

Penn Genetics Summer Short Course in Computational Genomics



Penn Medicine



VA
HEALTH CARE | Defining
EXCELLENCE
in the 21st Century

Educational Objectives

Data science

- Read, interpret, and critique the scientific literature
- Ingest, manipulate, summarize, and evaluate large datasets
- Develop reproducible workflows
- Utilize a high-performance compute cluster for data analysis

Genetics

- Understand key concepts in complex trait genetics
- Perform GWAS meta-analysis and downstream computational analyses
- Possess a framework for interpreting GWAS findings to biological insights

Deliverables

- GWAS meta-analysis
- Downstream analysis of summary statistics
 - Significant loci
 - Prioritized variants and genes
 - Biological meaning

Structure

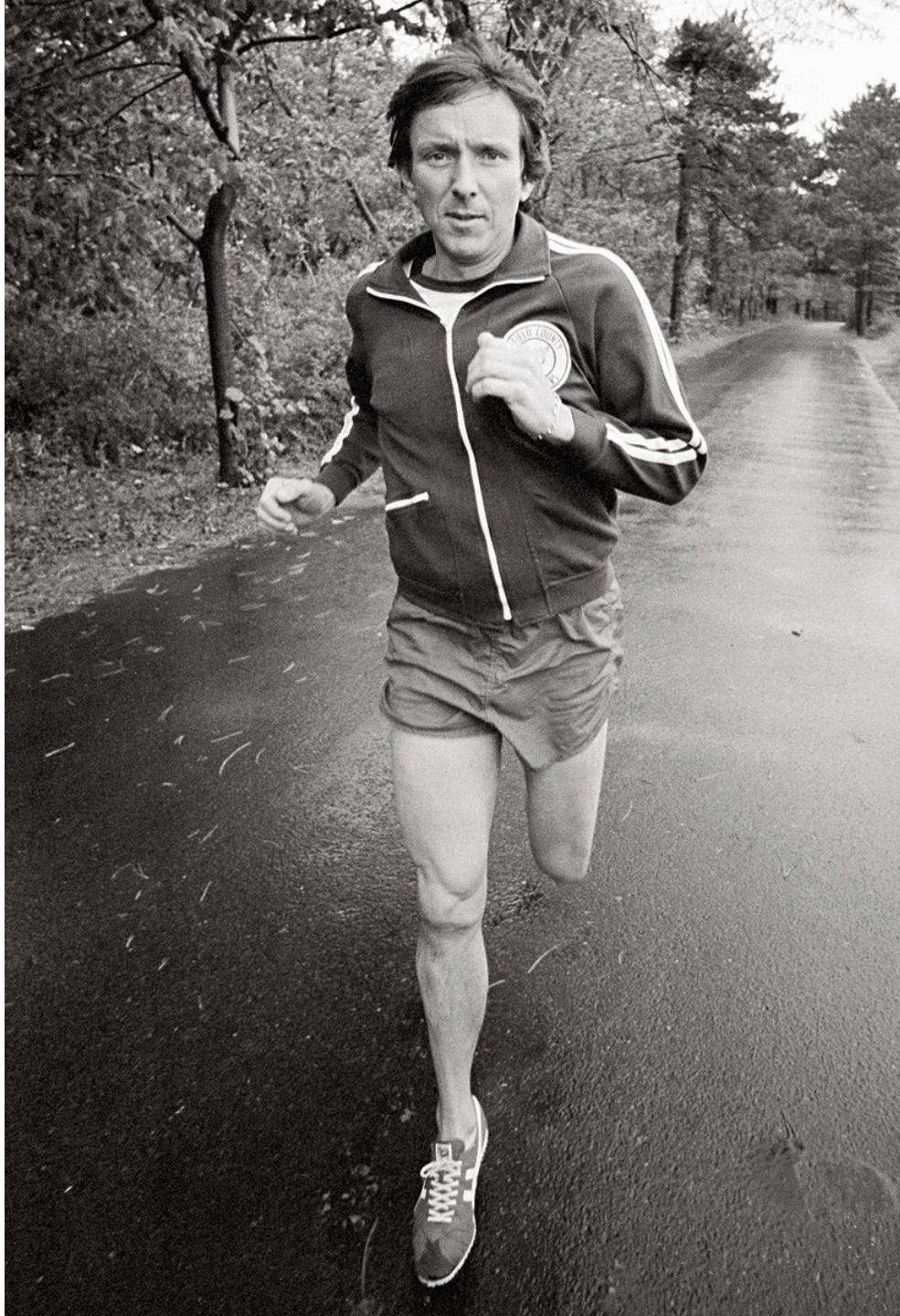
- Weekly meetings (Wednesday 1-3 [except Juneteenth – Tuesday June 18])
 - **Didactic**
 - Presentation
 - Theory / background
 - Practical advice
 - Coding demonstration
 - Alternating with **Discussion Section**
 - Project updates
 - Troubleshooting
- Literature review
- Projects
- Mentors

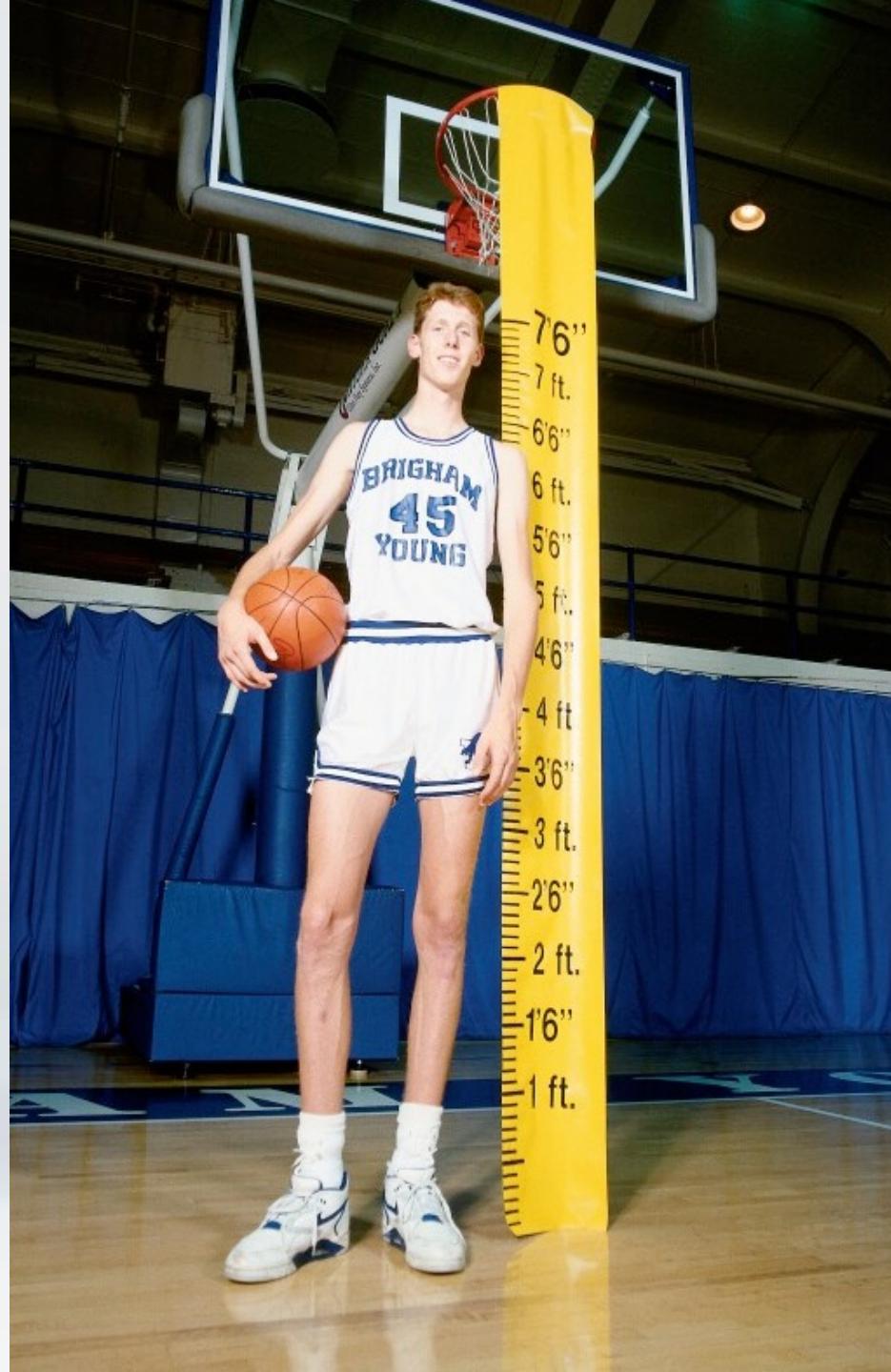
Genetic alchemy: From base data to noble

Scott M. Damrauer, MD

Vice-Chair of Clinical Research, Department of Surgery
William Maul Measey Associate Professor of Surgery II
Associate Professor of Genetics



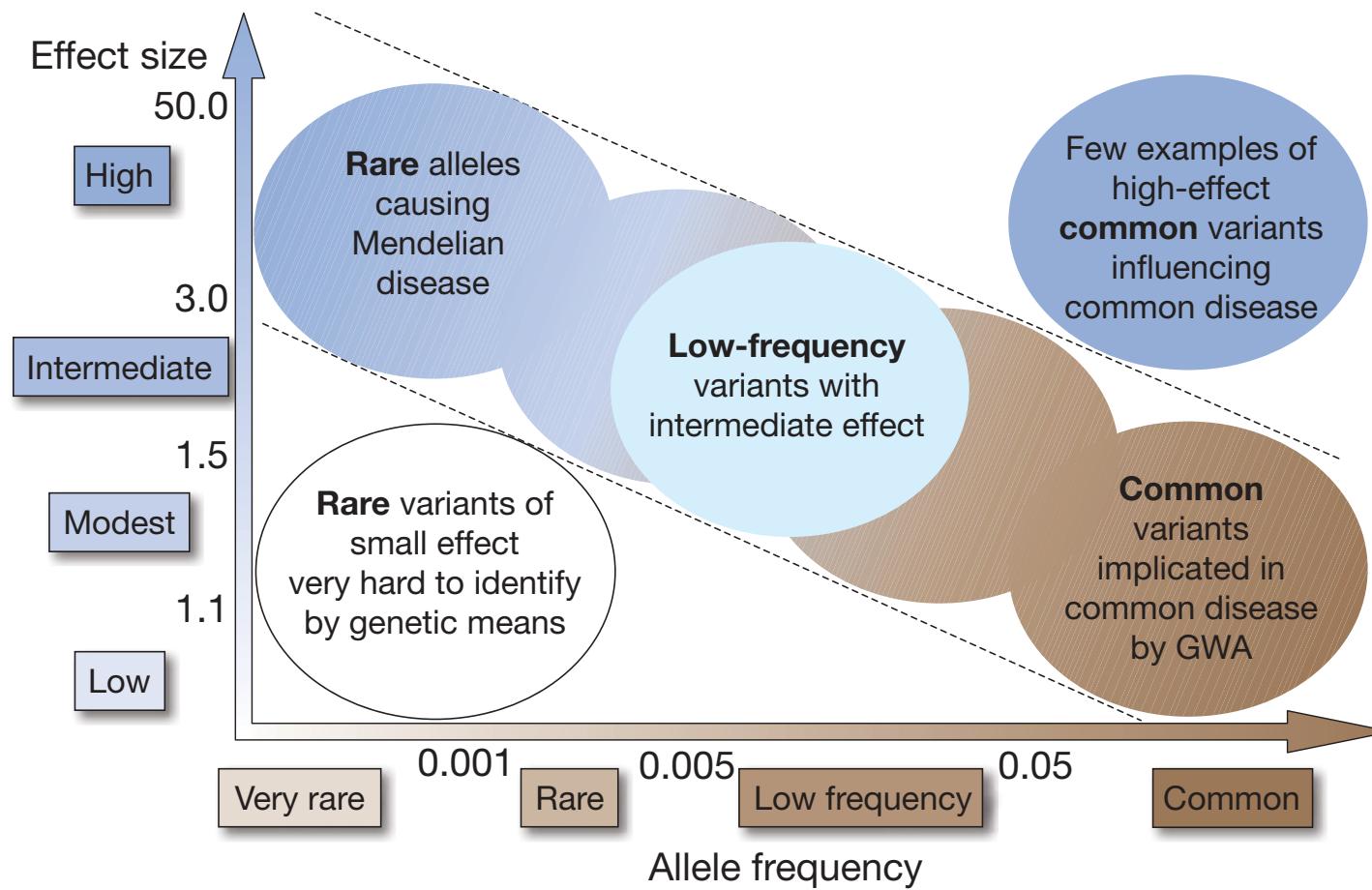






David Teniers the Younger (1610–1690), *The Alchemist*, ca. 1643–45. Oil on panel, 20 1/8 x 28 in (51 x 71 cm). Herzog Anton Ulrich Museum, Braunschweig (139)

Common variant genetics for risk prediction



All humans are at least 99.5% genetically similar

- Most common form of genetic variation occurs when a single nucleotide differs between individuals
 - This is a **Single Nucleotide Polymorphism (SNP)**
 - On average ~ 1/300 base pairs
- Each individual has 2 copies of each gene / “alleles”

The diagram illustrates four individuals' DNA sequences, each with three SNP markers indicated by arrows pointing to specific nucleotides. Individual 1 has a SNP at position 4 (C to G). Individual 2 has a SNP at position 4 (C to T). Individual 3 has a SNP at position 4 (A to T). Individual 4 has a SNP at position 4 (C to A).

Individual 1	AAC A C GCC A.... TT C G G GT C.... AGT C GACCG....
Individual 2	AAC A C GCC A.... TT C G A GGT C.... AGT C AACCG....
Individual 3	AAC A TGCC A.... TT C G GGT C.... AGT C AACCG....
Individual 4	AAC A C GCC A.... TT C G GGT C.... AGT C GACCG....

Testing genetic association

Individual 1	A A C A C G C C A.... T T C G G G G T C.... A G T C G A C C G....
Individual 2	A A C A C G C C A.... T T C G A G G T C.... A G T C A A C C G....
Individual 3	A A C A T G C C A.... T T C G G G G T C.... A G T C A A C C G....
Individual 4	A A C A C G C C A.... T T C G G G G T C.... A G T C G A C C G....

Does the “A” allele Increase the risk for Disease relative to “G”?

Testing genetic association in populations



Quantifying associations with linear outcomes

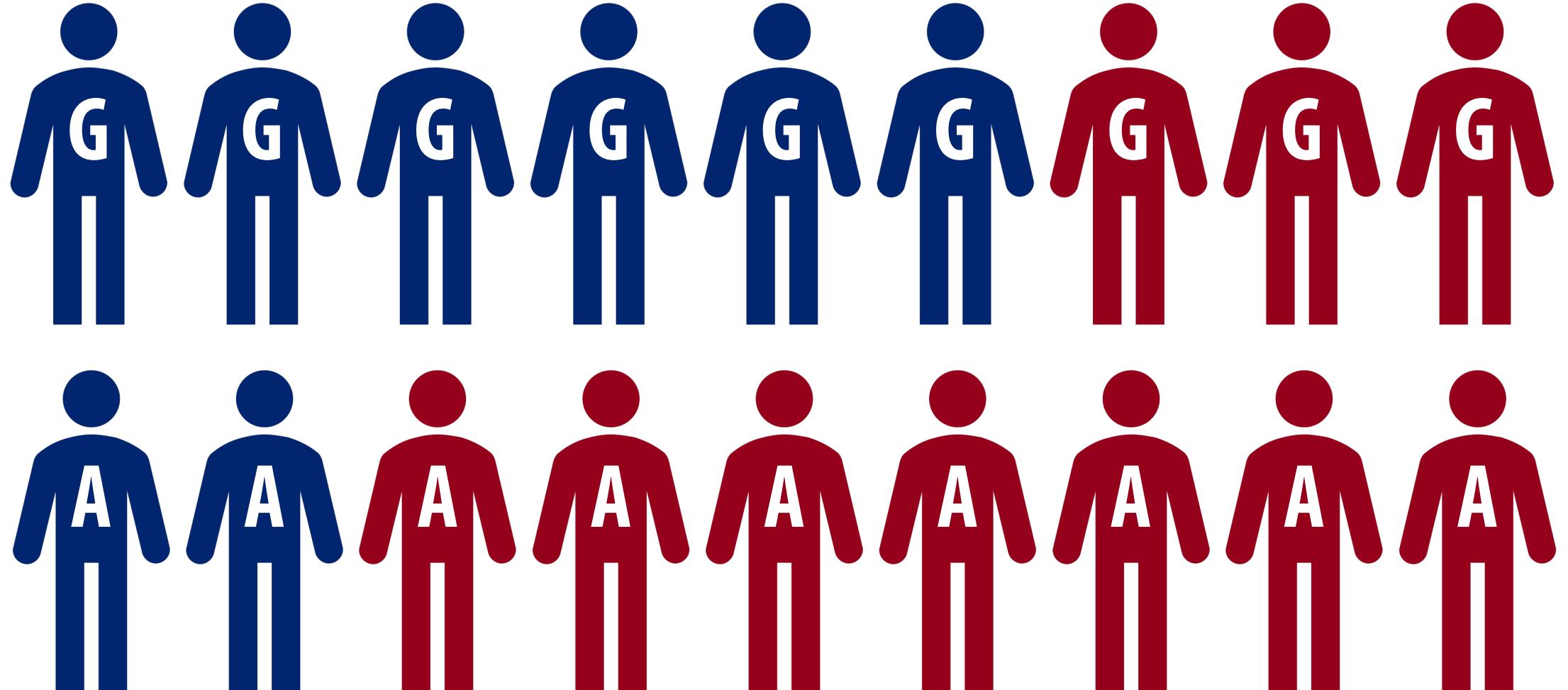
T-test or ANOVA to compare means

- Hypothesis testing

Linear regression to quantify association

- Hypothesis testing
- Modeling
- Multiple exposures (Multivariable)

Testing genetic association in populations



Quantifying associations with categorical outcomes

Chi Squared

- Unequal proportions from a $r \times c$ table
- Hypothesis testing
- Odds ratio

Can also use

- Fisher's Exact Test

Logistic regression

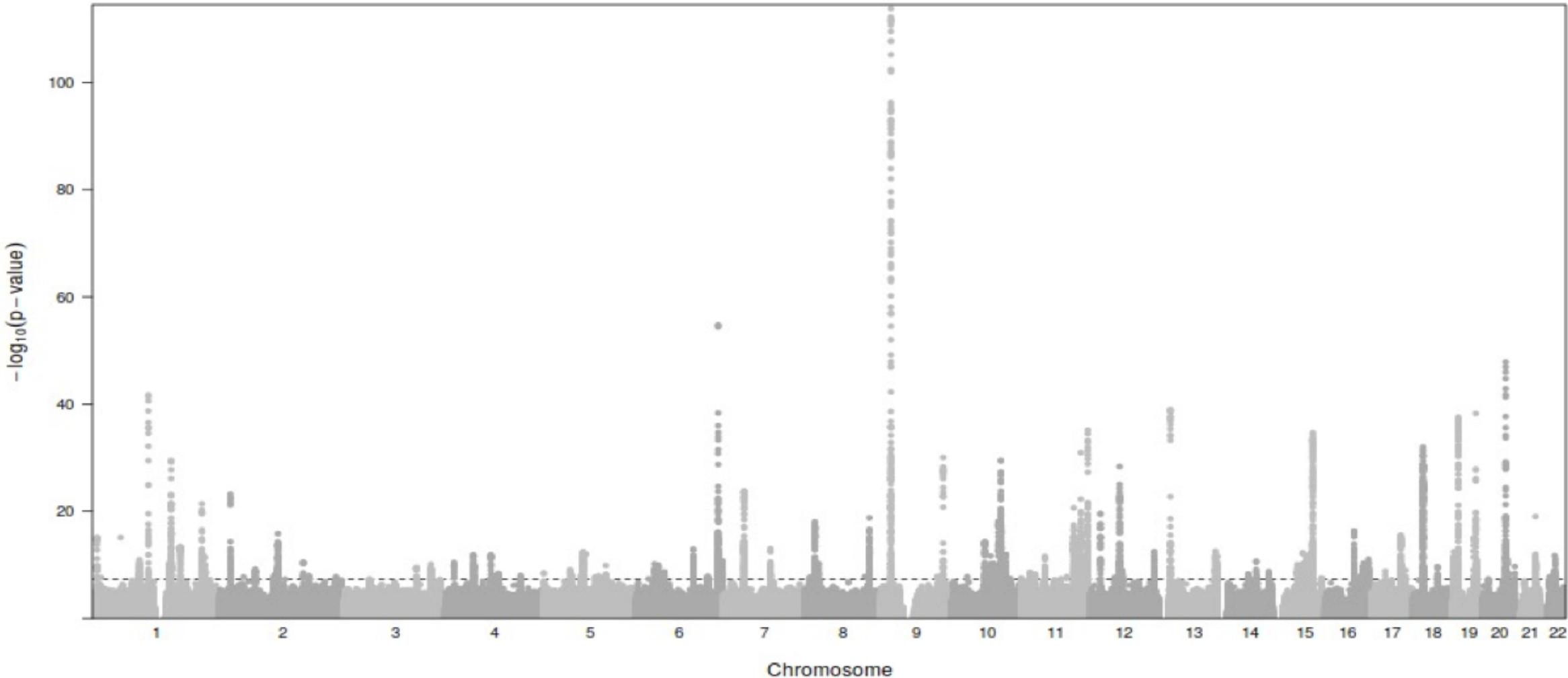
- Hypothesis testing
- Modeling
- Multivariable

Genome Wide Association Study (GWAS)

$$\text{logit}(P|X, G) = \beta_0 + \boldsymbol{\beta}_G \cdot \mathbf{G} + \beta_{age} \cdot age + \beta_{sex} \cdot sex + \sum_{k=1}^{10} \beta_{PC} \cdot PC_k$$

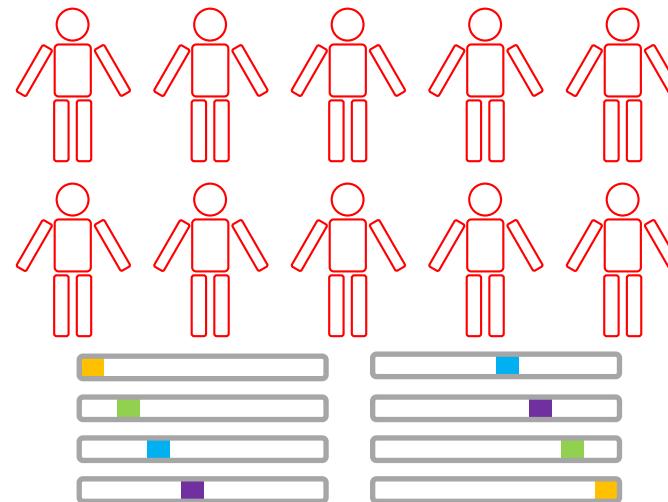
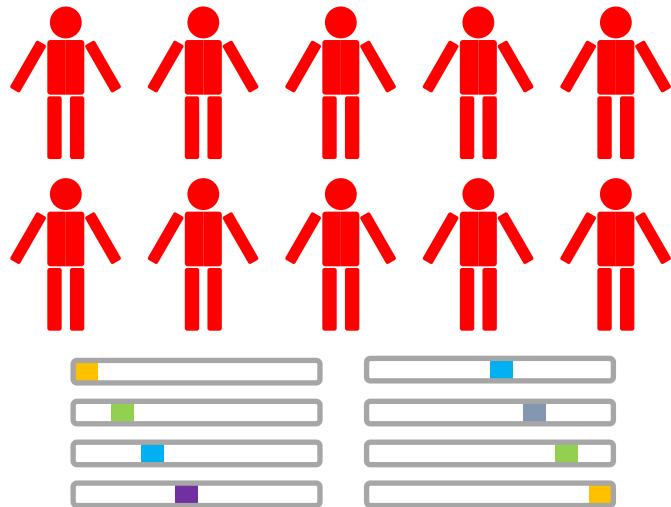
- Repeat the model for every genetic variant in the panel
 - ~300,000 to 400,000 for chip based GWAS
 - ~ 50 million times for imputed GWAS
- P-values adjusted for multiple testing
 - Genome-wide significance $P < 5 \times 10^{-8}$

AAAGen meta-analysis



ASCERTAINED case-control genetic association study

- Case control cohort based on disease status
- Perform genotyping on all participants
- Impute from directly genotyped data to more dense genetic data based on reference panels
- Perform association testing sequentially for each variant

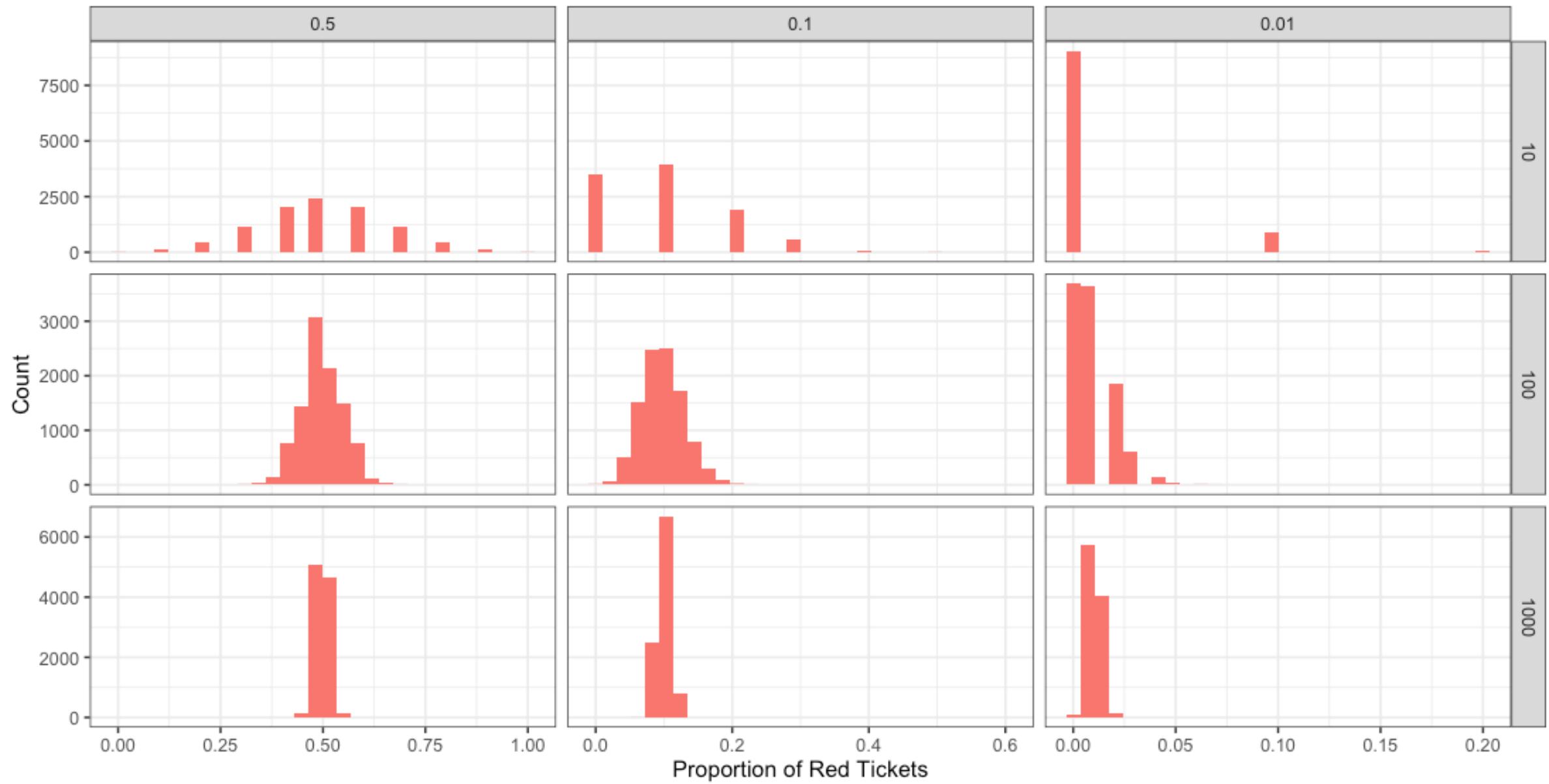


Sample size

- Traditionally thought of as
 - Acceptable risk of a false negative (power, beta, probability of a type 2 error)
 - Acceptable false positive rate (alpha, probability of a type 1 error)
 - Magnitude of effect
 - Variance (how precisely you can measure the effect)
- In the setting of an observational study, also depends on:
 - Rate of exposure in the population

Drawing lottery tickets





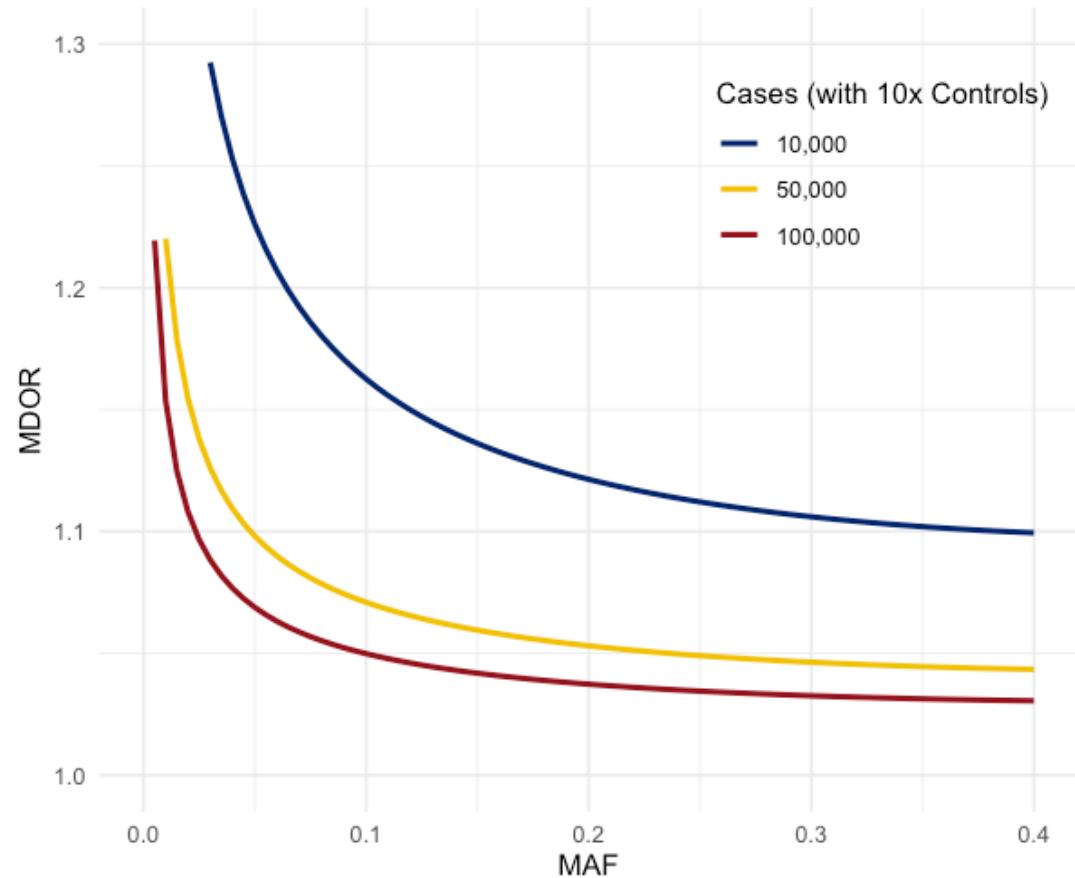
**What percentage of
smokers have lung
cancer?**

**Does smoking increase
the risk of lung cancer?**

GWAS power

- Bigger sample size is always better
 - This refers to cases more than controls
 - After 10 controls per case, the value of more controls becomes less
- For the “power” calculation, we solve for the minimum detectable odds ratio (MDOR)
 - genpwr package in R

MDOR for binary trait



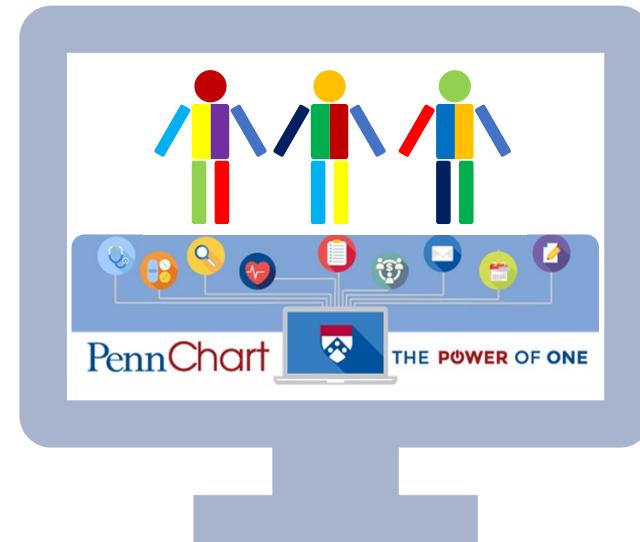
- 10 x controls per case
- 80% Power
- Alpha = GWS
- Additive genetic model

EHR based genetics

Genetic data is
acquired **ONCE** per
person

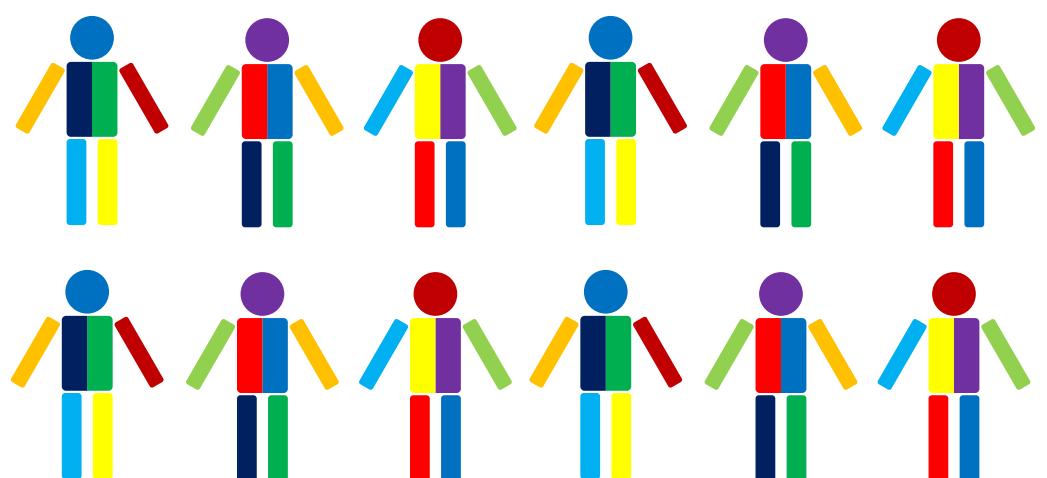
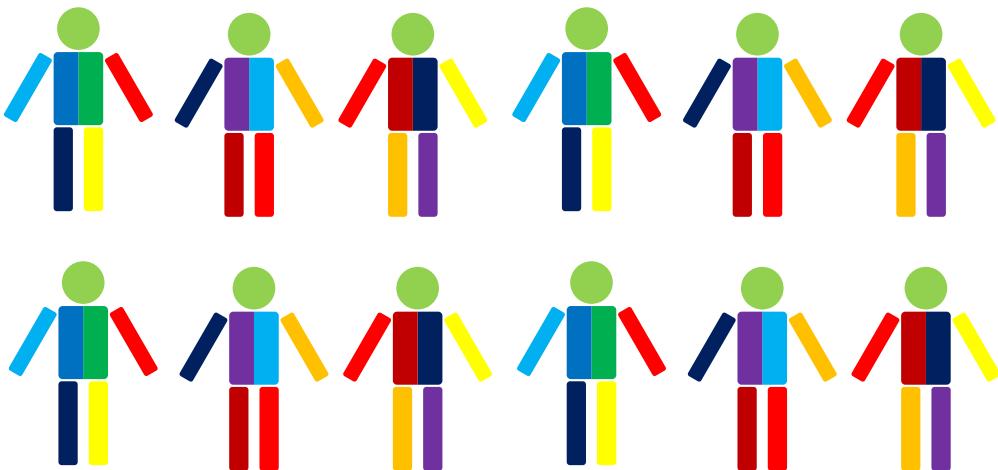


EHR captures the
full range of
biomedical
"traits"



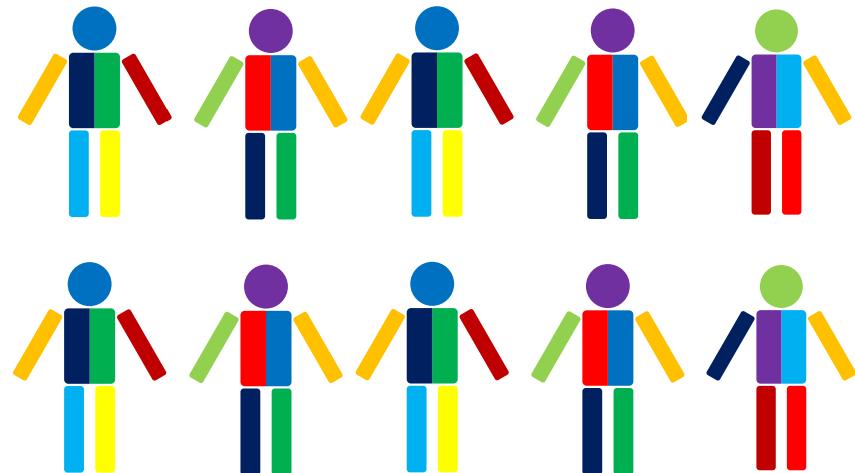
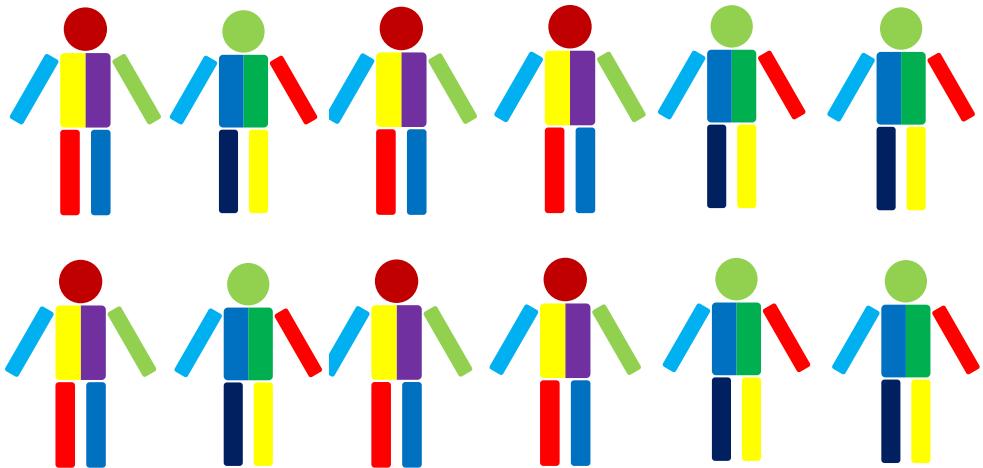
EHR based genetics

- “**Phenotyping**” performed on all aspects of data
- Structured: ICD-9/10 codes, CPT codes, Vital signs, labs
- Semi-structured: Text reports
- Unstructured: Office notes, Images



EHR based genetics

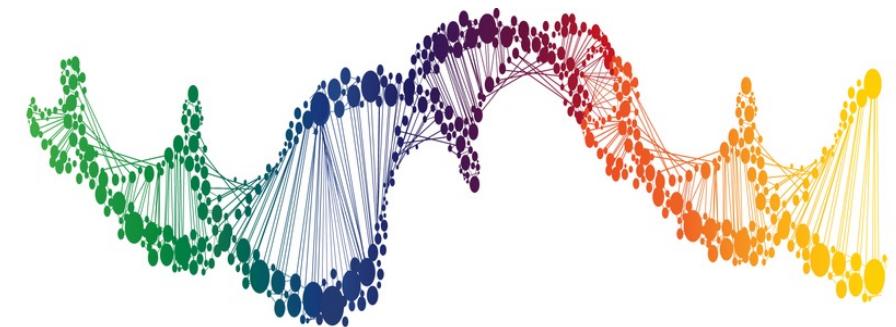
- “**Phenotyping**” performed on all aspects of data
- Structured: ICD-9/10 codes, CPT codes, Vital signs, labs
- Semi-structured: Text reports
- Unstructured: Office notes, Images



Phenome-wide association study



GWAS looks for associations with a specific phenotype across the entire genome



PheWAS looks for associations with a specific genotype across the entire phenotype

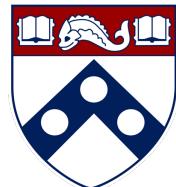
The Biobank Continuum



FINN GEN



estonian genome center
university of tartu



Penn Medicine

biobank^{uk}



Mount
Sinai



VANDERBILT
UNIVERSITY

Geisinger



Mass General Brigham

VA Million Veteran Program

 90%
Male

 10%
Female

 80K + (or 8%)
Hispanic

180K + (or 18%) Black

MVP has the largest cohort of people of African population of any research program in the world

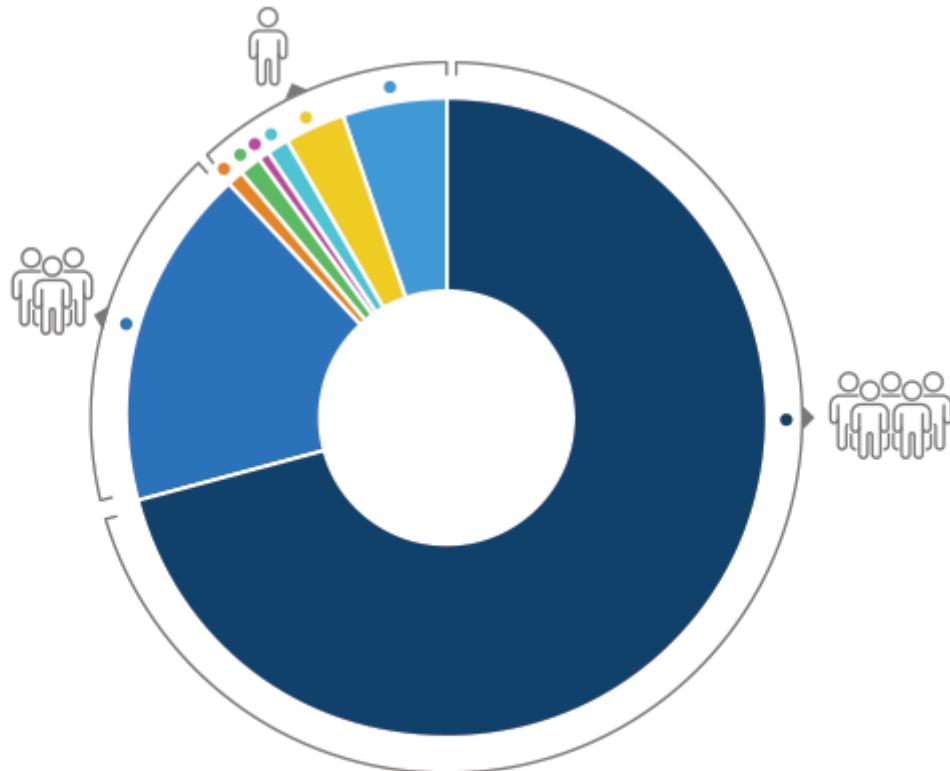
25% minority

Racial and ethnic groups represented

Average age is 66

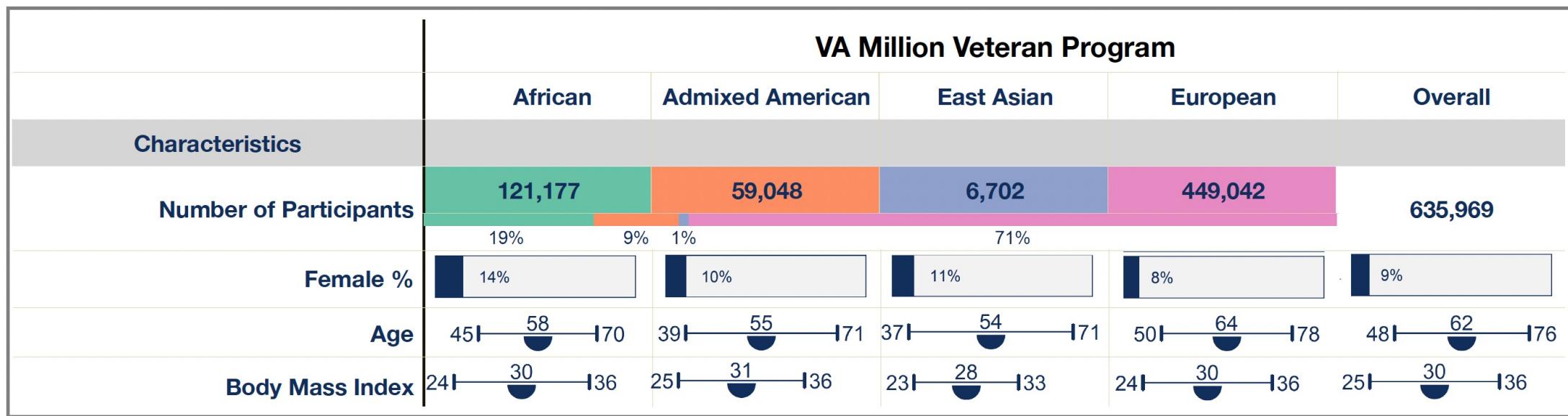
with largest cohort (33%) aged 70–79 (83% over 50)

1,023,500+ enrolled
(as of March 13, 2024)

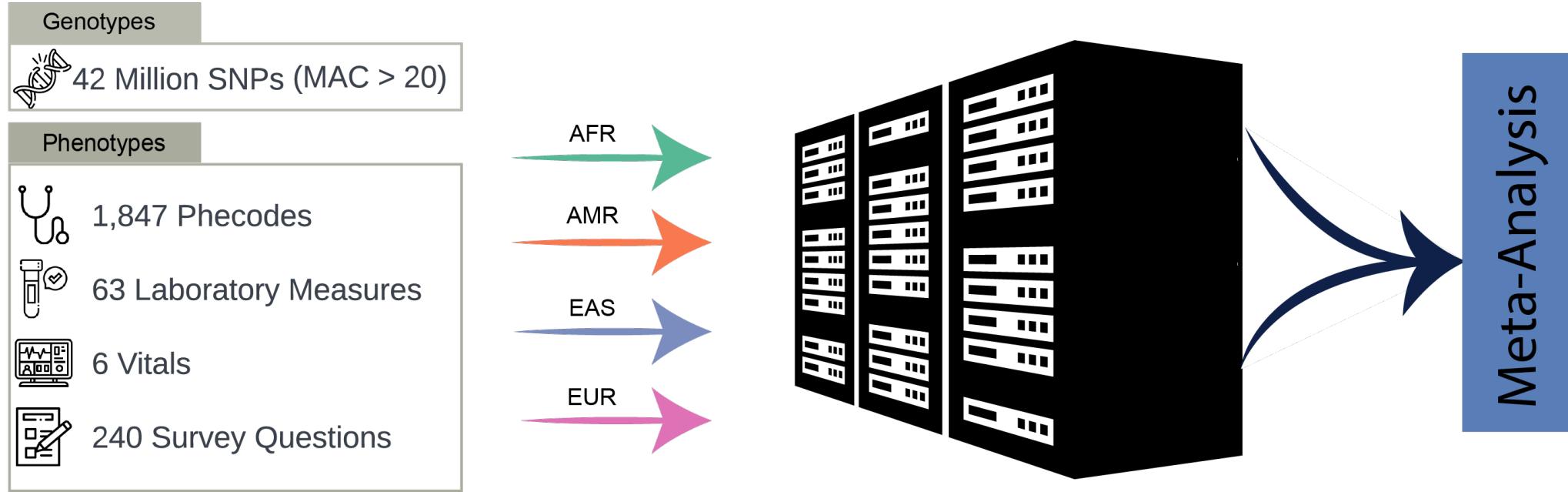


- White 70.99%
- Black or African American 17.32%
- American Indian or Alaskan Native 0.75%
- Asian 1.07%
- Native Hawaiian or Other Pacific Islander 0.47%
- Other 1.35%
- Multi-Racial 2.88%
- Missing 5.18%

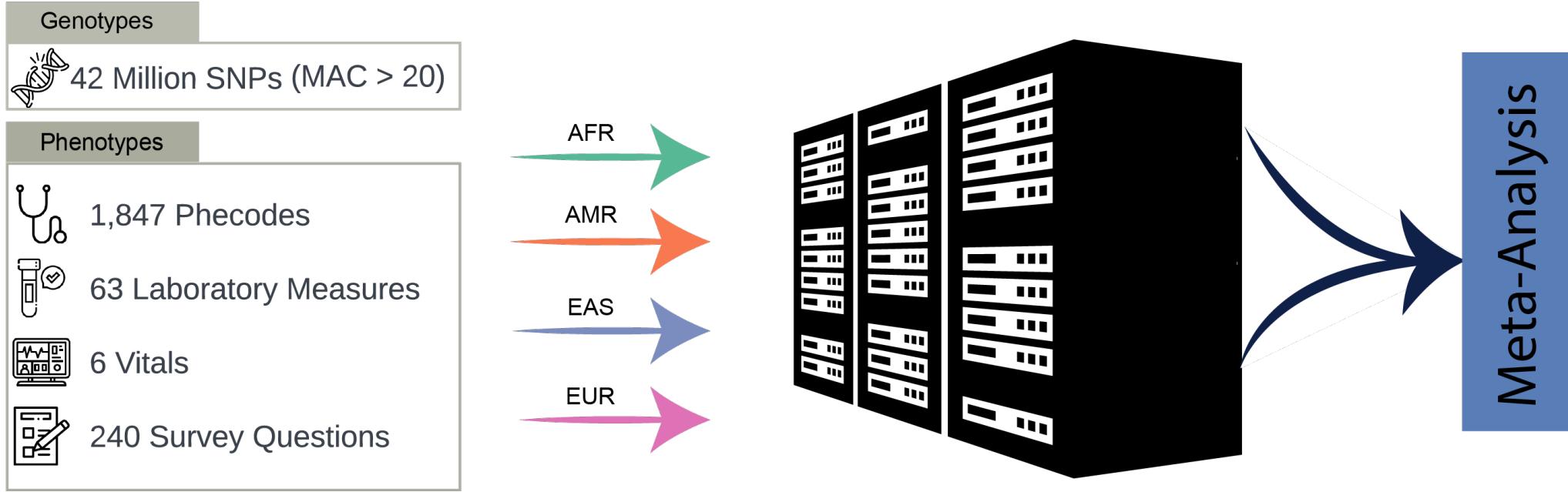
VA MVP – Study Participants



Analysis pipeline

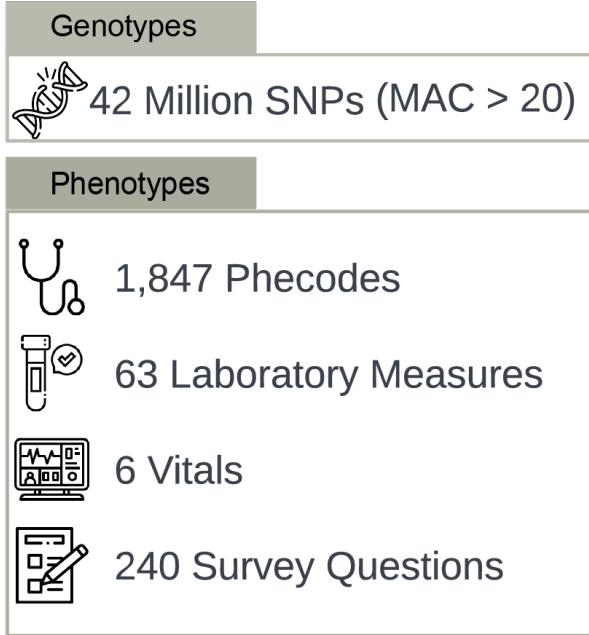


Analysis pipeline



Estimated wall time of 8 years on existing compute infrastructure

Analysis pipeline

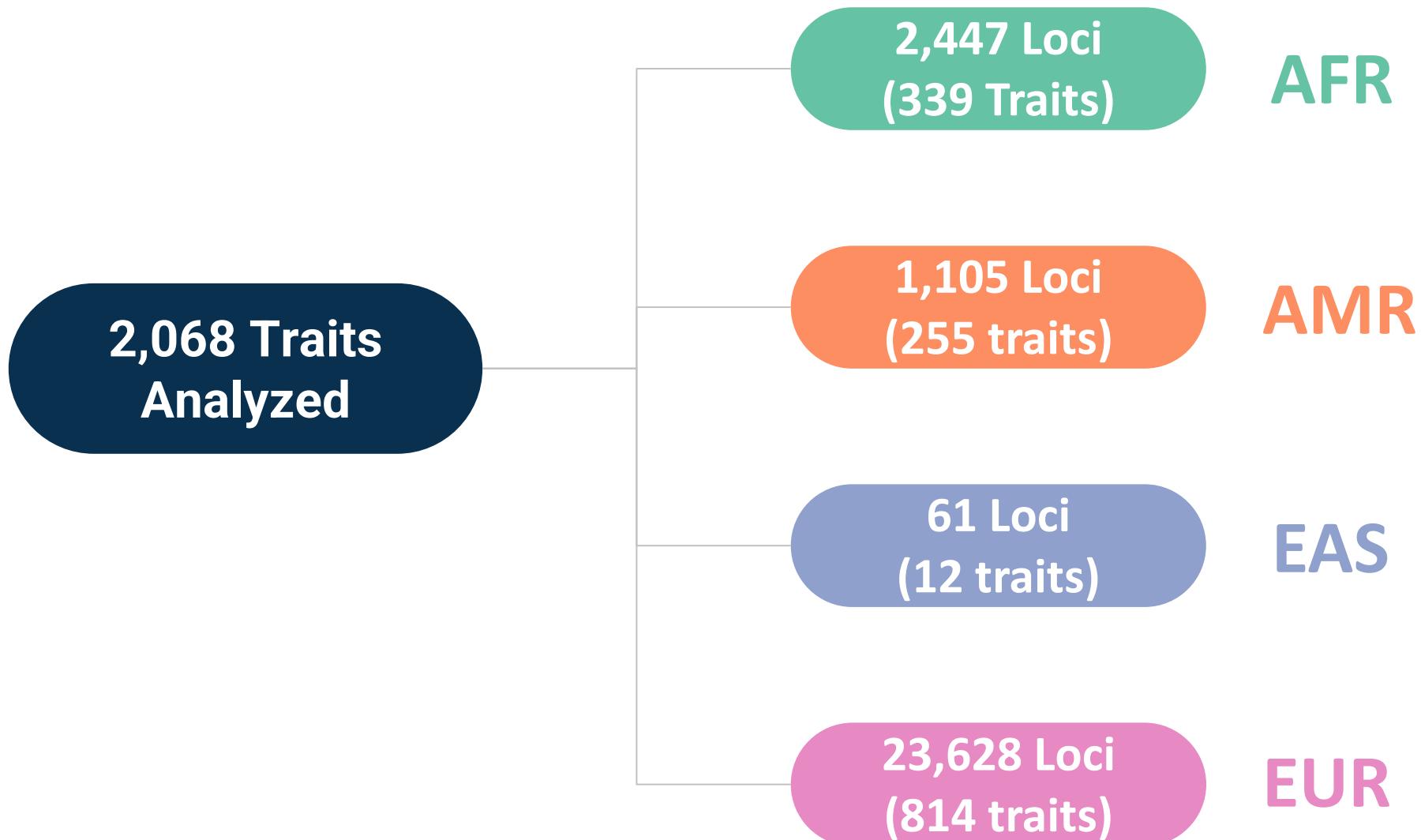


AFR
AMR
EAS
EUR



Actual wall time of 14 days on Summit

Genome-wide loci within each population group



Study wide p-value threshold 4.6×10^{-11}

Diversity and Scale: Genetic Architecture of 2,068 Traits in the VA Million Veteran Program

Anurag Verma^{†,1,2,3}, Jennifer E Huffman^{†,4,5,6}, Alex Rodriguez^{†,7}, Mitchell Conery^{†,8}, Molei Liu^{†,9}, Yuk-Lam Ho⁴, Youngdae Kim¹⁰, David A Heise¹¹, Lindsay Guare², Vidul Ayakulangara Panickan¹², Helene Garcon⁴, Franciel Linares¹³, Lauren Costa¹⁴, Ian Goethert¹⁵, Ryan Tipton¹⁶, Jacqueline Honerlaw⁴, Laura Davies¹⁷, Stacey Whitbourne^{6,14,18}, Jeremy Cohen¹¹, Daniel C Posner⁴, Rahul Sangar¹⁴, Michael Murray¹⁴, Xuan Wang¹², Daniel R Dochtermann¹⁹, Poornima Devineni¹⁹, Yunling Shi¹⁹, Tarak Nath Nandi⁷, Themistocles L Assimes²⁰, Charles A Brunette^{21,6}, Robert J Carroll²², Royce Clifford^{23,24}, Scott Duvall^{25,26}, Joel Gelernter^{27,28}, Adriana Hung²⁹, Sudha K Iyengar³⁰, Jacob Joseph^{31,32}, Rachel Kember^{33,34}, Henry Kranzler^{33,34}, Daniel Levey^{35,27}, Shiu-Wen Luoh^{36,37}, Victoria C Merritt²³, Cassie Overstreet²⁷, Joseph D Deak^{38,39}, Struan F A Grant^{40,41,42,43}, Renato Polimanti³⁸, Panos Roussos⁴⁴, Yan V Sun⁴⁵, Sanan Venkatesh⁴⁴, Georgios Voloudakis⁴⁴, Amy Justice^{35,46,47}, Edmon Begoli⁴⁸, Rachel Ramoni⁴⁹, Georgia Tourassi⁵⁰, Saiju Pyarajan¹⁹, Philip S Tsao^{20,51}, Christopher J O'Donnell⁵², Sumitra Muralidhar⁴⁹, Jennifer Moser⁴⁹, Juan P Casas⁴, Alexander G Bick⁵³, Wei Zhou^{54,55,56}, Tianxi Cai^{‡,12}, Benjamin F Voight^{‡,1,8,43,57}, Kelly Cho^{‡,6,14,18}, Michael J Gaziano^{‡,6,14,18}, Ravi K Madduri^{‡,7}, Scott M Damrauer^{‡,1,43,58,59}, Katherine P Liao^{‡,60,61}

Other sources of GWAS summary statistics

- Pan UK Biobank project -- <https://pan.ukbb.broadinstitute.org>
- FinnGen -- <https://r10.finngen.fi>
- All of US
- Biobank Japan
- Published, phenotype specific studies
 - Open GWAS project -- <https://gwas.mrcieu.ac.uk>
 - GWAS catalog -- <https://www.ebi.ac.uk/gwas/>

Penn Genetics Summer Short Course in Computational Genomics



Penn Medicine



VA
HEALTH CARE | Defining
EXCELLENCE
in the 21st Century