

# **PARKINSON'S PROGRESSIVE MARKERS INITIATIVE (PPMI)**

## **Data User Guide – updated February 26<sup>th</sup>, 2024**

### CONTENTS

1. INTRODUCTION .....	3
About this user guide .....	3
What is PPMI? .....	3
What is in the PPMI data? .....	3
What can I use the PPMI data for? .....	5
How do I access the data? .....	5
How is the data maintained? .....	5
2. SOFTWARE TOOLS NEEDED TO WORK WITH THE PPMI DATA .....	6
3. FUNDAMENTALS .....	7
3.1 Participant identifiers .....	7
3.2 Participant cohorts and status .....	8
3.3 Further subgrouping using the Analytic Dataset.....	9
3.4 Selecting data for download .....	9
3.5 Using the data dictionary .....	10
3.6 Getting started .....	11
3.7 Identifying genetic subgroups.....	12
4. STATIC PARTICIPANT DATA.....	14
4.1 Types of static data and how to look up codes.....	14
4.2 Creating a participant master table.....	14
5. STUDY VISITS AND LONGITUDINAL DATA .....	16
6. CLINICAL ASSESSMENTS.....	18
6.1 MDS-UPDRS.....	18
6.2 Non-motor assessment example - UPSIT .....	21
7. MEDICATION .....	23
7.1 Levodopa equivalent medication .....	23
7.2 Other medication .....	25
8. BIOSPECIMENS AND PROTEOMIC DATA .....	26
8.1 Overview of biospecimen and proteomic data.....	26

8.2	Querying biospecimen data .....	27
8.3	Requesting new biospecimen and proteomic data .....	28
8.4	Projects relating to the Amprion alpha-synuclein seeding amplification assay (αS-SAA) 28	
8.5	Further details on other significant projects .....	31
9.	IMAGE DATA.....	33
9.1	Image metadata.....	33
9.2	Viewing and downloading images .....	34
10.	GENETIC DATA .....	35
10.1	Genetic subgroups .....	35
10.2	Tests for known PD genetic risk factors .....	35
10.3	RNA-Sequencing data.....	37
10.4	Other genetic data .....	37
11.	CURATED DATA.....	40
12.	PPMI ONLINE.....	44
13.	PPMI REMOTE.....	45
14.	PPMI FOUND .....	46
15.	DIGITAL SENSOR DATA .....	47
15.1	Roche PD monitoring app: additional information .....	48
15.2	Verily study watch data: additional information .....	49
	APPENDIX A – SUMMARY OF KEY PPMI STUDY DATA .....	51
	APPENDIX B – SQL SCRIPTS.....	63
	APPENDIX C – R SCRIPTS .....	71
	APPENDIX D – IMPORTING DATA INTO A RELATIONAL DATABASE .....	83
	REFERENCES .....	85

# 1. INTRODUCTION

## About this user guide

This user guide provides an introduction to the PPMI data for first-time users as well as serving as a reference for more experienced users. The guide uses practical examples and assumes that most users come from a biomedical background and may have limited knowledge of software tools and coding.

***This version of the user guide is as of February 26<sup>th</sup>, 2024. A revision is in progress and will be posted, as a new version, to the same location. Note that all code/scripts within this guide are provided as is with no guarantee of accuracy or completeness. For questions, reach out to [resources@michaeljfox.org](mailto:resources@michaeljfox.org).***

## What is PPMI?

In 2010, The Michael J. Fox Foundation and a core group of academic scientists and industry partners launched the Parkinson's Progression Markers Initiative (PPMI) toward critically needed biological markers of Parkinson's onset and progression. PPMI has since engaged thousands of partners — more sites; scientists and clinicians; industry experts; philanthropic partners; and, most importantly, study volunteers — to build a cornerstone of Parkinson's research. Analysis from its open-access data set and available biosample library has deepened understanding of disease and informed design of dozens of therapeutic trials. For more background on PPMI, see <https://www.ppmi-info.org/about-ppmi>.

## What is in the PPMI data?

PPMI consists of three main collections of data:

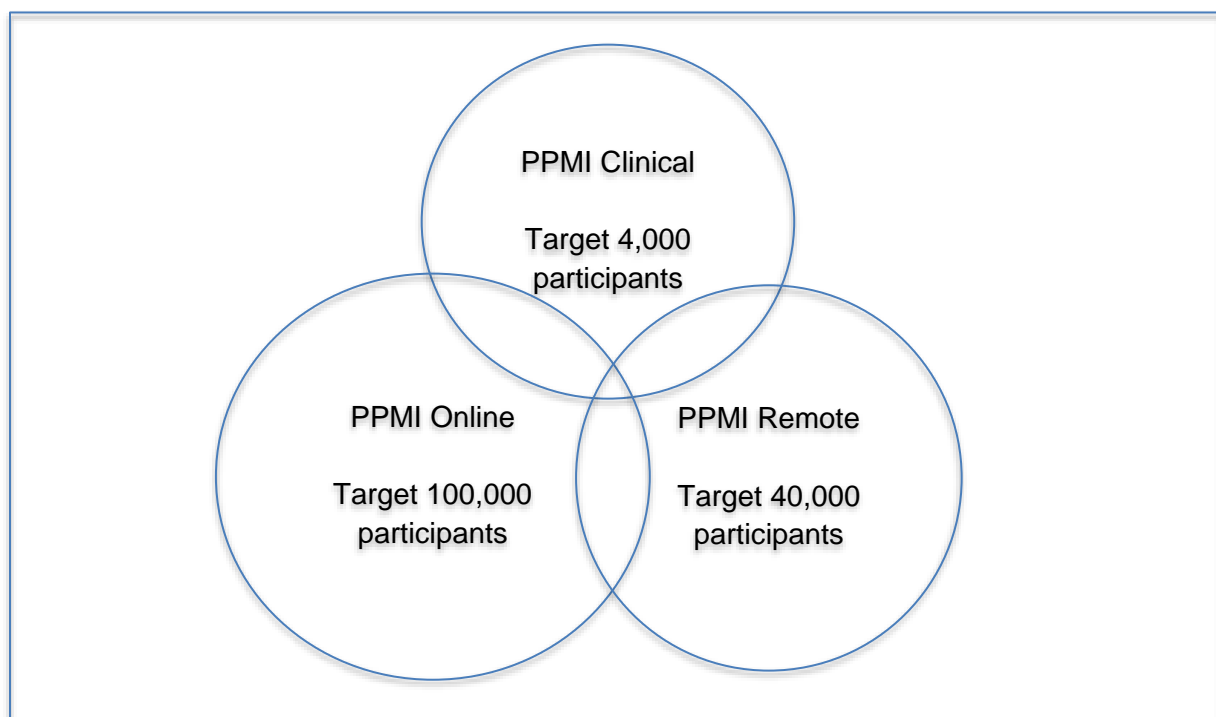
- **PPMI Clinical** contains data that has been captured by in-person clinical assessments, covering people with a confirmed diagnosis of PD, people who exhibit PD risk factors but have not yet been diagnosed with PD, and healthy controls. PPMI Clinical includes static data about the participants, history of clinical presentation using the MDS-UPDRS questionnaire,<sup>1</sup> various non-motor assessments such as the Montreal Cognitive Assessment, biospecimen analysis, genetic test results, MRI scans, DaTscan analysis, medical history, and more. Note that some of these datasets are only available for a subset of the participant population. PPMI is targeted to eventually cover 4,000 participants.

---

<sup>1</sup> Movement Disorder Society Unified Parkinson's Disease Rating Scale. See Goetz *et al.* (2008) and [https://www.movementdisorders.org/MDS-Files1/PDFs/Rating-Scales/MDS-UPDRS\\_English\\_FINAL.pdf](https://www.movementdisorders.org/MDS-Files1/PDFs/Rating-Scales/MDS-UPDRS_English_FINAL.pdf)

- **PPMI Remote** contains data that has been captured by remote assessment. The focus is on gathering data about the pre-diagnostic phase of PD. PPMI Remote collects data on olfactory symptoms and, in some cases, genetic testing, with the intention that this data can be used to identify early indicators of PD. PPMI Remote is targeted to eventually cover 40,000 participants.
- **PPMI Online** contains data that has been self-reported by participants and captured online via a web application. The data includes a subset of the data captured in PPMI Clinical. For instance, in PPMI Online only parts 1 and 2 of the MDS-UPDRS questionnaire are captured whereas in PPM Clinical, all four parts are recorded. PPMI Online contains data for both people with a confirmed diagnosis with PD, and people without a confirmed diagnosis of PD. PPMI Online is targeted to eventually cover 100,000 participants.

Note that individual participants may have data captured by one, two or all three of these methods as illustrated below.



The majority of this guide focuses on the PPMI Clinical Data. PPMI Online and PPMI Remote are discussed in Sections 12 and 13 respectively. There is also a supplementary dataset, PPMI FOUND that contains data captured from follow up telephone consultations for some PPMI Clinical participants; this is covered in Section 14.

Appendix A gives a summary of all the data that is available for immediate download.

## What can I use the PPMI data for?

The PPMI data can be used for a wide variety of research purposes, including studying different ways in which PD presents clinically, how it progresses, how it responds to therapeutic intervention, the effect of genetic factors, as well as studies relating to biomarkers, comorbidities, epidemiology and so on.<sup>2</sup> Any and all uses of the PPMI dataset and biospecimens must conform with the study's Data Use Agreement and the Biospecimen Use Agreement.

## How do I access the data?

You first need to have access approved, after which you can download the data. For further information on these steps, please refer to <https://www.ppmi-info.org/access-data-specimens/download-data>. Note that certain additional data such as raw 'omics, raw digital sensor, and sequestered data are not available from the standard download site but can be requested.

## How is the data maintained?

**PPMI data is updated continuously by the study team. In addition to new data being added, in some cases existing data may be updated, and there is no formal version control.** As such, it is important that you note the date when you took a download of the data. We also recommend developing a repeatable process for downloading and analyzing the data, so that you can easily update your analysis with the latest data if desired in the future.

## A note on examples in this user guide

In this user guide we show a number of example queries and result sets. Where these relate to aggregation of data you will likely get different results because our examples were taken at a point in time and new records may have been subsequently added. Where the examples relate to individual participants, we have used dummy data and you will again get different results.

---

<sup>2</sup> To reference PPMI in a research paper, refer to the [PPMI Publication Policy](#).



## 2. SOFTWARE TOOLS NEEDED TO WORK WITH THE PPMI DATA

The PPMI data is mostly stored as a series of relational tables. When downloaded, most of the data will be in the form of comma separated variable (CSV) files. These are suitable for importing into a number of tools including spreadsheets such as Microsoft **Excel**, statistical packages such as **SPSS**, or Integrated Development Environments (IDEs) such as **RStudio**.<sup>3</sup>

You may find that a single tool, such as R running under RStudio, is sufficient for your needs. A combination of tools can also work well, particularly if you want to do more complex data manipulation. In such cases, you may find it useful to first import the data into a Relational Database Management System (RDBMS) such as **Microsoft SQL Server**,<sup>4</sup> **Oracle**, or **PostgreSQL** where you can use Structured Query Language (SQL) scripts to join tables and manipulate the data efficiently, and then export the results to a tool like R or SPSS for analysis and visualization. Appendix D provides some guidance on importing PPMI data into an RDBMS.

Throughout this user guide we show example data extracts. We have provided scripts for these in both SQL (Appendix B) and R (Appendix C). You may find these to be a useful starting point for your own analysis.

---

<sup>3</sup> RStudio can be downloaded for free from <https://posit.co/download/rstudio-desktop/>

<sup>4</sup> As an example, Microsoft SQL Server Management Studio Express Edition which runs on Windows and requires 16 GB of memory to run effectively can be downloaded for free from <https://www.microsoft.com/en-GB/sql-server/sql-server-downloads>

### 3. FUNDAMENTALS

Before you can start working with the PPMI data you need to identify which participant<sup>5</sup> records to select and which tables to download, and you need to understand a few fundamentals about how the data is organized. In the following, we make the assumption that you will work mostly with PPMI Clinical. For information about PPMI Online and PPMI Remote, refer to Sections 12 and 13 respectively.

#### 3.1 Participant identifiers

Each participant is identified by a unique **PATNO**.<sup>6</sup> This is common across PPMI Clinical, PPMI Remote and PPMI Online. Below is an extract from the **Participant\_Status** table in the PPMI Clinical data which is likely to be your starting point (note: synthetic data for illustrative purposes; some columns not shown for clarity).

PATNO	COHORT	COHORT_DEFINITION	ENROLL_DATE	ENROLL_STATUS	STATUS_DATE	ENROLL_AGE
9001	2	Healthy Control	Jan-11	enrolled	May-21	64.0
9002	1	Parkinson's Disease	Feb-11	enrolled	Sep-21	47.1
9003	1	Parkinson's Disease	Mar-11	enrolled	Sep-21	67.5
9004	1	Parkinson's Disease	Apr-11	enrolled	Jan-22	55.2
9005	2	Healthy Control	Apr-11	enrolled	Jan-22	59.6
9006	1	Parkinson's Disease	Apr-11	Excluded	Apr-11	
9007	1	Parkinson's Disease	May-11	Withdrew	Oct-12	77.5
9008	1	Parkinson's Disease	May-11	Withdrew	Jun-12	61.2
9009	2	Healthy Control	Jun-11	enrolled	Apr-21	81.0
9010	2	Healthy Control	Jun-11	enrolled	May-21	79.7

**PATNO is an important field that you will use to link all data relating to an individual participant.**

<sup>5</sup> We use the term 'participant' in this guide, but you will sometimes see the term 'patient' in the underlying data. The two terms can be considered as synonymous.

<sup>6</sup> To protect participant confidentiality, the actual participant names or similar identifiers are not available to PPMI users.

## 3.2 Participant cohorts and status

PPMI Clinical participants have been enrolled into one of five **COHORTs**, namely:

1. **Parkinson's Disease**, i.e., people who have a formal diagnosis of Parkinson's disease (PD)
2. **Prodromal**, i.e., people who are at risk of developing PD based on clinical features, genetic variants or other biomarkers but have not been formally diagnosed
3. **Healthy Controls**, i.e., people with no neurologic disorder and no first-degree relative with PD
4. **SWEDD** (Scan without dopaminergic deficit). This is a small legacy cohort that you may wish to exclude, depending on your research purpose; for more details, see <https://www.ppmi-info.org/study-design/study-cohorts#legacy/>
5. **Early Imaging** is a cohort of participants with a confirmed diagnosis of PD who were untreated and underwent additional tests including DaTscan and AV-133 imaging.<sup>7</sup> For most analyses, this cohort can be combined with the Parkinson's Disease cohort, though data available for this cohort differs slightly from the Parkinson's Disease cohort.

PPMI Online has only two cohorts, Parkinson's and Non-Parkinson's, with the latter including prodromal cases and healthy controls.

Which cohorts you need will depend on the nature of your study. For further information on cohorts in general, see <https://www.ppmi-info.org/study-design/study-cohorts>. Note that cohorts for PPMI Clinical are further sub-divided as described in the next section.

There is also an **ENROLL\_STATUS** field which can take the following values relating to the status of the participant in the study:

- **Pending** – Is awaiting initial screening and, as such, does not yet have all the initial study data captured
- **Screened** – Has been screened for inclusion in the PPMI study but has not yet had all the initial study data captured
- **Declined** – After screening, the participant declined to give consent and no further study data was captured
- **Screen\_failed** and **Excluded** – After screening, the participant was deemed not suitable for the PPMI study and hence no further study data was captured
- **Enrolled** – Enrolled in the PPMI study and has data from one or more study visits recorded
- **Withdrew** – Participant withdrew from the study but had completed one or more study visits
- **Complete** – Participant completed all of their planned study visits and elected to not re-enroll in further study visits

---

<sup>7</sup> Details of the data captured for early imaging can be found here: <https://blackfynn-ppmi-schema.readthedocs-hosted.com/en/stable-documentation-branch/index.html>



Given the above definitions, you are most likely to want to include participants with a status of **Enrolled, Withdrew or Complete**.

### 3.3 Further subgrouping using the Analytic Dataset

In addition to the cohorts and status groups detailed above, **it is important to note the cohorts are also further organized into subgroups, described in a spreadsheet called the Consensus Committee Analytic Dataset**, which can be downloaded from the Quick Start section of the download page. In particular, this sub-divides participants with a PD diagnosis into those who have sporadic PD and those carrying disease-correlated variants of the LRRK2, GBA, and SNCA genes. Prodromal participants are similarly sub-divided into genetic groups and also into a general risk/hyposmia group (HPSM) and an REM behavior disorder (RBD) group. This data is also available in the Participant\_Status table which can be regarded as the latest source for this information.

### 3.4 Selecting data for download

As a starting point, tables you are likely to always need to download from the download page (<https://ida.loni.usc.edu/pages/access/studyData.jsp>) are:

- Study Docs
  - Data & Databases
    - Code List (Annotated)
    - Data Dictionary (Annotated)
- Subject Characteristics
  - Patient Status
    - Participant Status
  - Subject Demographics
    - Demographics
- Motor Assessments
  - Motor / MDS-UPDRS
    - MDS-UPDRS Part I Patient Questionnaire
    - MDS-UPDRS Part I Non-motor Aspects of Daily Living
    - MDS-UPDRS Part II Patient Questionnaire: Motor Aspects of Daily Living
    - MDS-UPDRS Part III Treatment Determination and Part III: Motor Examination

To download these, select the “Study Data” option from the Download drop down menu. Appendix A will help you determine which additional data files to download for your study.

In general, it is not recommended to use the download ALL option at the bottom of the list as this will take up a lot of storage, though you may find it helpful to use the ALL option for some of the data subgroups.

### 3.5 Using the data dictionary

The PPMI Clinical data dictionary is very useful for understanding both the content and format of tables. This can be downloaded from the main download page (see previous section) or directly from this link: <https://www.ppmi-info.org/access-data-specimens/data-dictionary>. There is a separate data dictionary for PPMI Online, also available on the main download page. The following notes will help you use these documents effectively:

1. In PPMI Clinical, there are two versions of the data dictionary, a version with the name suffix “\_Annotated”, and another with the suffix “\_Harmonized”. The two are identical except that the “\_Annotated” one provides some additional comments, hence in general, we would recommend using the “\_Annotated” version.
2. The field **MOD\_NAME** provides the table name as it exists in the PPMI database used by the data management team responsible for maintaining the data but when you download the file itself to your directory, it will generate a different filename. To see the mapping between the two, open the data dictionary in Excel and filter on **ITM\_NAME** = (Blanks). The **DSCR** column will then show the download filename. For example, the MOD\_NAME “AE” maps to the physical table name “Adverse\_Event\_Log” and “CLCKDRAW” maps to “Clock\_Drawing”.

	A	B	C	D
1	MOD_NAME	ITM_NAME	PAG_NAME	DSCR
2	AE		AE	Adverse Event Log
38	AVCNSNT		AVCNSNT	Early Imaging Documentation of Informed Consent
56	AVELIG		AVELIG	Early Imaging Eligibility
76	AVIMAG		AVIMAG	Early Imaging AV-133 Imaging
102	AVPREGNANC		AVPREGN	Early Imaging Pregnancy Test
119	AVREPPREG		AVREPPRE	Early Imaging Report of Pregnancy
131	AVSCRNF		AVSCRNF	Early Imaging Screen Fail
143	AVTEL		AVTEL	Early Imaging Adverse Event Telephone Assessment
218	CLCKDRAW		CLCKDRAW	Clock Drawing

3. The field **PAG\_NAME**, meaning “Page Name” is a generally used identifier indicating the version of electronic Case Reference Form (eCRF) used. It is typically the same as the MOD\_NAME for data collected via the electronic data capture system used for the clinical protocol.
4. **ITM\_NAME** refers to an individual column in a table.
5. **ITM\_TYPE** gives the datatype of a column (one of CHAR, DATE, NUMBER, TEXT or TIME). **FLD\_LEN** gives the size of the field, which may include decimal places given by **DECML**.

There are also codebooks for PPMI Clinical and PPMI Online, and these are covered in Section 4.

## 3.6 Getting started

Once you've set up your chosen tool and imported the data, you can quickly generate basic statistical information, for example the following SQL queries return the number of participants in each cohort and each status respectively.<sup>8,9,10</sup>

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Number of participants by cohort */
SELECT COHORT_DEFINITION, COUNT(*) AS 'Record Count'
FROM PPMI.dbo.Participant_Status
GROUP BY COHORT_DEFINITION
ORDER BY COUNT(*) DESC
```

COHORT_DEFINITION	Record Count
Parkinson's Disease	1097
Prodromal	731
Healthy Control	297
SWEDD	81
Early Imaging (original study participants only)	19

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Number of participants by status */
SELECT ENROLL_STATUS, COUNT(*) AS 'Record Count'
FROM PPMI.dbo.Participant_Status
GROUP BY ENROLL_STATUS
ORDER BY COUNT(*) DESC
```

ENROLL_STATUS	Record Count
enrolled	1357
Withdrew	307
Excluded	245

<sup>8</sup> In these examples, we assume you loaded the data into a database called 'PPMI' under user 'dbo'.

<sup>9</sup> As noted in the Introduction, you may get different results when executing these and other example queries in this user guide.

<sup>10</sup> Note that the files normally download with a date appended to the filename, e.g., Participant\_Status\_24Jul2023.csv, so in each case, you will need to either take this into account when importing the data or modify the script accordingly

ENROLL_STATUS	Record Count
Declined	114
Complete	102
screen_failed	53
screened	45
pending	2

To achieve the same as the above in RStudio, you could execute the following R script: <sup>11,12,13,14</sup>

```
# This script is provided as is with no guarantees of completeness or accuracy
# Participant counts by cohort and enroll status. NOTE: Requires tidyverse
package to be installed.
library (readr)
library (dplyr)

setwd ("C:\\PPMI")
Participant_Status <- read_csv ("Participant_Status.csv")

group_by (Participant_Status, COHORT_DEFINITION) %>% summarize(RECORD_COUNT =
n()) %>% arrange(desc(RECORD_COUNT))
group_by (Participant_Status, tolower(ENROLL_STATUS)) %>% summarize(RECORD_COUNT
= n()) %>% arrange(desc(RECORD_COUNT))
```

### 3.7 Identifying genetic subgroups

To get a count of participants in the Parkinson's Disease and Prodromal cohorts by genetic subgroup you can execute code (see Script 1 in appendices) that filters on ENROLL\_STATUS and COHORT, giving the following results.

COHORT_DEFINITION	Genetic subgroup	Participant count
Parkinson's Disease	SRDC	452
Parkinson's Disease	NULL	199
Parkinson's Disease	LRRK2	172
Parkinson's Disease	GBA	88
Parkinson's Disease	SNCA	29
Parkinson's Disease	PARKIN	8
Parkinson's Disease	PINK1	1
Parkinson's Disease	Multiple factors	1

<sup>11</sup> Change the C:\\PPMI to the location of your downloaded files

<sup>12</sup> Note that the files normally download with a date appended to the filename, e.g., Participant\_Status\_24Jul2023.csv, so in each case, you will need to either rename the file or modify the script accordingly

<sup>13</sup> The ENROLL\_STATUS has some statuses capitalized and some not capitalized, for example "Withdraw" and "withdrew", hence the use of the tolower() function in the R script to ignore the case

<sup>14</sup> Hereafter, please refer to Appendix B for SQL scripts and Appendix C for corresponding R scripts.



COHORT_DEFINITION	Genetic subgroup	Participant count
Prodromal	LRRK2	230
Prodromal	GBA	189
Prodromal	HPSM	115
Prodromal	RBD	106
Prodromal	Multiple factors	16
Prodromal	SNCA	9
Prodromal	PINK1	1
Prodromal	NULL	1

Note that we have introduced a 'Multiple factors' option as participants can sometimes test positive for both LRRK2 and GBA risk factors. You should consider whether you want to include these participants in your analysis.

You may also want to consider excluding participants with a NULL value for the subgroup as the reason for their allocation to the cohort may not be clear.

## 4. STATIC PARTICIPANT DATA

### 4.1 Types of static data and how to look up codes

Static data (i.e., data that does not change over time) about participants can be found in several tables, including the following in PPMI Clinical:

- **Demographics**
- **Family\_History**

Most fields in these tables use codes; you are likely to want to use the decode value at some point. For instance, SEX on the Demographics table is stored as a 0 for Female and 1 for Male. The decode values can all be found in the PPMI code book table which can be downloaded from <https://www.ppmi-info.org/access-data-specimens/data-dictionary>. In the following examples we refer to this table as the Codes table. There is a similar code book for PPMI Online also available in the same download site.

To determine the number of males and females in the PPMI data who have a status of enrolled, withdrew or complete, and either a confirmed PD diagnosis or are prodromal, you could use Script 2 in the appendices. This queries the Demographics table to retrieve SEX, joins this with the Codes table to provide the decoded value (Female or Male) and also joins with the Participant\_Status table using PATNO to filter on ENROLL\_STATUS and COHORT\_DEFINITION. It then performs a “group by” to aggregate the data.

COHORT_DEFINITION	SEX	Participant Count
Parkinson's Disease	Female	351
Parkinson's Disease	Male	536
Prodromal	Female	292
Prodromal	Male	271

### 4.2 Creating a participant master table

You may find it useful to create a Participant\_Master table that contains all the static participant data that you need for your analysis in one place by joining multiple tables using PATNO. You can then use this as a starting point to define the participant population for your study. Script 3 in the appendices joins the Demographics, Participant\_Status and Codes tables and selects a population of all participants who are in either the Parkinson's or Prodromal cohorts and who have a status of Enrolled, Withdrew or Complete. The result is a table like the following.



PATNO	BIRTHDT	COHORT_ DEFINITION	Genetic subgroup	ENROLL _AGE	ENROLL _DATE	ENROLL _STATUS	SEX	HANDED	PD diagnosis date
9001	Jan-56	Parkinson's Disease	SRDC	64.0	Mar-11	enrolled	Female	Right	03/01/20 10
9002	Aug-73	Parkinson's Disease	GBA	47.1	Mar-11	enrolled	Female	Right	02/01/20 10
9003	Aug-54	Parkinson's Disease	SRDC	67.5	Apr-11	enrolled	Female	Right	11/01/20 09
9004	Oct-63	Parkinson's Disease	LRRK2	55.2	May-11	Withdrew	Male	Right	09/01/20 10
9005	Nov-56	Parkinson's Disease	LRRK2	59.6	May-11	Withdrew	Male	Left	02/01/20 11
9006	Jun-54	Parkinson's Disease	SRDC	60.1	Jun-11	enrolled	Male	Mixed	03/01/20 11
9007	Jan-43	Parkinson's Disease	SRDC	77.5	Nov-11	Complete	Female	Right	01/01/20 11
9008	Jun-49	Parkinson's Disease	SRDC	61.2	Dec-11	Withdrew	Male	Right	03/01/20 11
9009	Sep-35	Parkinson's Disease	SRDC	81.9	Apr-12	enrolled	Male	Left	07/01/20 12
9010	May-38	Parkinson's Disease	NULL	79.7	Apr-12	Enrolled	Female	Right	08/01/20 12

(Data shown is synthetic for illustrative purposes).

This is just an example, and many more fields could be added to this.

## 5. STUDY VISITS AND LONGITUDINAL DATA

The PPMI data contains many measurements that change over time, for example, the results of MDS-UPDRS questionnaires. Once recorded this data does not change but it is important to note that additional measurements and other data may be added to participants' records in the future.

To illustrate this longitudinal feature of the data here is an extract of the MDS-UPDRS\_PART\_I table for participant #9001 (note: synthetic data for illustrative purposes; some columns not shown for clarity).

PATNO	EVENT_ID	INFODT	NP1COG	NP1HALL	NP1DPRS	NP1ANXS	NP1APAT	NP1DDS	NP1RTOT
9001	SC	01/2012	0	0	1	0	0	0	1
9001	BL	03/2012	0	0	0	0	0	0	0
9001	V01	05/2012	0	0	0	0	0	0	0
9001	V02	08/2012	0	0	1	0	0	0	1
9001	V03	11/2012	0	0	1	0	0	0	1
9001	V04	03/2013	0	0	1	1	0	0	2
9001	V05	09/2013	0	0	0	1	0	0	1
9001	V06	05/2014	0	0	0	0	0	0	0
9001	V07	09/2014	0	0	0	0	0	0	0
9001	V08	04/2015	0	0	0	1	0	0	1
9001	V09	11/2015	0	0	1	0	0	0	1
9001	V10	04/2016	0	0	0	1	0	0	1
9001	V11	06/2016	1	0	0	1	0	0	2
9001	V12	09/2016	1	0	1	1	0	0	3
9001	V14	03/2017	1	0	1	1	0	0	3
9001	V15	03/2018	2	0	1	1	1	1	6
9001	V17	09/2020	2	1	2	1	1	0	7
9001	R17	11/2021	2	1	2	2	2	0	9

The **EVENT\_ID** is the identifier of the PPMI visit at which the individual measurement was taken and the **INFODT** indicates the date on which the measurement was taken.<sup>15</sup> **EVENT\_ID is an important field that, in combination with PATNO, you will used to link test and measurement data relating to the same point in time for an individual participant.** *Please note, due to changes in the PPMI data collection infrastructure, as of the date of this version of the user guide, the recommended best practices for use of the EVENT\_ID for longitudinal analysis is under review. Once this review concludes, a new version of this user guide reflecting these changes will be posted to the same location.*

The columns to the right show the actual measurement values, in this example the scores for MDS-UPDRS Part 1 questions like cognitive issues (NP1COG), hallucinations (NP1HALL) etc. The full list of column definitions can be found in the data dictionary described in section 3.4.

<sup>15</sup> Dates in PPMI are rounded to protect participant privacy.



Note that there are some special values of EVENT\_ID. **BL** represents the baseline measurement when the participant formally entered the PPMI study. In some cases, this will have been preceded by **SC**, a screening measurement to assess suitability to join the study and capture static data such as demographics, typically 60 calendar days before the baseline. Note that, for confirmed cases of Parkinson's, the date of diagnosis will typically be different from either of these events.

The **V** numbers represent the subsequent study visits on which tests were performed and measurements recorded. The **R** numbers represent where measurements were recorded remotely. There are also **U** numbers which relate to unscheduled visits. The following document gives the planned schedule of tests performed at these points in time, though actual schedules and measurements recorded for individual participants may differ from this:

[https://www.ppmi-info.org/sites/default/files/docs/PPMI-2.0-AM-2-Protocol\\_SCHEDULE-OF-ACTIVITIES.pdf](https://www.ppmi-info.org/sites/default/files/docs/PPMI-2.0-AM-2-Protocol_SCHEDULE-OF-ACTIVITIES.pdf)

## 6. CLINICAL ASSESSMENTS

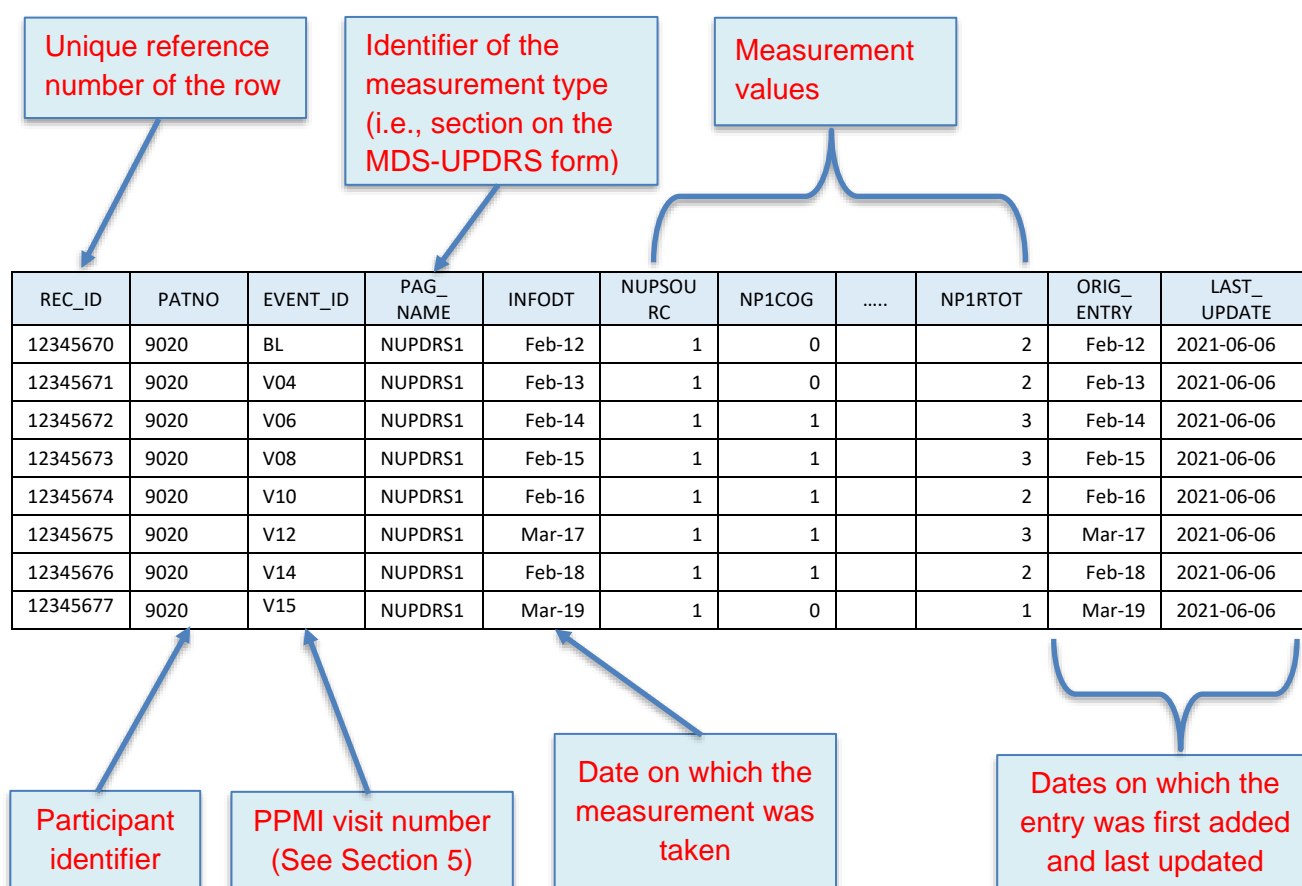
PPMI Clinical has captured a wide range of clinical assessment data, relating to both motor and non-motor symptoms of PD. Here we step through two examples. Other clinical assessments follow similar patterns.

### 6.1 MDS-UPDRS

Clinical measurements using the MDS-UPDRS questionnaire have been recorded at regular intervals for most of the PPMI participants and are stored in the following tables that mirror the different sections of the form, some of which are completed by a medical professional and some of which are completed by the participant.

- MDS-UPDRS\_Part\_I
- MDS-UPDRS\_Part\_I\_Patient\_Questionnaire
- MDS\_UPDRS\_Part\_II\_\_Patient\_Questionnaire
- MDS-UPDRS\_Part\_III
- MDS-UPDRS\_Part\_III\_ON\_OFF\_Determination\_\_\_Dosing
- MDS-UPDRS\_Part\_IV\_\_Motor\_Complications

These tables, and tables associated with other clinical measurements, all follow a similar structure:



The diagram illustrates the structure of the MDS-UPDRS tables. Callouts point to specific columns in the table below:

- Unique reference number of the row**: Points to the **REC\_ID** column.
- Identifier of the measurement type (i.e., section on the MDS-UPDRS form)**: Points to the **PAG\_NAME** column.
- Measurement values**: Points to the **NP1COG** and **NP1RTOT** columns.
- Participant identifier**: Points to the **PATNO** column.
- PPMI visit number (See Section 5)**: Points to the **EVENT\_ID** column.
- Date on which the measurement was taken**: Points to the **INFODT** column.
- Dates on which the entry was first added and last updated**: Points to the **ORIG\_ENTRY** and **LAST\_UPDATE** columns.

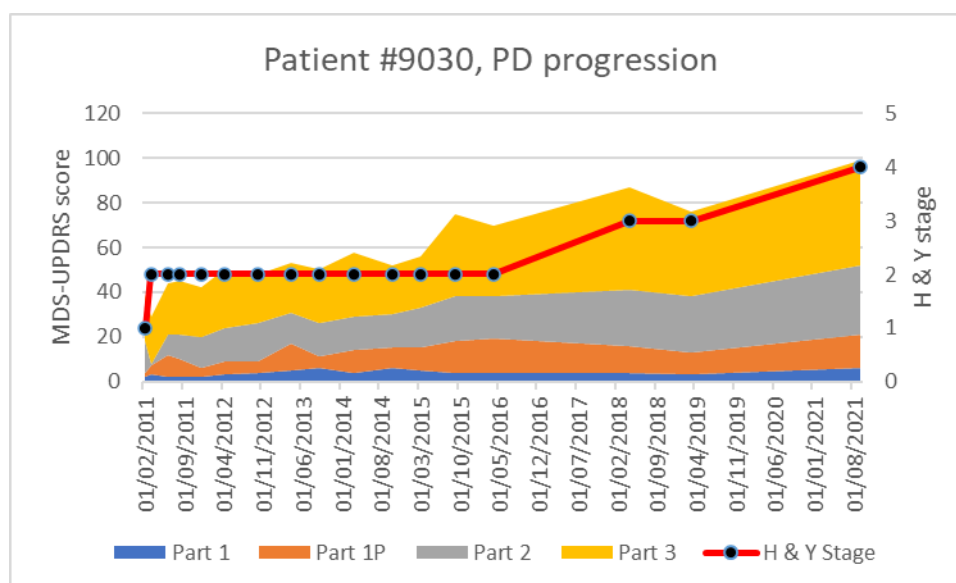
REC_ID	PATNO	EVENT_ID	PAG_NAME	INFODT	NUPSOURC	NP1COG	.....	NP1RTOT	ORIG_ENTRY	LAST_UPDATE
12345670	9020	BL	NUPDRS1	Feb-12	1	0		2	Feb-12	2021-06-06
12345671	9020	V04	NUPDRS1	Feb-13	1	0		2	Feb-13	2021-06-06
12345672	9020	V06	NUPDRS1	Feb-14	1	1		3	Feb-14	2021-06-06
12345673	9020	V08	NUPDRS1	Feb-15	1	1		3	Feb-15	2021-06-06
12345674	9020	V10	NUPDRS1	Feb-16	1	1		2	Feb-16	2021-06-06
12345675	9020	V12	NUPDRS1	Mar-17	1	1		3	Mar-17	2021-06-06
12345676	9020	V14	NUPDRS1	Feb-18	1	1		2	Feb-18	2021-06-06
12345677	9020	V15	NUPDRS1	Mar-19	1	0		1	Mar-19	2021-06-06

Script 4 in the appendices extracts the progression of MDS-UPDRS scores and also the Hoehn & Yahr stage<sup>16</sup> over time for a participant by taking the total scores from each section of the questionnaire.

INFODT	EVENT_ID	Part 1	Part 1P	Part 2	Part 3	H & Y Stage
2011-02-01	SC	2	2	14	21	1
2011-03-01	BL	3	4	1	21	2
2011-06-01	V01	2	10	9	23	2
2011-08-01	ST	2	8	11	24	2
2011-12-01	V03	2	4	14	22	2
2012-04-01	V04	3	6	15	26	2
2012-10-01	V05	4	5	17	22	2
2013-04-01	V06	5	12	14	22	2
2013-09-01	V07	6	5	15	24	2
2014-03-01	V08	4	10	15	29	2
2014-10-01	V09	6	9	15	22	2
2015-03-01	V10	5	10	18	23	2
2015-09-01	V11	4	14	20	37	2
2016-04-01	V12	4	15	19	32	2
2018-04-01	V14	4	12	25	46	3
2019-03-01	V15	3	10	25	38	3
2021-09-01	V17	6	15	31	47	4

(Note: synthetic data for illustrative purposes)

It is then straightforward to export this data to another tool for analysis and visualization. For example, the image below was plotted in Excel from the above data.



<sup>16</sup> This is the original rather than the modified scale. See Hoehn & Yahr (1967)

Note the following:

1. Occasionally there can be **missing MDS-UPDRS measurements**. These will normally appear as NULL values. In the above example visits with missing values were excluded by adding checks for nulls in the join statements. Alternative approaches are to set NULLs to a default value like zero (not recommended) or to impute a value such as a mean. Users should consider what missing value treatment is most appropriate based on their specific research question(s).
2. Sometimes there are **multiple MDS-UPDRS measurements on the same date**. These are because measurements may be taken when the participant is experiencing a medication “on” period and when the participant is experiencing an “off” period. In the above example we used the MAX function to take the maximum score on a given day, thereby ensuring just one value (the “worst case”) per date. Alternatives are either to take all values or filter on either “on” or “off” values. To filter on “on”/“off”, use column **PDSTATE** on MDS\_UPDRS\_Part\_III though note that in a small number of cases this field is missing data (NULL value), in which case you need instead query the PAG\_NAME field which will have a value of NUPDR3OF or NUPDR3ON<sup>17</sup>. See also the document “Methods for Defining PD Med Use” available from the download page under the folder Study Docs/Study Methods.
3. The **MDS-UPDRS scores can fluctuate over time**, i.e., they can go down as well as up, although the long-term trend is for them to increase. This can be caused by a number of factors including changes to levodopa and other medication, subjectivity of measurements and normal fluctuation of MDS-UPDRS scores, for example resulting from “on” and “off” periods.

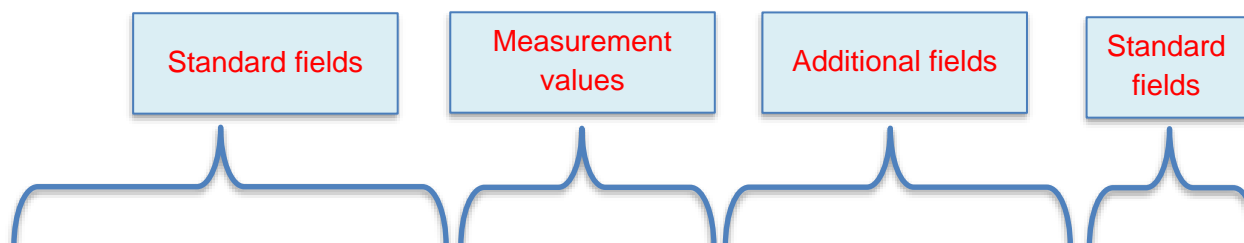
---

<sup>17</sup> In the future all new records will have PAG\_NAME = ‘NUPDRDOSE3’. This version of the user guide is as of April 25<sup>th</sup>, 2023. A revision is in progress and will be posted, as a new version, to the same location.

## 6.2 Non-motor assessment example - UPSIT

PPMI has captured a lot of data relating to clinical assessments other than MDS-UPDRS, but it is important to note that these are not always recorded at every study visit or for every participant.

These assessments are all stored in a similar way but have individual differences. Take the example of the file `University_of_Pennsylvania_Smell_Identification_Test__UPSIT_` as shown below (note: synthetic data for illustrative purposes).



REC_ID	PATNO	EVENT_ID	PAG_NAME	INFODT	SCENT_01_CORRECT	.....	TOTAL_CORRECT	UPSIT_FORM	UPSIT_PRCNTGE	UPSIT_PRCTVER	UPSIT_SOURCE	ORIG_ENTRY	LAST_UPDATE
12345681	9050	V04	UPSITCLINIC	Mar-21	0		32		55.5	Dec-20	Clinical Site	Mar-21	17/01/2022
12345682	9060	V04	UPSITCLINIC	Apr-21	0		38	2	59	Dec-20	Clinical Site	Apr-21	25/08/2021
12345683	9070	V04	UPSITCLINIC	Feb-21	1		35		35	Dec-20	Clinical Site	Feb-21	30/04/2021
12345684	9080	V04	UPSITCLINIC	Jun-21	0			2		Dec-20	Clinical Site	Jun-21	03/12/2021
12345685	9090	V04	UPSITCLINIC	Jul-21	1		34		40	Dec-20	Clinical Site	Jul-21	20/05/2021

As shown above there are standard fields, namely:

- REC\_ID: Uniquely identifies the individual row on the table
- PATNO: The unique identifier of the participant
- EVENT\_ID and INFODT: Capture longitudinal information about the data
- PAG\_NAME: Takes the value UPSITCLINIC for tests performed under PPMI Clinical, or UPSITPRO for tests performed under PPMI Remote (see also section 13)
- ORIG\_ENTRY: The date the record was added to the PPMI database
- LAST\_UPDATE: The date the record was last updated (will often be the same as ORIG\_ENTRY)

You won't normally need to be concerned with REC\_ID, ORIG\_ENTRY and LAST\_UPDATE.

The table also contains the actual measurement values and may also contain additional informational columns specific to the test. To understand these, open the data dictionary spreadsheet in Excel and filter on MOD\_NAME (in this case "UPSIT" – see also section 3.4). This will reveal, for example that:

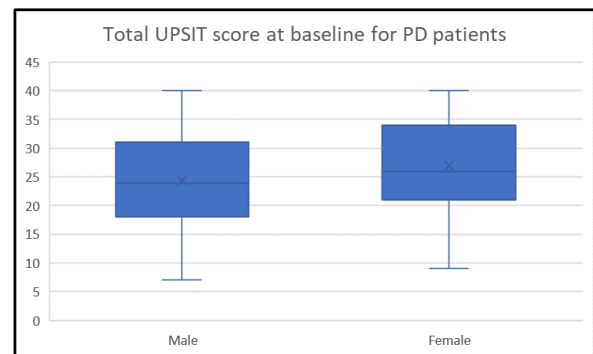
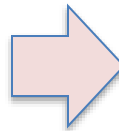
- UPSITFORM is the version of the UPSIT form administered
- UPSIT\_SOURCE is the source of the study where the results come from (which, in this case, can also be determined from the PAG\_NAME field described above)

Because these fields use codes, we can further open the Codes\_List spreadsheet in Excel and again filter on MOD\_NAME = "UPSIT", then filter on either the column name (ITM\_NAME in the spreadsheet) and/or a specific code value (CODE in the spreadsheet) to understand what the data means. This tells us the values that may appear in these fields:

- UPSITFORM: 1 = Original, 2 = Revised
- UPSIT\_SOURCE: Clinical site or Remote

Suppose we want to look at the relationship between performance in the UPSIT test and sex for participants with a diagnosis of Parkinson's disease. Script 5 in the appendices joins several tables to prepare the data, which can then be used in an application like Excel or R for analysis and visualization (data shown is synthetic for illustrative purposes).

PATNO	SEX	TOTAL_CORRECT
9041	Male	24
9042	Male	33
9043	Male	24
9044	Female	28
9045	Female	31
9046	Female	31
9047	Female	26
...	...	...



Details of the methods used to capture UPSIT data, and many other tests can be found in PDF documents which can be downloaded from the main download site. In this case, there is a document PPMI\_UPSIT\_Methods\_Document\_18NOV2016.pdf which can be found under the group 'Non-motor Assessments' and subgroup 'Olfactory Tests'.

## 7. MEDICATION

### 7.1 Levodopa equivalent medication

PPMI Clinical contains information about medications taken by the participant, both dopaminergic drugs such as levodopa, and other medications.

Dopaminergic medications, including levodopa, dopamine agonists and monoamine oxidase inhibitors, are stored as levodopa equivalent daily dosage (LEDD)<sup>18</sup> and this can be found on the **LEDD\_Concomitant\_Medication\_Log** table. Here is a sample of this table for participant #9052 generated by Script 6 in the appendices (note: synthetic data for illustrative purposes; some columns omitted for clarity):

PATNO	EVENT_ID	PAG_NAME	INFODT	LEDTRT	STARTDT	STOPDT	LEDD
9052	LOG	CMED	NULL	AMANTADINE	01/2012	01/2012	100
9052	LOG	CMED	NULL	AMANTADINE	01/2012	02/2013	200
9052	LOG	CMED	NULL	AMANTADINE	03/2013	02/2015	300
9052	LOG	CMED	NULL	CARBIDOPA/LEV ODOPA 25/100	02/2014	08/2014	400
9052	LOG	CMED	NULL	CARBIDOPA/LEV ODOPA 25/100	09/2014	09/2015	450
9052	LOG	CMED	NULL	PRAMIPEXOLE	11/2014	06/2020	25
9052	LOG	CMED	NULL	CARBIDOPA/LEV ODOPA 25/100	12/2014	08/2015	50
9052	LOG	CMED	NULL	CARBIDOPA/LEV ODOPA 25/100	09/2015	12/2017	750
9052	LOG	CMED	NULL	CARBIDOPA/LEV ODOPA ER 50/200	02/2016	06/2020	150
9052	V17	LEDDLOG	01/06/2022	SINEMET	12/2016	06/2020	25
9052	LOG	CMED	NULL	CARBIDOPA/LEV ODOPA 25/100	01/2018	06/2020	200
9052	LOG	CMED	NULL	RYTARY	01/2018	06/2020	600
9052	LOG	CMED	NULL	ENTACAPONE	04/2018	06/2020	LD x 0.33
9052	V17	LEDDLOG	01/06/2022	SINEMET	07/2020		1500
9052	V17	LEDDLOG	01/06/2022	SELEGILINE	02/2022		100

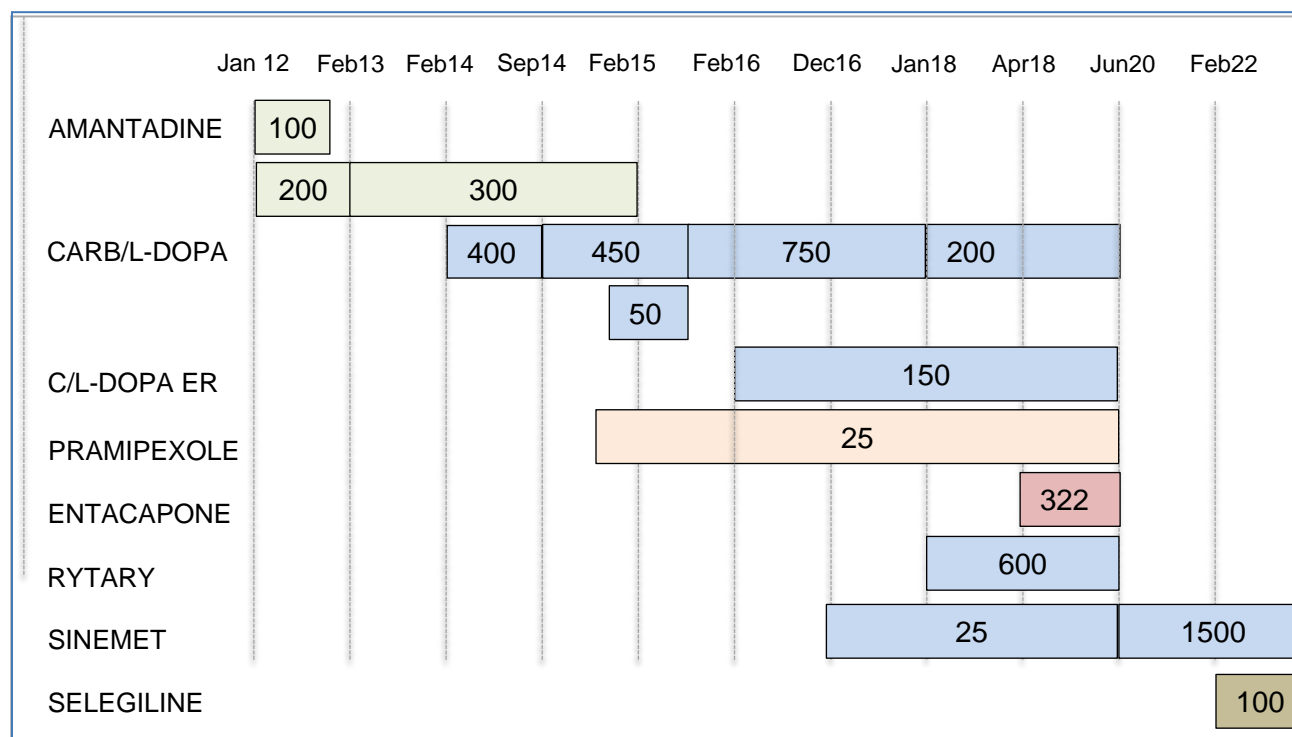
Each row represents a drug, dosage, and period of time when it was in effect. Where there is no end date, this represents the current prescription. Note that several drugs may be taken in parallel and that the table does not automatically come in date order.

Note also that there are two different values for PAG\_NAME: data collected prior to 2020 is labeled **CMED** and data captured from 2020 onwards is labeled **LEDDLOG**. The difference is that the later records additionally capture information about the visit (EVENT\_ID).<sup>19</sup>

<sup>18</sup> See Tomlinson (2010) for an explanation of how LEDD calculations are made. There is also documentation available from the PPMI download page under Study Docs/ALL Study Methods.

<sup>19</sup> In the old system EVENT\_ID = 'SC' or 'Log' and in the new system EVENT\_ID = 'ED'

A visual representation is useful for understanding what is stored in this table (note: dates *not* shown to scale):



The figures represent milligrams of levodopa equivalent. Certain values need to be calculated, for instance, Entacapone is expressed as  $0.33 \times \text{the levodopa dosage}$ , i.e.,  $0.33 \times (200 + 150 + 600 + 25) = 975 \times 0.33 = 321.75^{20}$ . Levodopa (which can also appear in the table under brand names such as Rytary, Madopar and Sinemet), is shown in blue in the above chart.

Note that the start and end dates are only measured to the month so we make an assumption that the medication is active for the full month. In this example, Amantadine with a LEDD value of 100 is active from January 2012 to January 2012 and we interpret this as the whole month of January 2012.

This table can be a little tricky to work with if you simply want to know at a point in time, what is the total LEDD value across all medications. Script 7 in the appendices creates a new table, LEDD, that contains this information:

If we query this new table for participant #9052 used as an example previously, we now have the following, which may be easier to work with as it simply records each time the LEDD value changes. The last value shows the current medication, which could also go to zero.

<sup>20</sup> In a very small number of cases of the drug Stalevo, which is a combination of levodopa and entacapone, there is not enough information in the LEDD log table to determine the correct LEDD value



PATNO	STARTDT	LDOPA	LEDD
9052	2012-01-01	0	300
9052	2012-02-01	0	200
9052	2013-03-01	0	300
9052	2014-02-01	400	700
9052	2014-09-01	450	750
9052	2014-11-01	450	775
9052	2014-12-01	500	825
9052	2015-04-01	500	525
9052	2015-09-01	750	775
9052	2016-02-01	900	925
9052	2016-12-01	925	950
9052	2018-01-01	975	1000
9052	2018-04-01	975	1321.75
9052	2020-07-01	1500	1500
9052	2022-02-01	1500	1600

Note that all prodromal and some Parkinson's participants (who have elected to delay the start of dopaminergic treatment) will not have LEDD records.

## 7.2 Other medication

As well as levodopa and equivalent medication, PPMI has captured data about other medications taken by the participant in the table **Concomitant\_Medication\_Log**. This has a similar structure to the LEDD table, using the same STARTDT and STOPDT fields to record the time period, but unlike the LEDD table each medication is recorded independently of other medications and there is no conversion to a common metric. One thing that makes this table more complicated is that for the dosage there are two fields: a numeric field, CMDOSE that records the dose, and a text field, CMDOSU, that specifies the unit of measurement, for example MG for milligrams.

## 8. BIOSPECIMENS AND PROTEOMIC DATA

### 8.1 Overview of biospecimen and proteomic data

PPMI contains data relating to many different biospecimens and proteomic analyses, including plasma, serum, blood, cerebrospinal fluid, DNA and RNA. This data is available under the Biospecimen folder on the main download site. There is general data such as blood test results and CSF test results available under the subfolder “Lab Collection Procedures” and there is data relating to specific research in the following subfolders:

- Biospecimen Analysis
- Metabolomic Analysis
- Proteomic Analysis

A list of research projects (both completed and in progress) can be found on the main PPMI site (<https://www.ppmi-info.org/access-data-specimens/ongoing-analysis/specimen-analysis>). This table contains information about the principal investigator, organization, current status, cohorts used and visit schedule. The methods for these projects can be found after registration in a series of PDF documents under the corresponding Methods subfolder.

For instance, in the link given above, we can look up project 181 which is “Adaptive Immune Markers for Predicting Cognitive Decline in PD”. We can then find details of the study methods in the document “181 Project Methods: Adaptive Immune Markers for Predicting Cognitive Decline” under the Biospecimen Methods subfolder, and the actual data in the file “Project 181 Adaptive Immune Markers for Predicting Cognitive Decline in PD” under the Biospecimen Sample Analysis subfolder. The data uses PATNO and EVENT\_ID as the key fields to join to the rest of the PPMI Clinical data in the usual way. Note that the results for many projects have been collated into the file “Current\_Biospecimen\_Analysis\_Results” in the Biospecimen analysis subfolder of the Biospecimen folder; filter this file on the Project ID that you are interested in.

There is also a table of neuropathology results in the subfolder of the same name that details the results of post-mortem analyses of PPMI participants. This gives information about the presence of Lewy Bodies, Braak Staging, hippocampal degeneration, substantia nigra depigmentation, the level of TDP43 and so on. This uses the field PATNO and records are labelled with an EVENT\_ID of AUT.

## 8.2 Querying biospecimen data

As a simple example of how to query this data, take the **Blood\_Chemistry\_\_Hematology** table, available under Lab Collection Procedures subfolder of the Biospecimen section on the main download site. This contains many blood test results, so to identify which results are captured for the majority of participants we can run the Script 8 in the appendices to obtain the following results.

LTSTCODE	LTSTNAME	Participant Count
RCT4	ALT (SGPT)	2194
RCT5	AST (SGOT)	2194
HMT12	Basophils	2194
HMT19	Basophils (%)	2194
RCT183	Calcium (EDTA)	2194
HMT11	Eosinophils	2194
HMT18	Eosinophils (%)	2194
HMT2	Hematocrit	2194
HMT40	Hemoglobin	2194
HMT9	Lymphocytes	2194
HMT16	Lymphocytes (%)	2194
HMT10	Monocytes	2194
HMT17	Monocytes (%)	2194
HMT8	Neutrophils	2194
HMT15	Neutrophils (%)	2194
HMT13	Platelets	2194
HMT3	RBC	2194
HMT71	RBC Morphology	2194
RCT17	Serum Bicarbonate	2194
RCT18	Serum Chloride	2194
RCT11	Serum Glucose	2194
RCT16	Serum Potassium	2194
RCT15	Serum Sodium	2194
RCT8	Serum Uric Acid	2194
RCT1	Total Bilirubin	2194
RCT12	Total Protein	2194

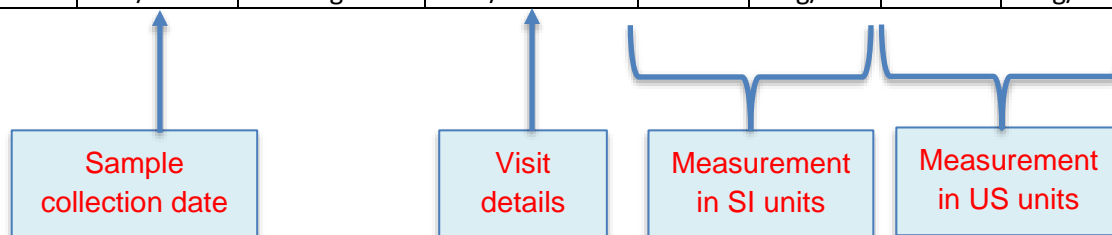
LTSTCODE	LTSTNAME	Participant Count
RCT6	Urea Nitrogen	2194
HMT7	WBC	2194
CGT283	Prothrombin Time	2186
RCT13	Albumin-QT	1631
RCT1407	Alkaline Phosphatase-QT	1631
RCT392	Creatinine (Rate Blanked)	1631
CGT626	APTT-QT	845
CGT284	APTT-QT	782
RCT3088	Albumin-BCG	576
RCT4447	Alkaline Phosphatase-QT	576
RCT408	Creatinine(Rate Blanked)-2dp	576
HMT4	MCV	576
CGT766	APTT-FSL	563
AMT7453	Which visit being performed?	82
BAT318	Serum beta hCG, Qualitative-QT	77
RCT3580	B-hCG, Qualitative	22
RCT3579	B-hCG, Quantitative	22
HMT95	Atypical Lymphocytes	9
HMT96	Atypical Lymphocytes	9
HMT20	Bands	6
HMT21	Bands (%)	6
AMT2917	Requisition Received ?	6
(Results truncated)		

This shows that there are 29 blood tests that have been captured for the majority of the participant population and a number of blood tests that have only partial coverage. It can be useful to understand data coverage in this way before embarking on detailed analysis.

We can then retrieve results for individual blood tests, for example Hemoglobin tests for participant #9002 using Script 9 in the appendices (note: example uses synthetic data for illustrative purposes):

EVENT_ID	LCOLLDT	LTSTNAME	LVISTYPE	LSIRES	LSIUNIT	LUSRES	LUSUNIT
SC	05/2011	Hemoglobin	Screening	145	g/L	14.5	g/dL
V04	06/2012	Hemoglobin	V04/Month 12	146	g/L	14.6	g/dL
V06	07/2013	Hemoglobin	V06/Month 24	150	g/L	15.0	g/dL

EVENT_ID	LCOLLDT	LTSTNAME	LVISTYPE	LSIRES	LSIUNIT	LUSRES	LUSUNIT
V08	09/2014	Hemoglobin	V08/Month 36	148	g/L	14.8	g/dL
V10	06/2015	Hemoglobin	V10/Month 48	153	g/L	15.3	g/dL
V12	11/2016	Hemoglobin	V12/Month 60	153	g/L	15.3	g/dL
V14	08/2018	Hemoglobin	V14/Month 84	142	g/L	14.2	g/dL



Several fields have been excluded from this query. In general, to look up the definition of a field, search on the field name in the Data Dictionary spreadsheet.

### 8.3 Requesting new biospecimen and proteomic data

If a particular biospecimen or proteomic analysis does not already exist in PPMI, a user can request biosamples and conduct a new analysis to derive a variable of interest see link here: <https://www.ppmi-info.org/access-data-specimens/request-specimens>. The resulting data is returned to the PPMI database for the benefit of all PPMI users; each approved biospecimen requested is given a project ID that is used to track the related data. Associated methods are documented in PDF files which can also be downloaded from the main download page.

### 8.4 Projects relating to the Amprion alpha-synuclein seeding amplification assay (αS-SAA)

The alpha-synuclein seed amplification assay, or SAA for short, is a recently developed diagnostic test for Parkinson's disease. Studies have shown that it has both high sensitivity and high specificity in identifying people with the condition and forms of Neuronal Synuclein Disease (NSD) (Siderowf *et al.* 2023). It is also being investigated for other uses, such as for differential diagnosis of dementia with Lewy Bodies and Multiple System Atrophy. A substantial subset of the PPMI participants have undergone this test and the results can be found in the **SAA Biospecimen Analysis Results** file, which is available for download from the Biospecimen / Biospecimen Analysis folder.

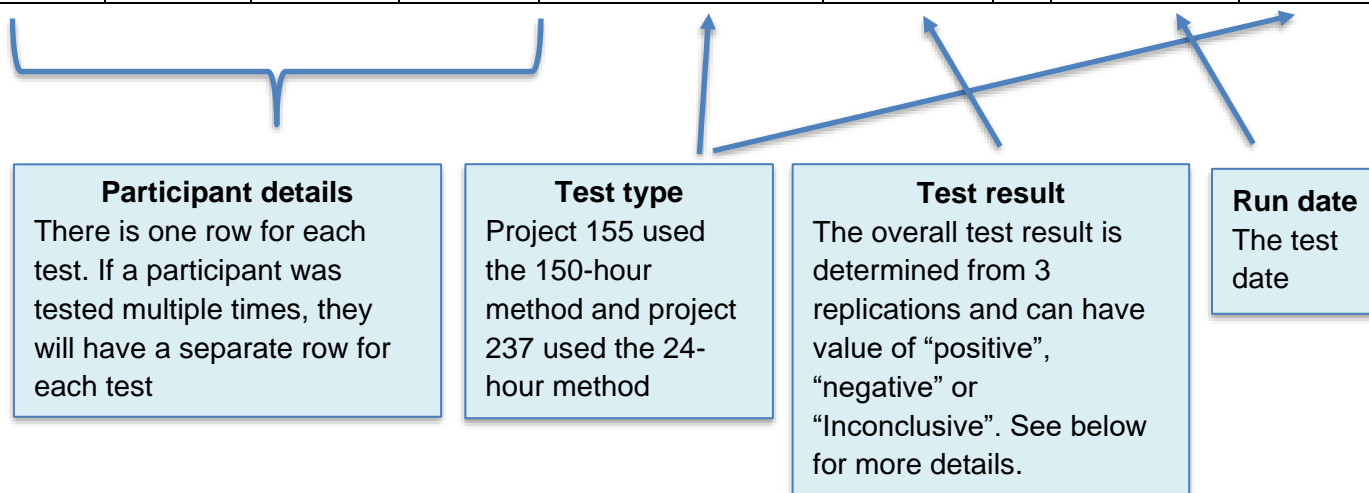
This file contains SAA results from two projects:

- Project 155 used the **Ampiron-αS-SAA** method where the assay process took 150 hours.

- Project 237 used the **Ampiron-24h- $\alpha$ S-SAA** method where the volume of recombinant  $\alpha$ -synuclein was reduced and the assay length shortened to 24 hours.

The file contains the following general columns (*note: some columns omitted for clarity*):

PATNO	SEX	COHORT	CLINICAL_EVENT	SAAMethod	SAA_Status	...	RUNDATE	PROJECT
90072	Male	PD	BL	Amprion-Alpha-synuclein-SAA	Positive		1/12/2023	155
90072	Male	PD	V03	Amprion-24h-alpha-synuclein-SAA	Positive		3/12/2023	237
90073	Female	Control	BL	Amprion-Alpha-synuclein-SAA	Negative		2/12/2023	155
90073	Female	Control	V02	Amprion-24h-alpha-synuclein-SAA	Negative		3/12/2023	237
90074	Male	PD	BL	Amprion-24h-alpha-synuclein-SAA	Positive		1/12/2023	237
90075	Female	PD	BL	Amprion-24h-alpha-synuclein-SAA	Positive		1/12/2023	237



Each SAA test comprises three replications. For each replication, there is a series of columns that capture the test details (i.e. each column will appear 3 times, suffixed with one of Rep1, Rep2 or Rep3). Note that many of the columns are specific to one of the two projects as detailed in the table below. Fluorescence is measured in relative fluorescence units (RFU) and there are thresholds for this value of 5,000 for project 155 and 3,000 for project 237 that determine the overall result of the test.

Kinetic Parameter	Column Name (suffixed with Rep1, Rep2 or Rep3)	Project 155	Project 237
Maximum Fluorescence in RFU after 24 hours	Fmax_24h		✓
Time to reach threshold in 24 hours	TTT_24h		✓
Area under the curve [RFU*hours]	AUC_24h		✓



Kinetic Parameter	Column Name (suffixed with Rep1, Rep2 or Rep3)	Project 155	Project 237
Time to reach the maximum slope value in hours	Tsmax_24h		✓
The maximum slope value obtained for a given assay reaction in RFU/h	SLOPEMax_24h		✓
Maximum Fluorescence in RFU after 150 hours	Fmax_150h	✓	
Time to reach threshold in 150 hours	TTT_150h	✓	
Area under the curve [RFU*hours]	AUC_150h	✓	
Time to reach 50% of FMax in hours	T50_150h	✓	
Slope in RFU/h	SLOPERep_150h	✓	
Identifier of the instrument used in the test (2, 4 or 6)	Instrument	✓	✓
Sample volume used in the test (takes the value NEV1 = 20-30µL, NEV2 = 30-40µL or NEV3 = <20µL)	SampleVol	✓	✓

See the corresponding methods documents for the individual projects for more details. These can be found in the Biospecimen Analysis Methods folder.

The overall test result is determined according to a set of decision criteria. For project 155, these are:

SAA_Status	Decision Criteria
Positive	<b>All 3</b> replicates have Fmax > <b>5,000</b> RFU
Negative	<b>0 or 1</b> replicate has Fmax > <b>5,000</b> RFU
Inconclusive	<b>Exactly 2</b> replicates have Fmax > <b>5,000</b> RFU

For project 237, the criteria are more complex:

SAA_Status	Decision Criteria
Positive	<b>All 3</b> replicates have Fmax > <b>3,000</b> RFU
	<b>Exactly 2</b> replicates have Fmax ≥ <b>3,000</b> RFU & < <b>45,000</b> RFU
	<b>Exactly 2</b> replicates have Fmax ≥ <b>45,000</b> RFU & <b>1</b> replicate has Fmax ≥ <b>3,000</b> RFU
Negative	<b>0 or 1</b> replicate has Fmax ≥ <b>3,000</b> RFU
Inconclusive	<b>Exactly 2</b> replicates have Fmax ≥ <b>45,000</b> RFU & <b>1</b> replicate has Fmax < <b>3,000</b> RFU

Where a participant has multiple test results, we recommend prioritizing the project 237 result with the most recent RUNDATE.

Project 237 also generates **Synucleinopathy Status** results. Based on the value of Fmax, a determination is made as to whether the case is “LBD-like” (Type 1), “MSA-like” (Type 2) or Undetermined. This is captured in the column SAA\_Type, according to the following criteria:

SAA_Type	Decision Criteria
Type 1	Positive samples where all 3 replicates have a Fmax =>45,000 RFU

SAA_Type	Decision Criteria
Type 2	Positive samples where 2 or 3 replicates have a Fmax <45,000 RFU
Undetermined	Positive sample has exactly 2 replicates with a Fmax =>45,000 RFU
NA	Not applicable, i.e. the SAA test result was Negative

Note that **these tables only contain data for Parkinson's and Healthy Control participants**; data for the prodromal participants is sequestered and accessible by submitting a written request (word document or PDF) to [ppmi@michaeljfox.org](mailto:ppmi@michaeljfox.org) that details who will be using the data and what analysis they will be performing in accordance with the [PPMI Data Access Guidelines](#).

## 8.5 Further details on other significant projects

In this section we give an overview of some key projects that have captured additional biospecimen, metabolomic and/or proteomic data.

**Project 151 (Identification of Proteins & Protein Networks in Human CSF That Differentiate Within PD Participants)** has conducted **cerebrospinal fluid tests** on Parkinson's, Prodromal and Healthy Control participants. The results are split across several files prefixed with "Project\_151". The test reference numbers (column TESTNAME) relate to the SOMA\_SEQ\_ID generated by the SomaScan software package<sup>21</sup> and can be looked up in column B of the spreadsheet "PPMI\_Project\_151\_pqtl\_Analysis\_Annotations\_20210210".

**Project 180 (Metabolomic Analysis of Penetrance, Prognosis, & Tracking Biomarkers of LRRK2 PD)** collected **plasma samples** from Parkinson's, Prodromal and Healthy Control participants and analyzed them using liquid chromatography with mass spectrometry (LC/MS) for a variety of metabolite and lipids including purines, lipids associated with lysosomal function like sphingolipids, bis(monoacyl)glycerophosphate (BMP), pro/antioxidant and pro/anti-inflammatory molecules, and environmental / dietary exposure markers. The results are split across 5 files prefixed with "Metabolomic\_Analysis\_of\_LRRK2\_PD". The values in the column TESTNAME take the form MZxxx.xx\_RTyyy.yy\_mmm where MZxxx.xx indicates a mass to charge ratio of xxx.xx, RTyyy.yy is the retention time in seconds, and mmm indicates either positive or negative ionization mode. For further information, please refer to the method document, "PPMI\_Methods\_Project\_180\_Metabolomics\_LRRK2\_20220705".

**Project 190 (Targeted & Untargeted MS Based Proteomics of Urine in PD)** analyzed **proteome profiles** of **urine samples** using mass spectrometry. The results are split across 5 files prefixed with "Targeted\_and\_untargeted\_MS-based\_proteomics\_of\_urine\_in\_PD". The test identifiers (column TESTNAME) are a concatenation of the gene identifier and the protein identifier. These can both be looked up in a protein database such as UniProt (<https://www.uniprot.org/>) or the Human Protein Atlas (<https://www.proteinatlas.org/>)

<sup>21</sup> <https://somalogic.com/>

**Projects 196 and 222 (Targeted Proteomics in AMP PD)** analyzed **cerebrospinal fluid and plasma** for a limited number of participants. The results are in a number of files prefixed with “PPMI\_Project\_196\_” and “PPMI\_Project\_222\_”. The columns in these tables are interpreted as related analyte values for each sample, and include protein name, Olink protein analyte target, UNIPROT ID, and the specific text values for each analyte included in the targeted assay panel. Note that results are broken into separate plasma and CSF files, and further divided into the four targeted assay panels used to make up the Olink Explore platform. For further information, see the related methods documents, as well as the Olink website for an overview of the assay and additional detail and guidance on interpreting these results. Projects 196 and 222 contain bridging samples allowing for the results to be combined and harmonized using an appropriate bridging protocol. For user convenience, a bridged copy of these two batches of Olink protein analysis are available as **Project 9000 (Bridged Results of Projects 196 and 222)** and is recommended for most users. This bridged file was created using the official Olink bridging protocol; users should carefully review the associated methods document for further details, and for interpretation notes.

Finally, the **iPSC Catalog Metadata** file (in the subfolder Biosample Inventory) contains metadata relating to multiple projects investigating induced pluripotent stem cells. The columns in this file are covered in the data dictionary but also explained in more detail in the Excel workbook “PPMI\_IPSC\_Catalog\_User\_Guide”.

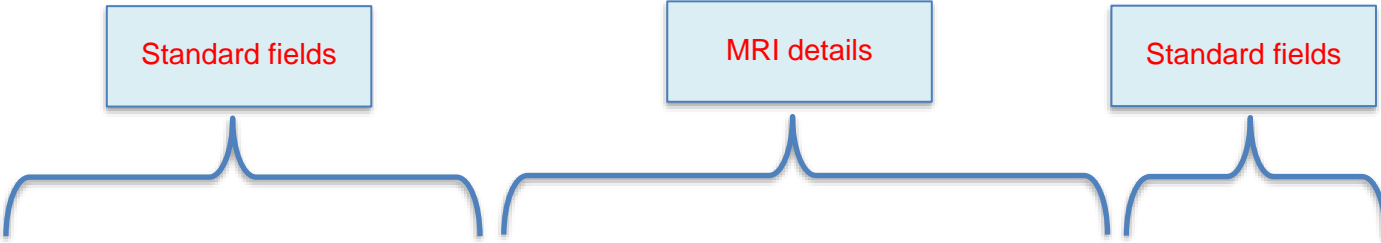


## 9. IMAGE DATA

PPMI contains a large number of images from various image modalities including Magnetic Resonance Imaging (MRI) and Dopamine Transporter scans (DaTscan). Metadata about these images is available for download in tabular format and the images themselves can be viewed and downloaded using a separate tool.

### 9.1 Image metadata

Image metadata is stored in a series of tables in a format similar to other measurement data. For example, for MRI there is the table `Magnetic_Resonance_Imaging__MRI_`:



REC_ID	PAT NO	EVEN T_ID	PAG_NAM E	INFOD T	MRI CMPL T	MRI WDT I	MRI WRS S	MRI RSLT	MRI RSSD F	PDM EDDT	PDM EDTM	ORIG_ENTRY	LAST_UPDAT E
12345681	9001	BL	MRI	Feb-11	1	0	NULL	2	NULL	NULL	NULL	Jan-12	2021-06-21
12345682	9001	V14	MRI	NULL	0	NULL	NULL	NUL L	NULL	NULL	NULL	Jan-19	2021-06-21
12345683	9003	BL	MRI	Mar-11	1	0	NULL	2	NULL	NULL	NULL	Jan-12	2021-06-21
12345684	9004	BL	MRI	Mar-11	1	0	NULL	2	NULL	NULL	NULL	Jan-12	2021-06-21
12345685	9005	BL	MRI	Apr-11	1	0	NULL	1	NULL	NULL	NULL	Jan-12	2021-06-21
12345686	9006	BL	MRI	Apr-11	1	0	NULL	2	NULL	NULL	NULL	Feb-12	2021-06-21
12345687	9007	BL	MRI	May-11	1	0	NULL	2	NULL	NULL	NULL	Feb-12	2021-06-21
12345688	9008	BL	MRI	May-11	1	0	NULL	2	NULL	NULL	NULL	Feb-12	2021-06-21
12345689	9009	BL	MRI	Jun-11	1	0	NULL	2	NULL	NULL	NULL	Feb-12	2021-06-21

The definitions of the MRI specific fields can be looked up in the Data Dictionary, and related permissible code values in the Code\_List:

Field name	Description	Permissible values
MRICMPLT	Brain MRI was completed	0 = Not Completed, 1 = Completed
MRICMPLTCM	Brain MRI not completed comment	(Text field)
MRIRSLT	Results	1 = Normal, 2 = Abnormal, not clinically significant, 3 = Abnormal, clinically significant
MRIRSSDF	Resting state seq-diff day PDmed use	0 = No, 1 = Yes
MRIWDTI	DTI MRI Scan	0 = No, 1 = Yes
MRIWRSS	With resting state sequences	0 = No, 1 = Yes
PDMEDDT	Last dopaminergic medication dose date	Date
PDMEDTM	Last dopaminergic medication dose time	Time

## 9.2 Viewing and downloading images

To view and download images, please refer to the walkthrough that provides a step-by-step guide on how to search, download, and view the desired neuroimaging data from PPMI in the following link:

<https://bit.ly/38A7tmt>

## 10. GENETIC DATA

Information about a participant's genetics can be found in the PPMI Clinical data in several places to different levels of detail.

### 10.1 Genetic subgroups

Genetic screening is performed against known highly penetrant genetic variants conferring elevated risk of developing PD.<sup>22</sup> As explained in Sections 3, this information is captured on both the **Consensus Committee Analytic Dataset** spreadsheet and as a series of columns on the **Participant\_Status** table.

This will tell you whether a participant was designated at the time of enrollment in the study as having a PINK1, LRRK2, Parkin, SNCA or GBA genetic variant (or some combination of more than one of these), as well as indicating whether the participant is a sporadic case, a prodromal RBD case or a prodromal HPSM case. Note that the RBD and HPSM cases might become sporadic cases after baseline if the participant subsequently develops Parkinson's, but the designation at the time of enrollment will not change.

Below are 4 example participants from this table: one sporadic case, one with a GBA variant, one with a LRRK2 variant and one with an SNCA variant.

PATNO	COHORT_DEFINITION	ENRLPINK1	ENRLPRKN	ENRLSRD C	ENRLHPS M	ENRLRB D	ENRLRRK 2	ENRLSNC A	ENRLGB A
9021	Parkinson's Disease	0	0	1	0	0	0	0	0
9022	Prodromal	0	0	0	0	0	0	0	1
9023	Parkinson's Disease	NULL	NULL	NULL	0	0	1	0	0
9024	Parkinson's Disease	0	0	0	0	0	0	1	0

If the value of one of these fields is NULL, then this indicates that the test was not performed, or the test result is not available.

### 10.2 Tests for known PD genetic risk factors

The table **iu\_genetic\_consensus** (identified in the download site as 'Consensus APOE SNPs and GBA and LRRK2 Coding Variant Summary') contains information about what genetic analysis has been performed on each participant, for example whether they have been subject to genome wide associated studies (GWAS) or RNA sequencing. Additionally, for LRRK2 and GBA, the table shows the specific alleles of these two genes for which the participant has been tested. This information is only sporadically updated.

<sup>22</sup> This information is commonly captured during screening as part of inclusion/exclusion criteria and verified through genetic testing prior to baseline.

Below is the content of this table for the same 4 example participants that we used previously (note: synthetic data for illustrative purposes; the table has been split into two sections for ease of viewing).

PATNO	CLIA	GWAS	WES	WGS	SANGE R	RNASE Q	RNASE Q _VIS	APOE	GBA_ POS	GBA_ PATHV AR
9021	-	X	X	X	X	X	5	E3/E3	0	0
9022	X	-	-	-	-	-	0	E2/E3	1	1
9023	X	X	-	X	-	-	0	E2/E3	0	0
9024	X	-	-	X	-	X	4	E3/E3	0	0

PATNO	GBA_VAR_ID	GBA_VAR_ CONF	LRRK2_ POS	LRRK2_ PATHVAR	LRRK2_VAR_ ID	LRRK2_VAR_ CONF
9021	none	none	1	0	S1647T/S164 7T / M2397T/M2 397T	S1647T/S164 7T / M2397T/M2 397T
9022	N409S/N409 S	none	0	0	none	none
9023	none	none	1	1	S1647T/S164 7T / G2019S / M2397T/M2 397T	S1647T/S164 7T / G2019S / M2397T/M2 397T
9024	none	none	1	0	N2081D / M2397T	N2081D / M2397T

The table **Genetic\_Testing\_Results** contains the genetic test results for GBA, LRRK2 and SNCA testing: A value of 1 in MUTRSLT indicates a positive match and the field GENECA has a value of 1 for LRRK2, 2 for SNCA and 3 for GBA. The field LRRKCD gives the specific LRRK2 allele match, in this example a value of 1 being G2019S|c.6055G>A|p.Gly2019Ser.

PATNO	EVENT_ID	PAG_NAME	INFODT	GENECAT	LRRKCD	MUTRSLT
9031	GMU	MUTRSLT	Feb-16	3		1
9032	GMU	MUTRSLT	Mar-19	1	1	1
9033	GMU	MUTRSLT	Apr-18	2		1

Additionally, the table **PPMI\_PD\_Variants\_Genetic\_Status** (identified in the download site as “Participant Genetic Status for Selected PD-Associated Variants”) contains information about matches to specific alleles on other genes known to be associated with PD, though it is important to note that this has not been updated since 2018.

## 10.3 RNA-Sequencing data

RNA-Sequencing was performed in Project 133, which produced 3 processed datasets, namely IR1, IR2 and IR3. IR1 and IR2 are earlier versions and IR3 is the latest version. These are all derived from the same the same sequencing data (the experiment was not repeated), but the structure and pipeline changed slightly to keep them up to date with progress in the field. The biggest difference was the genome build used for alignment. “Build 37” (used for IR1/IR2) is slightly older, and more widely used, whereas “Build 38” (used for IR3) is newer. We recommend using IR3, unless there is a specific reason to utilize the older build and annotation. The raw files are also available upon request in BAM<sup>23</sup> and FASTQ<sup>24</sup> formats, but note that they are very large with a total size of ~142 TB. Please note, while these files are available via web portal download, due to the file size, you may encounter difficulties. We recommend submitting a request for Aspera-based cloud data transfer using the Genetic Data Request Form described below.

Some technical notes about the RNA-Sequencing data:

- The downloaded .tar.gz files contain transcript abundance estimates presented as both featureCounts (Liao *et al.* 2014) and Transcripts per Million (TPM) calculated using the Salmon method (Patro *et al.* 2017). The featureCounts are stored in files under a sub-folder “counts” and the TPM figures are stored in files with the suffix “.sf” in a sub-folder “quant”. Note that these figures have been normalized and details of the normalization method can be found in Craig *et al.* (2021).
- In the file metaDataIR3.tsv can be found information about the phase it was sequenced in, e.g., “PPMI-Phase2-IR1.3174.V08.0000373111.5104-SL-0001”. Phase 1 and Phase 2 represent sequencing waves, or freezes, and are a point where we started analyzing while waiting for additional plates. Covariates are located here. Diagnostic and other clinical data are available elsewhere on ppmi-info.org.

## 10.4 Other genetic data

PPMI has several other processed genetic datasets that are available for download via the “Genetic Data” option on the main PPMI download page. Note that many of these files are very large and you will likely need specialist software to analyze them. We recommend reviewing the methods documentation to understand available data formats prior to downloading these files.

---

<sup>23</sup> See <https://github.com/samtools/hts-specs/blob/master/SAMv1.pdf> for the specification of SAM/BAM file formats.

<sup>24</sup> See Cock *et al.* 2010

**Whole Genome Sequencing (WGS)** data has been generated by **project 118** for most participants (for historical reasons). Processed files are available in Variant Call Format (VCF) and Genomic Variant Call Format (.gVCF). The raw files are also available at request in BAM and FASTQ formats but note that they are very large with a total size of ~184TB. There is also a method document available on the download page.

**Whole Exome Sequencing (WES)** data has been generated by **project 116** for most participants. Processed files are in Genomic Variant Call Format (.gVCF). The raw files are also available at request in BAM and FASTQ formats but note that they are very large with a total size of ~17.5TB. There is also a method document available on the download page.

The **Foundational Data Initiative for Parkinson's Disease (FOUNDIN-PD; project 150)** generated a series of 'omics datasets for 95 inducible pluripotent stem cell lines from PPMI participants as part of a study to identify progression markers in PD. For an overview of the study and further information, refer to the methods documents provided and see <https://www.foundinpd.org>. Processed files are in a variety of file formats corresponding to the distinct data types generated from these resources. The raw files are also available at request in FASTQ, .idat, .jpg, and .raw formats but note that they are very large with a total size of ~27TB.

**Gene Sequencing** refers to primer-targeted Sanger sequencing data that was generated by **project 115** for all PPMI participants, to specifically analyze genetic variability in the SNCA gene. The files are in VCF format; for further information, see the associated method document.

**Genotyping** refers to NeuroX array-based genotyping data generated by **project 107** which performed a series of genotyping studies on several hundred of the PPMI participants against a wide range of polymorphisms covering various autoimmune, inflammatory and neurological diseases as well as PD loci. The data is stored in the PLINK format (Purcell *et al.*, 2007) and further information can be found in the methods documents.

**Methylation Profiling** contains data relating to **projects 120 and 140**, which performed whole genome array-based methylation profiling on most participants. The files are in the Illumina .idat format and further information can be found in the methods documents.

If you require a genetic dataset in its raw format, then fill out and submit the related request form, both found under the "Download Genetic Data" section of PPMI@LONI. FASTQ/BAM users should submit the "Genetic Data Request Form"; Project 193 iPSC WGS users should submit the "PPMI Cell Lines Data Request Form."

## Download Genetic Data

Browse the available items within categories or use the Search feature to find items by keyword.

**Reminder: The PPMI Data Use agreement prohibits unauthorized sharing of these data, posting to public databases and any attempt data to identify individuals using these data. By downloading these data you acknowledge our [Terms and Conditions](#).**

Data Request Form  
Request Form  
Exome Sequencing  
FOUNDIN-PD  
Gene Sequencing  
Genotyping  
Methylation Profiling  
RNA Sequencing  
Whole Genome Sequencing  
ALL

Data Request Form:

*Form for requesting access to download data in FASTQ/BAM format*

Name of Dataset	Version	File Type	Last Download
Request Form			
<a href="#">Genetic Data Request Form</a>	2022-08-02	.DOC format	2022-10-27

## 11. CURATED DATA

Recognizing that to perform even a simple study using the PPMI Clinical data it is necessary to join several tables together, the PPMI data team has prepared three curated datasets where frequently used data from multiple tables has been merged (i.e., denormalized) into a single large table for ease of use.

For investigators that are interested in starting their PPMI data exploration with a dataset, **we recommend using dataset #3 (labeled v.2023-06-12 in the download data)**. Each of these curated data sets comes with a data dictionary containing descriptions of each of the variables, decode values where applicable, and supplementary information regarding the source of the variables. The data dictionaries can be found in the same folder as the data files.

The data dictionary for dataset #3 is labeled v.2023-12-14. Some key points to note from this data dictionary are:

- The data dictionary lists the variables in the same order in which they appear in the data file. For each variable (column B) it provides a category (column A - for instance ID for identifiers, Clinical for clinical data, Genetics for genetic test data, and so on) and a description (column C). The variables are grouped by categories.
- Where the variable relates to a coded value such as COHORT, EVENT\_ID, gender or race, the data dictionary provides the full list of permissible values and the corresponding decode values (columns D, E).
- Columns F, G, H and I of the data dictionary give the source of the data. If Derived Variable (column F) is set to "No", then columns G and H will detail the source column name and table name respectively from where the value was sourced. For example, the variable COHORT is taken from the STATUS table, column COHORT. To find the physical table name for download, recall from Section 3.5 that this can be looked up in the main data dictionary; in this example the source table name is PATIENT\_STATUS.
- If the Derived Variable value is set to "Yes", then the source variable(s) and table(s) are listed in columns G and H respectively and column I is populated with the details. For example, the variables relating to the MOCA (cognitive assessment) test have been aggregated into a single overall test result. About a third of the variable are derived, often to cut down on the overall number of columns.

We recommend you review the data dictionary for the curated dataset you are using because some of the exact field definitions may be different from what you may assume them to be.

The PPMI data team may create more curated cuts in the future if there is sufficient demand.



#	Participant population	Longitudinal profile	Summary of data attributes provided (not exhaustive)
1	683 PD participants	Baseline + 5 annual study visits	<p>159 attributes, including:</p> <ul style="list-style-type: none"> <li>• PATNO</li> <li>• EVENT_ID</li> <li>• Demographics (sex, race, age, etc.)</li> <li>• Family history of PD</li> <li>• Primary diagnosis</li> <li>• Initial symptoms at diagnosis</li> <li>• Symptomatic therapy start date</li> <li>• MDS-UPDRS Part 1 detailed scores</li> <li>• MDS-UPDRS Parts 1, 2, 3, 4 total scores</li> <li>• Tremor score</li> <li>• Rigidity score</li> <li>• Dominant side at diagnosis</li> <li>• TD/PIGD classification</li> <li>• Hoehn &amp; Yahr stage</li> <li>• PD medications</li> <li>• Levodopa Equivalent Daily Dose (LEDD)</li> <li>• QUIP disorder</li> <li>• REM sleep behavior disorder</li> <li>• HVLT results</li> <li>• SCOPA-AUT scores</li> <li>• Semantic fluency scores</li> <li>• STAi scores</li> <li>• MCI test scores</li> <li>• MOCA score</li> <li>• Symbol digit modalities score</li> <li>• Benton judgement of line score</li> <li>• Modified Schwab &amp; England ADL score</li> <li>• Epworth sleepiness score</li> <li>• UPSIT score</li> <li>• CSF A-beta and tau test results</li> <li>• Hemoglobin test results</li> <li>• Serum uric acid</li> <li>• DaTscan measurements for caudate, putamen and striatum</li> <li>• APOE test results</li> <li>• Geriatric depression score</li> <li>• Genetic screening (MAPT, SNCA)</li> </ul>

#	Participant population	Longitudinal profile	Summary of data attributes provided (not exhaustive)
2	86 Prodromal participants	Baseline + year 1 study visit	<p>115 attributes, including:</p> <ul style="list-style-type: none"> <li>• PATNO</li> <li>• EVENT_ID</li> <li>• Demographics (sex, race, age, etc.)</li> <li>• Family history of PD</li> <li>• Primary diagnosis</li> <li>• Symptomatic therapy start date</li> <li>• MDS-UPDRS Part 1 detailed scores</li> <li>• MDS-UPDRS Parts 1, 2, 3, 4 total scores</li> <li>• Tremor score</li> <li>• Rigidity score</li> <li>• TD/PIGD classification</li> <li>• Hoehn &amp; Yahr stage</li> <li>• PD medications</li> <li>• Levodopa Equivalent Daily Dose (LEDD)</li> <li>• QUIP disorder</li> <li>• REM sleep behavior disorder</li> <li>• HVLT results</li> <li>• SCOPA-AUT scores</li> <li>• Semantic fluency scores</li> <li>• STAi scores</li> <li>• MOCA score</li> <li>• Symbol digit modalities score</li> <li>• Benton judgement of line score</li> <li>• Modified Schwab &amp; England ADL score</li> <li>• Epworth sleepiness score</li> <li>• UPSIT score</li> <li>• CSF A-beta and tau test results</li> <li>• Hemoglobin test results</li> <li>• Serum uric acid</li> </ul>
3	2347 covering comprising: Sporadic PD – 803 GBA PD – 292 LRRK2 PD – 404 PINK1 PD – 1 Parkin PD – 10	Baseline + up to 12 annual study visits	<p>161 attributes, including:</p> <ul style="list-style-type: none"> <li>• PATNO</li> <li>• COHORT</li> <li>• subgroup</li> <li>• EVENT_ID</li> <li>• subgroup</li> <li>• Demographics (sex, race, age, etc.)</li> </ul>

#	Participant population	Longitudinal profile	Summary of data attributes provided (not exhaustive)
	SNCA PD – 39 Prodromal (RBD and/or hyposmia) – 528 Healthy controls - 260		<ul style="list-style-type: none"> <li>• DaTscan measurements for caudate, putamen and striatum</li> <li>• APOE test results</li> <li>• CSFSAA: CSF SAA combined results</li> <li>• Biologics, including tests for A-beta, A-synuclein, tau, urate, and neurofilament light</li> <li>• MDS-UPDRS Part 1 detailed scores</li> <li>• MDS-UPDRS Parts 1, 2, 3, 4 total scores</li> <li>• Tremor score</li> <li>• Rigidity score</li> <li>• Dominant side at diagnosis</li> <li>• TD/PIGD classification</li> <li>• Hoehn &amp; Yahr stage</li> <li>• Levodopa Equivalent Daily Dose (LEDD)</li> <li>• QUIP disorder</li> <li>• REM sleep behavior disorder</li> <li>• HVLT results</li> <li>• SCOPA-AUT scores</li> <li>• Semantic fluency scores</li> <li>• STAi scores</li> <li>• MCI test scores</li> <li>• MOCA score</li> <li>• Symbol digit modalities score</li> <li>• Benton judgement of line score</li> <li>• Modified Schwab &amp; England ADL score</li> <li>• Epworth sleepiness score</li> <li>• UPSIT score</li> <li>• Geriatric depression score</li> </ul>

## 12. PPMI ONLINE

PPMI Online is data from participant self-assessments that has been captured online. There is some participant overlap with both PPMI Clinical and PPMI Remote so that some participants may have data for all three, some for two of the three, and some for PPMI Online only. Data for PPMI Online is available from the same data download site as PPMI Clinical under the section “PPMI Online”.

The PATNO identifier is shared with PPMI Clinical and PPMI Remote data, but for longitudinal data, a different set of EVENT\_IDs is used, with the PPMI Online identifiers all having the form “OLxx” where xx is a two digit number.

Taking the example of MDS-UPDRS scores, in order to combine the data across PPMI Clinical and PPMI Remote, we need to take into account the different column names for the same fields, noting also that INFODT maps to CREATED\_AT. Script 10 in the appendices performs this mapping and results in a combined view like the following (note the “OL” EVENT\_IDs);

INFODT	EVENT_ID	NP1SLPN	NP1SLPD	NP1PAIN	NP1URIN	NP1CNST	....
03/01/2011	SC	3	1	1	3	1	
04/01/2011	BL	3	1	1	3	0	
06/01/2011	V01	2	1	0	3	1	
08/01/2011	V02	3	1	3	3	0	
02/01/2012	V03	2	1	1	1	0	
04/01/2012	V04	3	1	1	3	0	
09/01/2012	V05	2	2	0	3	0	
04/01/2013	V06	2	1	1	1	0	
10/01/2013	V07	2	1	1	2	1	
04/01/2014	V08	2	1	1	3	1	
10/01/2014	V09	2	1	1	2	1	
04/01/2015	V10	3	1	1	2	0	
10/01/2015	V11	3	1	1	2	0	
04/01/2016	V12	2	2	1	3	0	
04/01/2017	V13	4	2	1	2	1	
04/01/2019	V15	3	2	1	2	2	
04/01/2022	OL01	1	2	1	3	0	
10/01/2022	OL02	2	2	1	3	1	

For a summary of the all the other data available for PPMI Online, please refer to Appendix A.

## 13. PPMI REMOTE

The objective of PPMI Remote is to gather data about the pre-diagnosis phase of PD by conducting remote tests such as olfactory tests and genetic analysis to capture possible risk indicators of later developing PD. People are referred to PPMI Remote via both PPMI Online and an initiative called “Smell Test Direct”, because they have a known risk factor for developing PD.

Data for PPMI Remote is available from the same data download site as PPMI Clinical under the section “PPMI Remote”. The tables available for download are prefixed with “Remote\_”:

1. Remote\_Screening\_Smell\_Test\_Direct captures information about participants referred for screening from Smell Test Direct
2. Remote\_Screening\_Screener captures participant demographic data such as ethnicity
3. Remote\_Screening\_High\_Interest captures whether the participant has key prodromal markers of PD such as anosmia or acting out their dreams
4. Remote\_Screening\_Participant\_Progress records dates of key events of participants in PPMI Remote, including enrollment, UPSIT testing and DaTscans
5. Remote\_Screening\_UPSIT\_Screening records results of an abridged UPSIT questionnaire that the participants are given
6. Remote\_University\_of\_Pennsylvania\_Smell\_Identification\_TEST (UPSIT) contains the actual UPSIT results for participants who are not also enrolled in PPMI Clinical\*

\*Note that a *limited* set of UPSIT test results (from participants who have completed a DaTSCAN screening visit) are stored in the same table as used by PPMI Clinical (University\_of\_Pennsylvania\_Smell\_Identification\_Test\_\_UPSIT) rather than the PPMI Remote table.; filter on UPSIT\_SOURCE = Remote. See also Section 6.2.

## 14. PPMI FOUND

Data for PPMI FOUND is available from the same data download site as PPMI Clinical under the section “PPMI FOUND”. FOUND stands for “Follow up of persons with Neurologic Disease” and records additional information for some PPMI Clinical participants captured by telephone or other remote consultation methods, in parallel with the main study.

The FOUND tables contain a variety of data relating to known or hypothesized PD risk factors, from alcohol consumption to pesticides at work. Further details can be found in Appendix A. There is also a data dictionary in the file FOUND\_RFQ\_Dictionary available for download.

The FOUND tables are referenced using PATNO in the usual way, but note that they do not have EVENT\_IDs. The date of data capture is recorded in a field called “datacompXX” (standing for data completion) where “XX” is different for each table. For example, in the table FOUND\_RFQ\_Head\_Injury, there is a patno column (note the lower-case spelling) and a column called datacomphi that records the date when the questionnaire was completed.

The data is captured one time only as it relates to historical information. The exception to this is the table FOUND\_Self-Reported\_Dx, which records diagnoses of PD or other conditions and is updated every 6 months.

## 15. DIGITAL SENSOR DATA

PPMI contains two types of digital sensor data: (1) data from the Roche PD monitoring app and (2) data from the Verily study watch. The Roche data relates to a number of active motor and non-motor tests performed using a mobile app, and is stored in a single table. The Verily data has been captured by smartwatch and comprises biometric measurements in a number of separate tables and is much larger in volume than the Roche data.

Key points about the two datasets are summarized below.

	Roche PD monitoring app	Verily study watch
Nature of data	Active tests of motor and non-motor function	Biometric data relating to sleep, activity and cardio function
Method of capture	Mobile app (Android)	Smartwatch
Data tables and nature of data captures	<p><b>Roche_PD_Monitoring_App_v2_data</b>: consolidated table containing a range of tests covering:</p> <ul style="list-style-type: none"> <li>• <b>voice / speech</b></li> <li>• <b>dexterity</b></li> <li>• <b>bradykinesia</b></li> <li>• <b>tremor</b></li> <li>• <b>balance / postural stability</b></li> </ul> <p>During the study period, information relating to various questionnaires was also captured in the same table, including medication details and functional <b>patient reported outcome measures</b> such as EQ-5D-5L.</p>	<p><b>ambulatory</b>: minutes in the day when the participant is walking, running etc.</p> <p><b>inbedtimes</b>: time the participant went to bed and got out of bed</p> <p><b>onwrist</b>: details of when the Verily study watch was attached to the participant's wrist and gathering data</p> <p><b>prv</b>: pulse rate variability per hour</p> <p><b>pulserate</b>: mean hourly pulse rate</p> <p><b>sleepmetrics2</b>: metrics relating to timing, amount and quality of sleep</p> <p><b>sleepstage</b>: timing and duration of individual sleep stages (REM, NREM deep, NREM light, awake and unknown)</p> <p><b>stepcount</b>: hourly step count</p> <p><b>timezone</b>: timezone that each measurement was taken in</p>
Number of participants (as at July 31 <sup>st</sup> , 2023) – see Section 3.2 for definitions	32 (20 Parkinson's and 12 Prodromal)	353 (148 Parkinson's, 158 Prodromal, 35 Healthy Control, 2 SWEDD, 10 Early Imaging study)

	Roche PD monitoring app	Verily study watch
Approximate date range of data capture (as of July 31 <sup>st</sup> 2023)	May 2019 – December 2020	September 2017 – January 2021

## 15.1 Roche PD monitoring app: additional information

The tests/questionnaires in the Roche data are identified by the field QRSCATID. A full list of test codes and response options can be found in the Codes table (see Section 4) but note that at as the publication date of this guide, only some tests have data available, and these are:

Type	QRSCATID	Test description
Tests	DAS	Draw a shape test (note: the field QRSSCAT gives sub-test)
	FST	Speech test
	HT	Maximum hand turn speed
	PTT	Postural tremor test
	Q135	eSDMT (symbol digit modalities test)
	RTT	Resting tremor test
	SBT	Balance test
	SPT	Voice test
	ST12	Dexterity test
	UTT	U-Turn average turn speed
Questionnaires	DMT	Daily medication tracker (was medication taken and when?)
	HADS	HADS (anxiety and depression rating)
	MAT	Missed Active Test (reason for missing a test)
	PD-DDS	Daily diary (e.g., sleep quality, bowel movements)
	PD-DQS	Daily patient reported outcomes
	PD-HS	Health survey – ANS (autonomic function)



Type	QRSCATID	Test description
	PGIC	PGIC (patient global impression of change)
	Q008	EQ-5D-5L (self-care, activities, anxiety etc.)
	Q070	PDQ39 (standard 39-item PD questionnaire)
	Q226	PAC-SYM (constipation assessment)
	SETUP	Setup survey (handedness and home WiFi availability)

The corresponding answers to the questions and measurements can be found in the fields QRSRESN (for numeric values) or QRSRESC (for text values). Where applicable, the field QRSORESU gives the units of the measurement. QRSSTDTC and QRSSTDTC\_TIME give the start date and time of the test respectively; QRSENDTC and QRSENDTC\_TIME give the end date and time of the test.

For further information on the methodology behind each of the tests, please refer to the Methods documents provided. Assessment studies of the various tests are described in Lipsmeier *et al.* (2022).

## 15.2 Verily study watch derived data: additional information

In the Verily data, to protect participant confidentiality there are no absolute date and time markers. Instead, the timing of each measurement is recorded in three ways:

1. As a day of the week (time\_day field) and a local time provided in hh:mm:ss format (time\_local)
2. As the number of milliseconds since the start of the study watch usage (time\_ms)
3. As the participant's age in seconds (age\_seconds) where each whole year is defined as 31,556,952 seconds (i.e., 365 days, 5 hours, 4 minutes and 12 seconds)

To link the Verily data to the rest of PPMI Clinical, we need to do two things:

- Link PATNO to the field 'subject' in the Verily data
- To find the EVENT\_ID corresponding to the age\_seconds field on the Verily measurement, we take the ENROLL\_AGE field from the Patient\_Status table and use this to determine the age of the participant at each event (i.e. by comparing with the INFODT fields on a table containing event data such as one of the MDS UPDRS tables). We can subsequently multiply the age by 31,556,952 to convert to a value in seconds then find the value that is closest to the age\_seconds on the Verily table.

To help with the second point, script 11 has been provided in the appendices. This creates a new version of the pulserate table with the nearest EVENT\_ID appended to the end of the table. The script can be easily extended to the other tables by following the same pattern.

For further information on the Verily data, please also refer to the Methods documents provided.

## APPENDIX A – SUMMARY OF KEY PPMI STUDY DATA

This appendix gives a high-level summary of the main PPMI study data and supporting documents available for download that researchers may be interested in. Not every dataset of document is listed; to see the full list, visit the download page itself (<https://ida.loni.usc.edu/pages/access/studyData.jsp?project=PPMI>).

In addition to the datasets identified below, there are separate download options for image data (see Section 9) and certain genetic data (see Section 10).

Group	Subgroup	Filename	Comments
START HERE	Quick Start	Overview/Quick Start Guide	Deprecated by this user guide
		Consensus Committee Analytic Dataset	Defines participant subgroups. See Section 3.
		PPMI Analytic Dataset Guide	Explanation of participant subgroups. See also Section 3.
Study Docs	Data & Databases	Code List	List of all codes used in tables. "Annotated" version contains additional comments
		Data Dictionary	List of all tables and columns. "Annotated" version contains additional comments
		PPMI Data Sequester Guidelines	Information about additional data that has not been made available for download, for example because of confidentiality restrictions
	Study Methods	Method for calculating Levodopa Equivalent Dose in PPMI Data	See Section 7
		Methods for Defining PD Med Use	Details of how MDS-UPDRS on and off periods are defined and recorded
		Symptomatic Therapy Visits	Information about "ST" visits conducted for some participants prior to starting medication
	Study Protocol & CRFs	Study CRFs and Assessments	Details of all standard questionnaires used by PPMI, including PASE (Physical Activity Scale for the Elderly); screening & demographics; socio-economics; criteria for

Group	Subgroup	Filename	Comments
			classification as PD / Prodromal / Healthy Control; AV-133 eligibility, diagnostic features, medical history, family history, neurological exam, physical exam, vital signs, medication, Schwab & England daily living, etc.
Subject characteristics	Family History	Family History	Details of PD family history
	Genetic Status	Various	See Section 10
	Patient Status	Participant Status	Core table of participant status and classification. See Section 3.
	Subject Demographics	Age at visit	Participant age at each study visit
		Demographics	Basic demographic data such as sex, ethnicity, and date of birth
		Socio-Economics	Years of education
		Subject Cohort History	Mapping of participants to cohorts
Biospecimen	Biosample Inventory	Various	Various catalogs of available biospecimens. For further information on biospecimen data see Section 8.
	Biospecimen Sample Analysis	Biospecimen Analysis Results Overview	Overview of different biospecimen result sets
		Current Biospecimen Analysis Results	Details of current biosample tests performed including test name, test result, date and researcher
		Deprecated Biospecimen Analysis Results	Biospecimen analysis results that have been superseded.
		Pilot Biospecimen Analysis Results	Results of an earlier pilot project relating to SAA. See Section 8.4.
		Project 181 Adaptive Immune Markers for Predicting Cognitive Decline in PD	Results of project 181

Group	Subgroup	Filename	Comments
		SAA Biospecimen Analysis Results	Results of the SAA project: see Section 8.4.
	Biospecimen Analysis Methods	Various	Documents detailing the sample collection and analysis methods used for each biosample project
	Lab Collection Procedures	Various	Details of lab procedures used for blood chemistry, lumbar punctures, skin biopsies etc.
	Metabolomic Analysis	Various	Metabolomic tests and result
	Proteomic Analysis	Various	Proteomic tests and results
	Proteomic Analysis Methods	Various	Documents detailing the procedures used to perform proteomic analysis
Curated Data Cuts	Curated Data	Curated data cut – Original Cohort Baseline to Year 5 Data	A prepared dataset containing selected measurements from baseline and 5 annual study visits for 683 PD participants
		Curated data cut – Original Cohort Baseline to 5 Year Data Dictionary	Data Dictionary for the above
		Curated data cut – Prodromal Cohort Baseline to Year 1 Data	A prepared dataset containing selected measurements from baseline and a study visit after 1 year for 86 prodromal participants
		Curated data cut – Prodromal Cohort Baseline to Year 1 Data Dictionary	Data Dictionary for the above
		Curated data cut (v.2023-06-12)	A prepared dataset taken on June 12 <sup>th</sup> , 2023 containing selected measurements from 2347 Parkinson's (both genetic and idiopathic), prodromal and healthy control participants. See Section 11.
		Curated data cut dictionary (v.2023-06-12)	Data dictionary for the above

Group	Subgroup	Filename	Comments
	Supporting Documentation	Variable Definitions and Score Calculations	Supporting documentation for the curated data
Digital Sensor	Roche Smartphone App	Roche PD Monitoring App v2 data and methods	Measurements for selected participants from the Roche Smartphone app and associated method document. See Section 15.
	Verily Study Watch	Various	Measurements for selected participants from the Verily Study Watch and associated method documentation. See Section 15.
Enrollment	Subject Enrollment	Various	Files associated with participant study enrollment, capturing participant consent, recording participation in sub-studies, etc.
Imaging	DaTSCAN Imaging	DaTscan Analysis	DaTscan measurements of striatal binding ratio (SBR) of left and right caudate, putamen and anterior putamen regions with respect to the occipital lobe
		DaTscan Imaging	Metadata relating to DaTscans, e.g., data of scan, location
		DaTscan Metadata	Further metadata relating to DaTscans, e.g., time of scan, dosage of radiotracer
		DaTscan SPECT Image Processing SBR Calculation Methods	Method document for the SBR calculation
		DaTscan SPECT Visual Interpretation Assessment Methods	Method document for the visual interpretation
		DaTscan Visual Interpretation Results	Summary result of DaTscan visual interpretation, i.e., positive / negative
	Diffusion Imaging	DTI Preprocessing Manual	Documentation and results relating to a diffusion tensor imaging (DTI) study

Group	Subgroup	Filename	Comments
			on 15 participants (Schuff <i>et al.</i> , 2015)
	Magnetic Resonance Imaging	Grey Matter volume extraction from MRI T1	Grey matter volume estimated from MRI scan
		Magnetic Resonance Imaging (MRI)	Metadata relating to MRI scans, e.g., date, type of scan, comments
		MRI Metadata	Metadata relating to MRI results, e.g., were DTI, NM-MT, T1-weighted, T2-weighted FLAIR sequences acquired?
		Rs-fMRI Data Availability and Usage	Documentation and results relating to a sub-study of resting-state fMRI on 194 participants between 2012 and 2019
	PET Imaging	Various	Various technical documents, metadata and results relating to the tau sub-study, <sup>25</sup> and the early imaging sub-study that generated AV-133 <sup>26</sup> PET images
Medical History	Medical	Tau sub-study (various tables)	Data relating to the participant participating in the tau sub-study
		Adverse Event Telephone Assessment	Indicates if an adverse event occurred shortly after a lumbar puncture, skin biopsy, or DaTscan imaging procedure, and what contact was made with the participant about it
		Concomitant Medication Log	History of medication other than dopaminergic medication
		Determination of Freezing and Falls	Details of freezing and falls captured in later study visits
		Early Imaging (various tables)	Data relating to the Early Imaging cohort (see Section 3)

<sup>25</sup> For details of the tau sub-study see [https://www.ppmi-info.org/sites/default/files/docs/PPMI%20Tau%20PET%20Imaging%20Protocol\\_26Apr2021\\_Final.pdf](https://www.ppmi-info.org/sites/default/files/docs/PPMI%20Tau%20PET%20Imaging%20Protocol_26Apr2021_Final.pdf)

<sup>26</sup> For details of the early imaging sub-study which generated the AV-133 images see [https://www.ppmi-info.org/sites/default/files/docs/PPMI%202.0\\_Early%20Imaging%20Protocol\\_v2\\_19June2020\\_Final.pdf](https://www.ppmi-info.org/sites/default/files/docs/PPMI%202.0_Early%20Imaging%20Protocol_v2_19June2020_Final.pdf)

Group	Subgroup	Filename	Comments
		Features of Parkinsonism	Key clinical features captured in later study visits (rigidity, bradykinesia, tremor, postural instability)
		General Physical Exam	Results of general physical examination carried out in screening
		LEDD Concomitant Medication Log	History of dopaminergic medication; see Section 7
		Medical Conditions Log	Year of diagnosis of other medical conditions (arthritis, cancer, high cholesterol, depression etc.)
		Other Clinical Features	Other clinical features (cf. Features of Parkinsonism table) captured in later study visits (tremor, dystonia, micrographia, shuffling gait, etc.)
		PD Diagnosis History	Date of PD diagnosis and details of key motor symptoms
		Pregnancy Test	Pregnancy tests and results
		Primary Clinical Diagnosis	Details of diagnosis (idiopathic PD, MSA, PSP, MND with Parkinsonism, etc.)
		Prodromal History	Details of prodromal symptoms (REM behavior disorder, first degree family member with PD, hyposmia, etc.)
		Report of Pregnancy	Details of pregnancy during study period
		Surgery for PD Log	Details of PD related surgeries (e.g., Deep Brain Stimulation)
		Vital signs	Details of vital signs (blood pressure, temperature, weight, height, etc.)
	Neurological Exam	Neurological Exam	Summary of neurological tests performed at study visits (motor exam, gait assessment, reflexes, etc.)
		Prodromal Diagnostic Code Memo	Guidance for coding on the Prodromal History table



Group	Subgroup	Filename	Comments
	Safety Monitoring	Adverse Event Log	Register of adverse symptoms during study (headache, nausea, pain around lumbar puncture site, back pain, etc.)
Motor Assessments	Motor/MDS-UPDRS	Gait Data & Arm swing	Gait and arm swing measurements
		Gait Data & Arm swing Methods	Description of the method for gait and arm swing measurements
		MDS-UPDRS Part I	Part 1 of the MDS-UPDRS, capturing non-motor aspects of daily living
		MDS-UPDRS Part I Patient Questionnaire	Part 1 of the MDS-UPDRS, capturing non-motor aspects of daily living as documented by the participant
		MDS-UPDRS Part II Patient Questionnaire	Part 2 of MDS-UPDRS, capturing motor aspects of daily living as documented by the participant
		MDS-UPDRS Part III	Part 3 of the MDS-UPDRS capturing results of the motor examination
		MDS-UPDRS Part III ON/OFF Determination	Part 3 of the MDS-UPDRS, supplemental data about how the examination was performed
		MDS-UPDRS Part IV Motor Complications	Part 4 of the MDS-UPDRS capturing information about dyskinesia, on/off fluctuations etc.
		Modified Schwab & England Activities of Daily Living	Schwab & England total score
		Neuro QoL: Lower Extremity Function (Mobility)	Quality of life questionnaire (mobility) supplementary to MDS-UPDRS
		Neuro QoL: Lower Extremity Function (Fine Motor, ADL)	Quality of life questionnaire (fine motor control and activities of daily living) supplementary to MDS-UPDRS
		Participant Motor Function Questionnaire	Motor function questionnaire, supplementary to MDS-UPDRS
Non-motor Assessments	Autonomic Tests	SCOPA-AUT	Results of SCOPA-AUT assessment, covering 25 data points on: gastrointestinal, urinary,

Group	Subgroup	Filename	Comments
			cardiovascular, thermoregulatory, pupillomotor, and sexual
	Cognition	Various	Results of various cognitive tests, including clock drawing, trail marking and modified Boston naming test
	Neurobehavioral Tests	Various	Results of various neurobehavioral assessments including depression and anxiety
	Neuropsychological Tests	Various	Results of various neurophysiological test including the Montreal Cognitive Assessment (MoCA), modified semantic fluency, and the Hopkins verbal learning test
	Olfactory Tests	Olfactory UPSIT Methods	The method for the UPSIT test (Doty <i>et al.</i> , 1984)
		University of Pennsylvania Smell Identification Test (UPSIT)	Results of the UPSIT test
	Sleep Disorder Tests	Various	Results of the Epworth Sleepiness Scale questionnaire and the REM Sleep Behavior Disorder Questionnaire
PPMI Found	Follow up of persons with Neurologic Disease	FOUND Brief Summary (PDF)	A summary of FOUND. See Section 14 for additional information on PPMI FOUND
		FOUND RFQ Alcohol	History and current level of alcohol consumption
		FOUND RFQ Anti-inflammatory Meds	History and current usage of ibuprofen, aspirin and other NSAIDs
		FOUND RFQ Caffeine	History and current level of coffee consumption
		FOUND RFQ Calcium Channel Blockers	History of usage of calcium channel blockers and reason for usage
		FOUND RFQ Dictionary	Data dictionary for the FOUND data

Group	Subgroup	Filename	Comments
		FOUND RFQ Female Reproductive Health	Details of pregnancies, menstruation, menopause and hormone replacement therapy
		FOUND RFQ Head Injury	Details of historical head injuries
		FOUND RFQ Height and Weight	Details of body type, height and weight at ages 25, 40 and 60
		FOUND RFQ Occupation	Details of military service, current occupation and main jobs in each decade of life
		FOUND RFQ Pesticides at Work	Details of exposure to pesticides at work in each decade of life, including specific compounds
		FOUND RFQ Pesticides Non-Work	Details of exposure to pesticides in non-work settings
		FOUND RFQ Physical Activity	Summary of level of physical activity and hours of sleep per night at each decade in life
		FOUND RFQ Residential History	History of residential location and proximity to farms and crop spraying, and possible water contamination
		FOUND RFQ Smoking History	History and current level of cigarette and other tobacco usage
		FOUND RFQ Toxicant History	History of exposure to toxicants other than pesticides, e.g., adhesives, solvents and metals
		FOUND RFQ Self Reported Dx	Diagnosis of PD and other conditions, reported every 6 months by the participant
PPMI Online	Data & Databases Online	PPMI Online Codebook PPMI Online Dictionary	The codebook and data dictionary for the PPMI Online study. See Section 12 for an overview of PPMI Online.
	Cognition-Online	Cognitive Change	Yes/no record of cognitive change (can be captured by someone other than the participant)



Group	Subgroup	Filename	Comments
	Family History-Online	Family History of PD	First degree relatives with a diagnosis of PD
	Medical-Online	Assessment of Constipation Caffeine Consumption Chemical Exposure COVID-19 History Head Injuries Health History Annually Health History Quarterly History of Falls: Baseline History of Falls: Surveillance Medication History Pesticides at Work Physical Activity Smoking History	Participant medical history. In some cases, these tables capture information at regular intervals in the participant's lifetime. Please refer to the PPMI Online Dictionary for further details.
	Motor / MDS-UPDRS-Online	MDS-UPDRS Part 1 MDS-UPDRS Part 2 Participant Motor Function	The sections of the MDS-UPDRS questionnaires that are completed by participants plus some supplementary functional motor questions
	Neurobehavioral-Online	Geriatric Depression Scale Parkinson Anxiety Scale	The short form Geriatric Depression Scale (Yesavage <i>et al.</i> , 1983) and the Parkinson's Anxiety Scale (Leentjens <i>et al.</i> , 2004)
	Neuropsychological-Online	Penn Parkinson's Daily Activities Questionnaire	PDAQ-15 (Brennan <i>et al.</i> , 2016)
	Patient Status-Online	Age of PD Diagnosis Genetic Testing Results High Interest Questions for non-PD Cohort High Interest Questions for PD Cohort Hyposmia PD Return Study for non-PD PD Return Study for PD Participant Visit Information	Various tables associated with participant status. study visits, prodromal symptoms and PD diagnosis. Please refer to the PPMI Online Dictionary for further details.
	Sleep Disorder-Online	Epworth Sleepiness Scale PD Sleep Scale PDSS-2 PPMI RBD Sleep Questionnaire	The Epworth Sleepiness Scale (Johns, 1991), the PD Sleep Scale, revised version (Trenkwalder <i>et al.</i> , 2011), a PPMI specific sleep

Group	Subgroup	Filename	Comments
		RBD1Q Postuma Acting Out Dreams	questionnaire, and the RBD1Q test for RBD (Postuma <i>et al.</i> , 2012)
	Study Enrollment-Online	How you heard about PPMI Non-Completer Survey Participant Enrollment Status	Information about participant enrollment and which cohort they are in (Parkinson's / Non-Parkinson's)
	Demographics-Online	Occupation and Military Service Race and Ethnicity Registration Information Residential Location Socioeconomic Status	Various tables capturing participant demographic information. Please refer to the PPMI Online Dictionary for further details.
PPMI Remote	PPMI Remote Screening	Remote Screening High Interest Questionnaire	A list of questions that indicate whether a participant has one or more key prodromal indicators (e.g., olfactory problems or acting out dreams). See also Section 13 for an overview of PPMI Remote.
		Remote Screening Participant Progress	Administrative details of remote screening: dates, statuses etc.
		Remote Screening Screener	Demographic information about the participant being screened, e.g., gender, ethnicity
		Remote Screening Smell Test Direct Screener	Basic information about participants referred through the Smell Test Direct initiative, e.g., date, PD diagnosis status
		Remote Screening UPSIT Screening Questionnaire	Selected questions from the University of Pennsylvania Smell Identification Test used to screen participants
		Remote University of Pennsylvania Smell Identification Test (UPSIT)	Results of the full UPSIT questionnaire comprising 40 smell tests. (Doty <i>et al.</i> , 1984)  Note that this only captures data for participants who have not also enrolled in PPMI Clinical. For participants who are enrolled in both PPMI Clinical and PPMI Remote, the data can be found in the UPSIT file

Group	Subgroup	Filename	Comments
			on the non-motor assessments group, where the PPMI Remote records can be identified with PAG_NAME = 'UPSITPRO'. See also Section 6.2.
Archived PPMI Data	Data Collected in Prior EDC	Various	Not likely to be required
	Publication-associated Archives	Various	Not likely to be required



## APPENDIX B – SQL SCRIPTS

*Note that all code/scripts within this guide are provided as is with no guarantee of accuracy or completeness. For questions, reach out to [resources@michaeljfox.org](mailto:resources@michaeljfox.org).*

### Script 1: Count of enrolled/withdrawn/complete participants by cohort and genetic subgroup (Section 3.7)

Note that, by default, MS SQL Server, and certain other DBMSs such as MySQL, are case insensitive with respect to column values. However, if using a DBMS such as PostgreSQL that is case sensitive it will be necessary to match on both 'Enrolled' and 'enrolled' and both 'Withdraw' and 'withdrew'.

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Number of enrolled Parkinson's and Prodromal participants by genetic subgroup */

SELECT COHORT_DEFINITION,
       (CASE WHEN ENRLPINK1 + ENRLPRKN + ENRLSRDC + ENRLHPSM + ENRLRBD + ENRLRRK2 +
        ENRLSNCA + ENRLGBA > 1 THEN 'Multiple factors'
        WHEN ENRLPINK1 = 1 THEN 'PINK1'
        WHEN ENRLPRKN =1  THEN 'PARKIN'
        WHEN ENRLSRDC =1  THEN 'SRDC'
        WHEN ENRLHPSM = 1 THEN 'HPSM'
        WHEN ENRLRBD = 1  THEN 'RBD'
        WHEN ENRLRRK2 = 1 THEN 'LRRK2'
        WHEN ENRLSNCA =1  THEN 'SNCA'
        WHEN ENRLGBA = 1  THEN 'GBA'
        ELSE NULL END) AS 'Genetic subgroup',
       COUNT (*) AS 'Patient count'
FROM PPMI.dbo.Participant_Status
WHERE COHORT_DEFINITION LIKE 'P%' /* Filter on Parkinson's and prodromal cases */
AND ENROLL_STATUS IN ('Enrolled', 'Withdrew', 'Complete')
GROUP BY COHORT_DEFINITION,
         (CASE WHEN ENRLPINK1 + ENRLPRKN + ENRLSRDC + ENRLHPSM + ENRLRBD + ENRLRRK2 +
        ENRLSNCA + ENRLGBA > 1 THEN 'Multiple factors'
        WHEN ENRLPINK1 = 1 THEN 'PINK1'
        WHEN ENRLPRKN =1  THEN 'PARKIN'
        WHEN ENRLSRDC =1  THEN 'SRDC'
        WHEN ENRLHPSM = 1 THEN 'HPSM'
        WHEN ENRLRBD = 1  THEN 'RBD'
        WHEN ENRLRRK2 = 1 THEN 'LRRK2'
        WHEN ENRLSNCA =1  THEN 'SNCA'
        WHEN ENRLGBA = 1  THEN 'GBA'
        ELSE NULL END)
ORDER BY COHORT_DEFINITION, COUNT(*) DESC
```



## Script 2: Count of enrolled/withdrawn/complete participants by cohort and sex (Section 4.1)

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Number of PD and prodromal participants by sex */

SELECT S.COHORT_DEFINITION, C.DECODE AS 'SEX',
       COUNT(*) AS 'Patient Count'
FROM PPMI.dbo.Demographics D
     LEFT OUTER JOIN PPMI.dbo.Codes C ON C.ITM_NAME = 'SEX' AND C.CODE = D.SEX
     LEFT OUTER JOIN PPMI.dbo.Participant_Status S ON S.PATNO = D.PATNO
WHERE S.COHORT_DEFINITION LIKE 'P%'
     AND S.ENROLL_STATUS IN ('Enrolled', 'Withdrew', 'Complete')
GROUP BY S.COHORT_DEFINITION, C.DECODE
ORDER BY S.COHORT_DEFINITION, C.DECODE
```

## Script 3: Create participant master table (Section 4.2)

Note that the first time you run this, the DROP TABLE command will issue a warning because the table does not yet exist; you can ignore this.

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Create master table of selected participants demographic and status data */

DROP TABLE PPMI.dbo.Participant_Master

SELECT D.PATNO, D.BIRTHDT, S.COHORT_DEFINITION,
       (CASE WHEN S.ENRLPINK1 + S.ENRLPRKN + S.ENRLSRDC + S.ENRLHPSM + S.ENRLRBD +
S.ENRLLRK2 + S.ENRLSNCA + S.ENRLGBA > 1 THEN 'Multiple factors'
       WHEN S.ENRLPINK1 = 1 THEN 'PINK1'
       WHEN S.ENRLPRKN = 1 THEN 'PARKIN'
       WHEN S.ENRLSRDC = 1 THEN 'SRDC'
       WHEN S.ENRLHPSM = 1 THEN 'HPSM'
       WHEN S.ENRLRBD = 1 THEN 'RBD'
       WHEN S.ENRLLRK2 = 1 THEN 'LRRK2'
       WHEN S.ENRLSNCA = 1 THEN 'SNCA'
       WHEN S.ENRLGBA = 1 THEN 'GBA'
       ELSE NULL END) AS 'Genetic subgroup',
       ROUND (S.ENROLL_AGE,1) AS 'ENROLL_AGE',
       S.ENROLL_DATE, S.ENROLL_STATUS, C1.DECODE as 'SEX', C2.DECODE as 'HANDED',
       (SELECT MIN(PDDXDT)
        FROM PPMI.dbo.PD_Diagnosis_History DH
        WHERE DH.PATNO = D.PATNO) AS 'PD diagnosis date'
INTO PPMI.dbo.Participant_Master

FROM PPMI.dbo.Demographics D
     LEFT OUTER JOIN PPMI.dbo.Codes C1 ON C1.ITM_NAME = 'SEX' AND C1.CODE = D.SEX
     LEFT OUTER JOIN PPMI.dbo.Codes C2 ON C2.ITM_NAME = 'HANDED' AND C2.CODE = D.HANDED
     LEFT OUTER JOIN PPMI.dbo.Participant_Status S ON S.PATNO = D.PATNO

WHERE S.COHORT_DEFINITION LIKE 'P%'
     AND S.ENROLL_STATUS IN ('Enrolled', 'Withdrew', 'Complete')
ORDER BY D.PATNO
```





#### Script 4: Progression of MDS-UPDRS scores (Section 6.1)

Change the participant number highlighted below to the participant in which you are interested.

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Extract progression of total scores and Hoehn and Yahr stage from each section of MDS-
UPDRS for a single participant */

SELECT P1.INFODT, P1.EVENT_ID,
       MAX(P1.NP1RTOT) AS 'Part 1'
       ,MAX(P1P.NP1PTOT) AS 'Part 1P'
       ,MAX(P2.NP2PTOT) AS 'Part 2'
       ,MAX(P3.NP3TOT) AS 'Part 3'
       ,MAX(P3.NHY) AS 'H & Y Stage'
FROM PPMI.dbo.MDS_UPDRS_Part_1 P1
JOIN PPMI.dbo.MDS_UPDRS_Part_1P P1P
  ON P1.PATNO = P1P.PATNO AND P1.INFODT = P1P.INFODT AND P1P.NP1PTOT IS NOT NULL
JOIN PPMI.dbo.MDS_UPDRS_Part_2P P2
  ON P1.PATNO = P2.PATNO AND P1.INFODT = P2.INFODT AND P2.NP2PTOT IS NOT NULL
JOIN PPMI.dbo.MDS_UPDRS_Part_3 P3
  ON P1.PATNO = P3.PATNO AND P1.INFODT = P3.INFODT AND P3.NP3TOT IS NOT NULL
WHERE P1.PATNO = 9030
GROUP BY P1.INFODT, P1.EVENT_ID
ORDER BY P1.INFODT, P1.EVENT_ID
```

#### Script 5: Sex and total UPSIT score by participant (Section 6.2)

This assumes that the smell test results have been imported into a table called UPSIT.

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Extract sex and total UPSIT score at baseline for PD participants */
SELECT U.PATNO, DECODE AS 'SEX', U.TOTAL_CORRECT
FROM PPMI.dbo.UPSIT U, PPMI.dbo.DEMOGRAPHICS D, PPMI.dbo.Codes C,
PPMI.dbo.PARTICIPANT_STATUS S
WHERE U.EVENT_ID = 'BL' AND U.PATNO = D.PATNO
AND U.TOTAL_CORRECT IS NOT NULL
AND S.PATNO = D.PATNO AND S.ENROLL_STATUS IN ('Enrolled', 'Withdrew', 'Complete')
AND S.COHORT_DEFINITION LIKE 'Park%'
AND C.CODE = D.SEX AND C.ITM_NAME = 'SEX'
```



## Script 6: LEDD medication history for single participant (Section 7.1)

Change the participant number highlighted below to the participant in which you are interested.

```
/* This script is provided as is with no guarantees of completeness or accuracy */  
/* LEDD medication history for single participant */  
SELECT PATNO, EVENT_ID, PAG_NAME, INFODT, LEDTRT, STARTDT, STOPDT, LEDD  
FROM PPMI.dbo.LEDD_Concomitant_Medication_Log  
WHERE PATNO = 9052  
ORDER BY STARTDT
```



## Script 7: Convert LEDD records into total LEDD at each point in time (Section 7.1)

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Convert LEDD values into total LEDD at a point in time */
DROP TABLE PPMI.dbo.LEDD

CREATE TABLE PPMI.dbo.LEDD (
    PATNO INT,
    STARTDT DATE,
    LDOPA FLOAT, /* Total Levodopa daily quantity at a point in time */
    LEDD FLOAT) /* Total LEDD across all dopaminergic medication at a point in time */

/* Get all the dates when medication started or stopped, i.e. dates when changed */
INSERT INTO PPMI.dbo.LEDD
SELECT DISTINCT PATNO, STARTDT AS 'STARTDT', 0.0, 0.0 AS LEDD FROM
PPMI.dbo.LEDD_Concomitant_Medication_Log
UNION
SELECT DISTINCT PATNO, DATEADD("MONTH", 1, STOPDT) AS 'STARTDT', 0.0, 0.0 AS LEDD FROM
PPMI.dbo.LEDD_Concomitant_Medication_Log WHERE STOPDT IS NOT NULL

/* Get the base L-Dopa value and base LEDD value at each point in time */
UPDATE L SET LEDD = S.LEDD, LDOPA = S.LDOPA FROM PPMI.dbo.LEDD L,
    (SELECT L1.PATNO AS 'PATNO', L1.STARTDT AS 'STARTDT',
        SUM(CASE WHEN (L2.LEDTRT LIKE '%LEV%' OR L2.LEDTRT LIKE '%DOP%' OR
L2.LEDTRT LIKE '%RYTA%'
            OR L2.LEDTRT LIKE '%MADOPAR%' OR L2.LEDTRT LIKE '%SINEMET%' OR
L2.LEDTRT LIKE '%STANEK%'
            OR L2.LEDTRT LIKE '%STALEVO%' OR L2.LEDTRT LIKE '%ISICOM%' OR
L2.LEDTRT LIKE '%INBRIJA%'
            OR L2.LEDTRT LIKE '%NACOM%' OR L2.LEDTRT LIKE '%PROLOPA%') THEN
        (CASE WHEN ISNUMERIC (L2.LEDD)=1 THEN CONVERT (FLOAT, L2.LEDD)
ELSE
            L2.LEDDSTRMG * L2.LEDDOSE * L2.LEDDOSFRQ END) ELSE 0.0 END) AS
'LDOPA',
        SUM(CASE WHEN ISNUMERIC (L2.LEDD)=1 THEN CONVERT (FLOAT, L2.LEDD) ELSE
        (CASE WHEN (L2.LEDTRT LIKE '%LEV%' OR L2.LEDTRT LIKE '%DOP%' OR L2.LEDTRT
LIKE '%RYTA%'
            OR L2.LEDTRT LIKE '%MADOPAR%' OR L2.LEDTRT LIKE '%SINEMET%' OR
L2.LEDTRT LIKE '%STANEK%'
            OR L2.LEDTRT LIKE '%STALEVO%' OR L2.LEDTRT LIKE '%ISICOM%' OR
L2.LEDTRT LIKE '%INBRIJA%'
            OR L2.LEDTRT LIKE '%NACOM%' OR L2.LEDTRT LIKE '%PROLOPA%') THEN
            L2.LEDDSTRMG * L2.LEDDOSE * L2.LEDDOSFRQ
            ELSE 0.0 END) END) AS 'LEDD'
FROM PPMI.dbo.LEDD L1, PPMI.dbo.LEDD_Concomitant_Medication_Log L2
WHERE L2.PATNO = L1.PATNO AND L2.STARTDT <= L1.STARTDT
AND (L2.STOPDT >= L1.STARTDT OR L2.STOPDT IS NULL)
AND (L2.STOPDT >= L2.STARTDT OR L2.STOPDT IS NULL)
GROUP BY L1.PATNO, L1.STARTDT) S
WHERE S.PATNO = L.PATNO AND S.STARTDT = L.STARTDT
```

Continued on next page...



```
/* Update LEDD value for drugs with LEDD profile of the form LD x A */
UPDATE L SET LEDD = L.LEDD + S.LEDD FROM PPMI.dbo.LEDD L,
  (SELECT L1.PATNO AS 'PATNO', L1.STARTDT AS 'STARTDT', SUM(L1.LDOPA * CONVERT
(FLOAT, RIGHT(L2.LEDD, LEN(L2.LEDD) - 4))) AS 'LEDD'
  FROM PPMI.dbo.LEDD L1, PPMI.dbo.LEDD_Concomitant_Medication_Log L2
  WHERE L2.PATNO = L1.PATNO AND L2.STARTDT <= L1.STARTDT
  AND (L2.STOPDT >= L1.STARTDT OR L2.STOPDT IS NULL)
  AND (L2.STOPDT >= L2.STARTDT OR L2.STOPDT IS NULL) AND L2.LEDD LIKE 'LD x%'
  GROUP BY L1.PATNO, L1.STARTDT) S
WHERE S.PATNO = L.PATNO AND S.STARTDT = L.STARTDT

/* Update LEDD value for drugs with LEDD profile of the form (B + LD) x A */
UPDATE L SET LEDD = L.LEDD + S.LEDD FROM PPMI.dbo.LEDD L,
  (SELECT L1.PATNO AS 'PATNO', L1.STARTDT AS 'STARTDT', SUM(( CONVERT (FLOAT,
SUBSTRING (L2.LEDD, 2, CHARINDEX ('+ LD', L2.LEDD, 1)-3)) + L1.LDOPA ) * CONVERT
(FLOAT, RIGHT(L2.LEDD, 4))) AS 'LEDD'
  FROM PPMI.dbo.LEDD L1, PPMI.dbo.LEDD_Concomitant_Medication_Log L2
  WHERE L2.PATNO = L1.PATNO AND L2.STARTDT <= L1.STARTDT
  AND (L2.STOPDT >= L1.STARTDT OR L2.STOPDT IS NULL)
  AND (L2.STOPDT >=
    L2.STARTDT OR L2.STOPDT IS NULL) AND L2.LEDD LIKE '%+ LD%'
  GROUP BY L1.PATNO, L1.STARTDT) S
WHERE S.PATNO = L.PATNO AND S.STARTDT = L.STARTDT
```

### Script 8: Numbers of participants with each blood test (Section 8.1)

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Determine which blood tests have been captured for all participants */
SELECT LTSTCODE, LTSTNAME, COUNT(DISTINCT PATNO) AS 'Patient Count'
  FROM PPMI.dbo.Blood_Chemistry___Hematology
  GROUP BY LTSTCODE, LTSTNAME
  ORDER BY COUNT(DISTINCT PATNO) DESC, LTSTNAME
```

### Script 9: Blood test details for single participant (Section 8.1)

Change the participant number highlighted below to the participant in which you are interested.

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Hemoglobin test results for single participants */
SELECT EVENT_ID, LCOLLDT, LTSTNAME, LVISTYPE, LSIRE, LSIUNIT, LUSRES, LUSUNIT
  FROM PPMI.dbo.Blood_Chemistry___Hematology
  WHERE LTSTCODE = 'HMT40'
  AND PATNO = 9002
  ORDER BY LCOLLDT
```



## Script 10: Merge MDS\_UPDRS scores from PPMI Clinical and PPMI Online (Section 12)

Change the participant number highlighted below to the participant in which you are interested.

```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Merge MDS-UPDRS Part 1 and Part 2 patient scores from PPMI Clinical and PPMI Online
data for a single participant */

SELECT P1.INFODT, P1.EVENT_ID, NP1SLPN, NP1SLPD, NP1PAIN, NP1URIN, NP1CNST,
NP1LTHD, NP1FATG, NP2SPCH, NP2SALV, NP2SWAL, NP2EAT, NP2DRES, NP2HYGN, NP2HWRT,
NP2HOBBS, NP2TURN, NP2TRMR, NP2RISE, NP2WALK, NP2FREZ
FROM PPMI.dbo.MDS_UPDRS_Part_1P P1
JOIN PPMI.dbo.MDS_UPDRS_Part_2P P2
ON P1.PATNO = P2.PATNO AND P1.INFODT = P2.INFODT
WHERE P1.PATNO = 9030
UNION
SELECT P1.CREATED_AT AS INFODT, P1.EVENT_ID, NP1SLPN_OL AS NP1SLPN, NP1SLPD_OL
AS NP1SLPD, NP1PAIN_OL AS NP1PAIN, NP1URIN_OL AS NP1URIN, NP1CNST_OL AS
NP1CNST, NP1LTHD_OL AS NP1LTHD, NP1FATG_OL AS NP1FATG, NP2SPCH_OL AS NP2SPCH,
NP2SALV_OL AS NP2SALV, NP2SWAL_OL AS NP2SWAL, NP2EAT_OL AS NP2EAT, NP2DRES_OL
AS NP2DRES, NP2HYGN_OL AS NP2HYGN, NP2HWRT_OL AS NP2HWRT, NP2HOBBS_OL AS NP2HOBBS,
NP2TURN_OL AS NP2TURN, NP2TRMR_OL AS NP2TRMR, NP2RISE_OL AS NP2RISE, NP2WALK_OL
AS NP2WALK, NP2FREZ_OL AS NP2FREZ
FROM PPMI.dbo.MDS_UPDRS_Part_1_Online P1
JOIN PPMI.dbo.MDS_UPDRS_Part_2_Online P2
ON P1.PATNO = P2.PATNO AND P1.CREATED_AT = P2.CREATED_AT
WHERE P1.PATNO = 9030

ORDER BY INFODT, EVENT_ID
```

## Script 11: Find EVENT\_IDs corresponding to Verily measurements (Section 15)



```
/* This script is provided as is with no guarantees of completeness or accuracy */
/* Append EVENT_ID corresponding to age_seconds to the pulserate table */
/* Creates a new table pulserate2 but this can easily be extended to other tables */

DROP TABLE PPMI.dbo.pulserate2

SELECT P.subject, P.age_seconds, MIN(ABS((P.age_seconds - S.ENROLL_AGE * 31556952)
    - DATEDIFF(ss, convert(DATE, CONCAT('01/', S.ENROLL_DATE)), M.INFODT))) AS
MIN_TIME_DIFF
INTO PPMI.dbo.pulserate_temp
FROM PPMI.dbo.MDS_UPDRS_Part_1 M,
     PPMI.dbo.Participant_Status S,
     PPMI.dbo.pulserate P
WHERE S.PATNO = subject AND M.PATNO = subject AND S.ENROLL_DATE IS NOT NULL
GROUP BY P.subject, P.age_seconds
ORDER BY P.subject, P.age_seconds

SELECT P.*, M.EVENT_ID AS NEAREST_EVENT_ID
INTO PPMI.dbo.pulserate2
FROM PPMI.dbo.MDS_UPDRS_Part_1 M,
     PPMI.dbo.Participant_Status S,
     PPMI.dbo.pulserate_temp E,
     PPMI.dbo.pulserate P
WHERE MIN_TIME_DIFF = ABS((E.age_seconds - S.ENROLL_AGE * 31556952) -
    DATEDIFF(ss, convert(DATE, CONCAT('01/', S.ENROLL_DATE)), M.INFODT))
    AND E.subject = P.subject AND S.PATNO = P.subject AND M.PATNO = P.subject
    AND E.age_seconds = P.age_seconds
    AND S.ENROLL_DATE IS NOT NULL
ORDER BY P.subject, P.age_seconds

DROP TABLE PPMI.dbo.pulserate_temp

SELECT TOP 1000 * FROM PPMI.dbo.pulserate2 ORDER BY subject, age_seconds, time_ms
```



## APPENDIX C – R SCRIPTS

***Note that all code/scripts within this guide are provided as is with no guarantee of accuracy or completeness. For questions, reach out to [resources@michaeljfox.org](mailto:resources@michaeljfox.org).***

In the sample R scripts we use functions from the readr, dplyr and stringr libraries to perform data manipulation. You will need to install the tidyverse package<sup>27</sup> before you can use these (in RStudio select the Install Packages option under Tools to do this).

All scripts assume the data has been downloaded as csv files to C:\PPMI; you can change this as required.

### Script 1: Count of enrolled/withdrawn/complete participants by cohort and genetic subgroup (Section 3.7)

```
# This script is provided as is with no guarantees of completeness or accuracy
# Number of enrolled/completed/withdrawn Parkinson's and Prodromal participants by
# genetic subgroup

library(readr)
library(dplyr)

# Load table into data frame
setwd ("C:\\PPMI")
Participant_Status <- read_csv ("Participant_Status.csv")

# Filter on status and cohort, determine genetic subgroup, then group by, aggregate
# and order results
filter(Participant_Status, tolower(ENROLL_STATUS) %in% c("enrolled", "withdrew",
"complete") & substr(COHORT_DEFINITION,1, 1)== "P") %>%
transmute(PATNO, COHORT_DEFINITION, GENETIC_SUBGROUP = factor(case_when(
  ENRLPINK1 + ENRLPRKN + ENRLSRDC + ENRLHPSM + ENRLRBD + ENRLLRK2 + ENRLSNCA +
ENRLGBA > 1 ~ 'Multiple factors',
  ENRLPINK1 == 1 ~ 'PINK1',
  ENRLPRKN == 1 ~ 'PARKIN',
  ENRLSRDC == 1 ~ 'SRDC',
  ENRLHPSM == 1 ~ 'HPSM',
  ENRLRBD == 1 ~ 'RBD',
  ENRLLRK2 == 1 ~ 'LRK2',
  ENRLSNCA == 1 ~ 'SNCA',
  ENRLGBA == 1 ~ 'GBA',
  TRUE ~ '' ))) %>%
group_by(COHORT_DEFINITION, GENETIC_SUBGROUP) %>%
summarize(PATIENT_COUNT = n()) %>%
arrange(COHORT_DEFINITION, desc(PATIENT_COUNT))
```

<sup>27</sup> See <https://www.tidyverse.org/>.



## Script 2: Count of enrolled/withdrawn/complete participants by cohort and sex (Section 4.1)

```
# This script is provided as is with no guarantees of completeness or accuracy
# Number PD and Prodromal enrolled/withdrawn/complete participants by cohort and sex

library(readr)
library(dplyr)

# Load tables into data frames
setwd ("C:\\PPMI")
Participant_Status <- read_csv ("Participant_Status.csv")
Demographics <- read_csv ("Demographics.csv")
Codes <- read_csv ("Code_List_-_Harmonized.csv")

# Filter based on cohort and enrollment status
Participant_Status_Filtered <- filter(Participant_Status, tolower(ENROLL_STATUS) %in%
c("enrolled", "withdrew", "complete") & substr(COHORT_DEFINITION,1, 1)== "P")

# Extract the sex decode values from the Codes data frame, noting that a type
conversion to integer is needed
Sex <- filter(Codes,ITM_NAME == "SEX") %>% select(CODE,DECODE) %>% transmute(CODE =
as.numeric(as.character(CODE)), DECODE)

# Join the tables, group, aggregate and order the results
Participant_Status_Filtered %>%
  left_join (Demographics, by = "PATNO") %>%
  left_join (Sex,c("SEX"= "CODE")) %>%
  group_by (COHORT_DEFINITION, DECODE) %>%
  summarize(PATIENT_COUNT = n()) %>%
  arrange (COHORT_DEFINITION, DECODE) %>%
  rename (SEX = DECODE)
```





### Script 3: Create participant master table (Section 4.2)

```
# This script is provided as is with no guarantees of completeness or accuracy
# Create a "participants master" data frame

library(readr)
library(dplyr)

# Load tables into data frames
setwd ("C:\\PPMI")
Participant_Status <- read_csv ("Participant_Status.csv")
Demographics <- read_csv ("Demographics.csv")
Codes <- read_csv ("Code_List_-_Harmonized.csv")
PD_Diagnosis_History <- read_csv ("PD_Diagnosis_History.csv")

# Filter on status and cohort and determine genetic subgroup
Participant_Status_Filtered <- filter(Participant_Status, tolower(ENROLL_STATUS) %in%
c("enrolled", "withdrew", "complete") & substr(COHORT_DEFINITION,1, 1)== "P") %>%
  mutate(GENETIC_SUBGROUP = factor(case_when(
    ENRLPINK1 + ENRLPRKN + ENRLSRDC + ENRLHPSM + ENRLRBD + ENRLRRK2 + ENRLSNCA +
    ENRLGBA > 1 ~ 'Multiple factors',
    ENRLPINK1 == 1 ~ 'PINK1', ENRLPRKN == 1 ~ 'PARKIN', ENRLSRDC == 1 ~ 'SRDC',
    ENRLHPSM == 1 ~ 'HPSM', ENRLRBD == 1 ~ 'RBD', ENRLRRK2 == 1 ~ 'LRRK2',
    ENRLSNCA == 1 ~ 'SNCA', ENRLGBA == 1 ~ 'GBA', TRUE ~ '' )))

# Extract the sex and handedness decode values from the Codes data frame, noting that
# a type conversion to integer is needed
Sex <- filter(Codes, ITM_NAME == "SEX") %>% select(CODE, DECODE) %>% transmute(CODE =
as.numeric(as.character(CODE)), DECODE)
Handed <- filter(Codes, ITM_NAME == "HANDED") %>% select(CODE, DECODE) %>% transmute(CODE
= as.numeric(as.character(CODE)), DECODE)

# PD Diagnosis date - tidy up date format and select earliest date for each patient
PD_Diagnosis_History <- mutate(PD_Diagnosis_History, PD_Diagnosis_Date =
as.Date(paste("01/", as.character(PDDXDT)), "%d/%m/%Y")) %>%
  group_by(PATNO) %>%
  summarize (PD_Diagnosis_Date = min(PD_Diagnosis_Date))

# Create the Patient_Master data frame
Participant_Master <- Participant_Status_Filtered %>%
  left_join (Demographics, by = "PATNO") %>%
  left_join (Sex, c("SEX" = "CODE")) %>%
  left_join (Handed, c("HANDED" = "CODE")) %>%
  left_join (PD_Diagnosis_History, by = "PATNO") %>%
  select (PATNO, BIRTHDT, COHORT_DEFINITION, GENETIC_SUBGROUP, ENROLL_AGE,
ENROLL_DATE, ENROLL_STATUS, SEX = DECODE.x, HANDED = DECODE.y, PD_Diagnosis_Date) %>%
  arrange (PATNO)

#Fix date formats
Participant_Master <- mutate(Participant_Master, BIRTHDT = as.Date(paste("01/",
as.character(BIRTHDT)), "%d/%m/%Y"))
Participant_Master <- mutate(Participant_Master, ENROLL_DATE = as.Date(paste("01/",
as.character(ENROLL_DATE)), "%d/%m/%Y"))

#Display data frame
Participant_Master
```



#### Script 4: Progression of MDS-UPDRS scores (Section 6.1)

Change the participant number highlighted below to the participant in which you are interested.

```
# This script is provided as is with no guarantees of completeness or accuracy
# Extract progression of total scores and Hoehn and Yahr sate from each section of
MDS-UPDRS for a single participant

library(readr)
library(dplyr)

# Load tables into data frames
setwd ("C:\\PPMI")
MDS_UPDRS_Part_1 <- read_csv ("MDS-UPDRS_Part_I.csv")
MDS_UPDRS_Part_1P <- read_csv ("MDS-UPDRS_Part_I_Patient_Questionnaire.csv")
MDS_UPDRS_Part_2P <- read_csv ("MDS_UPDRS_Part_II_Patient_Questionnaire.csv")
MDS_UPDRS_Part_3 <- read_csv ("MDS-UPDRS_Part_III.csv")

#Fix date formats
MDS_UPDRS_Part_1 <- mutate(MDS_UPDRS_Part_1, INFODT = as.Date(paste("01/",
as.character(INFODT)), "%d/%m/%Y"))
MDS_UPDRS_Part_1P <- mutate(MDS_UPDRS_Part_1P, INFODT = as.Date(paste("01/",
as.character(INFODT)), "%d/%m/%Y"))
MDS_UPDRS_Part_2P <- mutate(MDS_UPDRS_Part_2P, INFODT = as.Date(paste("01/",
as.character(INFODT)), "%d/%m/%Y"))
MDS_UPDRS_Part_3 <- mutate(MDS_UPDRS_Part_3, INFODT = as.Date(paste("01/",
as.character(INFODT)), "%d/%m/%Y"))

# Join the data frames filtering on a single participant and take the maximum scores
across each INFODT/EVENT_ID
filter(MDS_UPDRS_Part_1, PATNO == 9300) %>%
  inner_join (MDS_UPDRS_Part_1P, by = c("PATNO","INFODT","EVENT_ID")) %>%
  inner_join (MDS_UPDRS_Part_2P, by = c("PATNO","INFODT","EVENT_ID")) %>%
  inner_join (MDS_UPDRS_Part_3, by = c("PATNO","INFODT","EVENT_ID")) %>%
  filter(!is.na(NP1RTOT), !is.na(NP1PTOT) & !is.na(NP2PTOT) & !is.na(NP3TOT) &
!is.na(NHY)) %>%
  select (INFODT, EVENT_ID, NP1RTOT, NP1PTOT, NP2PTOT, NP3TOT, NHY) %>%
  group_by (INFODT, EVENT_ID) %>%
  summarize (Part1 = max(NP1RTOT), Part1P = max(NP1PTOT), Part2 = max(NP2PTOT), Part3
= max(NP3TOT), HY_Stage = max(NHY)) %>%
  arrange (INFODT, EVENT_ID)
```



## Script 5: Sex and total UPSIT score by participant (Section 6.2)

```
# This script is provided as is with no guarantees of completeness or accuracy
# Extract sex and total UPSIT score at baseline for PD participants

library(readr)
library(dplyr)

# Load tables into data frames
setwd ("C:\\PPMI")
Participant_Status <- read_csv ("Participant_Status.csv")
Demographics <- read_csv ("Demographics.csv")
UPSIT <- read_csv
("University_of_Pennsylvania_Smell_Identification_Test__UPSIT_.csv")
Codes <- read_csv ("Code_List_-_Harmonized.csv")

# Filter based on cohort and enrollment status
Participant_Status_Filtered <- filter(Participant_Status,
tolower(ENROLL_STATUS) %in% c("enrolled", "withdrew", "complete") &
substr(COHORT_DEFINITION,1, 4)== "Park")

# Extract the sex decode values from the Codes data frame, noting that a type
conversion to integer is needed
Sex <- filter(Codes,ITM_NAME == "SEX") %>% select(CODE,DECODE) %>%
transmute(CODE = as.numeric(as.character(CODE)), DECODE)

# Join the data frames filtering out the patients with NULL values in the total
UPSIT score
Participant_Status_Filtered %>%
  inner_join (filter(UPSIT,!is.na(TOTAL_CORRECT)), by = "PATNO") %>%
  inner_join (Demographics, by = "PATNO") %>%
  inner_join (Sex,c("SEX"= "CODE")) %>%
```



## Script 6: LEDD medication history for single participant (Section 7.1)

Change the participant number highlighted below to the participant in which you are interested.

```
# This script is provided as is with no guarantees of completeness or accuracy
# LEDD medication history for participant

library(readr)
library(dplyr)

# Load tables into data frames
setwd ("C:\\PPMI")
LEDD <- read_csv ("LEDD_Concomitant_Medication_Log.csv")

# Fix date formatting
LEDD <- mutate(LEDD, STARTDT = as.Date(paste("01/", as.character(STARTDT)),
"%d/%m/%Y"), STOPDT = as.Date(paste("01/", as.character(STOPDT)), "%d/%m/%Y"))

# Select columns of interest from LEDD log for a single participant
LEDD %>% filter(PATNO == 9052) %>%
  select (PATNO, EVENT_ID, PAG_NAME, INFODT, LEDTRT, STARTDT, STOPDT, LEDD) %>%
  arrange (STARTDT)
```



## Script 7: Convert LEDD records into total LEDD at each point in time (Section 7.1)

```
# This script is provided as is with no guarantees of completeness or accuracy
# Convert LEDD values into total LEDD at a point in time

library(readr)
library(dplyr)
library(stringr)

# Load table into data frames
setwd ("C:\\PPMI")
LEDD_Concomitant_Medication_Log <- read_csv ("LEDD_Concomitant_Medication_Log.csv")

# Convert the start and stop dates from MM/YYYYY to 01/MM/YYYY
LEDD_Concomitant_Medication_Log <- mutate(LEDD_Concomitant_Medication_Log, STARTDT =
as.Date(paste("01/", as.character(STARTDT)), "%d/%m/%Y"))
# Add 1 month to the STOPDT because we assume that the dosage is taken until the end
of the month, hence the change takes effect from the start of the next month
LEDD_Concomitant_Medication_Log <- mutate(LEDD_Concomitant_Medication_Log, STOPDT =
as.Date(paste("01/", (as.integer(str_sub(STOPDT, 1, 2))%12+1), "/",
as.integer(str_sub(STOPDT, 4, 7))+(as.integer(str_sub(STOPDT, 1, 2))+1)%13, sep=""),
"%d/%m/%Y"))

# Get the set of unique dates when LEDD changes for each PATNO
LEDD <- unique(rbind(LEDD_Concomitant_Medication_Log %>% select(PATNO, STARTDT),
LEDD_Concomitant_Medication_Log %>% select(PATNO, STARTDT =
STOPDT) %>% filter(!is.na(STARTDT)))) %>% arrange (PATNO, STARTDT)

# Get the base L-Dopa value at each point in time
LDOPA_VALUES <-
LEDD_Concomitant_Medication_Log[str_detect(LEDD_Concomitant_Medication_Log$LEDTRT,regex
("LEV|DOP|RYTA|SINEMET|STANEK|STALEVO|INBRIJA|ISICOM|NACOM|PROLOPA", ignore_case=TRUE)
),] %>%
  select(PATNO, STARTDT, STOPDT, LEDD, LEDDSTRMG, LEDDOSE, LEDDOSFRQ) %>%
# Also calculate the LEDD value from base medication details, if available, in case
not populated in LEDD column
mutate(LEDD3 = case_when (is.numeric(LEDDSTRMG) ~ (LEDDSTRMG * LEDDOSE * LEDDOSFRQ) ,
TRUE ~ 0.0)) %>%
  mutate(LEDD2 = case_when (!is.na(as.double(as.character(LEDD)))) ~ as.double
(as.character(LEDD)) , TRUE ~ LEDD3))
```

Continued on next page...



```
LEDD <- left_join (LEDD,
  left_join(LEDD, LDOPA_VALUES, by="PATNO") %>%
    filter(STARTDT.y <= STARTDT.x & (is.na(STOPDT) | (STOPDT >
STARTDT.x & STOPDT > STARTDT.y))) %>%
    mutate (LEDD2 = case_when (is.na(LEDD2) ~ as.double (0.0), TRUE
~ LEDD2)) %>%

  group_by(PATNO,STARTDT.x) %>%
  summarize (LDOPA = sum(LEDD2), .groups = "drop") %>%
  rename (STARTDT = STARTDT.x),
  by = c("PATNO","STARTDT"))

# Get the base LEDD value at each point in time
# If the drug is levodopa and the LEDD value is not available, attempt to derive it
from base medication details
LEDD_VALUES <- mutate(LEDD_Concomitant_Medication_Log, LEDD3 = case_when
(str_detect(LEDD_Concomitant_Medication_Log$LEDTRT,regex("LEV|DOP|RYTA|SINEMET|STANEK|
STALEVO|INBRIJA|ISICOM|NACOM|PROLOPA",ignore_case=TRUE)) ~ (LEDDSTRMG * LEDDOSE *
LEDDOSFRQ) , TRUE ~ 0.0)) %>%
  mutate(LEDD2 = case_when (!is.na(as.double(as.character(LEDD)))) ~ as.double
(as.character(LEDD)) , TRUE ~ LEDD3))

LEDD <- left_join (LEDD,
  left_join(LEDD, LEDD_VALUES, by="PATNO") %>%
    filter(STARTDT.y <= STARTDT.x & (is.na(STOPDT) | (STOPDT >
STARTDT.x & STOPDT > STARTDT.y))
    & !is.na(LEDD2)) %>%
  group_by(PATNO,STARTDT.x) %>%
  summarize (LEDD = sum(LEDD2), .groups = "drop") %>%
  rename (STARTDT = STARTDT.x),
  by = c("PATNO","STARTDT"))

# Update base LEDD value for drugs with LEDD profile of the form LD * A
VALUES_TO_UPDATE <- LEDD_Concomitant_Medication_Log %>%
  filter (substr(LEDD,1,4) == "LD x") %>%
  select (PATNO, STARTDT,STOPDT, LEDD) %>%
  mutate (FACTOR = as.double(str_match(LEDD, "\\d+\\.\\d+$"))))

LEDD <- left_join (LEDD,
  left_join(LEDD, VALUES_TO_UPDATE, by="PATNO") %>%
    filter(STARTDT.y <= STARTDT.x & (is.na(STOPDT) | (STOPDT >
STARTDT.x & STOPDT > STARTDT.y))) %>%
  group_by(PATNO,STARTDT.x) %>%
  summarize (FACTOR1 = sum(FACTOR), .groups = "drop") %>%
  rename (STARTDT = STARTDT.x),
  by = c("PATNO","STARTDT"))
```

Continued on next page...



```
# Update base LEDD value for drugs with LEDD profile of the form (B + LD) x A
VALUES_TO_UPDATE <- LEDD_Concomitant_Medication_Log %>%
  filter (str_detect(LEDD, "\\+\\ LD")) %>%
  select (PATNO, STARTDT, STOPDT, LEDD) %>%
  mutate (FACTOR2 = as.double(str_match(LEDD, "\\d+\\.\\.\\d+")),
          FACTOR3 = as.double(str_match(LEDD, "\\d+\\.\\.\\d+$")),
          FACTOR4 = as.double(str_match(LEDD, "\\d+\\.\\.\\d+$")))

LEDD <- left_join (LEDD, left_join(LEDD, VALUES_TO_UPDATE, by="PATNO") %>%
  filter(STARTDT.y <= STARTDT.x & (is.na(STOPDT) | (STOPDT >
STARTDT.x & STOPDT > STARTDT.y)))) %>%
  group_by(PATNO, STARTDT.x) %>%
  summarize (FACTOR2 = sum(FACTOR2), FACTOR3 =
mean(FACTOR3), FACTOR4 = sum(FACTOR4), .groups = "drop") %>% rename (STARTDT =
STARTDT.x), by = c("PATNO", "STARTDT"))

# Tidy up null values and make final calculations
LEDD <- mutate (LEDD, LDOPA = case_when(is.na(LDOPA) ~ 0.0, TRUE ~ LDOPA), LEDD =
case_when(is.na(LEDD) ~ 0.0, TRUE ~ LEDD))

LEDD <- mutate (LEDD, LEDD = LEDD
  + case_when(!is.na(FACTOR1) ~ FACTOR1 * LDOPA, TRUE ~ 0.0)
  + case_when(!is.na(FACTOR2) ~ (FACTOR2 * FACTOR3) + (LDOPA * FACTOR4),
TRUE ~ 0.0)) %>%
  select (-FACTOR1, -FACTOR2, -FACTOR3, -FACTOR4)

# Display results
LEDD
```



## Script 8: Numbers of participants with each blood test (Section 8.1)

```
# This script is provided as is with no guarantees of completeness or accuracy
# Determine which blood tests have been captured for all participants

library(readr)
library(dplyr)

# Load table into a data frame
setwd ("C:\\PPMI")
Blood_Chemistry <- read_csv ("Blood_Chemistry___Hematology.csv")

# Count the number of distinct participants with each type of blood test
Blood_Chemistry %>% group_by (LTSTCODE, LTSTNAME) %>%
  summarize (PATIENT_COUNT = n_distinct(PATNO)) %>%
  arrange (desc(PATIENT_COUNT), LTSTNAME) %>%
  print (n = Inf)
```

## Script 9: Blood test details for single participant (Section 8.1)

Change the participant number highlighted below to the participant in which you are interested.

```
# This script is provided as is with no guarantees of completeness or accuracy
# Hemoglobin test results for a single participant

library(readr)
library(dplyr)

# Load table into a data frame
setwd ("C:\\PPMI")
Blood_Chemistry <- read_csv ("Blood_Chemistry___Hematology.csv")

# Fix date formatting
Blood_Chemistry <- mutate(Blood_Chemistry, LCOLLDT = as.Date(paste("01/",
as.character(LCOLLDT)), "%d/%m/%Y"))

# Count the number of distinct participants with each type of blood test
filter (Blood_Chemistry, PATNO == 9002 & LTSTCODE == "HMT40") %>%
  select (EVENT_ID, LCOLLDT, LTSTNAME, LVISTYPE, LSIRES, LSIUNIT, LUSRES, LUSUNIT) %>%
  arrange (LCOLLDT)
```





## Script 10: Merge MDS\_UPDRS scores from PPMI Clinical and PPMI Online (Section 12)

Change the participant number highlighted below to the participant in which you are interested.

```
# This script is provided as is with no guarantees of completeness or accuracy
# Merge MDS-UPDRS Part 1 and Part 2 patient scores from PPMI Clinical and PPMI Online
data for a single participant

library(dplyr)

# Load tables into data frames
setwd ("C:\\PPMI")
MDS_UPDRS_Part_1P <- read.csv ("MDS-UPDRS_Part_I_Patient_Questionnaire.csv",
stringsAsFactors = TRUE)
MDS_UPDRS_Part_2P <- read.csv ("MDS_UPDRS_Part_II_Patient_Questionnaire.csv",
stringsAsFactors = TRUE)
MDS_UPDRS_Part_1_Online <- read.csv("MDS-UPDRS_Part_I_Non-Motor_Aspects__Online_.csv",
stringsAsFactors = TRUE)
MDS_UPDRS_Part_2_Online <- read.csv("MDS-UPDRS_Part_II_Motor_Aspects__Online_.csv",
stringsAsFactors = TRUE)

# Fix date formats
MDS_UPDRS_Part_1P <- mutate(MDS_UPDRS_Part_1P, INFODT = as.Date(paste("01/",
as.character(INFODT)), "%d/%m/%Y"))
MDS_UPDRS_Part_2P <- mutate(MDS_UPDRS_Part_2P, INFODT = as.Date(paste("01/",
as.character(INFODT)), "%d/%m/%Y"))
MDS_UPDRS_Part_1_Online <- mutate(MDS_UPDRS_Part_1_Online, INFODT =
as.Date(paste("01/", as.character(CREATED_AT)), "%d/%m/%Y"))
MDS_UPDRS_Part_2_Online <- mutate(MDS_UPDRS_Part_2_Online, INFODT =
as.Date(paste("01/", as.character(CREATED_AT)), "%d/%m/%Y"))

# Join the data frames for PPMI Clinical
MDS_UPDRS_Clinical <- filter(MDS_UPDRS_Part_1P, PATNO == 9030) %>%
  inner_join (MDS_UPDRS_Part_2P, by = c("PATNO","INFODT","EVENT_ID")) %>%
  select (INFODT, EVENT_ID, NP1SLPN, NP1SLPD, NP1PAIN, NP1URIN, NP1CNST, NP1LTHD,
NP1FATG, NP2SPCH, NP2SALV, NP2SWAL, NP2EAT, NP2DRES, NP2HYGN, NP2HWRT, NP2HOBB,
NP2TURN, NP2TRMR, NP2RISE, NP2WALK, NP2FREZ) %>%
  arrange (INFODT, EVENT_ID)

# Join the data frames for PPMI Online and rename columns to match PPMI Clinical data
MDS_UPDRS_Online <- filter(MDS_UPDRS_Part_1_Online, PATNO == 9030) %>%
  inner_join (MDS_UPDRS_Part_2_Online, by = c("PATNO","INFODT","EVENT_ID")) %>%
  select (INFODT, EVENT_ID, NP1SLPN_OL, NP1SLPD_OL, NP1PAIN_OL, NP1URIN_OL,
NP1CNST_OL, NP1LTHD_OL, NP1FATG_OL, NP2SPCH_OL, NP2SALV_OL, NP2SWAL_OL, NP2EAT_OL,
NP2DRES_OL, NP2HYGN_OL, NP2HWRT_OL, NP2HOBB_OL, NP2TURN_OL, NP2TRMR_OL, NP2RISE_OL,
NP2WALK_OL, NP2FREZ_OL) %>%
  rename (NP1SLPN = NP1SLPN_OL, NP1SLPD = NP1SLPD_OL, NP1PAIN = NP1PAIN_OL, NP1URIN =
NP1URIN_OL, NP1CNST = NP1CNST_OL, NP1LTHD = NP1LTHD_OL, NP1FATG = NP1FATG_OL, NP2SPCH
= NP2SPCH_OL, NP2SALV = NP2SALV_OL, NP2SWAL = NP2SWAL_OL, NP2EAT = NP2EAT_OL, NP2DRES
= NP2DRES_OL, NP2HYGN = NP2HYGN_OL, NP2HWRT = NP2HWRT_OL, NP2HOBB = NP2HOBB_OL, NP2TURN
= NP2TURN_OL, NP2TRMR = NP2TRMR_OL, NP2RISE = NP2RISE_OL, NP2WALK = NP2WALK_OL,
NP2FREZ = NP2FREZ_OL) %>%
  arrange (INFODT, EVENT_ID)

# Combine the two data frames
union(MDS_UPDRS_Clinical, MDS_UPDRS_Online) %>% arrange (INFODT, EVENT_ID)
```



## Script 11: Find EVENT\_IDs corresponding to Verily measurements (Section 15)

```
# This script is provided as is with no guarantees of completeness or accuracy
# Verily smartwatch data - script to append the nearest EVENT_ID to the Verily data

library(dplyr)
# Load tables into data frames
setwd ("C:\\PPMI")
MDS_UPDRS_Part_1P <- read.csv ("MDS-UPDRS_Part_I_Patient_Questionnaire.csv",
stringsAsFactors = TRUE)
pulserate <- read.csv ("pulserate.csv", stringsAsFactors = TRUE)
Participant_Status <- read.csv("Participant_Status.csv", stringsAsFactors = TRUE)

#Fix date format on INFODT and ENROLL_DATE
MDS_UPDRS_Part_1P <- mutate(MDS_UPDRS_Part_1P, INFODT = as.Date(paste("01/",
as.character(INFODT), sep=""), "%d/%m/%Y"))
Participant_Status <- mutate(Participant_Status, ENROLL_DATE = as.Date(paste("01/",
as.character(ENROLL_DATE), sep=""), "%d/%m/%Y"))

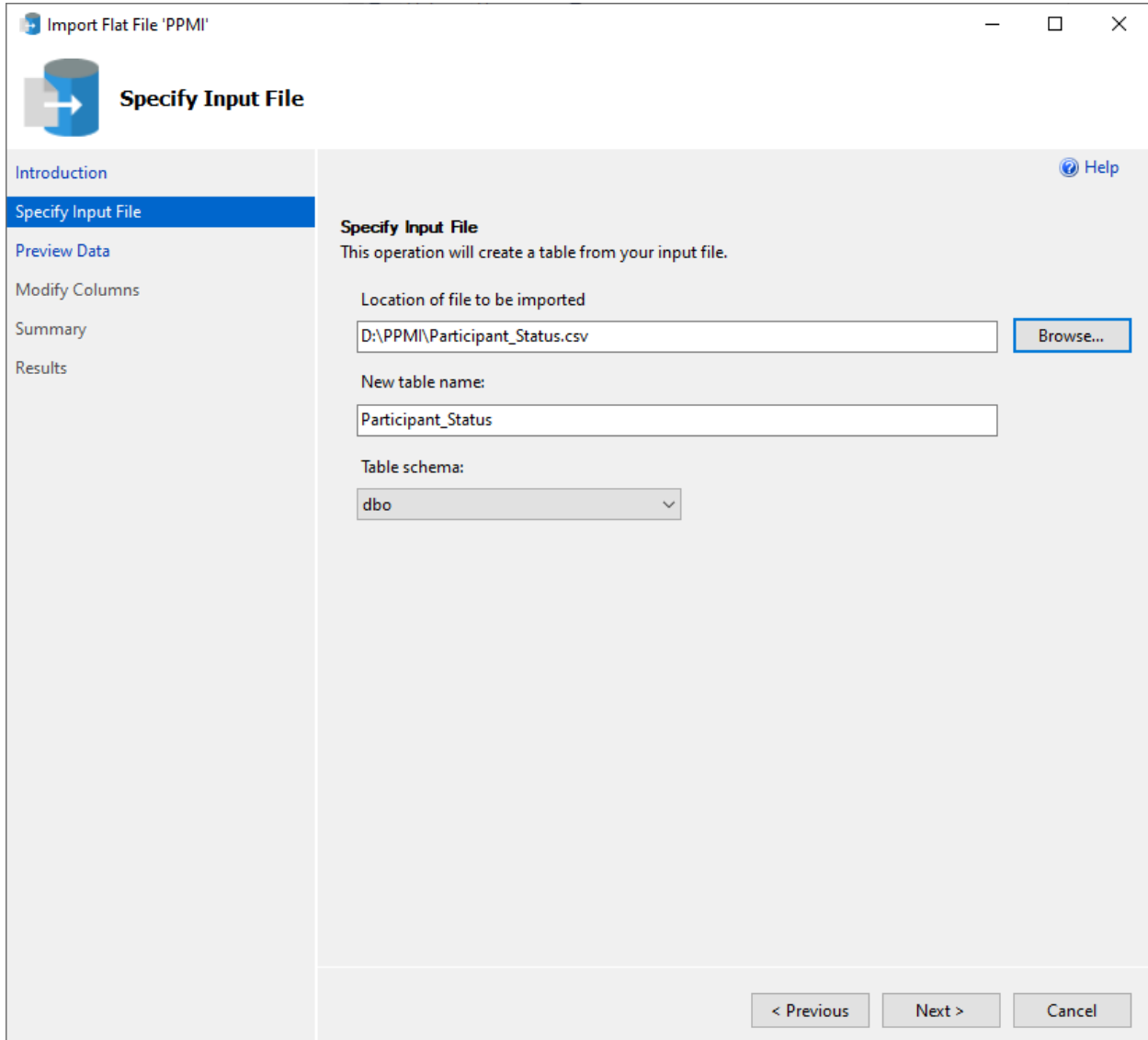
# Create a table that contains the minimum value of the time difference between
age_seconds and INFODT
temp <- select (pulserate %>% rename(PATNO = subject), PATNO, age_seconds) %>%
  left_join (Participant_Status, by = "PATNO") %>%
  select(PATNO, age_seconds, ENROLL_DATE, ENROLL_AGE) %>%
  mutate(diff_ages = (age_seconds - (ENROLL_AGE * 31556952))) %>%
  left_join (select (MDS_UPDRS_Part_1P,PATNO,EVENT_ID,INFODT),by='PATNO') %>%
  mutate (diff_events = difftime(INFODT, ENROLL_DATE, units = 'secs')) %>%
  mutate(delta = abs(diff_ages - diff_events)) %>%
  group_by(PATNO, age_seconds) %>%
  summarize(min_delta_seconds = min(delta))

# Now use this to append the nearest EVENT_ID to the Verily table
pulserate2 <- select (pulserate %>% rename(PATNO = subject), PATNO, age_seconds,
time_ms, time_local, time_day, hourly_mean_pulse_rate, sample_count) %>%
  left_join (Participant_Status, by = "PATNO") %>%
  select(PATNO, age_seconds, ENROLL_DATE, ENROLL_AGE, time_ms, time_local, time_day,
hourly_mean_pulse_rate, sample_count) %>%
  mutate(diff_ages = (age_seconds - (ENROLL_AGE * 31556952))) %>%
  left_join (select (MDS_UPDRS_Part_1P,PATNO,EVENT_ID,INFODT),by='PATNO') %>%
  mutate (diff_events = difftime(INFODT, ENROLL_DATE, units = 'secs')) %>%
  mutate(delta = abs(diff_ages - diff_events)) %>%
  inner_join (temp, by=c('PATNO' = 'PATNO','age_seconds' = 'age_seconds', 'delta' =
'min_delta_seconds')) %>%
  select (PATNO, age_seconds, time_ms, time_local, time_day, hourly_mean_pulse_rate,
sample_count, NEAREST_EVENT_ID = EVENT_ID) %>%
  arrange (PATNO, age_seconds, time_ms)

pulserate2 # The above example can easily be applied to the other Verily tables
```

## APPENDIX D – IMPORTING DATA INTO A RELATIONAL DATABASE

To import the PPMI data into an RDBMS you will have first downloaded the data you need into CSV files. You will then use a data import tool. In the following we use the example of the Import Flat File wizard in MS SQL Server Management Studio loading the **Participant\_Status** table:



The screenshot shows the 'Import Flat File 'PPMI'' wizard window. The 'Specify Input File' step is selected in the left-hand navigation pane. The main area displays the following fields and options:

- Location of file to be imported:** A text box containing 'D:\PPMI\Participant\_Status.csv' and a 'Browse...' button.
- New table name:** A text box containing 'Participant\_Status'.
- Table schema:** A dropdown menu with 'dbo' selected.

At the bottom right, there are three buttons: '< Previous', 'Next >', and 'Cancel'.

In an RDBMS it is important to pay attention to datatypes and you will likely need to override some of the default datatypes assigned, using the data dictionary described in section 3.4 as a guide. We recommend the following to minimize import errors and rejection of records:



Import Flat File 'PPMI'

## Modify Columns

Introduction  
Specify Input File  
Preview Data  
Modify Columns

**Modify Columns**  
This operation generated the following table schema. Please verify if schema is accurate, and if not, please make any changes.

Column Name	Data Type	Primary Key	<input type="checkbox"/> Allow Nulls
PATNO	smallint	<input type="checkbox"/>	<input type="checkbox"/>
COHORT	tinyint	<input type="checkbox"/>	<input type="checkbox"/>
COHORT_DEFINITION	nvarchar(50)	<input type="checkbox"/>	<input type="checkbox"/>
ENROLL_DATE	nvarchar(50)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
ENROLL_STATUS	nvarchar(50)	<input type="checkbox"/>	<input type="checkbox"/>
STATUS_DATE	nvarchar(50)	<input type="checkbox"/>	<input type="checkbox"/>
ENROLL_AGE	float	<input type="checkbox"/>	<input checked="" type="checkbox"/>
INEXPAGE	nvarchar(50)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
AV133STDY	tinyint	<input type="checkbox"/>	<input checked="" type="checkbox"/>
PPMI_ONLINE_ENROLL	bit	<input type="checkbox"/>	<input type="checkbox"/>
CONCOHORT	tinyint	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONCOHORT_DEFINITION	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONLRRK2	bit	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONGBA	bit	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONSNA	tinyint	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONHPSM	nvarchar(1)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONRBD	nvarchar(1)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
PHENOCNV	nvarchar(1)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
DIAG1	nvarchar(1)	<input type="checkbox"/>	<input checked="" type="checkbox"/>

Row granularity of error reporting (performance impact with smaller ranges)

< Previous Next > Cancel

**Annotations:**

- For numeric fields, avoid use of datatypes like bit, smallint and tinyint. Import all numeric fields as either int or float.
- Make all fields nullable
- For text fields, you may get truncation errors if the default fields sizes aren't big enough. Look at the data dictionary and adjust the datatypes accordingly.
- Import date fields using datatype date where appropriate

Note that in some cases, for example MDS-UPDRS measurements, there are missing values. In most cases these will be blank fields, and this is dealt with by making the field nullable as shown above. However, in a few cases you may find text values in numeric fields, typically indicating that a measurement was not possible. To deal with these cases, you have three options:

- Filter these rows out before you import the data (e.g., delete the rows in MS Excel) - Recommended
- Change them to some default value (e.g., using the find and replace function in MS Excel)
- Import them as text fields and deal with them in your analysis

## REFERENCES

- Brennan, L., Siderowf, A., Rubright, J., Rick, J., Dahodwala, N., Duda, J., Hurtig, H., Stern, M., Xie, S., Rennert, L., Karlawish, J., Shea, J., Trojanowski, J. & Weintraub, D. (2016). "The Penn Parkinson's Daily Activities Questionnaire-15: Psychometric properties of a brief assessment of cognitive instrumental activities of daily living in Parkinson's disease." *Parkinsonism Relat Disord.* 2016 Apr;25:21-6. <https://doi.org/10.1016/j.parkreldis.2016.02.020>
- Cock, P., Fields, C., Goto, N., Heuer, M., & Rice, P. (2010). "The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants". *Nucleic Acids Research* 38(6): 1767–1771, <https://doi.org/10.1093/nar/gkp1137>
- Craig, D., Hutchins, E., Violich, I. *et al.* (2021). "RNA sequencing of whole blood reveals early alterations in immune cells and gene expression in Parkinson's disease." *Nat Aging* 1:734–747. <https://doi.org/10.1038/s43587-021-00088-6>
- Doty, R., Shaman, P., Kimmelman, C. & Dann, M. (1984). "University of pennsylvania smell identification test; A rapid quantitative olfactory function test for the clinic." *The Laryngoscope* 94(2):176-178. <https://doi.org/10.1288/00005537-198402000-00004>
- Goetz, C., Tilley, B., Shaftman, S., Stebbins, G., Fahn, S., Martinez-Martin, P., Poewe, W., Sampaio, C., Stern, M., Dodel, R., Dubois, B., Holloway, R., Jankovic, J., Kulisevsky, J., Lang, A., Lees, A., Leurgans, S., LeWitt, P., Nyenhuis, D., Olanow, C., Rascol, O., Schrag, A., Teresi, J., van Hilten, J., LaPelle, N. (2008) "Movement Disorder Society-sponsored revision of the Unified Parkinson's Disease Rating Scale (MDS-UPDRS): Scale presentation and clinimetric testing results." *Movement Disorders* 23(15):2129-2170. <https://doi.org/10.1002/mds.22340>
- Hoehn M. & Yahr M. (1967). "Parkinsonism: onset, progression and mortality." *Neurology* 17:427-442. [http://info-centre.jenage.de/assets/pdfs/library/hoehn\\_yahr\\_NEUROLOGY\\_1967.pdf](http://info-centre.jenage.de/assets/pdfs/library/hoehn_yahr_NEUROLOGY_1967.pdf)
- Johns, M. (1991). "A New Method for Measuring Daytime Sleepiness: The Epworth Sleepiness Scale". *Sleep*, Volume 14, Issue 6, November 1991, Pages 540–545. <https://doi.org/10.1093/sleep/14.6.540>
- Leentjens, A.F.G., Dujardin, K., Pontone, G.M., Starkstein, S.E., Weintraub, D. & Martinez-Martin, P. (2014), "The Parkinson Anxiety Scale (PAS): Development and validation of a new anxiety scale." *Mov Disord.*, 29: 1035-1043. <https://doi.org/10.1002/mds.25919>
- Liao Y., Smyth, G. & Shi, W. (2014). "featureCounts: an efficient general purpose program for assigning sequence reads to genomic features." *Bioinformatics* 30(7):923–930. <https://doi.org/10.1093/bioinformatics/btt656>
- Lipsmeier, F., Taylor, K., Postuma, R. *et al.* (2022) "Reliability and validity of the Roche PD Mobile Application for remote monitoring of early Parkinson's disease". *Sci Rep* 12, 12081. <https://doi.org/10.1038/s41598-022-15874-4>
- Marek, K. *et al.* (2011). "The Parkinson Progressive Marker Initiative (PPMI)." *Progress in Neurobiology* 95(4):629-635. <https://doi.org/10.1016/j.pneurobio.2011.09.005>



- Patro, R., Duggal, G., Love, M. *et al.* (2017). "Salmon provides fast and bias-aware quantification of transcript expression." *Nat Methods* 14:417–419. <https://doi.org/10.1038/nmeth.4197>
- Postuma, R., Arnulf, I., Hogl, B., Iranzo, A., Miyamoto, T., Dauvilliers, Y, *et al.* (2012). "A single-question screen for rapid eye movement sleep behavior disorder: A multicenter validation study". *Movement Disorders* 27(7) 913-916. <https://doi.org/10.1002/mds.25037>
- Purcell S., Neale B., Todd-Brown K., Thomas L., Ferreira M., Bender D., Maller J., Sklar P., de Bakker P., Daly M. & Sham P. (2007) "PLINK: a tool set for whole-genome association and population-based linkage analyses". *Am J Hum Genet.* 2007 Sep;81(3):559-75. <https://doi.org/10.1086/519795>
- Schuff, N., Wu, I., Buckley, S., Foster, E., Coffey, C., Gitelman, D., Mendick, S., *et al.* (2015). "Diffusion imaging of nigral alterations in early Parkinson's disease with dopaminergic deficits." *Movement Disorders* 30(14) 1885-1892. <https://doi.org/10.1002/mds.26325>
- Siderowf, A., Concha-Marambio, L., Lafontant, D., Farris, C. Ma, Y., Urenia, P. *et al.* (2023). "Assessment of heterogeneity among participants in the Parkinson's Progressive Markers Initiative cohort using  $\alpha$ -synuclein see amplification: a cross-sectional study." *Lancet Neurology* 22(5): 407-417. [https://doi.org/10.1016/S1474-4422\(23\)00109-6](https://doi.org/10.1016/S1474-4422(23)00109-6)
- Tomlinson, C., Stowe, R., Patel, S., Rick, C., Gray, R. & Clark, C. (2010). "Systematic Review of Levodopa Equivalency Reporting in Parkinson's Disease." *Movement Disorders* 25(15):2649-2653. <https://doi.org/10.1002/mds.23429>
- Trenkwalder, C., Kohnen, R., Högl, B. Metta, V., Sixel-Döring, F., Frauscher, B., Hülsmann, J., Martinez-Martin, P. & Chaudhuri, K. (2011). "Parkinson's disease sleep scale—validation of the revised version PDSS-2". *Movement Disorders* 26(4) 644-652. <https://doi.org/10.1002/mds.23476>
- Yesavage JA, Brink TL, Rose TL *et al* (1983). "Development and validation of a Geriatric Depression Screening Scale: a preliminary report." *J Psychiatr Res* 17:37–49.