

# UK Biobank prospective cohort design and analytical considerations

Naomi E. Allen<sup>1,2</sup>, Ben Lacey<sup>1,2</sup>, Deborah A. Lawlor<sup>3,4</sup>, Jill P. Pell<sup>5</sup>, John Gallacher<sup>6,7</sup>, Liam Smeeth<sup>8</sup>, Paul Elliott<sup>9,10</sup>, Paul M. Matthews<sup>12</sup>, Ronan A. Lyons<sup>13</sup>, Anthony D. Whetton<sup>14</sup>, Anneke Lucassen<sup>15,16</sup>, Matthew E. Hurles<sup>17</sup>, Michael Chapman<sup>18</sup>, Andrew W. Roddam<sup>19</sup>, Natalie K. Fitzpatrick<sup>20</sup>, Anna L. Hansell<sup>21</sup>, Rebecca Hardy<sup>22</sup>, Riccardo E. Marioni<sup>23</sup>, Valerie B. O'Donnell<sup>24</sup>, Julie Williams<sup>25</sup>, Cecilia M. Lindgren<sup>26</sup>, Mark Effingham<sup>1</sup>, Jonathan Sellors<sup>1</sup>, John Danesh<sup>27,28,29</sup>, Rory Collins<sup>1,2</sup>

<sup>1</sup> UK Biobank Ltd, Stockport, UK

<sup>2</sup> Nuffield Department of Population Health, University of Oxford, Oxford, UK.

<sup>3</sup> Population Health Science, Bristol Medical School University of Bristol, Bristol, UK

<sup>4</sup> Medical Research Council Integrative Epidemiology Unit at the University of Bristol, Bristol, UK

<sup>5</sup> School of Health and Wellbeing, University of Glasgow, Scotland

<sup>6</sup> Department of Psychiatry, University of Oxford, Oxford, UK

<sup>7</sup> Dementias Platform UK

<sup>8</sup> London School of Hygiene and Tropical Medicine, London, UK

<sup>9</sup> MRC Centre for Environment and Health, School of Public Health, School of Public Health, Imperial College London, London, UK

<sup>10</sup> NIHR Biomedical Research Centre, Imperial College London, UK

<sup>11</sup> Health Data Research UK, Imperial College London, London, UK

<sup>12</sup> UK Dementia Research Centre Institute and Department of Brain Sciences, Imperial College London, UK

25 <sup>13</sup> Population Data Science, Swansea University Medical School, Swansea, Wales

26 <sup>14</sup> Veterinary Health Innovation Engine, University of Surrey, Guildford, UK

27 <sup>15</sup> Nuffield Department of Medicine, University of Oxford, Oxford, UK

28 <sup>16</sup> Faculty of Medicine, Southampton University, Southampton, UK

29 <sup>17</sup> Wellcome Trust Sanger Institute, Wellcome Genome Campus, Cambridge, UK

30 <sup>18</sup> NHS Digital, London, UK

31 <sup>19</sup> Our Future Health, London, UK

32 <sup>20</sup> Institute of Health Informatics, University College London, London, UK

33 <sup>21</sup> Centre for Environmental Health and Sustainability, University of Leicester, Leicester, UK

34 <sup>22</sup> School of Sport, Exercise and Health Sciences, Loughborough University, Loughborough,

35 UK

36 <sup>23</sup> Centre for Genomic and Experimental Medicine, University of Edinburgh, Edinburgh,

37 Scotland

38 <sup>24</sup> School of Medicine, Cardiff University, Cardiff, Wales

39 <sup>25</sup> UK Dementia Research Institute, Cardiff University, Cardiff, Wales

40 <sup>26</sup> Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of

41 Oxford, Oxford, UK

42 <sup>27</sup> British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public Health

43 and Primary Care, University of Cambridge, Cambridge, UK

44 <sup>28</sup> Health Data Research UK Cambridge, Wellcome Genome Campus and University of

45 Cambridge, Cambridge, UK.

46 <sup>29</sup> National Institute for Health Research Cambridge Biomedical Research Centre, University

47 of Cambridge and Cambridge University Hospitals, Cambridge, UK.

48

49 Corresponding author Emails:

50 [naomi.allen@ndph.ox.ac.uk](mailto:naomi.allen@ndph.ox.ac.uk)

51

52

53 **Overline: Biobanks**

54

55 **One sentence summary:** This article describes approaches to study design,  
56 resource access and data analysis in UK Biobank to facilitate health-related research

57

58 **Abstract**

59 Population-based prospective studies are valuable for generating and testing hypotheses  
60 about the potential causes of disease. We describe how the approach to UK Biobank's study  
61 design, data access policy, and statistical analysis can minimise error and improve the  
62 interpretability of research findings, with implications for other studies being established  
63 worldwide.

64

## 65 Introduction

66 Population health research has come a long way in the last few decades, with major  
67 advances in our understanding of the causes of disease. In particular, prospective studies  
68 that were initiated in the 1950s, such as the British Doctors Study (1) and the Framingham  
69 Heart Study (2), have been invaluable for understanding the association between lifestyle  
70 factors and disease risk as they overcome many of the biases inherent in case-control  
71 studies (most notably that exposures (i.e. risk factors for disease) are measured prior to  
72 disease onset). However, until recently, the conclusions that could be drawn from such  
73 studies were limited by small sample size, varying analytical approaches taken to define  
74 various risk factors and the relatively short duration of follow-up time to assess health  
75 outcomes. It was not until data from these different studies were integrated into large-scale  
76 individual-level meta-analyses that associations of exposures with disease risk were  
77 identified robustly. For example, it is now well established that circulating lipids and blood  
78 pressure are causally related to vascular disease (3), adiposity with cardiovascular disease  
79 (4), menopausal hormone therapy use and alcohol consumption with breast cancer (5, 6)  
80 and oral contraceptive use with a reduced risk of ovarian cancer (7).

81 More recently, there has been remarkable progress in research on the genetic  
82 determinants of disease. In the early 2000s, the literature was dominated by a plethora of  
83 genetic studies that focused on associations with particular conditions within specific  
84 “candidate” genes that were of *a priori* interest. Many of these studies involved small  
85 numbers of disease cases and yielded false-positive results that failed to replicate, often  
86 because of undue emphasis on *post hoc* selective reporting of the more extreme  
87 associations that were observed. Subsequently, improvements in assay technology led to  
88 genome-wide association studies (GWAS) that allowed hypothesis-free identification across  
89 the genome of variants associated with a particular phenotype. Much effort was typically  
90 spent on characterising the phenotype under investigation precisely in the belief that

91 outcome misclassification would have a substantial impact on the ability to detect  
92 associations. However, when meta-analyses of different studies were performed that yielded  
93 much larger numbers of individuals with the outcome of interest (albeit differently defined),  
94 small-to-moderate associations between genetic variants and outcomes began to be  
95 identified reproducibly after stringent adjustment for multiple testing (8).

96 Even larger sample sizes – of the order of hundreds of thousands of participants –  
97 are needed to study gene-environment interactions, especially where the genetic variant or  
98 environmental exposure of interest is rare or has a small effect on disease risk (9).  
99 Consequently, there is a strategic need to establish large-scale, well-characterised,  
100 population-based prospective cohorts in which biological samples are collected and health  
101 outcomes are followed long-term to facilitate research into the determinants of disease.

#### 102 **UK Biobank combines scale, depth, duration and accessibility**

103 UK Biobank is a population-based prospective cohort of 500,000 men and women designed  
104 to enable research into the genetic, lifestyle and environmental determinants of a wide range  
105 of diseases of middle-to-old age ([www.ukbiobank.ac.uk](http://www.ukbiobank.ac.uk)). It was established by the UK  
106 Medical Research Council (MRC) and Wellcome, which continue to fund it along with the  
107 British Heart Foundation (BHF), Cancer Research UK (CR-UK) and National Institute for  
108 Health and Care Research (NIHR). The key design features are its easy accessibility, large-  
109 scale prospective nature, depth and range of risk factor data, and comprehensive linkage to  
110 health outcomes, which together enable academic and industry researchers worldwide to  
111 perform discovery science (Supplementary Table 1).

112 UK Biobank was designed to promote innovative science by maximising access to  
113 the data in an equitable and transparent manner. All approved researchers (academic or  
114 commercial) can access all of the de-identified data in order to perform any type of health-  
115 related research that is in the public interest. This is the key criterion against which

116 applications to access the data are considered, with restrictions only placed on their use for  
117 potentially contentious research (for example, investigations that could lead to racial or  
118 sexual discrimination). Access to biological samples is currently largely restricted to assays  
119 that will be conducted on the whole cohort or large representative samples of the cohort.

120           Ready access to such a large-scale, in-depth resource has encouraged researchers  
121 from many disciplines across academia and industry to collaborate to ensure that different  
122 types of complex data (e.g., whole-exome and whole-genome sequencing data, magnetic  
123 resonance imaging (MRI) scans, accelerometer wave-form data, and electronic health  
124 records) are generated and analysed appropriately. The ready accessibility of the data at low  
125 cost without requiring collaboration with, or peer review from, the UK Biobank study  
126 investigators has led to an exponential increase in research output. By the end of 2023,  
127 there were more than 30,000 registered researchers (80% from outside the UK) and about  
128 9,000 publications (attracting 270,000 citations), with the number of publications increasing  
129 exponentially each year. In particular, the release to the worldwide research community of  
130 cohort-wide genome-wide genotyping and imputation data in 2017 has been hugely  
131 influential in advancing our understanding of the genetic determinants of disease.

132           The requirement that researchers publish their findings and make available any  
133 derived variables that have been generated as part of their research, together with the  
134 underlying code that generated the research output, enables the wider scientific community  
135 to critique, modify and build upon the work of others in a transparent manner (10). For  
136 example, research groups with expertise in signal processing have created derived variables  
137 related to the intensity and duration of physical activity from the raw accelerometer data (11,  
138 12). Similarly, academic and commercial research groups with expertise in image analysis  
139 have made available variables derived from the MRI scans related to body fat distribution  
140 (13), fat and iron content of specific organs (14, 15) and metrics of the structure and function  
141 of the brain (16) and heart (17). In this way, complex data that might otherwise only be of

142 use to specialists in a narrow field of research are turned into well-curated derived variables  
143 that are integrated with other UK Biobank data and can be used extensively by non-  
144 specialists to answer a range of research questions.

145 Easy access to such a wealth of data has led to new ways of presenting results. For  
146 example, summary statistics of all of the associations of individual genetic variants (18, 19)  
147 and polygenic risk scores (20) with a wide range of phenotypes are now available via online  
148 browsers. This move towards the publication of all summary results rather than publication of  
149 particular results in traditional scientific journals (where cherry-picking the most 'interesting'  
150 associations may introduce bias) is likely to accelerate scientific discovery and provide  
151 easier replication of associations across different studies. To help democratise access  
152 further, UK Biobank launched a cloud-based Research Analysis Platform in 2021 that allows  
153 streamlined access for researchers worldwide (in particular to the genome sequence data  
154 that are too large to transfer to researchers), as well as free computing and data storage for  
155 researchers from low- and middle-income countries and for early career researchers.

156 One consequence of researchers with different expertise accessing this wealth of data is the  
157 potential for unfamiliarity with various types of biases that are inherent in prospective studies  
158 that might influence results, as well as with the complexities associated with data that are  
159 outside of their areas of expertise. All researchers accessing biomedical resources to study  
160 the determinants of disease need to be aware of small sample size (that may produce  
161 imprecise estimates due to random error), incomplete or inadequate measurement of risk  
162 factors (that may lead to systematic under-estimation of disease associations), and health  
163 outcomes (that may lead to more imprecise estimates) and their potential confounding  
164 factors (that may obscure or lead to spurious associations between exposures and  
165 outcomes). Insufficient duration of follow-up may also lead to reverse causation bias,  
166 whereby the disease process influences potential risk factors (in particular, non-genetic  
167 ones), especially for conditions with a long prodromal phase, such as Alzheimer's disease.

168 UK Biobank has been set up to help minimise random and systematic error so that it  
169 can support reliable research into the determinants of disease (Supplementary Table 1),  
170 although the general principles of careful study design and appropriate data analysis apply  
171 equally to all large-scale, prospective studies. There are a number of trade-offs that need to  
172 be considered when designing a cohort study, which relate to the size and heterogeneity of  
173 the study population, and to the methods used for its recruitment, data collection and follow-  
174 up. UK Biobank has aimed to generate a large-scale, prospective biomedical resource that  
175 includes a wide range of exposure and health outcome measures collected as accurately as  
176 possible, with easy accessibility to the data. However, as with all prospective studies, it is  
177 important to consider, and if possible correct for, potential biases arising from the study  
178 design and collection of data.

#### 179 **The importance of a large-scale prospective design**

180 UK Biobank recruited 502,000 volunteers aged 40-69 years at recruitment between  
181 2006 and 2010 from across England, Wales and Scotland. This age group was selected to  
182 include individuals who were young enough that relatively few would have developed health  
183 conditions at the time of recruitment. As a prospective study, UK Biobank has many  
184 advantages for investigating the effects of genetic, lifestyle and environmental factors on  
185 disease outcomes (21). In particular, information on exposures to potential risk factors can  
186 be assessed before disease develops, which avoids bias caused by differential recall of  
187 information about past exposures depending on an individual's outcome status (recall bias).  
188 The prospective design also allows investigation of factors that might be affected by disease  
189 processes or their treatment, or by changes in an individual's behavior following the  
190 development of some condition (reverse causation bias). In addition, it can support studies  
191 of conditions that cannot readily be investigated retrospectively (e.g. fatal illnesses).  
192 Furthermore, by allowing a wide range of different conditions to be studied within the same  
193 study population, the full effects of a particular exposure on all aspects of health can be



194 better assessed (e.g. smoking on a wide range of different diseases). Likewise, the effects of  
195 many different exposures on a single disease can be determined, provided that sufficient  
196 numbers of cases have occurred to allow the separate and combined effects of exposures to  
197 be assessed reliably.

198           Prospective studies need to be large, as only a relatively small proportion of the  
199 participants will develop any given condition during follow-up. The rationale for recruiting  
200 500,000 adults into UK Biobank was that it would enable large numbers of cases of the most  
201 common diseases to develop within a reasonable follow-up period (while also allowing  
202 detailed exposure information to be collected within funding and organisational constraints).  
203 For example, after a median follow-up of 12 years (i.e. by end-2020), linkage to electronic  
204 healthcare record data indicated that there had been at least 30,000 incident cases of  
205 diabetes, 25,000 cases of depression, 15,000 cases of myocardial infarction, and 10,000  
206 cases of breast cancer (Table 1). For the reliable detection of risk ratios of about 1.3 for the  
207 main effects of different exposures (ranging from those that are dichotomous variables to  
208 those that are continuous measures), about 5,000-10,000 incident cases of a particular  
209 disease would be required (22). The need for a large sample size is even more evident when  
210 assessing combined effects. For example, when estimating the joint effect of blood pressure  
211 and age on the risk of coronary heart disease, the standard error of the estimates (and  
212 hence the 95% confidence intervals) are, on average, three times narrower with 500,000  
213 versus 50,000 participants (23). As the UK Biobank participants age, the number of incident  
214 cases of different diseases is increasing substantially, allowing a wider range of outcomes to  
215 be investigated more completely. For example, by 2032 there will be over 50,000 cases of  
216 diabetes and chronic obstructive pulmonary disease. The sheer size of the study also means  
217 that robust research into less common conditions will also be possible. For example,  
218 between 2020 and 2027, the number of cases of Alzheimer's disease, hip fracture and  
219 Parkinson's disease is expected to more than double to about 17,000, 13,000 and 10,000,  
220 respectively (Table 1).

## 221 **Comparing cohort characteristics with that of the wider population**

222 In UK Biobank, the well-defined sampling frame means that it is possible to assess  
223 not just the overall participation rate, but also the extent to which the study population differs  
224 from the wider population from which it was drawn. Postal invitations were sent to 9.2 million  
225 individuals aged 40–69, who were registered with the UK’s National Health Service (NHS)  
226 and lived within a short travelling time (typically about 25 miles) of one of 22 dedicated  
227 assessment centers. The choice of their location was determined by population density,  
228 ease of access, and potential to reach certain types of participants (e.g. ethnic minority  
229 groups and those living in more socio-economically deprived areas). During 2006-2010,  
230 502,000 participants were recruited (5.5% of those invited). Although the participation rate  
231 was low, and those who joined the study were somewhat healthier and wealthier than the  
232 UK population across the same age range (24), the cohort includes large numbers of  
233 individuals across a broad spectrum of risk factors (i.e. that vary from low to high exposure  
234 levels of a wide range of potential risk factors).

235 It is this heterogeneity across different levels of exposures (e.g., genetic, lifestyle,  
236 sociodemographic and environmental exposures) - and not the relatively low overall  
237 participation rate - that largely determines the generalisability of the findings to the broader  
238 UK population (25). For example, although individuals from more socio-economic deprived  
239 areas are under-represented in UK Biobank (16% versus 33% in the UK population), there  
240 are sufficiently large numbers of this group (82,000) to enable reliable assessment of the  
241 association of socio-economic deprivation with disease risk. By contrast, although UK  
242 Biobank is reasonably representative of the distribution for different ethnic groups, with  
243 29,000 participants recruited from Black and other ethnic minority groups (which was about  
244 the same proportion, ~5%, as the rest of the UK population at the time) (26), it is insufficient  
245 to examine reliably the differences in exposure-disease associations by ethnicity. Even  
246 though UK Biobank is currently the largest study in the world with whole-genome sequencing

247 data on individuals of African and South Asian ancestry (27), the numbers are still relatively  
248 small (with about 10,000 participants in each ethnic group).

249           Researchers who wish to present simple summary statistics (for example, means or  
250 proportions) using UK Biobank data that are representative of the underlying population  
251 could consider using sampling weights that reflect the population distribution of the variables  
252 under investigation, although such techniques have not been used widely. However, one  
253 research group found that standardisation of the prevalence of lifestyle factors with those  
254 derived from national survey data did not substantially alter the magnitude or direction of the  
255 association of lifestyle factors with mortality from cardiovascular disease or cancer (28). The  
256 one notable exception was an attenuation of the apparent protective association of alcohol  
257 with cardiovascular disease, which has been shown to be likely affected by bias (29).

258 There are circumstances where lack of representativeness may introduce bias, particularly if  
259 the risk factors of interest are related to study selection (an example of collider bias) (30).  
260 For example, UK Biobank participants are more likely to be non-smokers and to live in more  
261 affluent areas than the general population in the same age range. Given that area-level  
262 socio-economic deprivation is moderately inversely correlated both with participation in UK  
263 Biobank and lung cancer, this non-representativeness may attenuate the observed  
264 association of smoking with lung cancer if the effects of smoking and socio-economic  
265 deprivation are not independent or synergistic (31). Likewise, UK Biobank participants were  
266 more likely to use supplements and to have lower incident disease rates than the general  
267 population (at least in the early years of follow-up), leading to an apparent inverse  
268 association between glucosamine supplement usage and mortality (32). Analyses involving  
269 genetic variants that cluster by place of birth also have the potential to yield biased  
270 associations if standard variables such as assessment centre and ancestry-based principal  
271 components cannot completely correct for this latent structure (33). However, for most

272 genetic analyses where confounding from other risk factors is likely low, selection bias would  
273 typically be expected to be modest.

274 Consequently, in the interpretation of all research findings – whether they arise from the UK  
275 Biobank study or other prospective studies – it is important to consider the extent to which  
276 they might be affected by selective participation (i.e., selection bias). Given that traditional  
277 methods of identifying and controlling for selection bias (and other types of bias) may not be  
278 adequate, graphical tools such as directed acyclic graphs may provide a useful visual  
279 representation of the underlying assumptions about the relationships between exposures,  
280 potential confounders, mediators, and outcomes, and how they relate to study participation  
281 (34). Sensitivity analyses that include factors correlated with participation (and ongoing  
282 engagement) as covariates in the exposure-disease model can be performed; probability  
283 weighting, simulations and multiple imputation can be used to explore the potential impact of  
284 missing values related to participation on effect estimates (31, 35).

285 The general consistency of research findings in UK Biobank with those in other studies (36-  
286 38) – in particular, studies considered to be representative of the underlying population –  
287 suggest that many of the exposure-disease associations found in UK Biobank are largely  
288 generalizable to other populations. For example, the direction and magnitude of associations  
289 of genetic variants with osteoarthritis in UK Biobank are consistent with the associations  
290 observed in deCODE, which recruited more than half of Iceland's adult population (39).  
291 Likewise, although the frequency of genetic variants may vary substantially in studies  
292 conducted in different populations (resulting in differing statistical power to detect  
293 associations), the direction and magnitude of genetic associations are typically similar  
294 across populations e.g. the association of specific *GPR75* gene variants with obesity in UK,  
295 US and Mexico cohorts (40).

296 Nonetheless, there may be circumstances in which associations between an exposure and  
297 disease risk varies across different populations. For example, polygenic risk scores

298 developed and tested in populations of European ancestry often perform less well when  
299 applied to African and South Asian populations, owing to differences in allele frequencies  
300 and linkage disequilibrium patterns between the ethnic groups (41). As such, other large  
301 population cohorts with biological samples are needed around the world to increase the  
302 heterogeneity of genetic and non-genetic risk factors for disease (42) (Table 2). For  
303 example, studies established in Mexico (150,000 participants) and China (500,000  
304 participants) at about the same time as UK Biobank have enabled reliable investigation into  
305 the association between the risk of hypertension with body weight above and below the  
306 Western norm (43, 44). Large-scale studies established across Europe and China have also  
307 taken advantage of the heterogeneity of dietary and other exposures across different  
308 populations (45,46). Genetic and other assays of stored samples in these studies are  
309 extending the range of genomic risk factors that can now be investigated. New large-scale  
310 prospective studies are now established in the US e.g., All of Us (47) and the Million  
311 Veterans Program (48), and are also being established in Asia and parts of Africa e.g., Non-  
312 communicable Diseases Genetic Heritage Study in Nigeria (49, 50). This will further  
313 increase the ability to assess associations with disease risk across a broad range of genetic  
314 (and non-genetic) factors as long as there is sufficient duration of follow-up.

### 315 **Reliable assessment of a wide range of exposures**

316 The inclusion of participants exposed to different levels of risk factors (e.g. ranging from low  
317 to high intake of different dietary factors, smoking, sun exposure, etc.) is of value in  
318 assessing the generalisability of findings, which is enhanced further by analyses across  
319 studies established in different populations. However, all observational studies face  
320 challenges of exposure measurement error, in which risk factors and their potential  
321 confounders are measured imperfectly or incompletely, thereby introducing both random  
322 error (when measurements fluctuate randomly around their true value) and systematic error

323 (when measurements vary in the extent to which they are higher or lower than their true  
324 value).

325 As a result, UK Biobank has put significant effort into collecting comprehensive, accurate  
326 and high-quality data for many different types of exposures. Repeated measures have also  
327 been conducted in subsets of participants to address random error in exposure levels and  
328 thereby be able to correct for regression-dilution bias. However, there is real value in being  
329 able to perform cohort-wide repeat measures that would allow the relevance of individual  
330 changes in exposures over time to be assessed.

### 331 **Depth and breadth of exposure measurement**

332 In UK Biobank, a wide range of questionnaires and physical devices (e.g. spirometer to  
333 measure lung function, sphygmomanometer to measure blood pressure, bioimpedance  
334 device to measure body composition, dynamometer to measure hand grip strength, etc.)  
335 have been used (Fig. 1) to collect data that are reliable, valid and of high scientific value (26,  
336 51); such data continue to be collected and extended. During recruitment, UK Biobank used  
337 touch-screen and computer-assisted personal interview direct data-entry systems (instead of  
338 paper-based approaches that were routinely used at the time in such studies), as well as  
339 direct data transfer from measurement devices. This strategy enhanced data accuracy and  
340 completeness by supporting automated real-time consistency checks and data quality  
341 monitoring to identify and correct errors. Participants were also asked to bring certain  
342 information (e.g. medications, operations, family history, and birth details) to reduce errors  
343 associated with memory recall. However, perhaps the greatest benefit of using a touch-  
344 screen data entry model was that it reduced the time taken to collect data and thereby  
345 enabled a greater range of potential risk factors for disease to be collected. For example,  
346 data on sociodemographic factors (income, education, occupation), ethnicity, family history,  
347 lifestyle (diet, alcohol consumption, smoking history, physical activity, sleep, sun exposure,  
348 sexual history), early life factors, psychosocial factors, medical history, cognition and

349 environmental exposures were all collected via the touch-screen questionnaire in about fifty  
350 minutes.

351 A wide range of physical measurements were also taken for all 500,000 participants,  
352 comprising blood pressure, anthropometry (sitting and standing height, weight, waist and hip  
353 circumference, and bioimpedance measures), hand grip strength, vision and lung function.  
354 Blood and urine samples were also collected for long-term storage (Fig. 1). A proportion of  
355 the cohort also underwent a heel ultrasound for bone density, pulse wave velocity of arterial  
356 stiffness, a hearing test (180,000 participants), an eye examination (including refractive  
357 index), intraocular pressure measurements, a retinal photograph and optical coherence  
358 tomography (120,000 participants), a cardio-respiratory fitness test with a 4-lead  
359 electrocardiogram (ECG) (78,000 participants), and collection of a saliva sample (~85,000  
360 participants). Since the baseline assessment, UK Biobank continues to collect additional  
361 data from large subsets of the cohort. This has included data on physical activity using a 7-  
362 day accelerometer (in 100,000 participants, with 2,500 undergoing a repeat assessment), a  
363 multi-modal imaging assessment (in up to 100,000 participants, with 60,000 undergoing a  
364 repeat assessment over the next few years) and a series of web-based questionnaires that  
365 cover specific exposures in more depth (e.g. diet, cognition, occupational history).

366 Rigorous approaches have also been taken to sample collection, processing, retrieval and  
367 assay measurement. Prior to the start of UK Biobank, a series of pilot studies were  
368 conducted to determine the optimal method for sample collection and processing (52),  
369 followed by the development of a state-of-the-art robotic system and sample tracking  
370 software to ensure consistency of sample processing. Currently, genomic data (genome-  
371 wide genotyping and imputation, whole-exome and whole-genome sequence data, telomere  
372 length), as well as hematological and biochemical data are available for the whole cohort  
373 (Fig. 1). UK Biobank's general policy of performing cohort-wide assays supports research  
374 into a wide number of conditions and helps to avoid measurement errors that would

375 otherwise occur with different assay methods, reagents and equipment in different  
376 laboratories used in different subsets of the cohort at different times. To facilitate quality  
377 control, algorithms were developed to retrieve sample aliquots in a sequence that avoided  
378 clustering by recruitment location, date or time of day (53). Consequently, when assaying  
379 samples from participants in this quasi-random order, the mean biomarker concentration  
380 across batches and analysers should be constant, which allows correction for variation  
381 caused by laboratory drift. Throughout the assay process, the data are reviewed to identify  
382 issues and either address them in real time (e.g., if specific batches require re-  
383 measurement) or make any adjustments retrospectively. For example, following assay  
384 measurements of blood biochemistry markers, these data were corrected for systematic  
385 error caused by unexpected dilution that occurred in some aliquots during sample  
386 processing (53). Moreover, the use of highly efficient assay methods minimises sample  
387 depletion (with currently less than 10% of the baseline blood sample used so far), which will  
388 allow other types of assays (e.g., epigenetics, transcriptomics and proteomics) to be  
389 conducted on a cohort-wide basis when technological advances make this possible.

390 The collection of different types of data that describe the same (or highly related) exposures  
391 can also contribute to accuracy. In particular, a more precise assessment performed in a  
392 subset of participants could be used to correct for any random and systematic error inherent  
393 in the less precise baseline measures conducted in the full cohort (54). For example, data  
394 from an accelerometer device worn by 100,000 UK Biobank participants was used to  
395 calibrate self-reported physical activity estimates provided by all 500,000 participants at  
396 recruitment (55). Similarly, data on body fat composition available from dual-energy x-ray  
397 absorptiometry scans (56), which are being collected in up to 100,000 participants attending  
398 an imaging assessment, can be used to calibrate the bio-impedance measures available  
399 from the full cohort. Detailed dietary data from web-based questionnaires collected in over  
400 200,000 participants can also be used to predict food and nutrient intake for the entire  
401 cohort, as demonstrated in other studies (54).



402 The collection of data on a wide range of measures enables researchers to allow not only for  
403 more complete and accurate measurement of exposures, but also for potential confounders  
404 (extraneous factors that are associated with the exposure and outcome) and mediators  
405 (factors that are on the causal pathway between the exposure and the outcome). This is  
406 important, as random error in exposure measures can cause systematic attenuation of any  
407 true association, whereas random measurement error of confounders can result in an  
408 apparent exposure-disease association, where none really exists. For example, the  
409 observed inverse association of fruit and vegetable intake with cardiovascular disease risk in  
410 UK Biobank is likely to be due largely to residual confounding by socio-economic factors,  
411 which are difficult to assess and therefore subject to measurement error (57). The ability of  
412 UK Biobank to obtain more detailed information in the future about socio-economic factors  
413 (such as education, occupation and income via linkage to administrative datasets or specific  
414 web-based questionnaires) would enable more precise characterisation and, hence, even  
415 better adjustment for these important factors.

416 Because all epidemiological studies suffer, to a greater or lesser extent, from imperfect  
417 measurement of exposures and their potential confounders, various analytical methods have  
418 been developed to quantify and control for this. Perhaps the simplest approach is the  
419 comparison of likelihood ratio statistics associated with the exposure of interest in the  
420 models before and after adjustment for covariates. Generally speaking, a large proportional  
421 reduction in the likelihood ratio chi-square ( $LR\chi^2$ ) test after the addition to the model of  
422 covariates is indicative that the association likely remains affected by residual confounding,  
423 as adjustment for confounders that are perfectly measured would be expected to reduce the  
424  $\chi^2$  even further (6). An increasingly popular approach to distinguish the likely causal effect of  
425 an exposure (from that of extraneous confounders) is the use of Mendelian Randomisation –  
426 facilitated in analyses of UK Biobank by the extensive genetic information available on all of  
427 the participants – whereby specific genetic variants are used as proxies for exposures of  
428 interest or their confounders. For example, this approach has provided strong support for a

429 causal role of body fat mass and interleukin-6 in development of cardiovascular conditions  
430 (58, 59). Conversely, Mendelian Randomisation has not provided support for a protective  
431 effect of vitamin D against COVID-19 (60), cancer or cardiovascular outcomes (61),  
432 although it should be noted that Mendelian Randomisation analyses may also be affected by  
433 bias in some circumstances (62). When associations of genetic variants with the relevant  
434 non-genetic risk factors are weak (such that Mendelian Randomisation may not be effective),  
435 the likely impact of residual confounding due to imprecision in measured variables included  
436 in the model can be assessed using other analytical approaches such as probabilistic or  
437 multiple-bias analysis (34, 63). The use of different analytical strategies to triangulate  
438 evidence (for example, comparing results from models that include traditional observational  
439 variables with those that use genetic instrumental variables) will enable researchers to  
440 assess different biases and their potential impact on causal inference in a more robust  
441 manner.

#### 442 **Repeated exposure measurements**

443 Random errors in the measurement of risk factors can lead to substantial underestimation of  
444 the strength of their associations with subsequent health outcomes (regression dilution bias)  
445 (64, 65), as well as to substantial residual confounding when measurement error is present  
446 in confounders (66). These biases may be increased further through random error in risk  
447 factor measurements that occur during prolonged follow-up in prospective cohorts. For  
448 example, the true association of blood pressure and cholesterol with cardiovascular disease  
449 risk may be underestimated by about one-third in the first decade of follow-up and up to two-  
450 thirds in the third decade (64). However, despite regression dilution being one of the most  
451 important biases in exposure-disease associations, it is often overlooked in analyses of  
452 prospective studies, including UK Biobank (with some exceptions) (67-70). It is possible to  
453 correct for regression dilution bias by using repeat measures from a relatively small subset  
454 of the cohort. UK Biobank performed a repeat assessment on 20,000 participants in 2012-

455 2013 to allow researchers to address this issue specifically. Re-measures collected during  
456 the imaging assessment of up to 100,000 UK Biobank participants during 2014-2024 and a  
457 repeat assessment of up to 60,000 during 2019-2029 can be used to make appropriate time-  
458 dependent corrections for the effects of regression dilution bias.

459 In addition to addressing error caused (largely) by random error in baseline risk factors,  
460 repeated measures would also enable correction for systematic intra-individual changes in  
461 exposures over time, which may lead to either over-estimation or under-estimation of  
462 associations depending on the nature and magnitude of misclassification. For example,  
463 secular trends in the reduction of smoking or exposure to environmental pollutants may lead  
464 to an underestimation of their association with disease risk if solely based on baseline  
465 measures. To help address this issue, UK Biobank is exploring opportunities to collect  
466 information on longitudinal changes in environmental exposures (e.g. from existing data on  
467 changes in participants' residential location or future data collection using smartphone GPS  
468 tracking) to enable more accurate inferences to be made about how changes in  
469 environmental exposures affect health in the long-term. It is also intended to repeat previous  
470 web-based questionnaires in order to capture longitudinal changes in specific lifestyle factors  
471 such as diet and sleep.

472 Whereas repeated measures of the baseline assessment are being captured during the  
473 imaging assessments in a subset of the UK Biobank cohort, it would be desirable to perform  
474 a future repeat assessment of a wide range of exposures in as many of the participants as  
475 possible. This would allow investigation of how lifestyle, and physical and biochemical  
476 changes over time influence disease risk and progression, thereby helping to determine the  
477 temporality of associations and their underlying mechanisms. Data collection for as many  
478 surviving participants as possible would also reduce systematic error caused by differential  
479 participation rates that are related to the exposures and outcomes under investigation. UK  
480 Biobank generally has excellent participant engagement with an ongoing series of repeated

481 web-based questionnaires (with response rates of >50%), physical activity monitoring (45%  
482 for the first assessment, of whom 63% also performed a repeat assessment), and imaging  
483 assessments (24% for the first assessment and 65% for a repeat assessment). However,  
484 researchers should be aware that participants who engage in ongoing data collection  
485 activities (including repeat assessments) might not be representative of the cohort as a  
486 whole. For example, genetic variants associated with completing UK Biobank online  
487 questionnaires and activity monitoring are correlated with several metrics of better health  
488 (31). Attrition bias has been documented in other prospective studies (71-73), suggesting  
489 that similar factors affect ongoing participant engagement in many cohorts, regardless of  
490 their design, recruitment strategy or study population.

#### 491 **Reliable assessment of a wide range of health outcomes**

492 To minimise bias in exposure-disease associations, it is important that health outcomes are  
493 identified in a comprehensive manner and in as much detail as possible. Linkage to routine  
494 electronic health records, supplemented with information from self-reported questionnaires  
495 and other remote methods, and in-person assessments focused on specific outcomes (such  
496 as dementia), will help to deeply characterise health outcomes that are of high priority. The  
497 ability to combine these data from disparate sources to generate 'off-the-shelf' outcomes that  
498 can be easily interpreted by non-specialists will further increase the usability and  
499 reproducibility of research using these data.

#### 500 **Comprehensive ascertainment of health outcomes**

501 All cohort studies need a comprehensive and efficient way of following participants' health  
502 over the long-term to identify a wide range of disease outcomes. Unlike many countries  
503 (including the US and most low-to-middle income countries), the UK's National Health  
504 Service (NHS) collates and stores electronic health administrative records for clinical care.  
505 However, the data content, format and governance requirements may differ for England,

506 Wales and Scotland. To identify a wide range of health outcomes over a prolonged period,  
507 UK Biobank has linked to these health administrative records for all participants. This has  
508 the advantage of minimising ascertainment bias and reducing loss-to-follow-up or attrition  
509 bias by providing cohort-wide follow-up information without the need for active participant re-  
510 contact, which may be incomplete. Moreover, the low rate of UK Biobank participants  
511 requesting that all of their data and samples be withdrawn from the study (0.2%; most of  
512 which occurred soon after recruitment) also minimises systematic bias associated with loss  
513 to follow-up from non-random subgroups of the cohort.

514 To date, UK Biobank has linked NHS healthcare data from centralised national cancer and  
515 death registries and hospital inpatient admissions for all participants. In contrast, primary  
516 care data are not centralised but instead are held by commercial electronic system suppliers  
517 under the control of individual general practices, so it has been more challenging to obtain  
518 the agreements necessary to obtain these data for all participants. Primary care data are  
519 currently available for 45% of the UK Biobank cohort for general research purposes (which  
520 represents complete coverage from one primary care system supplier, up to 2016/2017) and  
521 for 80% of the cohort for COVID-19 research (complete coverage from two system suppliers  
522 in England, up to mid-2021, enabled by emergency legislation to facilitate COVID-19  
523 research). Whereas both subsets are broadly representative of the cohort with respect to the  
524 distribution of potential exposures, researchers should be encouraged to check these  
525 underlying assumptions prior to analysis. Incorporation of primary care data for all 500,000  
526 participants for all types of health-related research would be of enormous value as it will  
527 increase substantially the number of health outcomes that can be detected (thereby  
528 increasing statistical power) and their diagnostic accuracy (thereby increasing specificity).  
529 For example, whereas addition of primary care data would increase the numbers of  
530 myocardial infarction cases identified by less than 5%, the numbers of cases identified of  
531 diabetes and chronic obstructive pulmonary disease (COPD) would increase by about 40%  
532 (Fig. 2). Primary care data are also important for investigating risk factors associated with

533 disease severity, where associations may differ between milder disease subtypes generally  
534 captured in primary care records and more severe disease captured in hospital admission  
535 data.

536         Whereas linkage to health records ensures comprehensive coverage, there is the  
537 possibility of “collider bias” if health outcomes are differentially ascertained based on  
538 participant characteristics (e.g., ethnicity), as reported by some researchers in the context of  
539 COVID-19 research (74). However, there are a range of analytical approaches that can be  
540 used to investigate this type of bias (74-76) and the ascertainment of most health outcomes  
541 are not so strongly influenced by these characteristics.

#### 542 **Specificity of health outcomes**

543 Given that the prospective nature of cohort studies facilitates research into many diseases,  
544 the challenge is not only how to identify probable cases of disease but also to increase the  
545 precision and specificity of those diagnoses. The aim is to avoid a situation where insufficient  
546 data on health outcomes leads to misclassification of cases and non-cases, thereby  
547 reducing statistical power to detect an association. As such, UK Biobank’s aim is to ascertain  
548 as many cases as possible (i.e., to achieve adequate sensitivity) while minimising the  
549 number of false-positive cases (i.e., achieving a high positive predictive value). It is worth  
550 recognising that it is not necessary to identify all cases of a disease as false negatives will  
551 be diluted by the much larger number of ‘true’ controls (and so have limited impact). To help  
552 identify as many cases as possible, UK Biobank administers various web-based  
553 questionnaires, developed in close collaboration with relevant experts, to collect data on  
554 health outcomes that are incompletely recorded in health records, such as depression and  
555 anxiety (77), and neurodevelopmental and gastrointestinal conditions.

556 It is also important to characterise disease subtypes as low biological specificity can limit  
557 interpretation of results. To address this, UK Biobank (78-80) and other open-access

558 resources (81) have developed a number of algorithmically defined health outcomes based  
559 on inter-operable code lists from electronic healthcare records. Diagnostic codes contained  
560 in these records have also been mapped to a common standard (ICD-10) to facilitate broad-  
561 brush research. Whereas these coded health outcomes may be sufficient for most research  
562 purposes, they may lack specificity to identify disease subtypes. This could lead to materially  
563 biased estimates of associations if the determinants of these apparently similar, but  
564 etiologically different, disease subtypes differ. For example, while blood pressure is strongly  
565 positively associated with the risk of both ischaemic and haemorrhagic stroke (82), the  
566 association of cholesterol and certain genetic variants with stroke differ substantially by  
567 subtype (83, 84) providing clues to the underlying aetiology of this heterogeneous condition.  
568 To increase the specificity of health outcomes beyond the available coded data, UK Biobank  
569 intends to collect detailed data on disease sub-types over the next few years. For example,  
570 this could include disease-specific registers such as the National Diabetes Audit that collects  
571 data on diabetes subtypes, clinical scans to identify stroke sub-types, digitised  
572 histopathology slides to determine tumour morphological subtypes, and in-person  
573 assessments to characterise dementia subtypes.

574 It is possible to identify some disease sub-types using other data already available in the UK  
575 Biobank resource. For example, biochemistry measures have been used to ascertain  
576 chronic kidney disease (85), MRI scans collected in up to 100,000 participants have been  
577 used to define dilated cardiomyopathy (86) and non-alcoholic fatty liver disease (87), and  
578 genetic data have been used to distinguish diabetes subtypes (88). However, researchers  
579 should be aware of the potential for generating misleading associations where the exposure  
580 of interest (e.g. genetic variants or biochemistry measures) has, in part, been used to define  
581 the outcome.

## 582 **Long duration of follow-up**

583 For any prospective study, a long duration of follow-up (i.e. decades or more) is needed for  
584 sufficiently large numbers of health outcomes to accrue for reliable investigation. It also  
585 allows for the identification of recurring events and factors associated with disease  
586 progression. While the incidence of common health outcomes during the early years of  
587 follow-up in UK Biobank was somewhat lower than in the general population due to the  
588 'healthy volunteer' effect, which is typical of such studies (89), its impact is now reduced as  
589 the cohort has aged. With prolonged follow-up, large numbers of incident cases of a wide  
590 range of conditions have already occurred. Over the next five to ten years there will be  
591 thousands of incident cases of common outcomes (Table 1), enabling reliable investigation  
592 of their genetic, lifestyle and environmental determinants.

593 The rationale for recruiting middle-aged participants was to collect risk factor data many  
594 years before the development of any given condition, thereby minimising reverse causation  
595 bias. However, conditions that have a long prodromal phase (e.g. dementia or diabetes) or  
596 that can exist for years before a clinical diagnosis is made (such as prostate cancer) may  
597 affect the levels of risk factors measured at recruitment and create spurious associations.  
598 For example, associations observed between high insulin-like growth factor-I (IGF-I)  
599 concentrations and increased risks of cataract and diabetes were substantially attenuated  
600 after excluding the first five years of follow-up in UK Biobank (90), suggesting that baseline  
601 IGF-I concentrations may be altered as a result of early pathophysiological processes. Other  
602 large-scale longitudinal studies have also shown that apparent inverse associations between  
603 lifestyle factors and dementia risk are also likely to be due to reverse causation bias during  
604 the first 10-15 years of follow-up (91). Consequently, researchers should consider the impact  
605 of exclusion of participants with prevalent disease prior to analysis and perform sensitivity  
606 analyses to assess exposure-disease associations across different periods of follow-up to  
607 determine whether the first years of follow-up should be excluded (92).

## 608 **Conclusions**



609 The success of UK Biobank has been due, in large part, to the altruism of the 500,000  
610 volunteers, but also the global research community who have been – and continue to be –  
611 involved in expanding the range of exposures and outcomes available for research. Such  
612 enhancements (e.g. sample assays, linkage to specific healthcare datasets and  
613 environmental sources, etc.) help to ensure that the resource fulfils the needs of researchers  
614 and remains at the forefront of scientific discovery.

615 UK Biobank's large-scale prospective design and easy access to a wealth of genetic,  
616 phenotypic and health data provides a powerful resource to help address previously  
617 unanswerable questions of the determinants of incident disease, as well as enabling  
618 research into risk prediction and identification of early biomarkers of disease. Whereas the  
619 UK Biobank study has attempted to minimise random and systematic errors in the  
620 measurement of exposures and outcomes with good study design, researchers need to use  
621 the data in ways that best answer the questions posed, and to be aware of and, where  
622 necessary, use analytical methods to take account of potential biases when interpreting  
623 research findings.

624 Easy accessibility of UK Biobank data and research results (including the underlying  
625 analytical code) is enabling the community to directly peer review research by undertaking  
626 replication analyses, or to apply different methods to the same research question, to confirm  
627 or refute the findings of others. In particular, investigation of approaches used to identify and  
628 quantify the uncertainty of the results based on sensitivity analyses that examine systematic  
629 bias, will provide a level of transparency in the interpretation of findings that has, until now,  
630 generally been under-reported.

631 Whereas UK Biobank is well suited to address a wide range of health-related research  
632 questions, similar studies in other populations with different ranges of exposures and  
633 outcomes are needed. Taken together, they will enable a greater range of risk factors and

634 diseases to be analysed and allow for replication of associations, which is essential before  
635 determining the extent to which any specific research findings are generalizable to different  
636 populations. Scientific discoveries benefit from the availability of data from diverse  
637 populations that cover a wide range of the many different genetic, ancestral, ethnic, lifestyle  
638 and environmental factors that may influence risk of a broad range of diseases.

639

640

641

642

643

644

645

646

647

648

#### 649 **References and Notes**

650 (1) R. Doll, A.B. Hill, Lung cancer and other causes of death in relation to smoking; a second  
651 report on the mortality of British doctors, *BMJ*, **2**, 1071-1081 (1956).

652

653 (2) J. Truett, J. Cornfield, W. Kannel, A multivariate analysis of the risk of coronary heart  
654 disease in Framingham, *J Chronic Dis*, **20**, 511-524 (1967).

655

656 (3) Emerging Risk Factors Collaboration, Lipoprotein(a) concentration and the risk of  
657 coronary heart disease, stroke, and nonvascular mortality, *JAMA*, **302**, 412-423 (2009).

658

659 (4) Emerging Risk Factors Collaboration, Separate and combined associations of body-mass  
660 index and abdominal adiposity with cardiovascular disease: collaborative analysis of 58  
661 prospective studies, *Lancet*, **377**, 1085-1095 (2011).

662

663 (5) Collaborative Group on Hormonal Factors in Breast Cancer, Breast cancer and hormone  
664 replacement therapy: collaborative reanalysis of data from 51 epidemiological studies of  
665 52,705 women with breast cancer and 108,411 women without breast cancer. *Lancet*, **350**,  
666 1047-1059 (1997).

667

668 (6) Collaborative Group on Hormonal Factors in Breast Cancer, Alcohol, tobacco and breast  
669 cancer - collaborative reanalysis of individual data from 53 epidemiological studies, including  
670 58,515 women with breast cancer and 95,067 women without the disease, *Br J Cancer*, **87**,  
671 1234-1245 (2002).

672

673 (7) Collaborative Group on Hormonal Factors in Ovarian Cancer, Ovarian cancer and oral  
674 contraceptives: collaborative reanalysis of data from 45 epidemiological studies including  
675 23,257 women with ovarian cancer and 87,303 controls, *Lancet*, **371**, 303-314 (2008).

676

677 (8) E. Uffelmann, Q.Q. Huang, N.S. Munung, J. de Vries, Y. Okada, A.R. Martin, C.M.  
678 Martin, T. Lappalainen, D. Posthuma, Genome-wide association studies. *Nat Rev Methods*  
679 *Primers*, **1**, 59 (2021).

680

681 (9) J.A. Luan, M.Y. Wong, N.E. Day, N.J. Wareham, Sample size determination for studies  
682 of gene-environment interaction, *Int J Epidemiol*, **30**, 1035-1040 (2001).

683

684 (10) M. Conroy, J. Sellors, M. Effingham, T.J. Littlejohns, C. Boultonwood, L. Gillions, C.L.M.  
685 Sudlow, R. Collins, N.E. Allen, The advantages of UK Biobank's open access strategy for  
686 health research, *J Intern Med*, **286**, 389-397 (2019).

687

688 (11) S. Cassidy, H. Fuller, J. Chau, M. Catt, A. Bauman, M.I. Trenell, Accelerometer-derived  
689 physical activity in those with cardio-metabolic disease compared to healthy adults: a UK  
690 Biobank study of 52,556 participants. *Acta Diabetologica*, **55**, 975-979 (2018).

691

692 (12) A. Doherty, D. Jackson, N. Hammerla, T. Plötz, P. Olivier, M.H. Granat, T. White, V.T.  
693 van Hees, M.I. Trenell, C.G. Owen, S.J. Preece, R. Gillions, S. Sheard, T. Peakman, S.  
694 Brage, N.J. Wareham, Large Scale Population Assessment of Physical Activity Using Wrist  
695 Worn Accelerometers: The UK Biobank Study. *PloS One*, **12**, e0169649 (2017).

696

697 (13) M. Borga, J. West, J.D. Bell, N.C. Harvey, T. Romu, S.B. Heymsfield, O. Dahlqvist  
698 Leinhard, Advanced body composition assessment: from body mass index to body  
699 composition profiling. *J Invest Med*, **66**, 1-9 (2018).

700

701 (14) Y. Liu, N. Basty, B. Witcher, J.D. Bell, E.P. Sorokin, N. van Bruggen, E.L. Thomas, M.  
702 Cule, Genetic architecture of 11 organ traits derived from abdominal MRI using deep  
703 learning, *eLife*, **10**, e65554 (2021).

704

705 (15) A. McKay, H.R. Wilman, A. Dennis, M. Kelly, M.L. Gyngell, S. Neubauer, J.D. Bell, R.  
706 Banerjee, E.L. Thomas, Measurement of liver iron by magnetic resonance imaging in the UK  
707 Biobank population. *PloS One*, **13**, e0209340 (2018).

708

709 (16) F. Alfaro-Almagro, M. Jenkinson, N.K. Bangerter, J.L.R. Andersson, L. Griffanti, G.  
710 Douaud, S.N. Sotiropoulos, S. Jbabdi, M. Hernandez-Fernandez, E. Vallee, D. Vidaurre, M.  
711 Webster, P. McCarthy, C. Rorden, A. Daducci, D.C. Alexander, H. Zhang, I. Dragonu, P.M.

712 Matthews, K.L. Miller, S.M. Smith, Image processing and Quality Control for the first 10,000  
713 brain imaging datasets from UK Biobank. *NeuroImage*, **166**, 400-424 (2018).

714

715 (17) W. Bai, H. Suzuki, J. Huang, C. Francis, S. Wang, G. Tarroni, F. Guitton, N. Aung, K.  
716 Fung, S.E. Petersen, S.K. Piechnik, S. Neubauer, E. Evangelou, A. Dehghan, D.P. O'Regan,  
717 M.R. Wilkins, Y. Guo, P.M. Matthews, D. Rueckert, A population-based phenome-wide  
718 association study of cardiac and aortic structure and function. *Nature Med*, **26**, 1654-1662  
719 (2020).

720

721 (18) O. Canela-Xandri, K. Rawlik, A. Tenesa, An atlas of genetic associations in UK  
722 Biobank. *Nature Genet*, **50**, 1593-1599 (2018).

723

724 (19) G. McInnes, Y. Tanigawa, C. DeBoever, A. Lavertu, J.E. Olivieri, M. Aguirre, M.A.  
725 Rivas, Global Biobank Engine: enabling genotype-phenotype browsing for biobank summary  
726 statistics. *Bioinformatics*, **35**, 2495-2497 (2019).

727

728 (20) T.G. Richardson, S. Harrison, G. Hemani, G. Davey Smith, An atlas of polygenic risk  
729 score associations to highlight putative causal relationships across the human phenome.  
730 *eLife*, **8**, e43657 (2019).

731

732 (21) D.A. Grimes, K.F. Schulz, Cohort studies: marching towards outcomes. *Lancet*, **359**,  
733 341-345 (2002).

734

735 (22) P.R. Burton, A.L. Hansell, I. Fortier, T.A. Manolio, M.J. Khoury, J. Little, P. Elliott, Size  
736 matters: just how big is BIG?: Quantifying realistic sample size requirements for human  
737 genome epidemiology. *Int J Epidemiol*, **38**, 263-273 (2009).

738

739 (23) S. Lewington, personal correspondence (2022).  
740  
741 (24) A. Fry, T.J. Littlejohns, C. Sudlow, N. Doherty, L. Adamska, T. Sprosen, R. Collins, N.E.  
742 Allen, Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank  
743 Participants with the General Population. *Am J Epidemiol*, **186**, 1026-1034 (2017).  
744  
745 (25) K.J. Rothman, J.E. Gallacher, E.E. Hatch, Why representativeness should be avoided.  
746 *Int J Epidemiol*, **42**, 1012-1014 (2013).  
747  
748 (26) C. Sudlow, J. Gallacher, N. Allen, V. Beral, P. Burton, J. Danesh, P. Downey, P. Elliott,  
749 J. Green, M. Landray, B. Liu, P. Matthews, G. Ong, J. Pell, A. Silman, A. Young, T. Sprosen,  
750 T. Peakman, R. Collins, UK biobank: an open access resource for identifying the causes of a  
751 wide range of complex diseases of middle and old age. *PLoS Med*, **12**, e1001779 (2015).  
752  
753 (27) B.V. Halldorsson, H.P. Eggertsson, K.H.S. Moore, H. Hauswedell, O. Eiriksson, M.O.  
754 Ulfarsson, G. Palsson, M.T. Hardarson, A. Oddsson, B.O. Jansson, S. Kristmundsdottir, B.D.  
755 Sigurpalsdottir, O.A. Stefansson, D. Beyter, G. Holley, V. Tragante, A. Gylfason, P.I. Olason,  
756 F. Zink, M. Asgeirsdottir, S.T. Sverrisson, B. Sigurdsson, S.A. Gudjonsson, G.T. Sigurdsson,  
757 G.H. Halldorsson, G. Sveinbjornsson, K. Norland, U. Styrkarsdottir, D.N. Magnusdottir, S.  
758 Snorraddottir, K. Kristinsson, E. Sobech, H. Jonsson, A.J. Geirsson, I. Olafsson, P. Jonsson,  
759 O.B. Pedersen, C. Erikstrup, S. Brunak, S.R. Ostrowski, G. Thorleifsson, F. Jonsson, P.  
760 Melsted, I. Jonsdottir, T. Rafnar, H. Holm, H. Stefansson, J. Saemundsdottir, D.F.  
761 Gudbjartsson, O.T. Magnusson, G. Masson, U. Thorsteinsdottir, A. Helgason, H. Jonsson,  
762 P. Sulem, K. Stefansson, The sequences of 150,119 genomes in the UK Biobank. *Nature*,  
763 **607**, 732-740 (2022).  
764

765 (28) E. Stamatakis, K.B. Owen, L. Shepherd, B. Drayton, M. Hamer, A.E. Bauman, Is Cohort  
766 Representativeness Passé? Poststratified Associations of Lifestyle Risk Factors with  
767 Mortality in the UK Biobank, *Epidemiol*, **32**, 179-188 (2021).  
768

769 (29) J.R. Emberson, D.A. Bennett, Effect of alcohol on risk of coronary heart disease and  
770 stroke: causality, bias, or a bit of both?, *Vasc Health Risk Manag*, **2**, 239-249 (2006).  
771

772 (30) S. Ebrahim, G. Davey Smith, Commentary: Should we always deliberately be non-  
773 representative?, *Int J Epidemiol*, **42**, 1022-1026 (2013).  
774

775 (31) J. Tyrrell, J. Zheng, R. Beaumont, K. Hinton, T.G. Richardson, A.R. Wood, G. Davey  
776 Smith, T.M. Frayling, K. Tilling, Genetic predictors of participation in optional components of  
777 UK Biobank. *Nature Comms*, **12**, 886 (2021).  
778

779 (32) K. Suissa, M. Hudson, S. Suissa, Glucosamine and lower mortality and cancer  
780 incidence: Selection bias in the observational studies, *Pharmacoepidemiology Drug Saf*, **31**,  
781 1272-1279 (2022).  
782

783 (33) S. Haworth, R. Mitchell, L. Corbin, K.H. Wade, T. Dudding, A. Budu-Aggrey, D.  
784 Carslake, G. Hemani, L. Paternoster, G.D. Smith, N. Davies, D.J. Lawson, J.T. N, Apparent  
785 latent structure within the UK Biobank sample has implications for epidemiological analysis.  
786 *Nature Comms*, **10**, 333 (2019).  
787

788 (34) T.L. Lash, M.P. Fox, R.F. MacLehose, G. Maldonado, L.C. McCandless, S. Greenland,  
789 Good practices for quantitative bias analysis. *Int J Epidemiol*, **43**, 1969-1985 (2014).  
790

791 (35) M.R. Munafò, K. Tilling, A.E. Taylor, D.M. Evans, G. Davey Smith, Collider scope: when  
792 selection bias can substantially influence observed associations. *Int J Epidemiol*, **47**, 226-  
793 235 (2018).

794

795 (36) Emerging Risk Factors Collaboration, Association of Cardiometabolic Multimorbidity  
796 With Mortality, *JAMA*, **314**, 52-60 (2015).

797

798 (37) H.S. Dashti, S.E. Jones, A.R. Wood, J.M. Lane, V.T. van Hees, H. Wang, J.A. Rhodes,  
799 Y. Song, K. Patel, S.G. Anderson, R.N. Beaumont, D.A. Bechtold, J. Bowden, B.E. Cade, M.  
800 Garaulet, S.D. Kyle, M.A. Little, A.S. Loudon, A.I. Luik, F. Scheer, K. Spiegelhalter, J.  
801 Tyrrell, D.J. Gottlieb, H. Tiemeier, D.W. Ray, S.M. Purcell, T.M. Frayling, S. Redline, D.A.  
802 Lawlor, M.K. Rutter, M.N. Weedon, R. Saxena, Genome-wide association study identifies  
803 genetic loci for self-reported habitual sleep duration supported by accelerometer-derived  
804 estimates, *Nature Comms*, **10**, 1100 (2019).

805

806 (38) J. Deelen, D.S. Evans, D.E. Arking, N. Tesi, M. Nygaard, X. Liu, M.K. Wojczynski, M.L.  
807 Biggs, A. van der Spek, G. Atzmon, E.B. Ware, C. Sarnowski, A.V. Smith, I. Seppälä, H.J.  
808 Cordell, J. Dose, N. Amin, A.M. Arnold, K.L. Ayers, N. Barzilai, E.J. Becker, M. Beekman, H.  
809 Blanché, K. Christensen, L. Christiansen, J.C. Collerton, S. Cubaynes, S.R. Cummings, K.  
810 Davies, B. Debrabant, J.F. Deleuze, R. Duncan, J.D. Faul, C. Franceschi, P. Galan, V.  
811 Gudnason, T.B. Harris, M. Huisman, M.A. Hurme, C. Jagger, I. Jansen, M. Jylhä, M.  
812 Kähönen, D. Karasik, S.L.R. Kardia, A. Kingston, T.B.L. Kirkwood, L.J. Launer, T. Lehtimäki,  
813 W. Lieb, L.P. Lytykäinen, C. Martin-Ruiz, J. Min, A. Nebel, A.B. Newman, C. Nie, E.A. Nohr,  
814 E.S. Orwoll, T.T. Perls, M.A. Province, B.M. Psaty, O.T. Raitakari, M.J.T. Reinders, J.M.  
815 Robine, J.I. Rotter, P. Sebastiani, J. Smith, T.I.A. Sørensen, K.D. Taylor, A.G. Uitterlinden,  
816 W. van der Flier, S.J. van der Lee, C.M. van Duijn, D. van Heemst, J.W. Vaupel, D. Weir, K.  
817 Ye, Y. Zeng, W. Zheng, H. Holstege, D.P. Kiel, K.L. Lunetta, P.E. Slagboom, J.M. Murabito,



818 A meta-analysis of genome-wide association studies identifies multiple longevity genes.  
819 *Nature Comm*, **10**, 3669 (2019).

820

821 (39) U. Styrkarsdottir, S.H. Lund, G. Thorleifsson, F. Zink, O.A. Stefansson, J.K. Sigurdsson,  
822 K. Juliusson, K. Bjarnadottir, S. Sigurbjornsdottir, S. Jonsson, K. Norland, L. Stefansdottir, A.  
823 Sigurdsson, G. Sveinbjornsson, A. Oddsson, G. Bjornsdottir, R.L. Gudmundsson, G.H.  
824 Halldorsson, T. Rafnar, I. Jonsdottir, E. Steingrimsson, G.L. Norddahl, G. Masson, P. Sulem,  
825 H. Jonsson, T. Ingvarsson, D.F. Gudbjartsson, U. Thorsteinsdottir, K. Stefansson, Meta-  
826 analysis of Icelandic and UK data sets identifies missense variants in SMO, IL11, COL11A1  
827 and 13 more new loci associated with osteoarthritis. *Nature Genet*, **50**, 1681-1687 (2018).

828

829 (40) P. Akbari, A. Gilani, O. Sosina, J.A. Kosmicki, L. Khrimian, Y.Y. Fang, T. Persaud, V.  
830 Garcia, D. Sun, A. Li, J. Mbatchou, A.E. Locke, C. Benner, N. Verweij, N. Lin, S. Hossain, K.  
831 Agostinucci, J.V. Pascale, E. Dirice, M. Dunn, W.E. Kraus, S.H. Shah, Y.I. Chen, J.I. Rotter,  
832 D.J. Rader, O. Melander, C.D. Still, T. Mirshahi, D.J. Carey, J. Berumen-Campos, P. Kuri-  
833 Morales, J. Alegre-Díaz, J.M. Torres, J.R. Emberson, R. Collins, S. Balasubramanian, A.  
834 Hawes, M. Jones, B. Zambrowicz, A.J. Murphy, C. Paulding, G. Coppola, J.D. Overton, J.G.  
835 Reid, A.R. Shuldiner, M. Cantor, H.M. Kang, G.R. Abecasis, K. Karalis, A.N. Economides, J.  
836 Marchini, G.D. Yancopoulos, M.W. Sleeman, J. Altarejos, G. Della Gatta, R. Tapia-Conyer,  
837 M.L. Schwartzman, A. Baras, M.A.R. Ferreira, L.A. Lotta, Sequencing of 640,000 exomes  
838 identifies GPR75 variants associated with protection from obesity. *Science*, **373**, eabf8683  
839 (2021).

840

841 (41) L. Duncan, H. Shen, B. Gelaye, J. Meijssen, K. Ressler, M. Feldman, R. Peterson, B.  
842 Domingue, Analysis of polygenic risk score usage and performance in diverse human  
843 populations, *Nature Comms*, **10**, 3328 (2019).

844

845 (42) R. Collins, M.K. Balaconis, S. Brunak, Z. Chen, M. De Silva, J.M. Gaziano, G.S.  
846 Ginsburg, P. Jha, P. Kuri, A. Metspalu, N. Mulder, N. Risch, Global priorities for large-scale  
847 biomarker-based prospective cohorts, *Cell Genomics*, **2**, 100141 (2022).

848

849 (43) Z. Chen, M. Smith, H. Du, Y. Guo, R. Clarke, Z. Bian, R. Collins, J. Chen, Y. Qian, X.  
850 Wang, X. Chen, X. Tian, X. Wang, R. Peto, L. Li, Blood pressure in relation to general and  
851 central adiposity among 500 000 adult Chinese men and women. *Int J Epidemiol*, **44**, 1305-  
852 1319 (2015).

853

854 (44) L. Gnatiuc, J. Alegre-Díaz, J. Halsey, W.G. Herrington, M. López-Cervantes, S.  
855 Lewington, R. Collins, R. Tapia-Conyer, R. Peto, J.R. Emberson, P. Kuri-Morales, Adiposity  
856 and Blood Pressure in 110 000 Mexican Adults. *Hypertension*, **69**, 608-614 (2017).

857

858 (45) E. Riboli, R. Kaaks, The EPIC Project: rationale and study design. European  
859 Prospective Investigation into Cancer and Nutrition. *Int J Epidemiol*, **26** Suppl 1, S6-14  
860 (1997).

861

862 (46) Z. Chen, J. Chen, R. Collins, Y. Guo, R. Peto, F. Wu, L. Li, China Kadoorie Biobank of  
863 0.5 million people: survey methods, baseline characteristics and long-term follow-up. *Int J*  
864 *Epidemiol*, **40**, 1652-1666 (2011).

865

866 (47) J.C. Denny, J.L. Rutter, D.B. Goldstein, A. Philippakis, J.W. Smoller, G. Jenkins, E.  
867 Dishman, The "All of Us" Research Program. *New Engl J Med*, **381**, 668-676 (2019).

868

869 (48) K.M. Harrington, X.T. Nguyen, R.J. Song, K. Hannagan, R. Quaden, D.R. Gagnon, K.  
870 Cho, J.E. Deen, S. Muralidhar, T.J. O'Leary, J.M. Gaziano, S.B. Whitbourne, Gender  
871 Differences in Demographic and Health Characteristics of the Million Veteran Program  
872 Cohort. *Womens Health Issues*, **29** Suppl 1, S56-66 (2019).

873

874 (49) T. Chikowore, A.B. Kamiza, O.H. Oduaran, T. Machipisa, S. Fatumo, Non-  
875 communicable diseases pandemic and precision medicine: Is Africa ready? *EBioMedicine*,  
876 **65**, 103260 (2021).

877

878 (50) P. Song, A. Gupta, I.Y. Goon, M. Hasan, S. Mahmood, R. Pradeepa, S. Siddiqui, G.S.  
879 Frost, D. Kusuma, M. Miraldo, F. Sassi, N.J. Wareham, S. Ahmed, R.M. Anjana, S. Brage,  
880 N.G. Forouhi, S. Jha, A. Kasturiratne, P. Katulanda, K.I. Khawaja, M. Loh, M.K. Mridha, A.R.  
881 Wickremasinghe, J.S. Kooner, J.C. Chambers, Data Resource Profile: Understanding the  
882 patterns and determinants of health in South Asians - the South Asia Biobank. *Int J*  
883 *Epidemiol*, **50**, 717-718e (2021).

884

885 (51) T.J. Littlejohns, J. Holliday, L.M. Gibson, S. Garratt, N. Oesingmann, F. Alfaro-Almagro,  
886 J.D. Bell, C. Boultonwood, R. Collins, M.C. Conroy, N. Crabtree, N. Doherty, A.F. Frangi, N.C.  
887 Harvey, P. Leeson, K.L. Miller, S. Neubauer, S.E. Petersen, J. Sellors, S. Sheard, S.M.  
888 Smith, C.L.M. Sudlow, P.M. Matthews, N.E. Allen, The UK Biobank imaging enhancement of  
889 100,000 participants: rationale, data collection, management and future directions. *Nature*  
890 *Comms*, **11**, 2624 (2020).

891

892 (52) P. Elliott, T.C. Peakman, The UK Biobank sample handling and storage protocol for the  
893 collection, processing and archiving of human blood and urine. *Int J Epidemiol*, **37**, 234-244  
894 (2008).

895

896 (53) N.E. Allen, M. Arnold, S. Parish, M. Hill, S. Sheard, H. Callen, D. Fry, S. Moffat, M.  
897 Gordon, S. Welsh, P. Elliott, R. Collins, Approaches to minimising the epidemiological impact  
898 of sources of systematic and random variation that may affect biochemistry assay data in UK  
899 Biobank. *Wellcome Open Res*, **5**, 222 (2021).

900

901 (54) R. Kaaks, E. Riboli, Validation and calibration of dietary intake measurements in the  
902 EPIC project: methodological considerations. *European Prospective Investigation into*  
903 *Cancer and Nutrition, Int J Epidemiol*, **26**, S15-25 (1997).

904

905 (55) M. Pearce, T. Strain, Y. Kim, S.J. Sharp, K. Westgate, K. Wijndaele, T. Gonzales, N.J.  
906 Wareham, S. Brage, Estimating physical activity from self-reported behaviours in large-scale  
907 population studies using network harmonisation: findings from UK Biobank and associations  
908 with disease outcomes. *Int J Behav Nutr Phys Act*, **17**, 40 (2020).

909

910 (56) D. Malden, B. Lacey, J. Emberson, F. Karpe, N. Allen, D. Bennett, S. Lewington, Body  
911 Fat Distribution and Systolic Blood Pressure in 10,000 Adults with Whole-Body Imaging: UK  
912 Biobank and Oxford BioBank. *Obesity*, **27**, 1200-1206 (2019).

913

914 (57) Q. Feng, J.H. Kim, W. Omiyale, J. Bešević, M. Conroy, M. May, Z. Yang, S.Y. Wong,  
915 K.K. Tsoi, N. Allen, B. Lacey, Raw and cooked vegetable consumption and risk of  
916 cardiovascular disease: a study of 400,000 adults in UK Biobank. *Frontiers in Nutr*, **9**,  
917 831470 (2022).

918

919 (58) M.K. Georgakis, R. Malik, D. Gill, N. Franceschini, C.L.M. Sudlow, M. Dichgans,  
920 Interleukin-6 Signaling Effects on Ischemic Stroke and Other Cardiovascular Outcomes: A  
921 Mendelian Randomization Study, *Circ Genom Precis Med*, **13**, e002872 (2020).

922

923 (59) S.C. Larsson, M. Bäck, J.M.B. Rees, A.M. Mason, S. Burgess, Body mass index and  
924 body composition in relation to 14 cardiovascular conditions in UK Biobank: a Mendelian  
925 randomization study. *Europ Heart J*, **41**, 221-226 (2020).

926

927 (60) G. Butler-Laporte, T. Nakanishi, V. Mooser, D.R. Morrison, T. Abdullah, O. Adeleye, N.  
928 Mamlouk, N. Kimchi, Z. Afrasiabi, N. Rezk, A. Giliberti, A. Renieri, Y. Chen, S. Zhou, V.

929 Forgetta, J.B. Richards, Vitamin D and COVID-19 susceptibility and severity in the COVID-  
930 19 Host Genetics Initiative: A Mendelian randomization study. *PLoS Med*, **18**, e1003605  
931 (2021).

932

933 (61) X. Meng, X. Li, M.N. Timofeeva, Y. He, A. Spiliopoulou, W.Q. Wei, A. Gifford, H. Wu, T.  
934 Varley, P. Joshi, J.C. Denny, S.M. Farrington, L. Zgaga, M.G. Dunlop, P. McKeigue, H.  
935 Campbell, E. Theodoratou, Phenome-wide mendelian-randomization study of genetically  
936 determined vitamin D on multiple health outcomes using the UK Biobank study. *Int J*  
937 *Epidemiol*, **48**, 1425-1434 (2019).

938

939 (62) G.D. Smith, Mendelian randomisation and vitamin D: the importance of model  
940 assumptions, *Lancet Diabetes Endocrinol*, **11**, 14 (2023).

941

942 (63) S. Greenland, Multiple-bias modelling for analysis of observational data, *J Royal Stat*  
943 *Soc: Series A*, **168**, 267-306, (2005).

944

945 (64) R. Clarke, M. Shipley, S. Lewington, L. Youngman, R. Collins, M. Marmot, R. Peto,  
946 Underestimation of risk associations due to regression dilution in long-term follow-up of  
947 prospective studies. *Am J Epidemiol*, **150**, 341-353 (1999).

948

949 (65) S. MacMahon, R. Peto, J. Cutler, R. Collins, P. Sorlie, J. Neaton, R. Abbott, J. Godwin,  
950 A. Dyer, J. Stamler, Blood pressure, stroke, and coronary heart disease. Part 1, Prolonged  
951 differences in blood pressure: prospective observational studies corrected for the regression  
952 dilution bias. *Lancet*, **335**, 765-774 (1990).

953

954 (66) A.N. Phillips, G.D. Smith, How independent are "independent" effects? Relative risk  
955 estimation when correlated exposures are measured imprecisely. *J Clin Epidemiol*, **44**,  
956 1223-1231 (1991).

957

958 (67) V. Codd, Q. Wang, E. Allara, C. Musicha, S. Kaptoge, S. Stoma, T. Jiang, S.E. Hamby,  
959 P.S. Braund, V. Bountziouka, C.A. Budgeon, M. Denniff, C. Swinfield, M. Papakonstantinou,  
960 S. Sheth, D.E. Nanus, S.C. Warner, M. Wang, A.V. Khera, J. Eales, W.H. Ouwehand, J.R.  
961 Thompson, E. Di Angelantonio, A.M. Wood, A.S. Butterworth, J.N. Danesh, C.P. Nelson,  
962 N.J. Samani, Polygenic basis and biomedical consequences of telomere length variation.  
963 *Nature Genet*, **53**, 1425-1433 (2021).

964

965 (68) C.E. Rutter, L.A.C. Millard, M.C. Borges, D.A. Lawlor, Exploring regression dilution bias  
966 using repeat measurements of 2858 variables in up to 49,000 UK Biobank participants, *Int J*  
967 *Epidemiol*, **52**, 1545-1556 (2022).

968

969 (69) S. Tin Tin, G.K. Reeves, T.J. Key, Endogenous hormones and risk of invasive breast  
970 cancer in pre- and post-menopausal women: findings from the UK Biobank. *Br J Cancer*,  
971 **125**, 126-134 (2021).

972

973 (70) K.A. Wartolowska, A.J.S. Webb, Midlife blood pressure is associated with the severity of  
974 white matter hyperintensities: analysis of the UK Biobank cohort study. *Europ Heart J*, **42**,  
975 750-757 (2021).

976

977 (71) M.J. Adams, W.D. Hill, D.M. Howard, H.S. Dashti, K.A.S. Davis, A. Campbell, T.K.  
978 Clarke, I.J. Deary, C. Hayward, D. Porteous, M. Hotopf, A.M. McIntosh, Factors associated  
979 with sharing e-mail information and mental health survey participation in large population  
980 cohorts. *Int J Epidemiol*, **49**, 410-421 (2020).

981

982 (72) J. Beller, S. Geyer, J. Epping, Health and study dropout: health aspects differentially  
983 predict attrition. *BMC Med Res Methodol*, **22**, 31 (2022).

984

985 (73) A.E. Taylor, H.J. Jones, H. Sallis, J. Euesden, E. Stergiakouli, N.M. Davies, S. Zammit,  
986 D.A. Lawlor, M.R. Munafò, G. Davey Smith, K. Tilling, Exploring the association of genetic  
987 factors with participation in the Avon Longitudinal Study of Parents and Children. *Int J*  
988 *Epidemiol*, **47**, 1207-1216 (2018).

989

990 (74) G.J. Griffith, T.T. Morris, M.J. Tudball, A. Herbert, G. Mancano, L. Pike, G.C. Sharp, J.  
991 Sterne, T.M. Palmer, G. Davey Smith, K. Tilling, L. Zuccolo, N.M. Davies, G. Hemani,  
992 Collider bias undermines our understanding of COVID-19 disease risk and severity. *Nature*  
993 *Comms*, **11**, 5749 (2020).

994

995 (75) M. Chadeau-Hyam, B. Bodinier, J. Elliott, M.D. Whitaker, I. Tzoulaki, R. Vermeulen, M.  
996 Kelly-Irving, C. Delpierre, P. Elliott, Risk factors for positive and negative COVID-19 tests: a  
997 cautious and in-depth analysis of UK Biobank data, *Int J Epidemiol*, **49**, 1454-1467 (2020).

998

999 (76) L.A.C. Millard, A. Fernández-Sanlés, A.R. Carter, R.A. Hughes, K. Tilling, T.P. Morris,  
1000 D. Major-Smith, G.J. Griffith, G.L. Clayton, E. Kawabata, G. Davey Smith, D.A. Lawlor, M.C.  
1001 Borges, Exploring the impact of selection bias in observational studies of COVID-19: a  
1002 simulation study, *Int J Epidemiol*, **52**, 44-57 (2023).

1003

1004 (77) K.A.S. Davis, J.R.I. Coleman, M. Adams, N. Allen, G. Breen, B. Cullen, C. Dickens, E.  
1005 Fox, N. Graham, J. Holliday, L.M. Howard, A. John, W. Lee, R. McCabe, A. McIntosh, R.  
1006 Pearsall, D.J. Smith, C. Sudlow, J. Ward, S. Zammit, M. Hotopf, Mental health in UK  
1007 Biobank - development, implementation and results from an online questionnaire completed  
1008 by 157,366 participants: a reanalysis. *BJPsych Open*, **6**, e18 (2020).

1009

1010 (78) K. Rannikmäe, K. Ngoh, K. Bush, R. Al-Shahi Salman, F. Doubal, R. Flaig, D.E.  
1011 Henshall, A. Hutchison, J. Nolan, S. Osborne, N. Samarasekera, C. Schnier, W. Whiteley, T.  
1012 Wilkinson, K. Wilson, R. Woodfield, Q. Zhang, N. Allen, C.L.M. Sudlow, Accuracy of

1013 identifying incident stroke cases from linked health care data in UK Biobank. *Neurology*, **95**,  
1014 e697-e707 (2020).

1015

1016 (79) B. Rubbo, N.K. Fitzpatrick, S. Denaxas, M. Daskalopoulou, N. Yu, R.S. Patel, H.  
1017 Hemingway, Use of electronic health records to ascertain, validate and phenotype acute  
1018 myocardial infarction: A systematic review and recommendations. *Int J Cardiol*, **187**, 705-11  
1019 (2015).

1020

1021 (80) T. Wilkinson, C. Schnier, K. Bush, K. Rannikmae, D.E. Henshall, C. Lerpiniere, N.E.  
1022 Allen, R. Flaig, T.C. Russ, D. Bathgate, S. Pal, J.T. O'Brien, C.L.M. Sudlow, Identifying  
1023 dementia outcomes in UK Biobank: a validation study of primary care, hospital admissions  
1024 and mortality data. *Eur J Epidemiol*, **34**, 557-565 (2019).

1025

1026 (81) V. Kuan, S. Denaxas, A. Gonzalez-Izquierdo, K. Direk, O. Bhatti, S. Husain, S. Sutaria,  
1027 M. Hingorani, D. Nitsch, C.A. Parisinos, R.T. Lumbers, R. Mathur, R. Sofat, J.P. Casas,  
1028 I.C.K. Wong, H. Hemingway, A.D. Hingorani, A chronological map of 308 physical and  
1029 mental health conditions from 4 million individuals in the English National Health Service.  
1030 *Lancet*, **1**, e63-e77 (2019).

1031

1032 (82) S. Lewington, R. Clarke, N. Qizilbash, R. Peto, R. Collins, Age-specific relevance of  
1033 usual blood pressure to vascular mortality: a meta-analysis of individual data for one million  
1034 adults in 61 prospective studies. *Lancet*, **360**, 1903-1913 (2002).

1035

1036 (83) S. Lewington, G. Whitlock, R. Clarke, P. Sherliker, J. Emberson, J. Halsey, N. Qizilbash,  
1037 R. Peto, R. Collins, Blood cholesterol and vascular mortality by age, sex, and blood  
1038 pressure: a meta-analysis of individual data from 61 prospective studies with 55,000  
1039 vascular deaths. *Lancet*, **370**, 1829-1839 (2007).

1040



1041 (84) CHARGE and ISGC Consortium, Identification of additional risk loci for stroke and small  
1042 vessel disease: a meta-analysis of genome-wide association studies. *Lancet Neurol*, **15**,  
1043 695-707 (2016).  
1044

1045 (85) W. Luo, L. Gong, X. Chen, R. Gao, B. Peng, Y. Wang, T. Luo, Y. Yang, B. Kang, C.  
1046 Peng, L. Ma, M. Mei, Z. Liu, Q. Li, S. Yang, Z. Wang, J. Hu, Lifestyle and chronic kidney  
1047 disease: a machine learning modeling study, *Frontiers Nutr*, **9**, 918576 (2022).  
1048

1049 (86) R.A. Shah, B. Asatryan, G. Sharaf Dabbagh, N. Aung, M.Y. Khanji, L.R. Lopes, S. van  
1050 Duijvenboden, A. Holmes, D. Muser, A.P. Landstrom, A.M. Lee, P. Arora, C. Semsarian,  
1051 V.K. Somers, A.T. Owens, P.B. Munroe, S.E. Petersen, C.A.A. Chahal, Frequency,  
1052 penetrance, and variable expressivity of dilated cardiomyopathy-associated putative  
1053 pathogenic gene variants in UK Biobank participants. *Circulation*, **146**, 110-124 (2022).  
1054

1055 (87) D. Chahal, D. Sharma, S. Keshavarzi, F.A.Q. Arisar, K. Patel, W. Xu, M. Bhat,  
1056 Distinctive clinical and genetic features of lean vs overweight fatty liver disease using the UK  
1057 Biobank. *Hepatol Int*, **16**, 325-336 (2022).  
1058

1059 (88) N.J. Thomas, S.E. Jones, M.N. Weedon, B.M. Shields, R.A. Oram, A.T. Hattersley,  
1060 Frequency and phenotype of type 1 diabetes in the first six decades of life: a cross-sectional,  
1061 genetically stratified survival analysis from UK Biobank, *Lancet Diabetes Endocrinol*, **6**, 122-  
1062 129 (2018).  
1063

1064 (89) P.R. Burton, A.L. Hansell, UK Biobank: the expected distribution of incident and  
1065 prevalent cases of chronic disease and the statistical power of nested case-control studies,  
1066 *UK Biobank Technical Reports*, (2005).  
1067

1068 (90) K. Papier, A. Knuppel, A. Perez-Cornago, E.L. Watts, T.Y.N. Tong, J.A. Schmidt, N.  
1069 Allen, T.J. Key, R.C. Travis, Circulating insulin-like growth factor-I and risk of 25 common  
1070 conditions: outcome-wide analyses in the UK Biobank study. *Europ J Epidemiol*, **37**, 25-34  
1071 (2022).

1072

1073 (91) S. Floud, R.F. Simpson, A. Balkwill, A. Brown, A. Goodill, J. Gallacher, C. Sudlow, P.  
1074 Harris, A. Hofman, S. Parish, G.K. Reeves, J. Green, R. Peto, V. Beral, Body mass index,  
1075 diet, physical inactivity, and the incidence of dementia in 1 million UK women. *Neurology*, **94**,  
1076 e123-e132 (2020).

1077

1078 (92) T. Strain, K. Wijndaele, S.J. Sharp, P.C. Dempsey, N. Wareham, S. Brage, Impact of  
1079 follow-up time and analytical approaches to account for reverse causality on the association  
1080 between physical activity and health outcomes in UK Biobank, *Int J Epidemiol*, **49**, 162-172  
1081 (2020).

1082

1083 (93) K. Bleicher, R. Summerhayes, S. Baynes, M. Swarbrick, T. Navin Cristina, H. Luc, G.  
1084 Dawson, A. Cowle, X. Dolja-Gore, M. McNamara, Cohort Profile Update: The 45 and Up  
1085 Study. *Int J Epidemiol*, **52**, e92-e101 (2023).

1086

1087 (94) T.J.B. Dummer, P. Awadalla, C. Boileau, C. Craig, I. Fortier, V. Goel, J.M.T. Hicks, S.  
1088 Jacquemont, B.M. Knoppers, N. Le, T. McDonald, J. McLaughlin, A.M. Mes-Masson, A.M.  
1089 Nuyt, L.J. Palmer, L. Parker, M. Purdue, P.J. Robson, J.J. Spinelli, D. Thompson, J. Vena,  
1090 M. Zawati, The Canadian Partnership for Tomorrow Project: a pan-Canadian platform for  
1091 research on chronic disease prevention. *CMAJ*, **190**, e710-717 (2018).

1092

1093 (95) H.S. Feigelson, C.L. Clarke, S.K. Van Den Eeden, S. Weinmann, A.N. Burnett-Hartman,  
1094 S. Rowell, S.G. Scott, L.L. White, M. Ter-Minassian, S.A.A. Honda, D.R. Young, A.  
1095 Kaminen, T. Chinn, A. Lituev, A. Bauck, E.A. McGlynn, The Kaiser Permanente Research

1096 Bank Cancer Cohort: a collaborative resource to improve cancer care and survivorship. *BMC*  
1097 *Cancer*, **22**, 209 (2022).

1098

1099

1100

1101

1102

1103 **Acknowledgements**

1104 We thank Jenny Mills and Alicia Motley for constructing the figures and George Davey Smith  
1105 for helpful suggestions. Additional thanks to the UK Biobank Access team for their tireless  
1106 work on research registrations, applications and output. The authors would like to thank the  
1107 500,000 participants in the UK Biobank study for their enormous generosity and altruism and  
1108 their continued interest, support and involvement.

1109 **Funding:** UK Biobank has core funding from the Medical Research Council, Wellcome,  
1110 British Heart Foundation, Cancer Research UK and National Institute for Health Research.

1111 **Competing interests:** All authors are past or present members of the UK Biobank Strategic  
1112 Oversight Committee or the UK Biobank Senior Team. J.D. serves on scientific advisory  
1113 boards for AstraZeneca and Novartis and consults; British Heart Foundation Centre of  
1114 Research Excellence, University of Cambridge; the National Institute for Health and Care  
1115 Research Blood and Transplant Research Unit in Donor Health and Behaviour, University of  
1116 Cambridge; Health Data Research UK; Wellcome Genome Campus and University of

1117 Cambridge; Department of Human Genetics, Wellcome Sanger Institute, Hinxton, UK. R.C.  
1118 is named on US patent #9957563B2 regarding a statin-related myopathy genetic test but  
1119 any share in royalty and other payments has been waived in favour of the Nuffield  
1120 Department of Population Health, University of Oxford. Financial Relationships: UK Biobank  
1121 Consortium Funding and Enhancements (Novartis, Regeneron Pharmaceuticals, Merck,  
1122 AstraZeneca). R. E. M. is a scientific advisor to Optima Partners and the Epigenetic Clock  
1123 Development Foundation and has received a speaker fee from Illumina. P.M has received  
1124 consultancy or speaker fees from Roche, Merck, Biogen, Rejuveron, Sangamo, Nodthera,  
1125 Novartis and Biogen. P.M. has received research or educational funds from Biogen,  
1126 Novartis, Merck and GlaxoSmithKline.

1127





1130 **Table 1. Cumulative numbers of observed (2020) and predicted incident cases of**  
 1131 **various health conditions**  
 1132

Condition	Year of diagnosis		
	Observed <sup>1</sup>	Predicted	
	2020	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
Alzheimer's disease	5,000	17,000	37,000
Parkinson's disease	4,000	10,000	14,000

1133 <sup>1</sup> Observed values are based on incident events identified from linkage to records of deaths, hospitalisations, cancers, and  
 1134 primary care in the cohort to the end of 2020.

1135 **Table 2. Sampling characteristics of selected general population prospective studies with at least 250,000 participants, containing**  
 1136 **genomic, behavioural and health outcome data<sup>1</sup>**

Study name	Recruitment dates (range)	Location	Sample size	Sex; Age at recruitment	Population from which the sample was recruited	Participation rate
23andMe ( <a href="http://www.23andme.com">www.23andme.com</a> )	2007 - present	Global (but mainly USA)	6.8M	MF; 13+	Customers of a personal genetics company	not known
45 and Up (93)	2006 - 2009	Australia	267,000	MF; 45+	New South Wales residents enrolled in Medicare, recruited through direct invitations	19%
All of Us (47)	2018 - present	USA	Ongoing. Aim: 1M	MF; 18+	Varied approaches, many of which are targeted at underrepresented groups via direct and indirect means	not known
Canadian Partnership for Tomorrow's Health (CanPath) (94)	2008 - present	Canada	330,000	MF; 30-74	Residents across Canada recruited into 7 regional cohorts using varied approaches	not known
China Kadoorie Biobank (46)	2004 - 2008	China	510,000	MF; 30-70	Residents of 10 geographically defined regions across China,	30%



						recruited through direct invitations	
European Prospective Investigation into Cancer, Chronic Diseases, Nutrition and Lifestyle (EPIC) (45)	1992 - 2000	10 European countries	520,000	MF; 35-70	Residents from 23 centres located in 10 European countries	recruited through direct invitations	not known
Kaiser Permanente Research Bank (95)	2007 – 2013	USA	400,000	MF; 18+	Members of Kaiser Permanente health plan recruited through direct invitations, in-person communication and via website.	20-50% of each areas' insured population	
Million Veterans Program (48)	2011 - present	USA	Ongoing. Aim: 1M	MF; 18+	Members of the Veterans Health Administration System recruited through direct invitations and indirect (promotional materials) methods	14%	
UK Biobank (26)	2006 - 2010	UK	500,000	MF: 40-69	Residents living close to 22 assessment centres in the UK, recruited via direct invitations	5.5%	

1137 <sup>1</sup> The IHCC (<https://ihccglobal.org/>) has details of other prospective studies of less than 250,000 participants

1138

1139 **Figure legends**

1140 **Fig. 1. Illustration of the types of data collected in UK Biobank, which includes data**  
1141 **collected at in-person assessments (e.g. lifestyle factors, medical history, blood**  
1142 **pressure and other physical measures, imaging scans), data from online**  
1143 **questionnaires, data generated from biological samples and data from electronic**  
1144 **healthcare records over time**

1145 **Fig. 2. The proportion of incident cases (i.e. ascertained since recruitment into the**  
1146 **study) identified through hospital inpatient admissions, primary care and death data**  
1147 **for some common exemplar conditions (myocardial infarction, diabetes and chronic**  
1148 **obstructive pulmonary disease)**  
1149