

Whole Genome Sequencing release and multi-omics in UK Biobank

ASHG

November 2nd 2023

Lucy Burkitt-Gray, PhD
Lead Data Analyst
UK Biobank



UK Biobank is a large-scale **biomedical database** and research **resource**, containing in-depth genetic and health information from **half a million** UK participants. The database is **regularly augmented** with additional data and is **globally accessible** to approved researchers undertaking vital research into the most **common and life-threatening diseases**. It is a major contributor to the **advancement of modern medicine** and treatment and has enabled several scientific discoveries that **improve human health**.

Non-profit charity,
established by:



Ongoing core funding and additional
funding from:



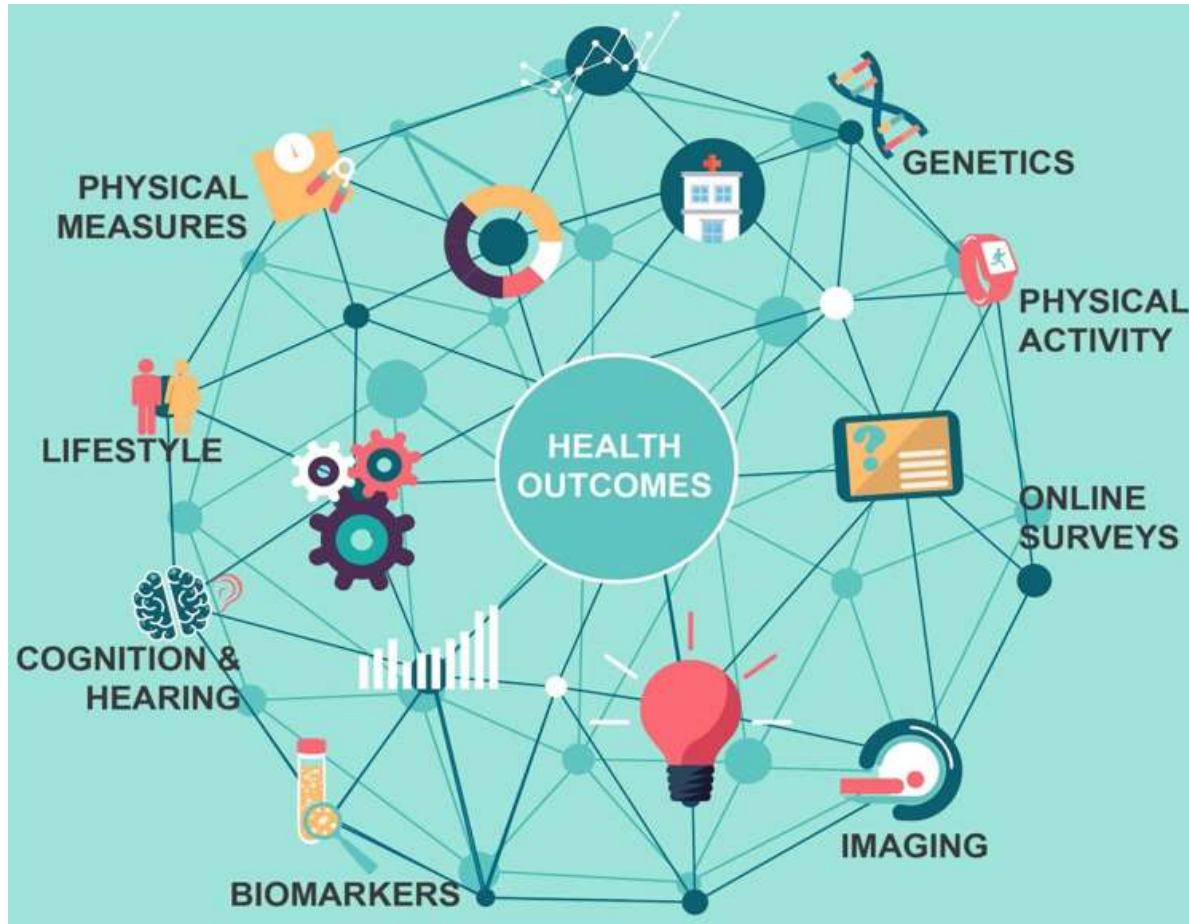
Link to 'The ground breaking UK Biobank Resource' video:
<https://www.youtube.com/watch?v=NGdegXRx8U0&t=152s>

Link to 'What is UK Biobank?' video:
<https://www.youtube.com/watch?v=66mol1ZHMYs>

Link to 'Celebrating 20 years of UK Biobank':
<https://www.ukbiobank.ac.uk/learn-more-about-uk-biobank/our-impact>



- 500,000 people aged 40 – 69 in 2006-10 from England (89%), Scotland (7%), and Wales (4%)
- 22 assessment centres located around the UK to enhance heterogeneity of the cohort
- Breadth and depth of data
 - Lifestyle and environmental exposures
 - Personal and family medical history
 - Cognitive function, hearing and vision tests
 - Physical measures (BP, lung function, body size)
 - Biological samples (blood, urine, saliva)
- Consent to access health-related records and to re-contact participants for further assessments



Regular updates on participant lifestyle, health and risk factors through focussed questionnaires

1-2 questionnaires per year on specific topic
Nutrition, sleep, mental health, pain

Linkage to NHS Digital records provide detailed health information on hospital admissions, GP records, prescriptions, and more

Linkage to death certificates for participant mortality and causes of death

At home physical activity monitoring, cardiac monitoring, COVID antibody testing

Genotyping of all 500,000 participants



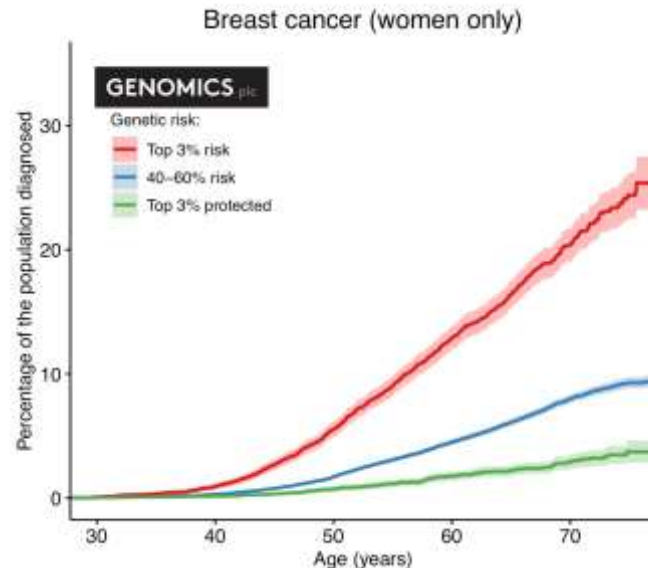
Custom-built genotyping array (850k variants) with imputed measures for 90M+ variants

Data made available for all 500,000 participants in 2017

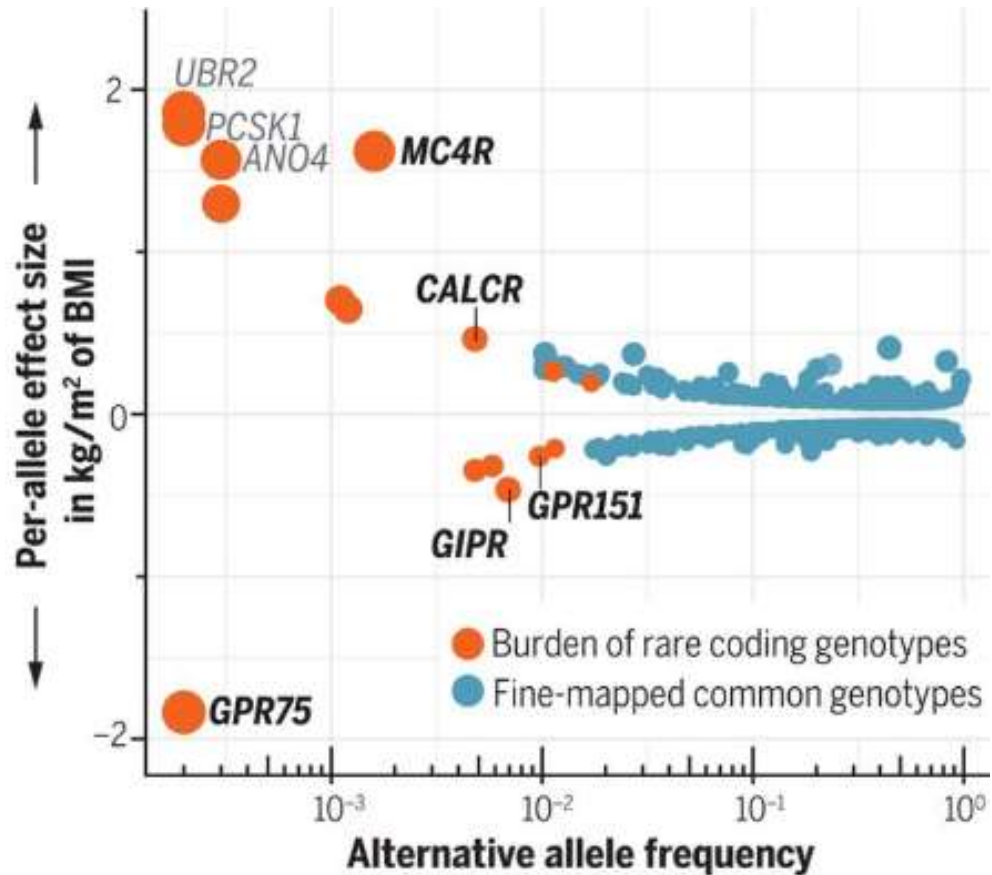
GWAS of hundreds of traits now publicly available

Polygenic risk scores generated for wide range of conditions

Made possible due to UK Biobank's large size, standardised assessment and health outcome follow-up through linkage to electronic records



Thompson *et al.*, MedRxiv 2022

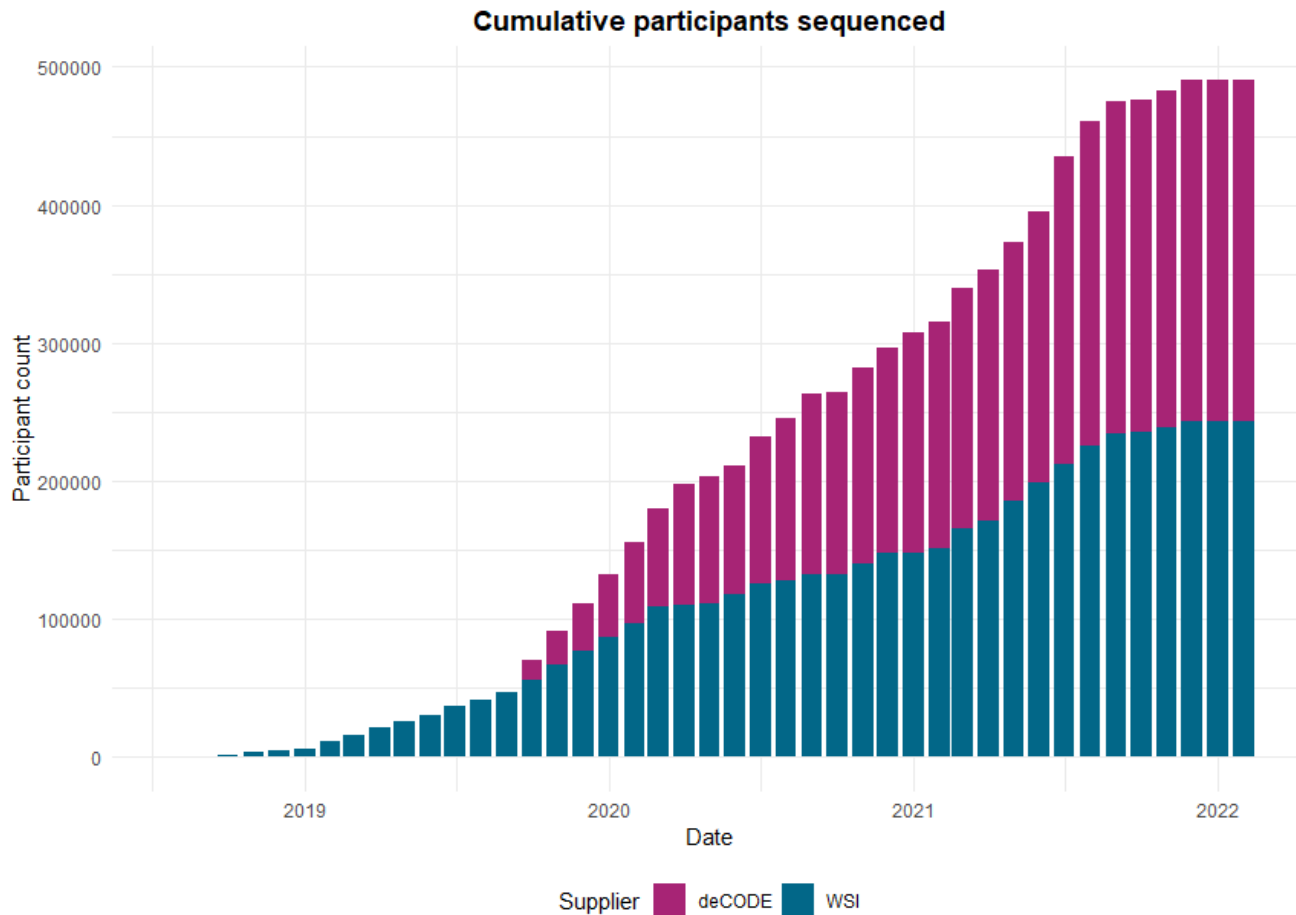


Power of massive-scale exome sequencing to enable discovery-based gene-burden analysis (e.g. GPR75 as obesity therapeutic target)

Akbari *et al.*, Science 2021

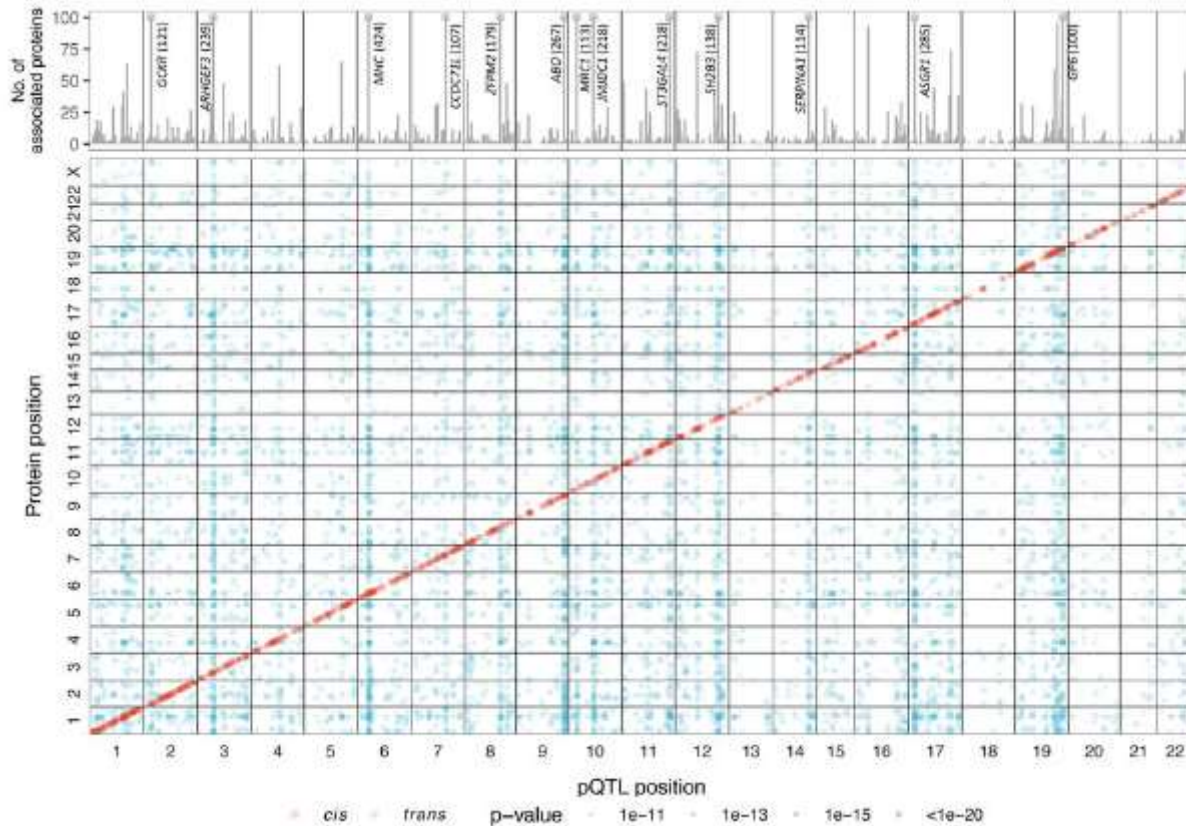
- Sequencing of all coding regions of the genome (~2%)
- Regeneron-led commercial funding
- Sequenced with Illumina NovaSeq 6000 platform using S2/S4 flow cells
- First 50,000 sequences released in 2019; final release of whole cohort in mid-2022





- 30X sequencing across the entire genome
- Government, charity, and industry funding
- Sequencing carried out between Wellcome Sanger Institute and deCODE Genetics using Illumina NovaSeq
- First 200,000 sequences publically available late 2021; full cohort to be released end 2023





Targeted proteomics on ~60,000 participants began in late 2020 supported by large pharmaceutical consortium

Using Olink platform to measure 3,000 circulating proteins from plasma samples

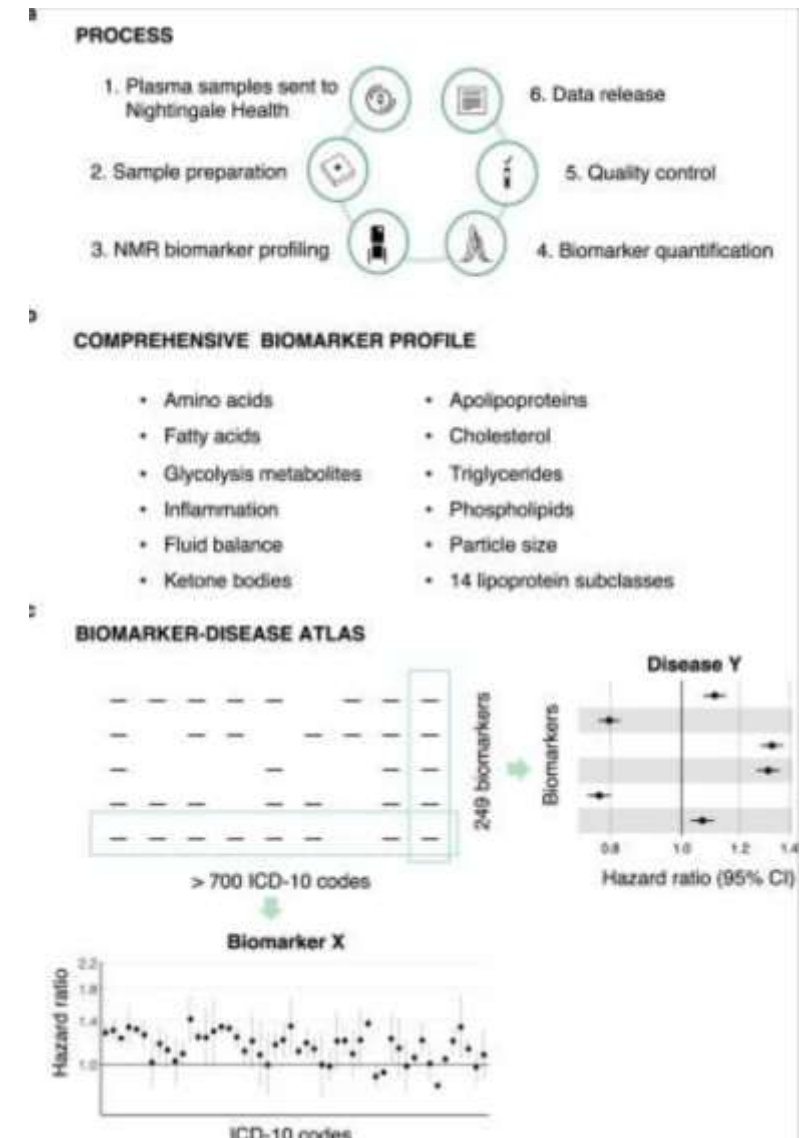
Already world's largest pQTL collection (10,000+)

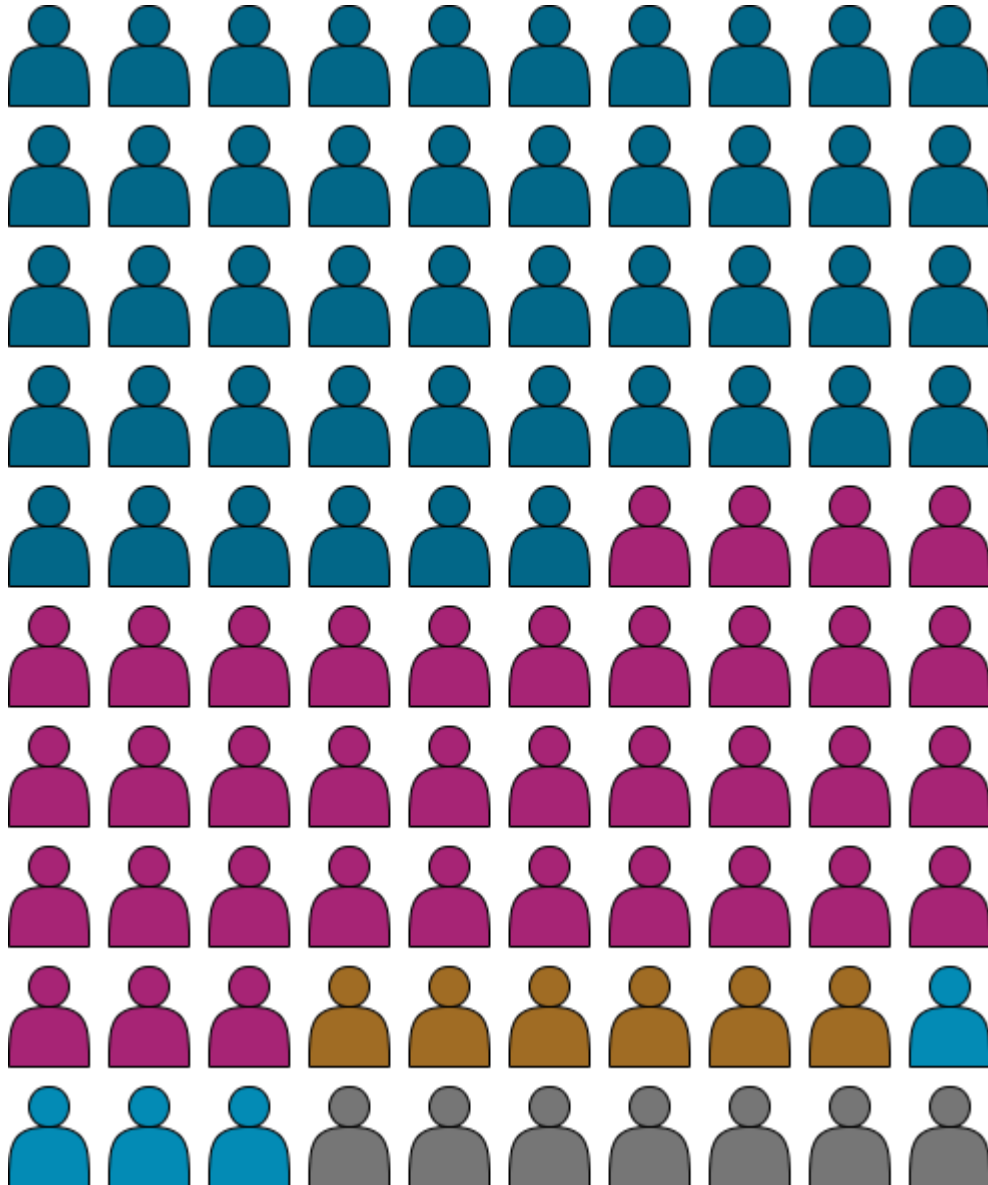
Data on 3,000 proteins for 56,000 participants now available

Approved mass spectrometry-based study commencing next year, for full-cohort analysis by Eliptica



- Metabolomic study carried out with Nightingale Health analysing 249 biomarkers in plasma
- Data for 275,000 participants made available in early 2023
- 16,000 repeat samples analysed, providing two timepoints 2-7 years apart
- Planned release of full cohort metabolomics late 2024/early 2025





As of the end of 2023:

- 46% of the cohort have whole genome, whole exome, genotyping, and metabolomic data
- A further 37% have whole exome, whole genome, and genotyping
- 6% of the cohort have data in all major omic datasets
- 4% have whole genome, whole exome, genotyping and proteomic data
- Overlaps will increase with future proteomic and metabolomic releases in 2024

Acknowledgements

Core funders



Medical
Research
Council



Collaborators

Lora Boteva

Oliver Gray

Daisy Vinter

Rachael Winkless

Dan Fry

Caroline Clark

Our 500k participants



Apply for access:

ukbiobank.ac.uk/enable-your-research/apply-for-access