

# When survey science met online tracking:

*Presenting an [error framework] for metered data*

**ORIOL J. BOSCH** | THE LONDON SCHOOL OF ECONOMICS / RECSM-UPF

**MELANIE REVILLA** | RECSM-UPF



*This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No849165), PI: Melanie Revilla*

# Tracking online behaviours using a meter

## Definition

**Metered data** is obtained from a meter willingly installed or configured by a sample of participants on their devices (PCs, tablets and/or smartphones).

A **meter** refers to a heterogeneous group of tracking technologies that allow sharing with the researchers, at least, **information about the URLs of the web pages visited by the participants.**

### Sample of participants

Collected from a designed sample of individuals

### Nonreactive

Collected by tracking the traces left by individuals when interacting with their devices online.

# Tracking online behaviours using a meter

## Benefits of metered data

- Objective and free of recall errors
- Continuously collected in real time
- Pre-designed sample of participants



BACKGROUND

# Metered data in past research



## The sources and correlates of exposure to vaccine-related (mis)information online<sup>☆</sup>

Andrew M. Guess<sup>a,\*</sup>, Brendan Nyhan<sup>b</sup>, Zachary O’Keefe<sup>c</sup>, Jason Reifler<sup>d</sup>

<sup>a</sup> Department of Politics, Princeton University, United States

<sup>b</sup> Department of Government, Dartmouth College, United States

<sup>c</sup> Department of Political Science, University of Michigan, United States

<sup>d</sup> Department of Politics, University of Exeter, United Kingdom

### ARTICLE INFO

Article history:  
Received 11 June 2020  
Received in revised form 1 October 2020  
Accepted 7 October 2020  
Available online 22 October 2020

Keywords:  
Vaccine hesitancy  
Vaccine skepticism  
Online  
Information  
Social media  
Search

### ABSTRACT

**Objectives:** To assess the quantity and type of vaccine-related information Americans consume online and its relationship to social media use and attitudes toward vaccines.  
**Methods:** Analysis of individual-level web browsing data linked with survey responses from representative samples of Americans collected between October 2016 and February 2019.

**Results:** We estimate that approximately 84% of Americans visit a vaccine-related webpage each year. Encounters with vaccine-skeptical content are less frequent; they make up only 7.5% of vaccine-related pageviews and are encountered by only 18.5% of people annually. However, these pages are more likely to be published by untrustworthy sources. Moreover, skeptical content exposure is more common among people with less favorable vaccine attitudes. Finally, usage of online intermediaries is frequently linked to vaccine-related information exposure. Google use is differentially associated with subsequent exposure to non-skeptical content, whereas exposure to vaccine-skeptical webpages is associated with usage of webmail and, to a lesser extent, Facebook.

**Conclusions:** Online exposure to vaccine-skeptical content is relatively rare, but vigilance is required given the potential for exposure among vulnerable audiences.

© 2020 Elsevier Ltd. All rights reserved.

Published in final edited form as:

*Nat Hum Behav.* 2020 March 02; 4(S): 472–480. doi:10.1038/s41562-020-0833-x.

## Exposure to untrustworthy websites in the 2016 U.S. election

Andrew M. Guess<sup>1</sup>, Brendan Nyhan<sup>2,\*</sup>, Jason Reifler<sup>3</sup>

<sup>1</sup>Department of Politics and Woodrow Wilson School, Princeton University, Princeton, NJ, USA

<sup>2</sup>Department of Government, Dartmouth College, Hanover, NH, USA

<sup>3</sup>Department of Politics, University of Exeter, Exeter, UK

### Abstract

Though commentators frequently warn about “echo chambers,” little is known about the volume or slant of political misinformation people consume online, the effects of social media and fact-checking on exposure, or its effects on behavior. We evaluate these questions for the websites publishing factually dubious content often described as “fake news.” Survey and web traffic data from the 2016 U.S. presidential campaign show that Trump supporters were most likely to visit these websites, which often spread via Facebook. However, these sites made up a small share of people’s informational diets on average and were largely consumed by a subset of Americans with strong preferences for pro-attitudinal information. These results suggest that widespread speculation about the prevalence of exposure to untrustworthy websites has been overstated.

# Predicting Voting Behavior Using Digital Trace Data

Ruben L. Bach<sup>1</sup>, Christoph Kern<sup>1</sup>, Ashley Amaya<sup>2</sup>, Florian Keusch<sup>1</sup>, Frauke Kreuter<sup>1,3,4</sup>, Ian Hecht<sup>5</sup> and Jonathan Heinemann<sup>6</sup>

## Abstract

A major concern arising from ubiquitous data is how to predict personal sensitive information. Although previous research on sociodemographic characteristics, little is known about how to predict voting behavior. Against this background, we investigate how to predict voting behavior, which is considered to be sensitive information. Using reconstructions of online users' browsing behavior, we find that the same individuals, who vote in the 2017 US midterm election, also visit the same information flows. These findings add to the literature on digital trace data.



journal homepage:

The sources and correlates of exposure to (mis)information online\*

Andrew M. Guess<sup>a,b</sup>, Brendan Nyhan<sup>b</sup>, Zachary

<sup>a</sup>Department of Politics, Princeton University, United States  
<sup>b</sup>Department of Government, Dartmouth College, United States  
<sup>c</sup>Department of Political Science, University of Michigan, United States  
<sup>d</sup>Department of Politics, University of Exeter, United Kingdom

## ARTICLE INFO

Article history:  
 Received 11 June 2020  
 Received in revised form 1 October 2020  
 Accepted 7 October 2020  
 Available online 22 October 2020

Keywords:  
 Vaccine hesitancy  
 Vaccine skepticism  
 Online information  
 Social media  
 Search

## ABSTRACT

**Objectives:** To assess the quantity and type of vaccine-related information Americans consume online and its relationship to social media use and attitudes toward vaccines.  
**Methods:** Analysis of individual-level web browsing data linked with survey responses from representative samples of Americans collected between October 2016 and February 2019.  
**Results:** We estimate that approximately 84% of Americans visit a vaccine-related webpage each year. Encounters with vaccine-skeptical content are less frequent; they make up only 7.5% of vaccine-related pageviews and are encountered by only 18.5% of people annually. However, these pages are more likely to be published by untrustworthy sources. Moreover, skeptical content exposure is more common among people with less favorable vaccine attitudes. Finally, usage of online intermediaries is frequently linked to vaccine-related information exposure. Google use is differentially associated with subsequent exposure to non-skeptical content, whereas exposure to vaccine-skeptical webpages is associated with usage of webmail and, to a lesser extent, Facebook.  
**Conclusions:** Online exposure to vaccine-skeptical content is relatively rare, but vigilance is required given the potential for exposure among vulnerable audiences.

© 2020 Elsevier Ltd. All rights reserved.

Social Science Computer Review

1-22

© The Author(s) 2019



Article reuse guidelines:  
[sagepub.com/journalsPermissions](http://sagepub.com/journalsPermissions)  
 DOI: 10.1177/0894439319882896  
[journals.sagepub.com/home/fssc](http://journals.sagepub.com/home/fssc)



International Journal of Public Opinion Research Vol. 31 No. 4 2019  
 © The Author(s) 2018. Published by Oxford University Press on behalf of The World Association for Public Opinion Research. All rights reserved.  
[doi:10.1093/ijpor/fdy035](https://doi.org/10.1093/ijpor/fdy035) Advance Access publication 15 December 2018

# Is Facebook Eroding the Public Agenda? Evidence From Survey and Web-Tracking Data

Ana S. Cardenal<sup>1</sup>, Carol Galais<sup>2</sup>, and Silvia Maió-Vázquez<sup>3</sup>

<sup>1</sup>Oberta de Catalunya, Spain;  
<sup>2</sup>Autònoma de Barcelona, Spain;  
<sup>3</sup>University of Oxford, UK

## The consequences of online partisan media

Andrew M. Guess<sup>a,b,1,2</sup>, Pablo Barberá<sup>c,1</sup>, Simon Munzert<sup>d,1</sup>, and JungHwan Yang (양정환)<sup>e,1,2</sup>

<sup>a</sup>Department of Politics, Princeton University, Princeton, NJ 08544; <sup>b</sup>School of Public and International Affairs, Princeton University, Princeton, NJ 08544; <sup>c</sup>Department of Political Science and International Relations, University of Southern California, Los Angeles, CA 90089; <sup>d</sup>Data Science Lab, Hertie School, 10117 Berlin, Germany; and <sup>e</sup>Department of Communication, University of Illinois at Urbana-Champaign, Urbana, IL 61801

Edited by Christopher Andrew Bail, Duke University, Durham, NC, and accepted by Editorial Board Member Margaret Levi February 17, 2021 (received for review June 29, 2020)

What role do ideologically extreme media play in the polarization of society? Here we report results from a randomized longitudinal field experiment embedded in a nationally representative online panel survey ( $N = 1,037$ ) in which participants were incentivized to change their browser default settings and social media following patterns, boosting the likelihood of encountering news with either a left-leaning (HuffPost) or right-leaning (Fox News) slant during the 2018 US midterm election campaign. Data on  $\approx 19$  million web visits by respondents indicate that resulting changes in news consumption persisted for at least 8 wk. Greater exposure to partisan news can cause immediate but short-lived increases in website visits and knowledge of recent events. After adjusting for multiple comparisons, however, we find little evidence of a direct impact on opinions or affect. Still, results from later survey waves suggest that both treatments produce a lasting and meaningful decrease in trust in the mainstream media up to 1 y later. Consistent with the minimal-effects tradition, direct consequences of online partisan media are limited, although our findings raise questions about the possibility of subtle, cumulative dynamics. The combination of experimentation and computational social science techniques illustrates a powerful approach for studying the long-term consequences of exposure to partisan news.

argues that media primarily reinforce existing predispositions (16). At the same time, more recent research strongly implies that newspapers and especially cable news can change people's voting behavior, especially those without strong partisan attachments (17–20). We propose an internet-age synthesis that views people's information environments through the lens of choice architecture (21): frictions, subtle design features, and default settings that structure people's online experience. In this view, small changes (or nudges) could disproportionately affect information consumption habits that have downstream consequences.

To that end, we designed a large, longitudinal online field experiment that subtly but naturalistically increased people's exposure to partisan news websites. Our choice of treatment is ecologically valid: Despite the importance of social media for agenda-setting (22) and public expression (23), more Americans continue to say that they get news from news websites or apps than social media sites (24). The intervention thus served as a nudge, boosting the likelihood that subjects encountered news framed with a partisan slant during their day-to-day web browsing experience, even if inadvertently. The powerful, sustained nature of the intervention and our ability to track participants with survey and behavioral data for months provided the opportunity to test a range of hypotheses about the long-term impact

of social integration, minimizing media are known for fragmenting attention to their effect on the timing news through Facebook of most important problems search design combines survey ferred news consumption influences people's information environments through the lens of when Facebook is a relevant news attachments (21): frictions, subtle design features, and default settings that structure people's online experience. In ; of our findings for the public this view, small changes (or nudges) could disproportionately affect information consumption habits that have downstream consequences.

## 2016 U.S. election

Princeton University, Princeton, NJ, USA  
 JSA

Though commentators frequently warn about “echo chambers,” little is known about the volume or slant of political misinformation people consume online, the effects of social media and fact-checking on exposure, or its effects on behavior. We evaluate these questions for the websites publishing factually dubious content often described as “fake news.” Survey and web traffic data from the 2016 U.S. presidential campaign show that Trump supporters were most likely to visit these websites, which often spread via Facebook. However, these sites made up a small share of people's information diets on average and were largely consumed by a subset of Americans with strong preferences for pro-attitudinal information. These results suggest that widespread speculation about the prevalence of exposure to untrustworthy websites has been overstated.

## How Much Time Do You Spend Online? Understanding and Improving the Accuracy of Self-Reported Measures of Internet Use

Theo Araujo, Anke Wonneberger, Peter Neijens, and Claes de Vreese

Amsterdam School of Communication Research (ASCoR), University of Amsterdam, Amsterdam, The Netherlands

### ABSTRACT

Given the importance of survey measures of online media use for communication research, it is crucial to assess and improve their quality, in particular because the increasingly fragmented and ubiquitous usage of internet complicates the accuracy of self-reported measures. This study contributes to the discussion regarding the accuracy of self-reported internet use by presenting relevant factors potentially affecting biases of self-reports and testing survey design strategies to improve accuracy. Combining automatic tracking data and survey data from the same participants (N = 690) confirmed low levels of accuracy and tendencies of over-reporting. The analysis revealed biases due to a range of factors associated with the intensity of (actual) internet usage, propensity to multitask, day of reference, and the usage of mobile devices. An anchoring technique could not be proved to reduce inaccuracies of reporting behavior. Several recommendations for research practice follow from these findings.



journal homepage:

The sources and correlates of exposure (mis)information online\*

Andrew M. Guess<sup>a,b</sup>, Brendan Nyhan<sup>b</sup>, Zachary

<sup>a</sup>Department of Politics, Princeton University, United States  
<sup>b</sup>Department of Government, Dartmouth College, United States  
<sup>c</sup>Department of Political Science, University of Michigan, United States  
<sup>d</sup>Department of Politics, University of Exeter, United Kingdom

### ARTICLE INFO

Article history:  
Received 11 June 2020  
Received in revised form 1 October 2020  
Accepted 9 October 2020  
Available online 22 October 2020

Keywords:  
Vaccine hesitancy  
Vaccine skepticism  
Online information  
Social media  
Search

### ABSTRACT

**Objectives:** To assess the quantity and its relationship to social media use  
**Methods:** Analysis of individual-level web site samples of Americans collected by Results: We estimate that approximate Encounters with vaccine-skeptical or related pageviews and are encountered likely to be published by untrustworthy among people with less favorable vacci linked to vaccine-related information e exposure to non-skeptical content, wht usage of webmail and, to a lesser exten  
**Conclusions:** Online exposure to vaccin given the potential for exposure among

resuring th  
wk. Greater  
short-lived i  
events. Afte  
find little ev  
results from  
duce a lastin  
media up to  
tion, direct  
although ou  
tle, cumulat  
and computi  
approach fo  
partisan nev

COMMUNICATION METHODS AND MEASURES  
2016, VOL. 10, NO. 1, 13–27  
<http://dx.doi.org/10.1080/19312458.2015.1118446>

## The Accuracy of Self-Reported Internet Use—A Validation Study Using Client Log Data

Michael Scharnow

University of Hohenheim

### ABSTRACT

The vast majority of empirical research on online communication, or media use in general, relies on self-report measures instead of behavioral data. Previous research has shown that the accuracy of these self-report measures can be quite low, and both over- and underreporting of media use are commonplace. This study compares self-reports of Internet use with client log files from a large household sample. Results show that the accuracy of self-reported frequency and duration of Internet use is quite low, and that survey data are only moderately correlated with log file data. Moreover, there are systematic patterns of misreporting, especially overreporting, rather than random deviations from the log files. Self-reports for specific content such as social network sites or video platforms seem to be more accurate and less consistently biased than self-reports of generic frequency or duration of Internet use. The article closes by demonstrating the consequences of biased self-reports and discussing possible solutions to the problem.

## Two Half-Truths Make a Whole? On Bias in Self-Reports and Tracking Data

Pascal Jürgens<sup>1</sup>, Birgit Stark<sup>1</sup>, and Melanie Magin<sup>2</sup>

### Abstract

The pervasive use of mobile information technologies brings new patterns of media usage, but also challenges to the measurement of media exposure. Researchers wishing to, for example, understand the nature of selective exposure on algorithmically driven platforms need to precisely attribute individuals' exposure to specific content. Prior research has used tracking data to show that survey-based self-reports of media exposure are critically unreliable. So far, however, little effort has been invested into assessing the specific biases of tracking methods themselves. Using data from a multimethod study, we show that tracking data from mobile devices is linked to systematic distortions in self-report biases. Further inherent but unobservable sources of bias, along with potential solutions, are discussed.

the view small changes for modest could disproportionate downstream



online field  
sed people's  
'treatment is  
al media for  
re Americans  
sites or apps  
served as a  
intered news  
y web brows-  
ul, sustained  
participants  
d the oppor-  
-term impact

## 2016 U.S. election

University, Princeton, NJ, USA

JSA

'shambers,' little is known about the volume online, the effects of social media and fact-evaluate these questions for the websites as "fake news." Survey and web traffic data Trump supporters were most likely to visit wever, these sites made up a small share of ly consumed by a subset of Americans with ese results suggest that widespread istworthy websites has been overstated.

# Inferences for finite populations

## Metered data can potentially suffer from different types of errors

Shared devices and observation of only part of the activity

- 60% of desktops, 40% of laptops and tablets, and 9% of smartphones shared to some degree (Revilla et al., 2017)
- 28% with the meter installed in all devices (Pew Research Center, 2020)

Technical issues and reactivity / social desirability bias (Jurgens et al., 2020; Toth and Trifonova, 2020)



# Inferences for finite populations

## Metered data can potentially suffer from different types of errors

Shared devices and observation of only part of the activity

- 60% of desktops, 40% of laptops and tablets, and 9% of smartphones shared to some degree (Revilla et al., 2017)
- 28% with the meter installed in all devices (Pew Research Center, 2020)

Technical issues and reactivity / social desirability bias (Jurgens et al., 2020; Toth and Trifonova, 2020)



Systematic **categorization** and **conceptualization** of metered data errors  
not available

## Total Error Framework for metered data

#1 **Summarize** the data collection and analysis process for metered data.

#2 **Conceptualize and categorize** all errors components (e.g. *measurement errors*) and causes (e.g. *social desirability*) that can occur when using metered data.

## Total Error Framework for metered data

#1 **Summarize** the data collection and analysis process for metered data.

#2 **Conceptualize and categorize** all errors components (e.g. *measurement errors*) and causes (e.g. *social desirability*) that can occur when using metered data.

- 1) Choose the best design options for metered data.
- 2) Make better informed decisions while planning when and how to supplement or replace survey data with metered data.
- 3) Help assess research using metered data.

## Total Error Framework for metered data

#1 **Summarize** the data collection and analysis process for metered data.

#2 **Conceptualize and categorize** all errors components (e.g. *measurement errors*) and causes (e.g. *social desirability*) that can occur when using metered data.

Bosch, O.J., and M. Revilla (2021). “**When survey science met online tracking: presenting an error framework for metered data.**” RECSM Working Papers Series, 62

## Adapting instead of reinventing

- Follow approach by Amaya et al (2020) with their **Total Error Framework** for **Big Data**
- 7 error components of the **TSE (Groves et al., 2009)** as starting point:
  - Coverage errors, sampling errors, *missing data errors*, adjustment errors, *specification errors*, measurement errors and processing errors

RESULTS

# Data collection and analysis process

**Define concept of  
interest**

*Average hours of consumption of  
online political news*

**Define concept of  
interest**



**Design  
measurement**

*Average hours of consumption of  
online political news*

*Average time recorded of the visits to  
online political outlets' URLs.*



**Define concept of interest**



**Design measurement**



**Develop/ choose the technology**

*Average hours of consumption of online political news*

*Average time recorded of the visits to online political outlets' URLs.*

*Proxy for IOS/ App for others*

**Define concept of interest**



**Design measurement**



**Develop/ choose the technology**

*Average hours of consumption of online political news*

*Average time recorded of the visits to online political outlets' URLs.*

*Proxy for IOS/ App for others*

**Define target inferential population**

*People living in the UK older than 18*

**Define concept of interest**

*Average hours of consumption of online political news*

**Design measurement**

*Average time recorded of the visits to online political outlets' URLs.*

**Develop/ choose the technology**

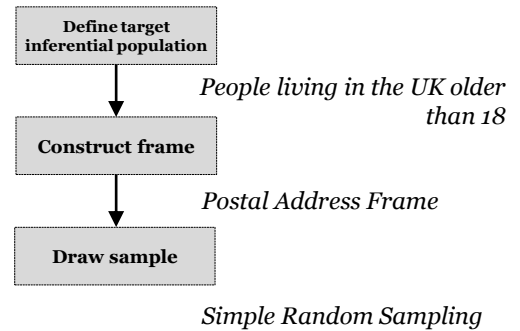
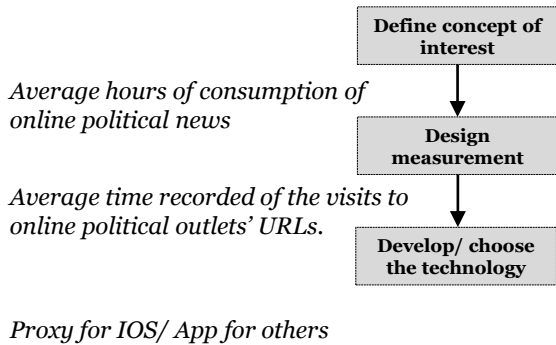
*Proxy for IOS/ App for others*

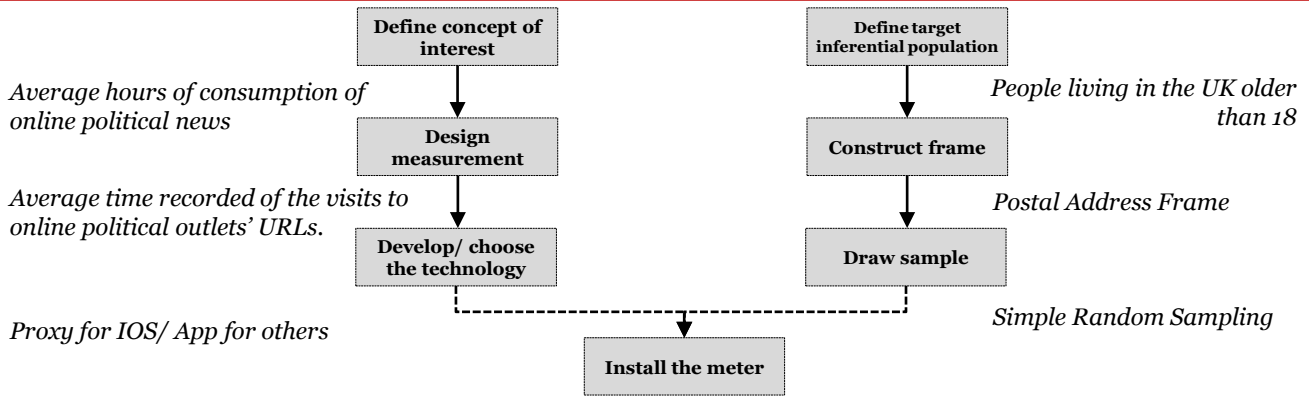
**Define target inferential population**

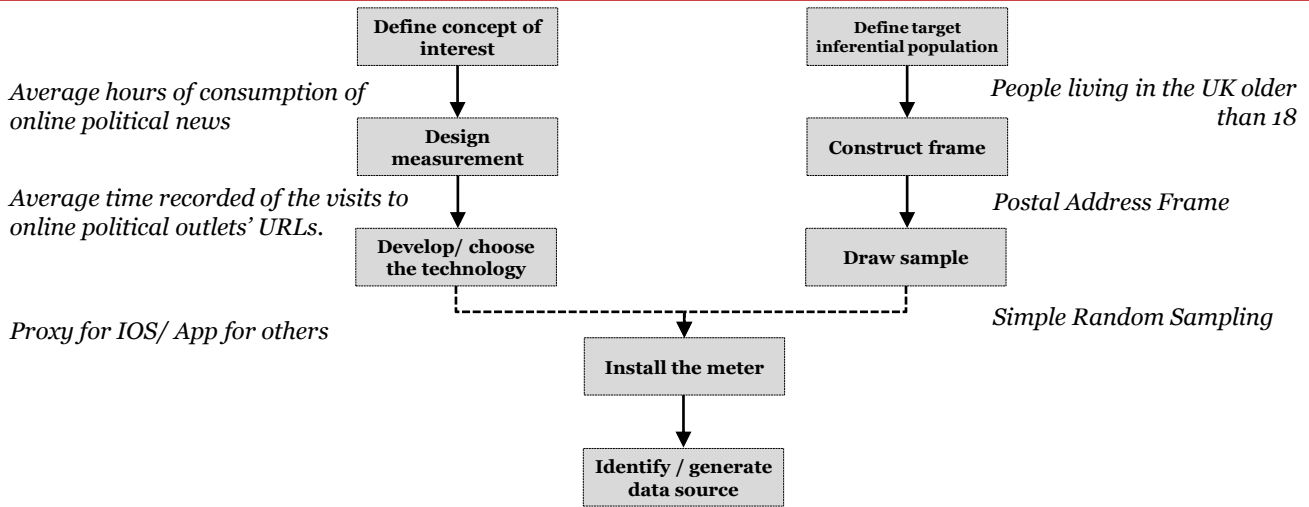
*People living in the UK older than 18*

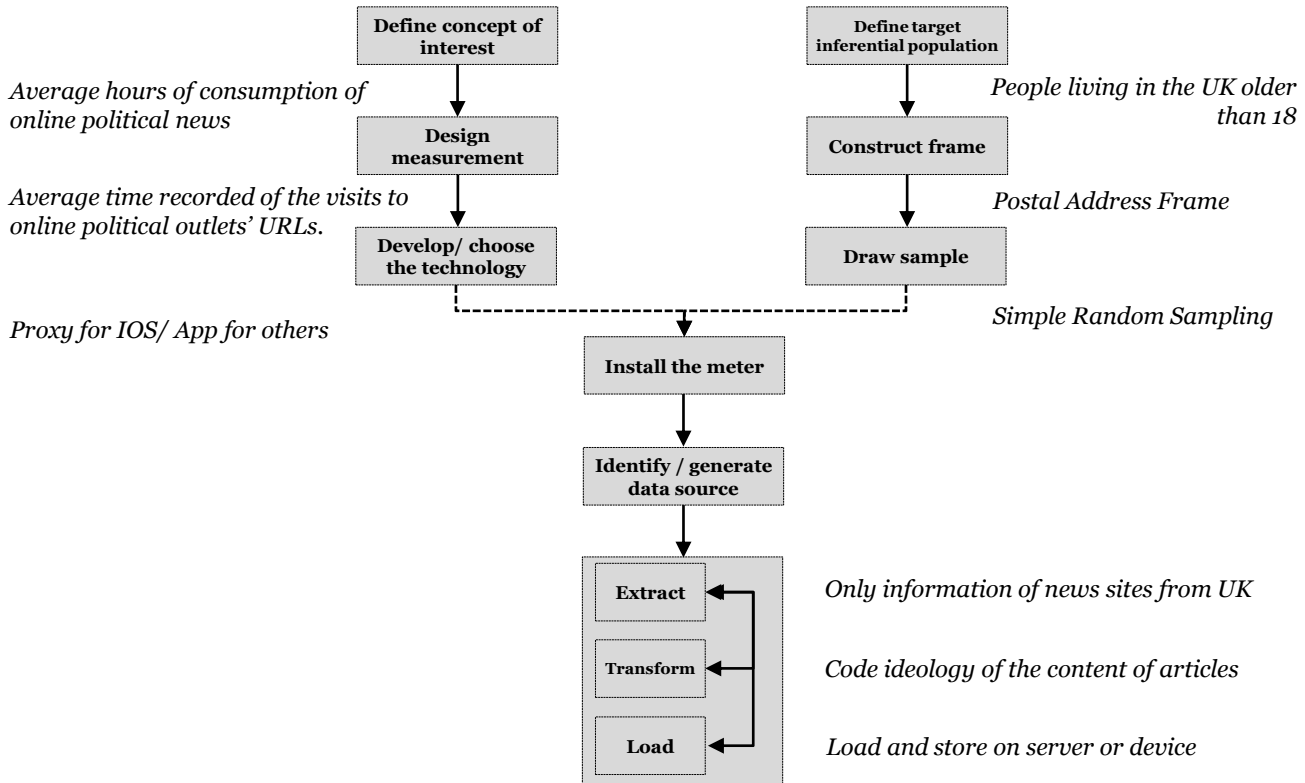
**Construct frame**

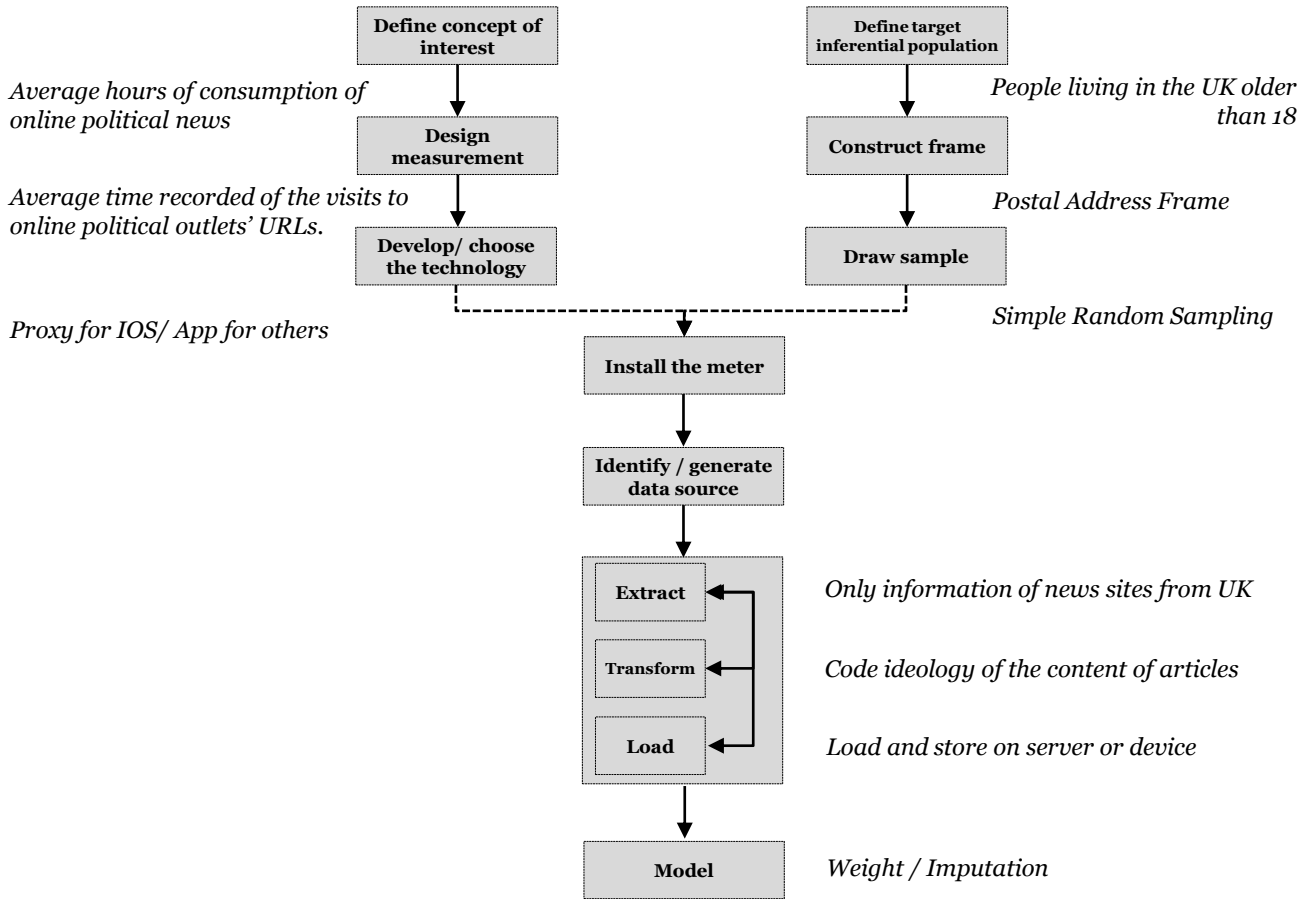
*Postal Address Frame*



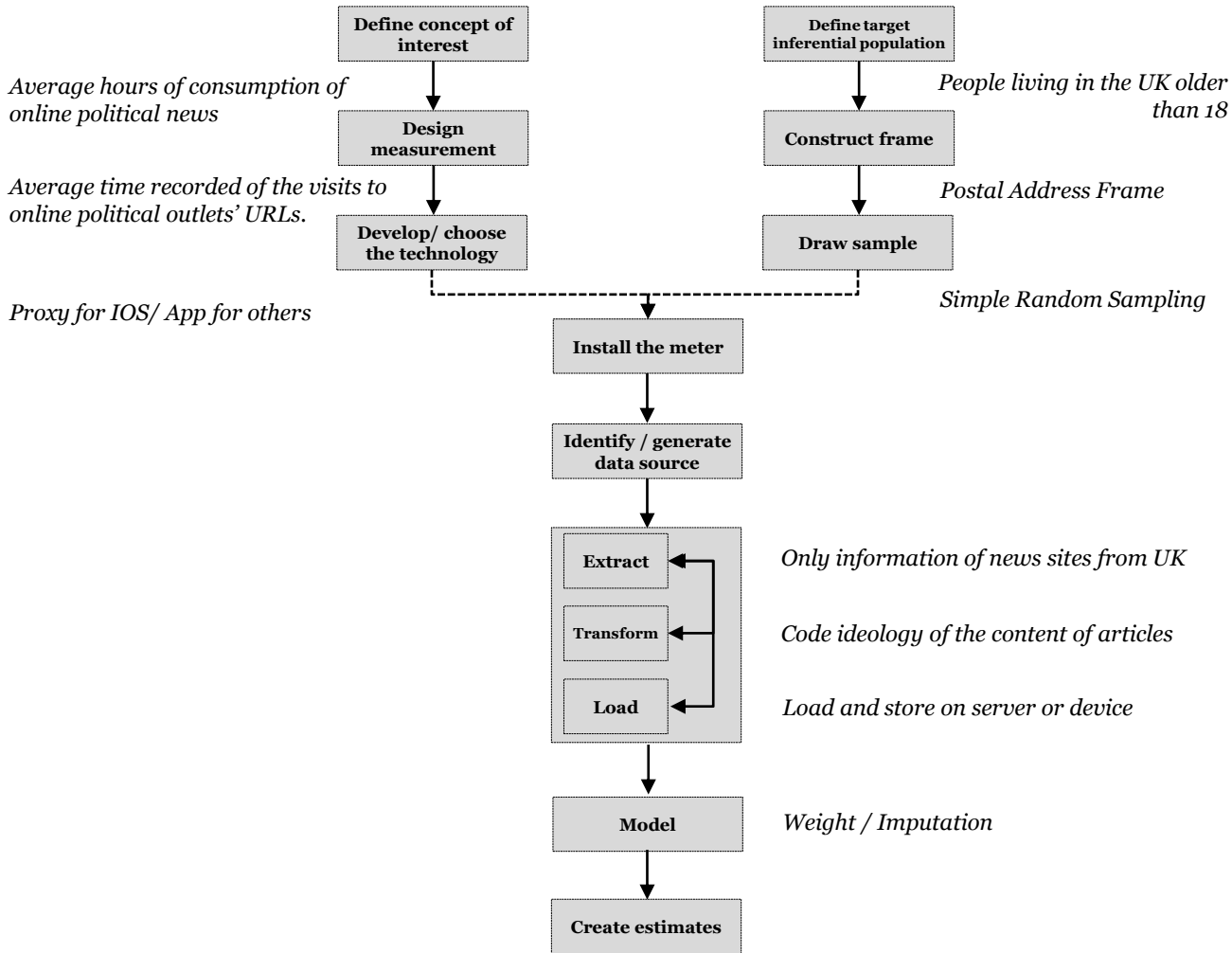


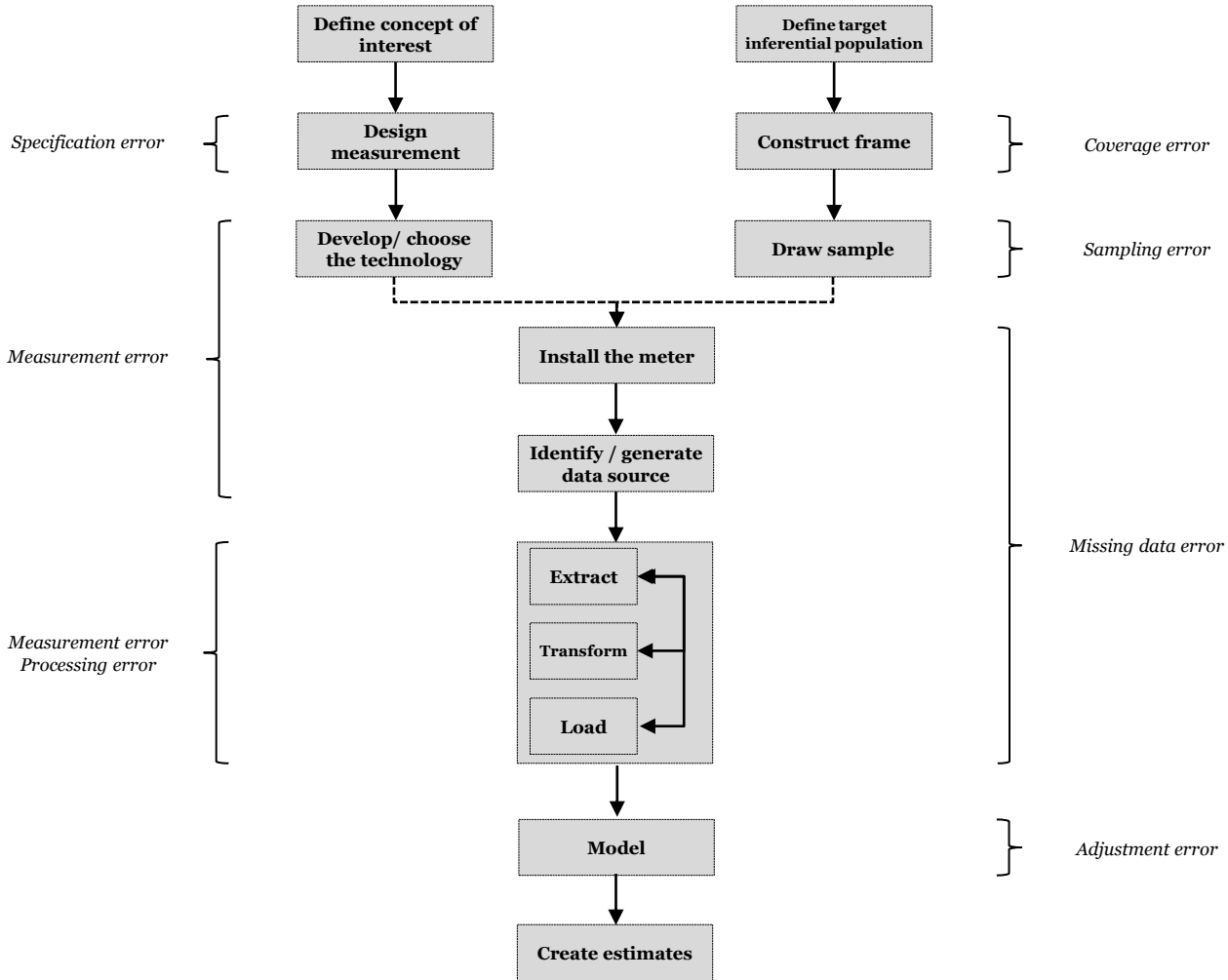












## RESULTS

## Error components and their causes

Error components	Specific error causes
Specification error	<ul style="list-style-type: none"> <li>– Measuring concepts from which not enough data is available</li> <li>– Inferring attitudes</li> <li>– Defining valid information</li> </ul>
Measurement error	<ul style="list-style-type: none"> <li>– Non-trackable target</li> <li>– Meter not installed</li> <li>– Uninstalling the meter</li> <li>– New non-tracked device</li> <li>– Technology limitations</li> <li>– Technology errors</li> <li>– Hidden behaviours</li> <li>– Shared device</li> <li>– Social desirability</li> <li>– Extraction error</li> </ul>
Processing error	<ul style="list-style-type: none"> <li>– Coding error</li> <li>– Aggregation at the domain level</li> <li>– Data anonymization</li> </ul>
Coverage error	<ul style="list-style-type: none"> <li>– Non-trackable individuals</li> </ul>
Sampling error	<ul style="list-style-type: none"> <li>– Same error causes than for surveys</li> </ul>
Missing data error	<ul style="list-style-type: none"> <li>– Noncontact</li> <li>– Non-consent</li> <li>– Non-trackable target</li> <li>– Meter not installed</li> <li>– Uninstalling the meter</li> <li>– New non-tracked device</li> <li>– Technology limitations</li> <li>– Technology error</li> <li>– Hidden behaviour</li> <li>– Social desirability</li> <li>– Extraction error</li> </ul>
Adjustment error	<ul style="list-style-type: none"> <li>– Same error causes than for surveys</li> </ul>

## RESULTS

# Error components and their causes

Error components	Specific error causes
Specification error	<ul style="list-style-type: none"> <li>- Measuring concepts from which not enough data is available</li> <li>- Inferring attitudes</li> <li>- Defining valid information</li> </ul>
Measurement error	<ul style="list-style-type: none"> <li>- Non-trackable target</li> <li>- Meter not installed</li> <li>- Uninstalling the meter</li> <li>- New non-tracked device</li> <li>- Technology limitations</li> <li>- Technology errors</li> <li>- Hidden behaviours</li> <li>- Shared device</li> <li>- Social desirability</li> <li>- Extraction error</li> </ul>
Processing error	<ul style="list-style-type: none"> <li>- Coding error</li> <li>- Aggregation at the domain level</li> <li>- Data anonymization</li> </ul>
Coverage error	<ul style="list-style-type: none"> <li>- Non-trackable individuals</li> </ul>
Sampling error	<ul style="list-style-type: none"> <li>- Same error causes than for surveys</li> </ul>
Missing data error	<ul style="list-style-type: none"> <li>- Noncontact</li> <li>- Non-consent</li> <li>- Non-trackable target</li> <li>- Meter not installed</li> <li>- Uninstalling the meter</li> <li>- New non-tracked device</li> <li>- Technology limitations</li> <li>- Technology error</li> <li>- Hidden behaviour</li> <li>- Social desirability</li> <li>- Extraction error</li> </ul>
Adjustment error	<ul style="list-style-type: none"> <li>- Same error causes than for surveys</li> </ul>

Most specific error causes on the side of measurement

## RESULTS

## Error components and their causes

Error components	Specific error causes
Specification error	<ul style="list-style-type: none"> <li>- Measuring concepts from which not enough data is available</li> <li>- Inferring attitudes</li> <li>- Defining valid information</li> </ul>
Measurement error	<ul style="list-style-type: none"> <li>- Non-trackable target</li> <li>- Meter not installed</li> <li>- Uninstalling the meter</li> <li>- New non-tracked device</li> <li>- Technology limitations</li> <li>- Technology errors</li> <li>- Hidden behaviours</li> <li>- Shared device</li> <li>- Social desirability</li> <li>- Extraction error</li> </ul>
Processing error	<ul style="list-style-type: none"> <li>- Coding error</li> <li>- Aggregation at the domain level</li> <li>- Data anonymization</li> </ul>
Coverage error	<ul style="list-style-type: none"> <li>- Non-trackable individuals</li> </ul>
Sampling error	<ul style="list-style-type: none"> <li>- Same error causes than for surveys</li> </ul>
Missing data error	<ul style="list-style-type: none"> <li>- Noncontact</li> <li>- Non-consent</li> <li>- Non-trackable target</li> <li>- Meter not installed</li> <li>- Uninstalling the meter</li> <li>- New non-tracked device</li> <li>- Technology limitations</li> <li>- Technology error</li> <li>- Hidden behaviour</li> <li>- Social desirability</li> <li>- Extraction error</li> </ul>
Adjustment error	<ul style="list-style-type: none"> <li>- Same error causes than for surveys</li> </ul>

Sampling and adjustment errors have no specific error causes

# Practical recommendations

- 1. Clearly define what your tracked data is measuring beforehand.**

# Practical recommendations

## 1. Clearly define what your tracked data is measuring beforehand.

**Concept:** *average hours of consumption of online political news*

**Measure:** *average time recorded of the visits to online political outlets' URLs.*

## 1. Clearly define what your tracked data is measuring beforehand.

**Concept:** *average hours of consumption of online political news*

**Measure:** *average time recorded of the **visits** to online political outlets' URLs.*

- What is considered a visit?



## 1. Clearly define what your tracked data is measuring beforehand.

**Concept:** *average hours of consumption of online political news*

**Measure:** *average time recorded of the **visits** to online political **outlets**' URLs.*

- What is considered a visit?
- Which online outlets?

## 1. Clearly define what your tracked data is measuring beforehand.

**Concept:** *average hours of consumption of online political news*

**Measure:** *average time recorded of the **visits** to online **political outlets**' URLs.*

- What is considered a visit?
- Which online outlets?
- Which URLs should be considered political?

## 1. Clearly define what your tracked data is measuring beforehand.

**Concept:** *average hours of consumption of online political news*

**Measure:** *average time recorded of the visits to online political outlets' URLs.*

- What is considered a visit?
- Which online outlets?
- Which URLs should be considered political?
- What time frame to use to compute an average?

# Practical recommendations

**2. Consider the impact of the chosen technologies on data quality.**

## 2. Consider the impact of the chosen technologies on data quality.

### Apps

**Where?**

Device

**Devices**

Not iOS

**Continuous?**

Yes

**Types of data**URLs, Time,  
Device, Search  
terms, Incognito,

### Plug-in A

**Where?**

Browser

**Devices**

Only PC &amp; MAC

**Continuous?**

Yes

**Types of data**URLs, Time,  
Device, Search  
terms, Incognito,  
HTML

### Plug-in B

**Where?**

Browser

**Devices**

Only PC &amp; MAC

**Continuous?**

No

**Types of data**URLs, Time,  
Device

### Proxy

**Where?**

Network

**Devices**

All

**Continuous?**

Yes

**Types of data**URLs, Time,  
Device

## 2. Consider the impact of the chosen technologies on data quality.

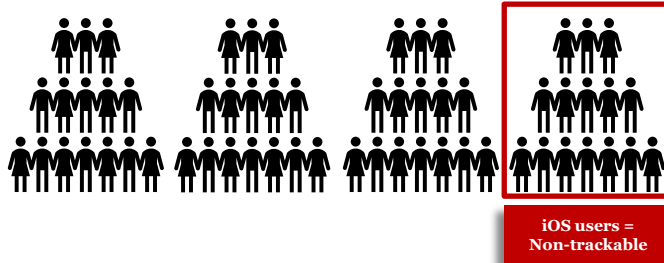
**Apps**

**Where?**  
Device

**Devices**  
Not iOS

**Continuous?**  
Yes

**Types of data**  
URLs, Time,  
Device, Search  
terms, Incognito



## 2. Consider the impact of the chosen technologies on data quality.

### Apps

Where?

Device

Devices

Not iOS

Continuous?

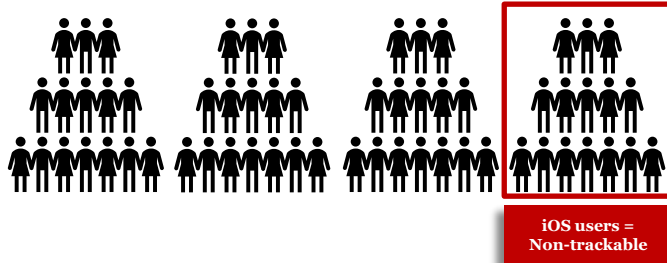
Yes

Types of data

URLs, Time,

Device, Search

terms, Incognito



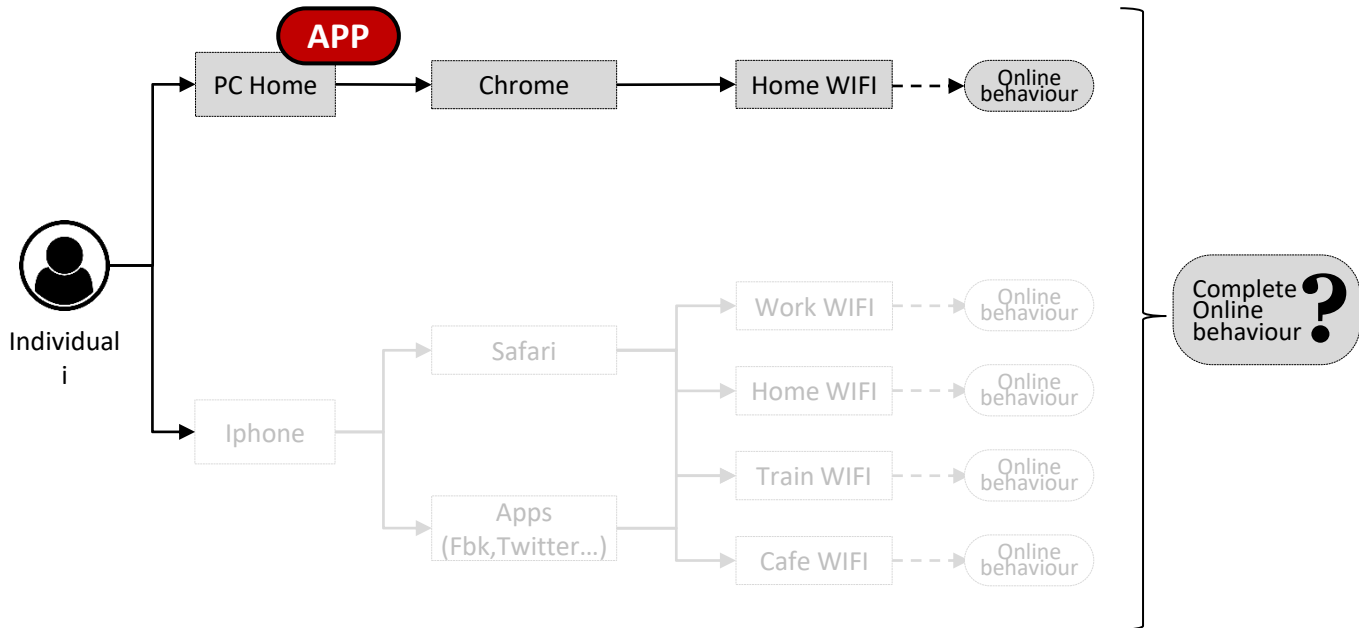
The screenshot shows the top of the Guardian website. At the top, there is a dark blue navigation bar with the text "Support the Guardian" and "Available for everyone, funded by readers". Below this are two yellow buttons labeled "Contribute ->" and "Subscribe ->". To the right of the navigation bar are links for "Search jobs", "Sign in", "Search", and "UK edition". The main navigation bar is dark blue with white text for "News", "Opinion", "Sport", "Culture", "Lifestyle", and "More". Below the navigation bar, there is a breadcrumb trail: "Environment > Climate change > Wildlife > Energy > Pollution". The main content area shows an "Opinion" article titled "Ignore the rhetoric: the UK government still fails to grasp the climate crisis" by "Chris Venables".

## Practical recommendations

**3. Explore strategies to increase the willingness of individuals to install the meter in all targets.**



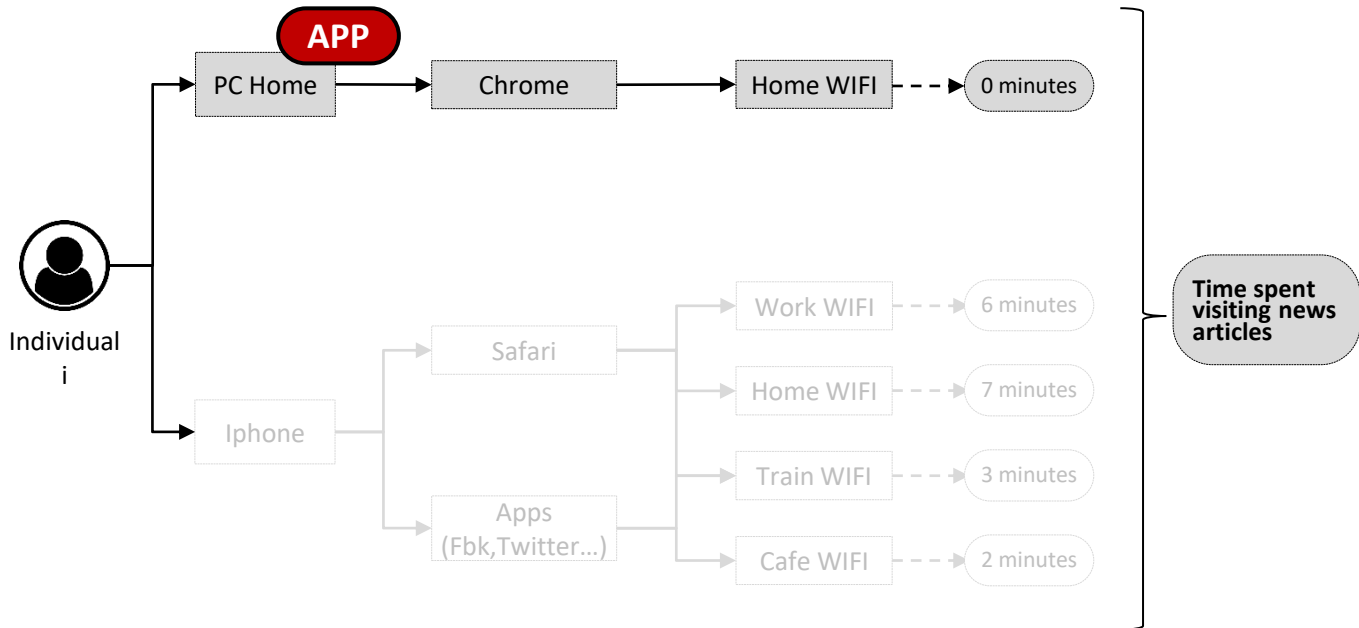
## 3. Explore strategies to increase the willingness of individuals to install the meter in all targets.



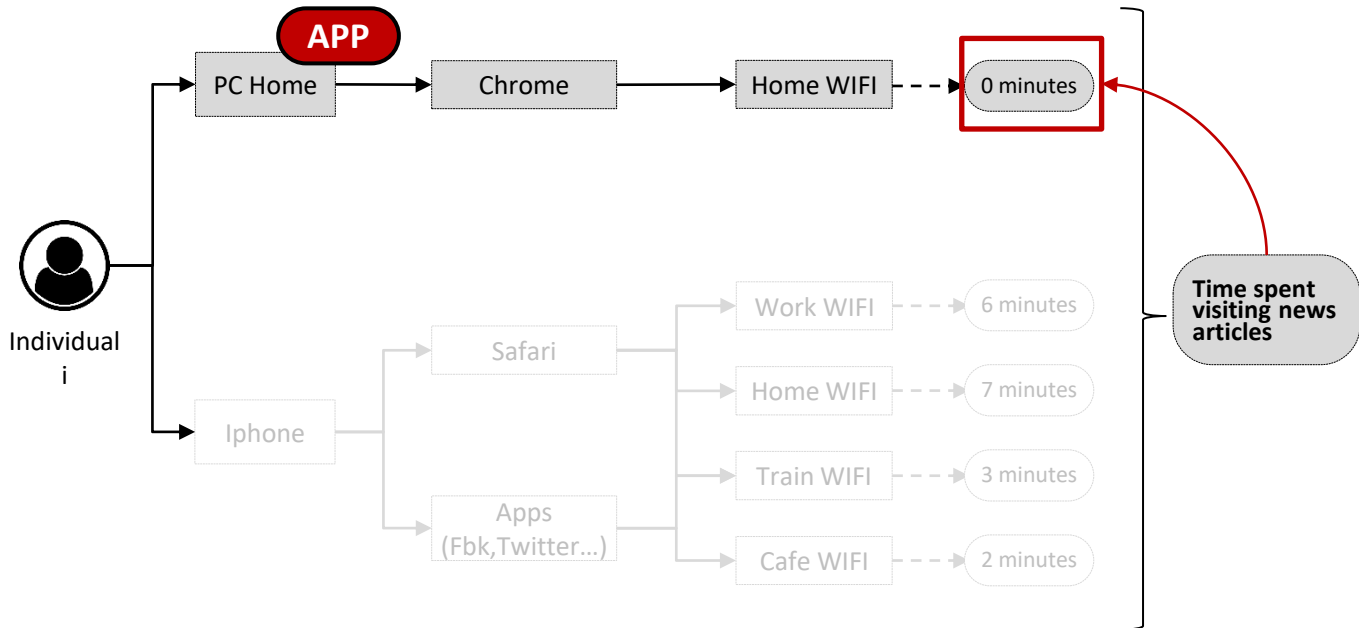
### **3. Explore strategies to increase the willingness of individuals to install the meter in all targets.**

- Tracking technologies present different installations processes.
- Multiple tracking technologies might need to be installed for the same participant.
- Targets (devices / browsers / networks used) are unknown.

## 4. Define strategies to maximise the information available to identify missing data.

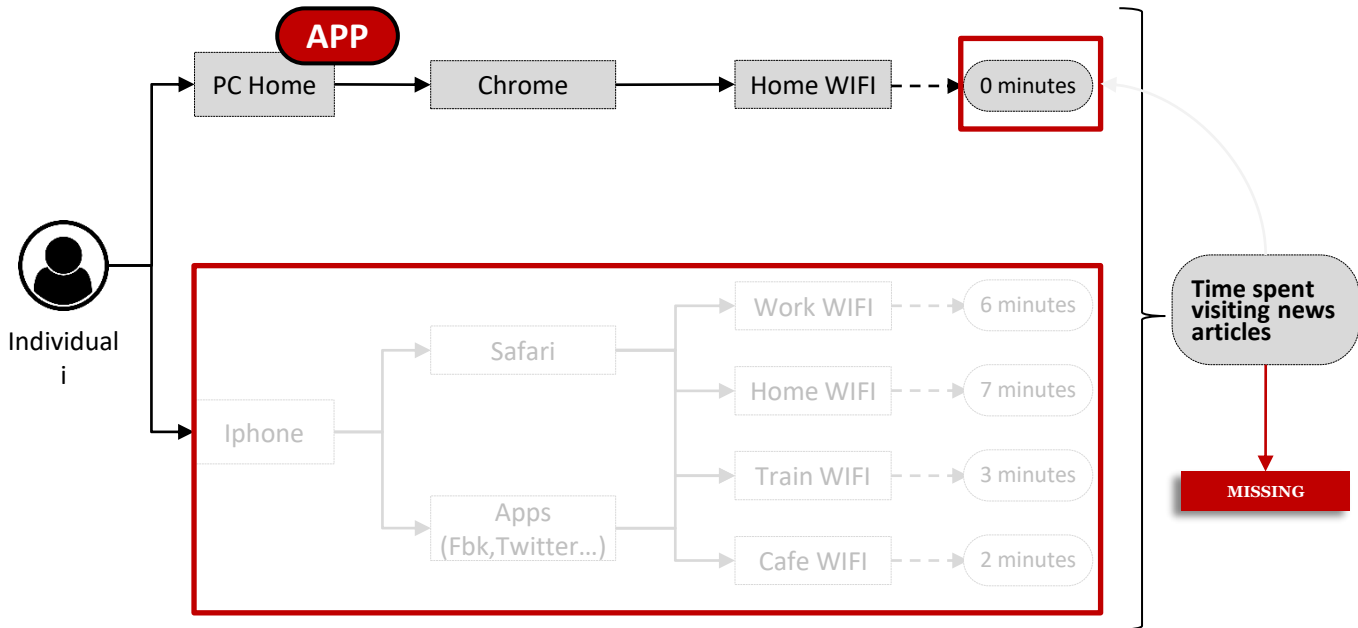


## 4. Define strategies to maximise the information available to identify missing data.



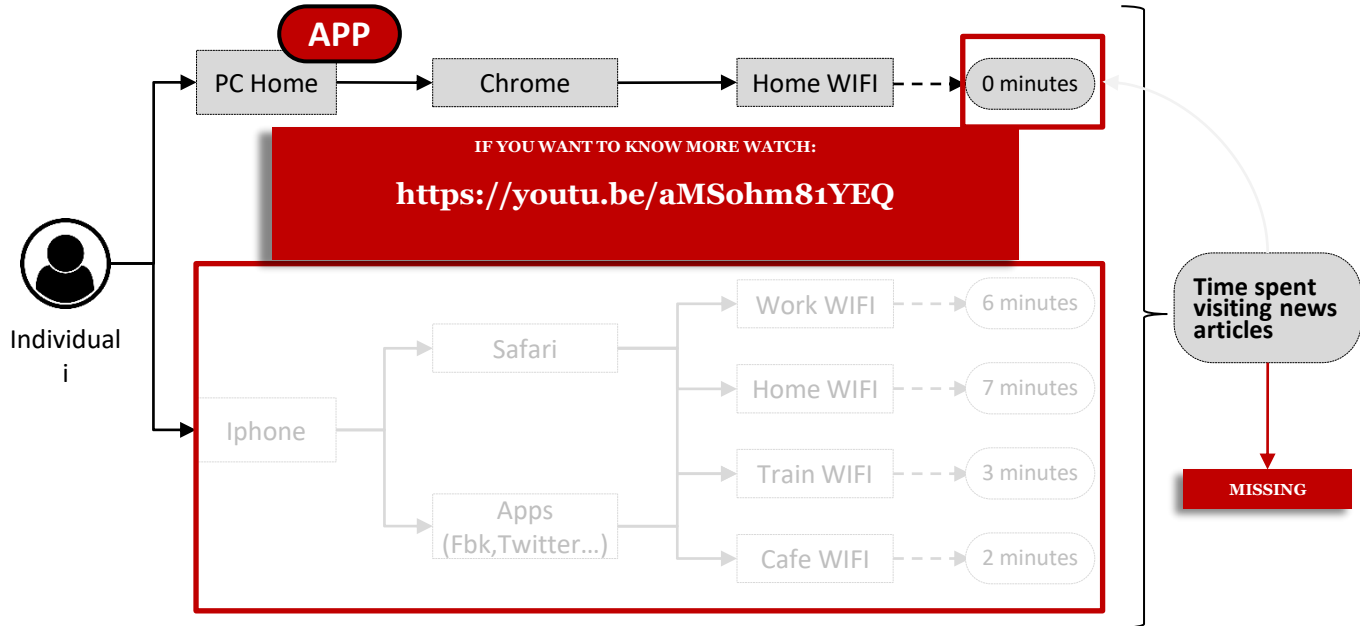


## 4. Define strategies to maximise the information available to identify missing data.



# Practical recommendations

## 4. Define strategies to maximise the information available to identify missing data.



## Limits

1. One specific definition of data quality.
2. Lack of previous empirical research.
3. Tracking technologies are constantly evolving.
4. Metered data errors are considered independently.

## Take-home messages

1. Using metered data is complex and many decisions must be taken.
2. Reporting these decisions and conducting robustness checks is necessary.
3. More empirical research is needed.
4. This framework can help on all these aspects.



# Thank you

## *Questions?*

ORIOI J. BOSCH | THE LONDON SCHOOL OF ECONOMICS / RECSM-UPF



[o.bosch-jover@lse.ac.uk](mailto:o.bosch-jover@lse.ac.uk)



[orioljbosch](https://twitter.com/orioljbosch)



<https://orioljbosch.com/>

Bosch, O.J., and M. Revilla (2021). **“When survey science met online tracking: presenting an error framework for metered data.”** RECSM Working Papers Series, 62

**LSE**

THE LONDON SCHOOL  
OF ECONOMICS AND  
POLITICAL SCIENCE ■

**upf.**

Universitat  
Pompeu Fabra  
Barcelona

**RECSM**

Research and Expertise Centre  
for Survey Methodology

web  
data  
opp