

UK Biobank: democratising access to large-scale genomic and phenotypic data for discovery science

Dr Ben Lacey

UK Biobank Epidemiology Group

University of Oxford

N8 CIR: Data Access for
Digital Health

8 Nov 2023



What is UK Biobank?

- Prospective cohort study, following the health of 500,000 people to understand the genetic, physiological, lifestyle, and environmental determinants of disease
- UK Biobank has a **unique** combination of scale, depth, and duration of follow-up
- Research database to enable scientific discoveries, which is **readily accessible** to researchers worldwide

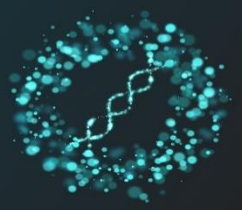




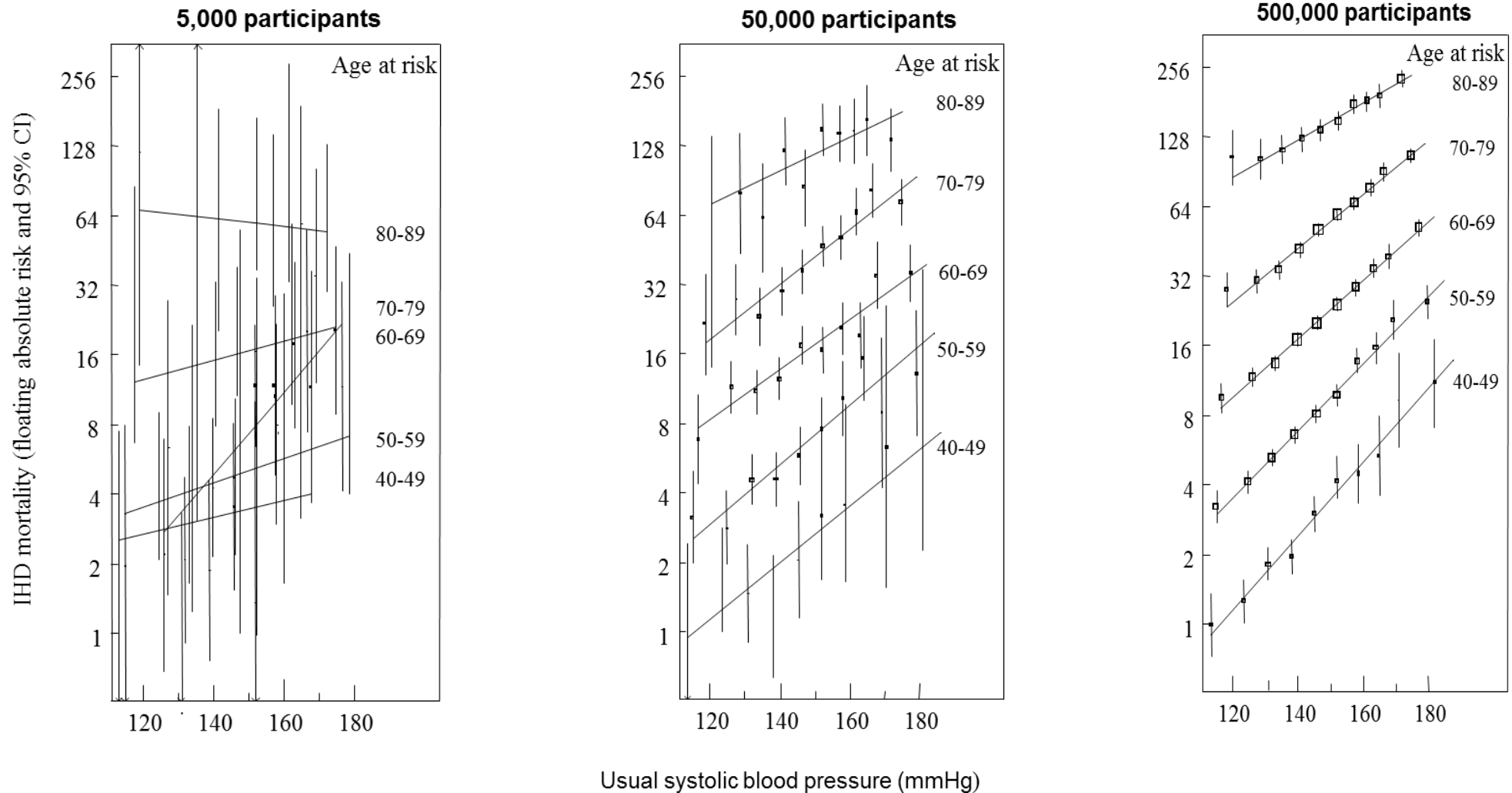
Why was UK Biobank established?

- The causes of most disease are due to a combination of many factors...genetics, physiological factors, lifestyle, environment (**detail** – we need a lot of data on participants)
- Moderate or small effects from common exposures, or large effects from rare exposures, are difficult to detect (**scale** - we need a lot of participants)

Need for a large-scale prospective cohort with deep characterization of participants and long follow-up



Why was UK Biobank established?





Overview of UK Biobank recruitment

- Between 2006-2010
- Aged 40-69 years old
- Registered with the NHS
- Living within ~25 miles of 1 of the 22 assessment centres





Data collected at UK Biobank Recruitment

**Touchscreen
questionnaire**

Demographics, environment, lifestyle, medical history, and cognitive function and hearing tests

**Verbal
interview**

Occupation, medical conditions, medications, operations, etc.

**Physical
measures**

Blood pressure, heart rate, anthropometry, spirometry etc.; and arterial stiffness, bone density, eye measures, fitness test (inc. ECG) in subset.

**Sample
collection**





Data collected at UK Biobank Recruitment

- Blood
- Urine
- Saliva

Total > 15 million aliquots



UK Biobank baseline characteristics



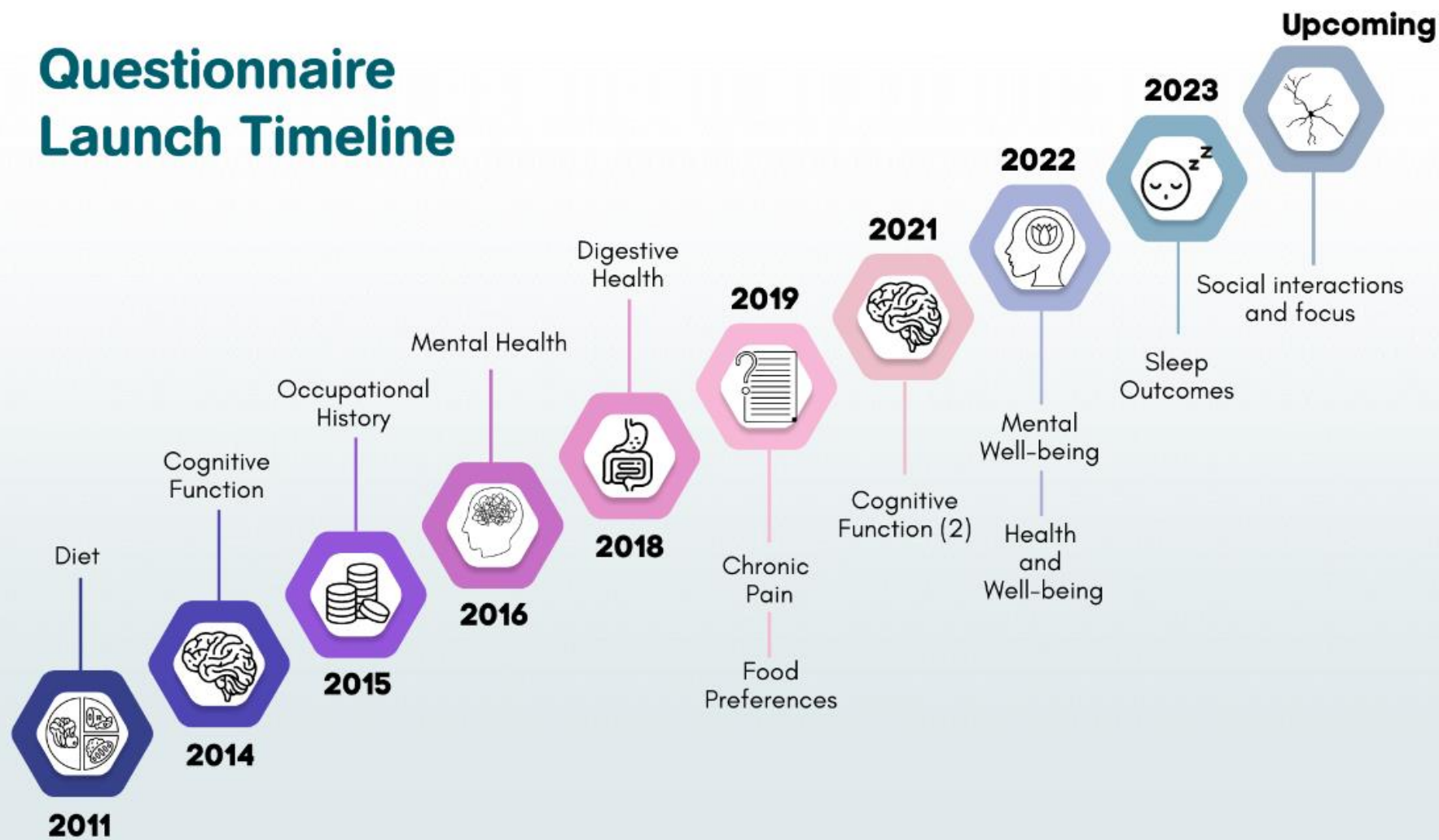
Characteristic	Category	Participants, n (%)
Age	40-49 yrs	119,000 (24%)
	50-59 yrs	168,000 (34%)
	60-69 yrs	213,000 (42%)
Sex	Male	230,000 (46%)
	Female	270,000 (54%)
Ethnicity	White	473,000 (95%)
	Other	27,000 (5%)
Deprivation	More	92,000 (18%)
	Average	166,000 (33%)
	Less	241,000 (46%)
Total		500,000

Wide range of backgrounds represented



Enhancements: Web-based questionnaires

Questionnaire Launch Timeline



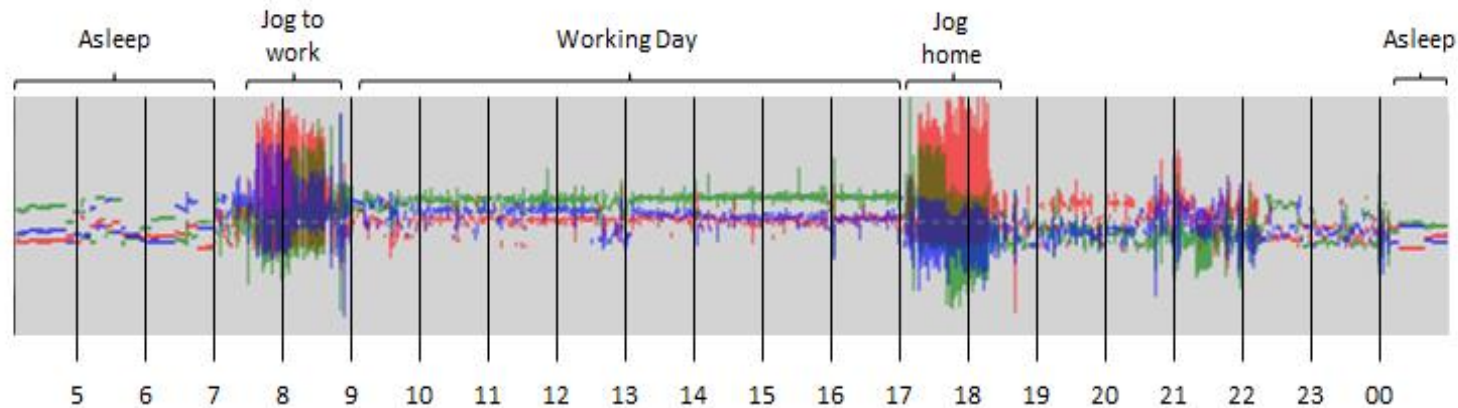
Among 330,000 participants for whom we have an email address.

Enhanced information on selected exposures and outcomes that is was not feasible to collect at baseline.



Accelerometer data

- 100,000 participants
- Worn continuously for 7 days
- 2,500 people repeated seasonally





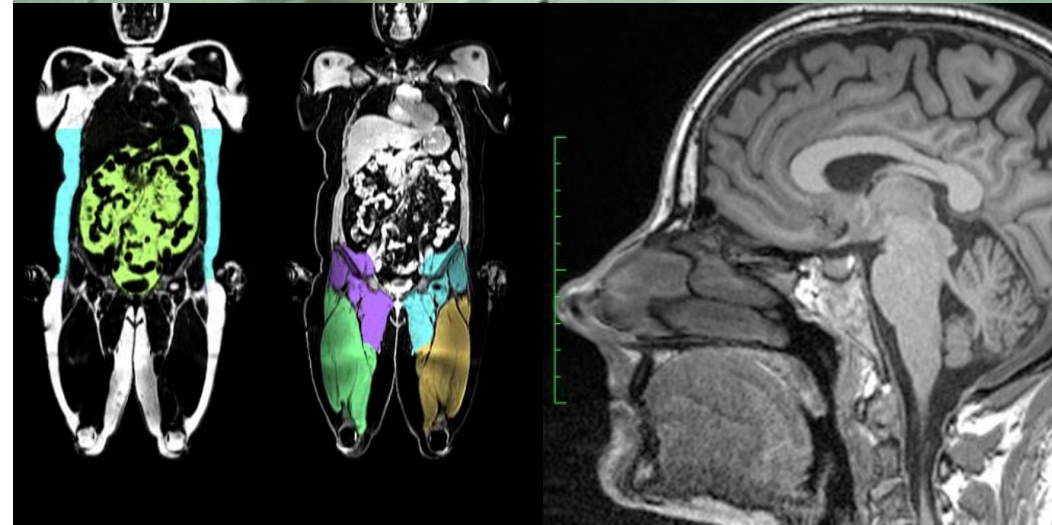
Repeat assessment in 20,000 participants (2012-13)

Multi-modal imaging (60,000 of 100,000 ppts; 2014-)

- MRI (heart, brain, abdomen)
- Full-body DEXA
- Carotid ultrasound
- 12-lead ECG

Repeat imaging underway (target of 60,000 ppts; 2019-)

The world's largest multi-modal imaging study





Genetics

Genome-wide genotyping

- 850k variants directly measured; >90M variants imputed
- 500,000 made available 2017



Whole Exome Sequencing

- 450,000 in Q4 2021; further 20,000 made available 2022

Whole Genome Sequencing

- First 200,000 made available Q4 2021
- Full cohort to be made available later in 2023





Enhancements: Samples into data

Biochemical measures in all 500,000

- 34 biomarkers in plasma, serum, red blood cells, and urine samples

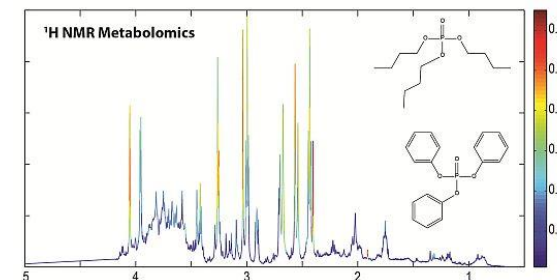
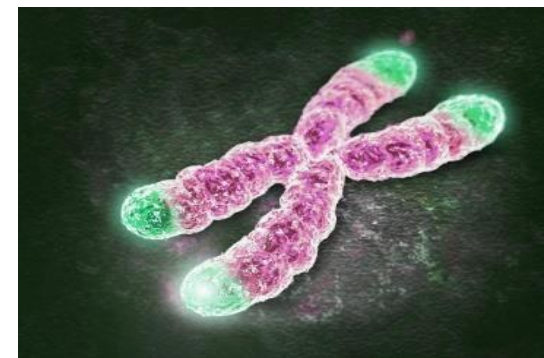
Telomere length in all 500,000

NMR-metabolomics in 300,000

- Data released for first 120,000 in 2021
- Second phase released Aug 2023

Proteomics in ~60,000

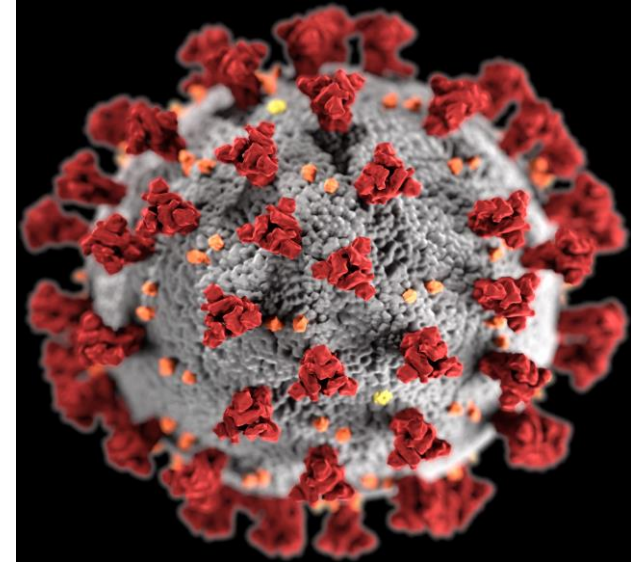
- Pharma consortium
- ~1500 plasma proteins using Olink's assay released
- Second phase release of further ~1500 plasma proteins later this year





Continuing data collection and high participant engagement enabled mobilisation of resources for COVID-19 research

- Enhanced COVID-related health outcome data
- Several COVID-related sub-studies
 - UK Biobank COVID-19 serology study
 - SARS-CoV-2 imaging study





700 research groups
300 papers so far

Gerontological Society of America
APOE e4 Genotype Predicts Severe COVID-19 in the UK Biobank Community Cohort
 Chia-Ling Kuo, PhD, Luke C. Milling, PhD, Janice L. Atkins, PhD, Jane A.H. Masolo, MSc, PhD, João Delgado, PhD, George A. Kuchel, MD, David Melzer, MBBCh, PhD
 Published: 26 May 2020
 Keywords: COVID-19, APOE, Genetic
 Issue Section: Letter to the Editor
 The novel respiratory disease COVID-19 produces varying symptoms, with fever, cough, and shortness of breath being common. In older adults, we found that pre-existing dementia is a major risk factor (odds ratio [OR] = 3.07, 95% CI: 1.71 to ...)

bioRxiv
 THE PREPRINT SERVER FOR BIOLOGY
 bioRxiv is receiving many new papers on coronavirus SARS-CoV-2. A reminder: these are preliminary reports that have not been peer-reviewed. They should not be regarded as conclusive, guide clinical practice/health-related behavior, or be reported in news media as established information.
Androgen Regulates SARS-CoV-2 Receptor Levels and Is Associated with Severe COVID-19 Symptoms in Men
 Zaniar Ghazizadeh, Homa Majd, Mikayla Richter, Ryan Samuel, Seyedeh Maryam Zekavat, Hosseinali Asgharian, Sina Farahvashi, Ali Kalantari, Jonathan Ramirez, Hongyu Zhao, Pradeep Natarajan, Hani Goodarzi, Faranak Fatahi
 doi: <https://doi.org/10.1101/2020.05.12.091082>
 This article is a preprint and has not been certified by peer review [what does this mean?]
 Abstract Full Text Info/History Metrics Preview PDF
 Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) infection has led to a global health crisis, and yet our understanding of the disease pathophysiology and potential treatment options remains limited. SARS-CoV-2 infection occurs through binding and internalization of the viral spike protein to angiotensin converting ...

ELSEVIER
 Brain, Behavior, and Immunity
 Available online 1 June 2020
 In Press, Corrected Proof
Ethnic disparities in hospitalisation for COVID-19 in England: The role of socioeconomic factors, mental health, and inflammatory and pro-inflammatory factors in a community-based cohort study
 Camille Lassale^{a, b, c, d, e}, Bamba Gaye^d, Mark Hamer^e, Catharine R. Gale^{f, g}, G David Batty^c
 show more

ELSEVIER
 Brain, Behavior, and Immunity
 Volume 87, July 2020, Pages 184-187
Lifestyle risk factors, inflammatory mechanisms, and COVID-19 hospitalization: A community-based cohort study of 387,109 adults in UK
 Mark Hamer^{a, b, c, d, e}, Mika Kivimäki^b, Catharine R. Gale^{f, g, h}, G. David Batty^b
 show more
<https://doi.org/10.1016/j.bbi.2020.05.059> Get rights and content

ELSEVIER
 Diabetes & Metabolic Syndrome: Clinical Research & Reviews
 Volume 14, Issue 4, July-August 2020, Pages 561-565
Vitamin D concentrations and COVID-19 infection in UK Biobank
 Claire E. Hastie^a, Daniel F. Mackay^a, Frederick Ho^a, Carlos A. Celis-Morales^{a, b}, Srinivasa Vittal Katikireddi^a, Claire L. Niedzwiedz^a, Bhautesh D. Jani^a, Paul Welsh^b, Frances S. Mair^a, Stuart R. Gray^a, Catherine A. O'Donnell^a, Jason MR. Gill^b, Naveed Sattar^{b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z}, Jill P. Pell^{a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z}
 show more
<https://doi.org/10.1016/j.dsx.2020.04.050> Get rights and content
 Highlights
 There is an urgent need to understand risk factors for contracting

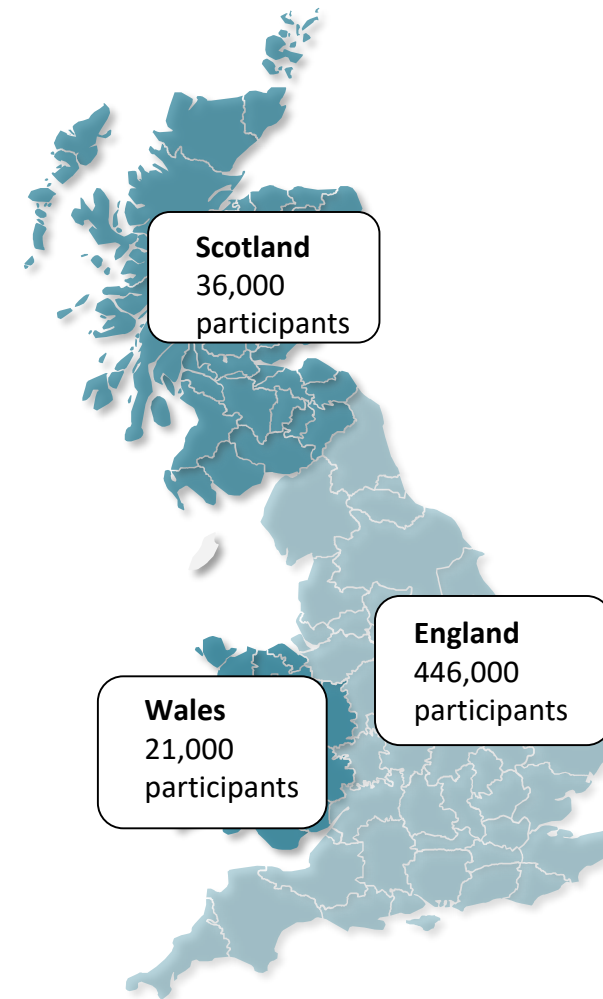
BMC Medicine
 Home About Articles Submission Guidelines
 We'd like to understand how you use our websites in order to improve our services for you.
 Research article | Open Access | Published: 29 May 2020
Ethnic and socioeconomic differences in SARS-CoV-2 infection: prospective cohort study using UK Biobank
 Claire L. Niedzwiedz, Catherine A. O'Donnell, Bhautesh Dinesh Jani, Evangelia Demou, Frederick K. Ho, Carlos Celis-Morales, Barbara L. Nicholl, Frances S. Mair, Paul Welsh, Naveed Sattar, Jill P. Pell & S. Vittal Katikireddi
 BMC Medicine 18, Article number: 160 (2020) | Cite this article
 2594 Accesses | 1 Citations | 176 Altmetric | Metrics



Follow-up of health outcomes

Regularly updated information on a wide range of diseases from NHS datasets in all 3 countries:

- **Deaths** (date and cause)
- **Cancers** (date and type)
- **Hospitalisations** (dates, diagnoses, procedures)
- **Primary care** (~45% of participants; dates, diagnoses)
- **COVID-related** (SARS-CoV-2 antigen tests)



Cumulative number of incident case over time

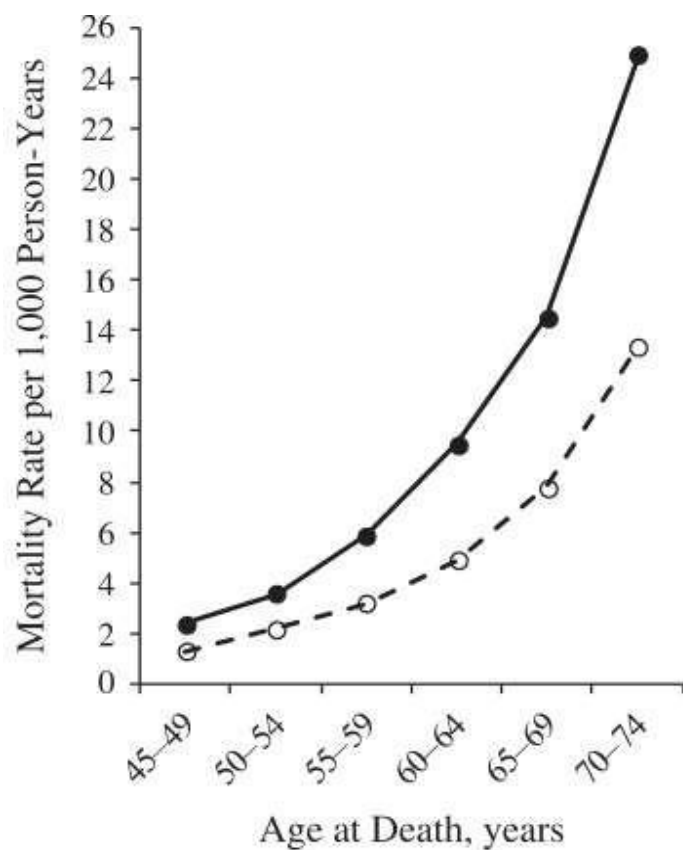


Condition	Year of diagnosis		
	Observed	Predicted	
	2022	2027	2032
Diabetes	31,000	54,000	70,000
Myocardial infarction	15,000	30,000	46,000
Stroke	12,000	25,000	37,000
COPD	25,000	47,000	65,000
Depression	25,000	39,000	47,000
Breast cancer	9,000	14,000	18,000
Colorectal cancer	5,000	8,000	11,000
Lung cancer	4,000	6,000	8,000
Prostate cancer	10,000	16,000	20,000
Hip fracture	5,000	13,000	22,000
Rheumatoid arthritis	4,000	6,000	8,000
Parkinson's disease	4,000	10,000	14,000
Alzheimer's disease	5,000	17,000	37,000

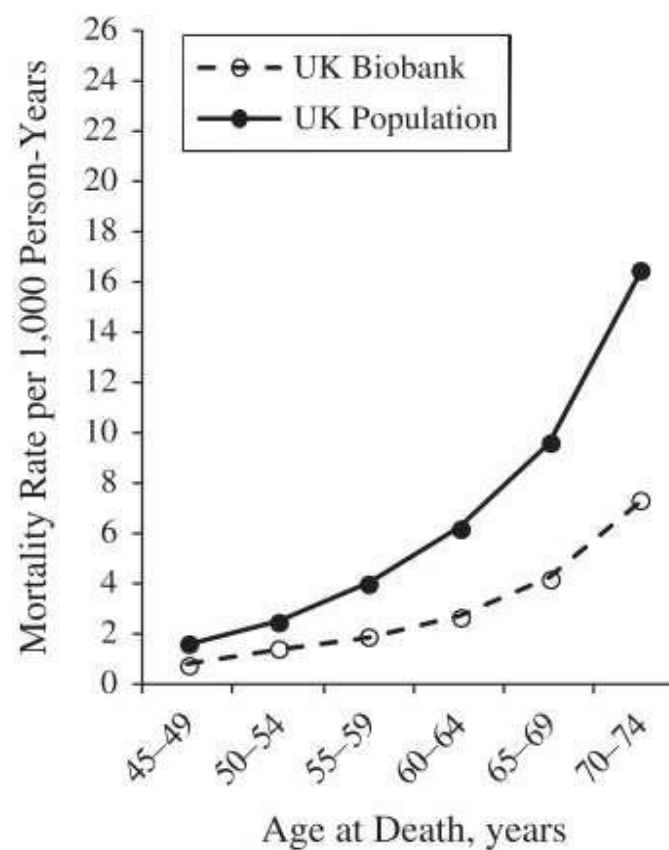


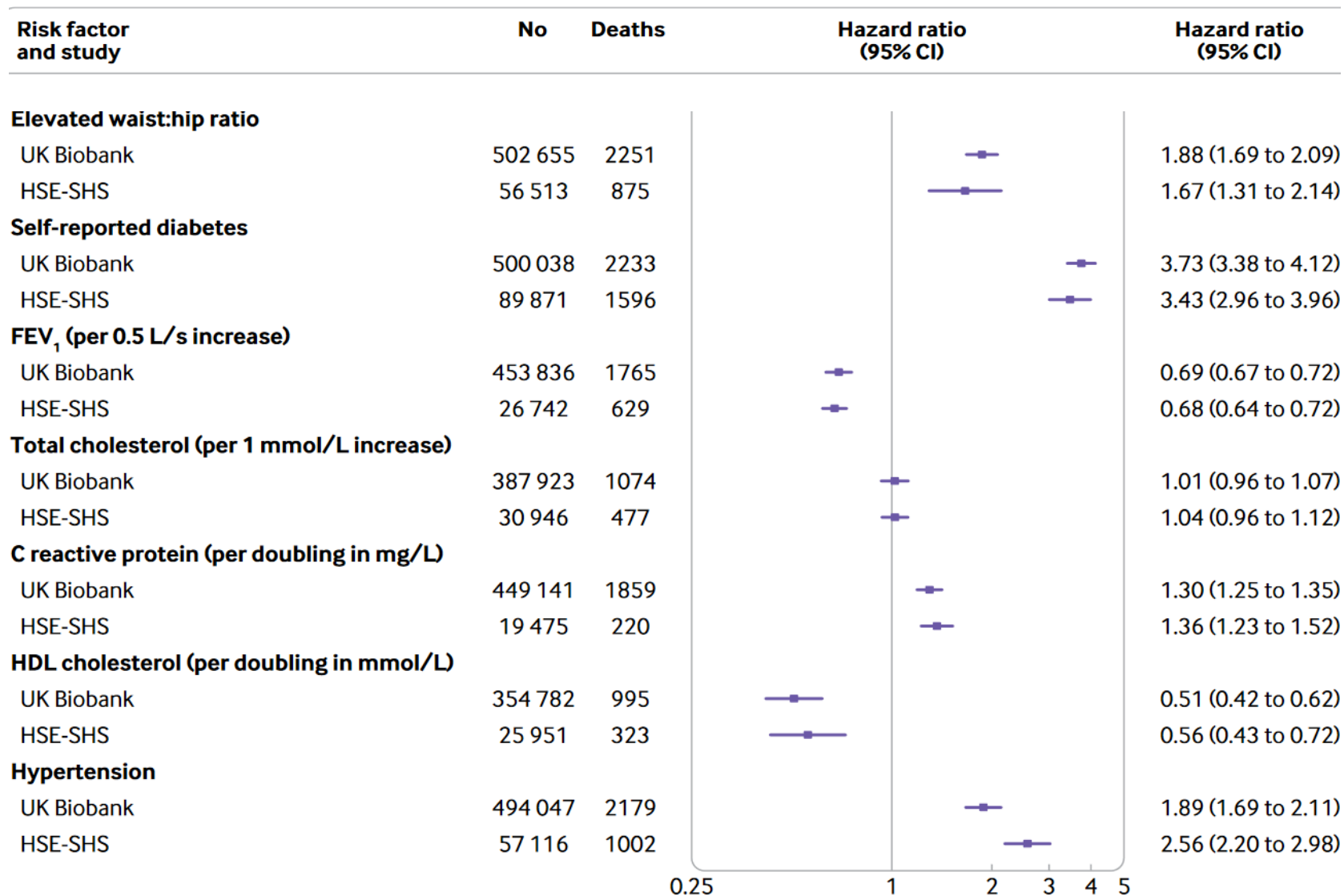
Mortality rates ~50% lower in UK Biobank

Men



Women





Hazard ratios for cardiovascular mortality.

HSE, Health Surveys for England, SHS, Scottish Health Surveys.

Batty et al. *BMJ* 2020;368:m131



Rationale for this work

- Important to identify as many cases as possible to increase the power of the study to investigate associations between risk factors and disease
- Reduce the number of false positive for disease, also important for study power
- Increase the range of diseases that can be studied
- To increase biological specificity of disease classification - increase the number of cases across the disease severity spectrum and across disease sub-types



Expanding health linkages with routine healthcare data



Potential new linkages:

- Mental health services
- Psychological therapy
- Microbiology
- Joint register
- Microbiology
- Clinical disease audits
- Ophthalmology datasets
- Govt. data (income, education..)



Based on:

- scientific added-value
- data quality
- comprehensiveness
- cost-effectiveness
- feasibility



Enhancing health outcomes

Use of online questionnaires to obtain self-reported data on health outcomes

Future web-based questionnaires: in development

- Neurodevelopment disorders (ADHD, autistic traits)
- Cognitive function outcomes
- Neurological disorders (lack of facial recognition and mental imagery)





Future efforts to identify dementia sub-types

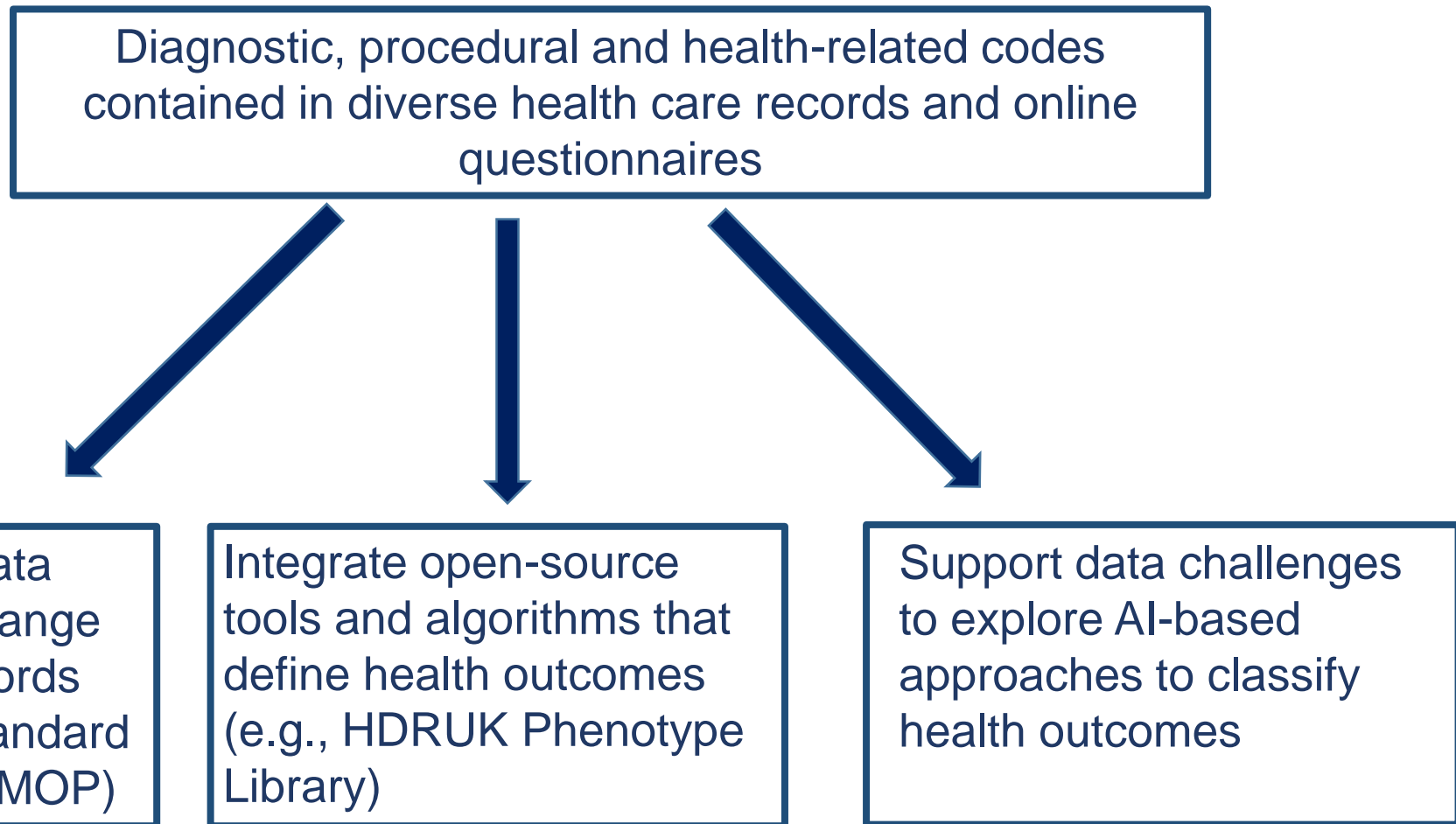
- Potential to invite participants with diagnosis of dementia to an assessment clinic for:
- **Blood-based neurodegenerative biomarkers**
 - E.g. Ptau-181 and -271, (A β)42, neurofilament light chain protein (NfL)
- **Imaging scans**
- **Clinical presentation**
 - Including possible use of apps





Enhancing health outcomes

Generation of 'research-ready' health outcomes





Accessing the data

Principles of accessing UK Biobank

- Open access resource, available for bona fide researchers to conduct health-related research that is in the public interest
- Available for use by academia and industry, both in the UK and overseas
- No preferential or exclusive access to the resource (*and limited exclusive access for data generated by researchers*)
- Researchers are obliged to return their results to UK Biobank so they can be shared with others



Accessing the data

Access procedure and options

Register and submit application



Application approved



Fees based on 'tier' of data requested

Tier 1: Clinic and follow-up data (£3k excl. VAT)

Tier 2: Above plus genotyping and other assay data (£6k excl. VAT)

Tier 3: Above plus sequencing and imaging data (£9k excl. VAT)



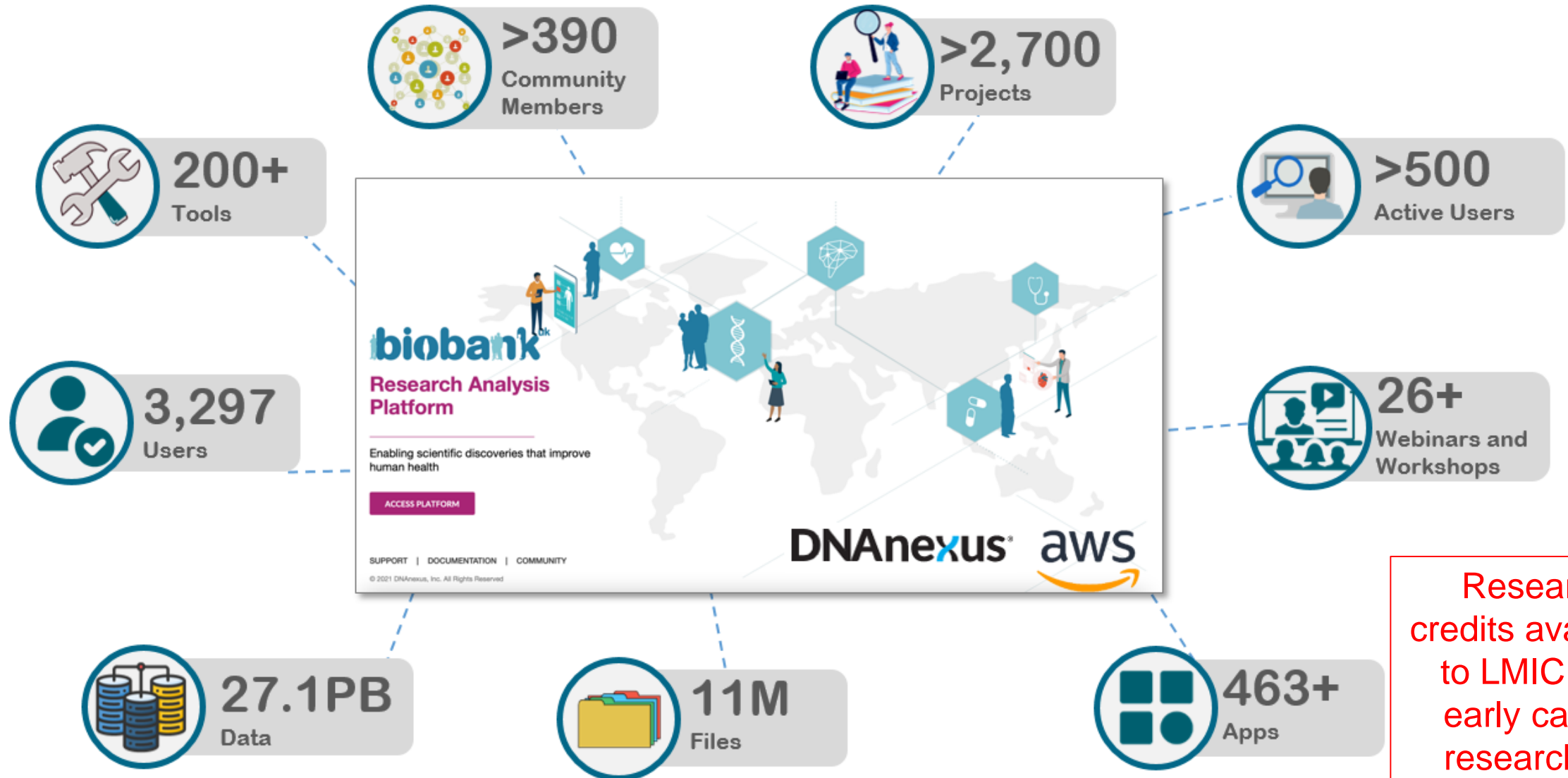
**Download dataset to
local computer**

**Access data in the Research
Analysis Platform**

- WES and WGS data only
available here



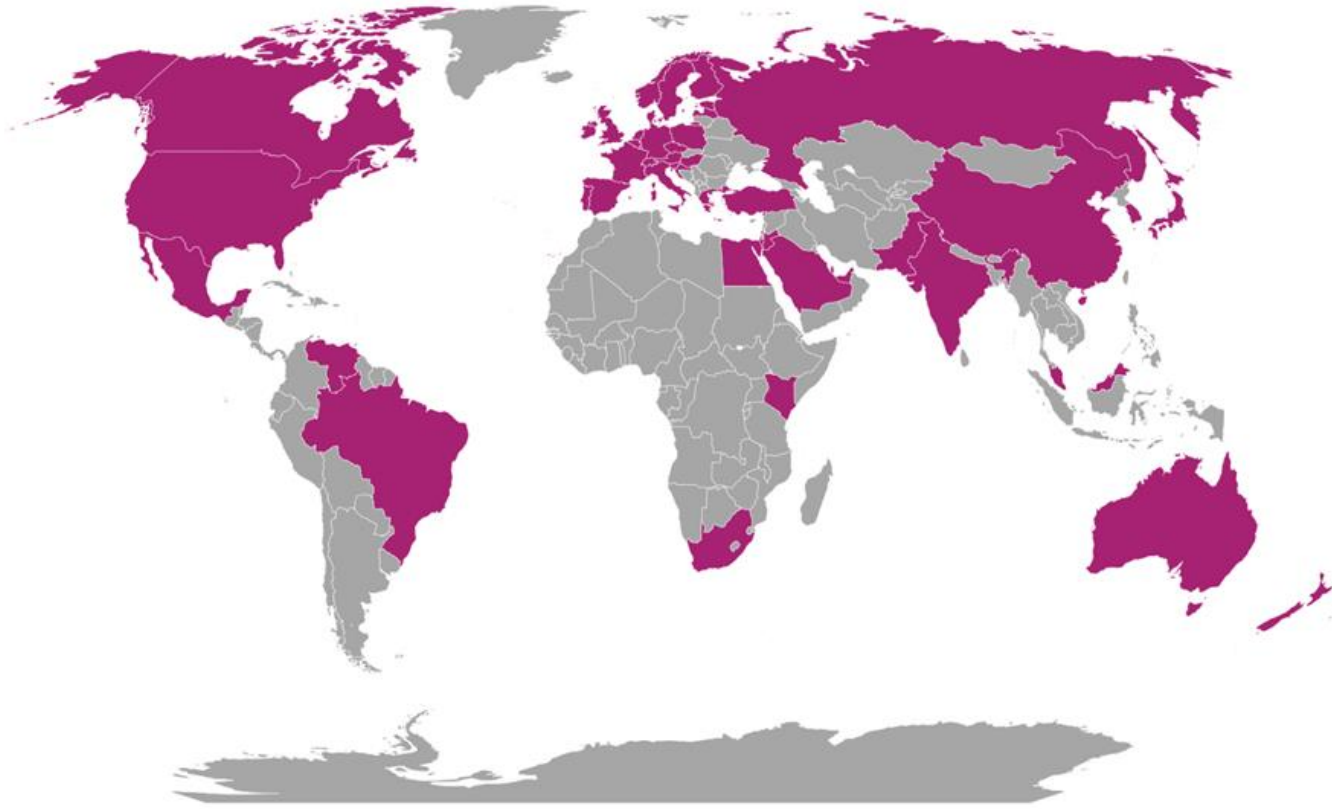
UKB Research Analysis Platform (RAP)



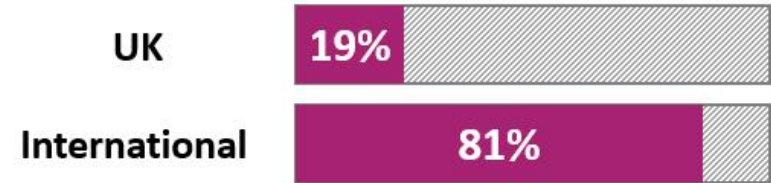
Research credits available to LMIC and early career researchers



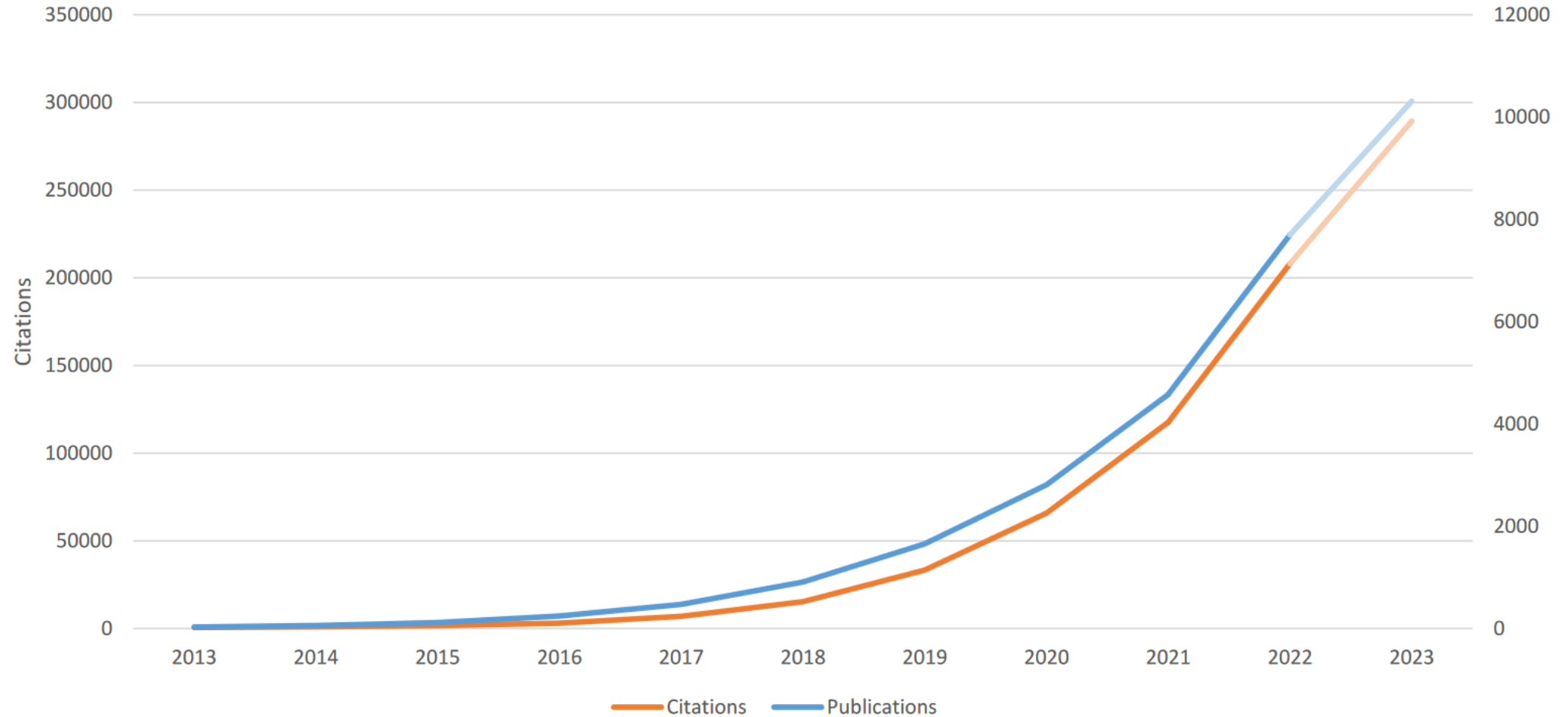
Who is using the UK Biobank dataset?



33,500 researchers
90 countries

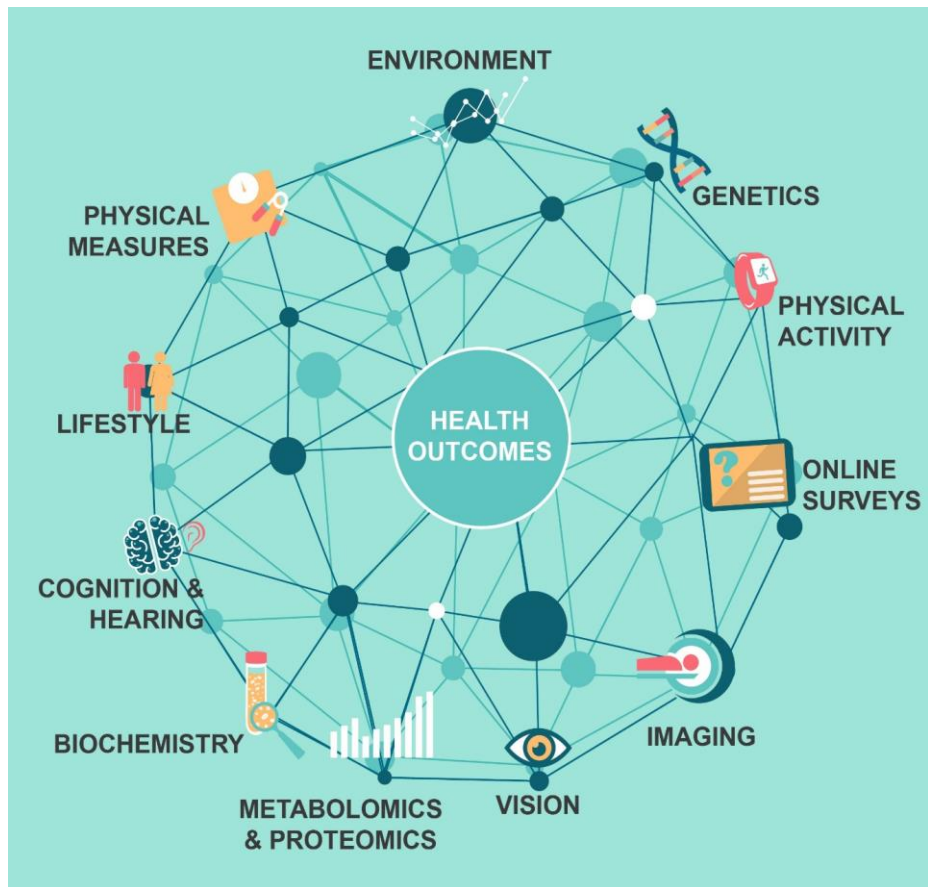


Publications and citations using UK Biobank





UK Biobank: Unique combination of 4 dimensions



Combination of size with increasing depth x duration x accessibility enabling cutting-edge science

- **SIZE:** 500,000 diverse individuals
- **DEPTH:** Genetics with extensive detail about lifestyle, environment and medical history, and other biological assays (biochemistry, genetics, -omics) and imaging
- **DURATION:** >10 years of follow-up has already yielded very large numbers of many different health outcomes
- **ACCESSIBILITY:** Very rapidly increasing number of different types of researcher globally using UK Biobank, facilitated by the UKB RAP

Acknowledgments

Core Funders



Medical
Research
Council



UK Biobank Team and Collaborators

Executive Team and Coordinating
Centre staff, Steering Committee,
International Scientific Advisory
Board, Scientific Working Groups,
Oxford University Team



Our 500k participants

