



# SURVEY FUTURES

SURVEY DATA COLLECTION  
METHODS COLLABORATION

# RS7: Data Integration

**Alexandru Cernat**

**Thomas O'Toole, Peter Barlow, Natalie Shlomo, Nikos Tzavidis, Joseph Sakshaug**



# What is Data Integration?

Data integration refers to the process of **bringing together information from multiple data sources in a coherent and consistent manner.**

- Data integration makes it possible to examine relationships between factors which might not be visible from any one data source alone.

***Research strand 7 of Survey Futures is concerned with how non-survey data can be used to enhance survey data.***

# What are our Research Themes?

## Practitioner Guide 1



Options for integrating non-survey and population survey data.

## Practitioner Guide 2



Using integrated non-survey data for evaluating and correcting non-response bias in surveys.

## Practitioner Guide 3



Using integrated non-survey data for monitoring and intervening in survey data collection.

# Accessing our work so far

## Practitioner guides: [surveyfutures.net/practice-guides/](https://surveyfutures.net/practice-guides/)

- Survey Practice Guide 1: Data Integration
- Survey Practice Guide 5: Using integrated non-survey data for evaluating and correcting for non-response
- *Using integrated non-survey data for monitoring and intervening in survey data collection*

## Events:

- [youtube.com/watch?v=Awjag8x4GiY](https://youtube.com/watch?v=Awjag8x4GiY)
- [youtube.com/watch?v=VxSC\\_obOO88](https://youtube.com/watch?v=VxSC_obOO88)
- [youtu.be/8Z\\_g-nb-XG0](https://youtu.be/8Z_g-nb-XG0)



**SURVEY  
FUTURES**  
SURVEY DATA COLLECTION  
METHODS COLLABORATION

# What's New in Data Integration?

## *A Scoping Review and Reporting Framework*



University of Essex



University of  
Southampton



Economic  
and Social  
Research Council

NCRM NATIONAL CENTRE FOR  
RESEARCH METHODS

Office for  
National Statistics

UCL

WARWICK  
THE UNIVERSITY OF WARWICK

MANCHESTER  
THE UNIVERSITY OF MANCHESTER

CITY  
UNIVERSITY OF THE FINANCIAL DISTRICT

National Centre  
for Social Research

LSE  
LONDON SCHOOL OF  
ECONOMICS AND  
POLITICAL SCIENCE

Ulster  
University

Unil  
UNIL | Université de Lausanne

Ipsos

KANTAR  
PUBLIC

# Research Questions



Which survey and non-survey data sources are commonly integrated?



Which methods are most frequently used to link survey and non-survey data?



For what research purposes are integrated survey and non-survey data used?



What is the quality of reporting for data integration?

# Methods

We conducted a PRISMA Scoping review (Tricco, 2018)

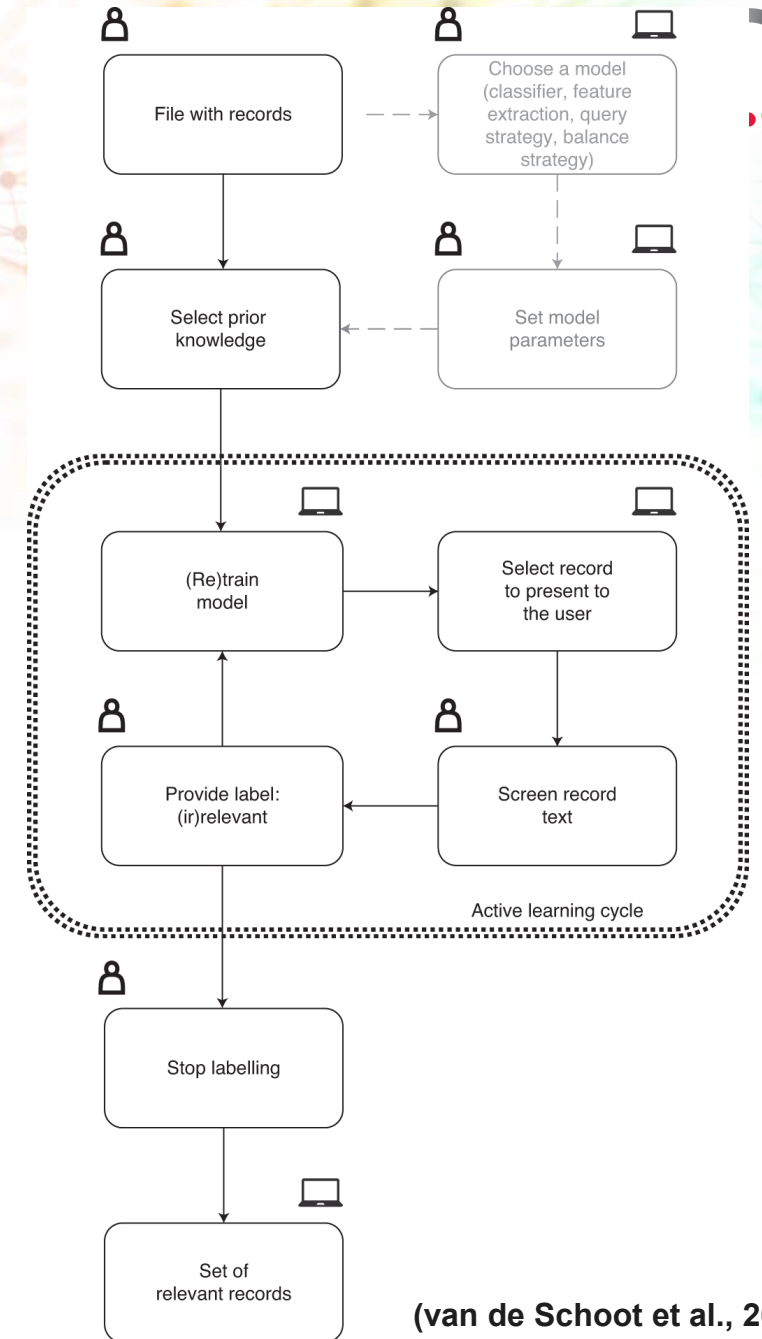
- Ovid (including the Cochrane Library, APAPsycInfo, Embase, Econlit and Medline)
- Web of Science (Core Collection)
- Scopus

Final searches conducted 16/05/2025

- *Screening was calibrated across two coders*

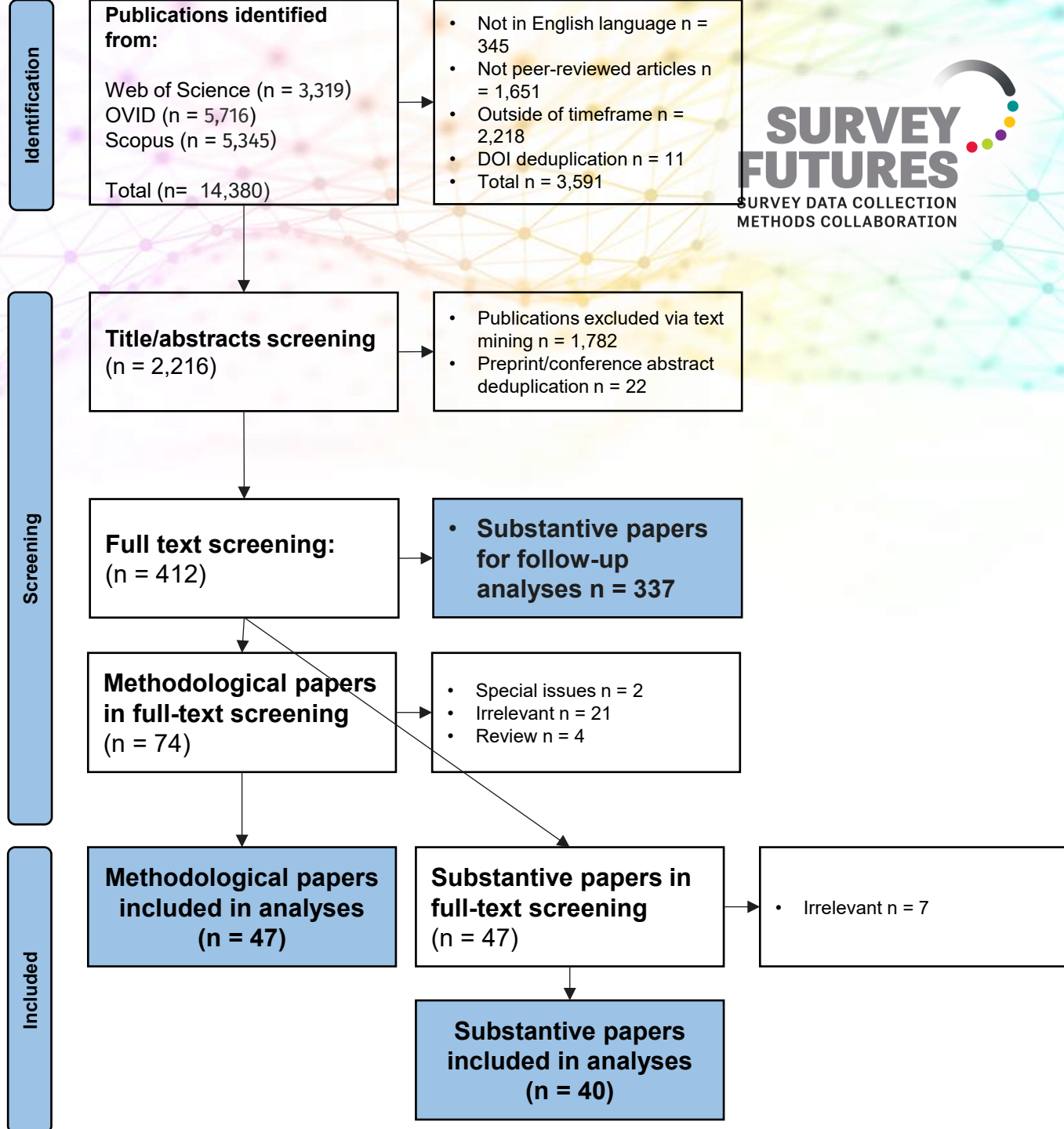
Active learning model (ASReview; van de Schoot et al., 2021)

- *Utrecht University's AI Lab "AI-aided Knowledge Discovery"*



(van de Schoot et al., 2021)

# The PRISMA Diagram



# RQ1: Which survey and non-survey data sources are commonly integrated in **methods literature**?

	Administrative records	Smart tracker app	Place-based characteristics	Sensor data	Commercial data	Social media data	Total
<b>Panel study</b>	18	4	1		2	1	26
<b>Birth cohort study</b>	3						3
<b>Repeated cross-sectional</b>	13		1	2			16
<b>Cross-sectional</b>	2		1				3
<b>Census</b>	1						1
<b>Total</b>	<b>37</b>	<b>4</b>	<b>3</b>	<b>2</b>	<b>2</b>	<b>1</b>	<b>49</b>

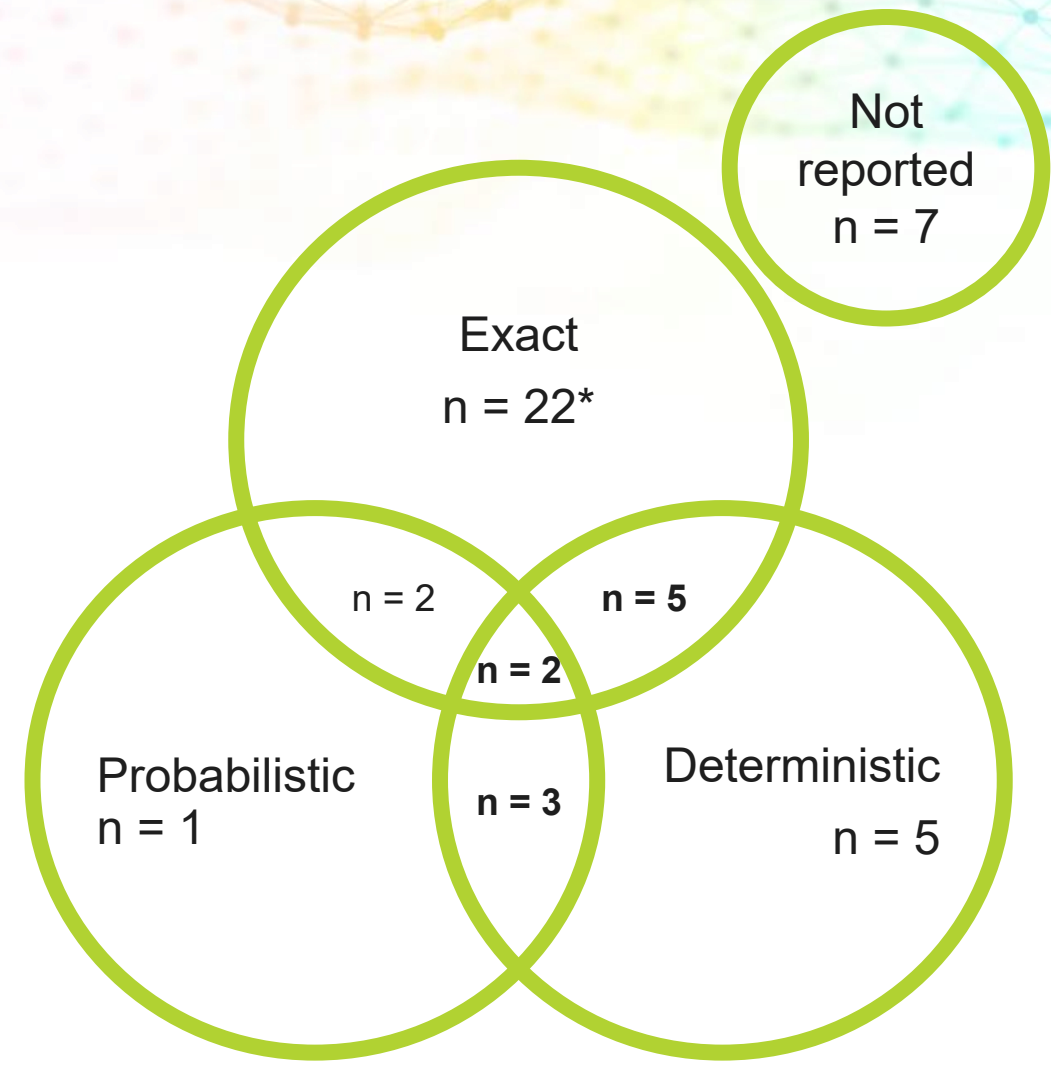
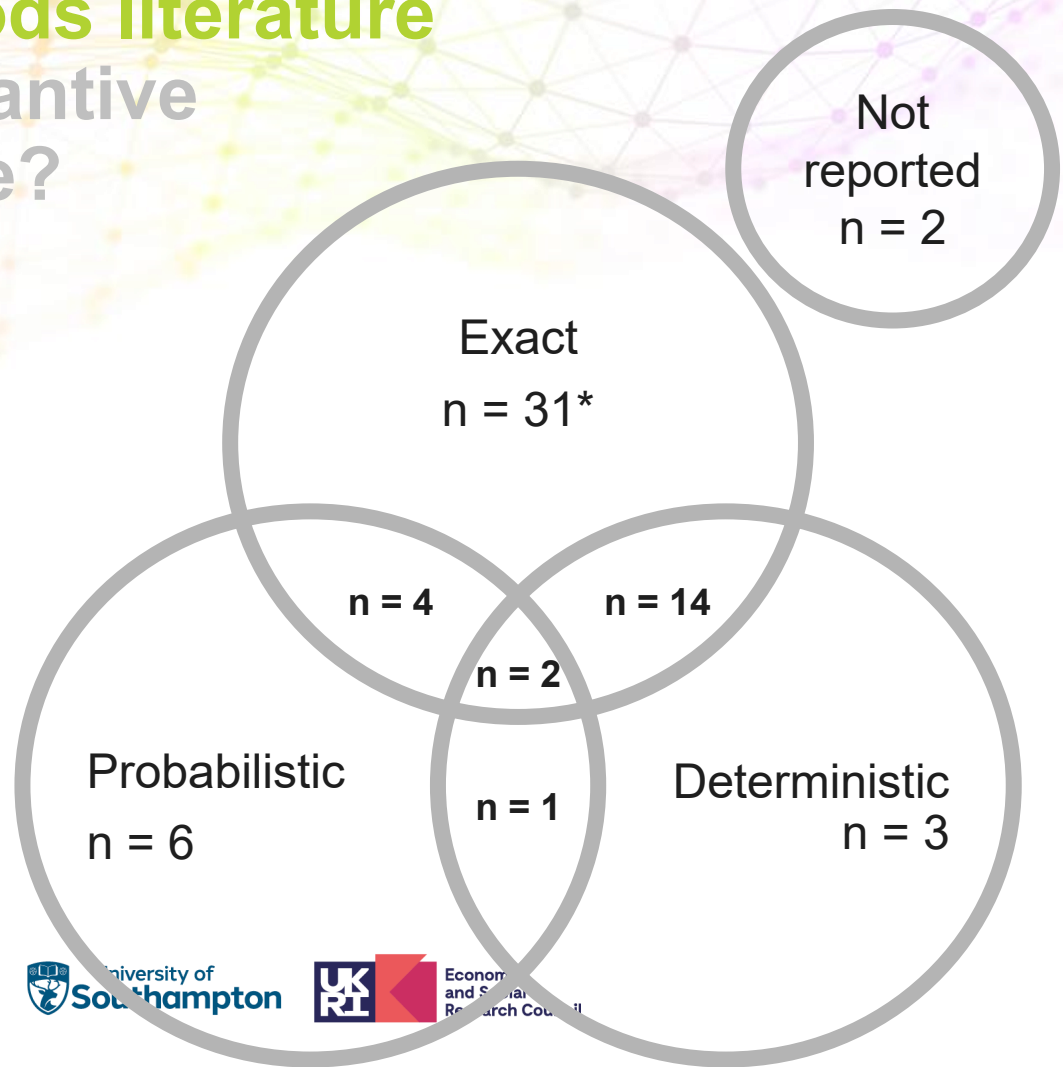


# RQ1: Which survey and non-survey data sources are commonly integrated in substantive literature?

	Administrative records	Smart tracker app	Place-based characteristics	Sensor data	Commercial data	Social media data	Total
<b>Panel study</b>							
<b>Birth cohort study</b>	22						22
<b>Repeated cross-sectional</b>	17	1	4				22
<b>Cross-sectional</b>	7				1		8
<b>Census</b>	28		1				29
<b>Total</b>	<b>74</b>	<b>1</b>	<b>5</b>		<b>1</b>		<b>81</b>

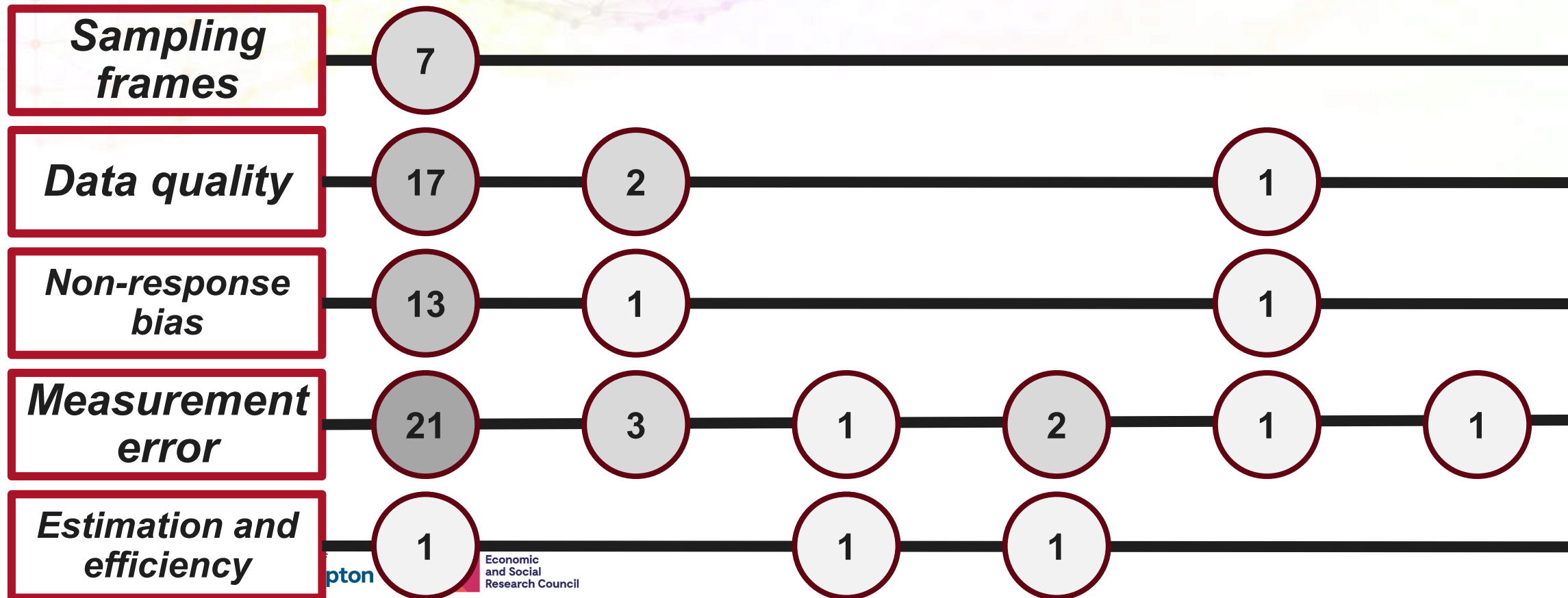


RQ2: Which methods are most frequently used to link survey and non-survey data in **methods literature** in substantive literature?



# RQ3: For what research purposes are integrated survey and non-survey data used?

	Administrative records	Smart tracker app	Place-based characteristics	Sensor data	Commercial data	Social media data
--	------------------------	-------------------	-----------------------------	-------------	-----------------	-------------------



# RQ4: What is the quality of reporting for data integration?

**Linkage and consent details were difficult to find and were inconsistently reported.**

**40.4%** Reported a response rate.

**31.9%** Reported a consent rate.

**55.3%** Mentioned consent to linkage.

**46.8%** Reported a linkage rate.

**42.5%** Mentioned whether consent was needed.

# Reporting guidelines for Integrated Data

## Survey data

(PRICSSA; Seidenberg,  
Moser & West, 2023)

## Non-survey data

- Title
- Type
- Sub-type
- Timeframe
- Population
- Citation
- DOI

## Integrated data information

- User guide/technical document
- Citation
- DOI

# *Reporting guidelines for Integrated Data*

## **Integration procedure**

- Exact matching
- Deterministic matching
- Probabilistic matching
- Statistical matching
- Linkage units

## **Integration quality**

- Linkage consent information
- Linkage consent rate
- Successful linkage rate
- Corrections applied

## **Integrated data access**

- Access conditions
- Access procedure



**SURVEY  
FUTURES**  
SURVEY DATA COLLECTION  
METHODS COLLABORATION

# Evaluating Data Quality in the Linked 1970 Birth Cohort Study



University of Essex



University of  
Southampton



Economic  
and Social  
Research Council

NCRM NATIONAL CENTRE FOR  
RESEARCH METHODS

Office for  
National Statistics

UCL

WARWICK  
THE UNIVERSITY OF WARWICK

MANCHESTER  
THE UNIVERSITY OF MANCHESTER

CITY  
UNIVERSITY OF BRISTOL

National Centre  
for Social Research

LSE  
LONDON SCHOOL OF  
ECONOMICS AND  
POLITICAL SCIENCE

Ulster  
University

Unil  
UNIL | Université de Lausanne

Ipsos

KANTAR  
PUBLIC

# Research Questions



What is the bias introduced by data linkage at each stage?

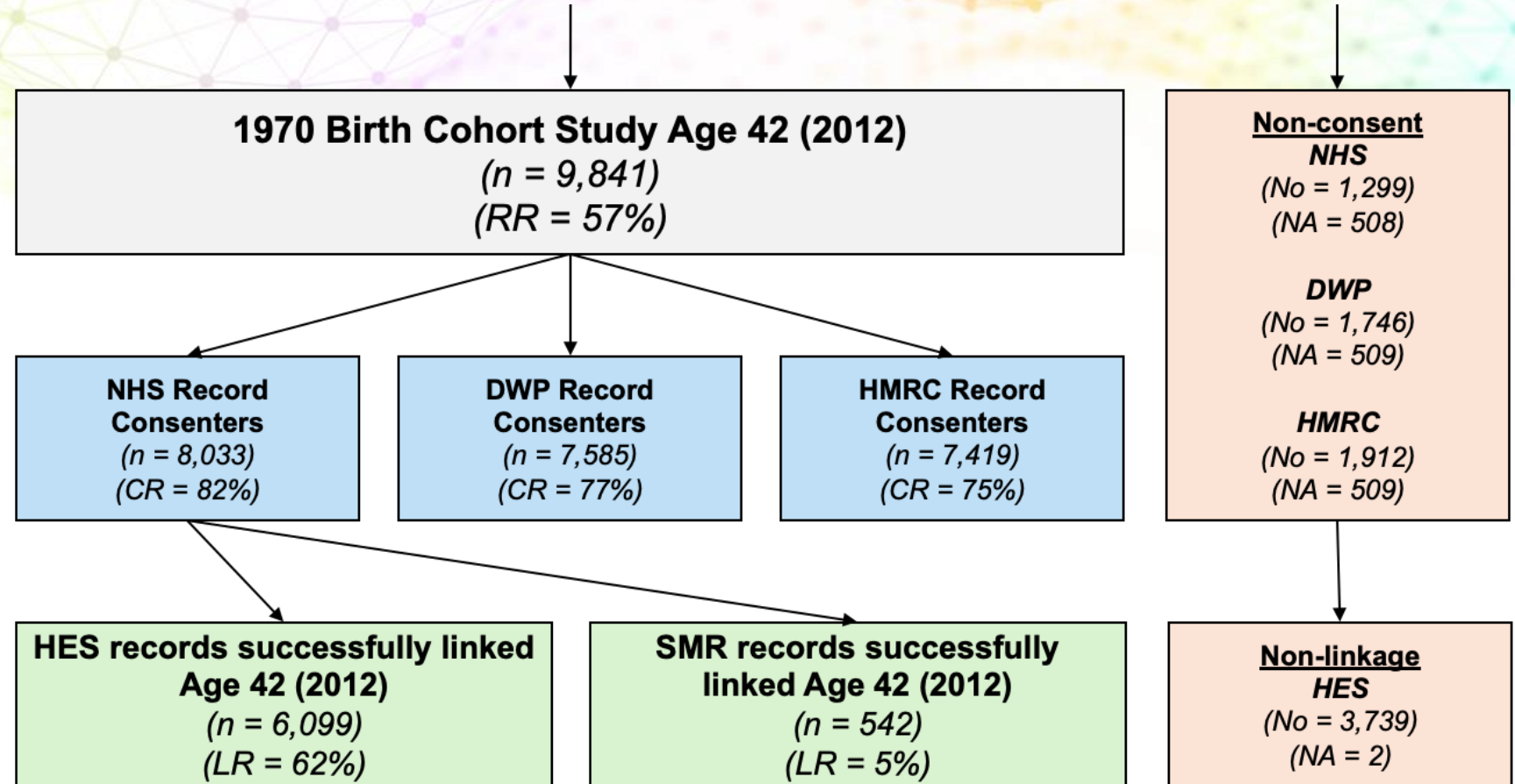
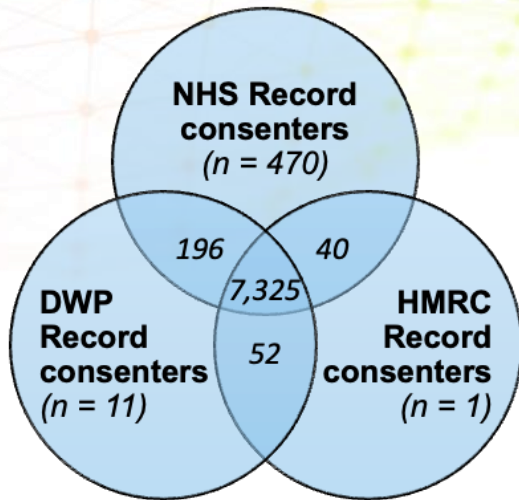


To what degree does this vary by domain/type of data?



Do differences in sample composition at each stage introduce bias to applied modelling?

# Data and selection process



# Predictors of consent

The strongest predictor of consent was **prior non-response**

**Renting** was positively associated with consent, as was lower **educational** attainment (GCSE/O-level).

**Former smoking status** was also positively associated with consent

Several **geographic regions** showed significantly higher odds of consent, particularly the East Midlands, North East, North West, South East, South West, West Midlands, Wales, and Yorkshire & Humber

The same set of variables were found to be associated with almost all linkage types

# Selection process by type of linkage

Dataset	R Indicator	Coefficient of Variation	n
Age 42 Sweep	0.82	0.11	9705
NHS Consenters	0.84	0.10	7957
HES Linkages	0.83	0.10	6041
SMR Linkages	<b>0.77</b>	<b>0.14</b>	<b>541</b>
DWP Consenters	0.84	0.09	7514
HMRC Consenters	0.84	0.09	7352

Note: Response propensities were estimated via binomial logistic regression modelling (Next wave response, as predicted by previous non-response, NHS consent, DWP consent, HMRC consent, sex, Socio-economic classification, household size, house ownership or renting, country of birth, highest educational qualification, marital status, smoking status, economic activity, benefit receipt, general health and government office region.



# SURVEY FUTURES

SURVEY DATA COLLECTION  
METHODS COLLABORATION

## *Does Additional Contact Information Improve Survey Outcomes?*

Hafsteinn Birgir Einarsson, Kim Backström & Alexandru Cernat



University of Essex



University of  
Southampton



Economic  
and Social  
Research Council

NCRM NATIONAL CENTRE FOR  
RESEARCH METHODS

Office for  
National Statistics

UCL

WARWICK  
THE UNIVERSITY OF WARWICK

MANCHESTER  
THE UNIVERSITY OF MANCHESTER

CITY  
UNIVERSITY OF THE CITY OF LONDON

National Centre  
for Social Research

LSE  
LONDON SCHOOL OF  
ECONOMICS AND  
POLITICAL SCIENCE

Ulster  
University

Unil  
UNIL | Université de Lausanne

Ipsos

KANTAR  
PUBLIC



What is the impact of additional contact information on participation?



How does the composition of the sample change due to extra contact info?

# Research Questions

# Other ongoing research



We investigate the effects of a **change in the Finnish Statistics Act** that granted Statistics Finland access to new contact information from the Finnish Social Insurance Institution.

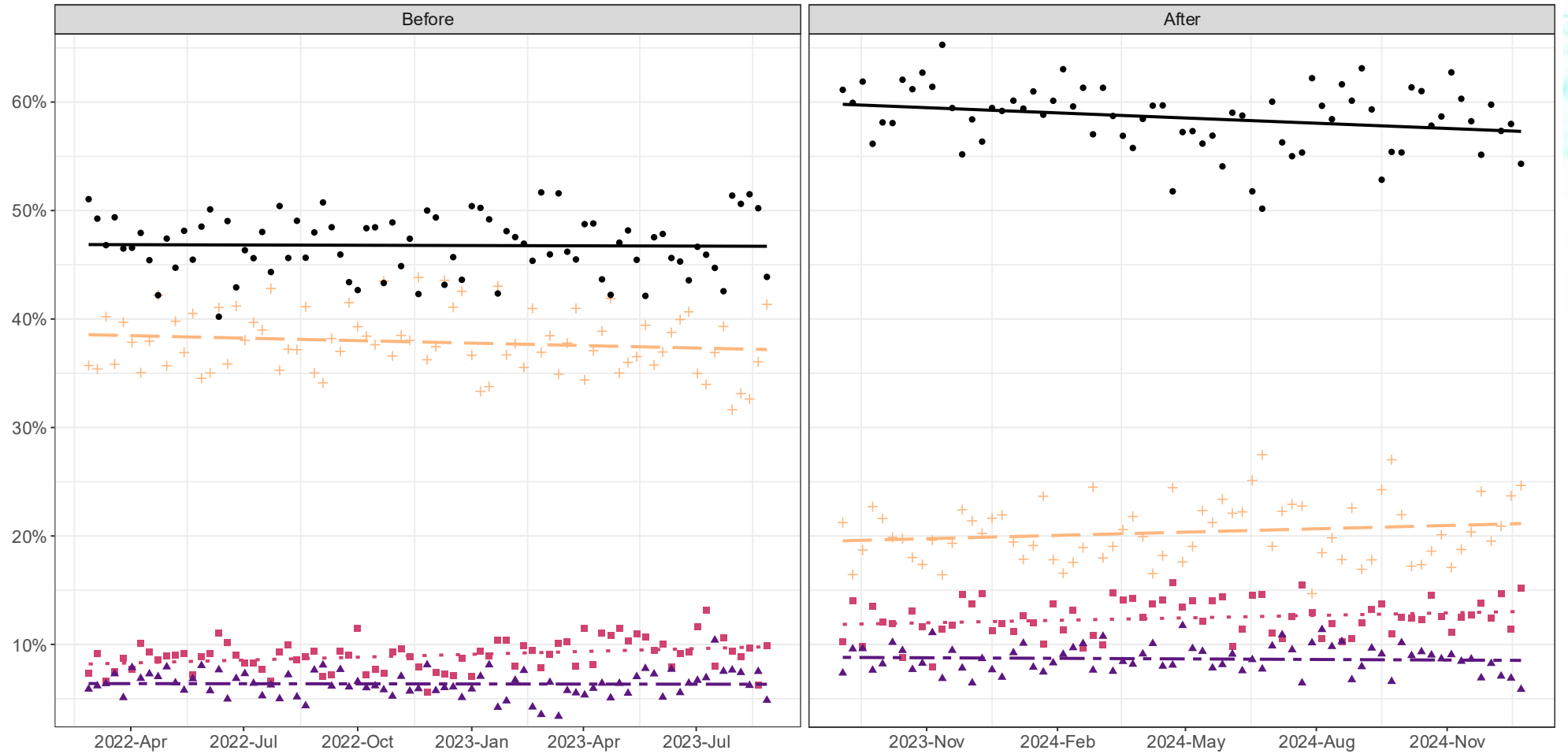
The additional contact information included

- phone numbers and
- email addresses

We treat **the legislative change as a quasi-experimental design** to investigate how extra contact information affected **response rates and selection bias in the Finnish LFS**

We look at **new samples of ~600 persons randomly selected from the population register each week before and after the policy change (148 weeks)**

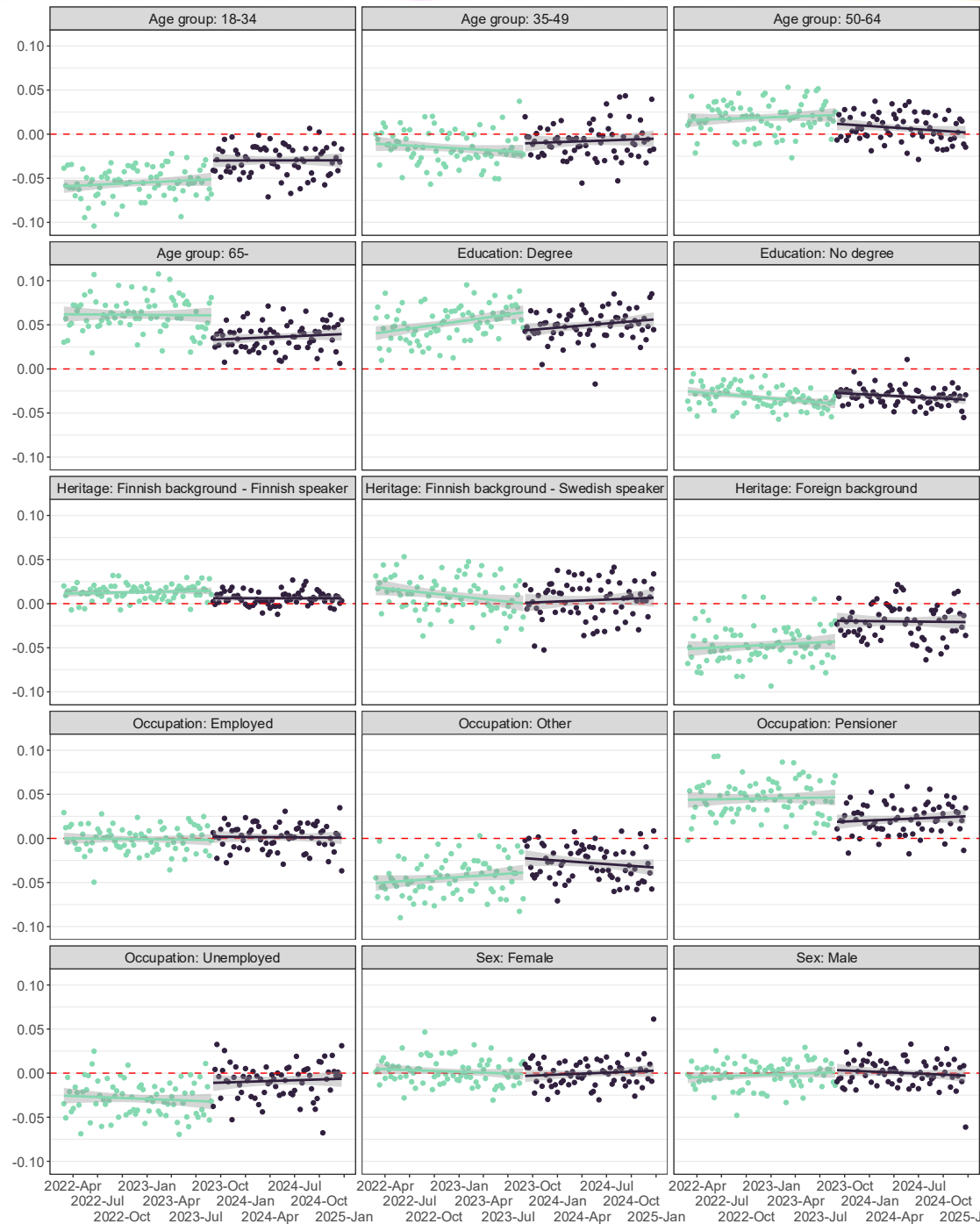
# Overall participation outcomes



# R-Indicators measuring the representativeness of respondent compositions in FI-LFS



# Partial R indicators show inconsistent changes in selection process



# Other ongoing research

**Peter** is leading a paper looking at how data at different levels (individuals, HHs, LSOA) predicts non-response in Understanding Society and how important it is to have up to date geographical data

**Natalie** is working on extending the classical Multi-Frame integration where both  $x_i$  and  $y_i$  are available in probability and nonprobability samples and same questionnaire applied to both samples

**Nikos** doing work around integrating survey and geospatial data



# SURVEY FUTURES

SURVEY DATA COLLECTION  
METHODS COLLABORATION

# RS7: Data Integration

**Alexandru Cernat**

**Thomas O'Toole, Peter Barlow, Natalie Shlomo, Nikos Tzavidis, Joseph Sakshaug**



# What search terms were used?

<b>Survey Data</b>	<b>AND</b>	<b>Non-Survey Data</b>	<b>AND</b>	<b>Data Integration</b>
<p>survey* OR "social survey*" OR "household survey*" OR "population representat*" OR population-based OR "nationally representat*" OR "random sampl*" OR "probability sampl*" OR questionair* OR interview* OR "cross-sectional stud*" OR "longitudinal stud*" OR "cohort stud*"</p>		<p>administrative OR operational OR "routinely collect*" OR "electronic record*" OR geospatial OR geographic OR spatial OR "gridded population*" OR "geographic information system*" OR satellite OR commercial OR "smart app*" OR "digital trace*" OR web-tracking OR "data donation" OR census</p>		<p>"data link*" OR "data integrat*" OR "record link*" OR "probabilistic link*" OR "deterministic link*" OR "data merg*" OR "data fus*" OR datalink* OR "multiple dataset*" OR "multiple data source*" OR "linkage stud*" OR "adaptive survey design*" OR "responsive survey design*" OR "total survey error*" OR "measurement error*"</p>

# Exclusion criteria

---

Not in English language

---

---

Not a peer-reviewed research  
paper

---

---

Not inside timeframe (2020 –  
current)

---

---

Not using both *survey and*  
non-survey data