



AN END-TO-END RESEARCH INFRASTRUCTURE FOR GENERATING AND SHARING DIVERSITY-AWARE DATA

DHAI WORKSHOP

Munich, 26th June 2023

MATTEO BUSSO¹, Ronald Chenu² and
Amalia de Götzen²

¹University of Trento (Italy), ²Aalborg University (Denmark)

WWW.INTERNETOFUS.EU

Index



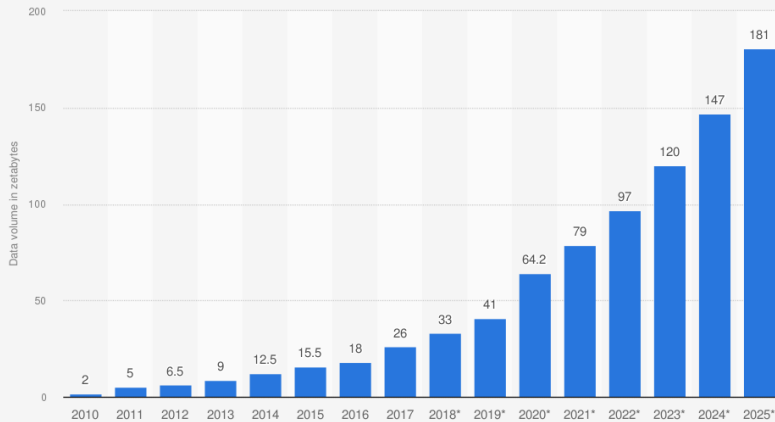
1. Introduction
2. Why a diversity aware RI?
3. What type of RI?
4. Current platforms
5. Benchmark
6. A former approach
7. Limits

1. Introduction

2. Why a diversity aware RI?
3. What type of RI?
4. Current platforms
5. Benchmark
6. A former approach
7. Limits

A world of data

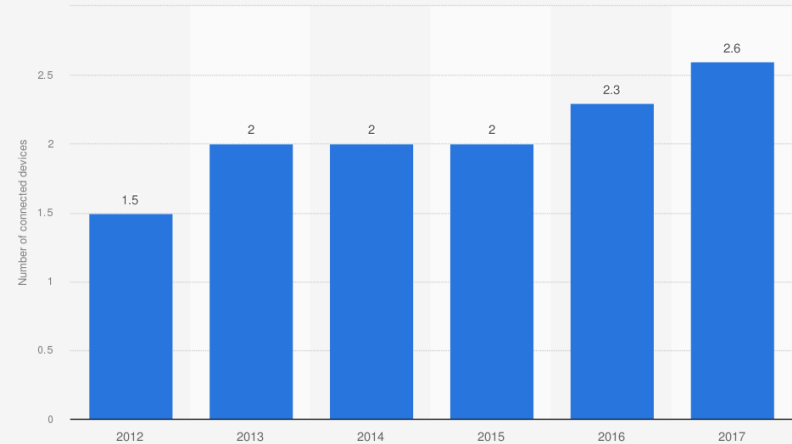
Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025 (in zettabytes)



Sources
IDC; Seagate; Statista estimates
© Statista 2022

Additional Information:
Worldwide; 2010 to 2020

Average number of connected devices used per person in Italy from 2012 to 2017



Sources
TNS Infratest; Google
© Statista 2022

Additional Information:
Italy; TNS Infratest; Google; 2012 to 2017; 1,000 respondents; 16 years and older; Computer-assisted telephone interview

1. Introduction
2. **Why a diversity aware RI?**
3. What type of RI?
4. Current platforms
5. Benchmark
6. A former approach
7. Limits

The computational turn



“[...] **digital technology is fundamentally changing** the way in which we engage in **the research process**. Many argue that this mediation is slowly beginning to change what it means to undertake research, affecting both the **epistemologies and ontologies** that underlie a research programme.” (Berry, 2011)

Bias in data analysis and reuse



- **COMPAS:** white offenders were identified to be labeled as lower risk more likely than black offenders despite their criminal history
- **Amazon's hiring algorithm:** reinforces the gender gap in hiring
- **Reproducibility crisis:** the results of many scientific studies are difficult or impossible to reproduce

Personal data in Google and Facebook (Asunciòn, 2017)

lack of valid consent given by their users

insufficient access and control given to users over their personal information

the risk of reidentification of anonymous personal data

1. Introduction
2. Why a diversity aware RI?
3. **What type of RI?**
4. Current platforms
5. Benchmark
6. A former approach
7. Limits

What is a Research Infrastructure?



According to EU Commission¹, RIs are "facilities that provide resources and services for the research communities to conduct research and foster innovation in their fields"

¹https://research-and-innovation.ec.europa.eu/funding/funding-opportunities/funding-programmes-and-open-calls/horizon-europe/research-infrastructures_en

What resources?



COLLECTION



MANAGEMENT



DISTRIBUTION



COMMUNITY

1. Introduction
2. Why a diversity aware RI?
3. What type of RI?
- 4. Current platforms**
5. Benchmark
6. A former approach
7. Limits

Data collection

Data from user interaction

- Psychlog
- Mobile Sensing Platform

Sensor data

- ResearchStack

Interaction and sensors

- AWARE
- CS Logger

Data management and distribution

UkDataArchive

GESIS

Open Science
Framework
(OSF)

ComputerOntario

BioBank

Crowd-sourcing vs. Communities



Pervasiveness of smart devices to collect data on large panel, e.g.:

- Participatory sensing
- Mobile Crowd Sensing

Limits

- Western countries panel
- Participant as a sensor

Focus on involvement and education, e.g.:

- Natural science: Zooniverse, CornellLab, iNaturalist
- Broader focus: SciStarter, EU citizen-science

Limits:

- Missing diversity aware data

1. Introduction
2. Why a diversity aware RI?
3. What type of RI?
4. Current platforms
5. **Benchmark**
6. A former approach
7. Limits

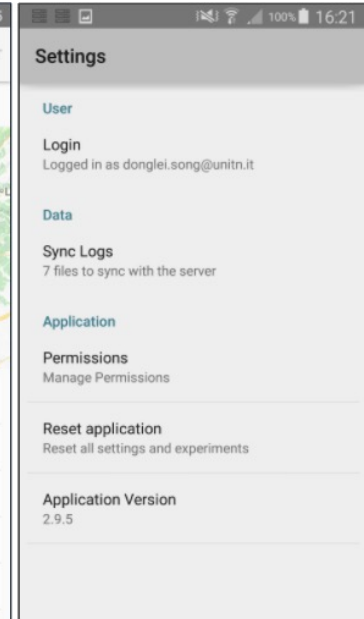
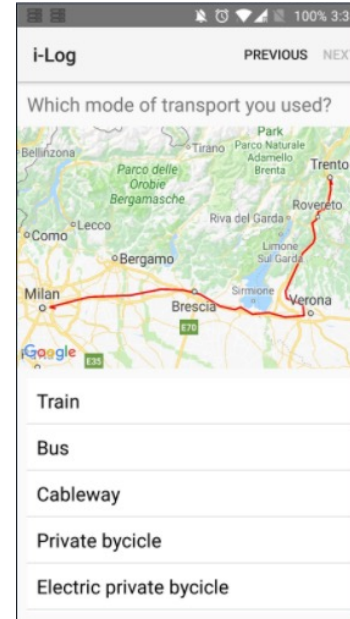
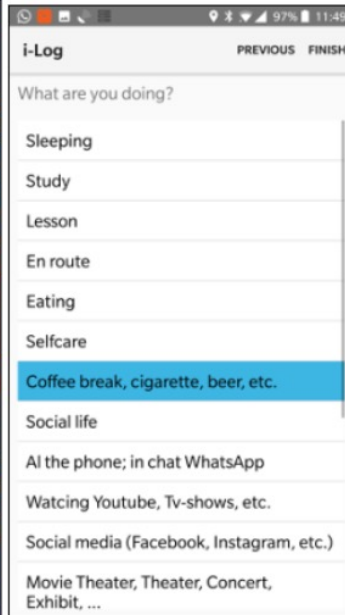
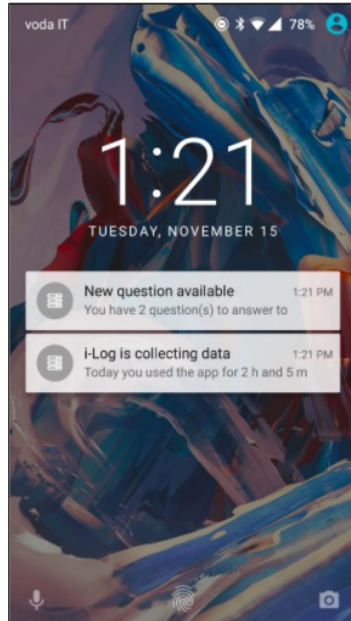
Benchmark



	Collection	Ethics & GDPR	Management	Distribution	Community
AWARE	✓		✓		
Citizen Science Logger	✓		✓		✓
UK Data Archive		✓	✓		✓
Zooniverse	✓			✓	✓

1. Introduction
2. Why a diversity aware RI?
3. What type of RI?
4. Current platforms
5. Benchmark
- 6. A former approach**
7. Limits

iLog App



Design: Support material



Design & Management

- Research design and management
- Research process Gantt
- iLog questions and sensors specifications

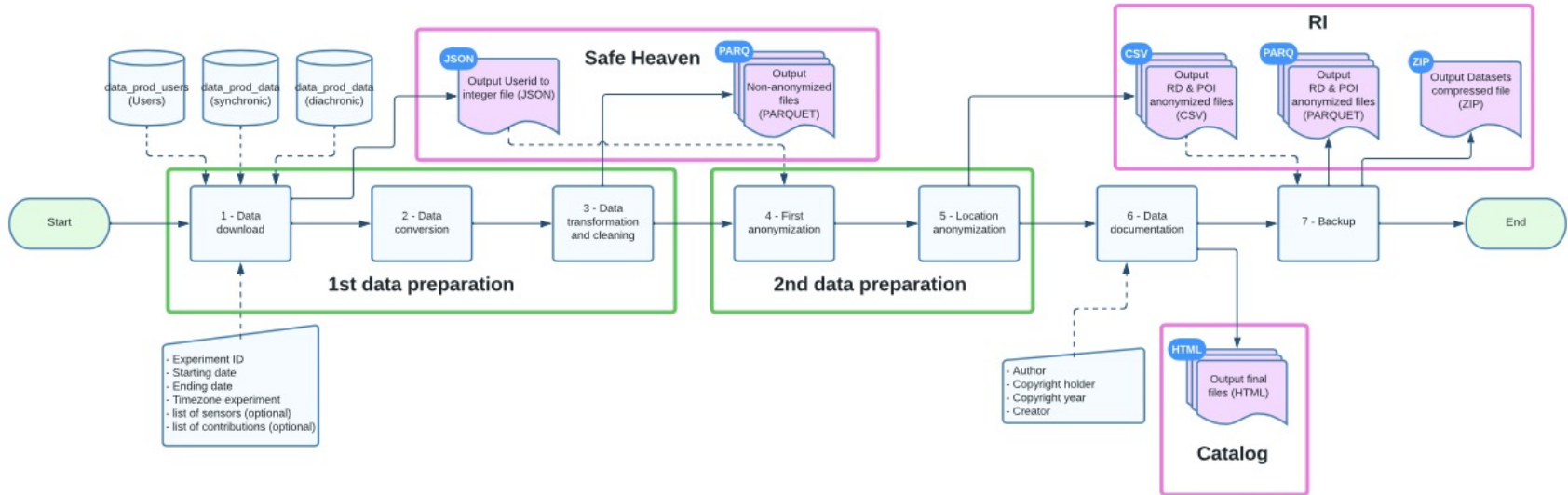
GDPR & Ethics

- Informed consent
- Privacy statement
- DPIA (Processor and Controller)
- Ethical committee approval

Guidelines

- Ethics Self Assessment Tool – UK Statistics Authority
- Hands on guide to research

Data preparation





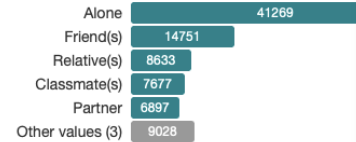
A5

Categorical

MISSING

With whom are you?

Distinct	8
Distinct (%)	< 0.1%
Missing	79840
Missing (%)	47.5%
Memory size	1.3 MiB



Overview Categories

Common Values

Value	Count	Frequency (%)
Alone	41269	24.6%
Friend(s)	14751	8.8%
Relative(s)	8633	5.1%
Classmate(s)	7677	4.6%
Partner	6897	4.1%
Roommate(s)	6265	3.7%
Colleague(s)	1470	0.9%
Other	1293	0.8%
(Missing)	79840	47.5%

Category Frequency Plot



- Contributionanswers
- Questionnaires
- Airplane mode event
- Application event
- Battery monitoring log
- Batterycharge event
- Cellular network
- Doze event
- Headset plug event
- Location event per time POI
- Location event per time RD
- Music event
- Notification event
- Proximity event
- Ring mode event
- Screen event
- Wifi event
- Wifi networks event

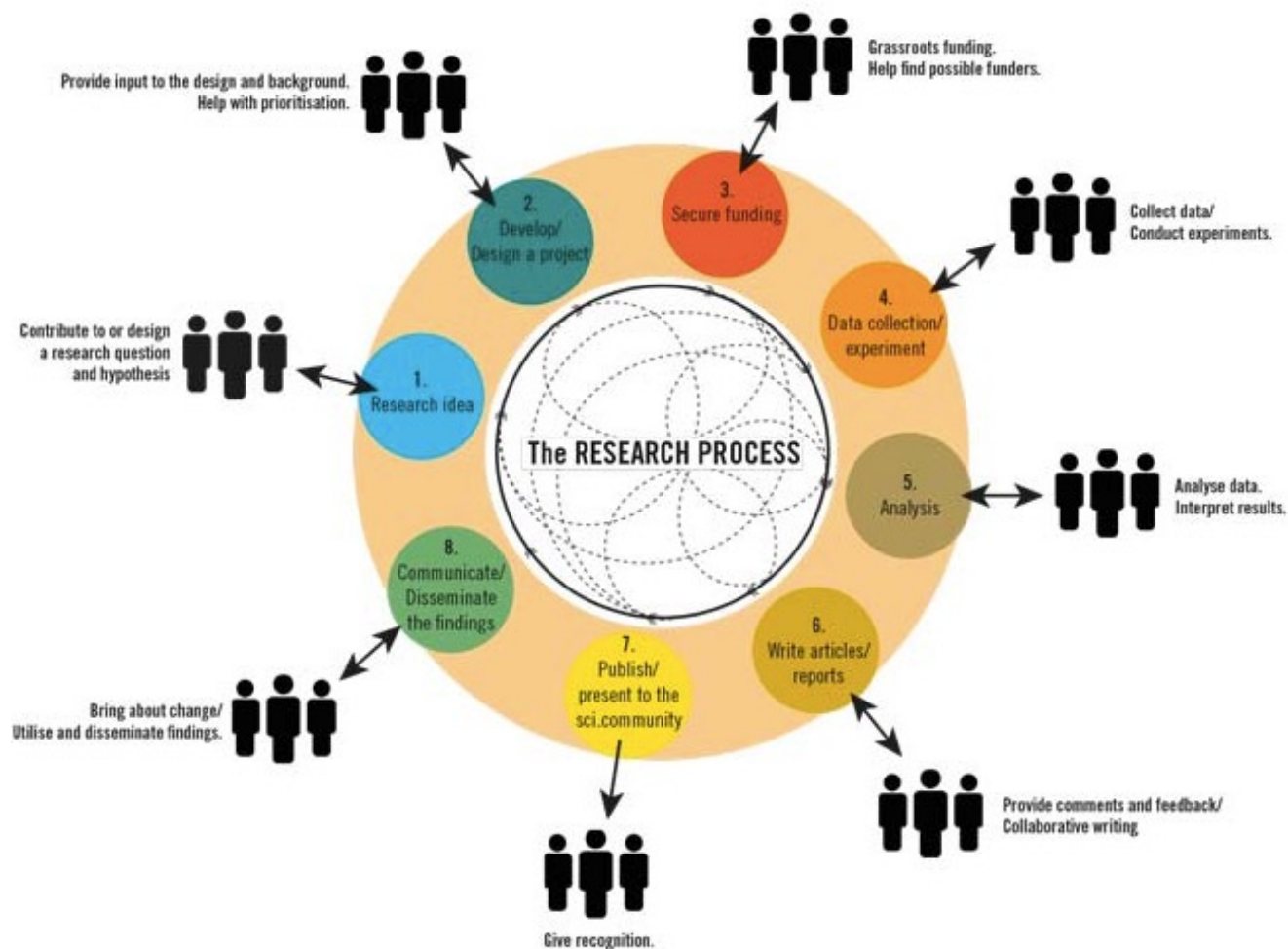
The LivePeople catalog

Community

Generating spaces for interactions:

among participants and researcher (e.g., between students and professors)

within participants or researcher (e.g., sharing knowledge on study design)



Community DB



Studies

Manage Profile

Boris Johnson



Studies

Studies / Create New

New Study

Title

Start date

End date

Funds

Status

CETS

Contacts

participant11@gmail.com,
Participant22@gmail.com,
participant33@gmail.com,
participant44@gmail.com

DESCRIPTION

Drugin this study user brain signal will be read while playing e-games
line 2

PURPOSE

Analyse decision making process.

TASKS

Task 1
Task 2
Task 3

ELIGIBILITY CRITERIA

Participant must not have been part of similar activities in last 3 months. Participant must also be Covid-19 vaccinated.

PRIVACY LINKS

[Privacy link 1, Privacy link 2](#)

Cancel Create

Details



David Miller

General info

Status Active

Institution University of Trento

Phone number (430) 065-7387

Email participant@gmail.com

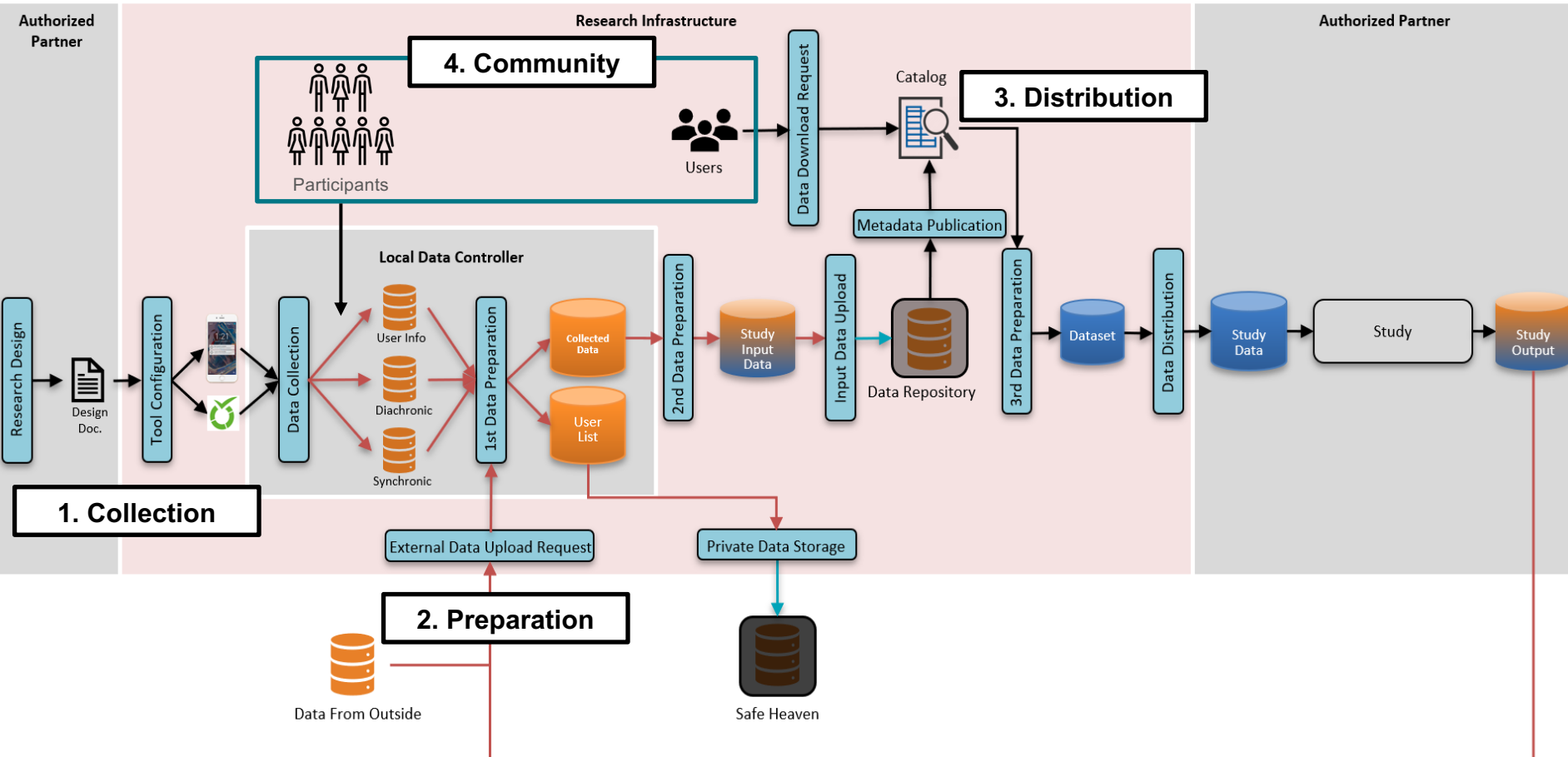
Address North Las Vegas, NV

[Data privacy policy](#)

Edit

Name <input type="text" value="David"/>	Surname <input type="text" value="Miller"/>
Gender <input type="text" value="Male"/>	Phone Number <input type="text" value="(430) 065-7387"/>
Place of Birth <input type="text" value="Las Vegas"/>	Date of Birth <input type="text" value="12/12/2000"/>
Address <input type="text" value="North Las Vegas, NV"/>	Nationality <input type="text" value="American"/>
Institution <input type="text" value="University of Trento"/>	Faculty <input type="text" value="Computer Science"/>
Course Year <input type="text" value="1"/>	<input checked="" type="checkbox"/> Is Erasmus
Codice Fiscale <input type="text" value="HMAA1212123434LK"/>	Status <input type="text" value="Active"/>
IBAN <input type="text" value="IT1210101010101010XY"/>	

Cancel Save



→ personal or sensitive data flow
 → data flow
 → storage flow
 ○ Personal data
 ○ Non-personal data
 ○ Possible Personal data
 GDPR compliant
 Ethical process

1. Introduction
2. Why a diversity aware RI?
3. What type of RI?
4. Current platforms
5. Benchmark
6. A former approach
7. **Limits**

Limits

- (i) the usability of iLog app by non-expert users, such as citizens;
- (ii) the validation of data management outcomes to effective reuse of resources;
- (iii) the community is still being tested and non-monetary incentive strategy is missing

GET IN TOUCH

-  <http://knowdive.disi.unitn.it/>
-  <http://datascientia.disi.unitn.it/>
-  [@knowdive](https://twitter.com/knowdive)
-  matteo.busso@unitn.it



WeNet project is funded by the EU's Horizon2020 programme under Grant Agreement number 823783.



THANK YOU!