

# UK Biobank data to fuel drug discovery

Jonathan Marchini

Executive Director, Head of Statistical Genetics and Machine Learning

**REGENERON**  
SCIENCE TO MEDICINE®

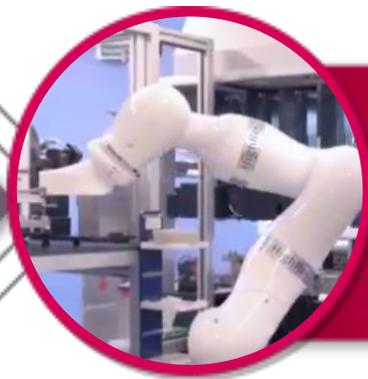
# Regeneron Genetics Center (RGC)

Established in 2014 and is now one of the largest operational human sequencing efforts

**SAMPLE  
BIOBANKING**



**LIBRARY PREPARATION  
AND EXOME CAPTURE**



**ILLUMINA-BASED  
SEQUENCING**



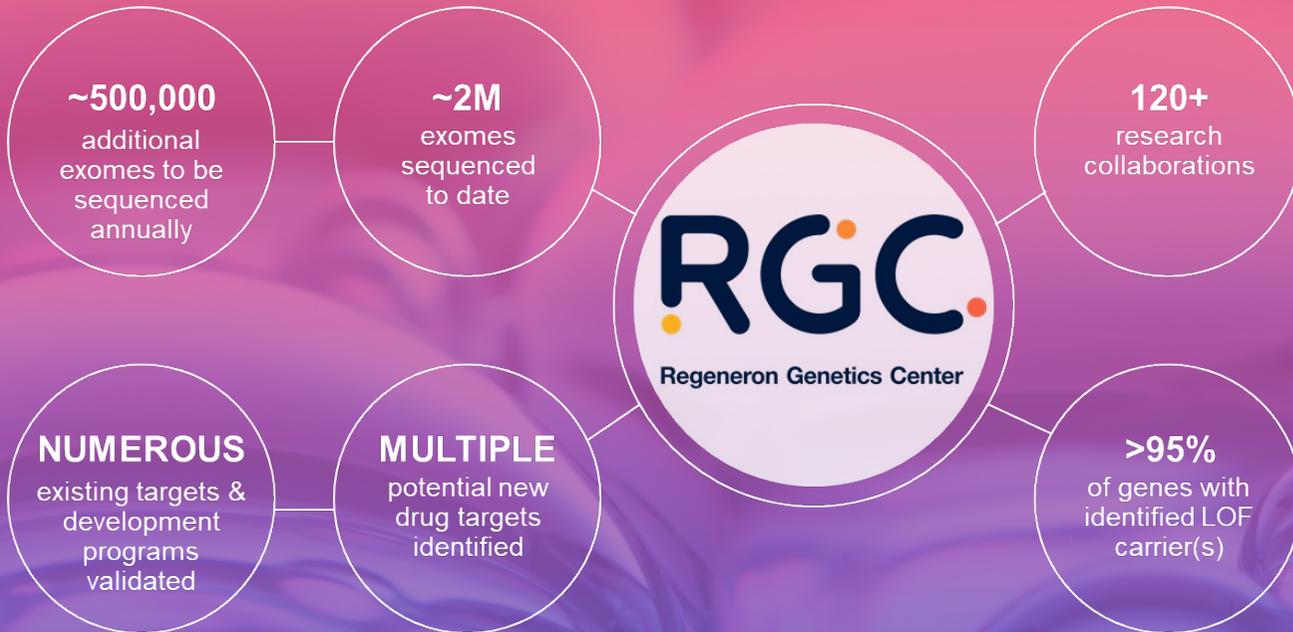
**CLOUD BASED INFORMATICS  
& ANALYSIS**



**Mission:**

Taking large scale human genetics to the next level for target discovery, support existing targets and identify novel indications

# Regeneron Genetics Center: Unprecedented Speed, Scale & Integration



All accomplished in just the first 9 years!

RGC has the most diverse collection and catalogue of human coding variation to date

The New York Times

## *Hospital and Drugmaker Move to Build Vast Database of New Yorkers' DNA*

Patients will be asked if their genetic sequence can be added to a database — shared with a pharmaceutical company — in a quest to cure a multitude of diseases.

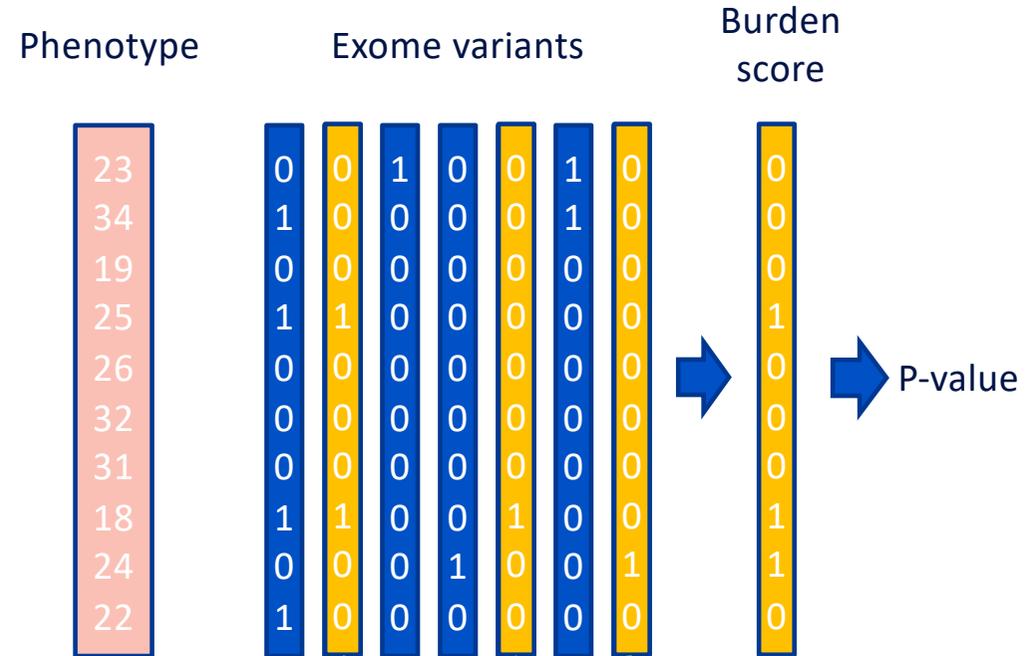
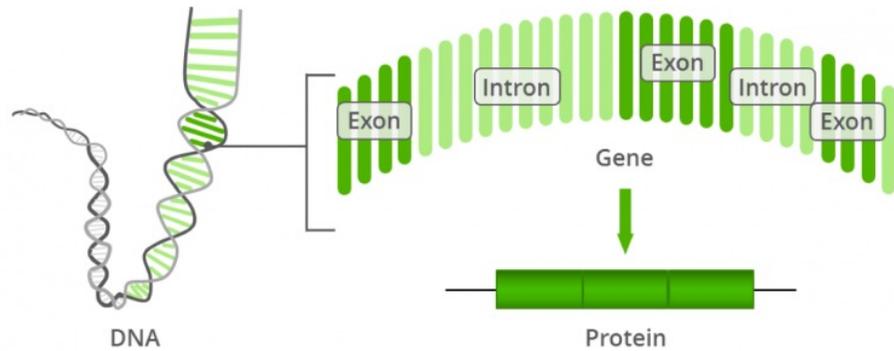
Give this article 97



Wilbert Gibson is a Mount Sinai patient who agreed to let the hospital system use his genetic information in research for treatment of a variety of diseases. Hiroko Masuike/The New York Times

# Exome sequencing of all 500,000 UK Biobank participants

The exons are the 1-2% of the genome that encode the proteins.



- Regeneron led consortium of commercial companies.
- All 500,000 sequences made available in mid-2022.



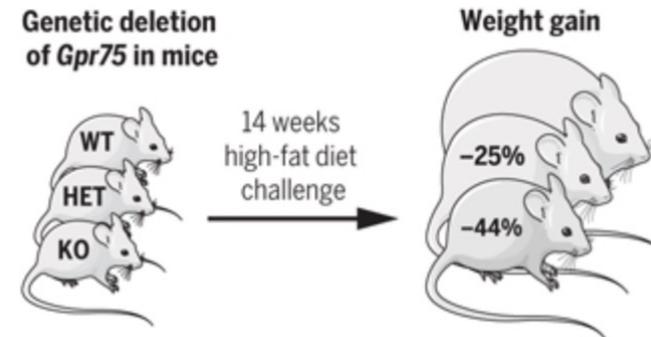
**Mask definitions**

- Annotations
- Allele frequencies

# Sequencing of 640,000 exomes identifies GPR75 variants associated with protection from obesity

PARSA ARBARI, ANKIT GILANI, OLUKAYODE SOSINA, JACK A. KOSMICKI, LORI KHRIMIAN, YI-YA FANG, TRIKALDARSHI PERSAUD, VICTOR GARCIA, DYLAN SUN, ALEXANDER LI, JOELLE MBATCHOU, ADAM E. LOCKE, CHRISTIAN BENNER, NIEK VERWELJ, NAN LIN, SAKIB HOSSAIN, KEVIN AGOSTINUCCI, JONATHAN V. PASCALE, ERCUMENT DIRICE, MICHAEL DUNN, REGENERON GENETICS CENTER<sup>†</sup>, DISCOVEREHR COLLABORATION<sup>‡</sup>, WILLIAM E. KRAUS, SVATHI H. SHAH, YI-DER I. CHEN, JEROME I. ROTTER, DANIEL J. RADER, OLLE MELANDER, CHRISTOPHER D. STILL, TOORAJ MIRSHAHI, DAVID J. CAREY, JAIME BERUMEN-CAMPOS, PABLO KURI-MORALES, JESUS ALEGRE-DIAZ, JASON M. TORRES, JONATHAN R. EMBERSON, BORY COLLINS, SUGANTHI BALASUBRAMANIAN, ALICIA HAWES, MARCUS JONES, BRIAN ZAMBROWICZ, ANDREW J. MURPHY, CHARLES PAULDING, GIOVANNI COPPOLA, JOHN D. OVERTON, JEFFREY G. REID, ALAN R. SHULDINER, MICHAEL CANTOR, HYUN M. KANG, GONCALO R. ABEGASIS, KATIA KARALIS, ARIS N. ECONOMIDES, JONATHAN MARCHINI, GEORGE D. YANCOPOULOS, MARK W. SLEEMAN, JUDITH ALTAREJOS, GIUSY DELLA GATTA, ROBERTO TAPIA-CONYER, MICHAL L. SCHWARTZMAN, ARIS BARAS, MANUEL A. R. FERREIRA, AND LUCA A. LOTTA

- Rare predicted loss of function coding variants in *GPR75* for heterozygous carriers found to be associated with
  - Lower BMI (-1.8 kg/m<sup>2</sup>)
  - Lower body weight (~5.3 kg or 11.6 lbs lower)
  - Protection against obesity (54% lower odds)



B

Study	Beta (95% CI) in SDs of BMI	Genotypes (RR RA AA)	p-value	Beta (95% CI) in kg/m <sup>2</sup> of BMI	Beta (95% CI) in kg of body weight
<b>Discovery Cohorts</b>					
MCPS	-0.48 (-0.82, -0.13)	95,816 30 0	7.1E-03	-2.6 (-4.4, -0.7)	-7.4 (-12.8, -2)
GHS_EUR	-0.27 (-0.52, -0.02)	121,010 51 0	3.6E-02	-1.4 (-2.8, -0.1)	-4.2 (-8, -0.3)
UKB_EUR	-0.34 (-0.49, -0.19)	428,572 147 0	6.6E-06	-1.8 (-2.6, -1)	-5.3 (-7.6, -3)
<b>Subgroup meta-analysis results</b>	-0.34 (-0.46, -0.22)	645,398 228 0	2.6E-08	-1.8 (-2.5, -1.2)	-5.3 (-7.1, -3.4)

# Diverse RGC cohorts: Novel genetic discoveries and therapeutic targets

nature communications



Article

<https://doi.org/10.1038/s41467-022-32398-7>

## Multiancestry exome sequencing reveals *INHBE* mutations associated with favorable fat distribution and protection from diabetes

Received: 8 February 2022

Accepted: 28 July 2022

Published online: 23 August 2022

Check for updates

Parsa Akbari<sup>1,13</sup>, Olukayode A. Sosina<sup>1,13</sup>, Jonas Bovijn<sup>1,13</sup>, Karl Landheer<sup>2</sup>, Jonas B. Nielsen<sup>1</sup>, Minhee Kim<sup>1</sup>, Senem Aykul<sup>1</sup>, Tanima De<sup>1</sup>, Mary E. Haas<sup>1</sup>, George Hindy<sup>1</sup>, Nan Lin<sup>1</sup>, Ian R. Dinsmore<sup>3</sup>, Jonathan Z. Luo<sup>3</sup>, Stefanie Hectors<sup>2</sup>, Benjamin Geraghty<sup>1</sup>, Mary Germino<sup>2</sup>, Lampros Panagis<sup>2</sup>, Prodromos Parasoglou<sup>2</sup>, Johnathon R. Walls<sup>2</sup>, Gabor Halasz<sup>2</sup>, Gurinder S. Atwal<sup>2</sup>, Regeneron Genetics Center<sup>4</sup>, DiscovEHR Collaboration<sup>4</sup>, Marcus Jones<sup>1</sup>, Michelle G. LeBlanc<sup>1</sup>, Christopher D. Still<sup>4</sup>, David J. Carey<sup>4</sup>, Alice Giontella<sup>5,6</sup>, Marju Orho-Melander<sup>5</sup>, Jaime Berumen<sup>7</sup>, Pablo Kuri-Morales<sup>7,8</sup>, Jesus Alegre-Díaz<sup>7</sup>, Jason M. Torres<sup>9,10</sup>, Jonathan R. Emberson<sup>9,10</sup>, Rory Collins<sup>10</sup>, Daniel J. Rader<sup>11</sup>, Brian Zambrowicz<sup>2</sup>, Andrew J. Murphy<sup>2</sup>, Suganthi Balasubramanian<sup>1</sup>, John D. Overton<sup>1</sup>, Jeffrey G. Reid<sup>1</sup>, Alan R. Shuldiner<sup>1</sup>, Michael Cantor<sup>1</sup>, Goncalo R. Abecasis<sup>1</sup>, Manuel A. R. Ferreira<sup>1</sup>, Mark W. Steeman<sup>2</sup>, Viktoria Gusarova<sup>2</sup>, Judith Altarejos<sup>2</sup>, Charles Harris<sup>2</sup>, Aris N. Economides<sup>1,2</sup>, Vincent Idone<sup>2</sup>, Katia Karalis<sup>1</sup>, Giusy Della Gatta<sup>1</sup>, Tooraj Mirshahi<sup>4</sup>, George D. Yancopoulos<sup>2</sup>, Olle Melander<sup>5,12</sup>, Jonathan Marchini<sup>1</sup>, Roberto Tapia-Conyer<sup>8,13</sup>, Adam E. Locke<sup>1,13</sup>, Aris Baras<sup>1,13</sup>, Niek Verweij<sup>1,13</sup> & Luca A. Lotta<sup>1,13</sup> ✉

Akbari et al. (2022) Nat Commun 13, 4844

Association with favorable fat distribution ( $p = 1.8 \times 10^{-09}$ ), favorable metabolic profile and protection from type 2 diabetes (~28% lower odds;  $p = 0.004$ ) for heterozygous protein-truncating mutations in *INHBE*

The NEW ENGLAND JOURNAL of MEDICINE

ORIGINAL ARTICLE

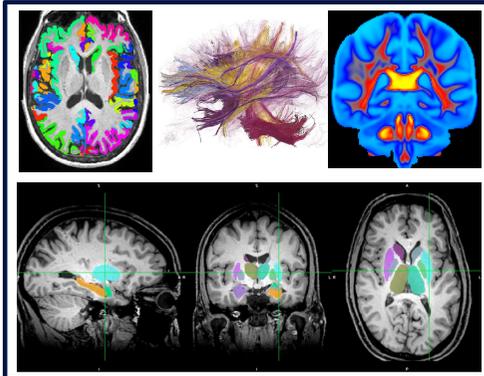
## Germline Mutations in *CIDEB* and Protection against Liver Disease

N. Verweij, M.E. Haas, J.B. Nielsen, O.A. Sosina, M. Kim, P. Akbari, T. De, G. Hindy, J. Bovijn, T. Persaud, L. Miloscio, M. Germino, L. Panagis, K. Watanabe, J. Mbatchou, M. Jones, M. LeBlanc, S. Balasubramanian, C. Lammert, S. Enhörning, O. Melander, D.J. Carey, C.D. Still, T. Mirshahi, D.J. Rader, P. Parasoglou, J.R. Walls, J.D. Overton, J.G. Reid, A. Economides, M.N. Cantor, B. Zambrowicz, A.J. Murphy, G.R. Abecasis, M.A.R. Ferreira, E. Smagris, V. Gusarova, M. Sleeman, G.D. Yancopoulos, J. Marchini, H.M. Kang, K. Karalis, A.R. Shuldiner, G. Della Gatta, A.E. Locke, A. Baras, and L.A. Lotta  
Verweij et al. (2022) N Engl J Med 387:332-344

Rare predicted loss-of-function variants plus missense variants in *CIDEB* associated with 33% lower odds of liver disease of any cause

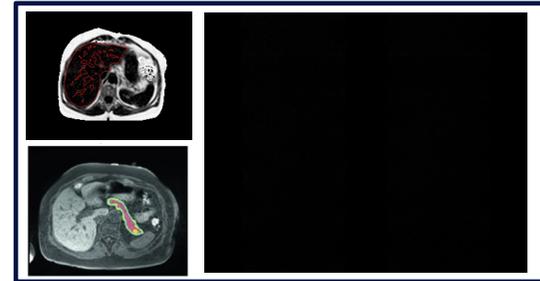
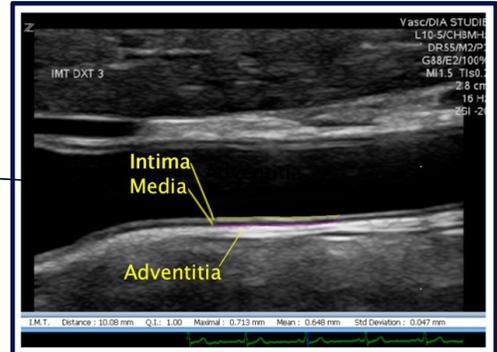
# The UK Biobank Imaging Study

- Aims to collect brain, heart, and abdomen scans from 100,000 participants.
- Repeat set of imaging on 60,000 participants.
- Raw images available, but derived phenotypes appear more slowly



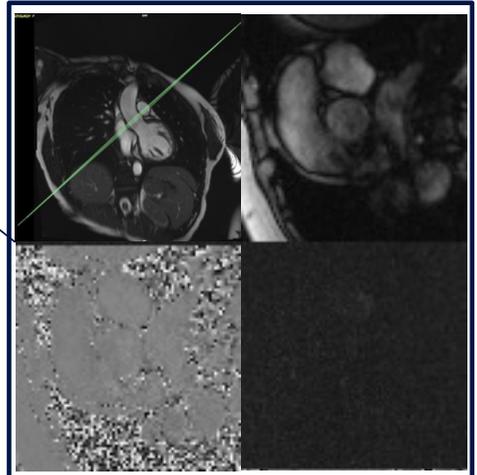
 Detailed MRI measures of the structure of the brain, including the connections between different parts of the brain.

 Ultrasound of the two arteries that take blood to the brain.



 Targeting the way fat is distributed in the abdomen, including around the liver and pancreas.

 Detailed assessment of the heart, including thickness of the heart wall, and how the heart changes as it pumps blood around the body.



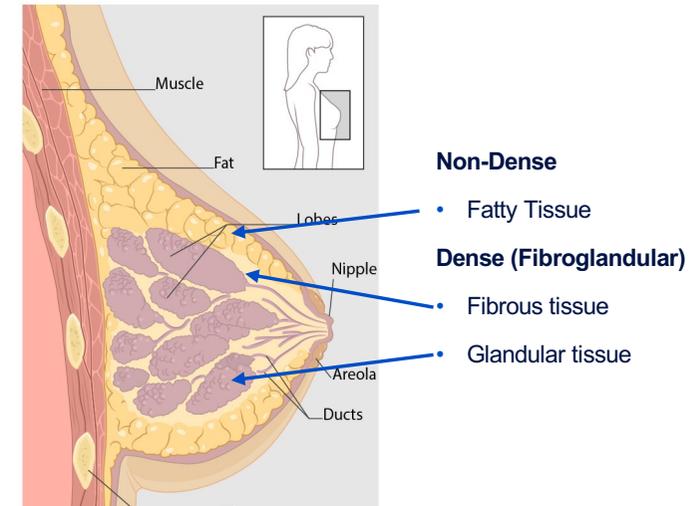
 Bone measurements, including fractures, with a focus on the spine, hips and knees.



# Can UKB whole body MRI be used to assess breast density?

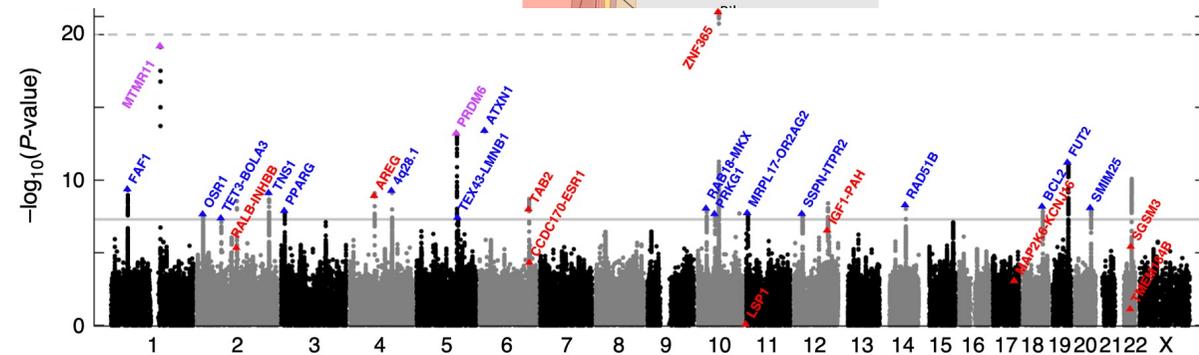
## Breast density is a risk factor for breast cancer, independent of other known risk factors <sup>1</sup>

- Women with density > 75% have an increased risk (range across 4 studies: 2.82-5.99) compared to women with < 10% density <sup>2</sup>
- Approximately 1/3 of breast cancer risk may be attributable to density (2 studies, attributable risk percent = 28% and 33%) <sup>2</sup>
- Breast density is modifiable—tamoxifen-induced reduction of density decreases the risk of subsequent breast cancer <sup>3</sup>



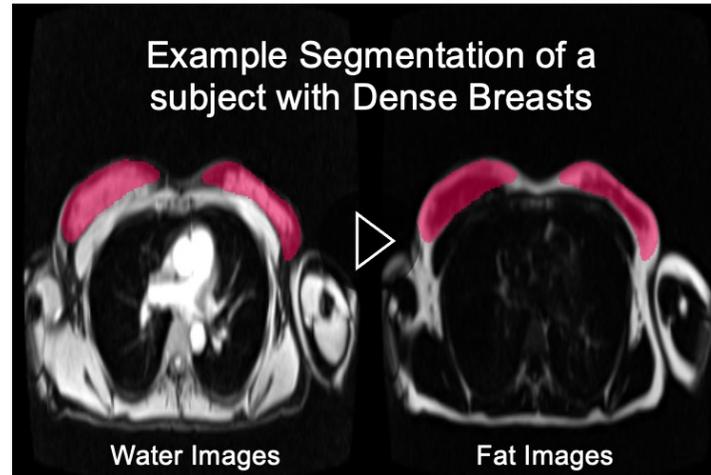
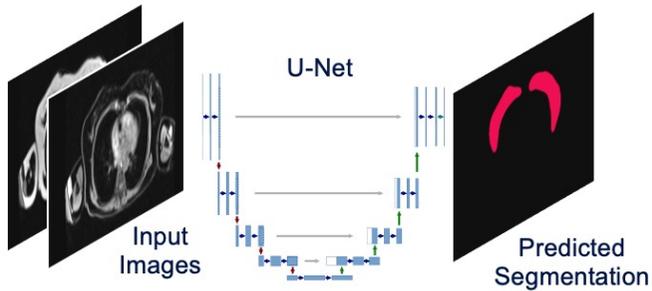
## Breast density has a strong genetic component

2 GWAS (BCAC and Kaiser) have identified 46 independent genome-wide significant breast density loci for three measures of breast density: dense area, non-dense area, and percent dense area, from mammograms. No ExWAS has been conducted to date



Sieh, W, et al. *Nat Comms.* 2020 11(1) :1-11 Mammography GWAS (N=24,192)

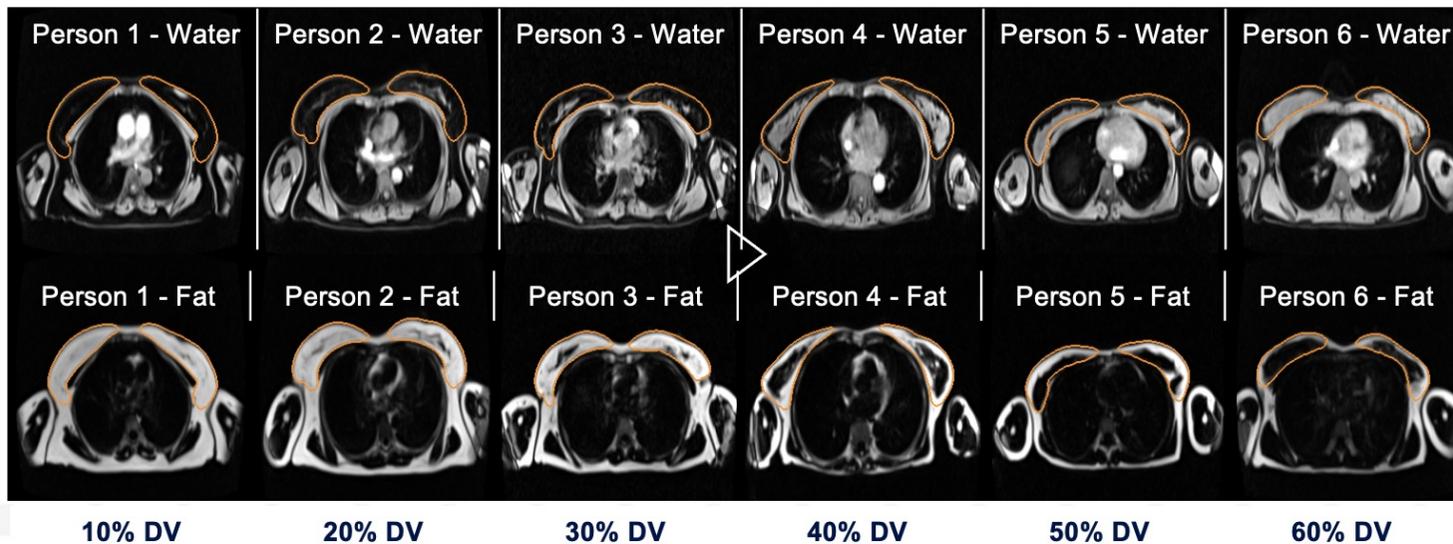
**Step 1 : Segment breasts** Ground truth breast segmentations from 122 subjects (~5,000 2D images) were used to train a U-Net model



**Step 2 : Measure density within breasts** Water and fat-fraction maps were generated and used to derive:

$$\text{Dense Volume (DV)} = \frac{\text{Water}}{\text{Water} + \text{Fat}}$$

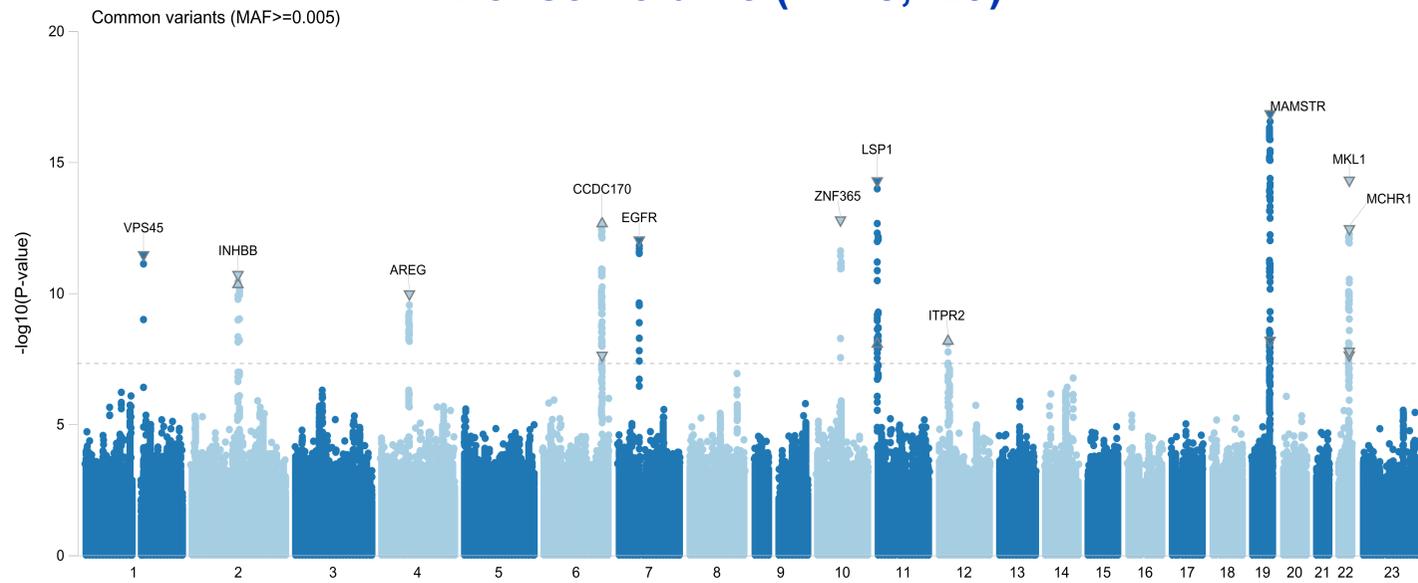
$$\text{Non-Dense Volume (NDV)} = \frac{\text{Fat}}{\text{Water} + \text{Fat}}$$



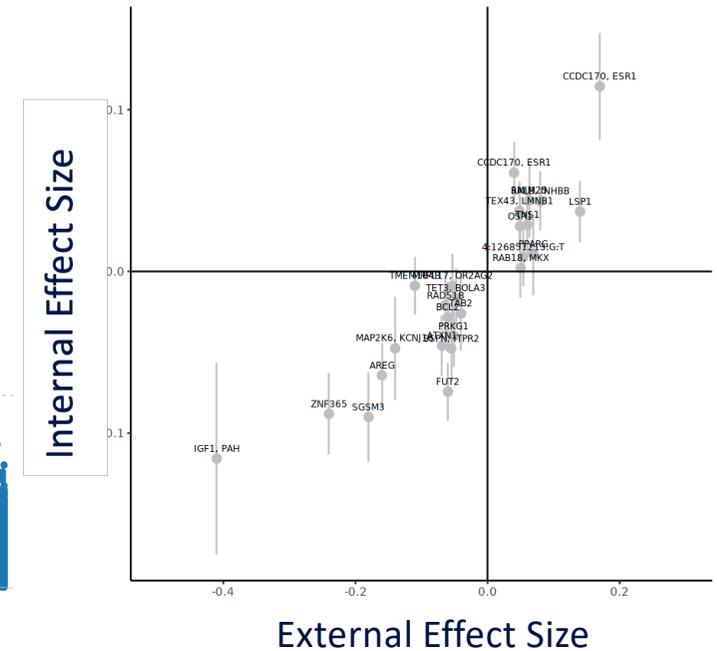
# GWAS with MRI-derived Breast Density Phenotypes

Genetic associations with MRI-derived dense volume recapitulate associations with known breast density loci

## Dense Volume (N=18,129)



Internal (y-axis) and external (x-axis) effect estimates are directionally consistent at 26 variants across previously published loci



# Data at this scale requires highly efficient and flexible analysis

Technical Report | [Published: 20 May 2021](#)

## Computationally efficient whole-genome regression for quantitative and binary traits

[Joelle Mbatchou](#), [Leland Barnard](#), [Joshua Backman](#), [Anthony Marcketta](#), [Jack A. Kosmicki](#), [Andrey Ziyatdinov](#), [Christian Benner](#), [Colm O'Dushlaine](#), [Mathew Barber](#), [Boris Boutkov](#), [Lukas Habegger](#), [Manuel Ferreira](#), [Aris Baras](#), [Jeffrey Reid](#), [Goncalo Abecasis](#), [Evan Maxwell](#) & [Jonathan Marchini](#) 

*Nature Genetics* **53**, 1097–1103 (2021) | [Cite this article](#)

<https://rgcgithub.github.io/regenie/>



- [regenie](#)
- [Citation](#)
- [License](#)
- [Contact](#)

## regenie

**regenie** is a C++ program for whole genome regression modelling of large [genome-wide association studies](#).

It is developed and supported by a team of scientists at the Regeneron Genetics Center.

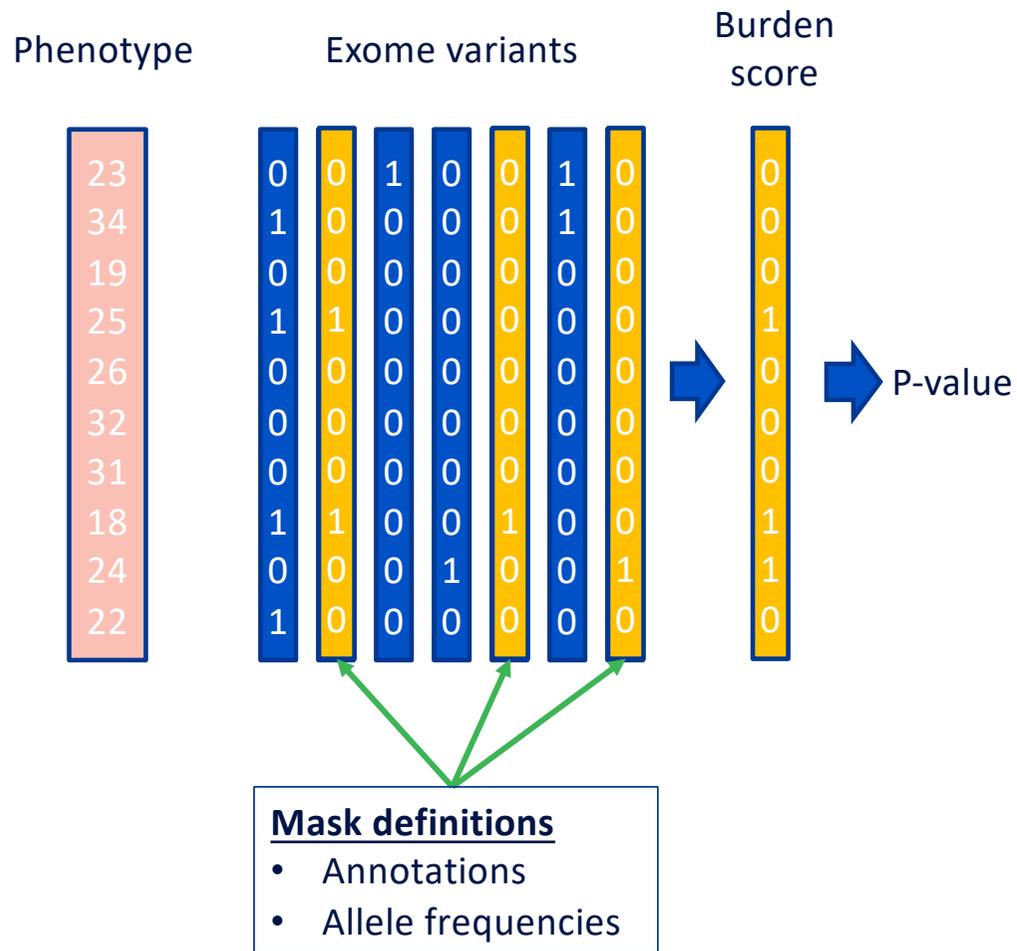
The method has the following properties

- It works on quantitative and binary traits, including binary traits with unbalanced case-control ratios
- It can handle population structure and relatedness
- It can process multiple phenotypes at once efficiently
- For binary traits, it supports Firth logistic regression and an SPA test
- It can perform gene/region-based tests (Burden, SKAT/SKATO, ACATV/ACATO)
- It can perform interaction tests (GxE, GxG) as well as conditional analyses
- It is fast and memory efficient 🔥
- It supports the [BGEN](#), [PLINK](#) bed/bim/fam and [PLINK2](#) pgen/pvar/psam genetic data formats
- It is ideally suited for implementation in [Apache Spark](#) (see [GLOW](#))
- It can be installed with [Conda](#)



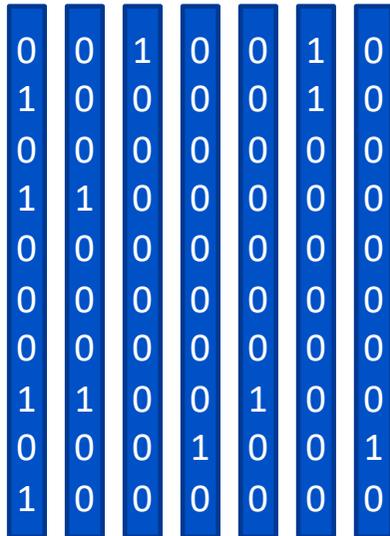
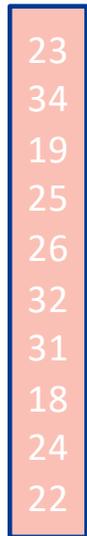
# Efficient meta-analysis of gene-based tests (REMETA)





Phenotype

Exome variants



### Advantages

- Works for Burden and SKAT tests [1]
- Easy to change mask definition
- Amenable to meta-analysis [2,3,4] i.e. metaSKAT and metaSTARR

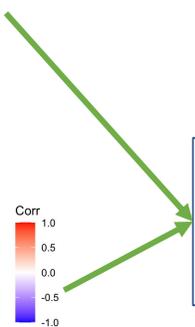
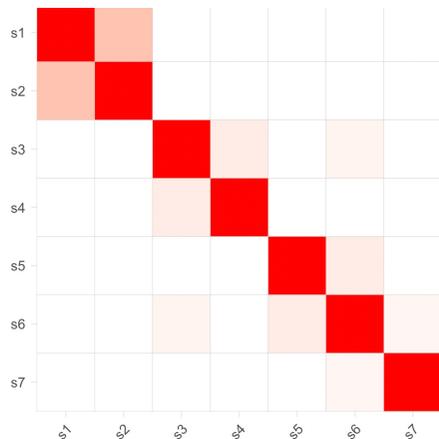
### Disadvantages

Calculating and keeping track of LD matrices per trait

Score stats

S<sub>1</sub> S<sub>2</sub> S<sub>3</sub> S<sub>4</sub> S<sub>5</sub> S<sub>6</sub> S<sub>7</sub>

SNP LD matrix



**Mask definitions**

- Annotations
- allele frequencies

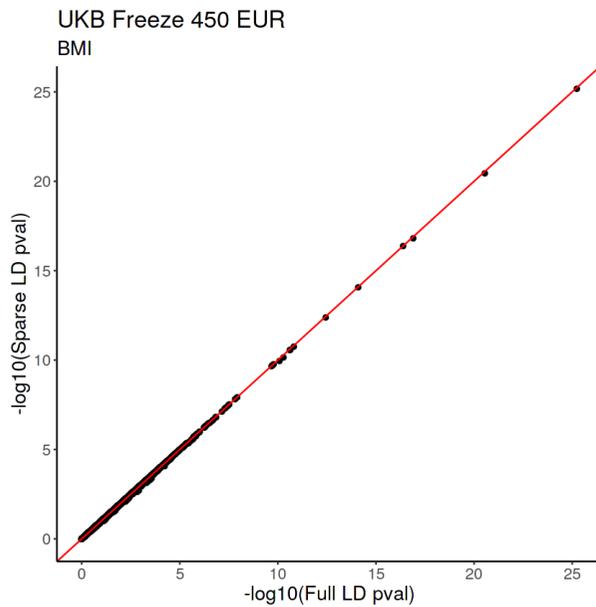


**P-value [1]**

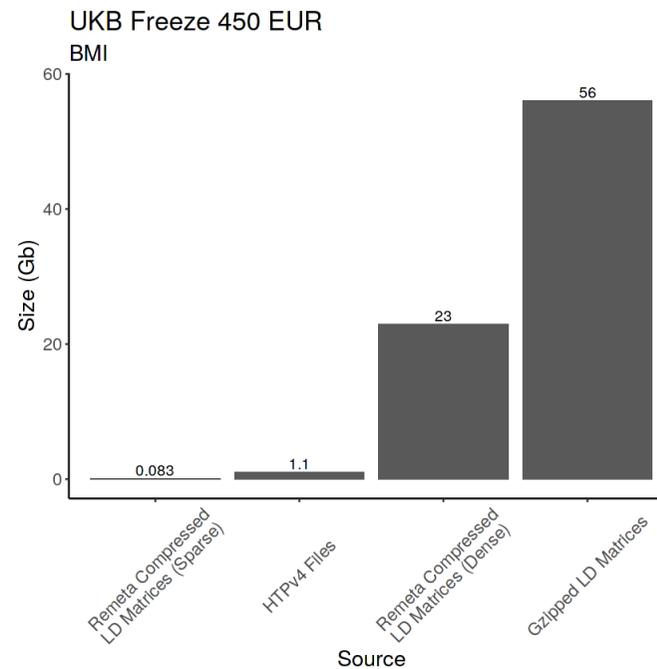
[1] Lee et al. (2012) AJHG 91(2): 224–237  
[2] Lee et al. (2013) AJHG 93(1):42–53  
[3] Liu et al. (2014) NatGen 46(2):200–4  
[4] Li et al. (2023) NatGen 55(1): 154–164.

# Gene-based meta-analysis on an industrial scale

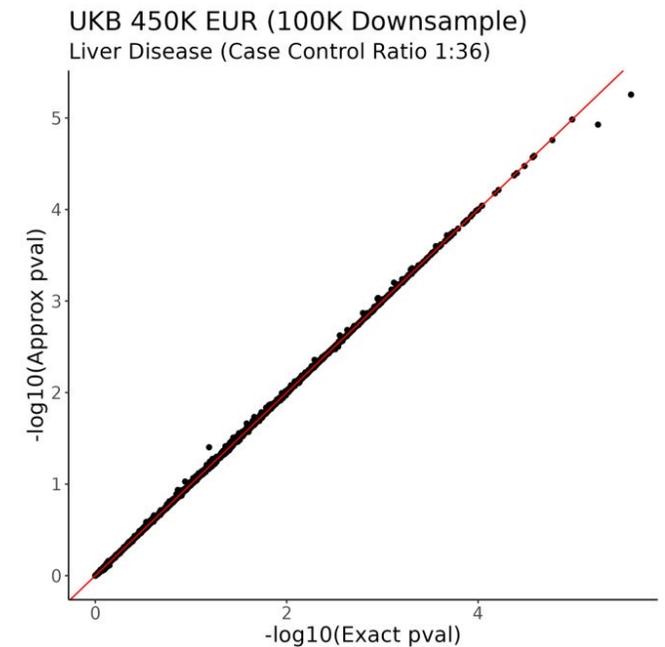
**Challenge** scale/ease-of-use when there are many cohorts, phenotypes/sub-phenotypes...



Sparse LD has no effect on gene p-values

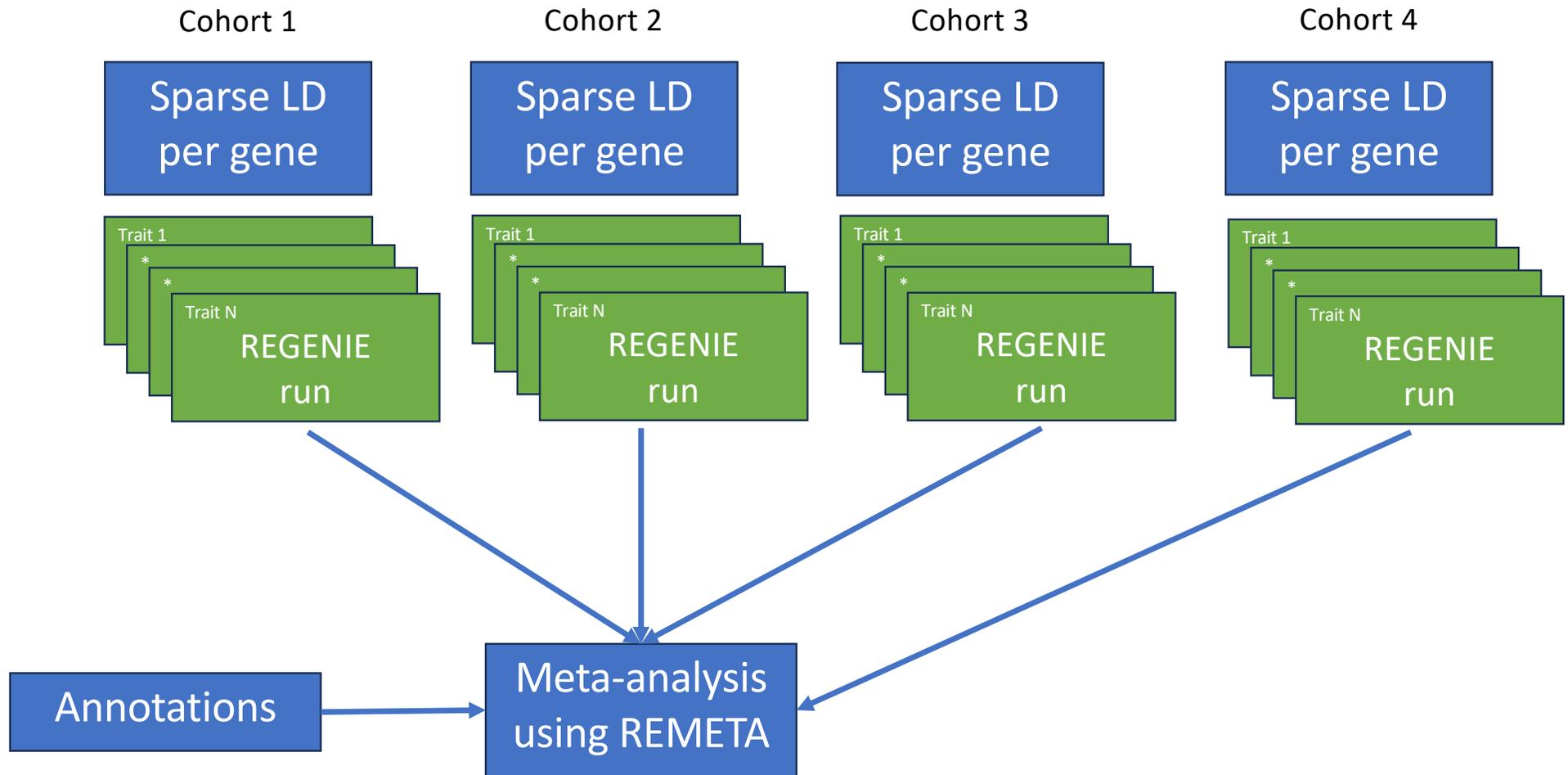


Sparse LD can be stored compactly



Using a reference LD file per cohort is accurate

# Meta-analysis workflow



# RGC presentations at ASHG

RGC Author	Assignment	Title	Presentation Date	Presentation Time	Location <sup>***</sup>
Kuan-Han Wu	Platform	Rare variants associated with prostate cancer risk discovered from 255,640 male exomes influence risk of prostate cancer metastasis	02 November 2023	9:15 a.m. – 9:30 a.m.	Conv Ctr/Room 146B/Level 1
Joelle Mbatchou	Poster	Using protein language model annotations to improve the power of exome-wide association studies (PB4424)	02 November 2023	3:00 p.m. – 5:00 p.m.	Exhibit & Poster Hall, Halls A/B
Vijay Kumar	Poster	Population-scale analysis of the trinucleotide repeat expansion in the Huntingtin gene (HTT) from 854,251 human exomes (PB1723)	02 November 2023	3:00 p.m. – 5:00 p.m.	Exhibit & Poster Hall, Halls A/B
Yuxin Zou	Poster	Joint fine-mapping of single variants and gene-based tests from exome sequencing and genotype imputation (PB4238)	02 November 2023	3:00 p.m. – 5:00 p.m.	Exhibit & Poster Hall, Halls A/B
Veera Rajagopal	Platform	Discovering genes linked to both cognition and psychiatric disorders through analysis of 888,052 exomes	03 November 2023	2:15 p.m. – 2:30 p.m.	Conv Ctr/Ballroom B/Level 3
Rujin Wang	Platform	A large-scale meta-analysis of genome-wide association studies reveals genetics underlying Parkinson's Disease leveraging electronic health records	03 November 2023	2:30 p.m. – 2:45 p.m.	Conv Ctr/Room 207A/Level 2
Arden Moscati	Poster	Genetic ancestry-based case-control matching to improve power to trait-specific association analysis (PB4158)	03 November 2023	3:00 p.m. – 5:00 p.m.	Exhibit & Poster Hall, Halls A/B
Ariane Ayer	Poster	Phenome-wide genetic associations of educational attainment with mental and behavioral disorders (PB1706)	03 November 2023	3:00 p.m. – 5:00 p.m.	Exhibit & Poster Hall, Halls A/B
Liron Ganel	Poster	Rare variant analysis of MRI-derived fat distribution phenotypes strengthens detected effects compared to larger meta-analysis (PB1754)	03 November 2023	3:00 p.m. – 5:00 p.m.	Exhibit & Poster Hall, Halls A/B
Sahar Gelfman	Poster	A large meta-analysis identifies genes associated with Anterior Uveitis (PB1313)	03 November 2023	3:00 p.m. – 5:00 p.m.	Exhibit & Poster Hall, Halls A/B
Tyler Joseph	Poster	REMETA: Efficient meta-analysis of gene-based tests in large-scale genetic studies (PB4350)	03 November 2023	3:00 p.m. – 5:00 p.m.	Exhibit & Poster Hall, Halls A/B
Jack Kosmicki	Platform	Exome sequencing of >1 million individuals identifies 209 genes associated with adult human height	04 November 2023	10:45 a.m. – 11:00 p.m.	Conv Ctr/Room 202A/Level 2
Blair Zhang	Poster	Genetic risk score in age-related macular degeneration subtypes across electronic health record cohorts (PB4168)	04 November 2023	2:15 p.m. – 4:15 p.m.	Exhibit & Poster Hall, Halls A/B
Manav Kapoor	Poster	Genome-wide exploration of positively selected loci and their association to disease phenotypes in 30,000 individuals from Sri Lanka and Bangladesh (PB3072)	04 November 2023	2:15 p.m. – 4:15 p.m.	Exhibit & Poster Hall, Halls A/B
Sophia Praggastis	Poster	A genome-wide meta-analysis connects iron homeostasis to metabolic disease through poly-unsaturated fatty acid synthesis (PB1308)	04 November 2023	2:15 p.m. – 4:15 p.m.	Exhibit & Poster Hall, Halls A/B
Kathy Burch	Platform	Leveraging ~937K exomes to estimate cancer risk conferred by rare deleterious germline variants in hereditary cancer risk genes	05 November 2023	10:00 a.m. – 10:15 a.m.	Conv Ctr/Room 202A/Level 2

\*\*\* Location: Walter E. Washington Convention Center