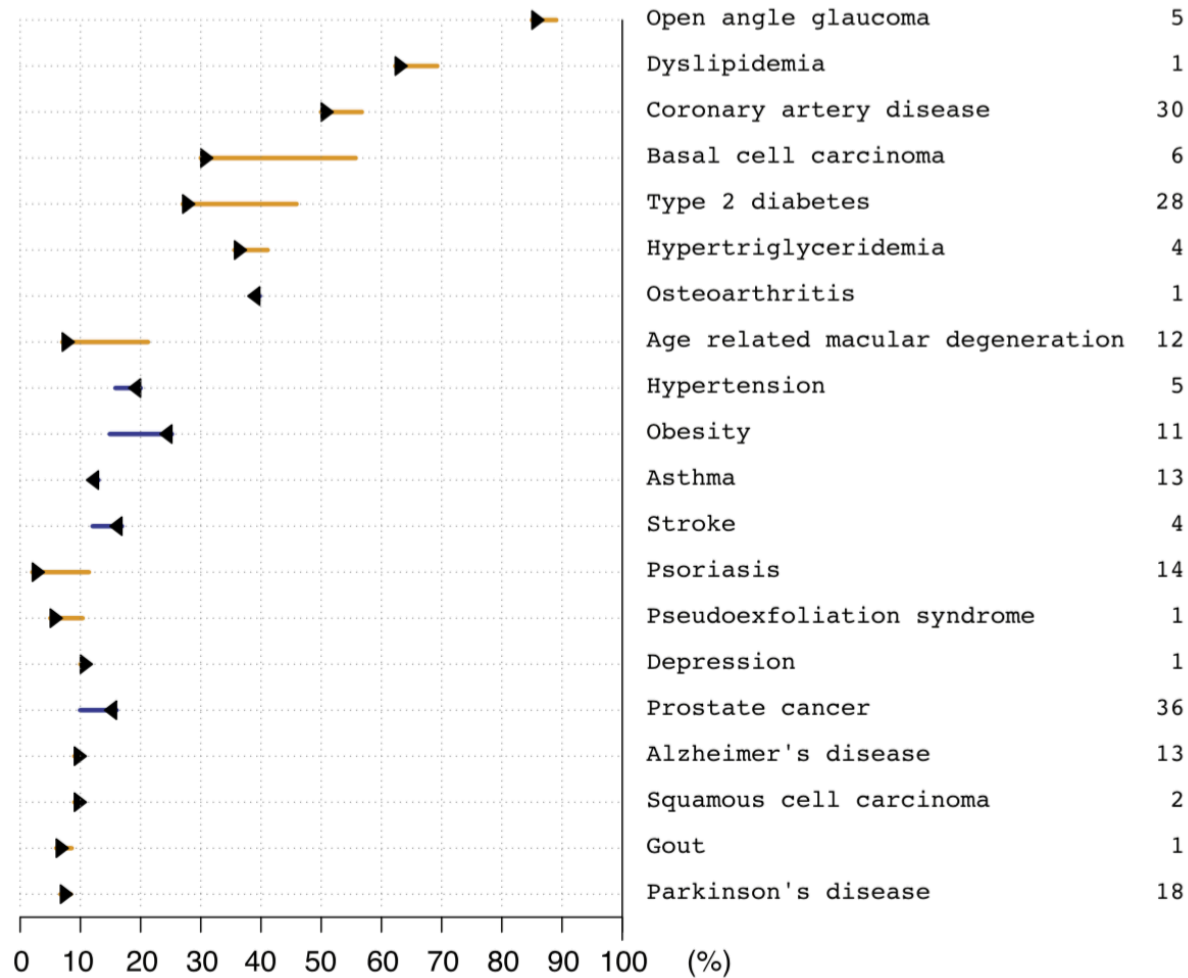
## High Interest Disease-Associated Rare Variants.

Gene	Position	Genotype	OMIM
SERPINA1	14:94844947	C/T	Emphysema due to AAT deficiency
TERT	5:1294397	C/T	Aplastic anemia
KCNJ11	11:17409571	T/T	Type 2 diabetes
GCKR	2:27730939	T/T	Hypertriglyceridemia
NUP54	4:77055431	G/A	Nuclear Pore Complex Protein

## High Interest Drug-Related Variants.

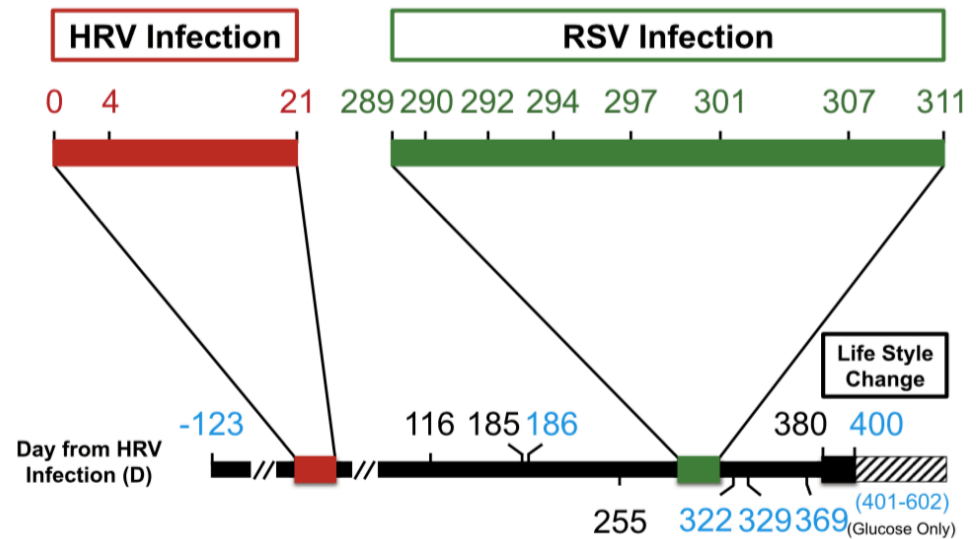
Gene	rsID	Genotype	Drug Response Affected
	rs10811661	C/T	Troglitazone (Increased Beta-Cell Function)
CYP2C19	rs12248560	C/T	Clopidogrel (Increased Activation)
LPIN1	rs10192566	G/G	Rosiglitazone (Increased Effect)
SLC22A1	rs622342	A/A	Metformin (Increased Effect)
VKORC1	rs9923231	C/T	Warfarin (Lower Dose Required)

# Genomic information helps refine disease risk estimates

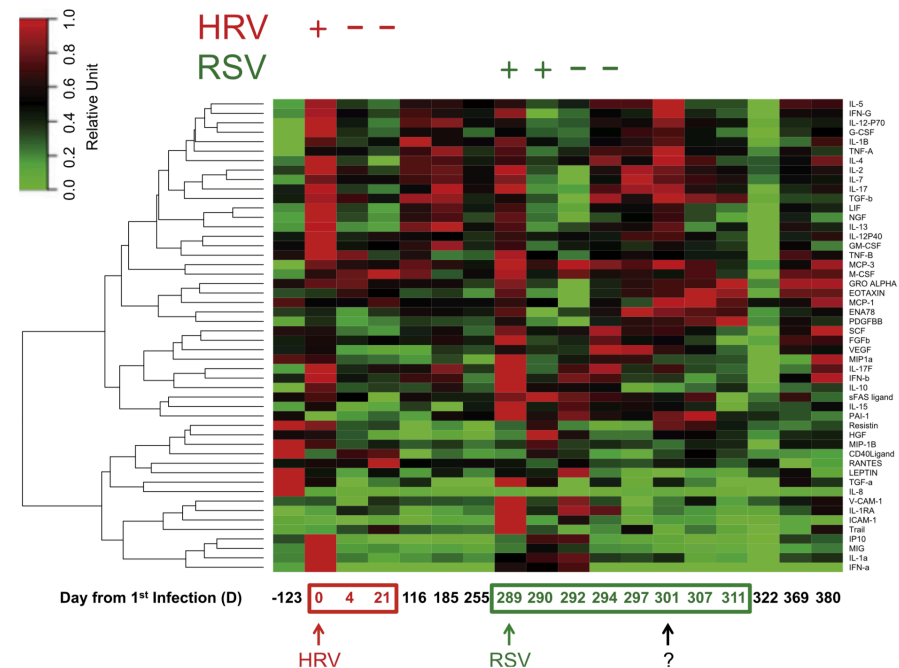
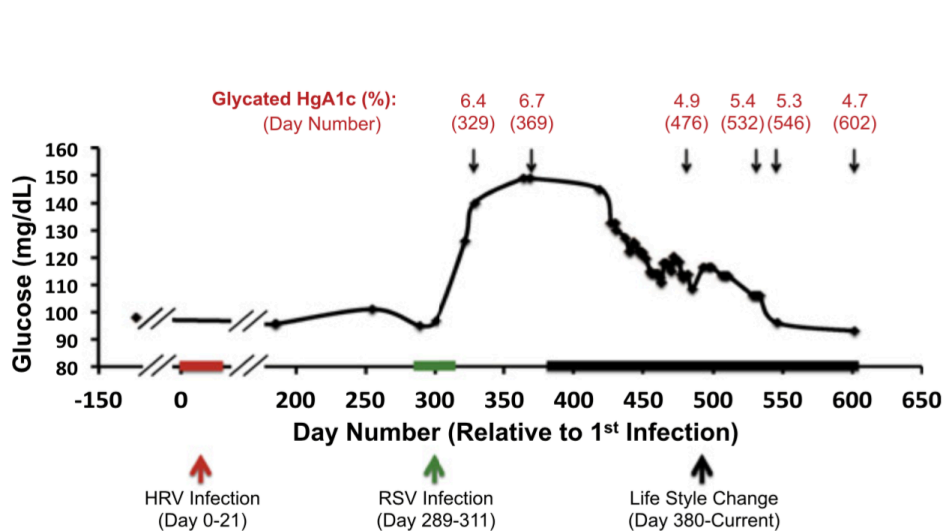


# Personal omics profiling time course

- Changing cellular state and environmental perturbations impact the genome.
- Longitudinal data collection tracks the dynamic regulation of the genome.



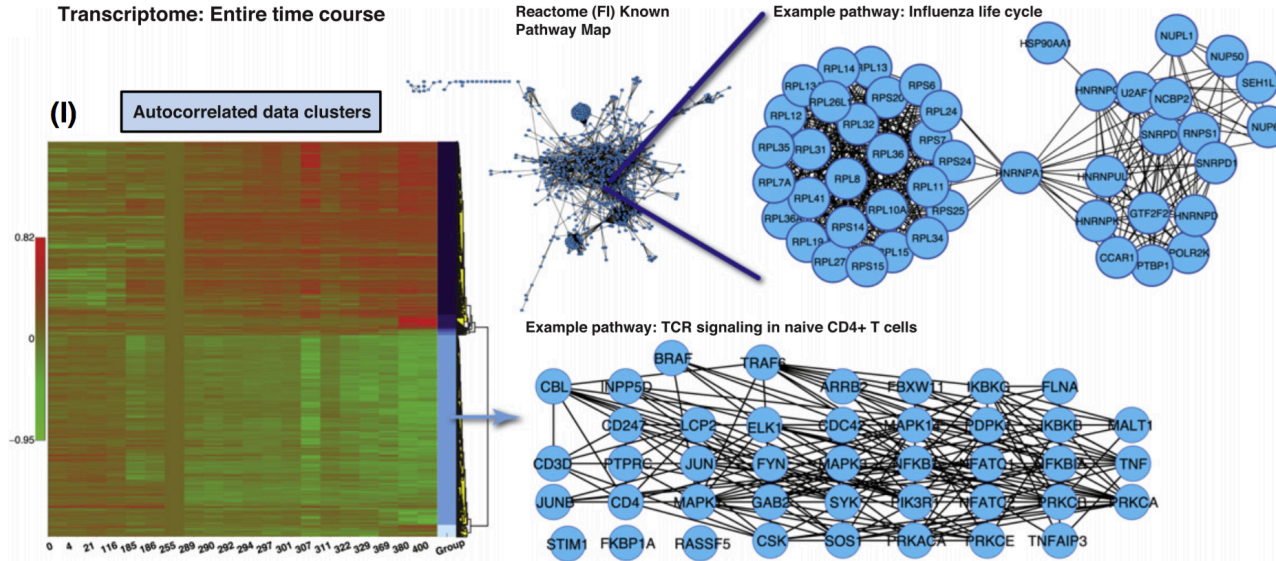
# Longitudinal medical data



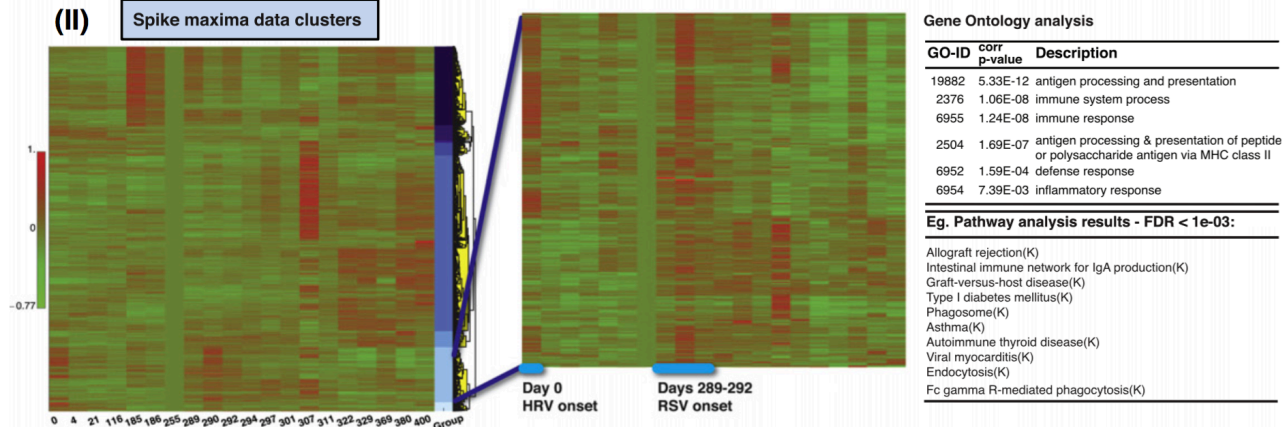
- Tracking relevant medical (e.g. blood glucose) data over time helps link phenotypic changes with changes at the molecular level.

# Transcriptomic time course

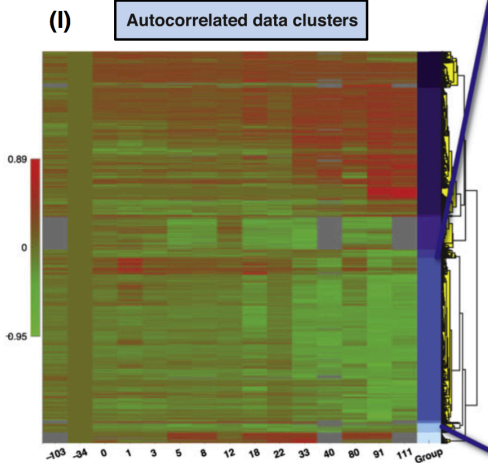
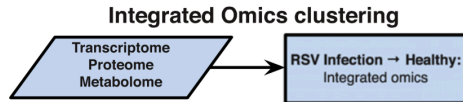
Transcriptome: Entire time course



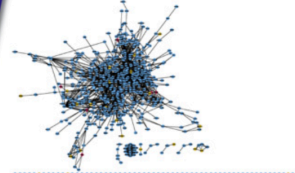
**(II) Spike maxima data clusters**



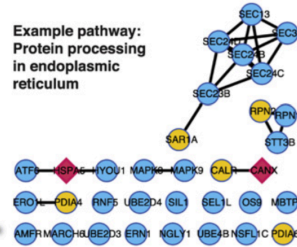
# Integration of multiple omics datasets



Full Reactome (FI) known pathway map for cluster:



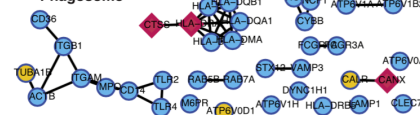
**Example pathway: Protein processing in endoplasmic reticulum**



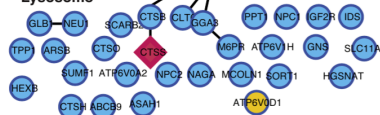
**Dynamic expression pattern observed in:**

- RNA (Blue circle)
- Protein (Yellow circle)
- Both RNA + Protein (Red diamond)

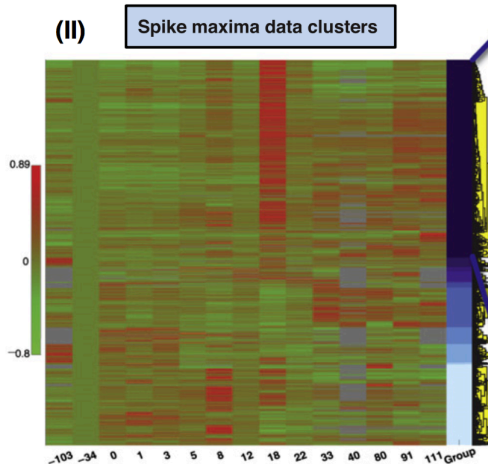
**Example pathway: Phagosome**



**Example pathway: Lysosome**



**Example pathway: Insulin**



**Eg. Pathway Analysis Results - FDR < 5e-02:**

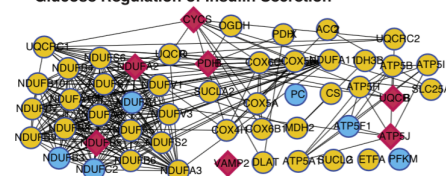
- Spliceosome(K)
- Glucose Regulation of Insulin Secretion(R)
- Formation and Maturation of mRNA Transcript(R)
- Oxidative phosphorylation(K)
- Electron Transport Chain(R)
- Parkinson's disease(K)
- Huntington's disease(K)
- Influenza Life Cycle(R)
- Metabolism of non-coding RNA(R)
- Transport of Mature Transcript to Cytoplasm(R)
- Protein export(K)
- Pyruvate metabolism and TCA cycle(R)

GO-ID	corr p-value	Description
8380	3.59E-71	RNA splicing
6396	2.53E-57	RNA processing
16070	5.93E-54	RNA metabolic process
16071	6.10E-50	mRNA metabolic process
10467	1.74E-48	gene expression
90304	1.78E-45	nucleic acid metabolic process

**Eg. Metabolites in Cluster**

- 3R-hydroxy-5Z-dodecenoic acid
- 5,6-DIHEHE-EA
- 7-Ethoxycoumarin
- Lauric acid
- 1-O-(1Z-tetradecenyl)-2-(9Z-octadecenyl)-sn-glycerol (23R)-1alpha,23,25-trihydroxy-24-oxovitamin D3 / (23R)-1alpha,23,25-trihydroxy-24-oxocholecalciferol
- 1alpha-hydroxy-26,27-dinorvitamin D3 25-carboxylic acid / 1alpha-hydroxy-26,27-dinorcholecalciferol
- 12-oxo-9-octadecynoic acid
- GPCho(O-16:0/O-4:0[U])
- 19-hydroxy-17-oxoandrost-5-en-3-beta-yl sulfate - 11.538899

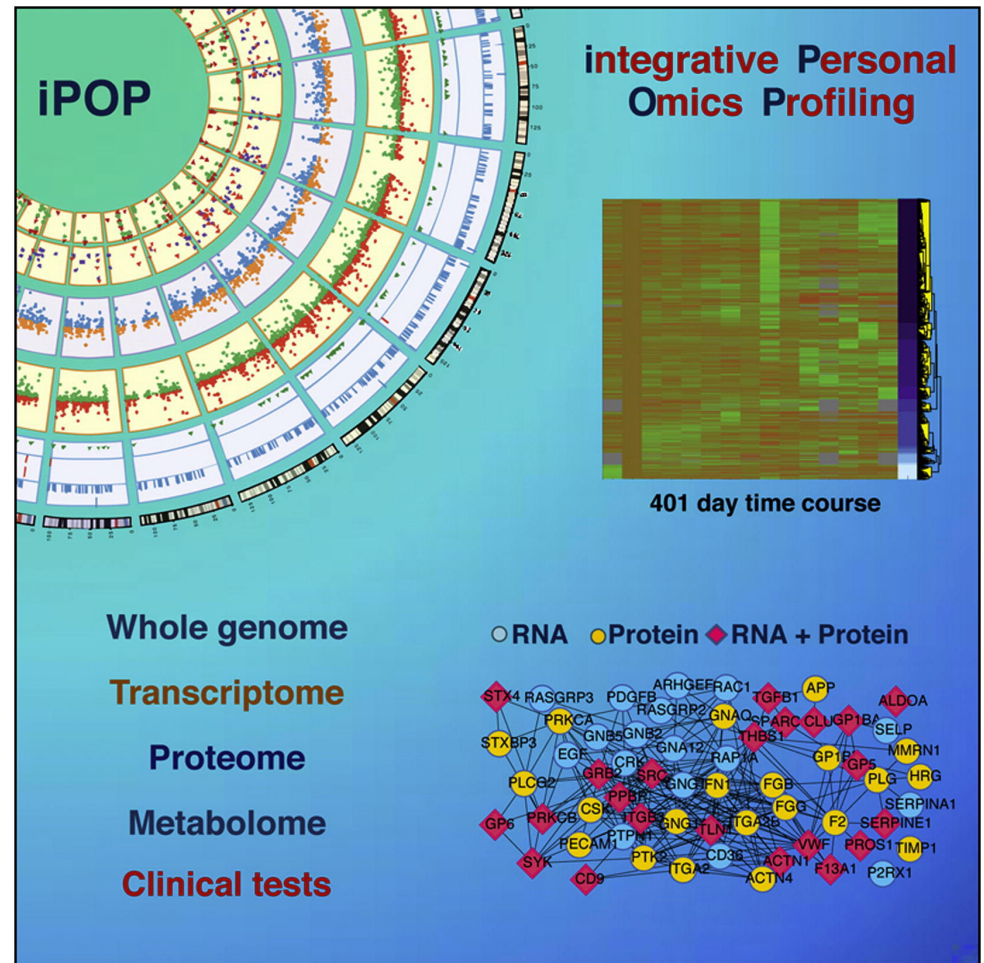
**Example Pathway: Glucose Regulation of Insulin Secretion**





# Integrated personal omics profile

- iPOP: Longitudinal study integrating multiple large-scale datasets.
- Recording medical and molecular data helps reveal molecular underpinnings of health and disease.



# Personal Genome Project

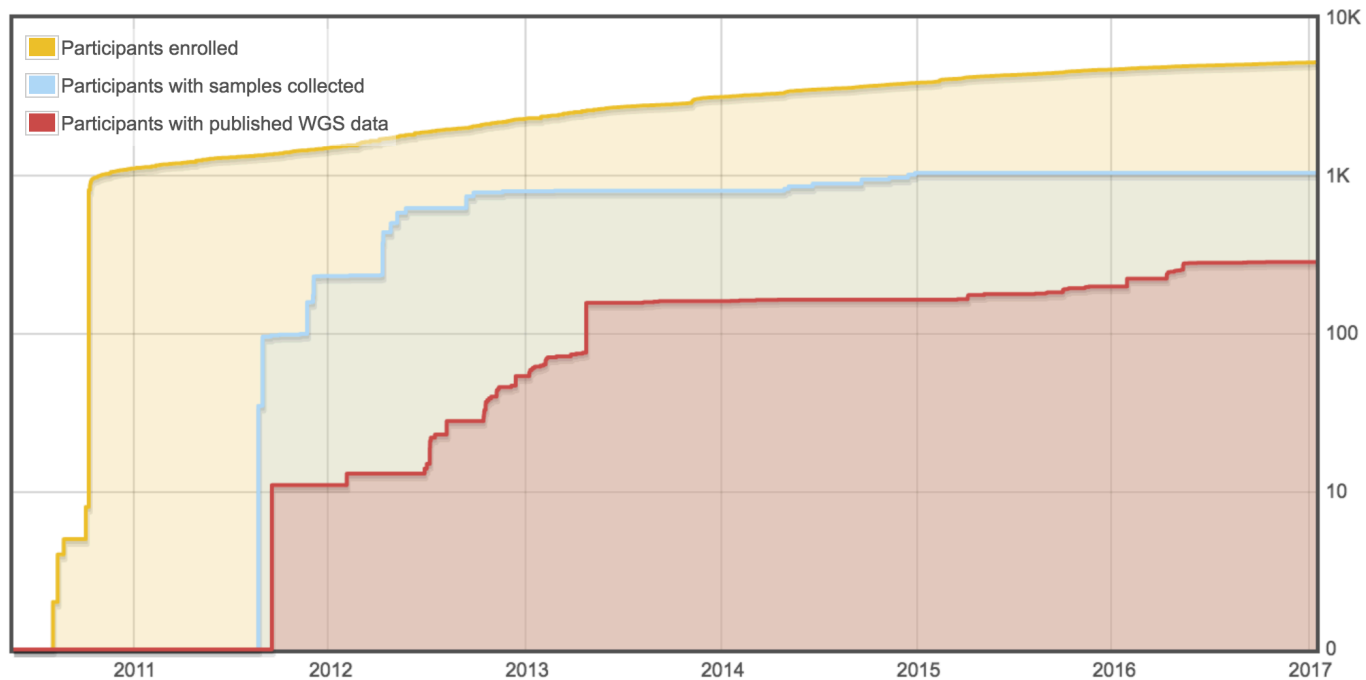
## Sharing Personal Genomes

The Personal Genome Project was founded in 2005 and is dedicated to creating public genome, health, and trait data. Sharing data is critical to scientific progress, but has been hampered by traditional research practices—our approach is to invite willing participants to publicly share their personal data for the greater good.



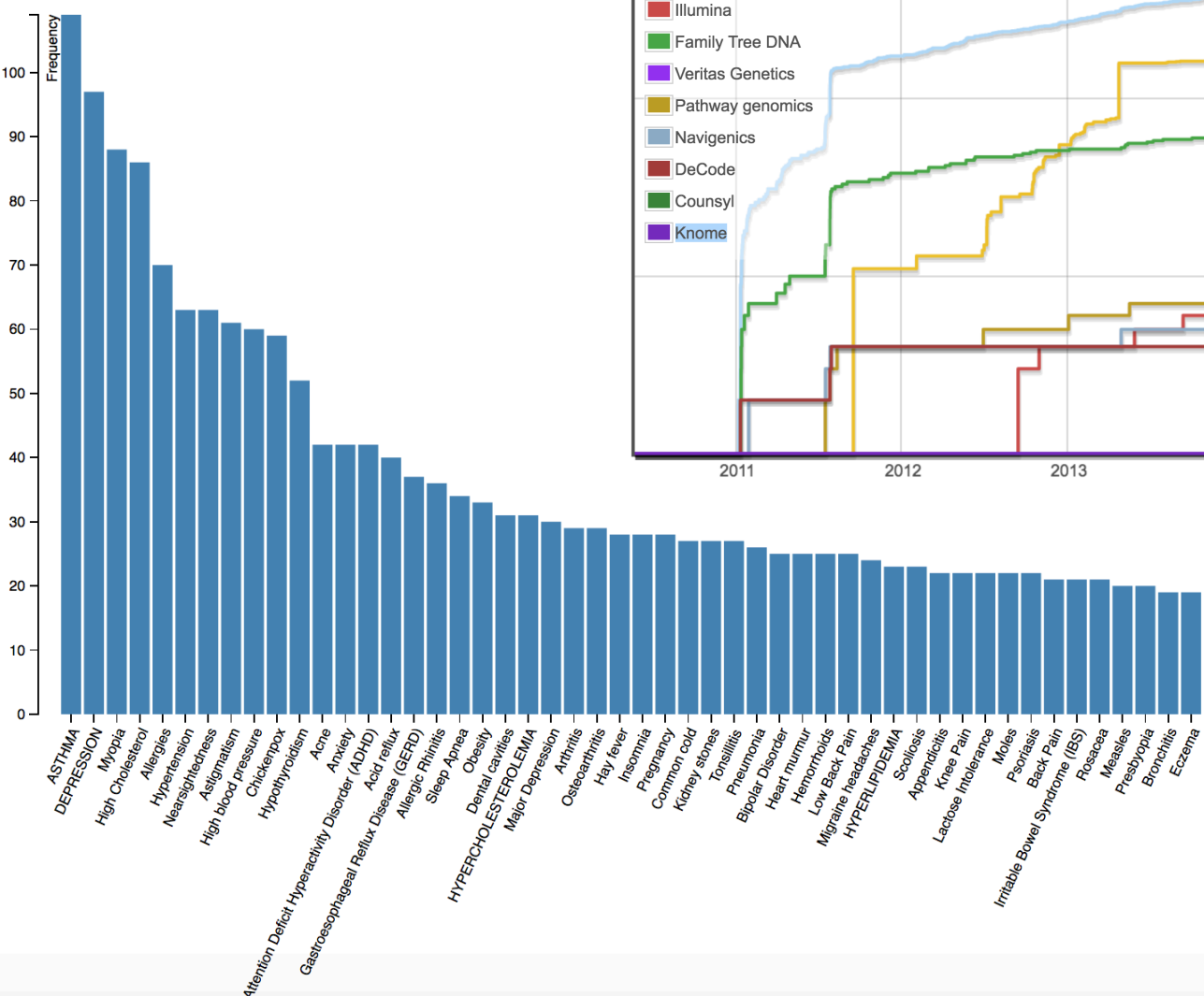
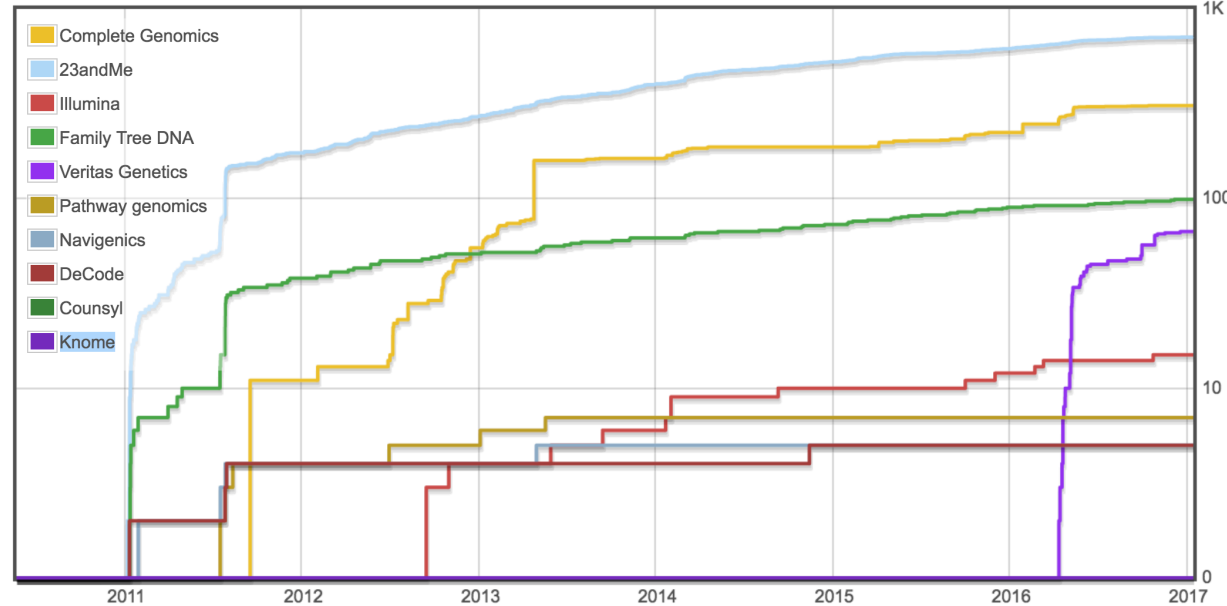
[Learn more >](#)

Pipeline: enrolled → samples collected → WGS data published



# Data Types in the Personal Genome Project

Number of participants per data type



Conditions

# Precision medicine in the clinic

- Increasingly genomic information is playing a role in the clinic.
  - Targeted therapeutics
  - Pharmacogenomics
    - informing treatment options based on patient drug sensitivity

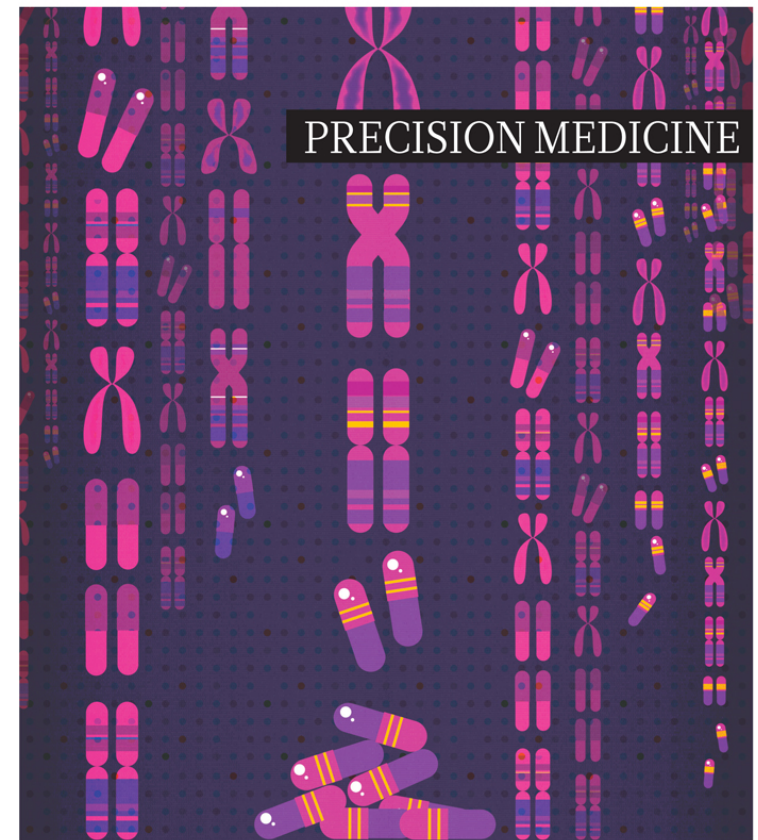
Table 1 | Examples of precision medicine

Condition	Gene	Action
<b>Mendelian disease</b>		
Cystic fibrosis	CFTR	Specific therapies such as ivacaftor and a combination of lumacaftor and ivacaftor
Long QT syndrome	KCNQ1, KCNH2 and SCN5A	Specific therapy for patients with SCN5A mutations
Duchenne muscular dystrophy	DMD	Ongoing phase III clinical trials of exon-skipping therapies
Malignant hyperthermia susceptibility	RYR1	Avoid volatile anaesthetic agents; avoid extremes of heat
Familial hypercholesterolaemia (FH)	PCSK9, APOB and LDLR	<ul style="list-style-type: none"> <li>Heterozygous FH (HeFH): eligible for PCSK9 inhibitor drugs</li> <li>Homozygous FH (HoFH): eligible for PCSK9 inhibitor drugs in addition to lomitapide and mipomersen</li> </ul>
Dopa-responsive dystonia	SPR	Therapy with dopamine precursor L-dopa and the serotonin precursor 5-hydroxytryptophan
Thoracic aortic aneurysm	SMAD3, ACTA2, TGFB1, TGFB2 and FBN1	Customization of surgical thresholds based on patient genotype
Left ventricular hypertrophy	MYH7, MYBPC3, GLA and TTR	Sarcomeric cardiomyopathy, Fabry disease and transthyretin cardiac amyloid disease have specific therapies
<b>Precision oncology</b>		
Lung adenocarcinoma	EGFR and ALK	Targeted kinase inhibitors, such as gefitinib and crizotinib
Breast cancer	HER2	HER2 (also known as ERBB2)-targeted treatment, such as trastuzumab and pertuzumab
Gastrointestinal stromal tumour	KIT	Targeted KIT kinase activity inhibitors, such as imatinib
Melanoma	BRAF	BRAF inhibitors, such as vemurafenib and dabrafenib
<b>Pharmacogenomics</b>		
Warfarin sensitivity	CYP2C9 and VKORC1	Adjust dosage of warfarin or consider alternative anticoagulant
Clopidogrel sensitivity, post-stent procedure	CYP2C19	Consider alternative antiplatelet therapy (for example, prasugrel or ticagrelor)
Thiopurine sensitivity	TPMT	Reduce thiopurine dosage or consider alternative agent
Codeine sensitivity	CYP2D6	Avoid use of codeine; consider alternatives such as morphine and non-opioid analgesics
Simvastatin sensitivity	SLCO1B1	Reduce dose of simvastatin or consider an alternative statin; consider routine creatine kinase surveillance

## Precision medicine in the clinic

- Precision medicine is leading to better defining and treating disease at the molecular level.
- It is both changing the prescription of existing medications and inspiring new targeted therapies.
- Precision medicine requires high quality patient genome sequences be obtained at reasonable cost.

nature**OUTLOOK**



Produced with support from:

illumina®

A personal approach  
to health care