# OPEN RESEARCH SOFTWARE INFRASTRUCTURE IN NEURO-MEDICINE

Adina Wagner

mas.to/@adswa

Psychoinformatics lab,
Institute of Neuroscience and Medicine, Brain & Behavior (INM-7)
Research Center Jülich

JÜLICH
Forschungszentrum

Slides: DOI 10.5281/zenodo.10149349 (Scan the QR code)
files.inm7.de/adina/talks/html/zimannheim.html

# OPEN SCIENCE AND OPEN SOFTWARE GO HAND IN HAND

- Science has **specific requirements**; research software from within science ("from scientists, for scientists") can fulfill them. Open formats, protocols, and code allow re-use, interoperability, and customization.

- Open and reproducible science has specific needs for **transparency**: Open source software provides the necessary auditability.

- Creating software becomes **increasingly possible for scientists**: The San Francisco Declaration on Research Assessment (DORA; signed by FZJ), the Agreement on Reforming Research Assessment (CoARA), and the DFG recognize software as academic output.

# THE INSTITUTE FOR NEUROSCIENCE & MEDICINE (INM-7)

# THE INSTITUTE FOR NEUROSCIENCE & MEDICINE (INM-7)

- Interdisciplinary institute with 11 research groups
- Research foci:
  - Infrastructure and method development: Digital biomarker, machine learning, meta analysis, research data management
  - Basic research in human brain mapping: Connectomics, genetic gradients, in-vivo brain mapping, multimodal integration
  - AI Applications in medical research: Cognition, Personality, Aging & neurodegenerative disease, Schizophrenia
  - Ethical implications of medical AI: Bias in AI applications, medical AI and society, individualized predictions

# SOFTWARE @ INM-7

- The institute has a history of open source software, starting with the SPM Anatomy Toolbox (Eickhoff, 2005)
- Multiple groups develop and maintain open source research software for their respecitve subdomain
- Recent integration efforts connect our open software stack to open research software infrastructure for neuro-medicine

![DataLad logo]

- Domain-agnostic data management tool **(command-line** + **graphical user interface**), built on top of Git & Git-annex
- 10+ year open source project (100+ contributors), available for all major OS
- Born from rethinking data:
  - Just like code, **data is not static**.
  - Just like code, **data is subject to collaboration**. Stream-lined workflows for sharing and collaborating should be possible, mirroring those in software development.
  - **Provenance** of data is essential for reproducible, trustworthy, and FAIR science
  - Flexibility and **interoperability with existing tools** is the key to sustainability and ease of use

- Domain-agnostic **command-line tool** (+ **graphical user interface**), built on top of Git & Git-annex
- 10+ year open source project (100+ contributors), available for all major OS
- Major features:
  **Version-controlling arbitrarily large content**
    Version control data & software alongside to code!
  **Transport mechanisms for sharing, updating & obtaining data**
    Consume & collaborate on data (analyses) like software
  **(Computationally) reproducible data analysis**
    Track and share provenance of all digital objects
  **(... and much more)**

# EXHAUSTIVE TRACKING OF RESEARCH COMPONENTS

text document
source code
binary file
regular file
file reference:
identity and
availability

Well-structured datasets (using community standards), and portable computational environments — and their evolution — are the precondition for reproducibility

```
# turn any directory into a dataset
# with version control

% datalad create <directory>
```

```
# save a new state of a dataset with
# file content of any size

% datalad save
```

# CAPTURE COMPUTATIONAL PROVENANCE

text document

source code

binary file

regular file

file reference:
identity and
availability

Which data was needed at which version, as input into which code, running with what parameterization in which compuational environment, to generate an outcome?

```
# execute any command and capture its output
# while recording all input versions too

% datalad run --input ... --output ... <command>
```

# EXHAUSTIVE CAPTURE ENABLES PORTABILITY



Precise identification of data and computational environments combined with provenance records form a comprehensive and portable data structure, capturing all aspects of an investigation.

```
# transfer data and metadata to other sites and services
# with fine-grained access control for dataset components


% datalad push --to <site-or-service>
```

# REPRODUCIBILITY STRENGTHENS TRUST



Outcomes of computational transformations can be validated by authorized 3rd-parties. This enables audits, promotes accountability, and streamlines automated "upgrades" of outputs
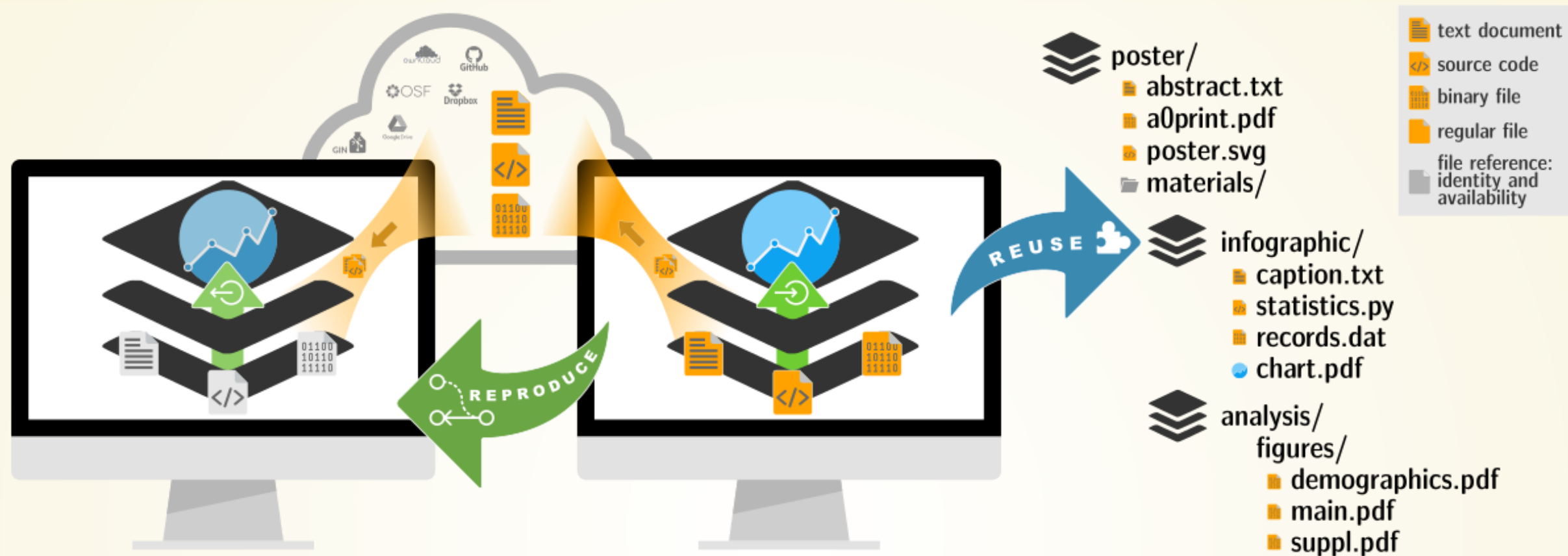
```
# obtain dataset (initially only identity,
# availability, and provenance metadata)


% datalad clone <url>
```

```
# immediately actionable provenance records
# full abstraction of input data retrieval


% datalad rerun <commit|tag|range>
```
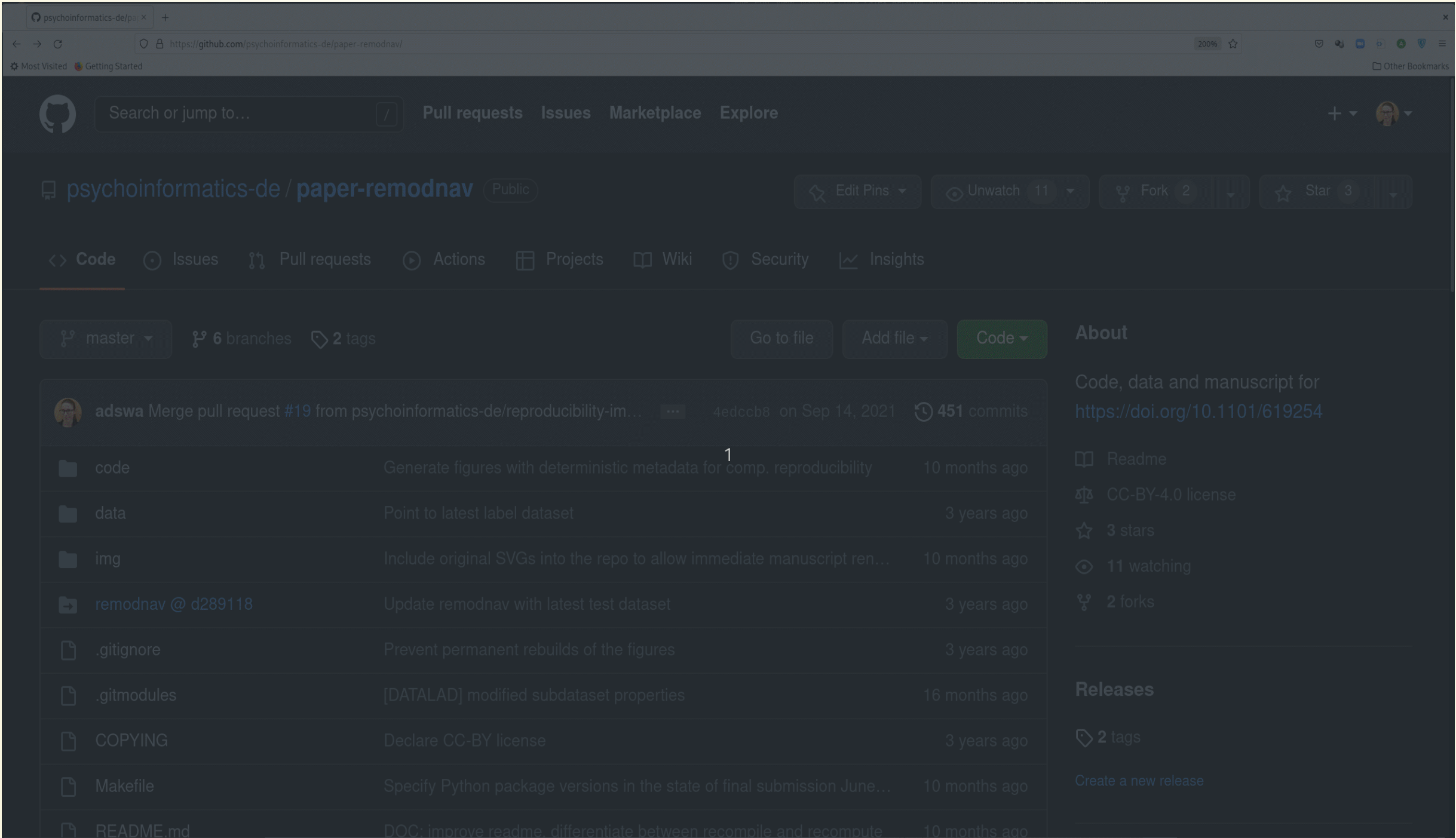
# ULTIMATE GOAL: (RE-)USABILITY



Verifiable, portable, self-contained data structures that track all aspects of an investigation exhaustively can be (re-)used as modular components in larger contexts — propagating their traits

```
# declare a dependency on another dataset and
# re-use it a particular state in a new context


% datalad clone -d <superdataset> <url> <path-in-dataset>
```

# DATALAD USECASES

# ACKNOWLEDGEMENTS

## Funders



NSF 1429999

germany-usa

SPONSORED BY THE

Federal Ministry of Education and Research

BMBF 01GQ1411

JÜLICH Forschungszentrum

DFG

EUROPEAN UNION
European Regional Development Fund

cbbs
center for behavioral brain sciences

SACHSEN-ANHALT
Ministerium für
Wissenschaft und Wirtschaft

**DataLad software & ecosystem**

- Psychoinformatics Lab, Research center Jülich
- Center for Open Neuroscience, Dartmouth College
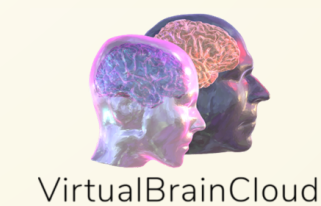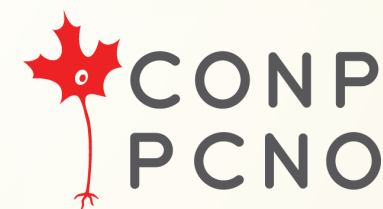- Joey Hess (git-annex)
- *>100 additional contributors*

## Collaborators



Human Brain Project

CONP PCNO

OpenNEURO

eBRAIN Health

MOTOR SFB

brainlife.io

cbrain

VirtualBrainCloud

# JTRACK: DIGITAL BIOMARKERS FROM YOUR SMARTPHONE

- **Objective**: Close monitoring of patients/participants in non-clinical settings
- Modern smartphones contain a variety of sensors for passive monitoring and active acquisition:
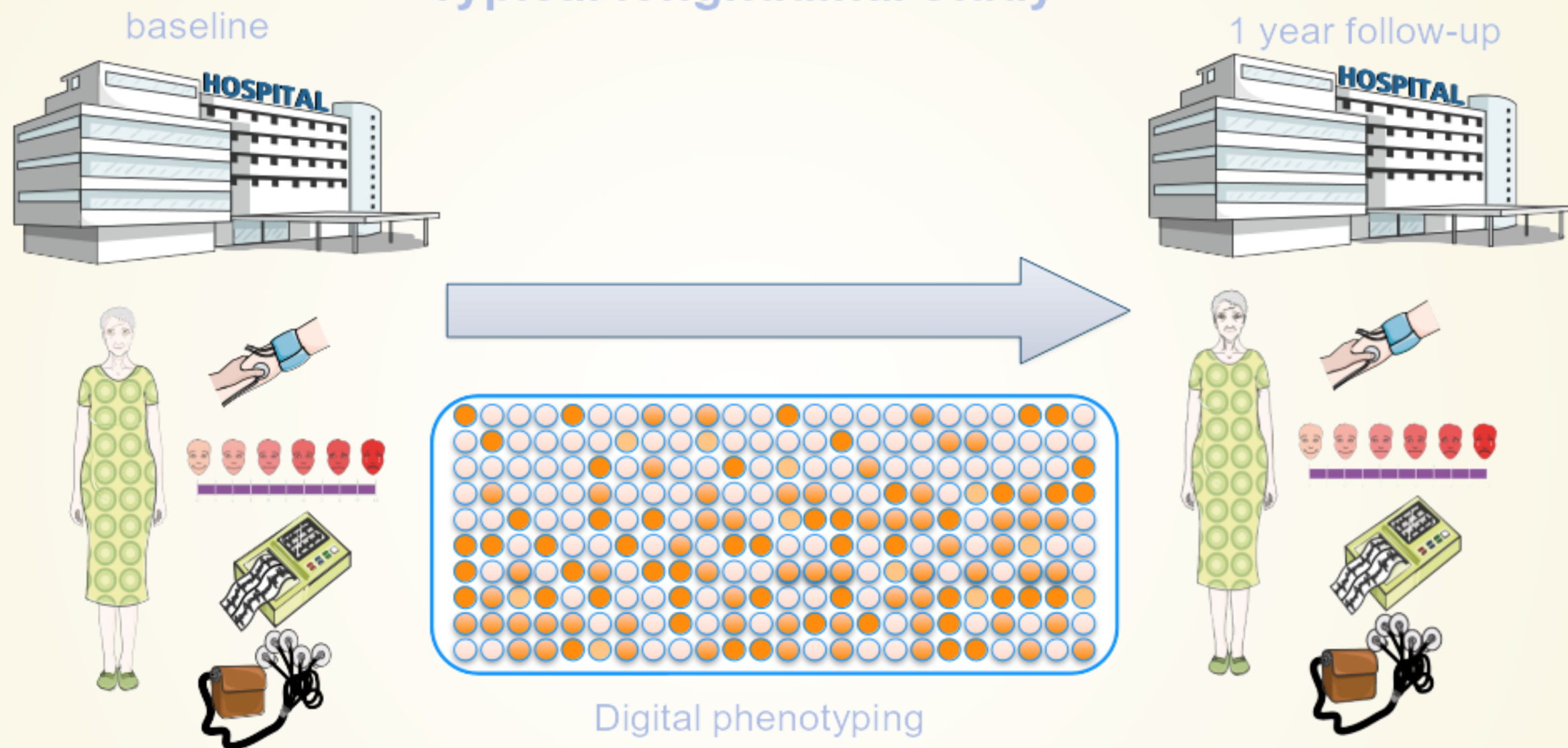
Any modern smart device



- gyroscope
- accelerometer
- location
- human activity recognition
- application usage
- screen time
- microphone

# JTRACK



Typical longitudinal study

baseline

1 year follow-up

HOSPITAL

HOSPITAL

Digital phenotyping

Flexible components for different users:
- **JTrack Social**: Smartphone app for participants
- **JTrack EMA**: Smartphone app for participants
- **JDash**: Monitoring and analytics tool for study owners

# JTRACK COMPONENTS: JTRACK SOCIAL

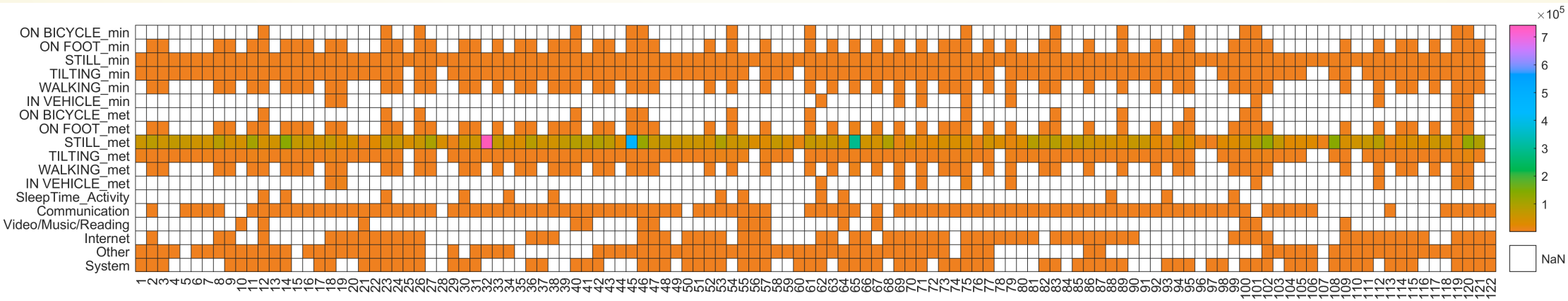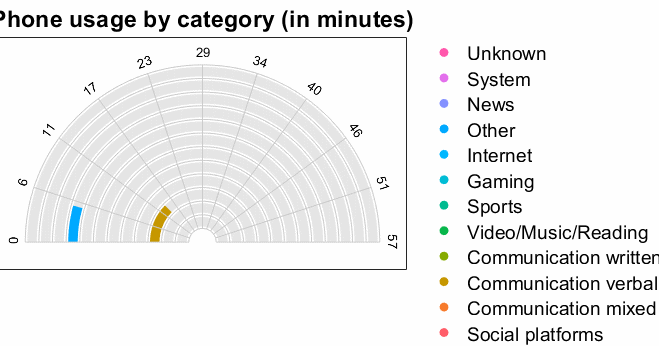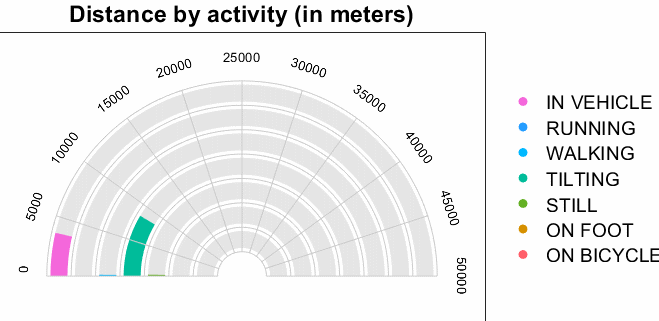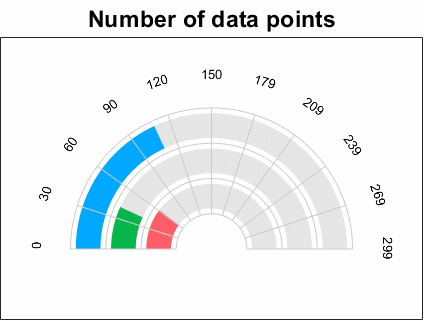Smartphone App for active labeling and passive monitoring

Available for Android + IPhone

- Sensor data (passive collection mode default: Accelerometer & Gyroscope)
- Application usage statistics
- Human activation recognition (e.g., walking, running, driving)
- Location information (anonymized)
- Active recording, e.g., free-speech generation tasks

# JTRACK COMPONENTS: JTRACK SOCIAL

# JTRACK COMPONENTS: JTRACK EMA

Smartphone App for Ecological Momentary Assessment <sup>Available for Android + IPhone</sup>
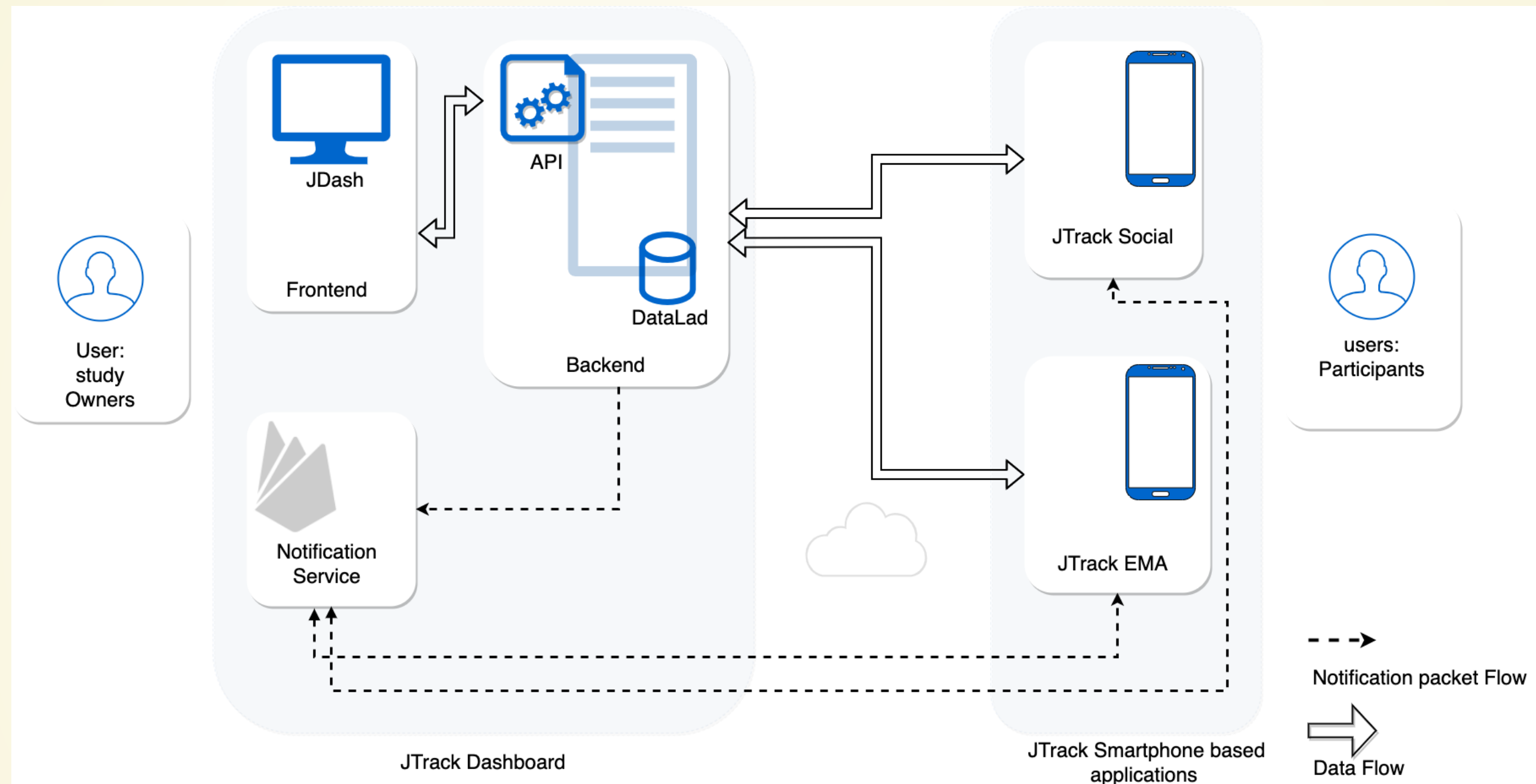
- Binary Questions
- Date and Time Questions
- Sliding Questions
- Multiple/Single coice questions

# JTRACK COMPONENTS: JDASH

Dashboard for Study Administration

- Investigator's study & user management
- Data Quality Control
- Notification Center

# BEHIND THE SCENES



- Servers in Germany
- Data versioning via DataLad
- Authenticated data access via JDash
- Data transfer via HTTPS

# INVESTIGATOR'S / PARTICIPANT'S POINT OF VIEW

- Install JTrack
- Scan QR code and give permissions to App
- Resume daily life

- Install JTrack with Participant
- Provide study- and subject-specific QR code from JDash
- Monitor study and communicate with participants via JDash

# ADVANTAGES

- Easy to deploy and free environment for collection of real world data (RWD) basically at no cost
- Standardized data collection across centers
- High-density longitudinal data with fully customizable data collection
- Opportunity for citizen science

# ACKNOWLEDGEMENTS

**Publications**:
- JTrack Social: Sahandi Far et al. 2021
- JTrack EMA: Sahandi Far et al. 2023

**Contact**:
- via JDash: jdash.inm7.de

- **JTrack Hour** (open to everyone)
  Every second Tuesday at 1PM (even-numbered calendar weeks) -
  Join the Zoom meeting!

**Team**:

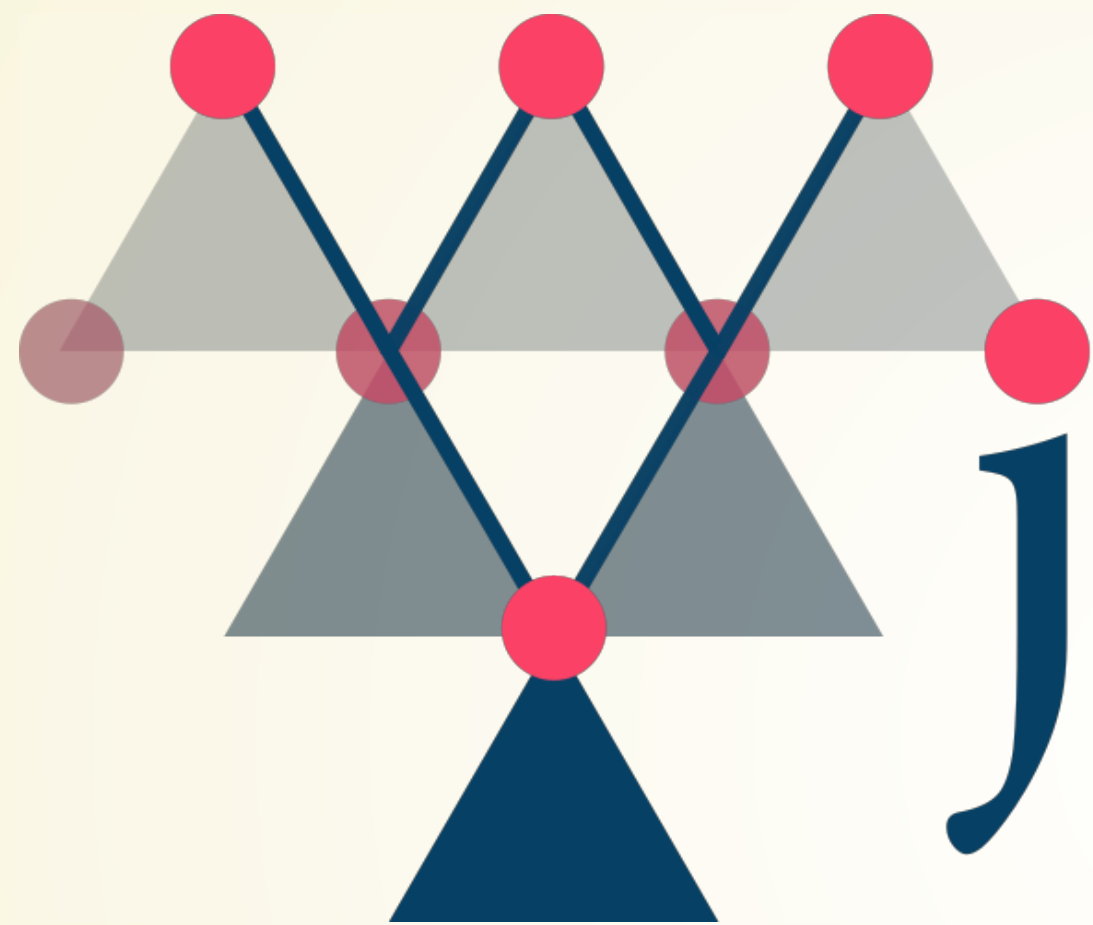**JTrack is a team effort**

Mamaka Narava   Jona Fischer   Michael Stolz
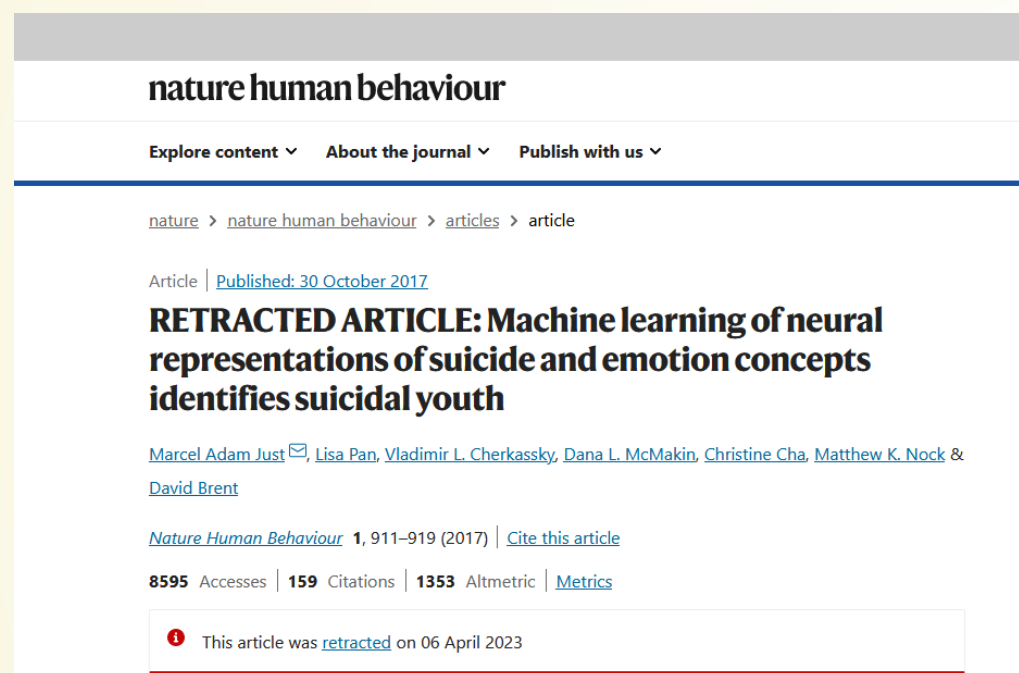
Mehran Sahandi Far   Juergen Dukart

keep calm and
**julearn**
run_cross_validation

- Open source Python library for easy-to use ML-pipelines, built upon scikit-learn
- Domain-general, but aims to simplify entry into ML for domain scientists with built-in guarantees against most common pitfalls:
  - Data leakage
  - Overfitting of hyperparameters
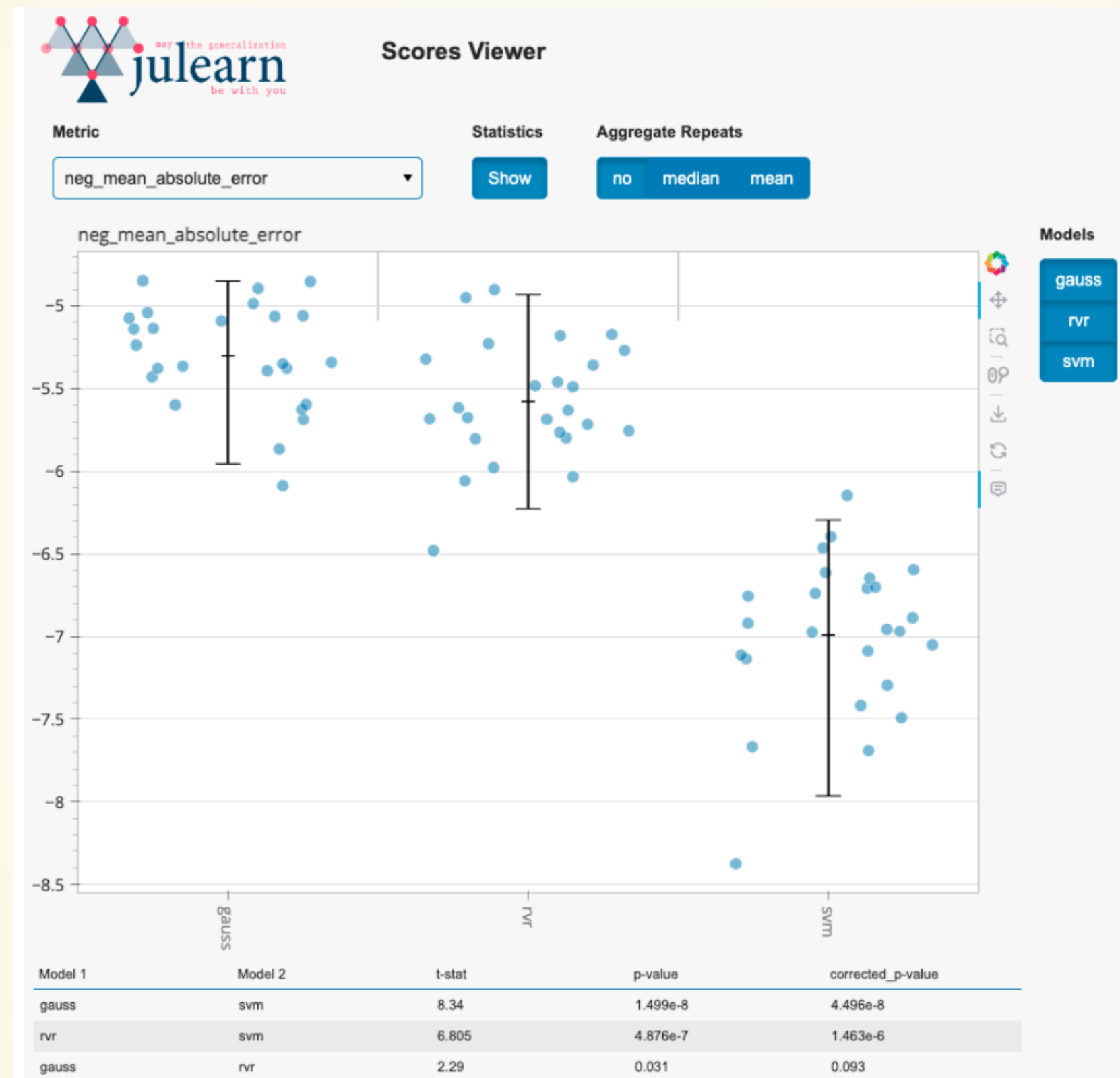
**The problem**: Expensive AI mistakes

**A solution**: User-friendly solutions to common complex use cases



- Simplifies common use cases for supervised ML pipelines, with feature such as:
    - Automatic usage of nested cross-validation for proper evaluation in hyperparameter tuning
    - Preprocessing based on feature types, incl. confound removal
    - Built-in visualization for model inspection and comparison
- Plug-and-play with scikit-learn transformers

# VISUALIZATION

Interactive "Scores Viewer" for easier model comparison

# JULEARN VS SCIKIT-LEARN

## Simple CV pipeline

```python
from julearn import run_cross_validation
run_cross_validation(
    X=X, y=y, data=data,
    preprocess=["zscore"], model="svm",
    problem_type="classification",
    X_types={"continuous": X} # X_types optional here
```
copy

```python
from sklearn.model_selection import cross_validate
from sklearn.svm import SVC # SVR in case of regression
from sklearn.preprocessing import StandardScaler
from sklearn.pipelines import make_pipeline

pipeline = make_pipeline(StandardScaler(), SVC())
cross_validate(X=data.loc[:,X], y=data.loc[:,y], estimator=pipeline)
```
copy

# JULEARN VS SCIKIT-LEARN

## Nested CV with hyperparameter tuning

```python
from julearn import run_cross_validation, PipelineCreator
creator=PipelineCreator(problem_type="classification")
creator.add("zscore", with_mean=[True, False])
creator.add("pca", n_components=2)
creator.add("svm", C=[1,2], degree=[3,4])

# X_types optional
run_cross_validatoin(
    X=X, y=y, data=data, model=creator, X_types={"continuous": X})
```

`copy`

```python
from sklearn.model_selection import cross_validate, GridSearchCV
from sklearn.svm import SVC # SVR in case of regression
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.pipelines import make_pipeline

pipeline = make_pipeline(StandardScaler(), PCA(), SVC())
param_grid = {
    "standardscaler__with_mean": [True, False],
    "pca__n_components": [2],
    "scv__C": [1,2],
    "svc__degree": [3, 4]
}
grid_pipeline = GridSearchCV(estimator=pipeline, param_grid=param_grid)
cross_validate(X=data.loc[:,X], y=data.loc[:,y], estimator=pipeline)
```
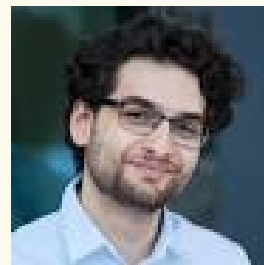
`copy`

# ACKNOWLEDGEMENTS

- Preprint: Hamdan et al., 2023
- Documentation: juaml.github.io/julearn
- Source Code: github.com/juaml/julearn
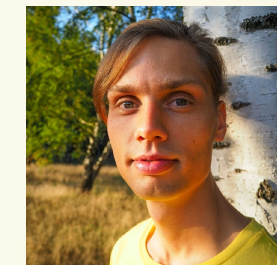


Fede Raimondo  Sami Hamdan  Kaustubh Patil  Shammi More  Vera Komeyer  Synchon Mandal  Leonard Sasse

**ML Hours** (open to everyone)
Consultancy on Machine-Learning, every second Thursday, 2-4pm
Chat: https://matrix.to/#/#ml:inm7.de

# EVEN BETTER TOGETHER

# THE ABCD-J PLATFORM

## AN OPEN SOURCE PLATFORM FOR DIGITAL BIOMARKER FOR NEURO-MEDICINE IN NRW

A collaboration between clinical, academic, and industry partners:
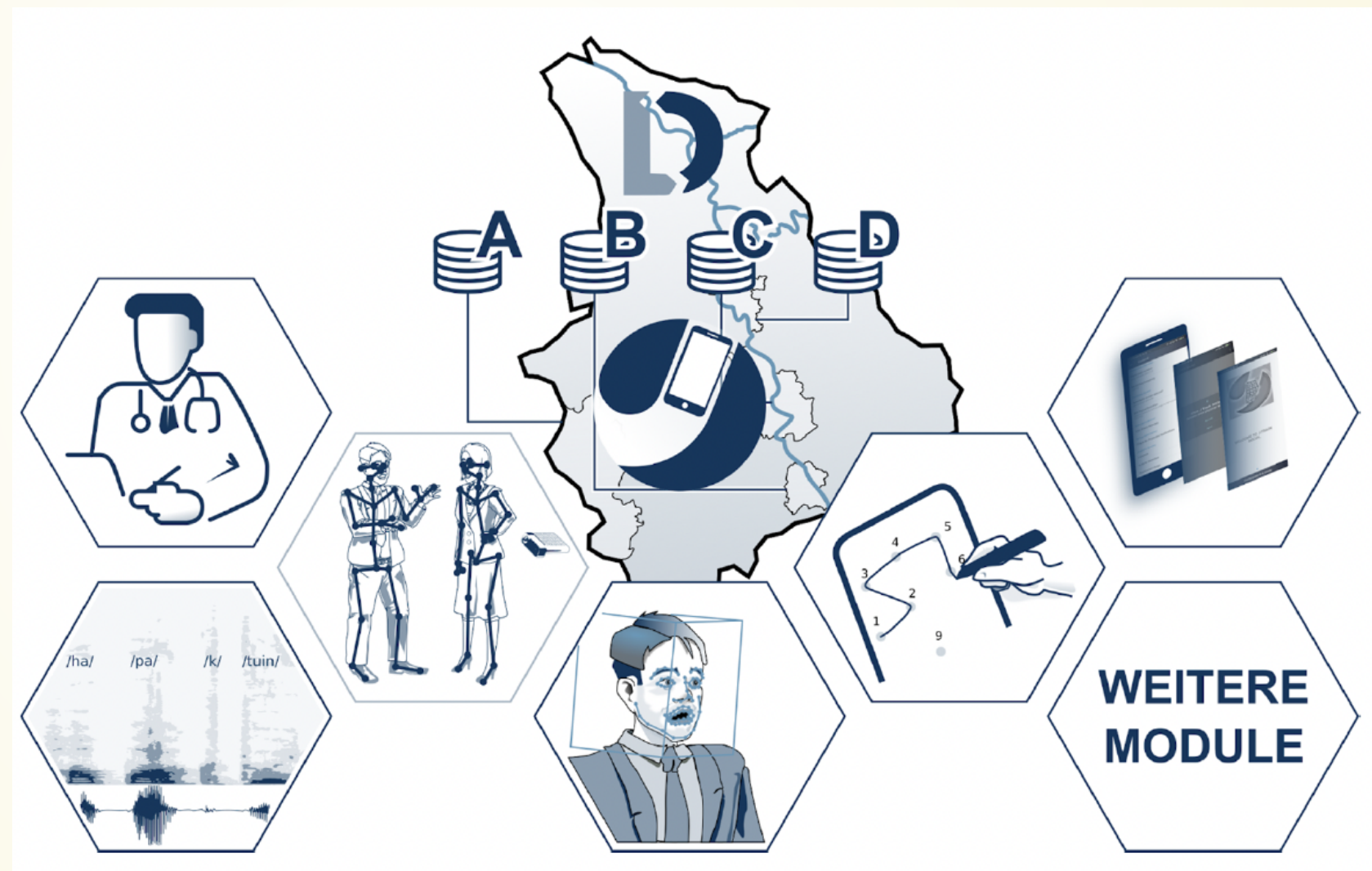
Research Center
Jülich

RWTH Aachen
University Bonn
University Cologne
HHU Düsseldorf

LVR Clinics
DZNE

PeakProfiling
CanControl
IXP

(open to future
additions)

# GOALS

- **Social:** Promote and facilitate collaboration between multiple centers
- **Technical:** Accelarate research through homogenization of workflows and processes, with emphasis on digital biomarker development; Elevate existing open technical solutions for research practice adoption
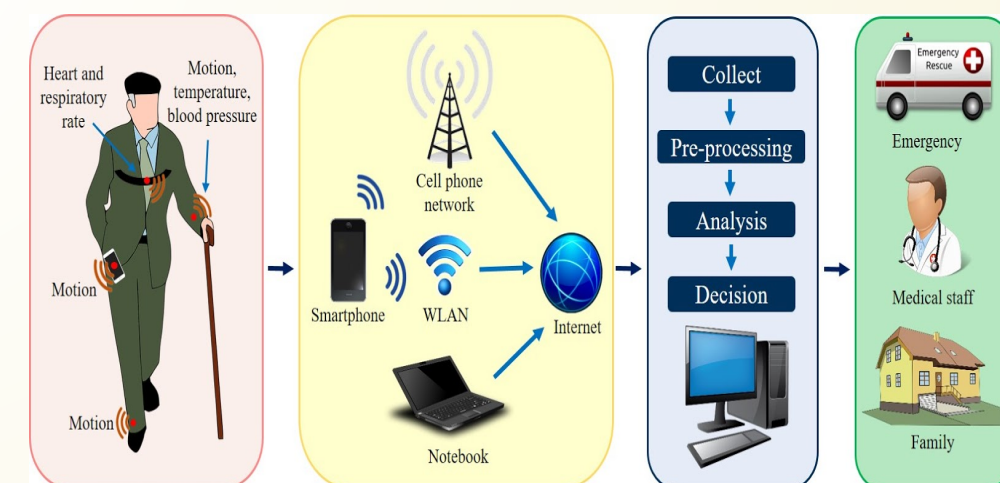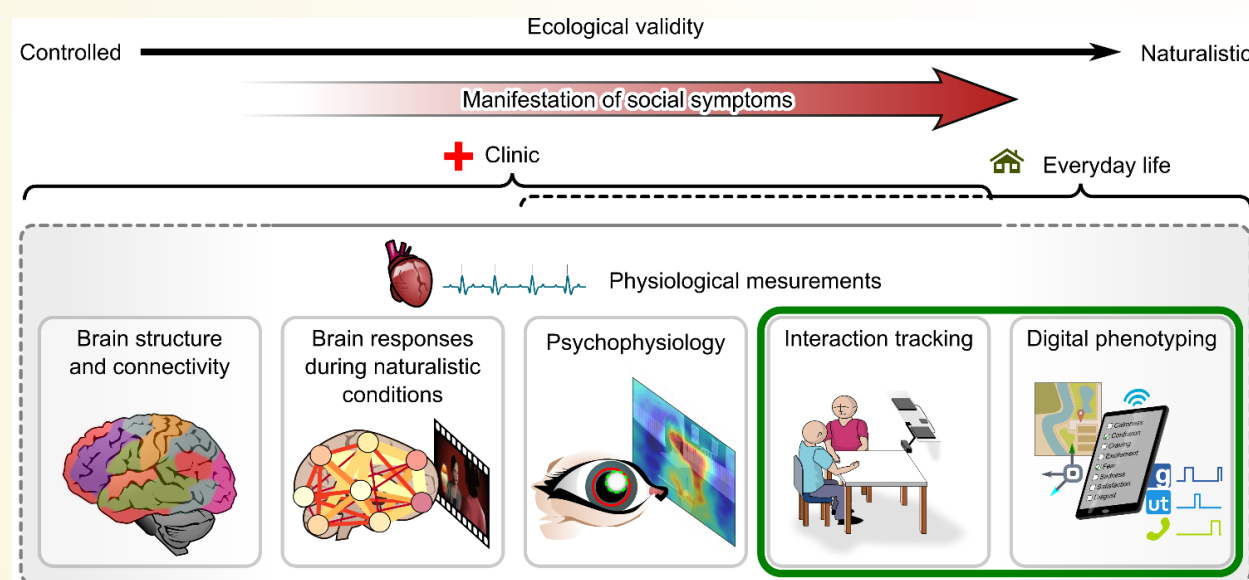
# OPEN RESEARCH INFRASTRUCTURE

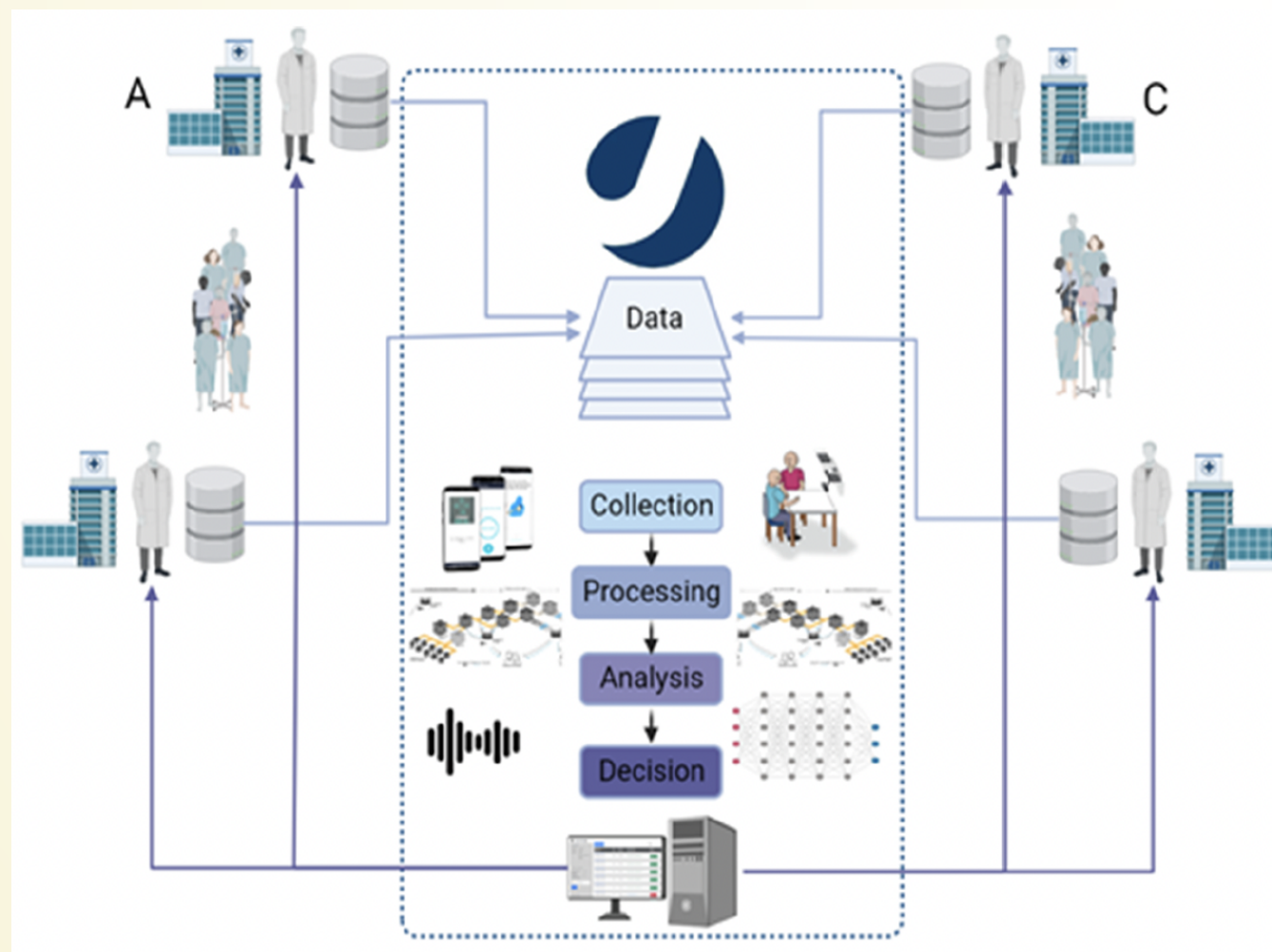| Clinicians' point of view | Patients' point of view |
|---|---|
| • "Deep phenotyping": ecologically valid, multimodal data<br>• Decentral data acquisition, standardized and reproducible<br>• Focus on patient well-being and optimal treatment | • Accurate diagnosis and optimal treatment<br>• Strict data protection<br>• Individual patient is central<br>• Minimal disturbance in daily life |

# OPEN RESEARCH INFRASTRUCTURE

**Research data management's point of view**

- Technical solutions exists, need elevation for research practice adoption
- Homogenization of workflows and processes fosters collaboration across sites
- Decentralized approach with centralized services and web-based multi-center integration

# FRONT-END AND BACK-END

- JTrack for decentral, ecologically valid acquisitions, complementing in-clinic assessments
- JDash for study management, participant management, and analytics overview (derived study data at subject & group level)
- Central data overview and analytics at FZJ
  - Provenance-tracked analysis and modeling
  - Automated meta-data extraction for data discoverability
  - Result overview for clinical decision making

# OPPORTUNITIES

**Software improves with its use cases**

- JTrack integration into different types wearables
- JTrack integration of cognitive tasks and feedback to participants
- Julearn integration into JDash
- More meta-data extractors for DataLad
- ...

# CURRENT (FIRST) STEPS

**Data cataloging**
- Leveraging legacy data via data census and meta-data catalog
  - improved discovery without direct data transfer
  - homogenization of access request procedures
  - establishing a legal basis for (re-)use
  - Example: data.sfb1451.de
- Demonstrator for data infrastructure based on §21 data (standardized and anonymized performance data of hospitals, legally required, submitted yearly to InEK by all hospitals)

**Feasibility/Proof-of-concept study**
- Recommendations for common digital tools and workflows for common tasks and processes
- Selection of digital measures from clinical routine for first trials

# SUMMARY

- Open Source (Research) Software aids in various domain-general or -specific applications.
- Open Science needs Open Source: For transparency and reproducibility, for science-specific requirements, for open formats, for re-use, and to enable interoperability across tools.
- Collaboration across clinical and research settings is a technical, social, and legal challenge. Technical solutions won't save us alone, but they are a good first step.

# THANKS!

Questions?

Inputs, Resources, Synergies