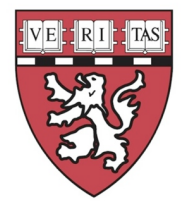# Software tools for genetic analysis

Benjamin Neale, Ph.D.
Analytic and Translational Genetics Unit, MGH
Stanley Center for Psychiatric Research & Program in Medical and
Population Genetics, Broad Institute

# First looks at the CCDG freeze 1 callset
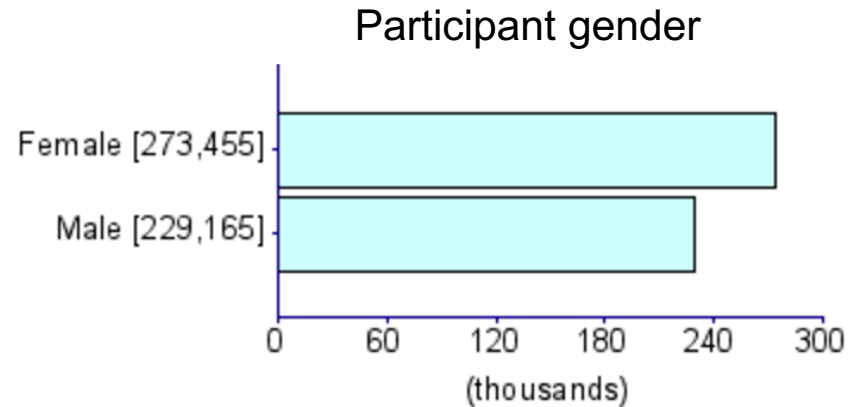
# UK Biobank!

# The team



- Liam Abbott
- Verneri Anttila
- Krishna Aragam
- Jon Bloom
- Sam Bryant
- Claire Churchhouse
- Joanne Cole
- Mark J. Daly

- Andrea Ganna
- Steven Gazal
- Jackie Goldstein
- Mary Haas
- Joel Hirschhorn
- Daniel Howrigan
- Sekar Kathiresan
- Dan King

- Duncan Palmer
- Tim Poterba
- Manuel Rivas
- Cotton Seed
- Sailaja Vedantam
- Raymond Walters

# UK Biobank Overview

- Prospective cohort of 502,620 participants
- Recruited across the UK 2006-2010
- Healthy volunteer effect - healthier, leaner, smoke less than population average

Age at recruitment

Mean 56.5 years (S.D. 8.1)

Participant gender

Female [273,455]

Male [229,165]

(thousands)

# Phenotyping

- Self-reported demographics (500,00), diet and exercise habits (211,000)

- Physical (500,00) and cognitive measurements (190,000)

- Imaging of brain, heart, abdomen (1,000 - 13,700)

- Blood, saliva and urine assays (500,00)

- Medical records (392,00) and cancer registers (80,000)

# Data Showcase http://biobank.ctsu.ox.ac.uk/crystal/

# Data Showcase http://biobank.ctsu.ox.ac.uk/crystal/

# Genotypic Data



UK Biobank Axiom array content
~825K markers total

**Markers relevant to specific phenotypes**
~45K

Alzheimer's disease
Autoimmune/inflammatory phenotypes
Blood phenotypes
Cancer (common and rare variants)
Neurological disease
Pharmacogenetics (ADME)
Cardiac disease
Cardiometabolic phenotypes
Lung function phenotypes

**Markers within genomic regions of interest**

Known GWAS loci (NHGRI GWAS catalog)
Expression quantitative trait loci (eQTLs)
Mitochondria
Y chromosome
Human leukocyte antigen (HLA) system
Killer-cell immunoglobulin-like receptor (KIR)
Apolipoprotein E (ApoE) gene
Neanderthal ancestry markers
~47K

**Genome-wide coverage for improved performance of array-based imputation**
~630K

Common variants
(MAF ≥ 5% in a European sample)

Low frequency variants
(1% < MAF < 5% in a European sample)

**Coding variation**
~125K

Protein truncating variants
Other rare coding variants
Rare, possibly disease causing, mutations

825,000 markers genotyped

Imputation to HRC
Plus 1KG + UK10K in works
96 million variants

Bycroft et al. on biorxiv

# Genome-wide association analysis of 2,400 traits!

How did we do it?

# Team hail

Alex Bloemendal Comp Bio

Jon Bloom
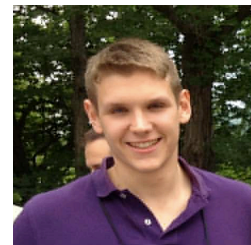Comp Bio



Cotton Seed
Software Eng
Team lead

Jackie Goldstein
Software Eng

Amanda Wang
Software Eng

Tim Poterba
Software Eng

Daniel King
Software Eng

# Scalability



10hr

10 compute hours

# Scalability

# Scalability

# Creating a Resource of GWAS Results

- **Enhance** the value of UK Biobank

- **Public** and **easily accessible**
  - www.nealelab.is/blog

- **Blog posts and code** available on GitHub
  - Remove time lag seen with publication
  - Better reflect updates and developments
  - Share how we performed analyses
  - Publish on novel methods and downstream analyses



RAPID GWAS OF THOUSANDS OF PHENOTYPES FOR 337,000 SAMPLES IN THE UK BIOBANK

# Association results for taking cholesterol lowering meds



**LDL Cholesterol**

GLGC Nat Genet

# Online resources

About   Downloads   Terms   Contact   Jobs   FAQ

## gnomAD browser beta | genome Aggregation Database

Search for a gene or variant or region

Example - Gene: PCSK9, Variant: 1-55516888-G-G.

### Ab

The Genome Aggregation Database (gnomAD investigators, with the goal of aggregating an from a wide variety of large-scale sequencing scientific community.

The data set provided on this website spans sequences from unrelated individuals sequen genetic studies. The gnomAD Principal Inves current release are listed here.

All data here are released for the benefit of th

---

**IBD Exomes Browser**

About   Downloads   Terms   Contact   FAQ

## Inflammatory Bowel Disease Exomes Browser

Search for a gene or variant or region

Examples - Gene: NOD2, Transcript: ENST00000407236, Variant: 16-50745656-G-A, RS ID: rs104895438, Region: 16:50727514-50766988

### Meta-an

This browser presents IBD case-cont cohorts, each of which contains samp tallies of samples included in the curr

- Ashkenazi Jews (AJ): 2641 IBD
- Non-Finnish Europeans (NFE): 4
- Finland (FINN): 696 IBD (210 CD

A full listing of collaborative partners

---

**Global Biobank Engine**

About   Downloads   Terms   Contact   HLA Alleles   Power   Genetic correlation   Decomposition

G Select Language ▼

## Global Biobank Engine (pre-alpha)

Search for a gene or variant or region or phenotype coding (coming soon)

Examples - Gene: F5, Transcript: ENST00000367797, Variant: 1:169519049, RS ID: rs6025, Region: 10:114686614-114786614

Recent News   Lab manuscripts

## Genetic Association Results

**Note:** We present summary statistic results from the UK Biobank hospital in-patient health-related outcomes summary information data (Data-Field 41202); computational grouping of phenotypes with

# Genetic parameters page



## Genetic Parameters: Genetic correlation

Version 0.01

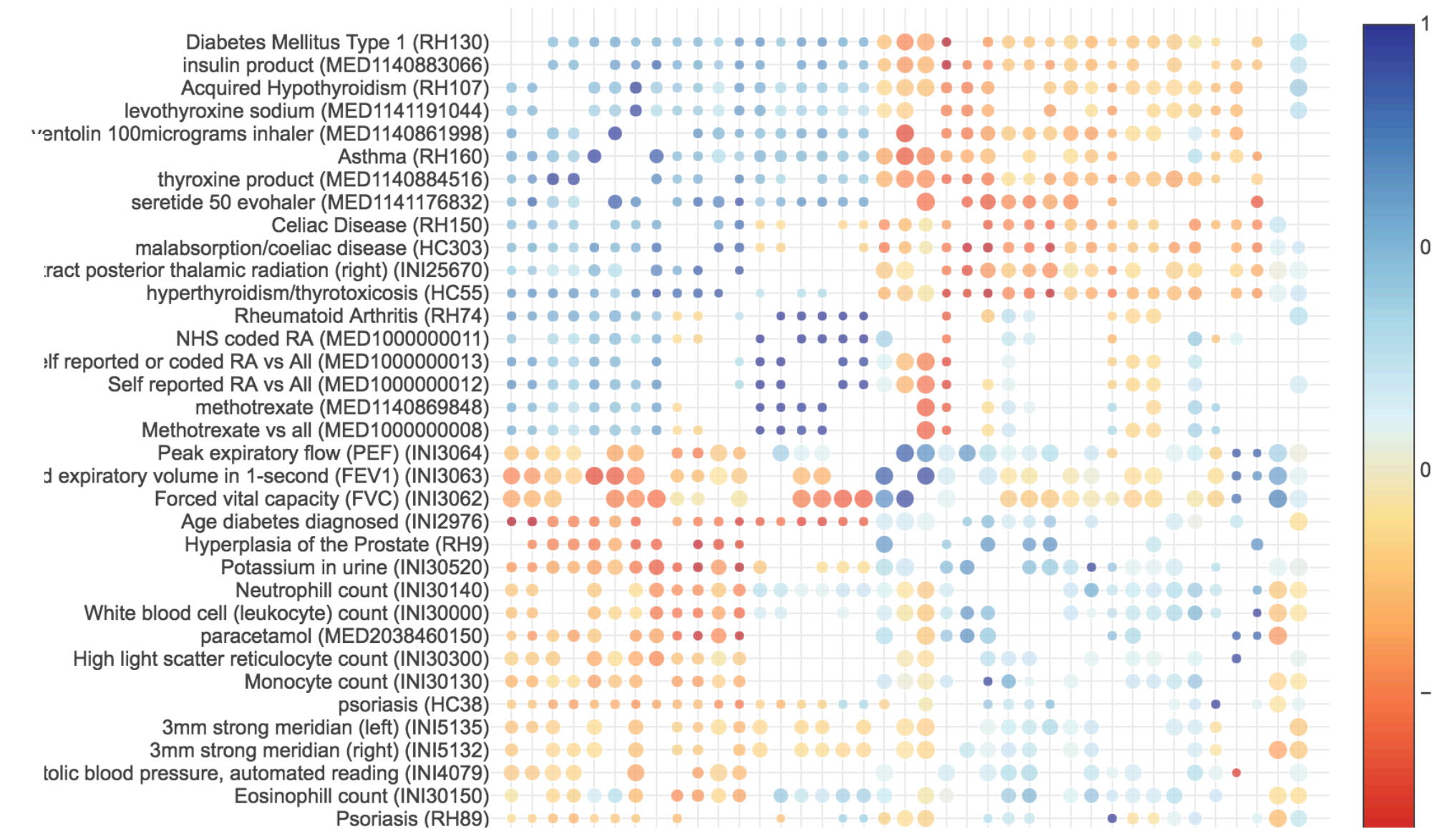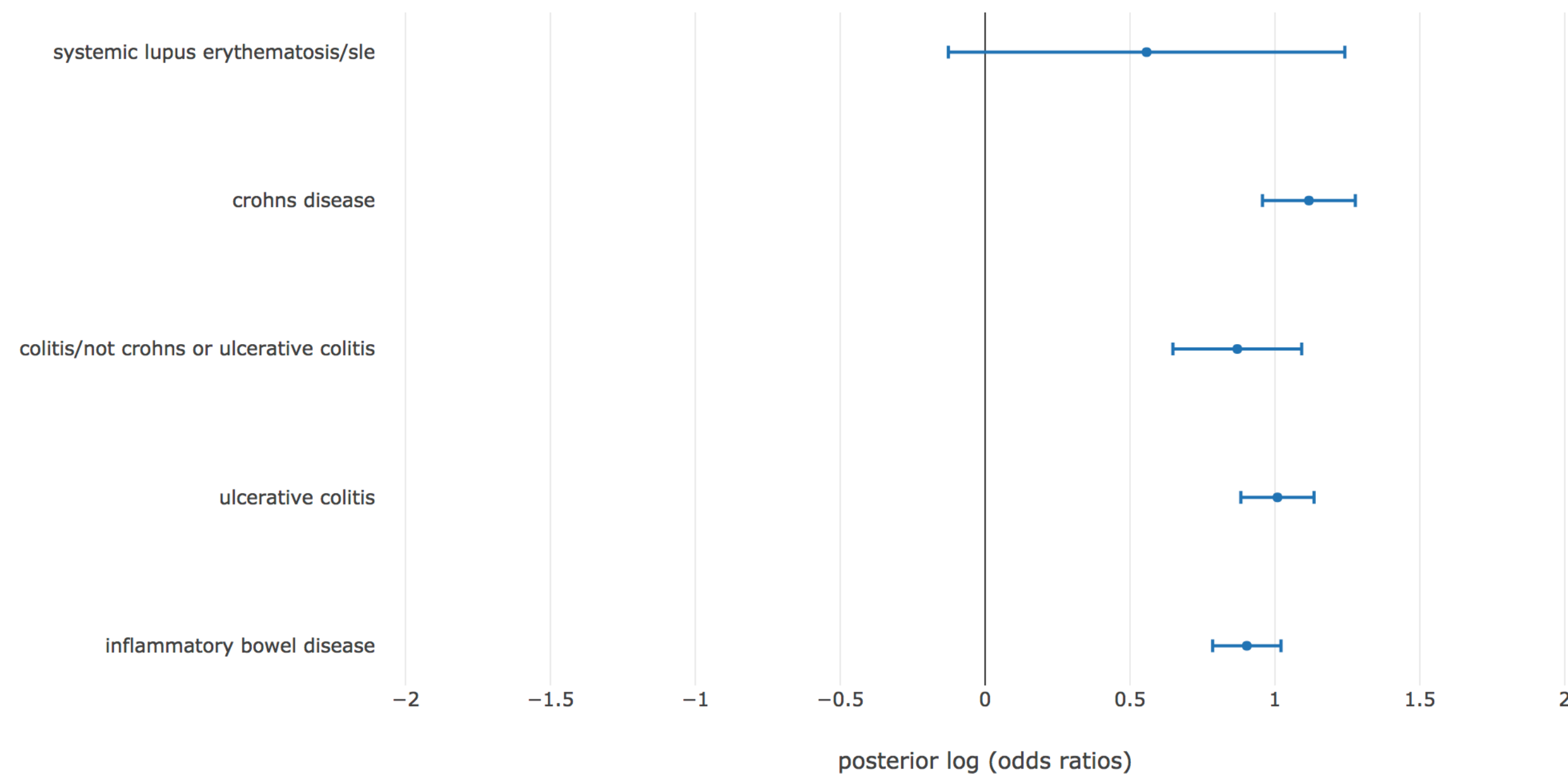Christopher DeBoever and Manuel A. Rivas

C.D. developed the web application and M.A.R designed the study.

This app displays the genetic correlation and other genetic parameters for various phenotypes from the UK Biobank.

**HLA Map**

Posterior log(odds ratios) for phenotypes with posterior probability >0.7 for DRB1*01:03

# PheWAS page



**"PheWAS" page**

log(OR) or BETA of PheWAS (chr2:163124051)

# Rare variant aggregate analysis



Code Phenotype: asthma

| | |
|---:|:---|
| **Code:** | HC382 |
| **Case Count or N:** | 43626 |
| **Single variant results:** | HC382 |
| **Phenotype distribution:** | HC382 |

Note: only genes with a log10BF greater than 1 are included in the manhattan plot with a hyperlink.

gene: IL33; log10(Bayes Factor): 9.41

gene: IL18RAP; log10(Bayes Factor): 3.65

**Rare variant aggregate analysis**

# Goals of Annotation

- **Annotate the genome comprehensively to**
  - Boost the power of rare variant association analysis
  - Assist with fine mapping to identify possible causal variants
  - Study population genetics
  - Improve the interpretability of whole-genome data.

- **Annotations at different levels:  SNV, SV, Indel, gene**

- **Annotate CCDG and CMG data, e.g. CCDG Freeze 1 data (n=24K)**

# GSP and TOPMed Collaboration on Annotation

- GSP Annotation WG (ACs, CCDGs, CMGs)

- TOPMed Annotation Interest WG

- Collaboration between GSP and TOPMed , with the goal of consistent annotations across both data sets.

- NHLBI RFI on Strategically Critical Resources or Infrastructures Using R24/U24

- Make annotation resources to the whole GSP community, as well as the general research community, .e.g., the NIH Data Commons Use Cases.

# Types of Raw Annotation Useful for Disease Mapping in WGS

- **MAFs**
- **Variant types**
  - VEP annotation using Genecode, e.g., LOF, missense, synonymous, 5", 3"
  - Promoters using FANTOM5
  - Enhancers using GeneHancer
- **Protein scores (coding variants only),** e.g., Polyphen and SIFT
- **Evolutionary/conservation scores**, e.g, Priphylop, GerpN, GerpS
- **Epigenetic scores** (ENCODE, ROADMAP)(cell specific), e.g. H3AK27ac,  H3K4me1
- **ChromHMM States**

# Types of Raw Annotation Useful for Disease Mapping in WGS: Continues

- **LD/heritability-related scores**, e.g., bStatistics, recombination rates
- **Higher-order chromatin interactions**, e.g., Hi-C
- **Transcpriptome scores (eQTLs)** (GTEx, n too small for rare variant eQTLs)
- **Methylation scores (mQTLs)**
- **Post-translational/protein (pQTL) scores** (UniProt)
- **Phenotype-specific heritability**
- **Structural Variant and Indel annotation**

# Composite Annotation Scores

- Existing (mainly driven by conservation scores for non-coding variants)
  - CADD/EIGEN (similar)
  - LINSIGHT (not defined for coding variants)
  - FATHMM-XF
  - GenoCanyon and GenoSkyline (cell/tissue specific)
- Newer scores
  - PINES (epigenetic-score based)
  - LACE (two dimension scores: epigenetics and conservation)
  - Annotation PCs

# Population Genetics annotation

- **Variant level -** additional population and evolutionary scores (i.e. Fst; recurrent mutations; allele age estimates)

- **Locus level -** local ancestry inference; local IBD inference; local recombination rate; haplotype selection scores (i.e. B-statistics)

- **Individual level -** self-reported ancestry, global genetic ancestry, inbreeding coefficient

- **Population level -** global genetic ancestry; IBD community; Founder effects (IBDNe, ROH), admixture;