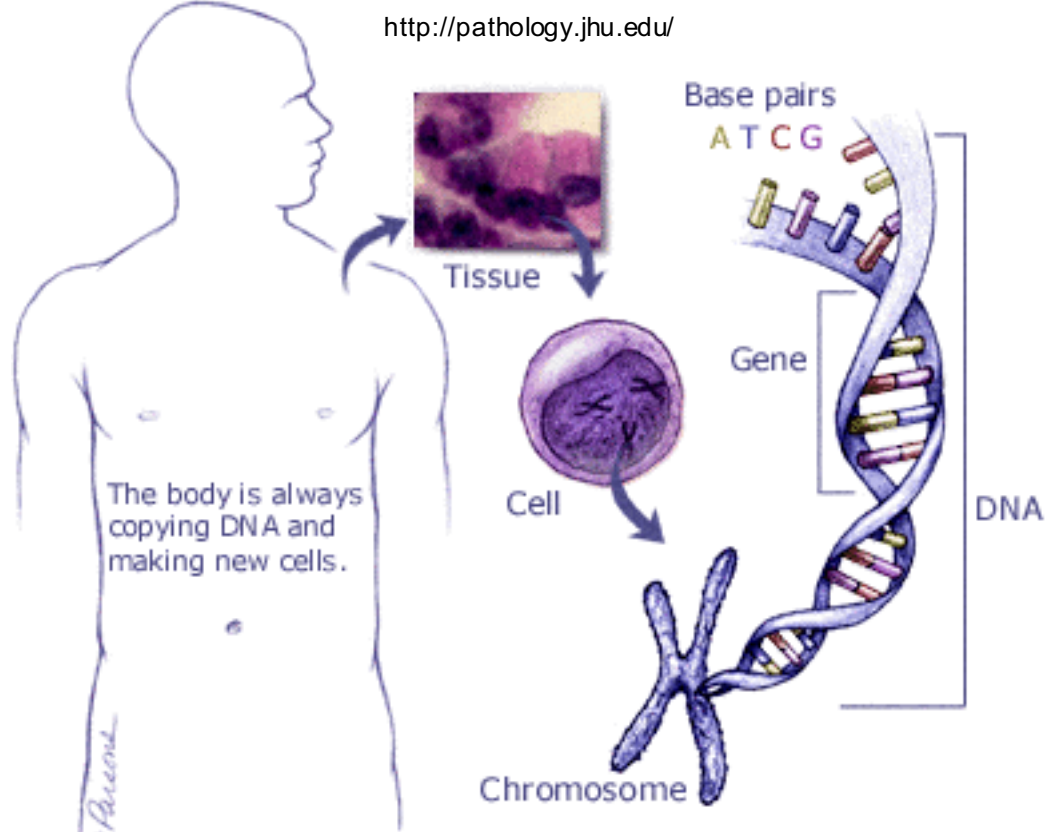


# What we learn from “clinically”-deep 10,000 genomes?

Ahmed Moustafa

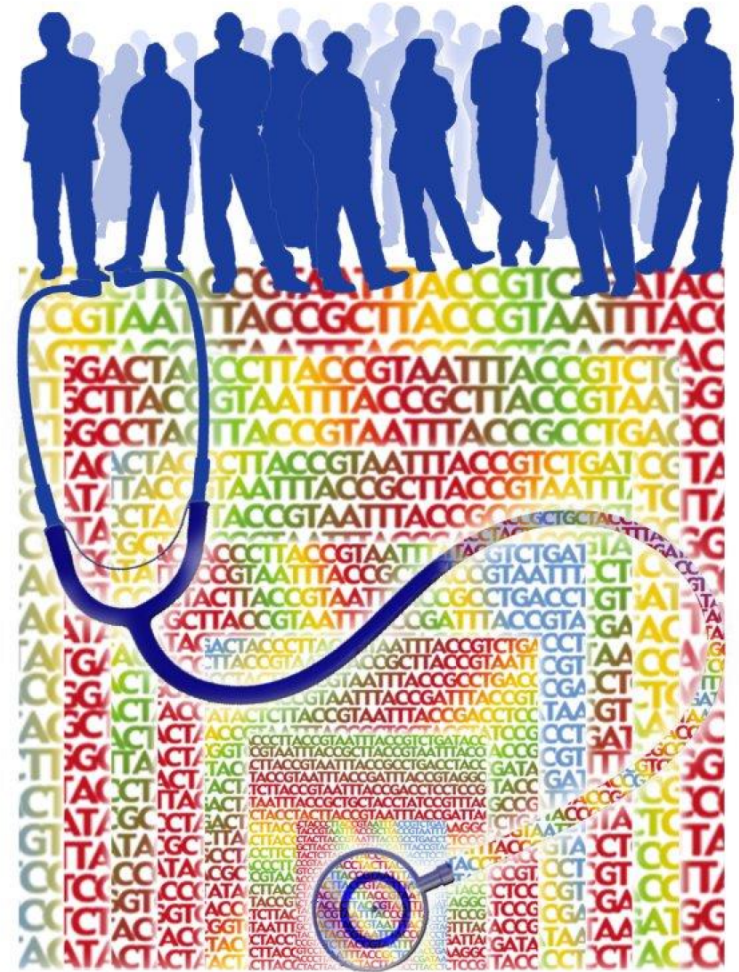
Human Longevity, Inc. (HLI)  
American University in Cairo (AUC)

BioVision Alexandria 2016  
Bibliotheca Alexandrina

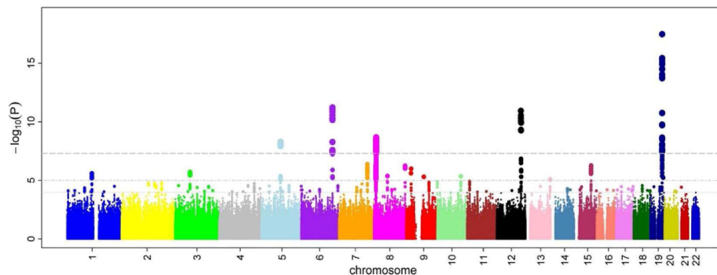


To identify such genetic changes (mutations, variants, polymorphism) that affect our “**phenotype**”, we need to “**genotype**”.

To associate w/ **statistical power** the genotype to the phenotype, we need to genotype “**populations**”.



<https://www.genome.gov/>



GWAS standard:  
 $p\text{-value} < 5e-8$

[https://en.wikipedia.org/wiki/Manhattan\\_plot](https://en.wikipedia.org/wiki/Manhattan_plot)

# Genotyping Approaches

Illumina  
HumanOmni

- **SNP Array** – with ~2.5 million markers (e.g., Illumina) covers only < 0.01% of the human genome. SNP arrays are designed based on known markers from published specific population studies → SNP bias
- **Whole-Exome Sequencing (WES)** – covers only the coding component of the genome (~2%). However, most disease associated SNPs from genome-wide association studies (GWAS) are non-coding
- **Whole-Genome Sequencing (WGS)**



# Public Genome Projects

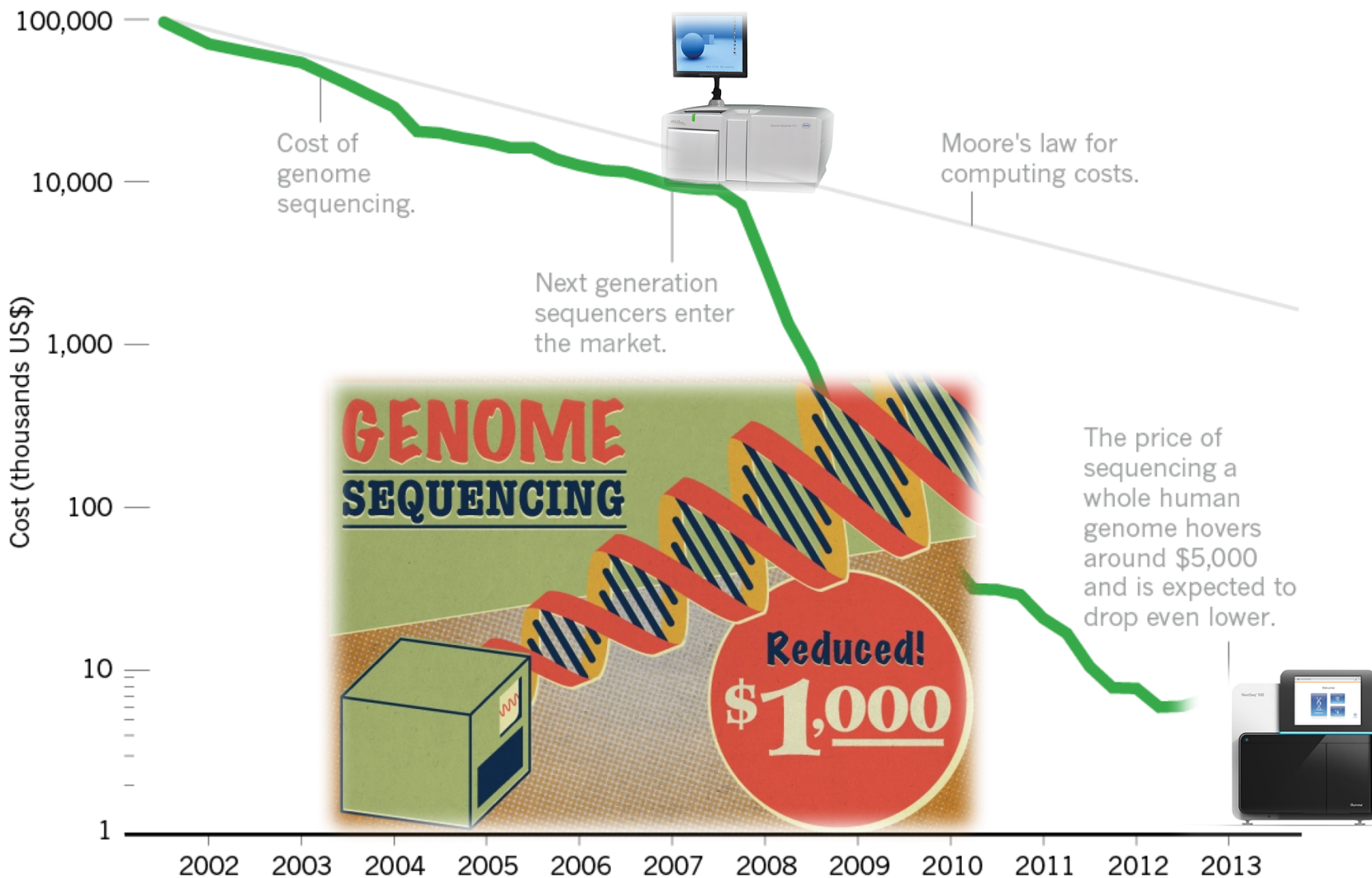
<b>Study</b>	<b>No.</b>	<b>Average Coverage</b>	<b>No. Deep Coverage</b>
<b>1000 Genomes 2015</b>	2,504	7.4x	453
<b>Japan Genome 2015</b>	1,070	32.4x	1,070
<b>Iceland Genome 2015</b>	3,545	20x	909
<b>UK Genome 2015</b>	3,781	7x	0
<b>African Genome 2015</b>	320	4x	0
<b>Sardinian Genome 2015</b>	2,120	4x	0
<b>Netherland Genome 2014</b>	750	13x	0
<b>Total</b>	<b>14,090</b>		<b>2,432</b>

# **The Clinical Genome**

- In Depth Coverage
- In Breadth Coverage
- Quality and Reproducibility

# Falling fast

In the first few years after the end of the Human Genome Project, the cost of genome sequencing roughly followed Moore's law, which predicts exponential declines in computing costs. After 2007, sequencing costs dropped precipitously.



# Workflow



HiSeq X Ten



72hrs

Sequencing (40 flowcells)

4

Packaging (9 TB/day)

4

Upload (650 GB/sample)



BCL to FASTQ

12

FASTQ to BAM (aligner)

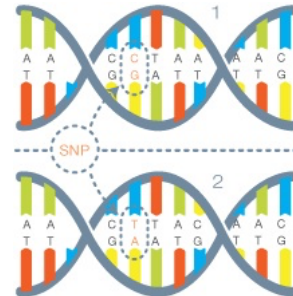
25hrs

∞

BAM to (g)VCF (variant caller)

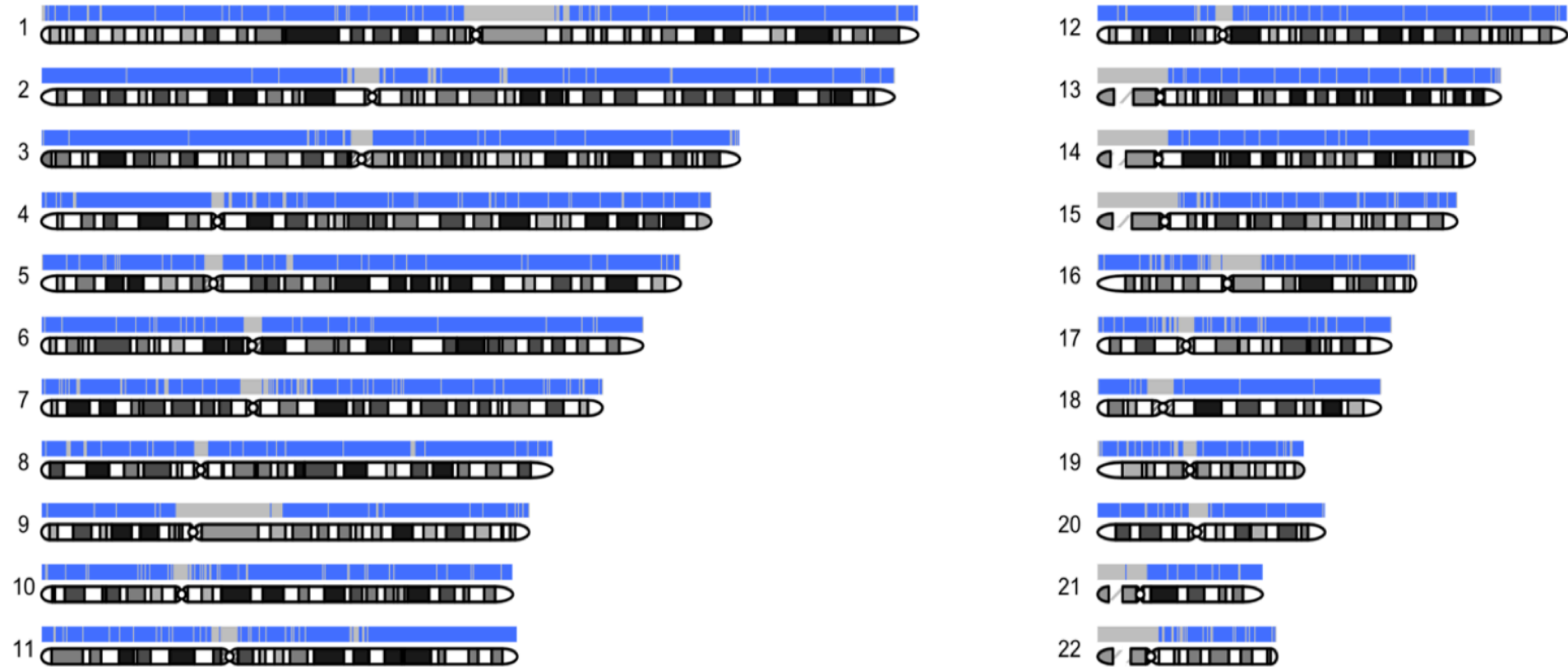
Permanent storage (S3)

Downstream analysis



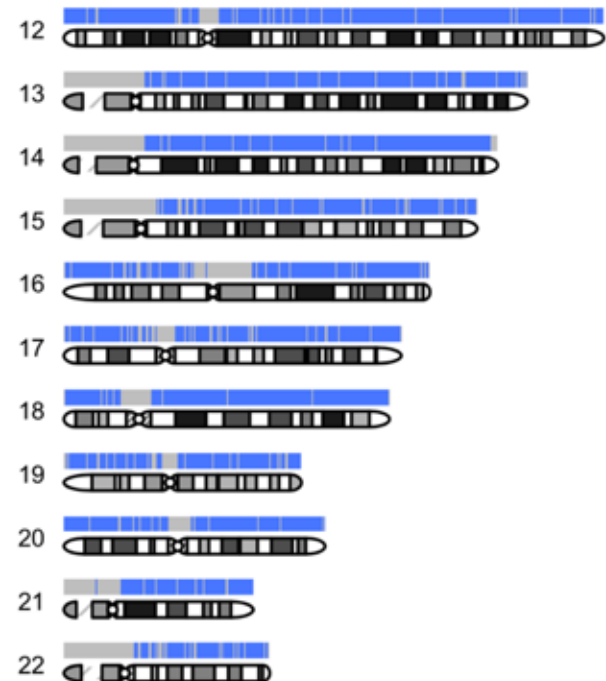
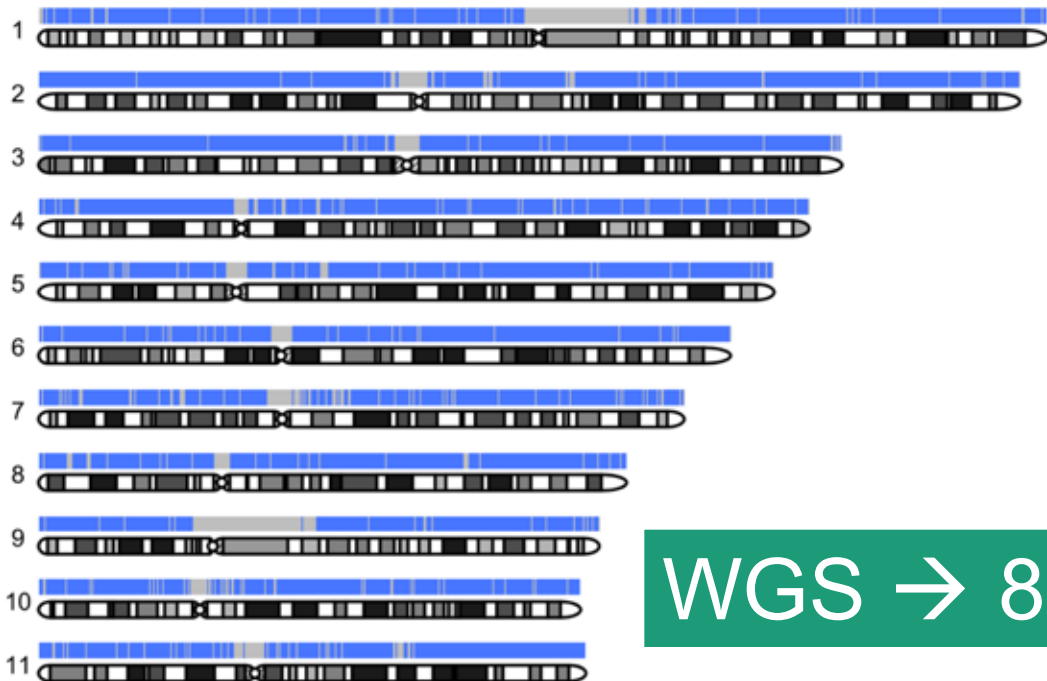


# Overall Quality Metrics



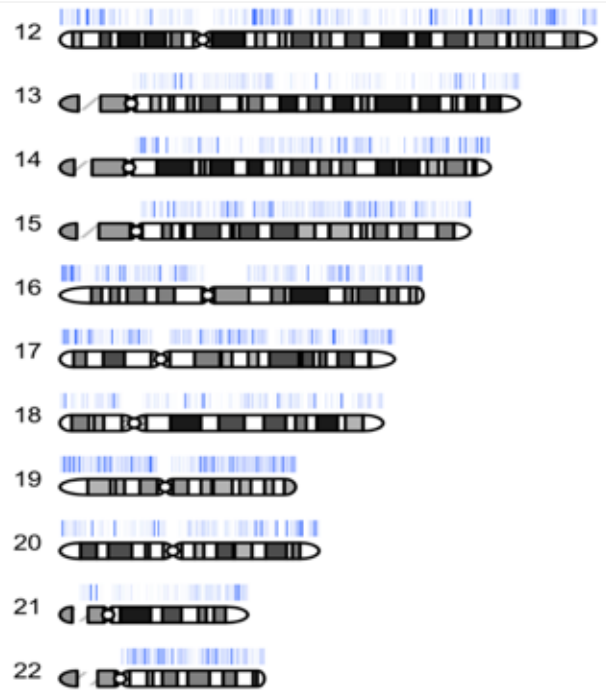
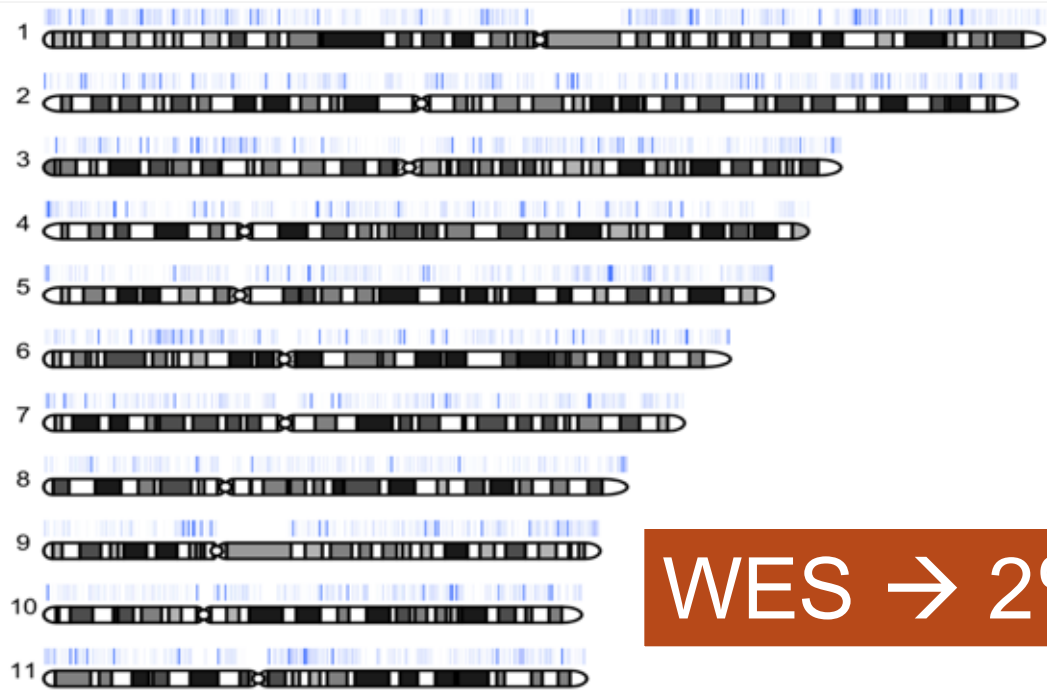
30x-40x Depth | 85% Coverage (Blue) |  
92% of Exome | 96% of Clinical Variants

Whole-Genome Sequencing



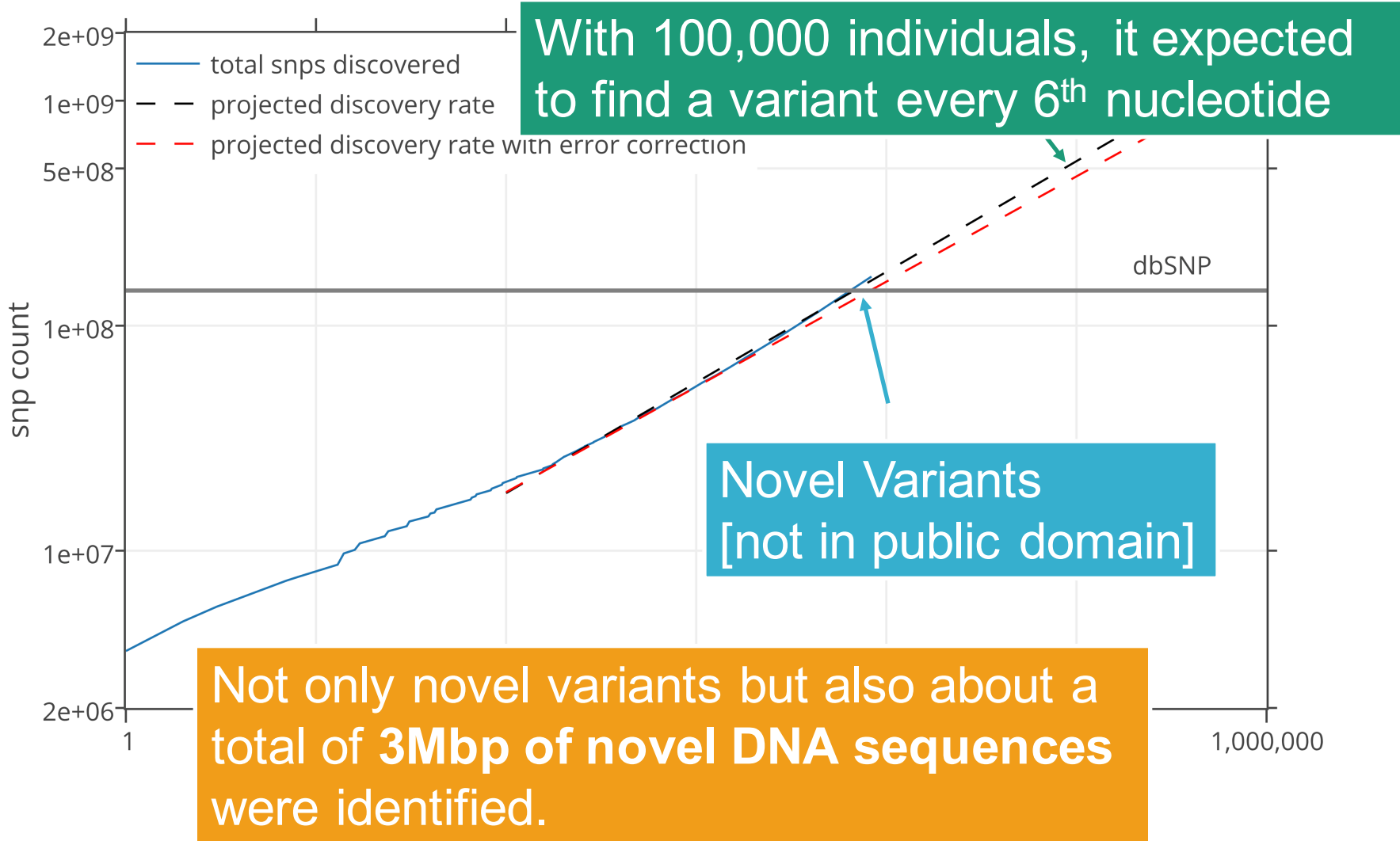
WGS → 85%

Whole-Exome Sequencing

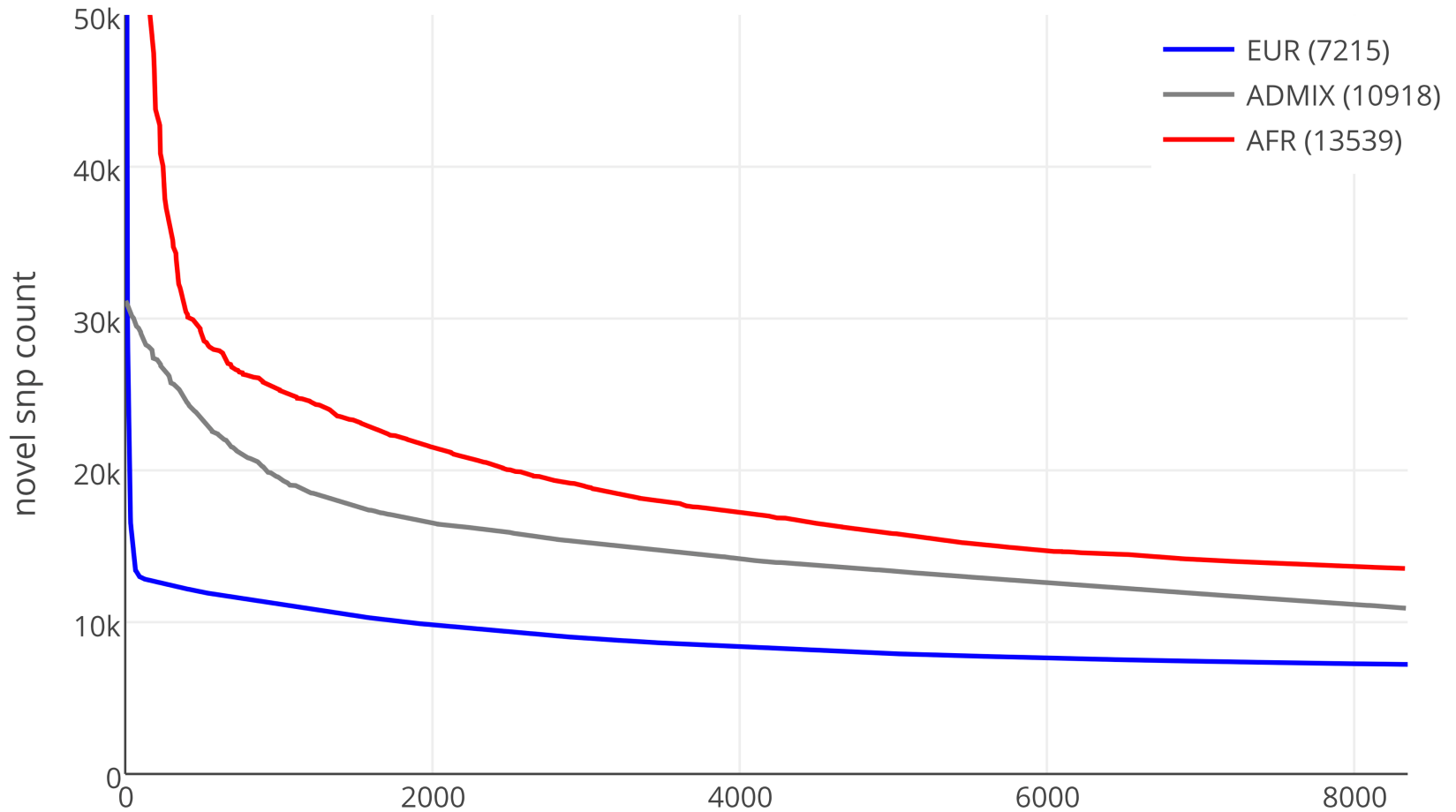


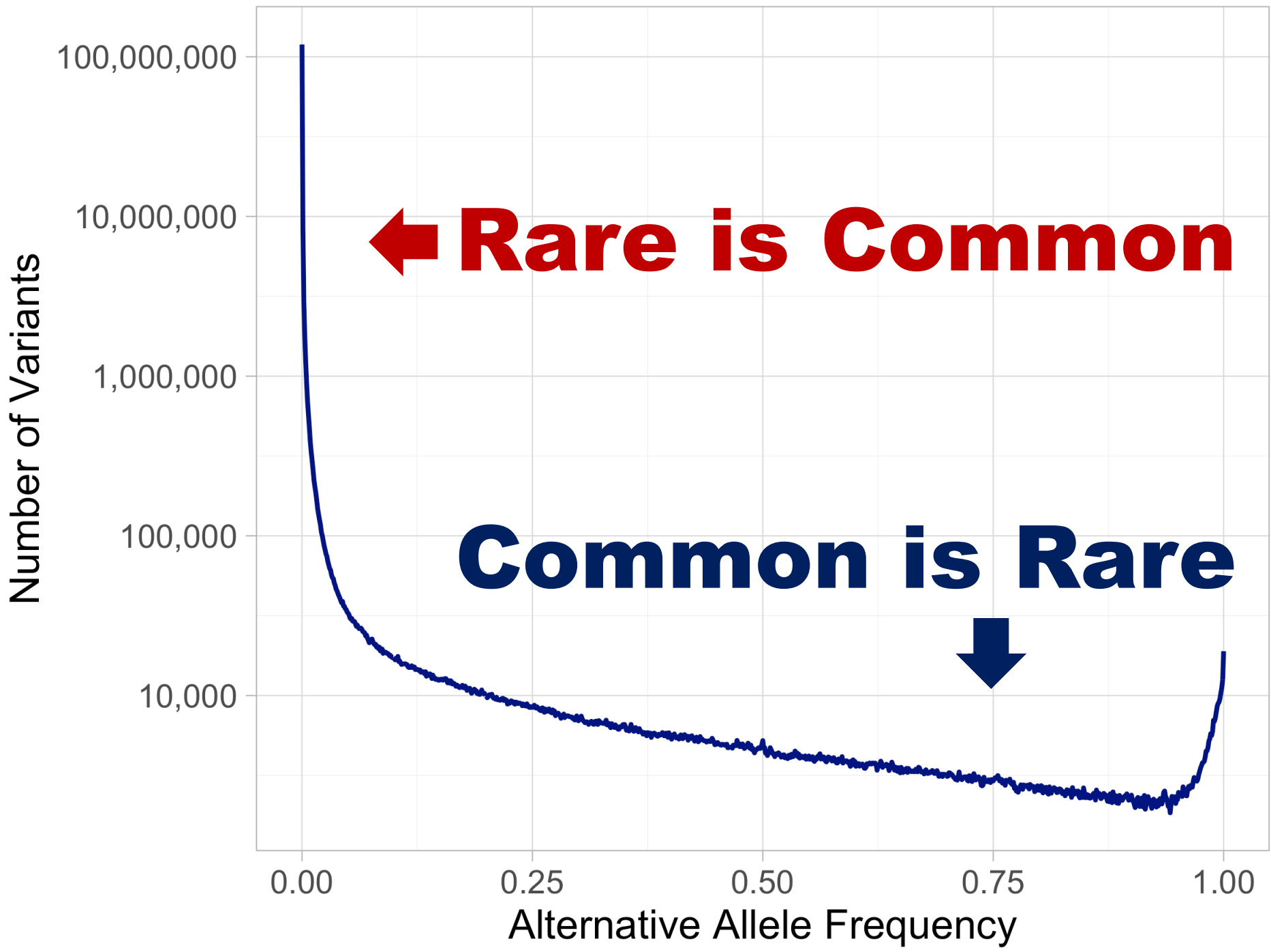
WES → 2%

# 150 Million Variants



# Novel Variants by Individual

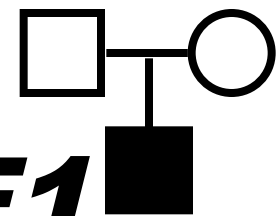




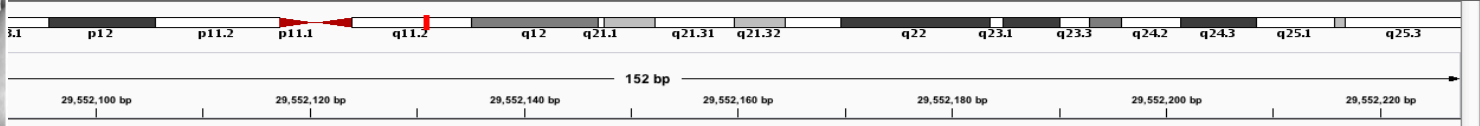
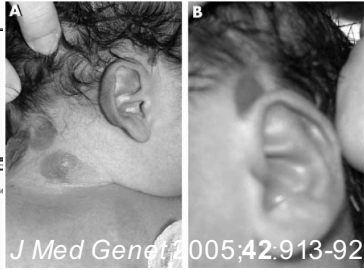
**← Rare is Common**

**Common is Rare**

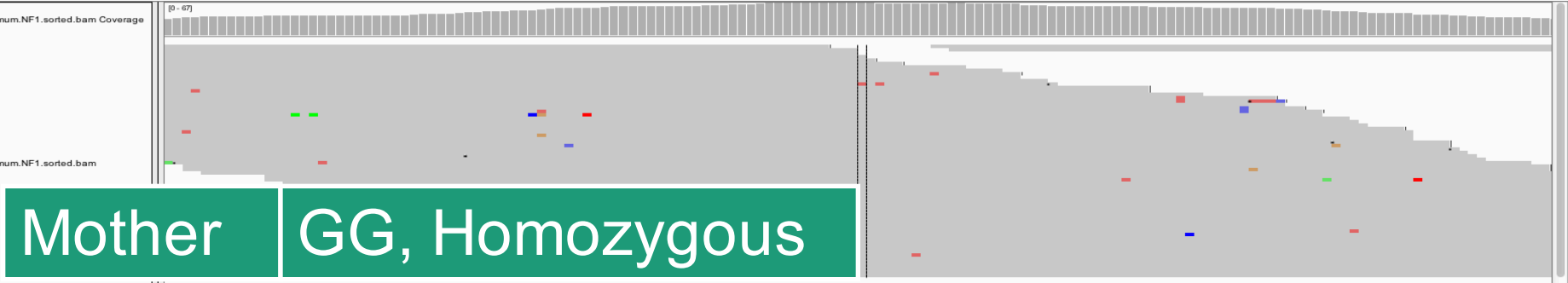




# De novo mutation in *NF1*



**Affected** GA, Heterozygous



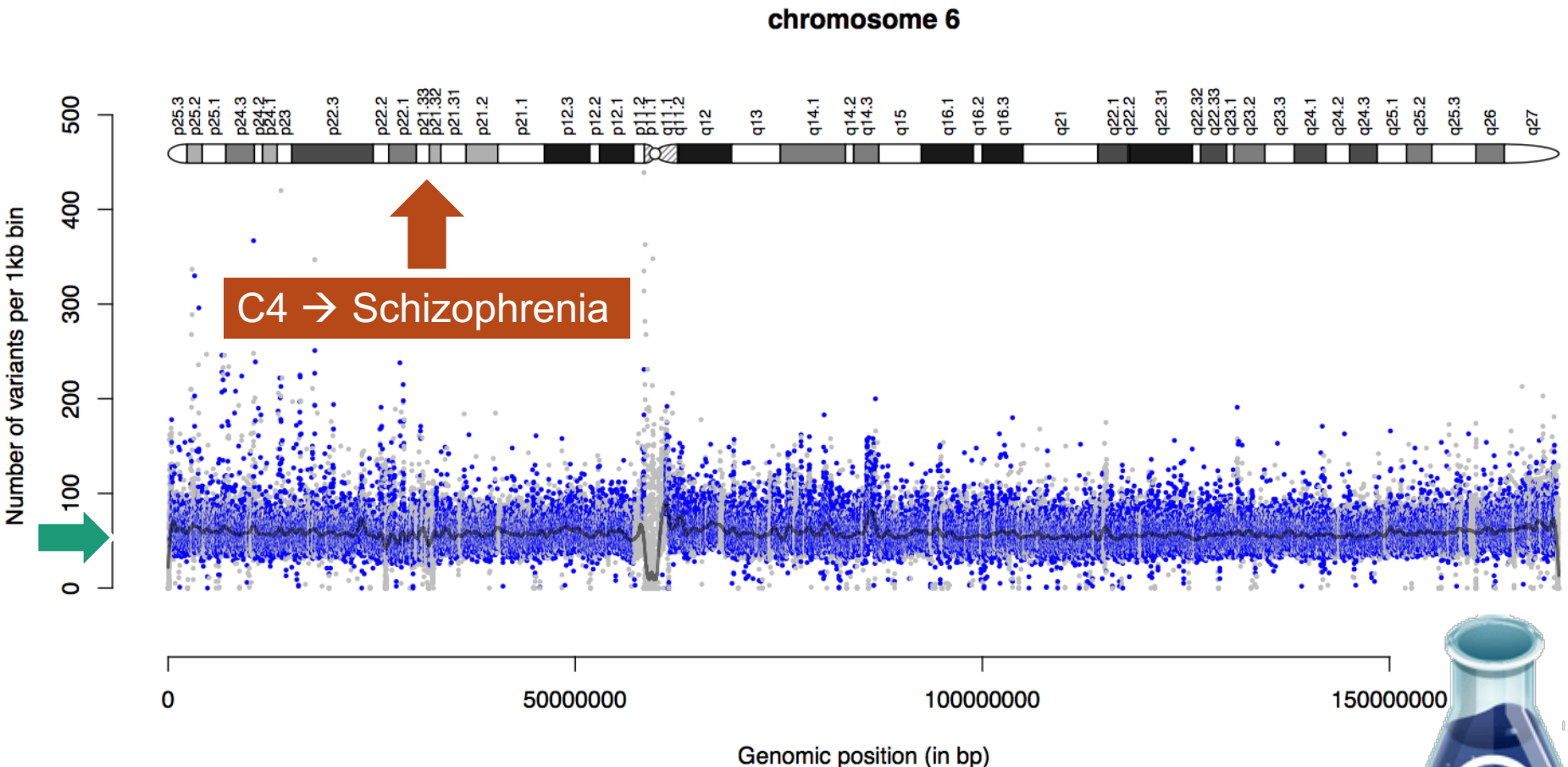
Sequence  
 RefSeq Genes  
 ...GTGCTTCAGTAAAGCTTATTTATTTATTTTCTAGCAGGCAGATAGAAGTTCCGTGCACTTCTCCTTTTATACGGGGTAGGATGATGATATTCCTTCTAGTGGAATACCAAGTCAAATGTCATGGATCATGAAGAATTACTACGTACTC...  
 ...QADRSSCHFLLLFYGVGCDIPSSGNTSQMSMDHEELLRT...

*NF1* → neurofibromin (tumor suppressor)  
 Delleman syndrome (congenital)

Mutation from G → A, only in the son, introduces a splice acceptor leading to skipping the beginning of an exon.

# **Rates of Variation**

# Rates of Variation in 10,000 Genomes

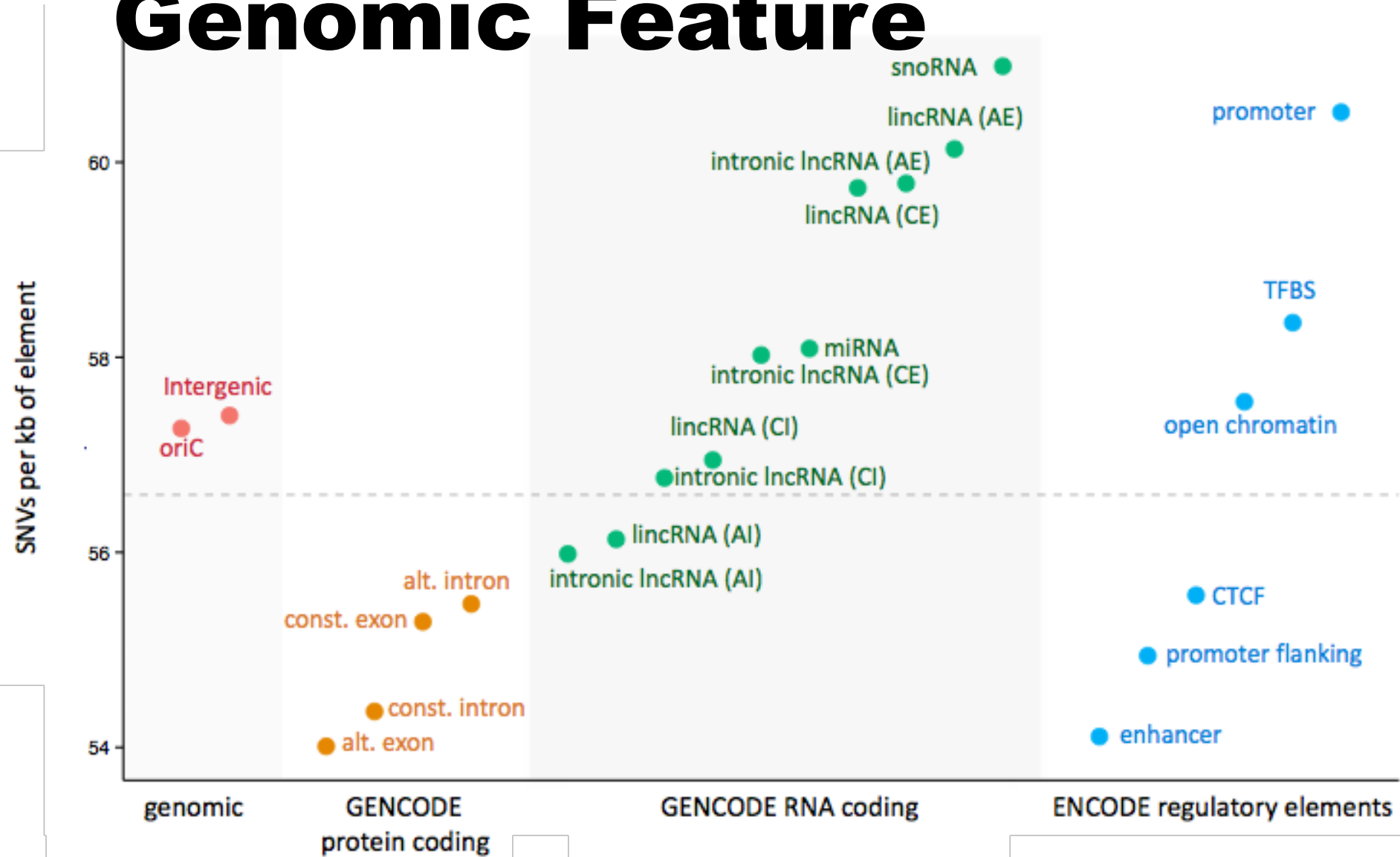


Blue → Overlap w/ Genome-in-a-Bottle (GIAB) high-confidence  
Grey → Outside of GIAB

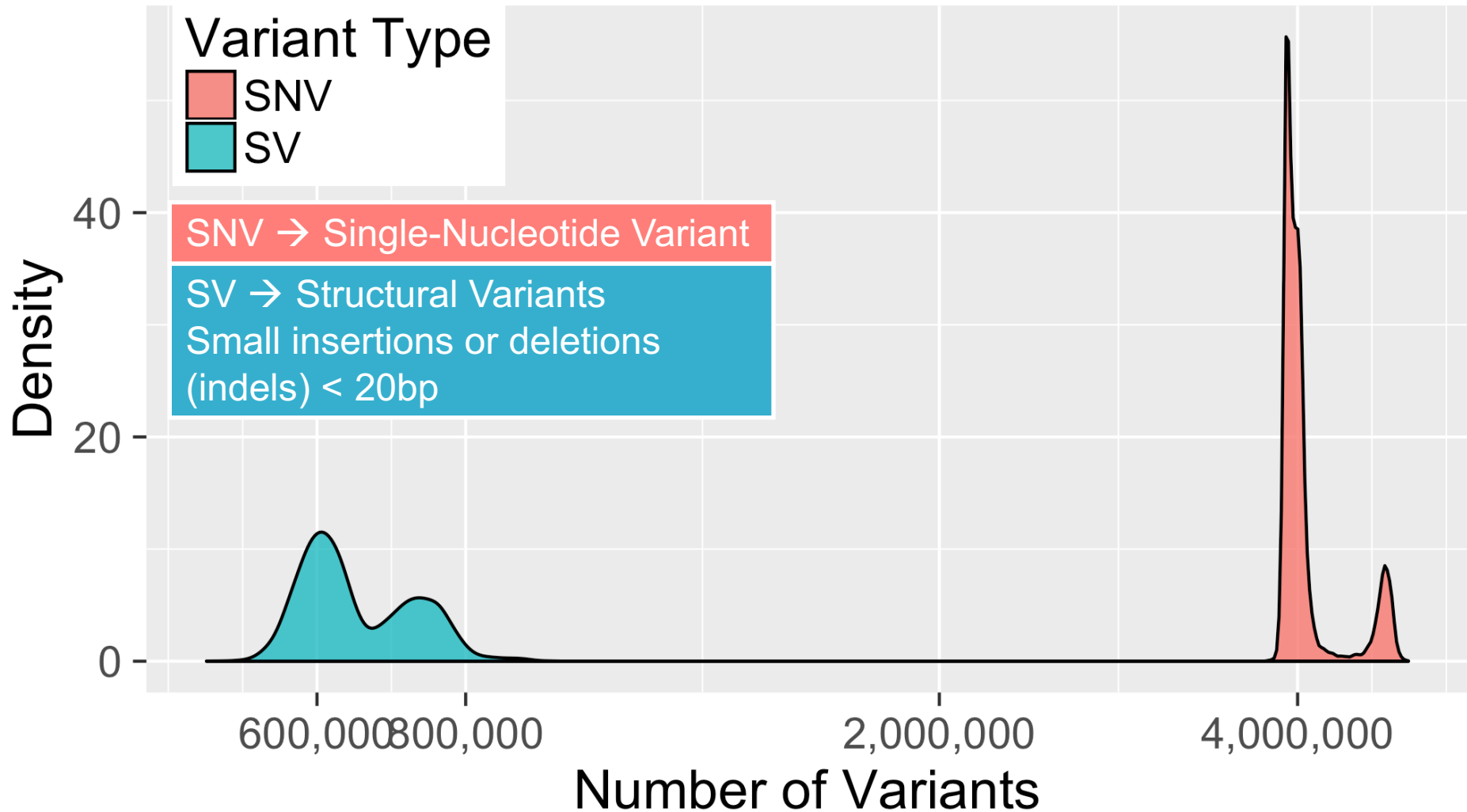




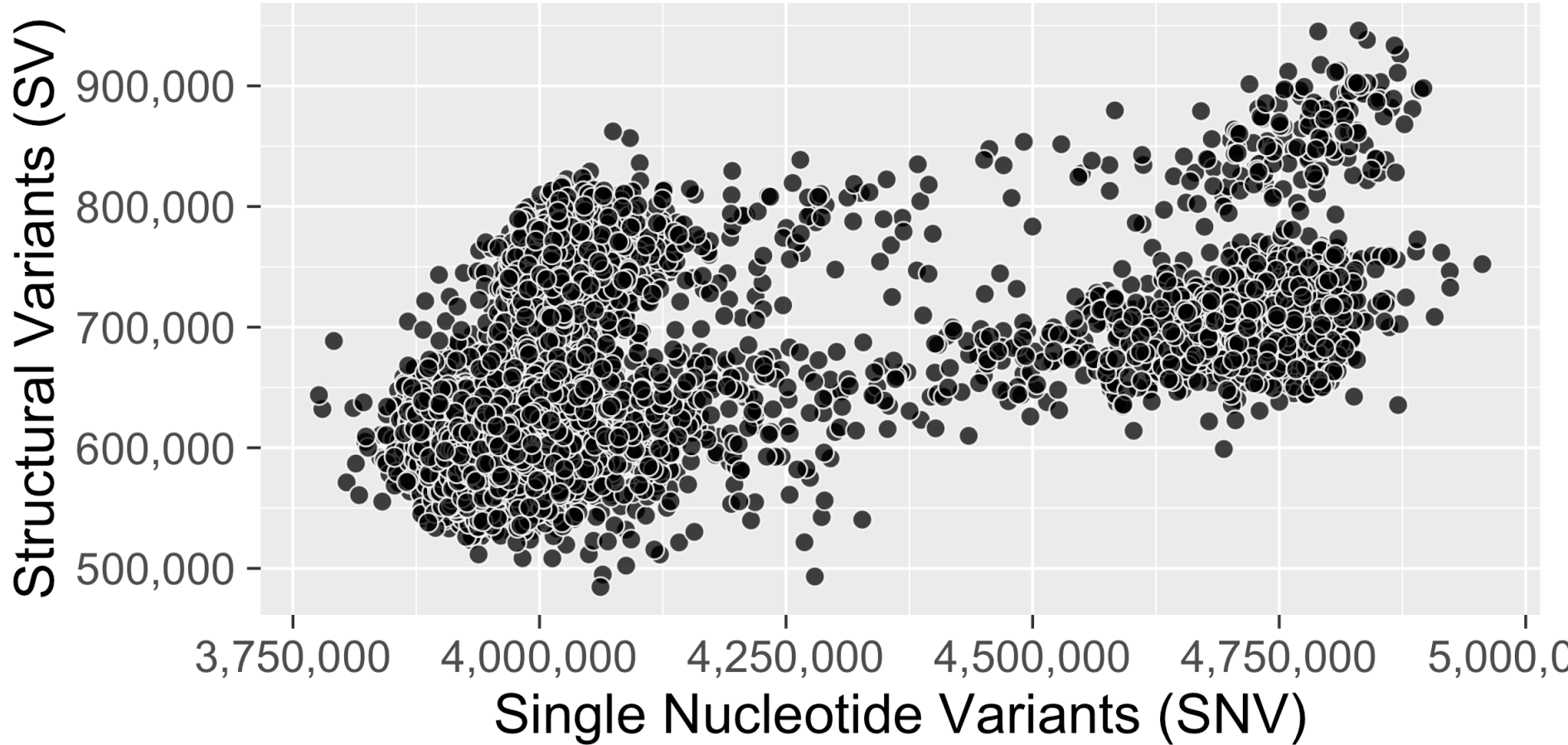
# Rates of Variation by Genomic Feature



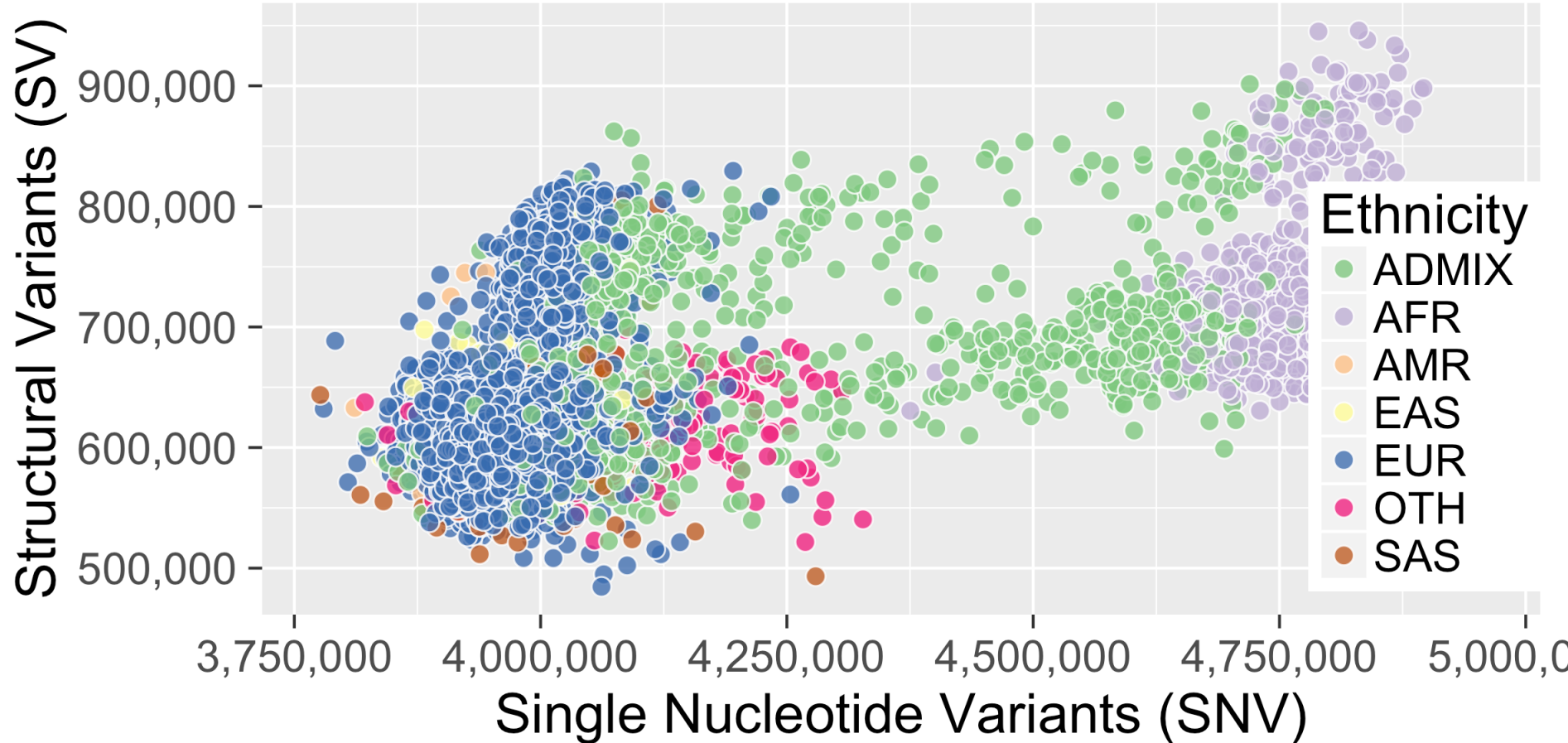
# Variants per Individual



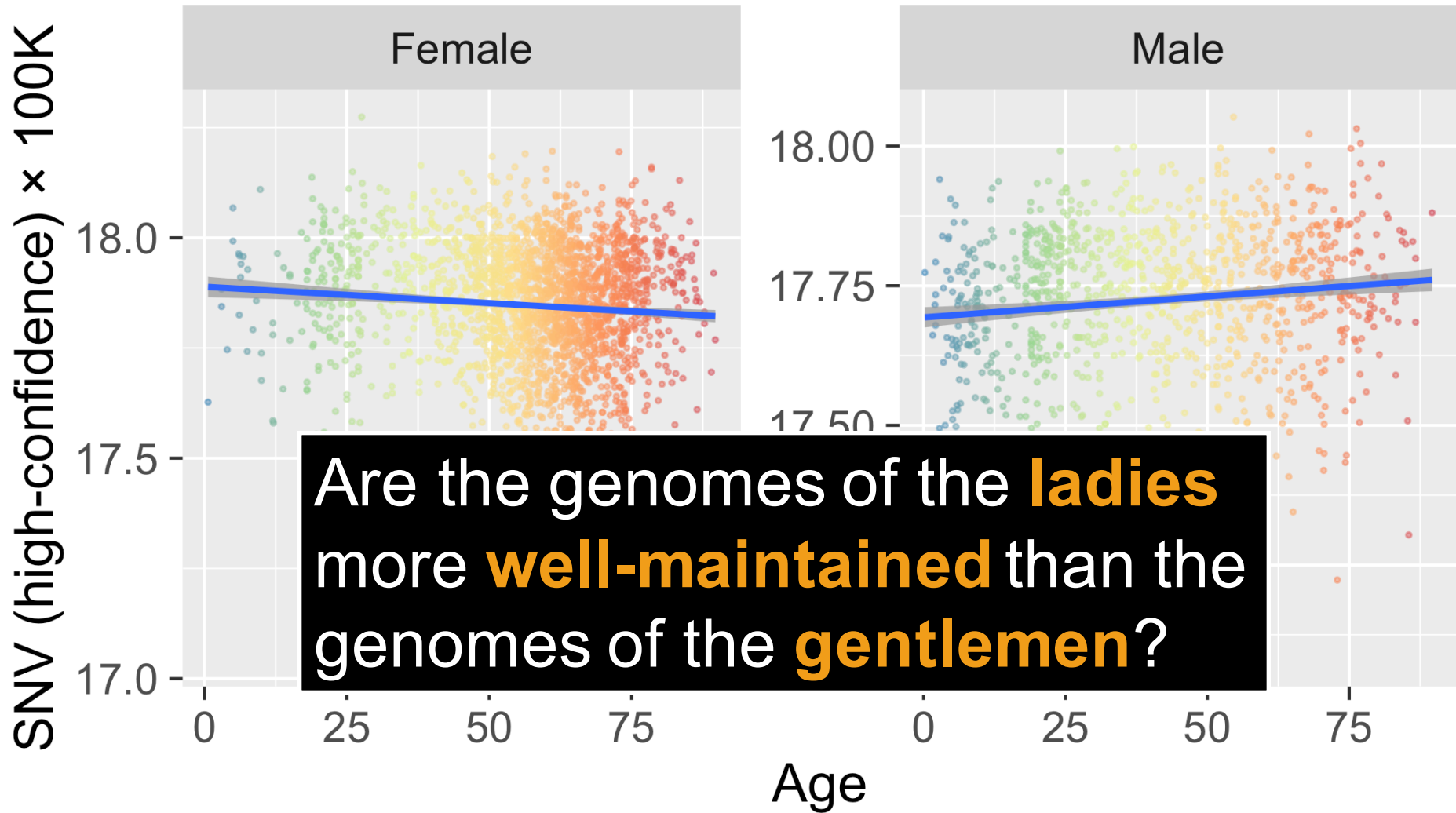
# Variants by Ancestry



# Variants by Ancestry

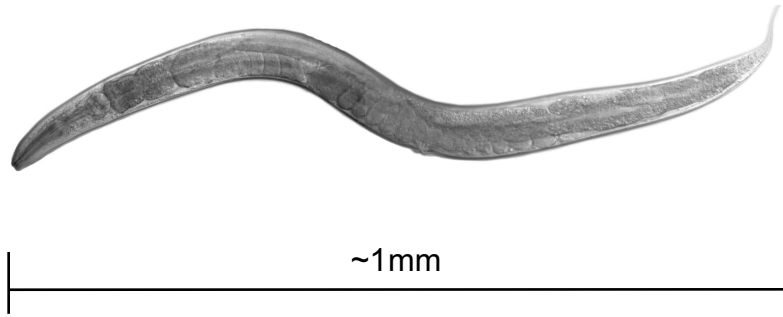


# Variants by Gender EUR



# **The 3D [Non-Coding] Genome**

# Gene Number $\neq$ Organismal Complexity



**Protein-coding genes**  
**~ 19,000**

**Protein-coding genes**  
**~ 19,000**

**Where does the complexity come from?**

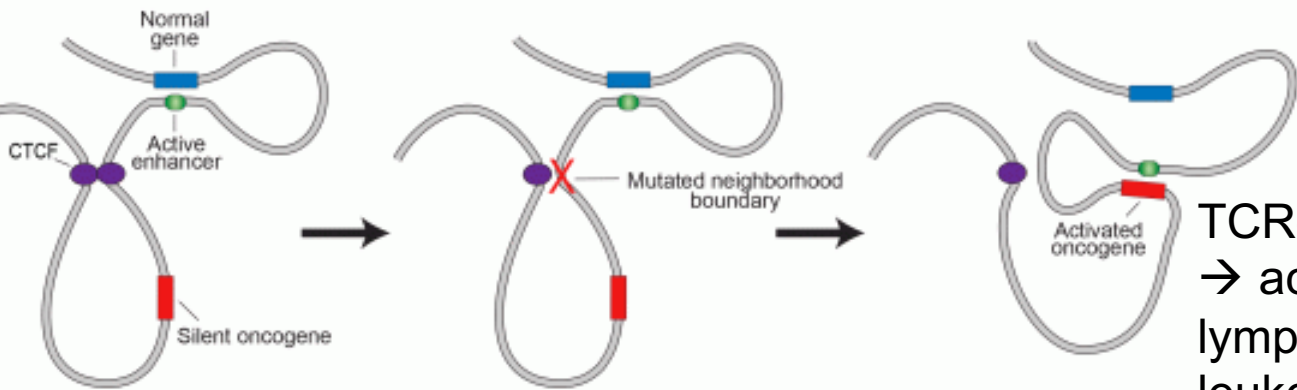
# “Gene regulation”

by region (e.g., breast versus kidney)  
in response to environmental signals  
in development (e.g., embryo versus adult)





# Topologically Associating Domains (TADs) via CTCF

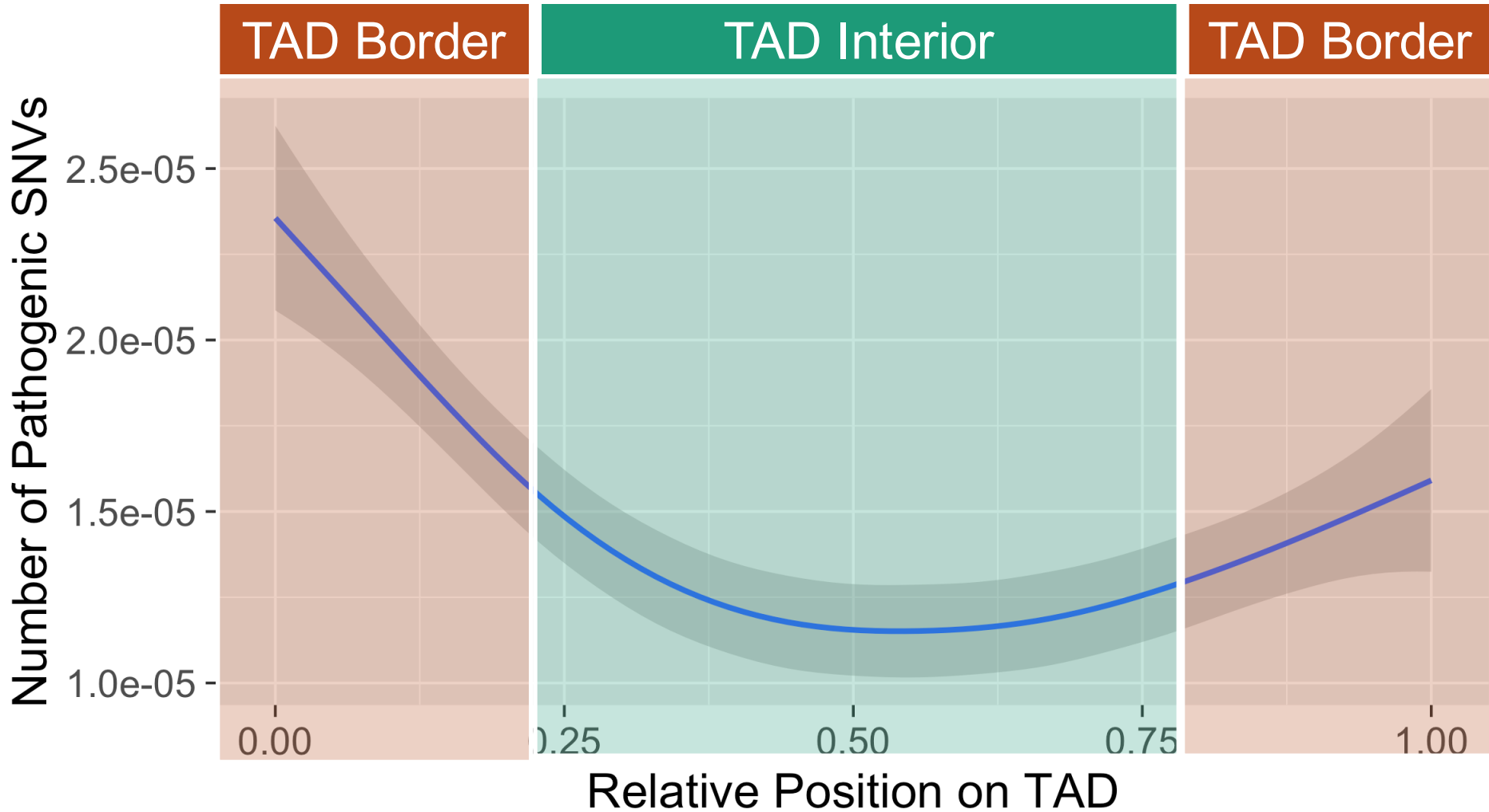


Normal Cell

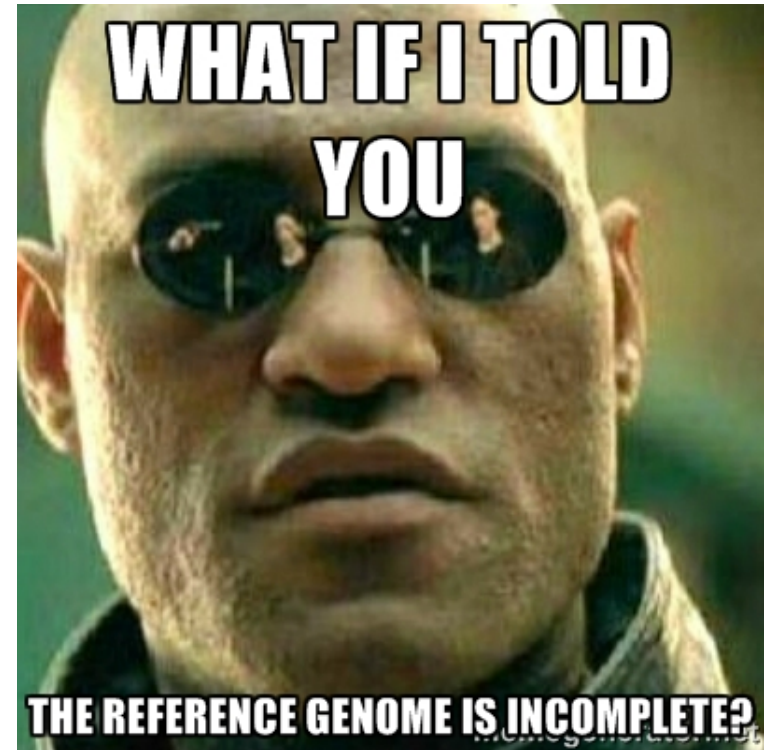
Cancer Cell

TCR-LMO2  
→ acute lymphoblastic leukemia (ALL)

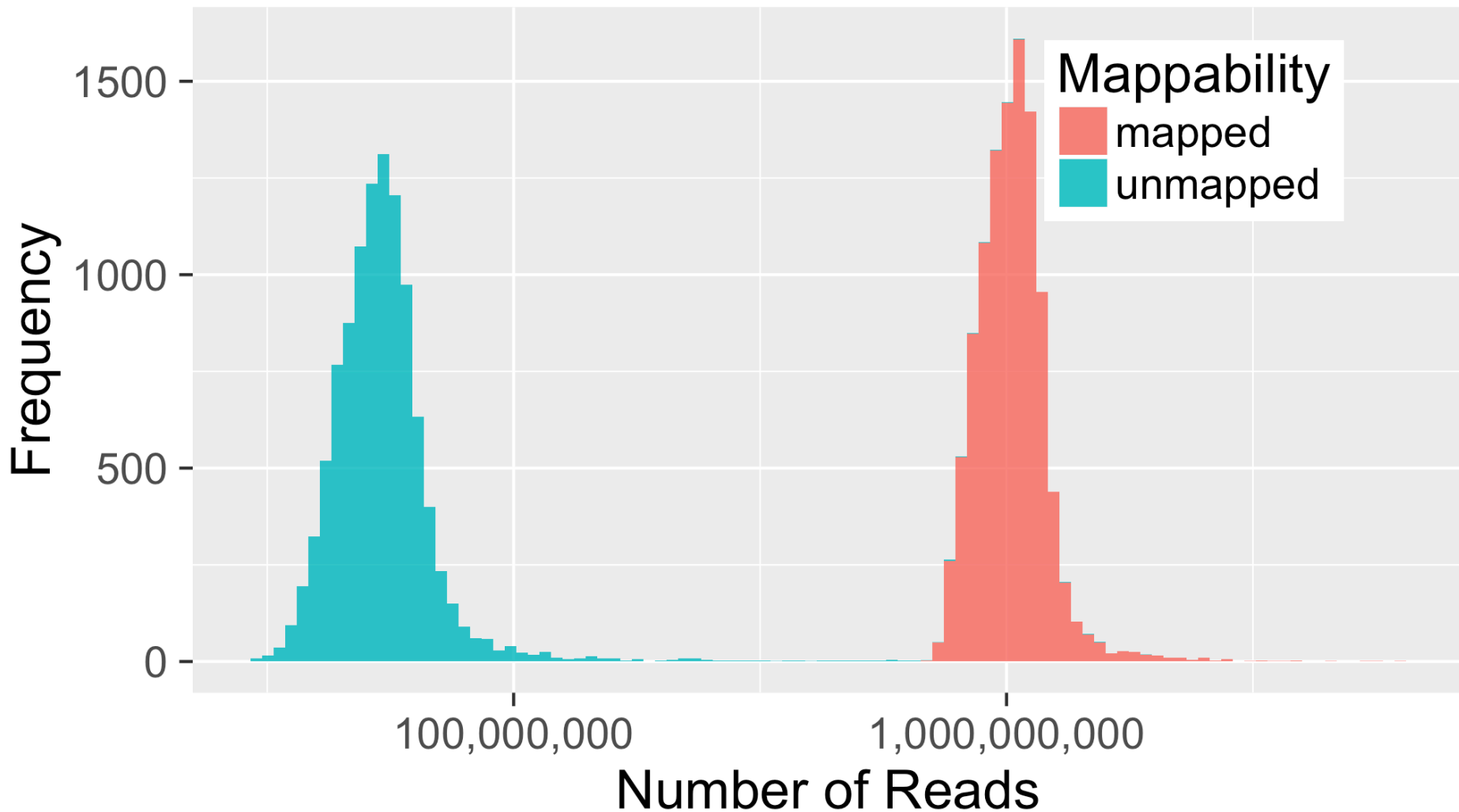
# Pathogenicity over TADs



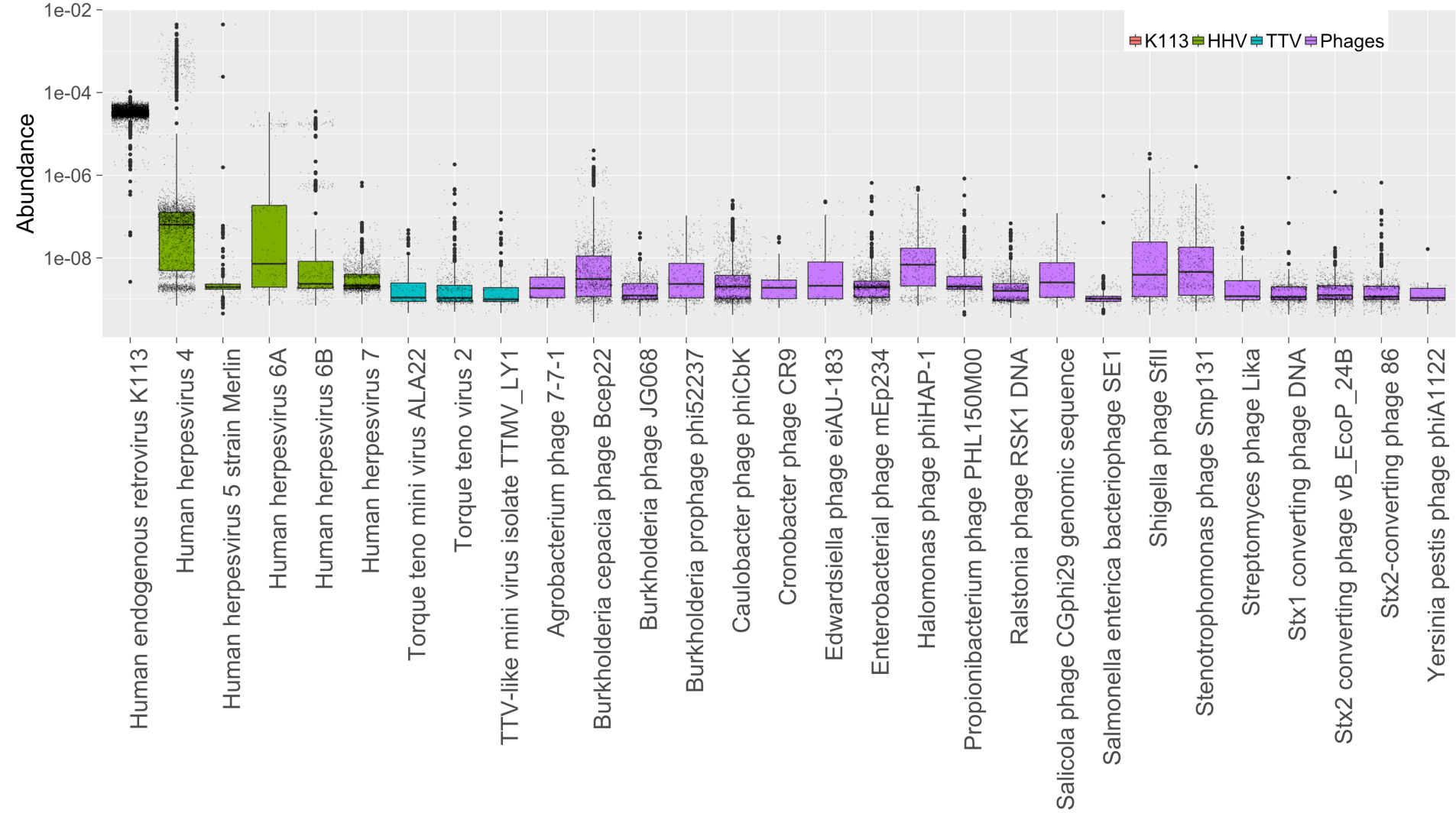
# The Unmapped Genome



# Mapped & Unmapped (per Individual)

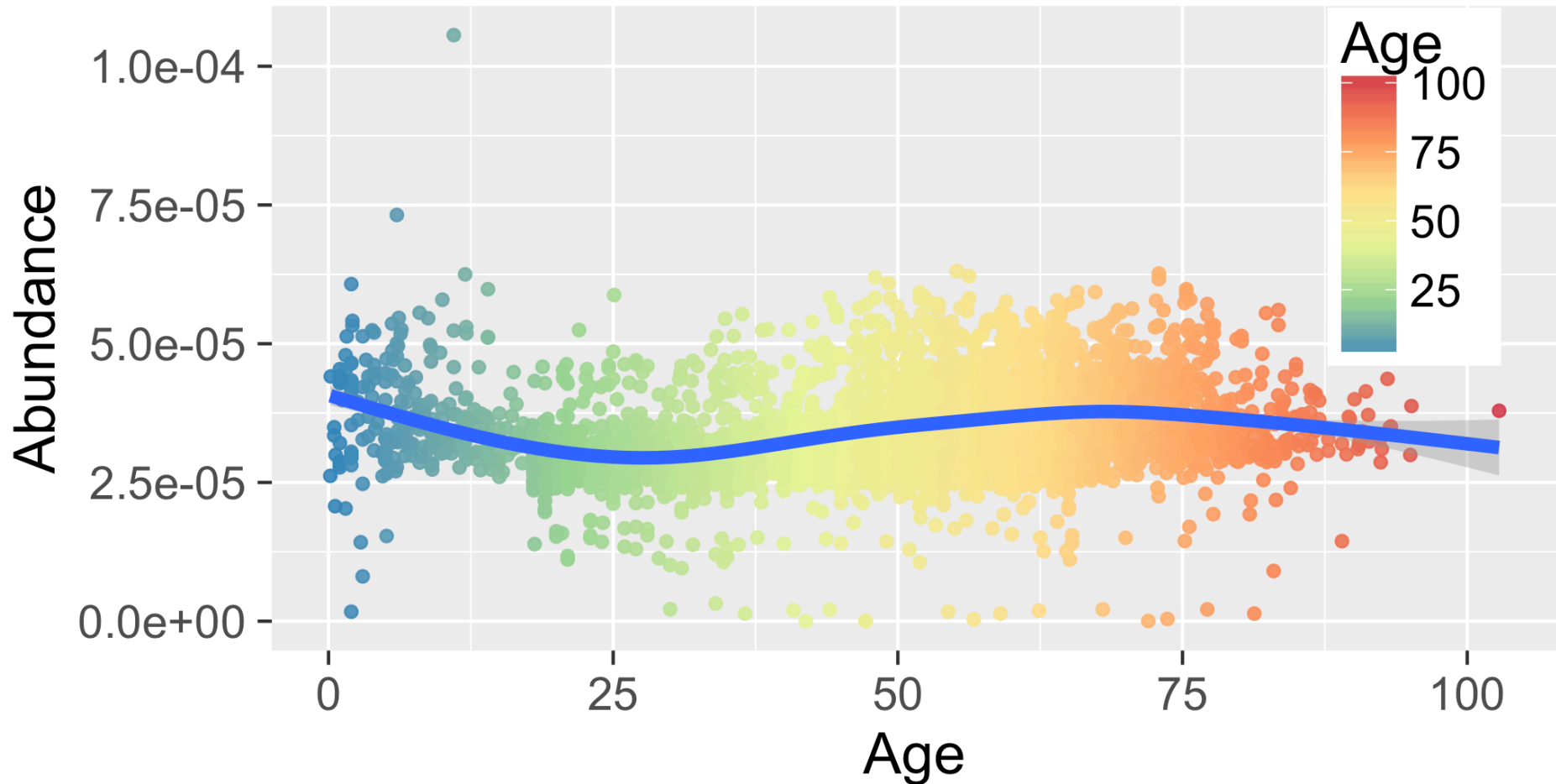


# Viral Load in Unmapped



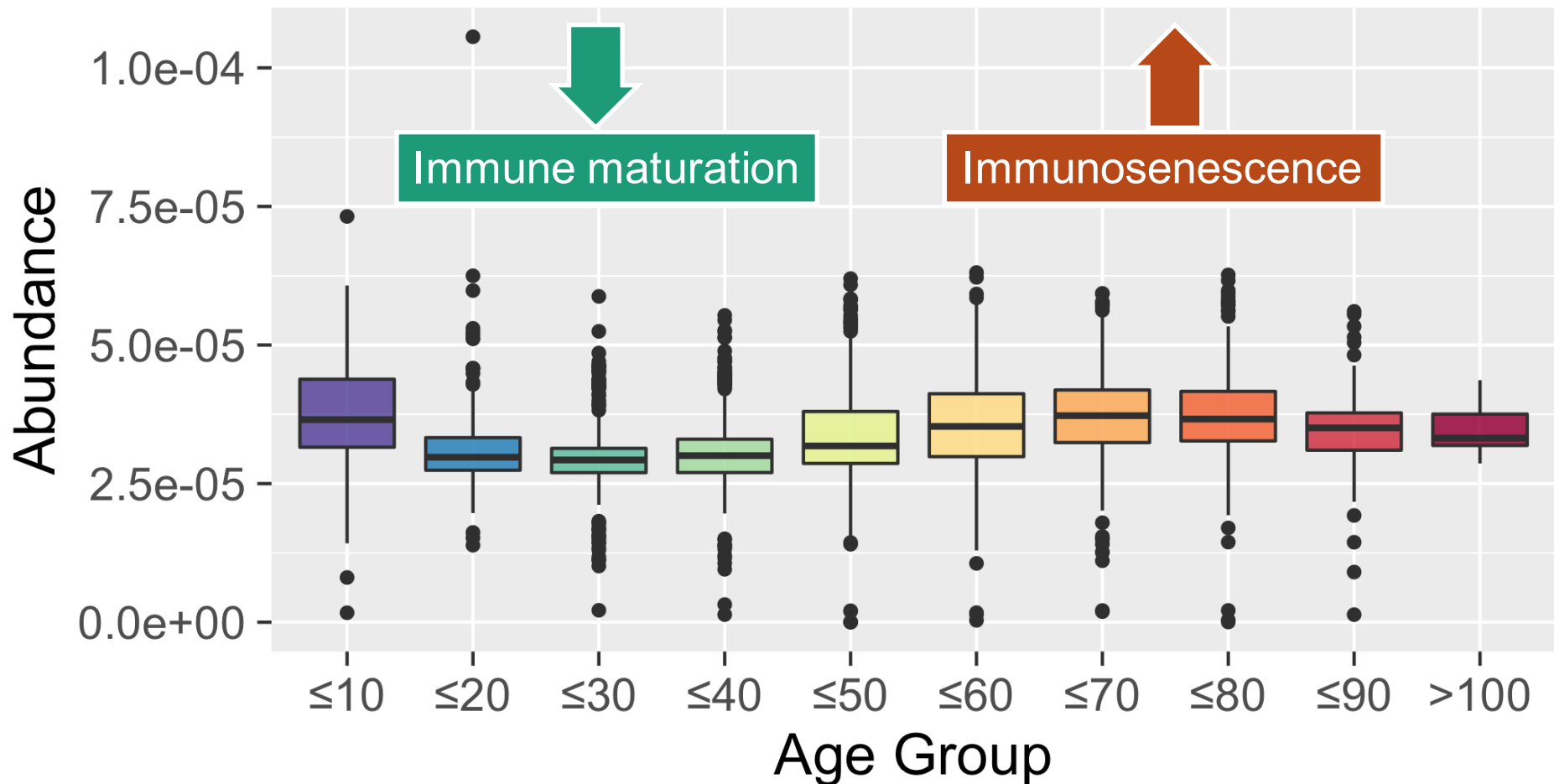
# Viral Load in Unmapped

Human endogenous retrovirus HERV-K113



# Viral Load in Unmapped

Human endogenous retrovirus HERV-K113



# Summary

- Generated first 10,000 deep coverage genomes with clinical standards (coverage, quality, reproducibility)
- Identified 150 million SNVs
- Identified an individual contribution of an average of 8,000 novel variants



# Summary cont'd

- Formulated high-resolution profiles in coding sites with tolerance score
- Identified association between pathogenicity of SNVs and TADs
- Detected dynamics in abundance of Human Endogenous Retrovirus K113 associated with age

Ischaemic  
heart  
disease?

Stroke?

Hypertension ?

# Egyptian Genome!

Cirrhosis?

Cancer?

Kidney  
Disease?

