# Collective Privacy Recovery

**Data-sharing Coordination via Decentralized Artificial Intelligence**

**Evangelos Pournaras**

**Trustworthy Distributed Intelligence**
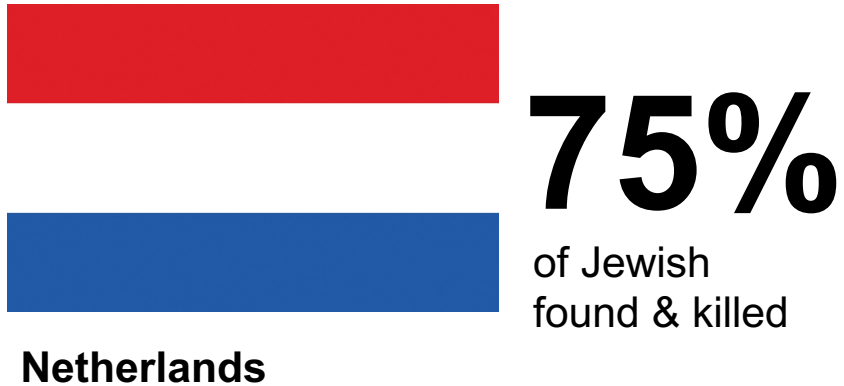
PNAS nexus

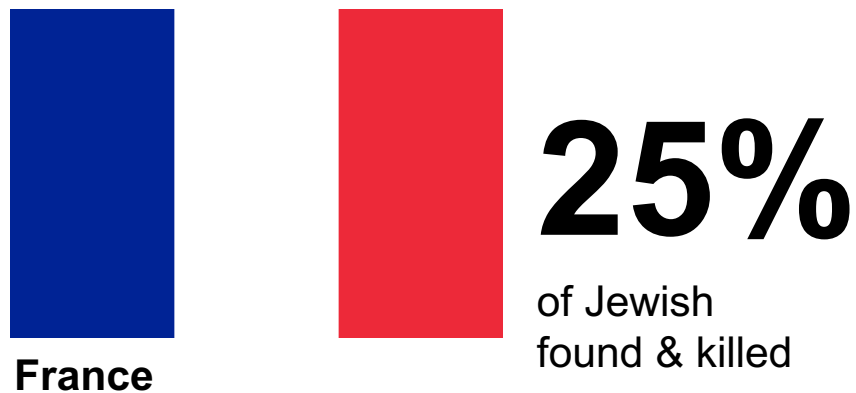# An Unforgiving Race of Power

## Privacy

# World War II

**75%**
of Jewish
found & killed

**Netherlands**

*What made such huge difference?*

**25%**
of Jewish
found & killed

**France**

# World War II



**75%**

of Jewish
found & killed

**Netherlands**

*What made such huge difference?*



PRIVACY IS POWER

WHY AND HOW YOU SHOULD
TAKE BACK CONTROL
OF YOUR DATA

CARISSA VÉLIZ

France had excluded sensitive
information from census for
privacy reasons

**25%**

of Jewish
found & killed

**France**

# Risks of Privacy Loss & the Privacy Paradox

*How many installed apps are needed to identify 91.2% of individuals?*

**?**

*How many spatio-temporal GPS records are needed to identify 95% of individuals?*

**?**

*From 90% of individuals who give up privacy, how many intend to protect it?*

**?**

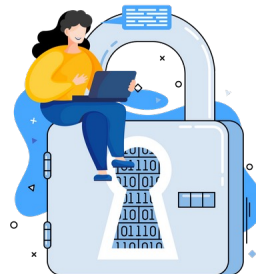# Risks of Privacy Loss & the Privacy Paradox

*How many installed apps are needed to identify 91.2% of individuals?*
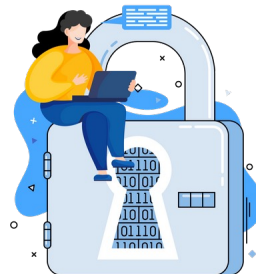


4

*How many spatio-temporal GPS records are needed to identify 95% of individuals?*



4

*From 90% of individuals who give up privacy, how many intend to protect it?*



76%

*See [4,5]*

# Implications of Collective Privacy Loss

**Environmental impact**

Data centers consume too much energy: faster growth of unprocessed data than Moore's law predictions

*Privacy loss resembles an ecological disaster with the global significance of climate change*

**Health impact**

Surveilance stress & anxiety [7]

**Social impact**

Algorithmic biases, discrimination, censorship, loss of freedoms

**Political impact**
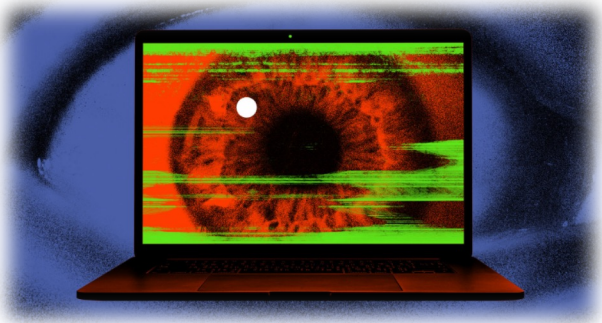
Influence of election results

# Privacy is not only an individual right…

**… it is also a shared value in the digital era!**

GDPR

**What are we missing here?**

Collective arrangements for sharing data that provide a *minimum quality of services* for *maximum privacy*

ELINOR OSTROM
2009 Nobel Laureate in Economic Sciences
Nobel medal © ® The Nobel Foundation

*who* is sharing to **whom**, **when**, **how much** of **what data** & for **what purpose?**

**Data as a scarce resource? Minimizing both excessive & insufficient levels of data**

Share data under the doctrine "*as little as possible, as much as necessary*"
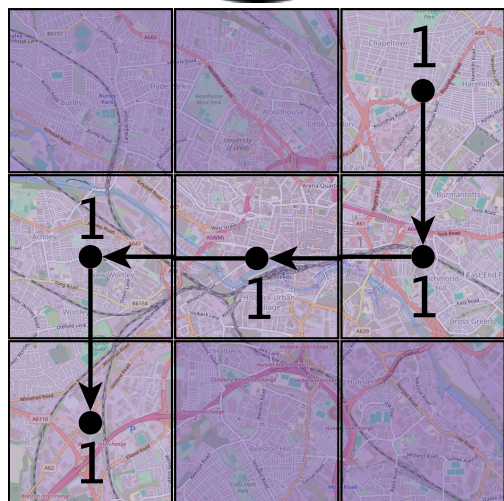
# Data collectives

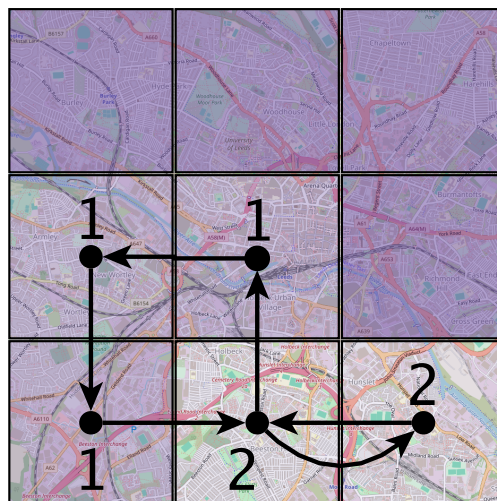# Privacy Loss is Coordination Deficit

A Toy Example

# Existing Status Quo of Data Sharing
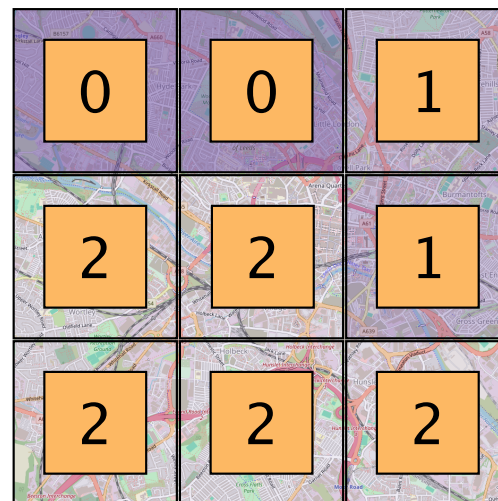
no coordination

Total 🗄 Data



5 ≥ 4 7

12 records

*Risk of identity inference [4]*

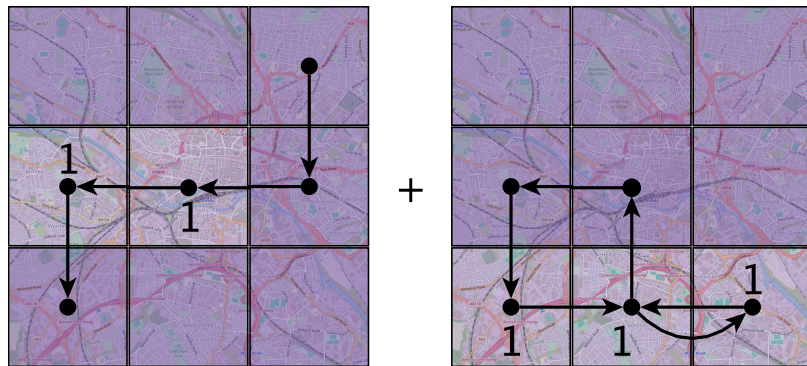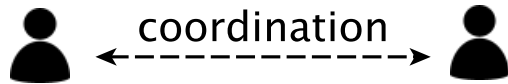*Is not this data (far most times) excessive?*

**Turn on your GPS?**

**Default**: Share all your personal data

# Coordinated Data Sharing

*Fairer data-sharing contributions*



no coordination

Total 🗄 Data

5

7

12 records

coordination

Lower 🗄 Data

2

3

5 records

**< 4**

*Reduced risk of identity inference [4]*

**> 50% ↓**

**Scenario:** Determine the highest traffic density areas

**Or.. selectively turning on & off your GPS?**

**Collective arrangement**: Share `as little as possible, as much as necessary'

# Coordinated Data Sharing

*Fairer data-sharing contributions*

no coordination — Total 🗄 Data

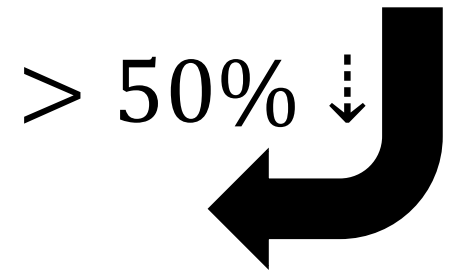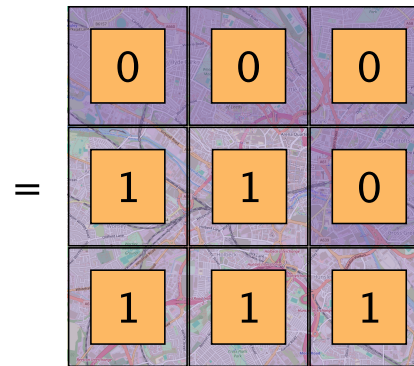5 + 7 = 12 records

coordination

Lower 🗄 Data

3 + 3 = 6 records

50% ↓

**Scenario:** Accurate traffic density estimation in the city center over periphery

< 4

*Reduced risk of identity inference [4]*

**Or.. selectively turning on & off your GPS?**

**Collective arrangement**: Share `as little as possible, as much as necessary'

UNIVERSITY OF LEEDS

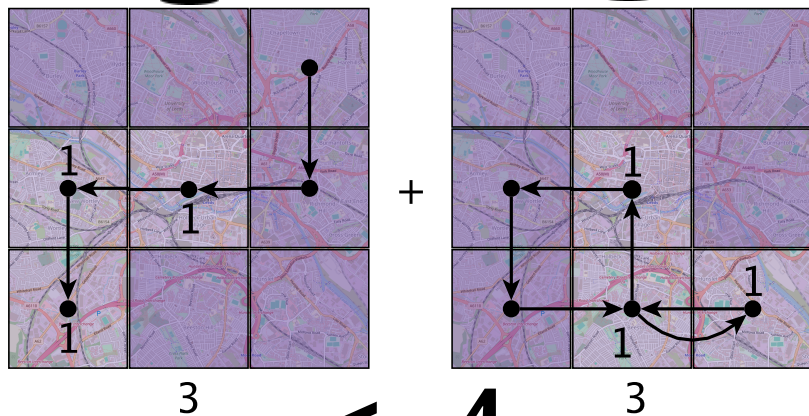# A Very Simple but Hard Idea to Materialize in Practice

*How to **automate & scale up** such
collective arrangements of data sharing?*

**Coordinated data sharing:**

A techno-socio-economic problem of computational complexity

**Modeling** as a *multi-agent discrete-choice optimization problem*

**Solving** using *decentralized, privacy-preserving & efficient AI*

# Related Work

**Security & cryptography: differential privacy, multi-party computation, k-anonymization**

Limited use of shared data

**Federated learning**

No coordination element for data-sharing optimization

**Personalized privacy assistants**

Privacy-intrusive themselves

**Methodological limitations**

Survey studies, limited realism, no causal inference

# A Living-lab Real-world Experiment

An inter-disciplinary study on coordinated data sharing

# Data Sharing Conditions & Hypotheses

A **novel & complete** spectrum for an in-depth understanding of data sharing choices



EXCESSIVE
DATA SHARING

AS LITTLE AS POSSIBLE
AS MUCH AS NECESSARY    010010
110101
INSUFFICIENT
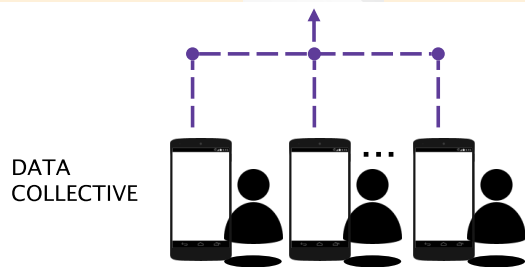DATA SHARING    101111

DATA
COLLECTIVE

...

DATA SHARING:         ATTITUDINAL

# Data Sharing Conditions & Hypotheses

A **novel & complete** spectrum for an in-depth understanding of data sharing choices



EXCESSIVE
DATA SHARING

AS LITTLE AS POSSIBLE
AS MUCH AS NECESSARY

010010
110101
101111

INSUFFICIENT
DATA SHARING

0010
0101

DATA
COLLECTIVE

DATA SHARING:          ATTITUDINAL                    INTRINSIC

# Data Sharing Conditions & Hypotheses

A **novel & complete** spectrum for an in-depth understanding of data sharing choices
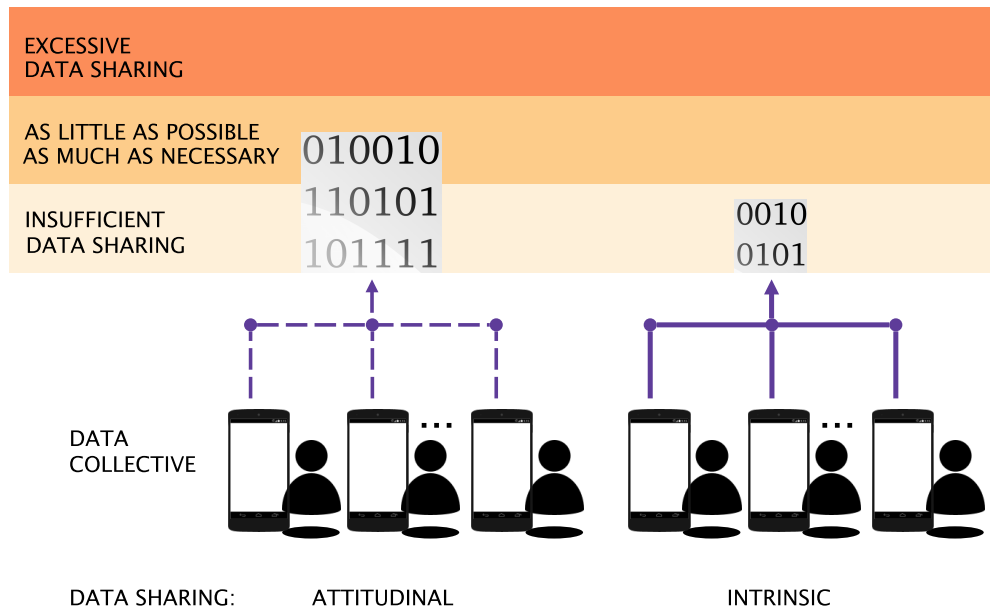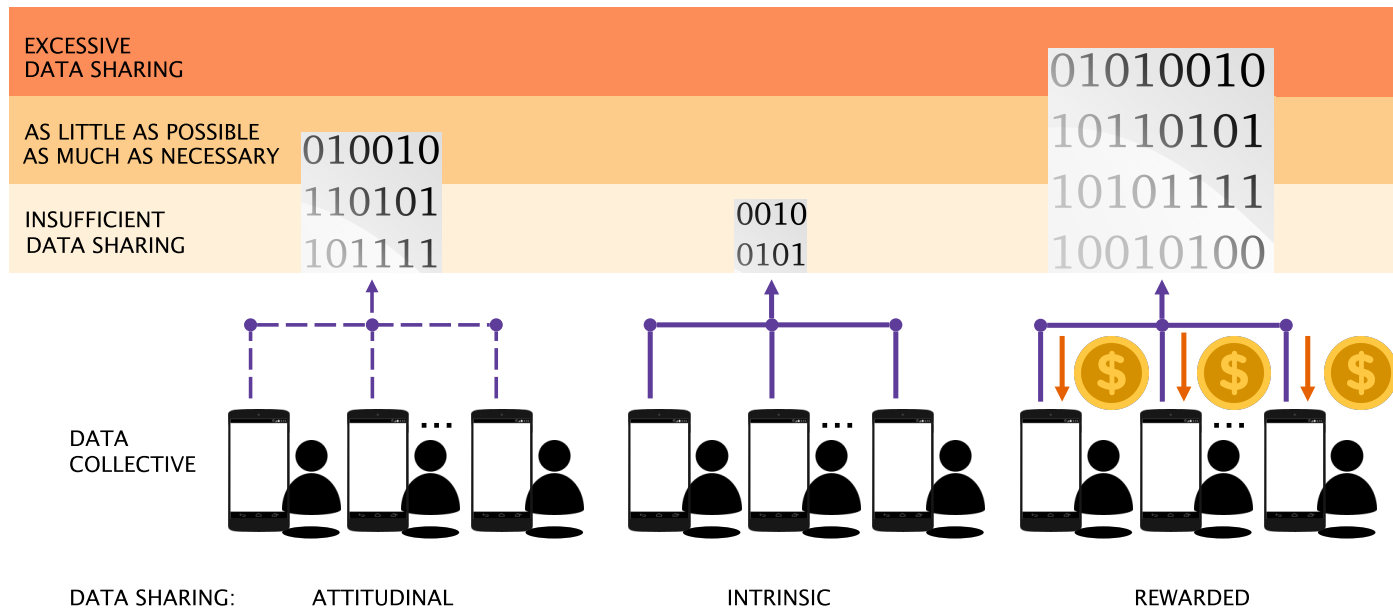
# Data Sharing Conditions & Hypotheses

A **novel & complete** spectrum for an in-depth understanding of data sharing choices



| | |
|---|---|
| EXCESSIVE DATA SHARING | COLLECTIVE PRIVACY RECOVERY |
| AS LITTLE AS POSSIBLE AS MUCH AS NECESSARY | |
| INSUFFICIENT DATA SHARING | |

010010
110101
101111

0010
0101

01010010
10110101
10101111
10010100

010010
110101
101111

DATA COLLECTIVE

AI–ASSISTED

DATA SHARING:   ATTITUDINAL   INTRINSIC   REWARDED   COORDINATED

# Data Sharing Model

**Data sharing criteria**: theory on *trust & risk* in data sharing [2]

**Mobile sensor data**: *ultimate killer app!*

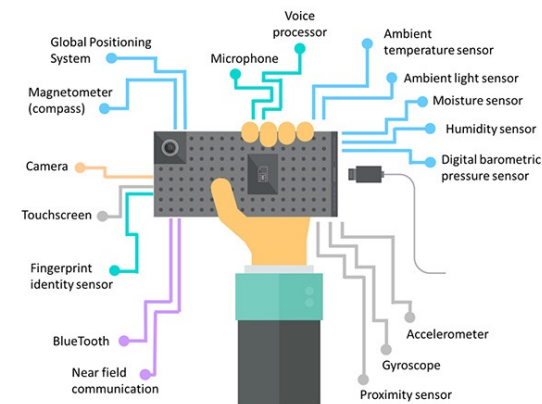**A full 4x4x4 factorial design:** 64 combinations to study!



Data-sharing Criteria

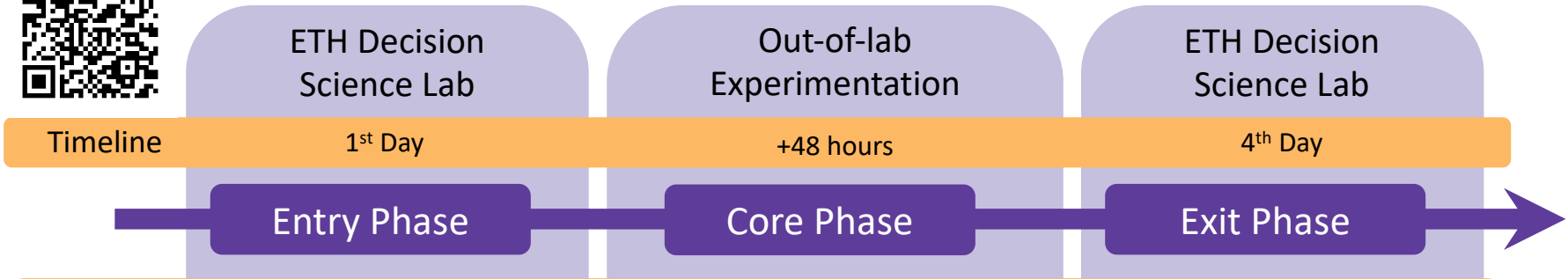| Sensor Type | Data Collector | Context |
| --- | --- | --- |
| Global Position System [gps] | Corporation [cor] | Social networking [soc] |
| Noise [noi] | Non-gov. organization [ngo] | Environment [env] |
| Accelerometer [acc] | Educational Institute [edu] | Transportation [tra] |
| Light [lig] | Gov. Organization [gov] | Health [hea] |

Data-sharing Elements

Data-sharing Scenarios

# A Novel Living-lab Experiment

>27,000 **real data disclosures** studied! **Open data** [6]



| Timeline | ETH Decision Science Lab | Out-of-lab Experimentation | ETH Decision Science Lab |
|---|---|---|---|
| | 1st Day | +48 hours | 4th Day |
| | Entry Phase | Core Phase | Exit Phase |

# A Novel Living-lab Experiment

>27,000 **real data disclosures** studied! **Open data** [6]



| | ETH Decision Science Lab | Out-of-lab Experimentation | ETH Decision Science Lab |
|---|---|---|---|
| Timeline | 1st Day | +48 hours | 4th Day |
| | **Entry Phase** | **Core Phase** | **Exit Phase** |
| Lab pool participants | 1. Instructions & consent<br>2. App installation<br>3. Entry app survey | 1. Daily app use<br>2. Data-sharing choices<br>3. Sensor data sharing | 1. Exit web survey<br>2. Interview<br>3. Compensation |

(1) Attitudinal & (2) intrinsic data sharing

(3) Rewarded data sharing (2x)

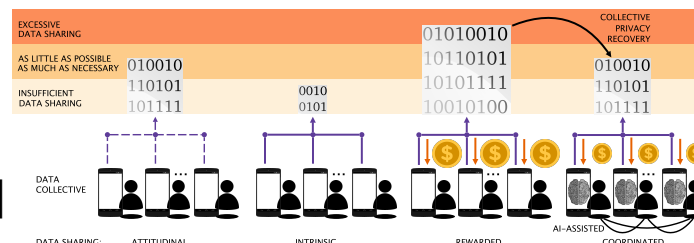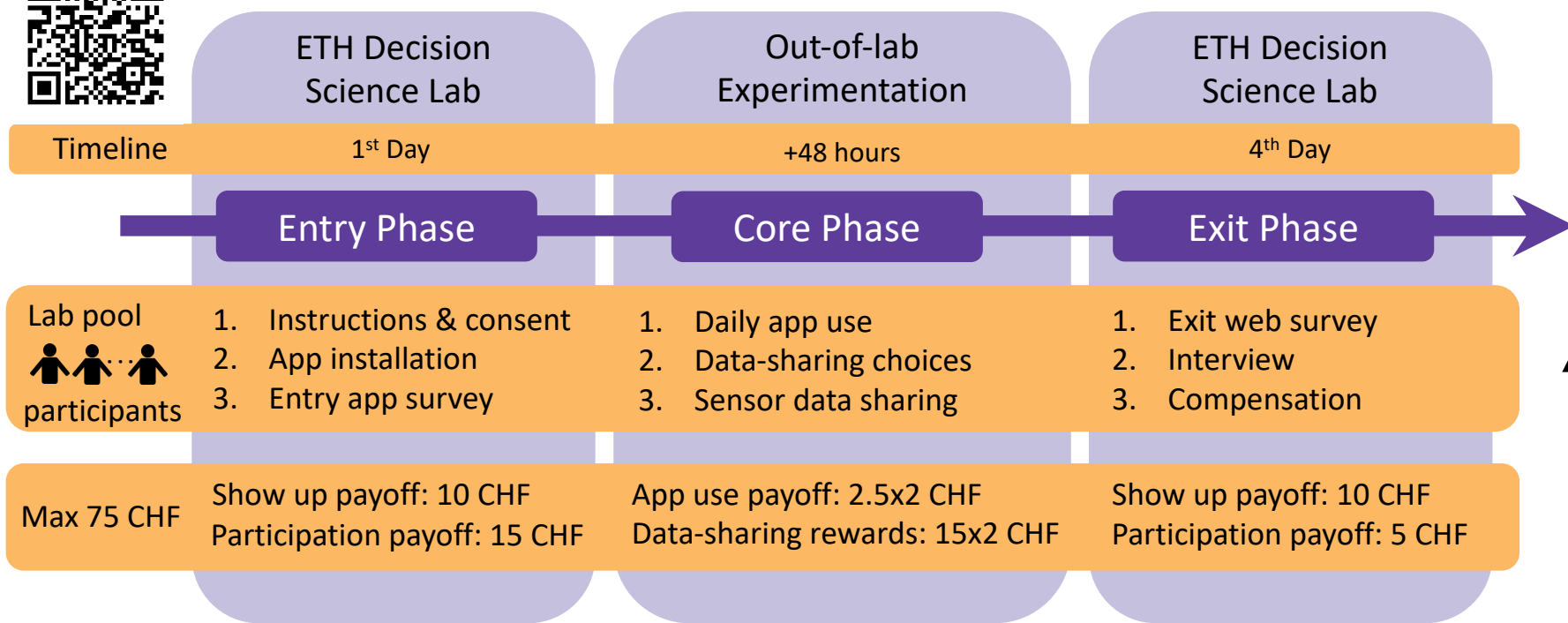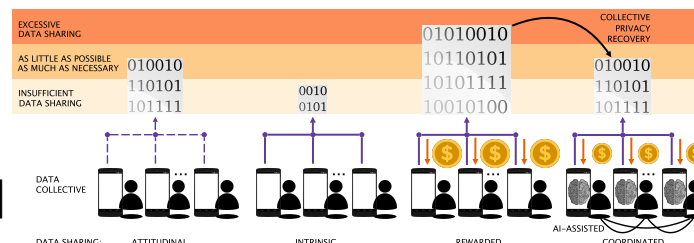(4) Coordinated data sharing

Three options to choose from:

- One intrinsic data sharing

- Two rewarded data sharing

# A Novel Living-lab Experiment

>27,000 **real data disclosures** studied! **Open data** [6]



| | ETH Decision Science Lab | Out-of-lab Experimentation | ETH Decision Science Lab |
|---|---|---|---|
| Timeline | 1st Day | +48 hours | 4th Day |
| | **Entry Phase** | **Core Phase** | **Exit Phase** |
| Lab pool participants | 1. Instructions & consent<br>2. App installation<br>3. Entry app survey | 1. Daily app use<br>2. Data-sharing choices<br>3. Sensor data sharing | 1. Exit web survey<br>2. Interview<br>3. Compensation |
| Max 75 CHF | Show up payoff: 10 CHF<br>Participation payoff: 15 CHF | App use payoff: 2.5x2 CHF<br>Data-sharing rewards: 15x2 CHF | Show up payoff: 10 CHF<br>Participation payoff: 5 CHF |

(1) Attitudinal & (2) intrinsic data sharing

(3) Rewarded data sharing (2x)

(4) Coordinated data sharing

Three options to choose from:
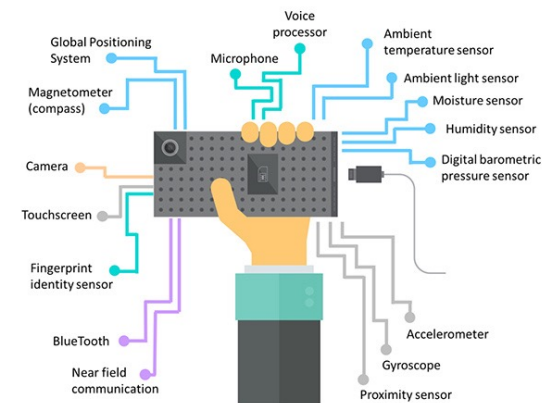
- One intrinsic data sharing
- Two rewarded data sharing

# Data Sharing Model

**Data sharing criteria**: theory on *trust & risk* in data sharing [2]

**Mobile sensor data**: *ultimate killer app!*

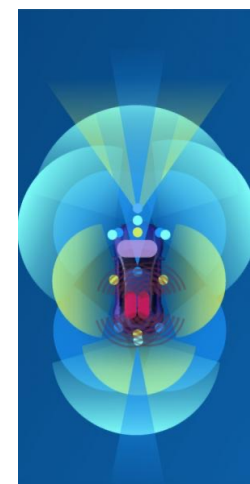**A full 4x4x4 factorial design:** 64 combinations to study!



Data-sharing Criteria

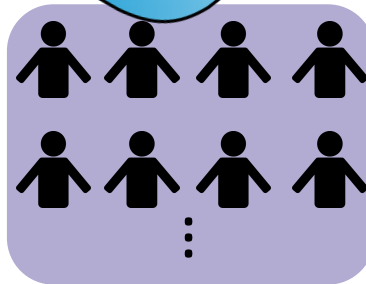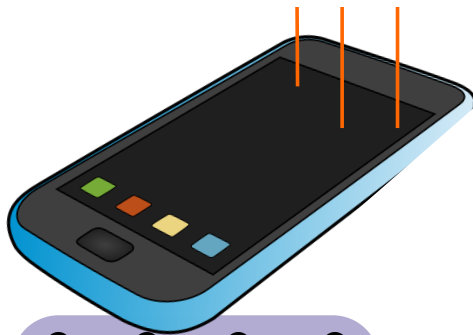| Sensor Type | Data Collector | Context |
|---|---|---|
| Global Position System [gps] | Corporation [cor] | Social networking [soc] |
| Noise [noi] | Non-gov. organization [ngo] | Environment [env] |
| Accelerometer [acc] | Educational Institute [edu] | Transportation [tra] |
| Light [lig] | Gov. Organization [gov] | Health [hea] |

Data-sharing Elements

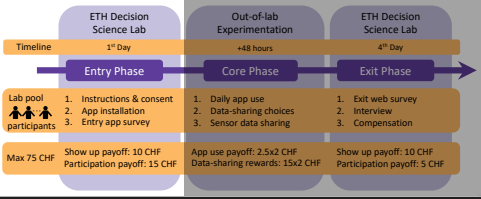Data-sharing Scenarios

# Data Collection Infrastructure



Sensor Data

Participants

Data Collectors

Private Newspaper

Software NGO

Technical University

Governmental Confederation

# Data Collection Infrastructure

Local Data Management System

Sensor Data

Experimental Data

Mobile App

What do you prefer to improve?

$ 🔒

Rewards    Privacy

0.0    0.0%

CHF    Privacy

Remore Data Management System

Shared Sensor Data

Participants

Data Access Portal

Private Newspaper    Software NGO    Technical University    Governmental Confederation

Data Collectors

UNIVERSITY OF LEEDS

| Timeline | ETH Decision Science Lab 1st Day | Out-of-lab Experimentation +48 hours | ETH Decision Science Lab 4th Day |
|---|---|---|---|
| | Entry Phase | Core Phase | Exit Phase |
| Lab pool participants | 1. Instructions & consent 2. App installation 3. Entry app survey | 1. Daily app use 2. Data-sharing choices 3. Sensor data sharing | 1. Exit web survey 2. Interview 3. Compensation |
| Max 75 CHF | Show up payoff: 10 CHF Participation payoff: 15 CHF | App use payoff: 2.5x2 CHF Data-sharing rewards: 15x2 CHF | Show up payoff: 10 CHF Participation payoff: 5 CHF |

UNIVERSITY OF LEEDS

# 1. Attitudinal Data

How intrusive are the following features of information sharing?
Sensors

Data collectors

Context/Purpose

**Data sharing criteria**

How privacy intrusive is the data sharing of the following sensors?
Accelerometer
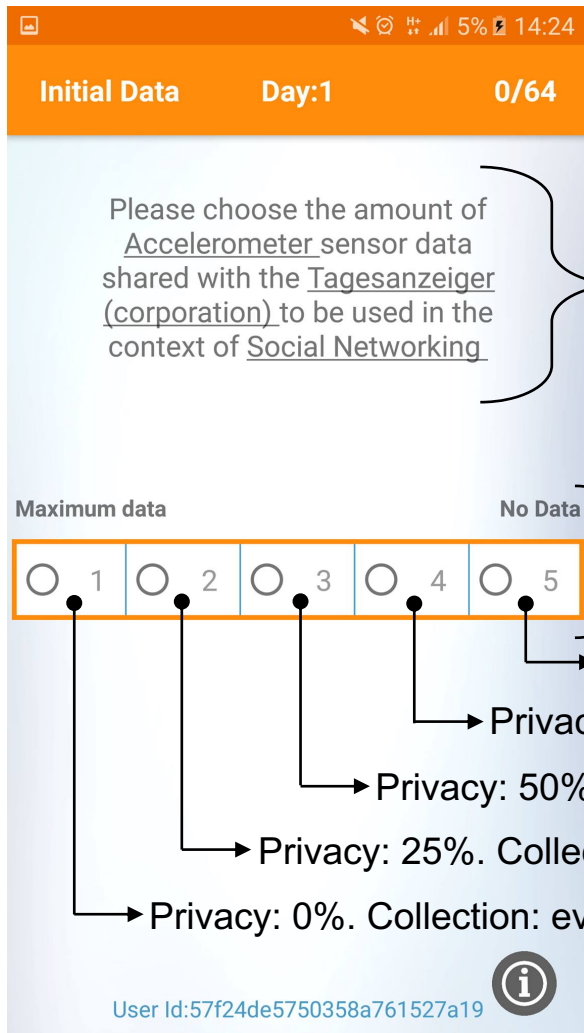
Location

Light

Noise

Sensors

How privacy intrusive are the following data collectors of your mobile sensor data?
Corporations

Non-governmental Organizations

Governments

Educational Institutes

Data collectors

**Studied data-sharing criteria**

How privacy intrusive are the following contexts under which sensor data is used by stakeholders?
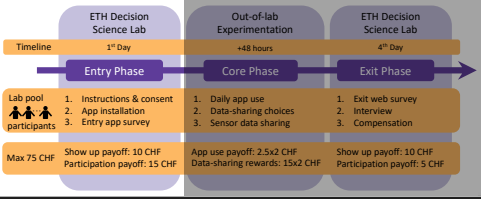Health/Fitness

Social Networking

Environment

Transportation

Contexts

*Survey questions to derive the privacy intrusion level*

| | ETH Decision Science Lab | Out-of-lab Experimentation | ETH Decision Science Lab |
|---|---|---|---|
| Timeline | 1st Day | +48 hours | 4th Day |
| | Entry Phase | Core Phase | Exit Phase |
| Lab pool participants | 1. Instructions & consent 2. App installation 3. Entry app survey | 1. Daily app use 2. Data-sharing choices 3. Sensor data sharing | 1. Exit web survey 2. Interview 3. Compensation |
| Max 75 CHF | Show up payoff: 10 CHF Participation payoff: 15 CHF | App use payoff: 2.5x2 CHF Data-sharing rewards: 15x2 CHF | Show up payoff: 10 CHF Participation payoff: 5 CHF |

UNIVERSITY OF LEEDS

# 2. Intrinsic Data Sharing of Participants



Initial Data    Day:1    0/64

Please choose the amount of Accelerometer sensor data shared with the Tagesanzeiger (corporation) to be used in the context of Social Networking

Question expressing a data sharing scenario

Maximum data                No Data

○ 1  ○ 2  ○ 3  ○ 4  ○ 5

Regulates the frequency of data collection

Privacy: 100%. No collection

Privacy: 75%. Collection: every 120s

Privacy: 50%. Collection: every 90s

Privacy: 25%. Collection: every 60s

Privacy: 0%. Collection: every 30s

User Id:57f24de5750358a761527a19

# 2. Intrinsic Data Sharing of Participants



Question expressing a data sharing scenario

Regulates the frequency of data collection

Privacy: 100%. No collection

Privacy: 75%. Collection: every 120s

Privacy: 50%. Collection: every 90s

Privacy: 25%. Collection: every 60s

Privacy: 0%. Collection: every 30s

*Privacy-utility trade-offs are also possible to make with differential privacy settings*

# 3. Rewarded Data Sharing of Participants

| Timeline | ETH Decision Science Lab 1ˢᵗ Day | Out-of-lab Experimentation +48 hours | ETH Decision Science Lab 4ᵗʰ Day |
| --- | --- | --- | --- |
| | Entry Phase | Core Phase | Exit Phase |
| Lab pool participants | 1. Instructions & consent 2. App installation 3. Entry app survey | 1. Daily app use 2. Data-sharing choices 3. Sensor data sharing | 1. Exit web survey 2. Interview 3. Compensation |
| Max 75 CHF | Show up payoff: 10 CHF Participation payoff: 15 CHF | App use payoff: 2.5x2 CHF Data-sharing rewards: 15x2 CHF | Show up payoff: 10 CHF Participation payoff: 5 CHF |

# 3. Rewarded Data Sharing of Participants

# 4. Coordinated Data Sharing

**A multi-agent discrete-choice**
**combinatorial optimization problem**

**3 options to choose from for each agent:**

*intrinsic vs. two rewarded data sharing*

# 4. Coordinated Data Sharing

**A multi-agent discrete-choice**
**combinatorial optimization problem**

**3 options to choose from for each agent:**
*intrinsic vs. two rewarded data sharing*

**Quality of service:**
<u>Global cost function</u>: *min root mean square error*
Matching indicator between shared & required data

**Privacy:**
<u>Local cost function</u>: data sharing level

# 4. Coordinated Data Sharing

**A multi-agent discrete-choice combinatorial optimization problem**

**Options to choose from for each agent:**

*intrinsic vs. two rewarded data sharing*

**Quality of service:**

Global cost function: *min root mean square error*

Matching indicator between shared & required data

**Privacy:**

Local cost function: data sharing level

***Collective learning heuristic of EPOS:***

*Decentralized*  *Unsupervised*  *Efficient*

*Privacy-preserving*  *Resilient*  *Scalable*



IRCAI GLOBAL TOP 100 OUTSTANDING PROJECT 2022
AI

**Open-source**
Github

# Three Key Results!

# Three Key Results

## 1. Coordinated data sharing is efficient

It <u>recovers privacy</u> for people & <u>reduces costs</u> for service providers by accessing less but better quality of data

# Three Key Results

**1. Coordinated data sharing is efficient**

It <u>recovers privacy</u> for people & <u>reduces costs</u> for service providers by accessing less but better quality of data

**2. Data collector & context are the most important criteria with which individuals makes data-sharing choices**

For <u>rewarded choices with privacy loss</u> though, the <u>type</u> of shared data becomes the most important criterion

# Three Key Results

**1. Coordinated data sharing is efficient**

It <u>recovers privacy</u> for people & <u>reduces costs</u> for service providers by accessing less but better quality of data

**2. Data collector & context are the most important criteria with which individuals makes data-sharing choices**

For <u>rewarded choices with privacy loss</u> though, the <u>type</u> of shared data becomes the most important criterion

**3. Individuals exhibit five key group-behavior changes from intrinsic to rewarded data sharing.**

They are <u>stable</u>, yet <u>reinforcing</u>

# 1. Coordinated data sharing is efficient

It <u>recovers privacy</u> for people & <u>reduces costs</u> for service providers by accessing less but better quality of data

# Privacy

**Significant privacy recovery via coordination**

High privacy-preservation choices involve
**data with low privacy sensitivity**

**Intrinsic vs. attitudinal**: correlated

**Reward-intrinsic vs. attitudinal**: correlated

# Privacy Goal Signals

Extracted "*easy*" & "*hard*" scenarios for the data collective to respond

**Very high**: Probability of sharing "5" at each data sharing scenario

…

**Very Low**: Probability of sharing "1" at each data sharing scenario

# Quality of Service

Rewards "**spoil**" data quality – Implications:

More data, more risks, more costs:

Financial, legal, environmental

Coordination "**mines**" data quality – Implications:

Less but more purposeful data

Minimizing excessive & insuffiecient data

# Data Sharing Cost

**Win-win for all**: higher privacy for people, <u>lower costs for service providers</u>

# 2. Data collector & context are the most important criteria with which individuals makes data-sharing choices

For <u>rewarded choices with privacy loss</u> though, the <u>type</u> of shared data becomes the most important criterion

# A Conjoint Analysis: Prediction Models

Type, collectors & contexts explain well privacy choices

# A Conjoint Analysis: Importance

**Rewards change the importance** of the data sharing criteria



**Data collector** & **context** determine privacy preservation

**Data type** determines rewarded choices with privacy loss

# 3. Individuals exhibit five key group-behavior changes from intrinsic to rewarded data sharing.

They are <u>stable</u>, yet <u>reinforcing</u>

# Data Sharing Behaviors

All possible behavioral changes

observed & unobserved:

| Data Sharing: | Without Rewards | | | With Rewards | | |
|---|---|---|---|---|---|---|
| | *Low* | Moderate | High | Low | Moderate | High |
| Privacy ignorants | | | ✓ | | | ✓ |
| Privacy neutrals | | ✓ | | | ✓ | |
| Privacy preservers | ✓ | | | ✓ | | |
| Rewards seekers | | ✓ | | | | ✓ |
| Rewards opportunists | ✓ | | | | | ✓ |
| Privacy sacrificers | ✗ | | | | ✗ | |
| Reward opposers (sharer) | | | ✗ | ✗ | | |
| Reward opposers (neutral) | | ✗ | | ✗ | | |
| Reward sacrificer (sharer) | | | ✗ | | ✗ | |

# Data Sharing Behaviors

All possible behavioral changes

observed & unobserved:

| Data Sharing: | Without Rewards | | | With Rewards | | |
|---|---|---|---|---|---|---|
| | *Low* | Moderate | High | Low | Moderate | High |
| Privacy ignorants | | | ✓ | | | ✓ |
| Privacy neutrals | | ✓ | | | ✓ | |
| Privacy preservers | ✓ | ✓ | | ✓ | | |
| Rewards seekers | | ✓ | | | | ✓ |
| Rewards opportunists | ✓ | | | | | ✓ |
| Privacy sacrificers | ✗ | | | | ✗ | |
| Reward opposers (sharer) | | | ✗ | ✗ | | |
| Reward opposers (neutral) | | ✗ | | ✗ | | |
| Reward sacrificer (sharer) | | | ✗ | | ✗ | |



| Clustering algorithms | k-means | hierachical | pamkCBI |
|---|---|---|---|
| Privacy ignorants | 0.79 (8) | 0.67 (41) | 0.58 (48) |
| Privacy neutrals | 0.93 (0) | 0.88 (1) | 0.7 (31) |
| Privacy preservers | 0.89 (7) | 0.76 (16) | 0.7 (31) |
| Rewards seekers | 0.83 (1) | 0.75 (17) | 0.61 (37) |
| Rewards opportunists | 0.84 (6) | 0.76 (14) | 0.56 (51) |

High bootstrap values, same
clusters among different algorithms

# Data Sharing Behaviors

All possible behavioral changes

observed & unobserved:



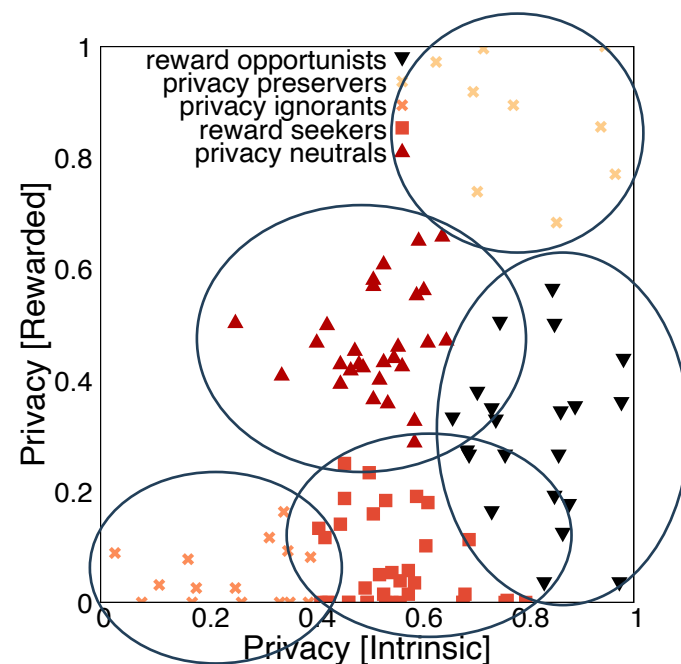| Data Sharing: | Without Rewards | | | With Rewards | | |
|---|---|---|---|---|---|---|
| | *Low* | Moderate | High | Low | Moderate | High |
| Privacy ignorants | | | ✓ | | | ✓ |
| Privacy neutrals | | ✓ | | | ✓ | |
| Privacy preservers | ✓ | | | ✓ | | |
| Rewards seekers | | ✓ | | | | ✓ |
| Rewards opportunists | ✓ | | | | | ✓ |
| Privacy sacrificers | ✗ | | | | ✗ | |
| Reward opposers (sharer) | | | ✗ | ✗ | | |
| Reward opposers (neutral) | | ✗ | | ✗ | | |
| Reward sacrificer (sharer) | | | ✗ | | ✗ | |

| Westin's population categories [7, 8] | | Data-sharing Groups ($n = 84$). | |
|---|---|---|---|
| Privacy fundamentalists | 25% | Privacy preservers / Reward opportunists | 26.2% |
| Privacy pragmatists | 57% | Privacy neutrals / Reward seekers | 57.14% |
| Privacy unconcerned | 18% | Privacy ignorants | 16.7% |

| Clustering algorithms | k-means | hierachical | pamkCBI |
|---|---|---|---|
| Privacy ignorants | 0.79 (8) | 0.67 (41) | 0.58 (48) |
| Privacy neutrals | 0.93 (0) | 0.88 (1) | 0.7 (31) |
| Privacy preservers | 0.89 (7) | 0.76 (16) | 0.7 (31) |
| Rewards seekers | 0.83 (1) | 0.75 (17) | 0.61 (37) |
| Rewards opportunists | 0.84 (6) | 0.76 (14) | 0.56 (51) |

Significant match to Westin's
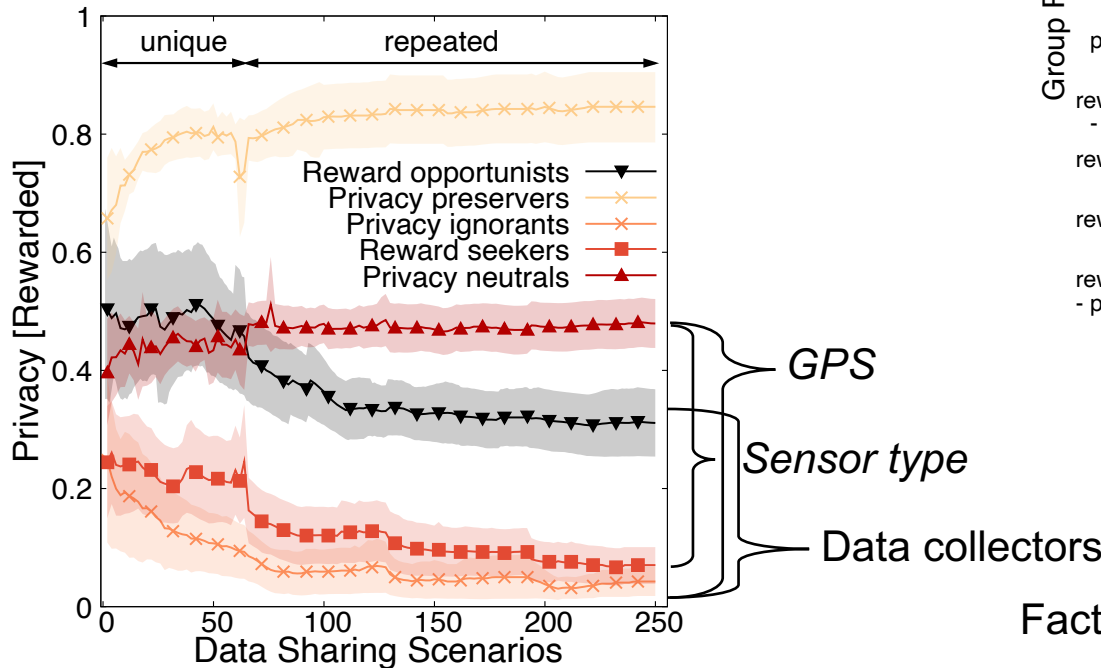general population categories [8]

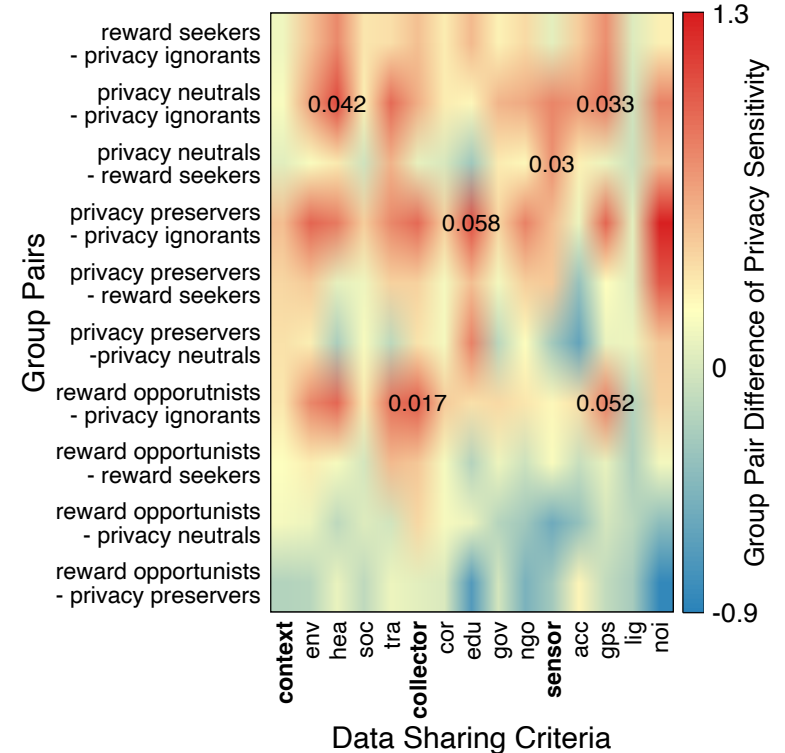High bootstrap values, same
clusters among different algorithms

# Data Sharing Polarization

Repititve data sharing dillemas **create polarization**

*Privacy preservers & ingorants tend to preserve & ignore further privacy*

ANOVA posthoc analysis



Factors explaining group differences

# Discussion, Lessons Learnt & Future Work

**Data collectives**: A win-win modus operandi for privacy
recovery & quality of service: <u>less & better data</u>

**Policy interventions**: Tailored campaigns based on the importance of data
sharing (i) **criteria** & (ii) **groups** for higher <u>privacy awareness & engagement</u>

**Generative AI**: An opportunity to build large language models **ethically aligned** to values
of communities sharing their data

**Temporal coordination** as an implementation of the "*right to be forgotten*"

# Questions?

E-mail: [e.pournaras@leeds.ac.uk](mailto:e.pournaras@leeds.ac.uk)



# References

[1] Pournaras, E., Ballandies, M.C., Bennati, S. and Chen, C.F., 2024. Collective privacy recovery: Data-sharing coordination via decentralized artificial intelligence. PNAS nexus, 3(2), p.pgae029.

[2] Adams, A. and Angela Sasse, M., 2001. Privacy in multimedia communications: Protecting users, not just data. In People and computers XV—interaction without frontiers: Joint Proceedings of HCI 2001 and IHM 2001 (pp. 49-64). Springer London.

[3] Sekara, V., Alessandretti, L., Mones, E. and Jonsson, H., 2021. Temporal and cultural limits of privacy in smartphone app usage. *Scientific reports*, *11*(1), p.3861.

[4] De Montjoye, Y.A., Hidalgo, C.A., Verleysen, M. and Blondel, V.D., 2013. Unique in the crowd: The privacy bounds of human mobility. *Scientific reports*, *3*(1), pp.1-5.

[5] Rose, J., Rehse, O. & Röber, B. The value of our digi- tal identity. URL https://www.bcg.com/en-gb/publications/2012/ digital-economy-consumer-insight-value-of-our-digital-identity, https://www.bcg.com/en-gb/publications/2014/ data-privacy-numbers.

[6] Data Collectives open data: https://doi.org/10.6084/ m9.figshare.21750158

[7] Oulasvirta, A., Pihlajamaa, A., Perkiö, J., Ray, D., Vähäkangas, T., Hasu, T., Vainio, N. and Myllymäki, P., 2012, September. Long-term effects of ubiquitous surveillance in the home. In Proceedings of the 2012 ACM Conference on Ubiquitous Computing (pp. 41-50).

[8] Kumaraguru, P. & Cranor, L. F. Privacy indexes: a survey of Westin's studies (Carnegie Mellon University, School of Computer Science, 2005)