## **Plan Overview**

A Data Management Plan created using DMPonline

Title: UK Biobank

Creator: Sander W. van der Laan

**Principal Investigator:** Sander W. van der Laan, Kristel R. van Eijk, Folkert W. Asselbergs, Charlotte N. Onland-Moret, Jessica van Setten , Vinicius Tragante do Ó

**Data Manager:** Sander W. van der Laan, Kristel R. van Eijk, Charlotte N. Onland-Moret, Jessica van Setten , Vinicius Tragante do Ó

**Project Administrator:** Sander W. van der Laan, Kristel R. van Eijk, Folkert W. Asselbergs, Charlotte N. Onland-Moret, Jessica van Setten , Vinicius Tragante do Ó

Affiliation: Other

Funder: European Commission

**Template:** UMC Utrecht DMP

**ORCID iD:** 0000-0001-6888-1404

ORCID iD: 0000-0002-1692-8669

**ORCID iD:** 0000-0002-2360-913X

**ORCID iD:** 0000-0002-4934-7510

ORCID iD: 0000-0002-8223-8957

**ID:** 80411

**Start date:** 01-01-2021

**End date:** 01-01-3000

**Last modified:** 20-04-2022

## **Copyright information:**

The above plan creator(s) have agreed that others may use as much of the text of this plan as they would like in their own plans, and customise it as necessary. You do not need to credit the creator(s) as the source of the language used, but using any of the plan's text does not imply that the creator(s) endorse, or have any relationship to, your project or proposal

# **UK Biobank**

#### 1. General features

## 1.1. Please fill in the table below. When not applicable (yet), please fill in N/A.

We have a project specific number from the UK Biobank team: 24711. The project name is "On the quest for tissue-specific contributions to disease" and the lead investigator is Folkert Asselbergs.

DMP template version	29 (don't change)
ABR number (only for human-related research)	n/a
METC number (only for human-related research)	n/a
DEC number (only for animal-related research)	n/a
Acronym/short study title	UK Biobank
Name Research Folder	smb://ds.umcutrecht.nl/data/LAB/lab_research/RES-Folder-LKCH/UK_Biobank
Name Division	Laboratories, Pharmacy, and Biomedical genetics
Name Department	Central Diagnostic Laboratory
Partner Organization	
Start date study	2021-01-01
Planned end date study	3000-01-01
Name of datamanager consulted*	Saskia Haitjema
Check date by datamanager	January 6th 2022

#### 1.2 Select the specifics that are applicable for your research.

- Fundamental / translational study
- WMO
- Use of Questionnaires
- Observational study
- Retrospective study

We will use data from the UK Biobank (UKB, <a href="https://www.ukbiobank.ac.uk">https://www.ukbiobank.ac.uk</a>) as reference in many studies or as part of a course curriculum in practicals to learn genetic analyses methods, as well as in study projects of researchers in different UMC groups. UK Biobank is a large-scale biomedical database and research resource, containing in-depth genetic and health information from half a million UK participants. The database is regularly augmented with additional data and is globally accessible to approved researchers undertaking vital research into the most common and life-threatening diseases. It is a major contributor to the advancement of modern medicine and treatment and has enabled several scientific discoveries that improve human health.

## 2. Data Collection

#### 2.1 Give a short description of the research data.

Subjects			Data Capture Tool	, .	IFORMAT	Storage space
Human	1	Genotype data		PLINK-format, Oxford- format	.vcf, .bed/.bim/.fam, .gen/.sample	±15Tb
Human		Phenotype data	R, python, etc	flat text file	.tab	±50gb

### 2.2 Do you reuse existing data?

· Yes, please specify

#### Existing data from the UKB:

- Genotype data
- 'Clinical' data, i.e. extensive phenotype information, see below for more information

#### Discover UK Biobank

UK Biobank is a world-leading biomedical database to enable scientific discoveries that improve human health. Our goal is to inspire the imaginations of health researchers around the world to meet the challenge of greater understanding, prevention, and treatment of a range of serious illnesses. UK Biobank has Research Tissue Bank Approval till 2026 and through access to an unmatched amount of biological and medical data on half a million people living in the UK, we can enable your vision of improving human health.

- 1. UK Biobank is a longitudinal study; it follows the health of 500,000 volunteer participants.
- 2. Participants were aged between 40-69 years when they joined UK Biobank between 2006-2010.
- 3. Each participant attended a baseline assessment at a centre in England (89%), Scotland (7%) and Wales (4%).
- 4. Participants provided their consent for long-term follow-up.
- 5. Participants answered lots of questions about their health & lifestyle.
- 6. Participants donated samples of blood, urine and saliva for long-term storage and analysis.
- 7. Physical measurements were also taken (e.g. height, weight, spirometry, blood pressure, heel bone density).
- 8. Many participants have undertaken MR brain & heart imaging, activity monitoring and online follow-up questionnaires.
- 9. We have genetic data on all 500,000 participants.
- 10. UK Biobank is not representative of the general population with evidence of a 'healthy volunteer' selection bias, details of which are published online.

#### 2.3 Describe who will have access to which data during your study.

Please note, that the data have been de-identified for the purpose of public sharing.

,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Who has access	
Pseudonymized data	Research team, Datamanager	

## 2.4 Describe how you will take care of good data quality.

#	Question	Yes	No	N/A
1.	Do you use a certified Data Capture Tool or Electronic Lab Notebook?			х
2.	Have you built in skips and validation checks?			х
3.	Do you perform repeated measurements?			х
4.	Are your devices calibrated?			х
5.	Are your data (partially) checked by others (4 eyes principle)?			х
6.	Are your data fully up to date?	Х		
7.	Do you lock your raw data (frozen dataset)	Х		
8.	Do you keep a logging (audit trail) of all changes?		X	
9.	Do you have a policy for handling missing data?			х
10.	Do you have a policy for handling outliers?	х		

Please note, that the data have been collected and parsed by the UK Biobank. We merely use it as-is.

## 2.5 Specify data management costs and how you plan to cover these costs.

#	Type of costs	<b>Division</b> ("overhead")	Funder	Other (specify)
1.	Archiving	n/a		
2.	Storage			Groups by Onland, Veldink, Pasterkamp, Asselbergs, De Ruijter
3.	Maintenance Dataset			Groups by Onland, Veldink, Pasterkamp, Asselbergs, De Ruijter
4.	Datamanager	х		
5.	Data analysis tool	х		

There is no data-archiving as these data are 'owned' by the UK Biobank, we are only allowed to use it as long as the project runs. The original data remains at the UK Biobank and can be downloaded/obtained from them at any time when needed.

Therefore, there are only costs for data access/right ('Maintenance Dataset'), and the datamanager ('Division') and data analysis tool/storage ('Other').

# 2.6 State how ownership of the data and intellectual property rights (IPR) to the data will be managed, and which agreements will be or are made.

UK Biobank is a large-scale biomedical database and research resource that is enabling new scientific discoveries to be made that improve public health. The resource provides accredited researchers access to medical and genetic data from half a million volunteer participants to improve our understanding of the prevention, diagnosis and treatment of a wide range of serious and life-threatening illnesses.

Through the long-term commitment of 500,000 participants, together with the support of our funders, we are enabling new scientific discoveries that benefit people's health.

Thus the UK Biobank is the main manager of the data and the study collection.

## 3. Personal data (Data Protection Impact Assessment (DPIA) light)

Will you be using personal data (direct or indirect identifying) from the Electronic Patient Dossier (EPD), DNA, body material, images or any other form of personal data?

· Yes, go to next question

These data are obtained from the UK Biobank. They provide access to these data under their informed consent and only to people they deem certified to access these data.

We obtain these data through a secure and encrypted link. The data are pseudonymised by the UK Biobank and are unique to each requested project; in other words the sampleIDs are random for each project number.

For further details please refer to https://www.ukbiobank.ac.uk.

## 3.1 Describe which personal data you are collecting and why you need them.

Which personal data?	Why?	
Clinical and questionnaire data, e.g. age, gender (sex), body measurements, medical history, etc.	To answer the research questions.	
Genetic data	To answer the research questions.	

### 3.2 What legal right do you have to process personal data?

#### · Other, please explain

These data are obtained from the UK Biobank. They provide access to these data under their informed consent and only to people they deem certified to access these data.

We obtain these data through a secure and encrypted link. The data are pseudonymised by the UK Biobank and are unique to each requested project; in other words the sampleIDs are random for each project number.

For further details please refer to https://www.ukbiobank.ac.uk.

#### 3.3 Describe how you manage your data to comply to the rights of study participants.

As noted before: we have full rights to these data on the basis of our project application. We receive periodically a list of samples to be removed for whatever reason (sometimes withdrawn informed consent). This list is send to all project-members. They are subsequently obliged to remove these samples during analyses. Again, this is the way the UK Biobank has setup things, which we follow.

# 3.4 Describe the tools and procedures that you use to ensure that only authorized persons have access to personal data.

First of all you need to be a project-member of our project. People wanting access need to be either a project-member which you will only get after approval by our team (Folkert Asselbergs) and the UK Biobank.

Second we use the data on the HPC. This is only accessible to those with a HPC account. You can only get this after obtaining an UMC account and approval by our team (Folkert Asselbergs).

Thirdly, the data is stored in a separate folder on the HPC with limited access - not all people having access to the HPC can access these data.

Again, this is the way the UK Biobank has setup things, which we follow.

#### 3.5 Describe how you ensure secure transport of personal data and what contracts are in place for doing that.

We do not share these data as this is not allowed. People wanting access need to be either a project-member which you will only get after approval by our team (Folkert Asselbergs) and the UK Biobank.

#### 4. Data Storage and Backup

#### 4.1 Describe where you will store your data and documentation during the research.

For purposes of analyses digital files are partly and temporarily stored on the high-performance computer cluster (HPC) facilitated by the UMC Utrecht.

Data storage is only accessible to authorized personnel.

Phenotype data is accessible based on an application number as given out by the UK Biobank as part of an approved application. The phenotype data should be stored per-group and only made accessible to named in the respective application.

These data are not stored in the RFS associated with this DMP: in the event the data are deleted we can re-download these from the UK Biobank servers.

### 4.2 Describe your backup strategy or the automated backup strategy of your storage locations.

We do not make backups of the data: in the event the data are deleted we can re-download these from the UK Biobank servers.

### 5. Metadata and Documentation

#### 5.1 Describe the metadata that you will collect and which standards you use.

We do not collect anything else, but the data we can obtain through a download.

#### 5.2 Describe your version control and file naming standards.

We will use GitHub as version control with a specific GitHub repository for the each individual project.

We will use the release-system native of GitHub and where possible link it to Zenodo (code only!). Summary statistics of analyses may be shared through Zenodo, DataverseNL or anything else following the requirements of UK Biobank.

## 6. Data Analysis

#### 6 Describe how you will make the data analysis procedure insightful for peers.

We will write analysis plans for each individual project in which we state why we will use which data and which statistical analysis we plan to do in which software. The analysis plan will be stored at GitHub or potentially through a pre-registration server, e.g. OSF. This way this will be findable for our peers. This is the responsibility of each individual project-member.

## 7. Data Preservation and Archiving

#### 7.1 Describe which data and documents are needed to reproduce your findings.

The data is not backed up or stored for any other purpose then analyses. The original data can be obtained from the UK Biobank. We will only store analysis plans/protocol describing the methods and materials, the script to process the data, the scripts leading to tables and figures in the publication, a codebook with explanations on the variable names, and a 'read\_me.txt' file with an overview of files included and their content and use. All at the responsibility of the individual project-members.

These data are stored at the individual RFS' for each research group to which an individual project-member belongs.

#### 7.2 Describe for how long the data and documents needed for reproducibility will be available.

Documentation needed to reproduce findings from this study will be stored for at least 10 years. The original data can be obtained through UK Biobank.

# 7.3 Describe which archive or repository (include the link!) you will use for long-term archiving of your data and whether the repository is certified.

We do not 'own' the data, it is controlled/managed by the <u>UK Biobank</u>. We will only keep copies for local use, and potentially archive projects through Archivemetica and share codes used publications etc through DataverseNL according to the principles of FAIR. However, the original and copy of data will not be stored by the UMC Utrecht, and remains at the UK Biobank.

### 7.4 Give the Persistent Identifier (PID) that you will use as a permanent link to your published dataset.

We will not publish the original data. It can be obtained through UK Biobank as described in previous sections.

#### 8. Data Sharing Statement

# 8.1 Describe what reuse of your research data you intend or foresee, and what audience will be interested in your data.

Specifically the methods and codes developed for the use of this data will be of interest to our peers. Since the data is managed by the <u>UK Biobank</u> we refrain from stating anything regarding data re-use, other than that in general these data make for an excellent population reference for multiple purposes.

# 8.2 Are there any reasons to make part of the data NOT publicly available or to restrict access to the data once made publicly available?

• Yes (please specify)

As the data is privacy-sensitive, and managed by the <u>UK Biobank</u> we will refrain from sharing these data publicly; this should go through UK Biobank.

# 8.3 Describe which metadata will be available with the data and what methods or software tools are needed to reuse the data.

Publications will be open access. The study protocol and this Data Management Plan will also be available. Along with the publication, the codebook of the data and scripts of analyses will be available through GitHub. Data (raw or processed) will be accessible under conditions set forward by the Biobank.

#### 8.4 Describe when and for how long the (meta)data will be available for reuse

• Other (please specify)

Meta data will be accessible under conditions set forward by the UK Biobank.

### 8.5 Describe where you will make your data findable and available to others.

We will publish and archive publication, codes, etc as described above through Archivemetica (local archiving) and DataverseNL (public) with a note that the data will be accessible under conditions set forward by the <u>UK Biobank</u>.

к