# Human Genetics Zeitgeist

Goncalo Abecasis

University of Michigan School of Public Health

@gabecasis

# Goal of Human Genetic Studies

Find biological processes that,
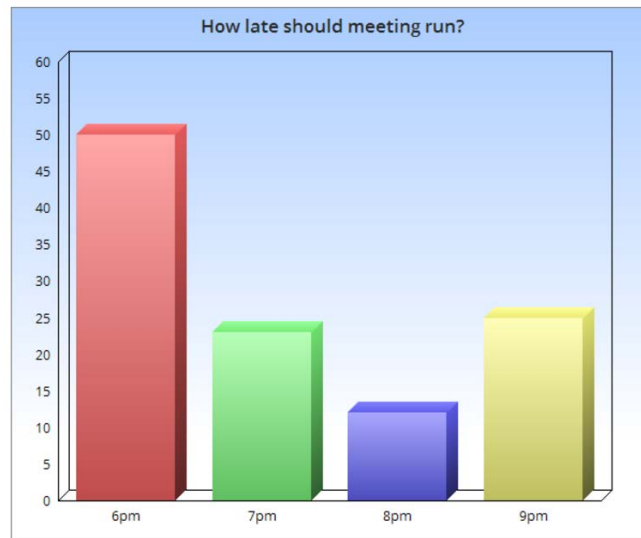when changed, alter disease course

**Understand Disease:**

Enable new treatments

**Predict disease:**

Enable early prevention, decision making

# Something to keep you awake.

- For the next bit of the talk, you will see some interspersed polls.

- To participate, go to:
  - On a browser:          Go to    **pollev.com/topmed**
  - On a phone:            Text      **topmed to 22333**

# Your poll will show here

**1**

**Install the app from**
**pollev.com/app**

**2**

**Make sure you are in**
**Slide Show mode**

Still not working? Get help at pollev.com/app/help
*or*
Open poll in your web browser

# Human Genetics, Sample Sizes over My Time

| Year | No. of Samples | No. of Markers | Publication |
|------|----------------|----------------|-------------|
| Ongoing | 120,000 | 600 million | NHLBI Precision Medicine Cohorts / TopMed |
| 2016 | 32,488 | 40 million | Haplotype Reference Consortium (Nature Genetics) |
| 2015 | 2,500 | 80 million | The 1000 Genomes Project (Nature) |
| 2012 | 1,092 | 40 million | The 1000 Genomes Project (Nature) |
| 2010 | 179 | 16 million | The 1000 Genomes Project (Nature) |
| 2010 | 100,184 | 2.5 million | Lipid GWAS (Nature) |
| 2008 | 8,816 | 2.5 million | Lipid GWAS (Nature Genetics) |
| 2007 | 270 | 3.1 million | HapMap (Nature) |
| 2005 | 270 | 1 million | HapMap (Nature) |
| 2003 | 80 | 10,000 | Chr. 19 Variation Map (Nature Genetics) |
| 2002 | 218 | 1,500 | Chr. 22 Variation Map (Nature) |
| 2001 | 800 | 127 | Three Region Variation Map (Am J Hum Genet) |
| 2000 | 820 | 26 | T-cell receptor variation (Hum Mol Genet) |

# A comprehensive review of genetic association studies

Joel N. Hirschhorn, MD, PhD[1–3], Kirk Lohmueller[1], Edward Byrne[1], and Kurt Hirschhorn, MD[4]

"… of the 166 associations which have been studied 3 or more times, only six have been consistently replicated."

Hirschhorn et al (2002)

# Your poll will show here

**1**

## Install the app from
## pollev.com/app

**2**

## Make sure you are in
## Slide Show mode
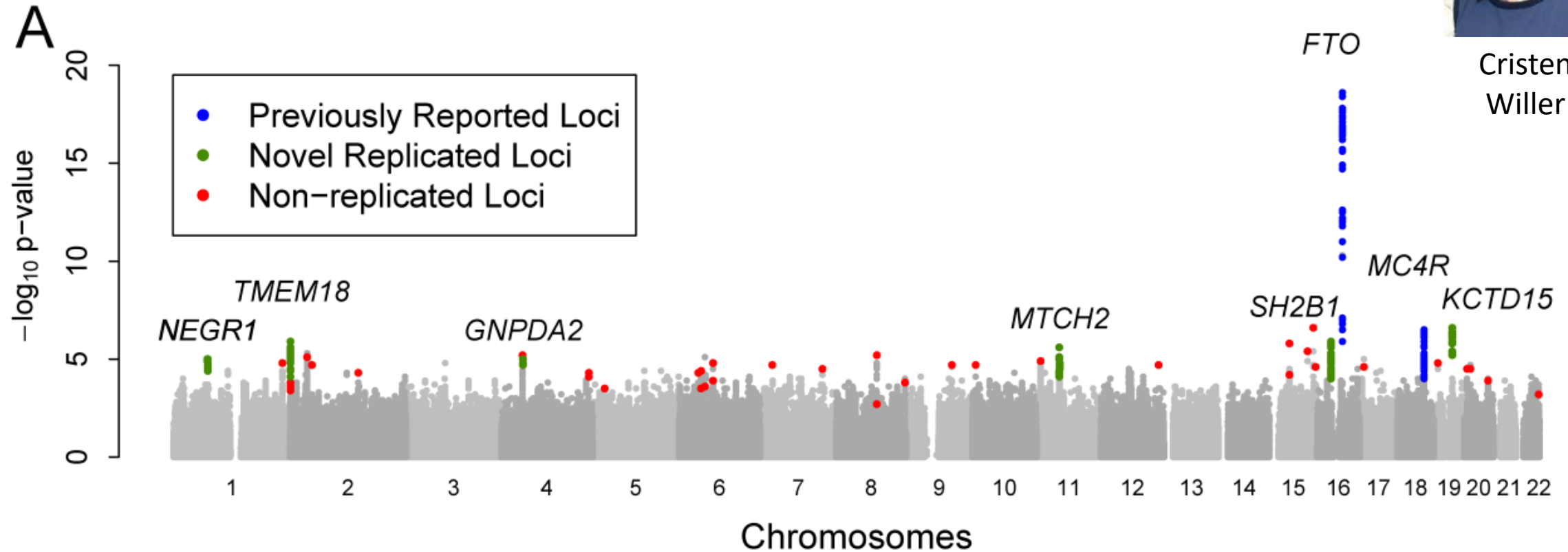
Still not working? Get help at pollev.com/app/help
*or*
Open poll in your web browser
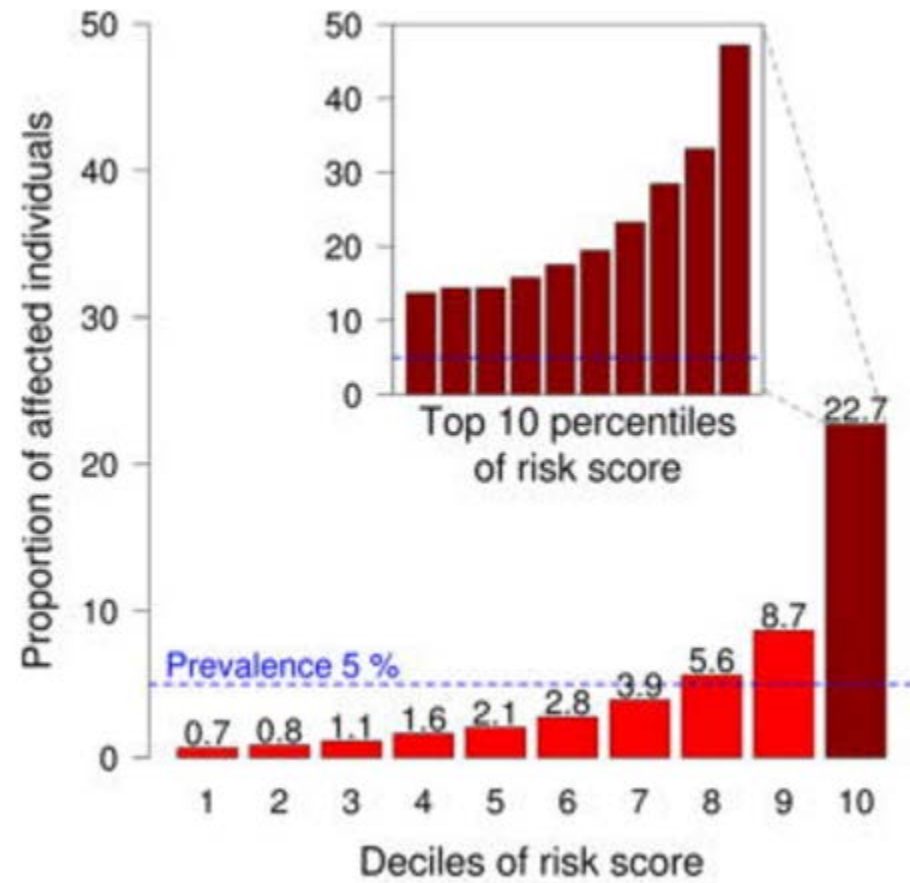
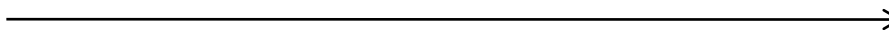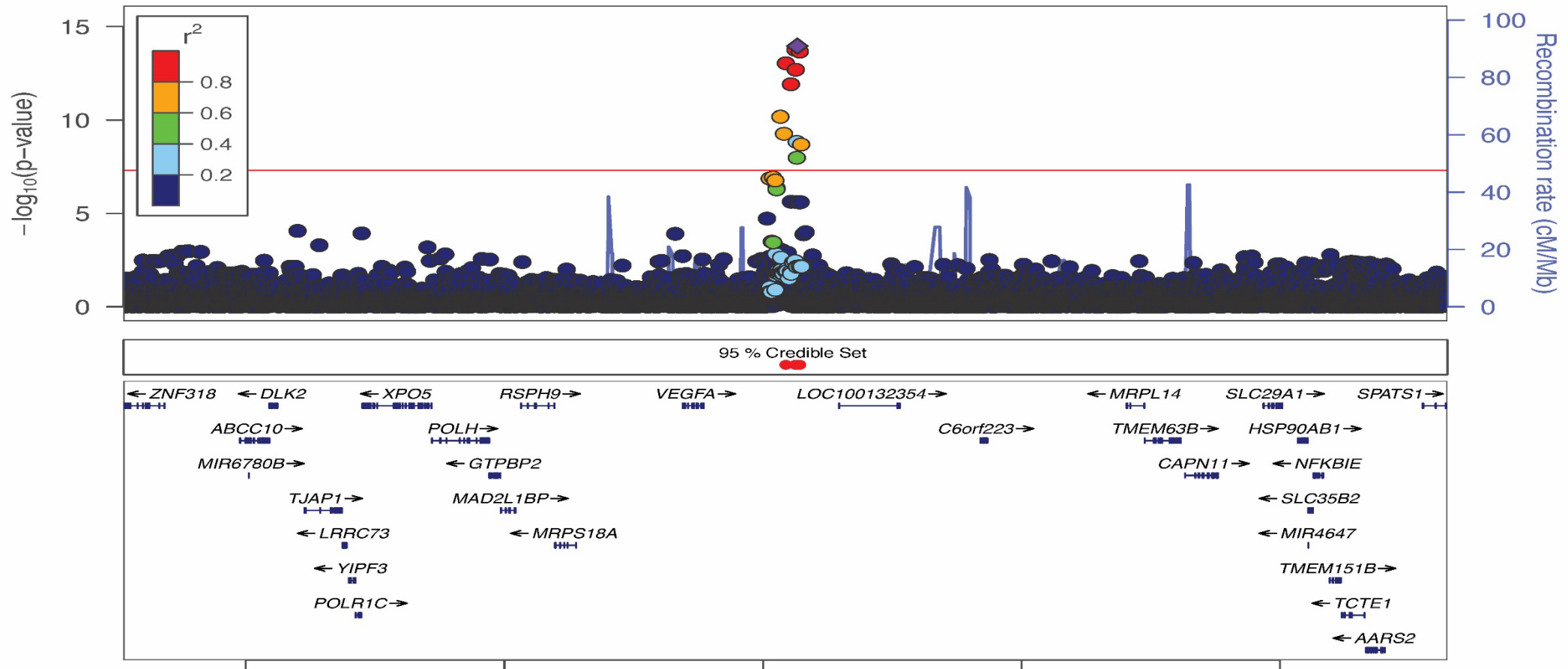# Search for Genetic Variants Influencing Body Mass



Cristen Willer

Seven of eight confirmed BMI loci show strongest expression in the brain…

Willer et al, *Nature Genetics,* 2009

# Macular Degeneration, 2010 - 2015



2010


Summary of Genome–Wide Scan Results for ~2.5 Million Imputed SNPs

2013



2015



As sample sizes grow, we have more loci.

How do we enable more scientists to explore data and explore big questions?

By default, easy to spend a lot of energy on the basics of data aggregation and initial analysis.

# Combined Effects of Many Alleles Strongly Predict Risk (2015)



Low risk → High risk

# Your poll will show here

**1**

**Install the app from
pollev.com/app**

**2**

**Make sure you are in
Slide Show mode**

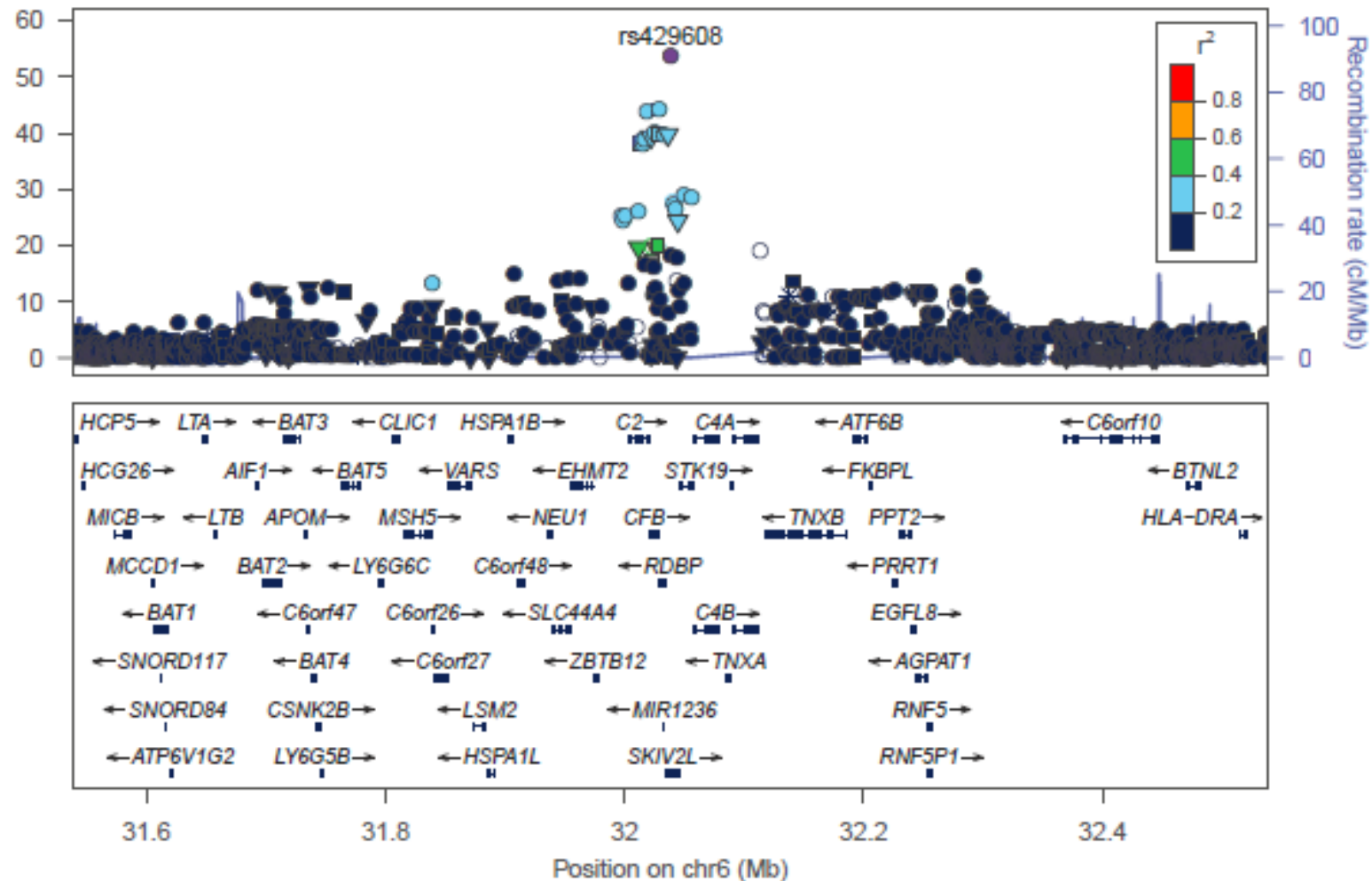Still not working? Get help at pollev.com/app/help
*or*
Open poll in your web browser

# Age-Related Macular Degeneration:
## Close-up near *VEGFA*

# Age Related Macular Degeneration:
## Close-Up of Specific Region
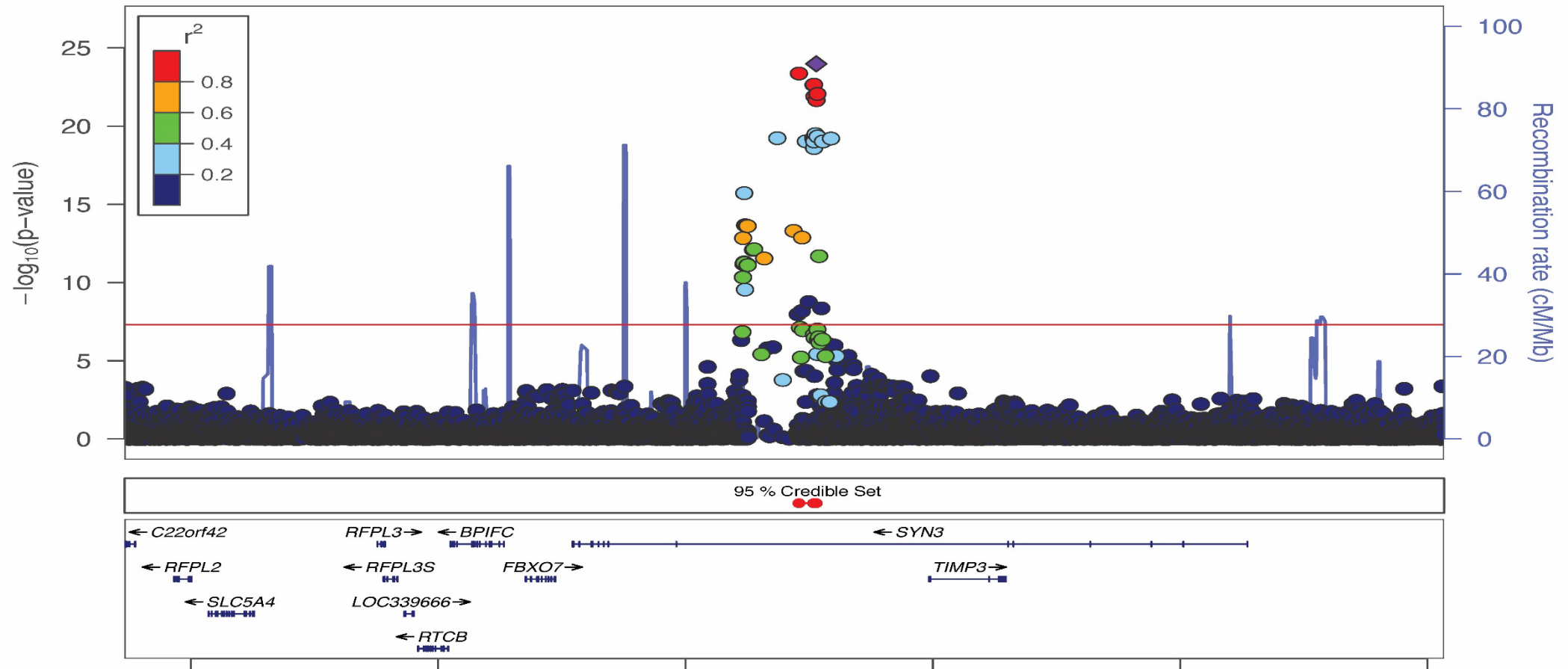
# What do AMD Associated Variants Have in Common?

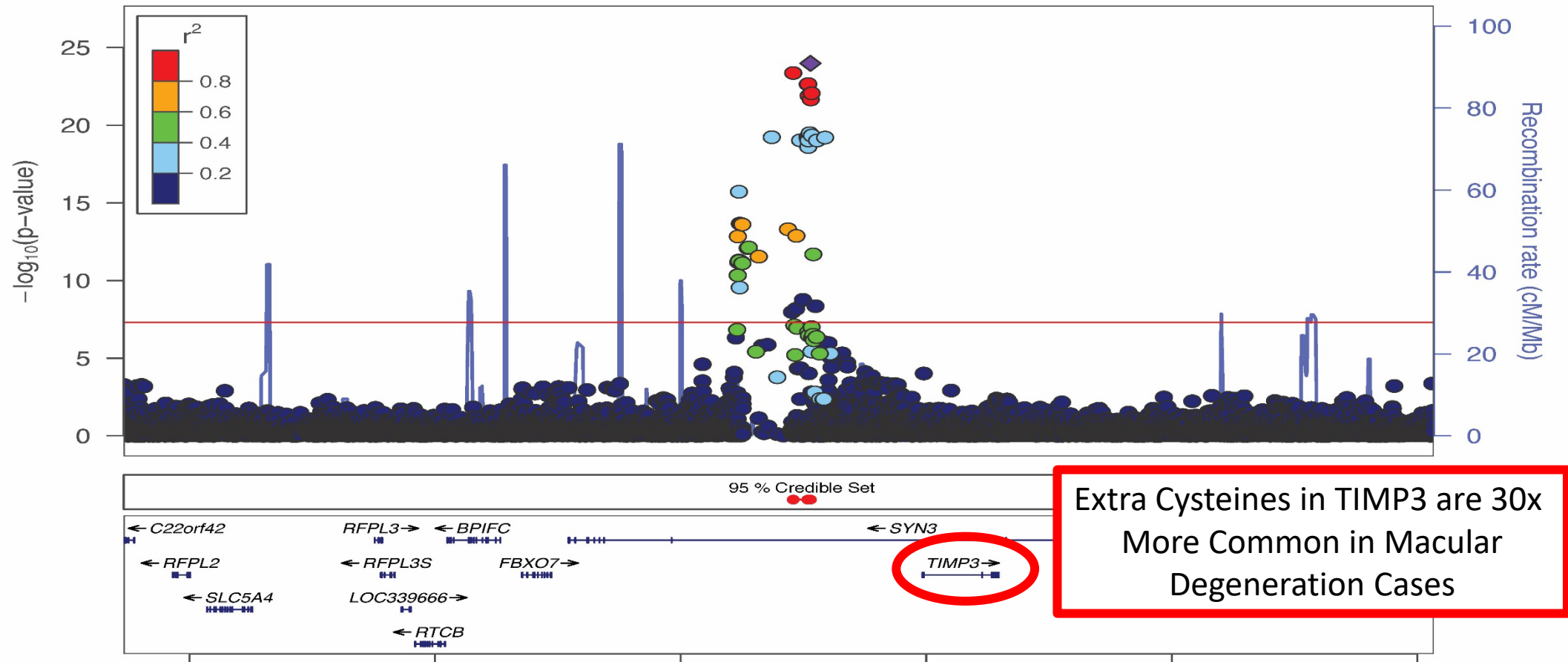|                        | Enrichment |
| ---------------------- | ---------- |
| Protein Coding         | 21x        |
| Near Complement Gene   | 183x       |
| Deleterious (CADD >20) | 255x       |

# Age-Related Macular Degeneration:
## Another Close-up, Chromosome 22

# TIMP3 Variants on AMD Array



8 of 10 Sorby's mutations cause unpaired cysteine residues
→ Polymerization of TIMP3 protein
→ Accumulation in extracellular matrix

# Your poll will show here

**1**

**Install the app from
pollev.com/app**

**2**

**Make sure you are in
Slide Show mode**

Still not working? Get help at pollev.com/app/help
*or*
Open poll in your web browser

# Age-Related Macular Degeneration:
## Another Close-up, Chromosome 22



Extra Cysteines in TIMP3 are 30x More Common in Macular Degeneration Cases

# Challenges

- How do we move faster from cataloguing loci to advancing biology?

- Engaging populations at the scale of 10,000s of individuals

- Sequencing at the scale of 10,000s of genomes

- Explore new technologies that accelerate functional analyses

- Make sure we don't get bogged down with basics
  - Simplify processes for running analyses we are good at
  - Simplify processes for trying new ideas on data

# How Great Analysts Contribute …

- Carry out top-notch analyses that point biology in the right direction

- New analysis tools and methods that scale, add value and meaning to data

- Enable new paradigms for collecting and sharing research data

- Expose data and analysis tools to broad community, including non-experts
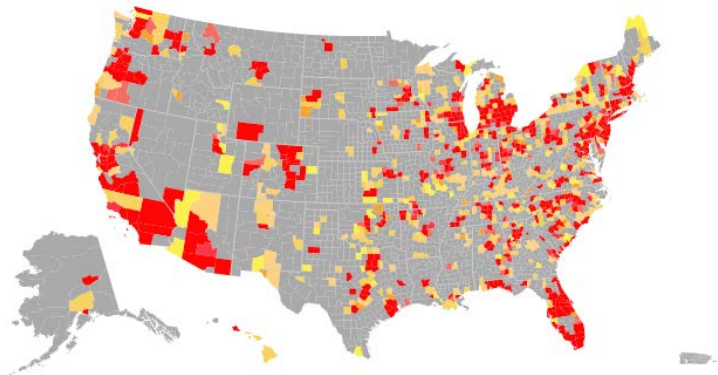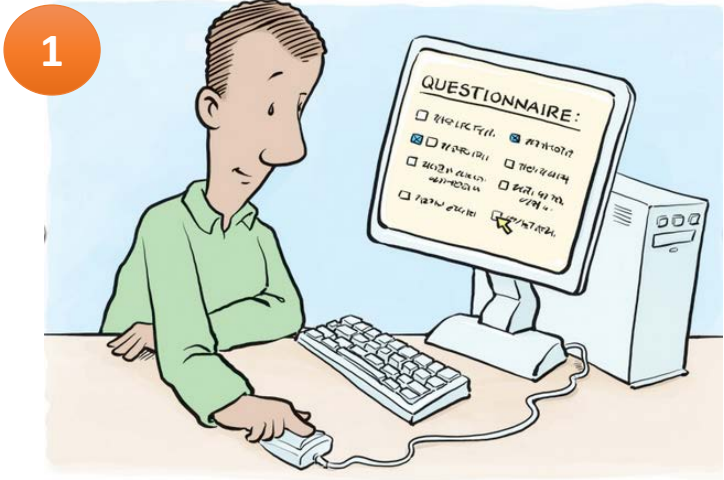  - Infuse high-quality health and genetic data in all research

# Making Data Available: Three Use Cases

- Non-technical users …
  - Need to access and understand results of large scale genomic studies
  - Interact with analysis and results on demand, often don't need individual data

- Statisticians & method developers
  - Need to compute over data without becoming big data computer scientists!
  - Use APIs to remix interesting analyses based on sufficient statistics

- Hard core method developers
  - Will want to access raw sequence data and develop low level computational methods
  - Need access to individual level data

- Three different needs and expectations
  - We should try to make some of these uses as close to zero friction as possible

# How Can We Engage 10,000s of Research Participants?

Part I – Genes for Good

# GENES for GOOD

1. QUESTIONNAIRE:
2. (DNA sample collection)
3. PROTECTIONS / DNA SAMPLE #273610
4. DAILY NEWS — ONLINE STUDY LEADS TO BIG BREAKTHROUGH MEDICAL FINDING

- Exploring new ways to engage populations in research

- Continuous Engagement, Web, Mobile Devices

- Currently, >25,000 participants

- www.genesforgood.org
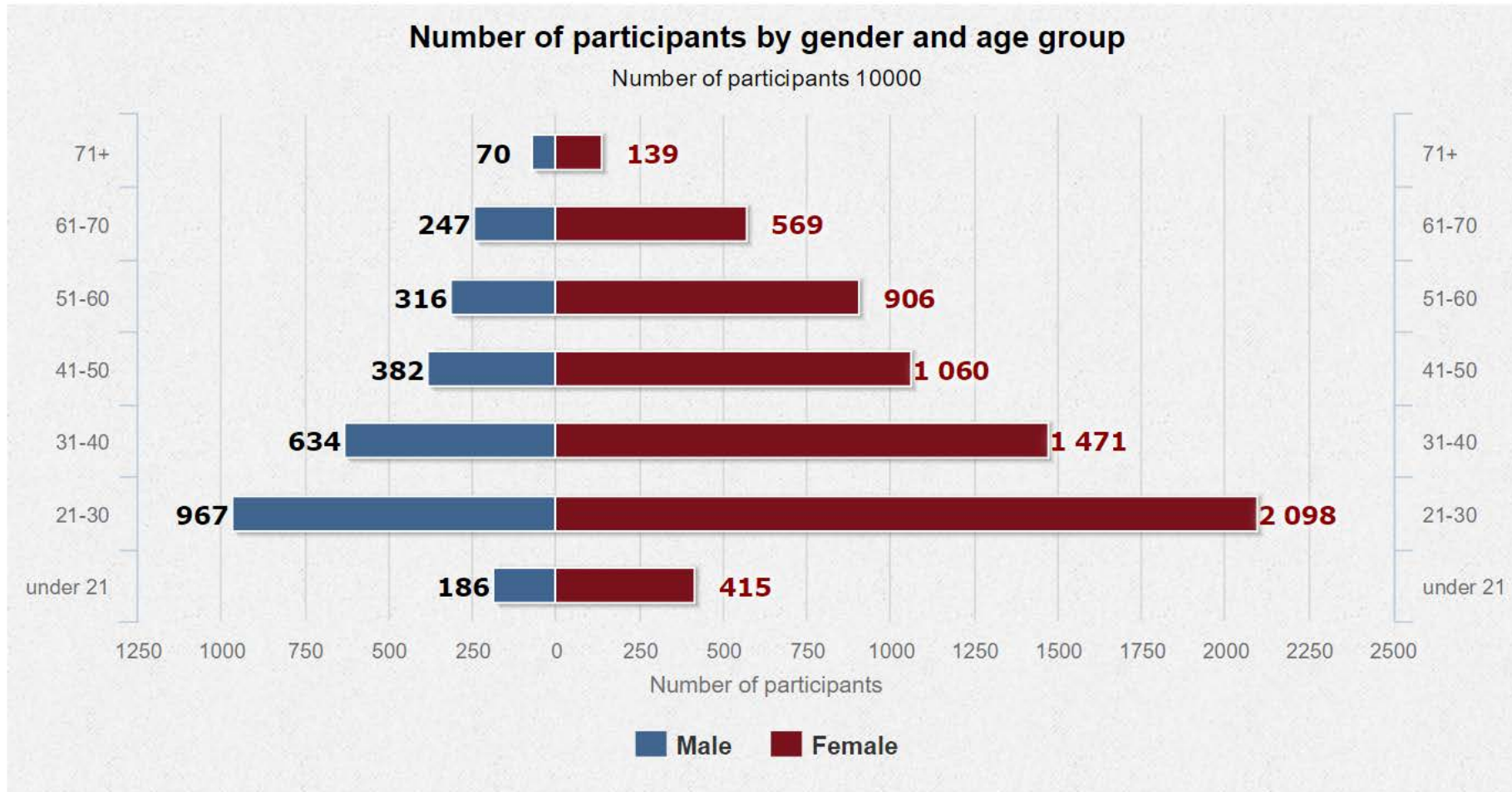
Genes for Good - Geographical Distribution of App Usage

Activity:
High

Low

2015-01-01          Participants: 0          Completed Surveys: 0

# 10,000 Participants...



Number of participants by gender and age group

# Return of Results

# Average Reported Sleep Hours Over One Year



2016

# BMI, Age & Diabetes



**Relationship of BMI with Diabetes Type 1 or 2**

Legend:
- GfG
- SHIELD
- NHANES

Y-axis: % with Diabetes (0, 5, 10, 15, 20, 25)

X-axis: BMI (< 18.5, 21.75, 26, 28.5, 32.5, 37.5, > 40.0)

Bays et al. (2007) *International Journal of Clinical Practice*

# Results: BMI GWAS

| Pheno | n | Chr:Pos | SNP | Gene | Our P | Other P* |
|-------|-----|---------|-----|------|-------|----------|
| BMI | 2,851 | 16:53803574 | rs1558902 | *FTO* | $5 \times 10^{-5}$ | $5 \times 10^{-120}$ |



*Speliotes et al. (2010) *Nature Genetics*

# Your poll will show here

**1**

## Install the app from
### pollev.com/app

**2**

## Make sure you are in
### Slide Show mode

Still not working? Get help at pollev.com/app/help
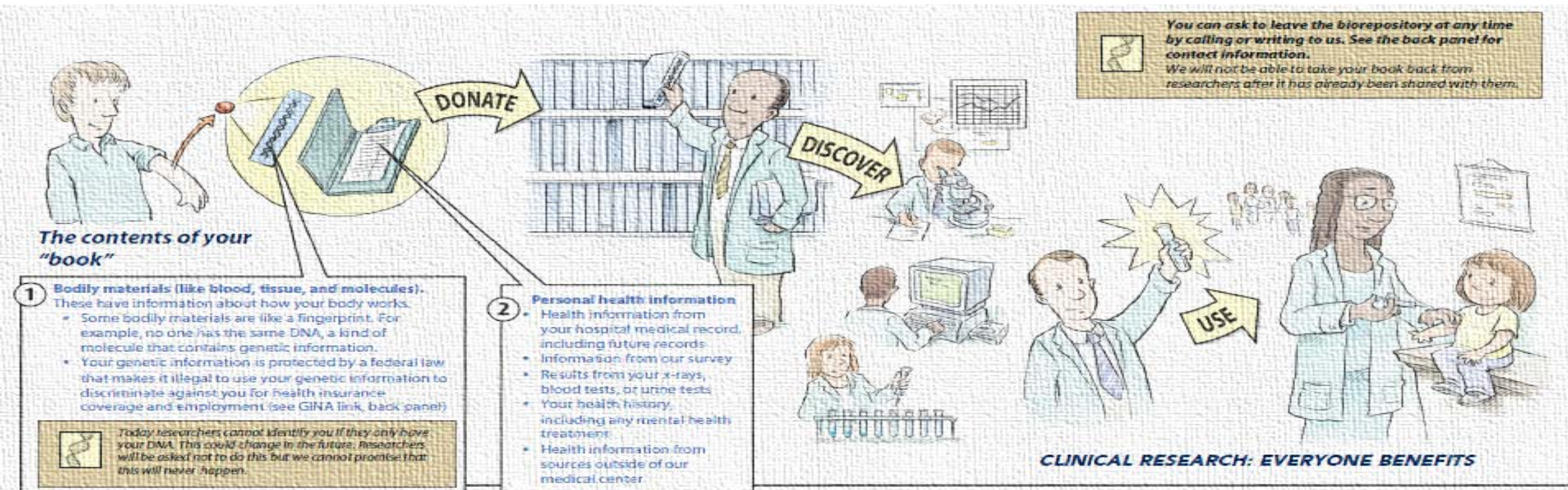*or*
Open poll in your web browser

# How Can We Engage 10,000s of Research Participants?

Part II – Michigan Genomics Initiative
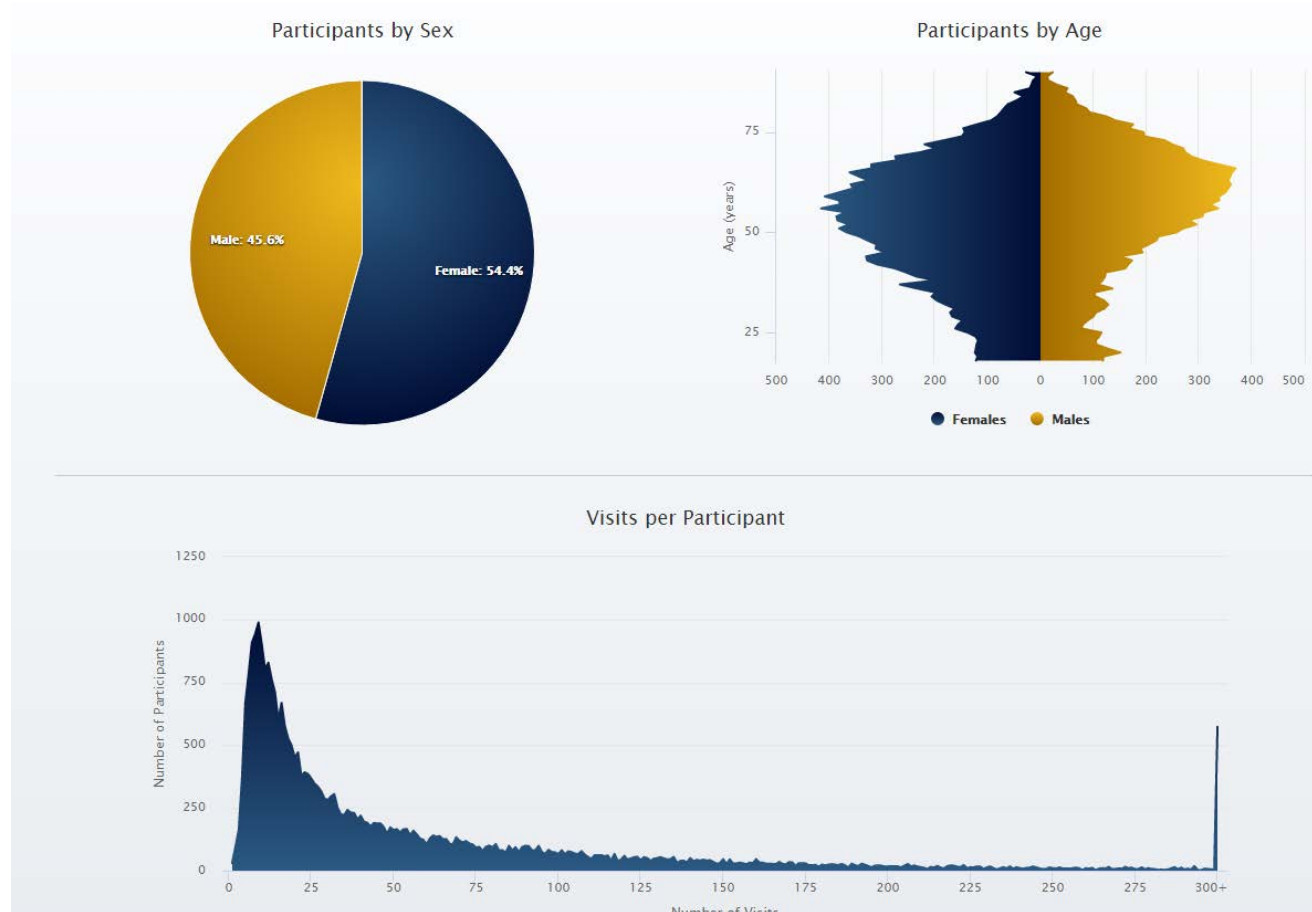
# Michigan Genomics Initiative

- Combine genetic and electronic health information on 40,000+ patients    50 new participants per day

- Use genetic information study many traits and diseases    Diverse traits – 40% w/cancer

- Build catalog of naturally occurring human knockouts    Speed and improve translation

- Clear, easy to understand consent – full participant buy-in.    Key for long term success

- **Team effort: Schmidt (Analysis), Ketherpal (Electronic Health Records), Brummett (Recruitment)**

# MGI demographics
http://www.michigangenomics.org/



Participants by Sex — Male: 45.6%, Female: 54.4%

Participants by Age

Visits per Participant

| Disease | % |
|---|---|
| Hypertension | 28% |
| Obesity | 24% |
| Arrhythmias | 22% |
| Sleep Apnea | 12% |
| Skin Cancer | 12% |
| Asthma | 8% |
| Cystic Fibrosis | 0.1% |

# Michigan Genomics Initiative (Freeze 1)
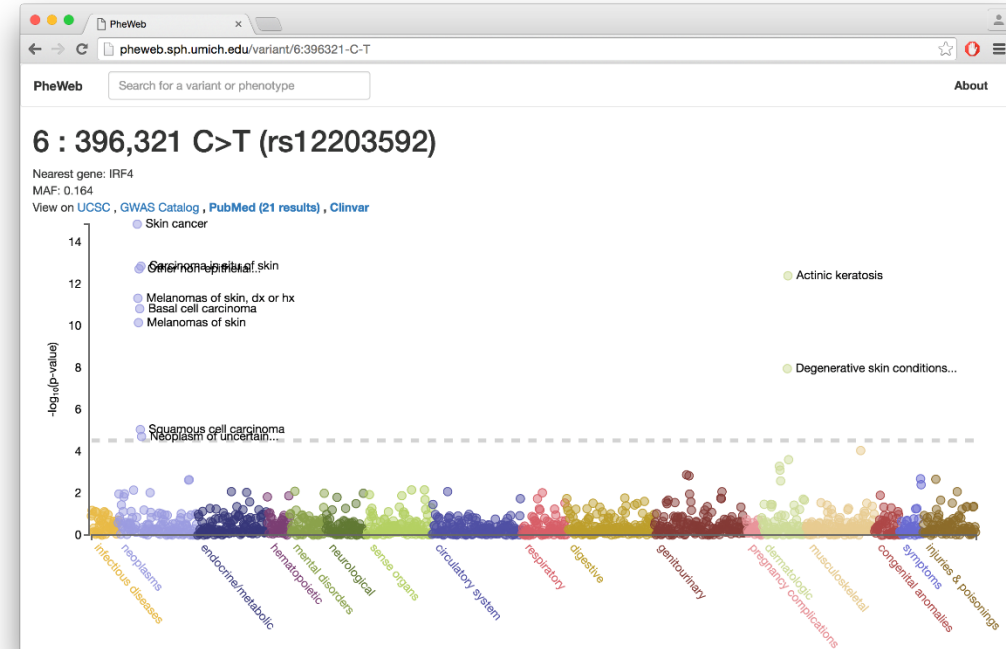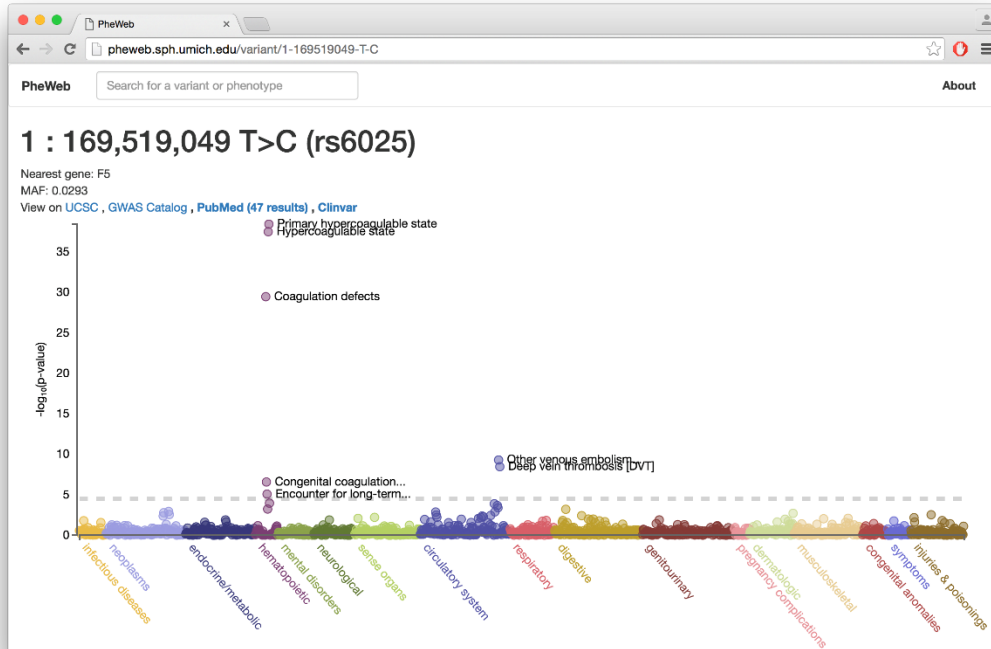## 20,000 individuals
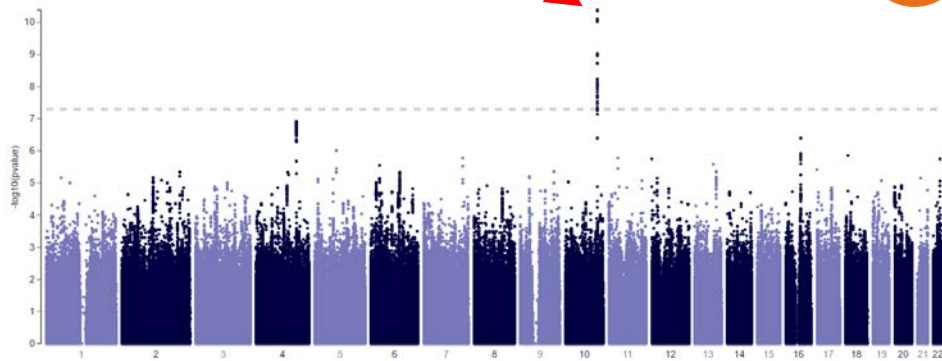## 7.5 million variants x 1,500 phenotypes

Ellen Schmidt

Peter VandeHaar

# Michigan Genomics Initiative Association Statistics
## http://pheweb.sph.umich.edu

# I heard rs10490924 in ARMS2 is associated with macular degeneration ...

# I heard rs738409 in PNPLA3 is associated with liver disease …

# I heard rs12203592 in IRF4 is associated with freckling, skin color …

# What signals map near IRF4?

Strongest association is with "Skin Cancer"

Other signals in this locus ...

| Top p-value | Phenotype | Top Variant |
|---|---|---|
| 6.7e-18 | Skin cancer | rs12203592 |
| 9.8e-15 | Actinic keratosis | rs12203592 |
| 3.1e-11 | Basal cell carcinoma | rs12203592 |
| 6.7e-11 | Melanomas of skin | rs12203592 |
| 3.3e-10 | Degenerative skin conditions and other dermatoses | rs12203592 |

# Federate!



Wouldn't it be nice to combine analysis without data use agreements and exchanging individual level data?

# PheWeb Goals

- Enable researchers to easily federate data

- Enable remixing interesting analysis without accessing individual data
  - Compute a novel association statistic
  - Retrieve association results for all variants in a set
  - Compute a new burden test for a gene or coding element
  - Carry out a Mendelian randomization analysis

- How?
  - Enable APIs to deliver intermediate algebra results that go into analyses

# Your poll will show here

**1**

**Install the app from pollev.com/app**

**2**

**Make sure you are in Slide Show mode**

Still not working? Get help at pollev.com/app/help
*or*
Open poll in your web browser

# The Scale of Genetic Data

# The 1000 Genomes Project (2008 – 2015)

# Shotgun Sequence Data

TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT

ATAGCTAG**A**TAGCTGATGAGCCCGATCGCTGCTAGCTC

ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC

AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG

GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

**A/C**

Predicted Genotype

# Variants per genome

| Type | Variant sites / genome |
|---|---|
| SNPs | ~3,800,000 |
| Indels | ~570,000 |
| Mobile Element Insertions | ~1000 |
| Large Deletions | ~1000 |
| CNVs | ~150 |
| Inversions | ~11 |

# Population histories

# Your poll will show here

**1**

**Install the app from pollev.com/app**

**2**

**Make sure you are in Slide Show mode**

Still not working? Get help at pollev.com/app/help

*or*

Open poll in your web browser

# Optimal Model for Analyzing 1000 Genomes?

| 1000 Genomes Call Set (CEU) | Reference Errors | Heterozygote Errors | Homozygous Non-Reference Errors |
|---|---|---|---|
| Broad | 0.66 | 4.29 | 3.80 |
| Michigan | 0.68 | 3.26 | 3.06 |
| Sanger | 1.27 | 3.43 | 2.60 |

# Optimal Model for Analyzing 1000 Genomes?

| 1000 Genomes Call Set (CEU) | Reference Errors | Heterozygote Errors | Homozygous Non-Reference Error |
|---|---|---|---|
| Broad | 0.66 | 4.29 | 3.80 |
| Michigan | 0.68 | 3.26 | 3.06 |
| Sanger | 1.27 | 3.43 | 2.60 |
| **Majority Consensus** | **0.45** | **2.05** | **2.21** |

- "Ensemble" outperforms the best method

# Challenges with the basic approach ...

⭐

ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG

5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

# Challenges with the basic approach ...

⭐

```
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGATGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG
```

5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

# Challenges with the basic approach …

⭐

```
                    CTAGATGATGAGCCCGATCGCTGCTAGCTC
                      AGATGATGAGCCCGATCGCTGCTAGCTCGA
                        GATGATGAGCCCGATCGCTGTTAGCTCGAC
                      AGATGATGAGCCCGATCGCTGCTAGCTCGA
                        ATGATGAGCCCGATCGCTGCTAGCTCGACG
                       GATGATGAGCCCGATCGCTGCTAGCTCGAC
                     AGATGATGAGCCCGATCGCTGCTAGCTCGA
                      GATGATGAGCCCGATCGCTGCTAGCTCGAC
              GCTAGCTAGCTGATGAGCCCGATCGCTGCT
           GATAGCTAGCTAGCTGATGAGCCCGCTCGC
                AGCTAGCTGATGAGCCCGATCGCTGCTAGC
              CTAGCTGATGAGCCCGATCGCTGCTAGCTC
          GCTGATAGCTAGCTAGCTGATGAGCCCGAT
        GATGCTAGCTGATAGCTAGCTAGCTGATGA
       GTCGATGCTAGCTGATAGCTAGCTAGCTGA
             TAGCTAGCTAGCTGATGAGCCCGATCGCTG
```

```
5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'
```

# Challenges with the basic approach …

# Variant Filtering

- Modern callers start with a candidate list of sites and annotate these …
  - Likely good sites: variants in HapMap or Omni 2.5M arrays
  - Likely problematic sites: variants that deviate from HWE or don't segregate in families

- Then, build a model that separates likely good sites from likely bad ones …
  - SVM, VQSR, self-organizing maps, ….

- Possible features …
  - What is the mapping quality of reads with the variant?
  - How many other differences in reads with the variant?
  - How many individuals are heterozygotes and homozygotes?
  - How many reads with the variant are on the forward and reverse strand?
  - What fraction of reads  have the variant in heterozygotes?
  - …

# The NHLBI TOPMed Program

- Trans-Omics for Precision Medicine

- Advance knowledge of heart, lung and blood disorders

- Add high-quality 'omics' data to high-priority studies
  - Whole genome sequencing currently executed at scale
  - Gene expression and metabolomics in pilot phases

- Data deposited in national databases, available for others to analyze

# TOPMed Sequencing as of February 15, 2017

http://nhlbi.sph.umich.edu/

- 68,503 genomes
  - 67,317 pass quality checks    (98.3%)
  - 823 flagged for low coverage   (  1.2%)
  - 358 fail quality checks        (  0.5%)

- Mean depth:        38.3x

- Genome covered: 98.7%

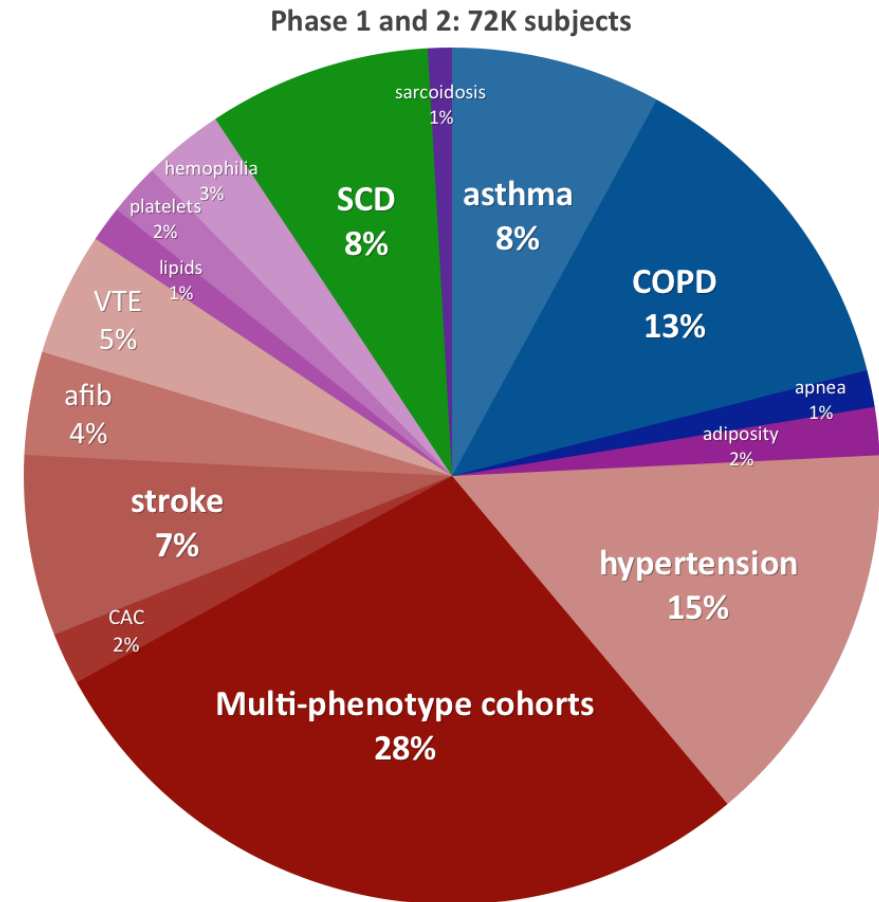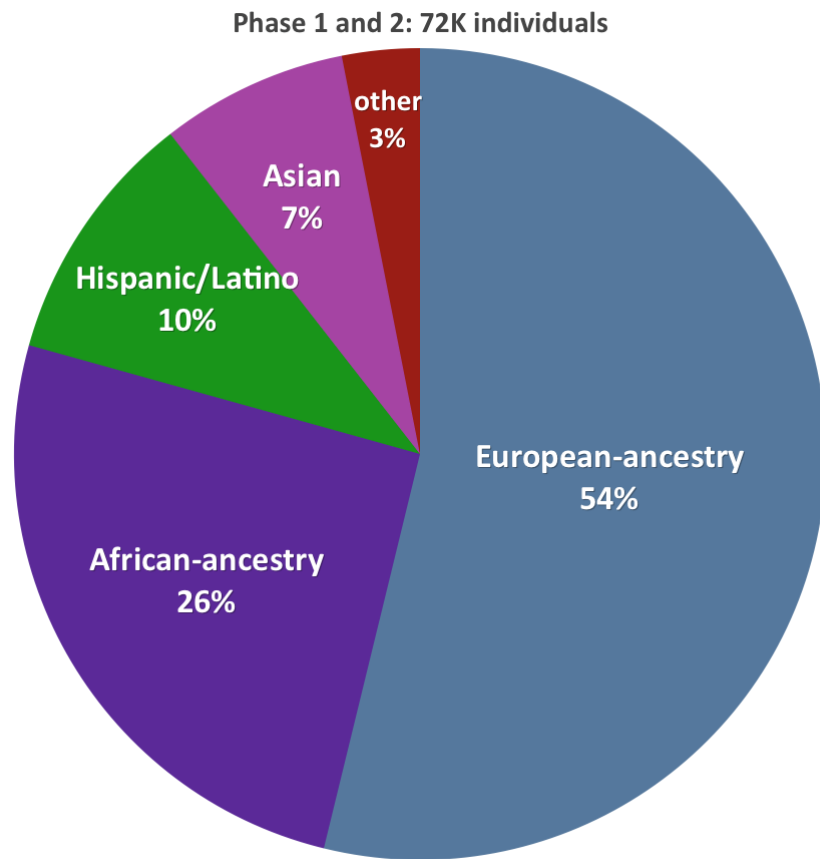- Contamination:     0.29%

- $9 \times 10^{15}$ sequenced bases

**Overall Genome Counts**

● Pass   ● Flag   ● Fail

67,317

# $9 \times 10^{15}$ sequenced bases



Number of snowflakes covering ~9 square miles in a 10-inch deep snowstorm.
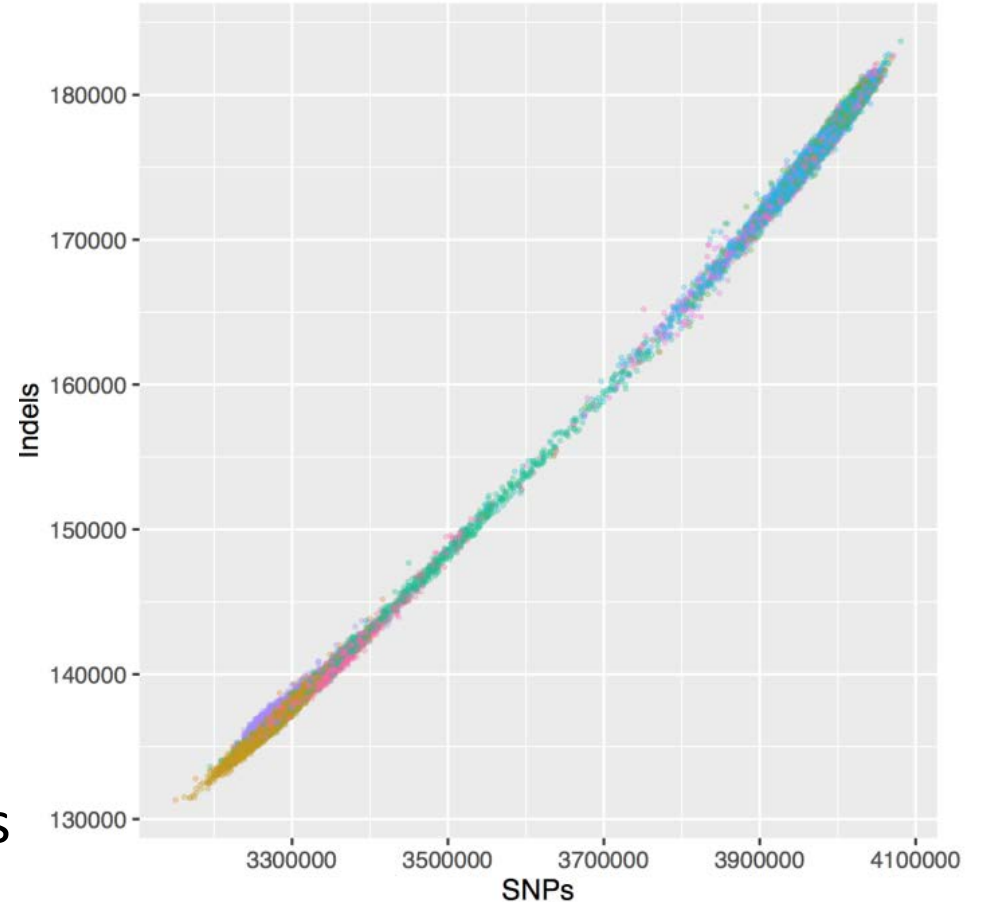100x more data than the 1000 Genomes Project.

# $9 \times 10^{15}$ sequenced bases



US corn production in 2014: $1.3 \times 10^{15}$ kernels

Image: Patrick Porter @ Smug Mug

# Ancestry and Focus Phenotypes in TOPMed



Phase 1 and 2: 72K individuals

Phase 1 and 2: 72K subjects

Quenna Wong and Cathy Laurie, University of Washington

# Current TOPMed Data Freeze

- 18,877 samples sequenced by early May
  - 4,047 individuals in 1,301 nuclear families

- 191 million SNPs
- 10.1 million indels

- Genotype VCF is very cumbersome
  - Extracting subsets of individuals can take days!
  - Post VCF formats are more compact
  - Post VCF formats not supported by analysis tools

# 1.6M Coding Variants

| Category | Count | Singletons |
|---|---|---|
| All SNPs | 191M | 43.5% |
| -- Missense SNPs | 1.5M | 47.9% |
| -- LoF SNPs | 39K | 55.5% |
| | | |
| All Indels | 10.1M | 43.2% |
| -- Inframe Coding Indels | 21K | 49.3% |
| -- Frameshift Indels | 31K | 59.2% |

# Browse All Variations Online
## http://bravo.sph.umich.edu

Peter VandeHaar

## KMT2D



## PCSK9



496 missense, 26 inframe indels, 0 stop or frameshifts
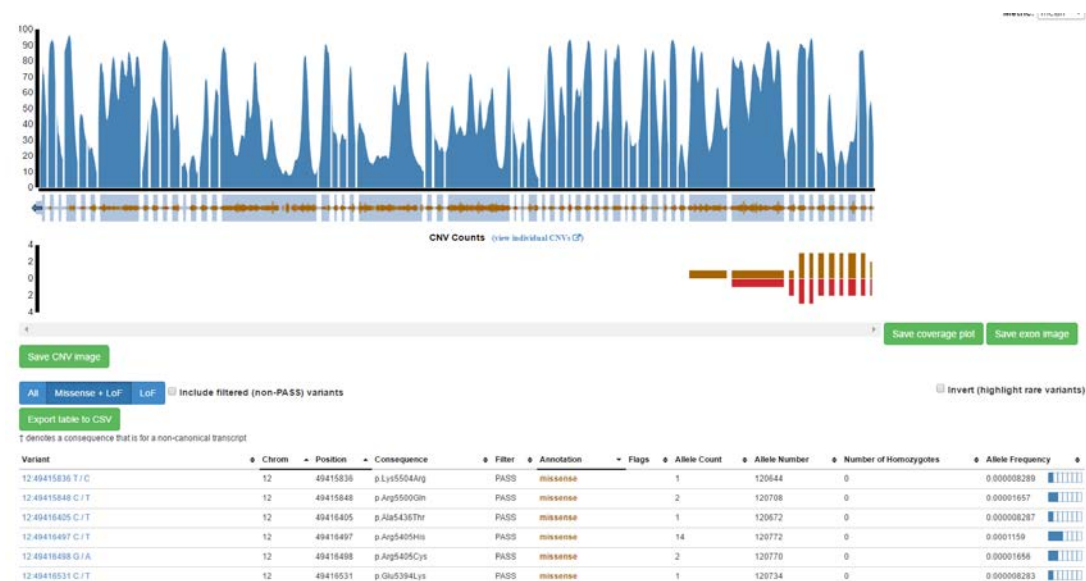
91 missense, 4 inframe indels, 7 stop or frameshifts

Federate!

KMT2D - BRAVO

KMT2D - ExAC

496 missense, 26 inframe indels, 0 stop or frameshifts

1842 missense, 23 inframe indels, 11 stop or frameshifts

# Functional Variants in Non-coding Regions?

| CADD Score | Variants | % Singletons |
|---|---|---|
| 0 – 9 | 149,808,517 | 40% |
| 10 – 14 | 9,481,664 | 42% |
| 15 – 19 | 3,168,122 | 43% |
| 20 – 24 | 885,525 | 45% |
| >= 25 | 334,964 | 50% |

- Starting genome-wide explorations for signatures of selection and function

- Strongest outlier regions (>100kb) in singleton proportion are near *HLA* and *ABO.*
- A few large examples of regions (>100kb) with high singleton proportions (e.g. *TP53BP1*)

# Functional Variants in Non-coding Regions?

| CADD Score | Variants | % Coding | % Singletons |
|:---:|:---:|:---:|:---:|
| 0 – 9 | 149,808,517 | 0.5% | 40% |
| 10 – 14 | 9,481,664 | 4% | 42% |
| 15 – 19 | 3,168,122 | 12% | 43% |
| 20 – 24 | 885,525 | 39% | 45% |
| >= 25 | 334,964 | 100% | 50% |

- Starting genome-wide explorations for signatures of selection and function

- Strongest outlier regions (>100kb) in singleton proportion are near *HLA* and *ABO.*
- A few large examples of regions (>100kb) with high singleton proportions (e.g. *TP53BP1*)

# How to help TOPMed advance discoveries?

- Genomewide analyses at scale are challenging

- Even simple analysis can require 1,000s of CPU days to complete

- Need to engage diverse teams in analysis and interpretation

```
snp,pvalue
rs1234,0.05
rs4343,0.0002
rs51101,0.61
rs981,0.000018
rs2223,0.72
```

# Plasma Lipids and Whole Genome Sequences

- Total Cholesterol, LDL, HDL, Triglycerides

- 8,394 TOPMed Participants
  - Jackson Heart Study
  - Framingham Heart Study
  - Amish Heart Study

- TOPMed Lipids Working Group
  - Leads: S. Kathiresan, C. Willer
  - PIs: A. Correa, A. Cupples, J. O'Connell, J. Wilson

Pradeep Natarajan

May Montasser

Maryam Zekavat

Gina Peloso

# How ENCORE works …

Matthew Flickinger    Jonathon LeFaive

# LDL Genomewide Analysis in ENCORE

# Mendelian Lipid Loci with LoF signals

| Phenotype | Gene | Cumulative LOF Frequency | Association P-value |
|-----------|------|--------------------------|---------------------|
| LDL | LDLR | 0.00024 | 0.009 |
| | APOB | 0.00061 | 0.00005 |
| | PCSK9 | 0.011 | $9 \times 10^{-28}$ |
| HDL | LCAT | 0.0004 | 0.035 |
| | ABCA1 | 0.0006 | 0.012 |
| | CETP | 0.001 | 0.001 |
| Triglycerides | APOC3 | 0.008 | $2 \times 10^{-19}$ |

Loci examined:

      LDL (**LDLR, APOB, PCSK9**, *LDLRAP1, ABCG5, ABCG8*)

      HDL (*APOA1*, **ABCA1, LCAT, CETP**, *LIPC, LIPG, SCARB1*)

      Triglycerides (*LPL, APOC2, APOA5*, **APOC3**, *GPIHBP1, LMF1, ANGPTL3, ANGPTL4*)

# TOPMed Production and Processing

# Your poll will show here

**1**

**Install the app from
pollev.com/app**

**2**

**Make sure you are in
Slide Show mode**

Still not working? Get help at pollev.com/app/help
*or*
Open poll in your web browser

# Sequencing on the Cheap: Imputation

**Observed GWAS Genotypes**

. . . . A . . . . . . A . . . . A . . . .
. . . . G . . . . . . C . . . . A . . . .

**Reference Haplotypes (e.g. 1000G)**

C G A G A T C T C C T T C T T C T G T G C
C G A G A T C T C C G A C C T C A T G G
C C A A G C T C T T T C T T C T G T G C
C G A A G C T C T T T C T T C T G T G C
C G A G A C T C T C C G A C C T T A T G C
T G G G A T C T C C G A C C T C A T G G
C G A G A T C T C C G A C C T T G T G C
C G A G A C T C T T T T C T T T T G T A C
C G A G A C T C T C C G A C C T C G T G C
C G A A G C T C T T T T C T T C T G T G C

# How long does it take to impute one genome?

- Depends on the reference size …

- 2007:       60 samples, 2.5M SNPs            14 min
- 2009:       60 samples, 7.3M SNPs            41 min
- 2011:       283 samples, 11.6M SNPs         1,287 min
- 2012:       381 samples, 37.4M SNPs         7,800 min
- 2015:       33,000 samples, 37M SNPs        63,000,000 min (estimated)
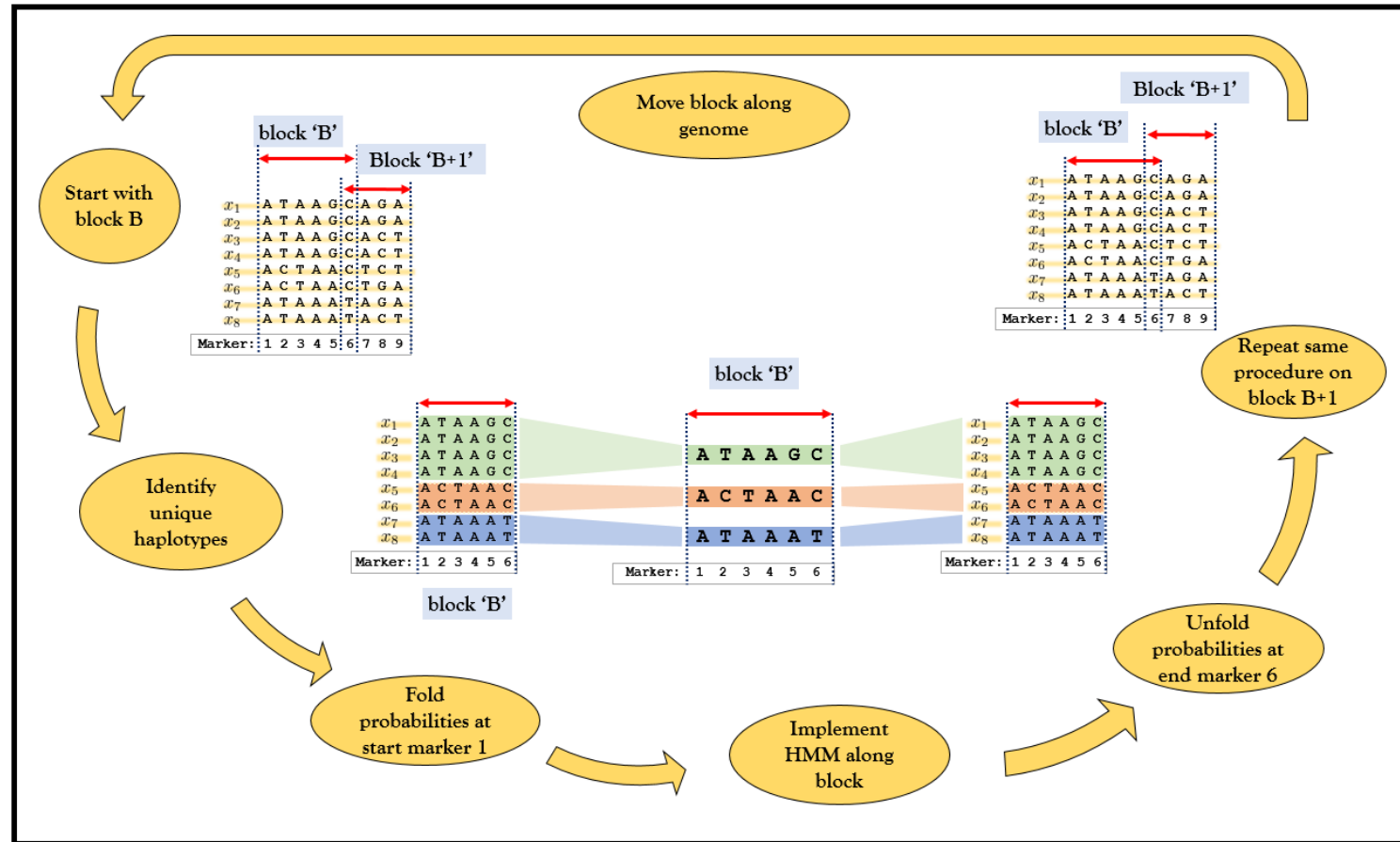
# How long does it take to impute one genome?

- Depends on the computational methods …

- 2007 software:    381 samples, 37.4M SNPs        7,800 min    O(MH$^2$)
- 2010 software:    381 samples, 37.4M SNPs        512 min    O(MH)
- 2012 software:    381 samples, 37.4M SNPs        24 min    O(MH)
- 2015 software:    381 samples, 37.4M SNPs        1 min    <O(MH)
- 2016 software:    381 samples, 37.4M SNPs        <5 secs    <O(MH)

# Most Recent Imputation Improvements Minimac3

# Imputation Servers

https://imputationserver.sph.umich.edu

# Imputation Accuracy w/TOPMed

Sayantan Das

- Imputation Accuracy in EUR
  - $r^2 > 0.85$ at 1% frequency
  - $r^2 > 0.75$ at 0.2% frequency

- Outperforms alternatives, including
  - 1000 Genomes (2,500 diverse genomes)
  - HRC (33,000 mainly European genomes)

- Working to define set of TOPMed samples that can be included in imputation panel

# TOPMed Data Resources

- 1st release of sequenced genomes and phenotypes now in dbGAP/SRA
  - >8,000 genomes now available without embargo

- Browsable Variant Catalog at http://bravo.sph.umich.edu
  - 169,454,024 variants in 14,559 individuals
  - Variant lists will also be deposited in dbSNP

- Track sequence data production at http://nhlbi.sph.umich.edu

- Imputation resource that will improve ability to reconstruct diverse genomes affordably planned

# Many new directions and opportunities …

- New techniques for exploring genomic function at scale
  - Use sequencing to measure enhancer activity after mass mutagenesis
  - Patwardhan et al (*Nature Biotechnology*, 2012)

- New techniques for dissecting biology of single cells
  - Use sequencing to profile expression in individual cells
  - Macosko et al (*Cell*, 2015)

- Open access resources like UK Biobank

# Your poll will show here

**1**

**Install the app from**
pollev.com/app

**2**

**Make sure you are in
Slide Show mode**

Still not working? Get help at pollev.com/app/help
*or*
Open poll in your web browser

# TOPMed Phase 1 and 2 – Acknowledgments

**Sickle Cell Disease (Boston-Brazil)**
Vijay G. Sankaran

**Asthma in African Descent Populations**
Kathleen C. Barnes

**Samoan Family Obesity Study**
Stephen T. McGarvey

**Atherosclerosis Risk Study**
Rasika Mathias

**PharmHU**
Eric Boerwinkle

**Cardiovascular Disease Risk**
Joanne E. Curran, David C. Glahn

**Cleveland Family Study**
Susan Redline

**Atherosclerosis Risk Study and VTE**
Eric Boerwinkle

**GOLDN**
Donna K. Arnett

**Race and Ethnic Disparity in Asthma**
Esteban González Burchard

**San Antonio Family Heart Study**
John Blangero

**Taiwan Study of Hypertension**
D. C. Rao

**Severe COPD Gene**
Edwin K. Silverman

**Early Atrial Fibrillation**
Patrick T. Ellinor

**Genotyping for Hemophilia**
Barbara Konkle

**GenNet Study of Salt Sensitivity**
Jiang He

**UW Northwest Genomics Center**
Debbie Nickerson

**Genes Influencing LDL Cholesterol**
Braxton D. Mitchell

**Severe Asthma Research Program**
Deborah A. Meyers

**Sarcoidosis in African Americans**
Courtney Montgomery

**Baylor HGSC**
Richard Gibbs

**Framingham Heart Study**
Vasan S. Ramachandran

**Women's Health Initiative**
Charles Kooperberg

**Walk-PHaSST**
Mark Gladwin

**Broad Institute**
Stacey Gabriel

**Genetic Epidemiology of COPD**
Edwin K. Silverman

**MESA Atherosclerosis Study**
Jerome I. Rotter, Stephen S. Rich

**REDS-III Brazil**
Brian Custer, Shannon Kelly

**NY Genome Center**
Soren Germer

**Asthma in Costa Rica**
Scott T. Weiss

**AA Coronary Artery Calcification**
Kent Taylor

**Sickle Cell Disease (OMG)**
Allison Ashley-Koch

**DCC**
Bruce Weir

**Jackson Heart Study**
Adolfo Correa, James Wilson

**HyperGEN and GENOA**
Donna K. Arnett

**Sickle Cell Disease (HW)**
Sergei Nekhai

**IRC**
Gonçalo Abecasis

# Acknowledgements

Thank you!
Michigan Team

- National Institutes of Health
  - **NHLBI**
  - NHGRI
  - NCBI

- TOPMed Study Investigators
- TOPMed Data Coordinating Center
- TOPMed Sequencing Centers

- TOPMed Lipids Working Group