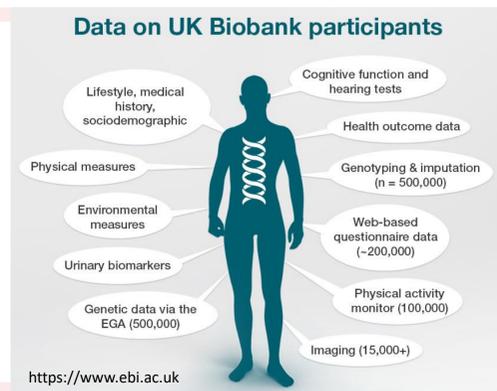# Curation of UK Biobank phenotype data

Emily HM Wong[1], David Ficenec[2], David Merberg[1], Yulia Skovpen[3], Alena Fedarovich[3], Tyler Kolisnik[3], Oksana Tyurina[3], Christine Loreth[3], Mahyar Sabripour[1], Bin Li[1], Yan Ge[1], Tania Khasanova[3], Julie Bryant[3], Sandor Szalma[1]

[1]Data Science Institute and [2]Information Technology, Takeda Pharmaceuticals, USA; [3] Rancho Biosciences, USA

## Introduction

- UK Biobank collects vast amounts of phenotype data and provides a valuable resource for studying the complex relationships between the human genome and phenome.

- As part of the UKB-Pharma consortium, we are building a pipeline to process available UK Biobank data. To facilitate downstream analyses, we and Rancho Biosciences first extensively curated phenotype data from all UK Biobank participants.


Data on UK Biobank participants
https://www.ebi.ac.uk

## Goal

Curate the UK Biobank data by aggregating and integrating the disparate datasets into a single comprehensive resource
- Flat files (PheWAS ready)
- tranSMART ready files

## Curation Process: Overview

**Beginning of 2018:** 502,616 subjects
3,390 data fields

**Data Inspection** → **Harmonization** → **Ontology Mapping** → **Preparation of PheWAS & tranSMART files**

- Data inspection
- UKB coding files update
- Creation of decoded version
- Preparation of supporting files for PheWAS
  - e.g. code reassignment file, field ID–coding file mapping, summary statistics file, etc.

- Harmonization of inconsistent data in multiple instances
- ~100 harmonization rules have been proposed and implemented

- RXNORM
- SNOMED
- ICD10
- MeSH

- Data preparation for PheWAS analysis
- Data conversion into a human-readable format to assist the translational biomedical research
- Upload of UKB phenotype data onto tranSMART

## UK Biobank in tranSMART

- The tranSMART platform is a knowledge management architecture for translational medicine.
- It allows for mining and searching multi-model data, including clinical and biomarker data.



## Identifying drug targets and pleiotropy

- Genomics data and the curated phenotype data were loaded onto a Spark and Hail framework on Amazon Databricks.
- This allows us to efficiently perform PHEWAS and GWAS to identify drug targets and pleiotropy in three major therapeutic areas – gastroenterology, neuroscience, and oncology.