

THE PSYCOURSE STUDY –  
FAIR LONGITUDINAL, MULTI-LEVEL GENOMIC AND PHENOMIC DATA  
FOR MACHINE LEARNING



Urs Heilbronner<sup>1</sup>; Monika Budde<sup>1</sup>; Daniela Reich-Erkelenz<sup>1</sup>; Heike Anderson-Schmidt<sup>2</sup>; Janos Kalman<sup>1,3,4</sup>; Fanny Senner<sup>1,3</sup>; Kristina Adorjan<sup>1,3</sup>; Eva C. Schulte<sup>1,3</sup>; Ashley L. Comes<sup>1,4</sup>; Sergi Papiol<sup>1,3</sup>; Till F. M. Andlauer<sup>5</sup>; Marcella Rietschel<sup>6</sup>; Markus M. Nöthen<sup>7,8</sup>; Peter Falkai<sup>3</sup>; Thomas G. Schulze<sup>1</sup>

<sup>1</sup>Institute of Psychiatric Phenomics and Genomics (IPPG), University Hospital, LMU Munich, Munich, Germany; <sup>2</sup>Department of Psychiatry and Psychotherapy, University Medical Center Goettingen, Goettingen, Germany; <sup>3</sup>Department of Psychiatry and Psychotherapy, University Hospital, LMU Munich, Munich, Germany; <sup>4</sup>International Max Planck Research School for Translational Psychiatry, Max Planck Institute of Psychiatry, Munich, Germany; <sup>5</sup>Department of Neurology, Klinikum rechts der Isar, School of Medicine, Technical University of Munich, Munich, Germany; <sup>6</sup>Department of Genetic Epidemiology in Psychiatry, Central Institute of Mental Health, Medical Faculty Mannheim, University of Heidelberg, Mannheim, Germany; <sup>7</sup>Institute of Human Genetics, University of Bonn School of Medicine & University Hospital Bonn, Bonn, Germany; <sup>8</sup>Department of Genomics, Life & Brain Center, University of Bonn, Bonn, Germany



INTRODUCTION

The integration of multiple levels of genetic information, longitudinal deep phenotyping, and machine learning holds promise to substantially increase our knowledge on mental disorders. However, well-annotated longitudinal datasets on comprehensively characterized individuals are rare. Also, the workload involved in pre-processing such large datasets can be substantial. Here, we present a phenotype dataset of the longitudinal PsyCourse Study (Budde et al., 2018), that has been extensively annotated and formatted for machine-learning analyses. Metadata of this dataset are publicly available (Heilbronner et al., 2020). Biological OMICS data are available for analyses (genomics, transcriptomics, proteomics, lipidomics) in subsets of participants (see below), these will be extended for additional cross-sectional and longitudinal analyses in the future.

METHODS

The latest version of the phenotype dataset (PsyCourse 4.1) contains n=1,320 clinical (affective-to-psychotic spectrum) and n=466 control participants, who were assessed at the first of four planned study visits (visit 2: n=788/288, visit 3: n=661/280, visit 4: n=589/251), each separated by approximately six months. A comprehensive test battery surveying demographics, psychiatric history, medication, substance abuse, diagnosis, cognition, psychiatric rating scales, and several questionnaires, was assessed. At the time of writing, the following biological data exist: SNP array (n=1,446, all diagnoses and controls, PsychChip), epigenetics (n=96 bipolar disorder, two measurement points, EPIC array), exome sequencing (n=104, bipolar disorder), small-RNAome sequencing (n=1,322, first study visit, all diagnoses and controls), mRNA sequencing (n=538, first study visit, schizophrenia spectrum), plasma proteome profiling (n=220, multiple diagnoses and study visits), serum protein profiling (n= 222, multiple diagnoses and study visits), and lipidomics (n=242 bipolar, n=186 schizophrenia, n=192 controls, multiple study visits). Importantly, missing phenotype data, ubiquitous in longitudinal studies, have been re-coded to identify structurally missing data. Researchers can access the data by filling out a short research proposal (see URL in „Results“).

ABSTRACT

- The PsyCourse Study is a longitudinal investigation of the affective-to-psychotic spectrum, including healthy controls, that combines deep phenotyping and collection of DNA, RNA, plasma, and serum (qr\_1)
- Multiple OMICS datasets of PsyCourse participants already exist
- The PsyCourse phenotype dataset has been pre-formatted for machine learning analyses, and metadata are publicly available (qr\_2)
- Data can be accessed by submitting a short research proposal (qr\_3)
- We encourage researchers to collaborate with us!

RESULTS

The framework of FAIR data (Wilkinson et al., 2016) describes a measurable set of principles that have been created to foster the re-use of scientific data (Figure 1). The PsyCourse dataset fulfills these criteria by being:

- *Findable* (public metadata),
- *Accessible* (via research proposals),
- *Interoperable* (available in wide and long format R files, also as a text file), and thus
- *Reusable*

For details, please refer to the PsyCourse Open Science website: [psycourse.de/openscience-en.html](https://psycourse.de/openscience-en.html).

DISCUSSION

While not in the public domain, PsyCourse data can be re-used by bona-fide researchers world-wide. Accepted research proposals are posted on our website, giving a transparent overview on ongoing research projects. We encourage researchers to collaborate with us!

REFERENCES

Budde et al. (2018) A longitudinal approach to biological psychiatric research: The PsyCourse study. Am J Med Genet B Neuropsychiatr Genet. DOI: 10.1002/ajmg.b.32639  
Heilbronner et al. (2020). The PsyCourse Codebook, Version 4.1. Open Data LMU. DOI: 10.5282/ubm/data.199.  
Wilkinson et al. (2016) The FAIR Guiding Principles for scientific data management and stewardship. Scientific Data. DOI: 10.1038/sdata.2016.18

GRANTS

TGS: DFG Grants SCHU 1603/5-1 and SCHU 1603/7-1; BMBF Grants IntegraMent and BipoLife; Dr. Lisa-Oehler-Foundation (Kassel, Germany).  
HB: DFG Grants BI 576/5-1 and Research Training Group “Scaling Problems in Statistics” RTG 1644.

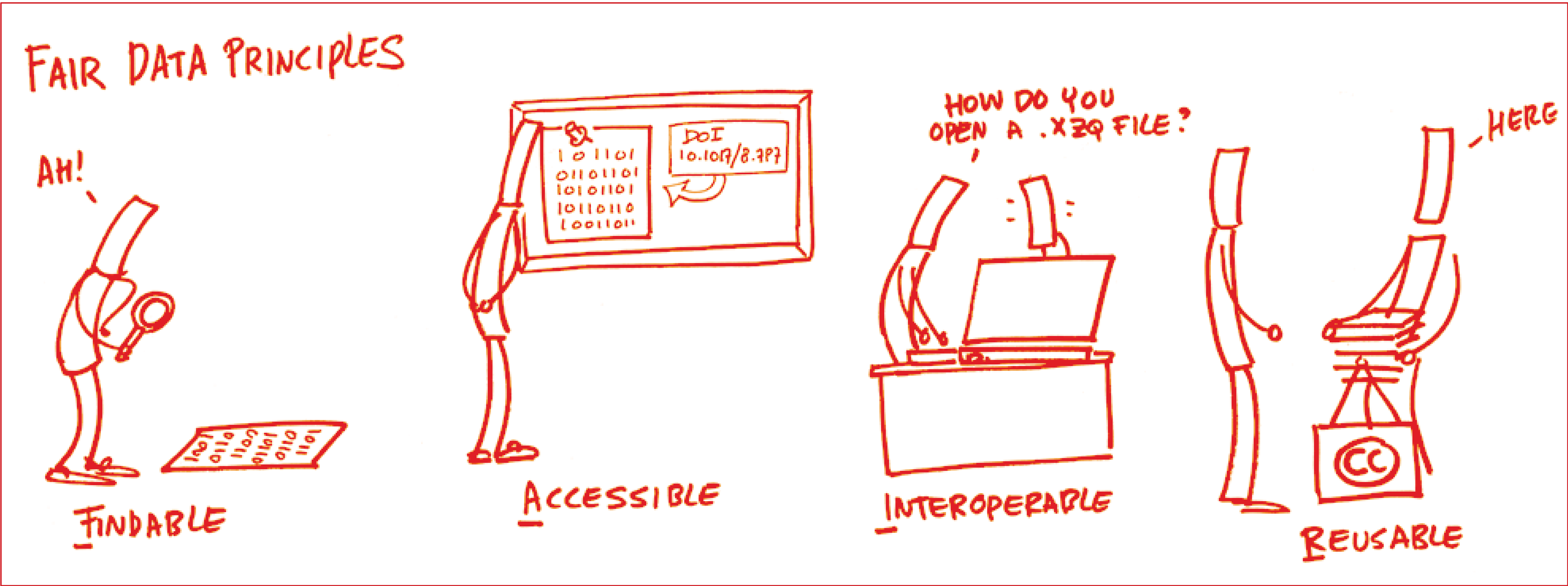


Figure 1. The FAIR Data Principles (<https://book.fosteropenscience.eu/>).

