

Recent Developments in Federal Statistics – Surveys and Beyond

Frauke Kreuter

JPSM – Uni Mannheim – IAB

Melbourne June 26 '18



AAPOR Report on Big Data

AAPOR Big Data Task Force
February 12, 2015

Prepared for AAPOR Council by the Task Force, with Task Force members including:

Lilli Japec, Co-Chair, Statistics Sweden
Franke Kreuter, Co-Chair, JPSM at the U. of Maryland, U. of Mannheim & IAB
Marcus Berg, Stockholm University
Paul Biemer, RTI International
Paul Decker, Mathematica Policy Research
Cliff Lampe, School of Information at the University of Michigan
Julia Lane, American Institutes for Research
Cathy O'Neil, Johnson Research Labs
Abe Usher, HumanGeo Group

Acknowledgement: We are grateful for comments, feedback and editorial help from Eran Ben-Porath, Jason McMillan, and the AAPOR council members.

The National Academies of
SCIENCES · ENGINEERING · MEDICINE

REPORT

INNOVATIONS IN FEDERAL STATISTICS

Combining Data Sources While
Protecting Privacy

The National Academies of
SCIENCES · ENGINEERING · MEDICINE

CONSENSUS STUDY REPORT

FEDERAL STATISTICS, MULTIPLE DATA SOURCES, AND PRIVACY PROTECTION

Next Steps



SPONSORED BY THE

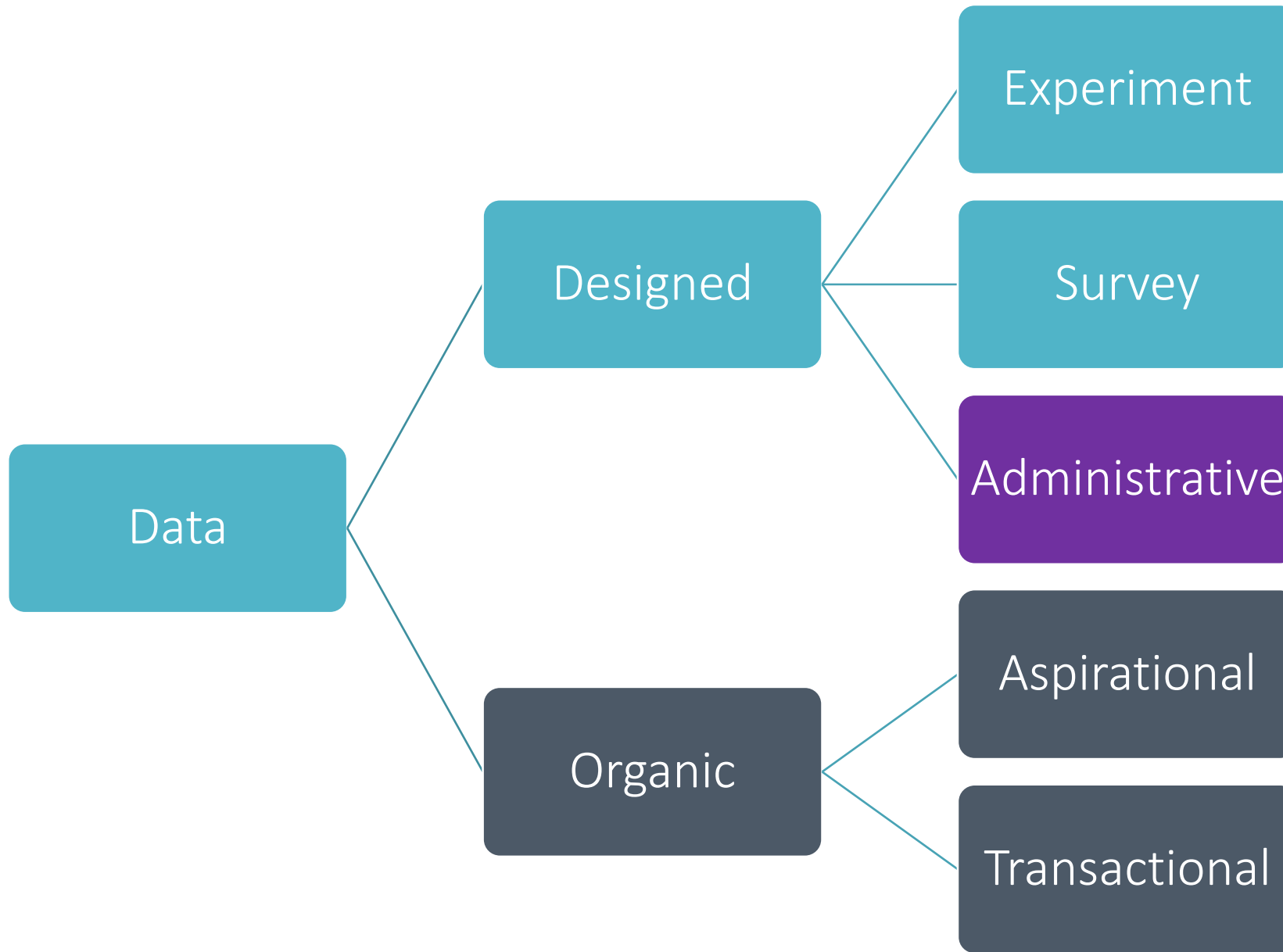


Federal Ministry
of Education
and Research



The project on which this report is based was funded by the Federal Ministry of Education and Research under the number [16OH22064]. Responsibility for the contents of this publication lies with the author.





Source: Roberto Rigobon, [Discussion on Applications and Issues with Using Commercial Data in Research](#), BEA Expert Meeting on Exploiting Commercial Data for Official Economic Statistics November 19, 2015

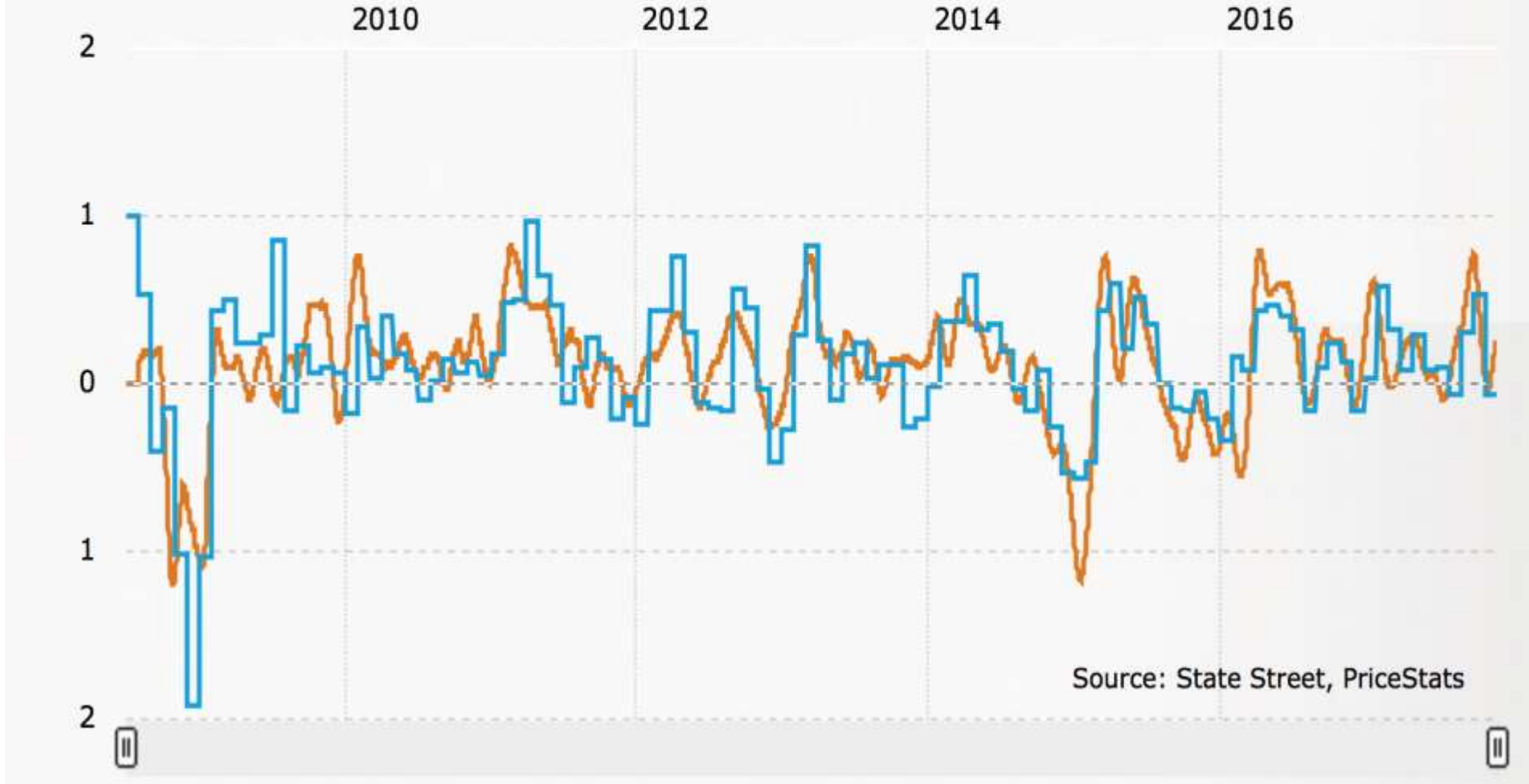
Economic Indicators

Examples for online data collection (and analysis)

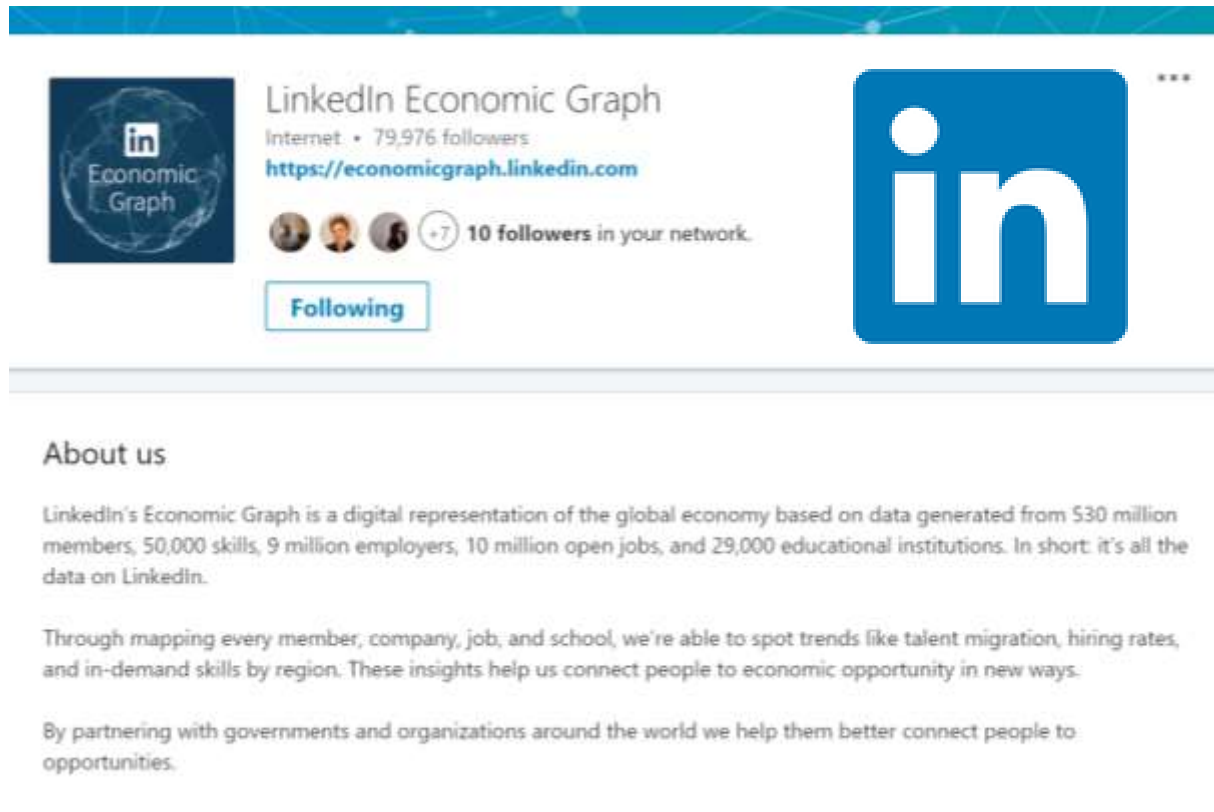
US Aggregate Inflation Series

(Monthly Rate, 2008 - Present)

- Official CPI
- PriceStats Index



Employment Websites + Aggregators



The image shows a LinkedIn profile for 'LinkedIn Economic Graph'. The profile picture is a globe with the LinkedIn logo and the text 'Economic Graph'. The name is 'LinkedIn Economic Graph', with 'Internet' as the industry and '79,976 followers'. The URL is 'https://economicgraph.linkedin.com'. There are 10 followers in the user's network, with a '+7' icon indicating more. A 'Following' button is visible. Below the profile is an 'About us' section with the following text:

About us

LinkedIn's Economic Graph is a digital representation of the global economy based on data generated from 530 million members, 50,000 skills, 9 million employers, 10 million open jobs, and 29,000 educational institutions. In short: it's all the data on LinkedIn.

Through mapping every member, company, job, and school, we're able to spot trends like talent migration, hiring rates, and in-demand skills by region. These insights help us connect people to economic opportunity in new ways.

By partnering with governments and organizations around the world we help them better connect people to opportunities.

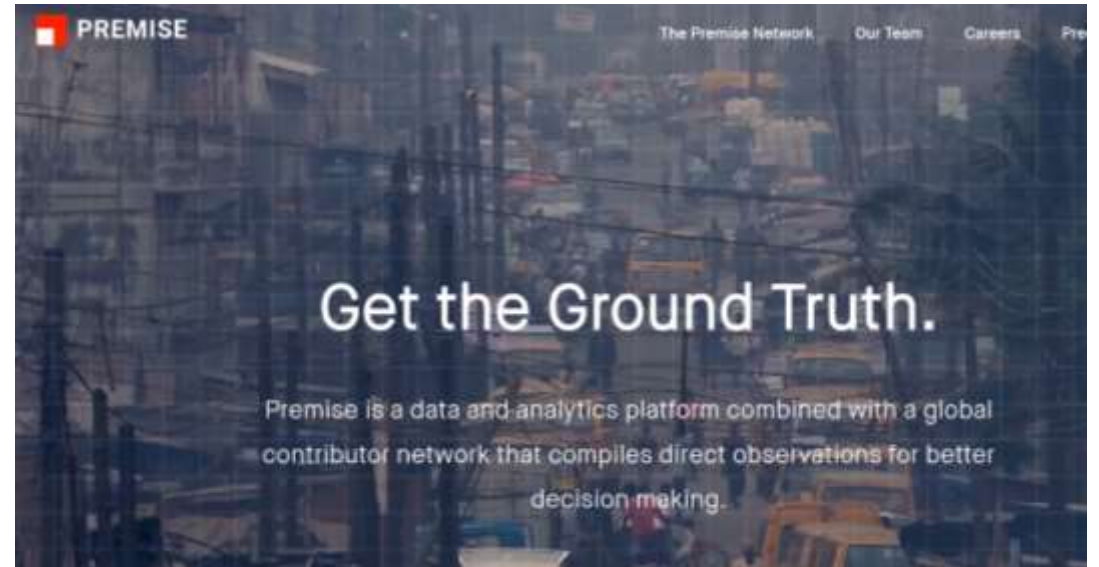


The image shows the homepage of burningglass TECHNOLOGIES. The header includes the logo and navigation links for HOME, ABOUT, INDUSTRIES, and PRODUCTS. The main headline is 'BUILD A BETTER WORKFORCE' with a sub-headline 'HUMAN CAPITAL MANAGEMENT'. A 'FIND OUT HOW' button is present. At the bottom, there are three circular icons representing 'HUMAN CAPITAL MANAGEMENT', 'HIGHER EDUCATION', and 'RECRUITING AND STAFFING'.



The image shows the logo for WANTEDANALYTICS, a CEB Company. The logo features an orange square with a white line graph icon, followed by the text 'WANTEDANALYTICS' in orange and 'a CEB Company' in black.

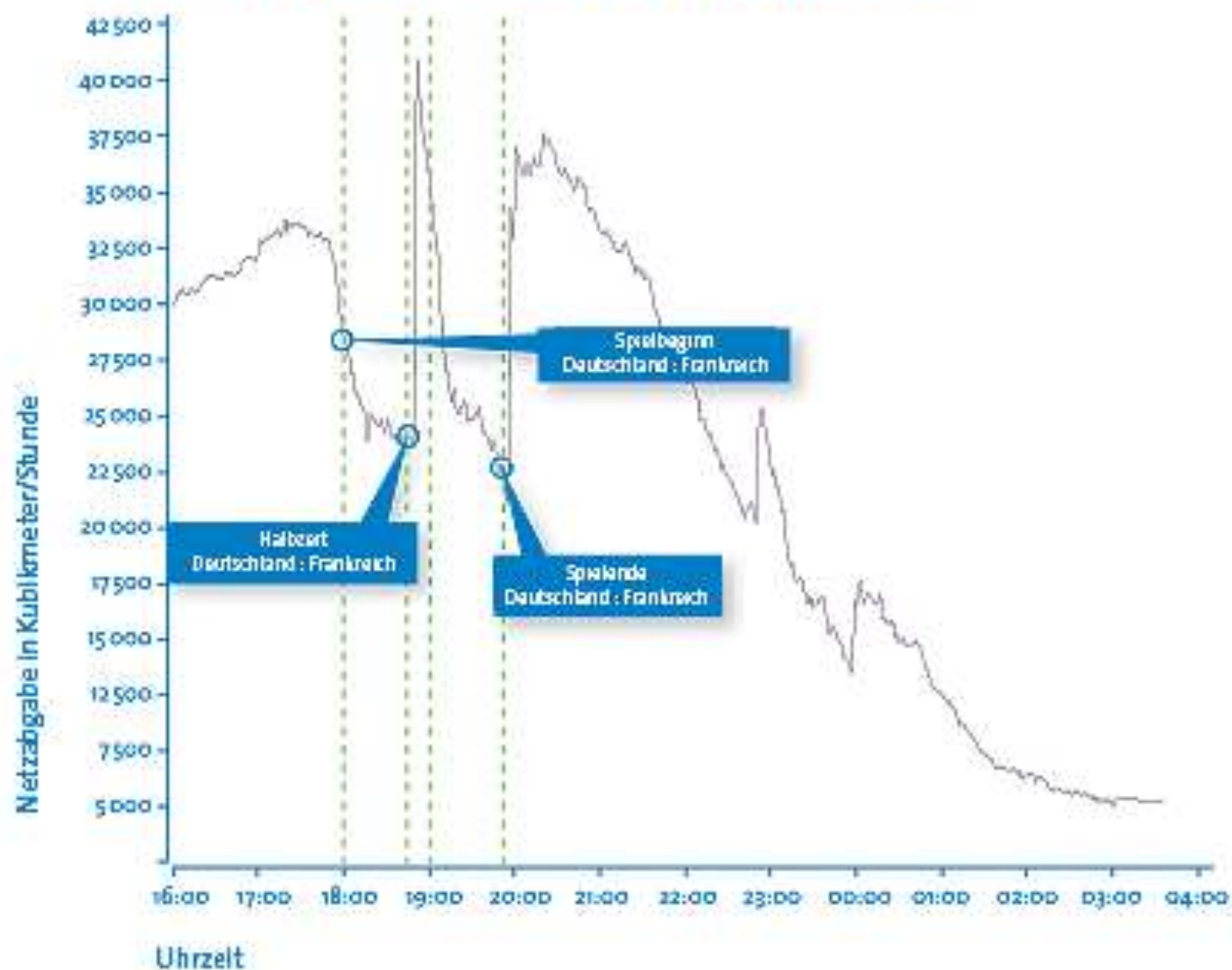
Platforms - Crowdsourcing



Observations



Netzabgabe für den 4. Juli 2014

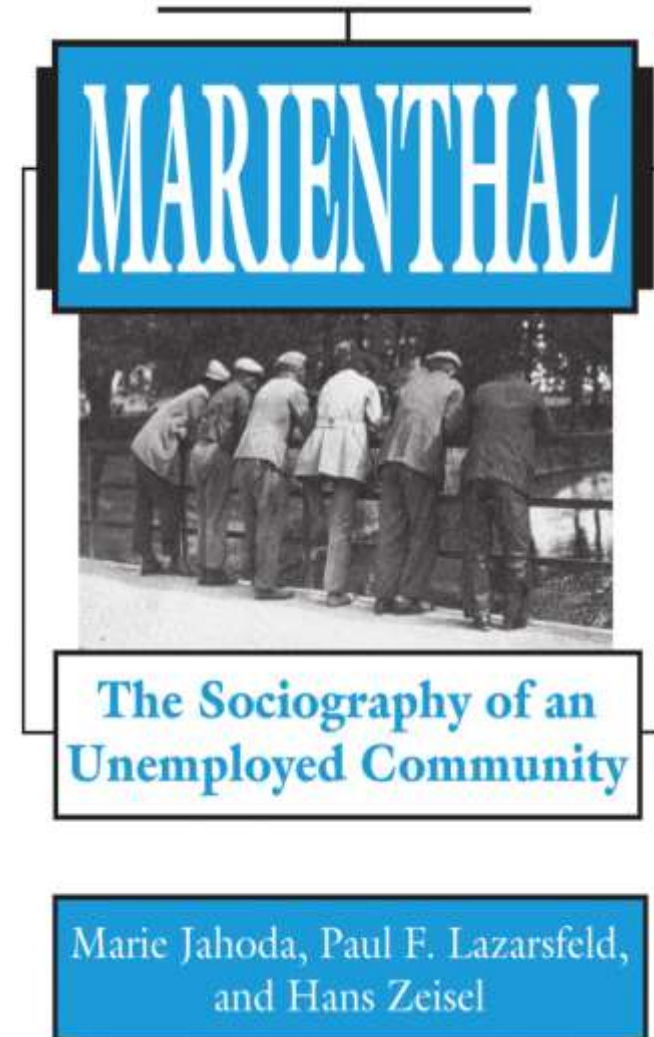
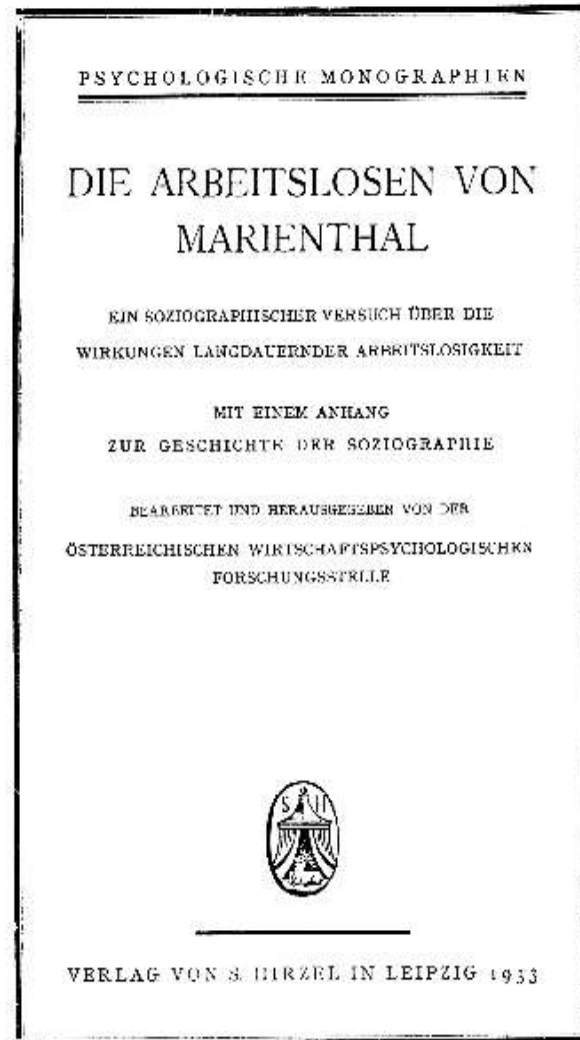


1. IAB – Research Examples

1st Refugees



2nd Unemployment





9. Januar 2018 | Projekte

IAB-SMART-Studie: Mit dem Smartphone den Arbeitsmarkt erforschen



Autoren

Sebastian Bähr

Georg-Christoph Haas

Florian Keusch

Frauke Kreuter

Mark Trappmann

PASS – Panel (10 years) + Administrative Data

Sample of households with at least one welfare benefit recipient (at reference date)

Refreshed annually
Surveyed annually

Random household sample of resident population

Refreshed annually

Surveyed annually



Meldung zur Sozialversicherung

Personalauswahl

Versicherungsnummer Personalnummer (freiwillige Angabe)

Name Vorsatz Zusatz Titel

Vorname

Straße und Hausnummer (Anschrift nur bei Anmeldung und Anschriftenänderung)

(Land) Postleitzahl Wohnort

Grund der Abgabe Entgelt in Gleitzone Namensänderung

Beschäftigungszeit

von bis Betriebsnummer des Arbeitgebers Personengruppe

Mehrfachbeschäftigung Ost West

Beitragsgruppen KV RV ALV PV Angaben zur Tätigkeit Aktuelle Staatsangehörigkeit

Beitragspflichtiges Bruttoarbeitsentgelt (in DM ohne Pfennige / Euro ohne Cent) DM Euro Statuskennzeichen

Wenn keine Versicherungsnummer angegeben werden kann:

Geburtsname Vorsatz Zusatz Geburtsort

Geburtsdatum Geschlecht männlich weiblich

General Data Processing Regulation (GDPR)

Opt-In

Consent in GDPR

Consent needs to be freely given.

Consent needs to be specific, per purpose.

Consent needs to be informed.

Consent needs to be an unambiguous indication.

Consent is an act: it needs to be given by a statement or by a clear act.

Consent needs to be distinguishable from other matters.

Consent request needs to be in clear, plain language; intelligible and easily accessible.

Informed Consent – Privacy - Technology

Instruction booklet with IAB-Smart app screen shots – separate and active opt-in required for all use cases

Bitte klicken Sie auf „AKZEPTIEREN“, um mit der Installation fortzufahren.



5 6

Damit alle Funktionen der App genutzt werden können, müssen Sie jeweils zustimmen:



Sie haben aber auch die Möglichkeit die App zu installieren und Punkte für die Beantwortung von Fragen zu erhalten, ohne den Nutzungsdatenzugriff zu erlauben.

1

Um die Erfassung von Nutzungsdaten zu erlauben, verschieben Sie den Schieber bitte nach rechts.



13 14

Consent to Linkage by Framing and Mode in %

Phone	Front	Back	Total n
Gain	90.8	78.7	598
Loss	90.5	81.2	610
Total n	613	595	1208

Web	Front	Back	Total
Gain	82.6	62.4	520
Loss	86.3	75.4	489
Total	511	498	1009

Lack of understanding

Phone	Consenters %correct	Non- consenters %correct
Answers sent to IAB	88.3	57.8
Merged with IAB	93.3	36.7
Name/Adress saved	68.3	38.8
Result lead to you	63.4	--
IAB only access	85.6	--
Public access to identifiable data	87.5	--

2. Role of methodologists

Data Output/Access

Data Analysis

Data Curation/Storage

Data Generating Process

Research Question



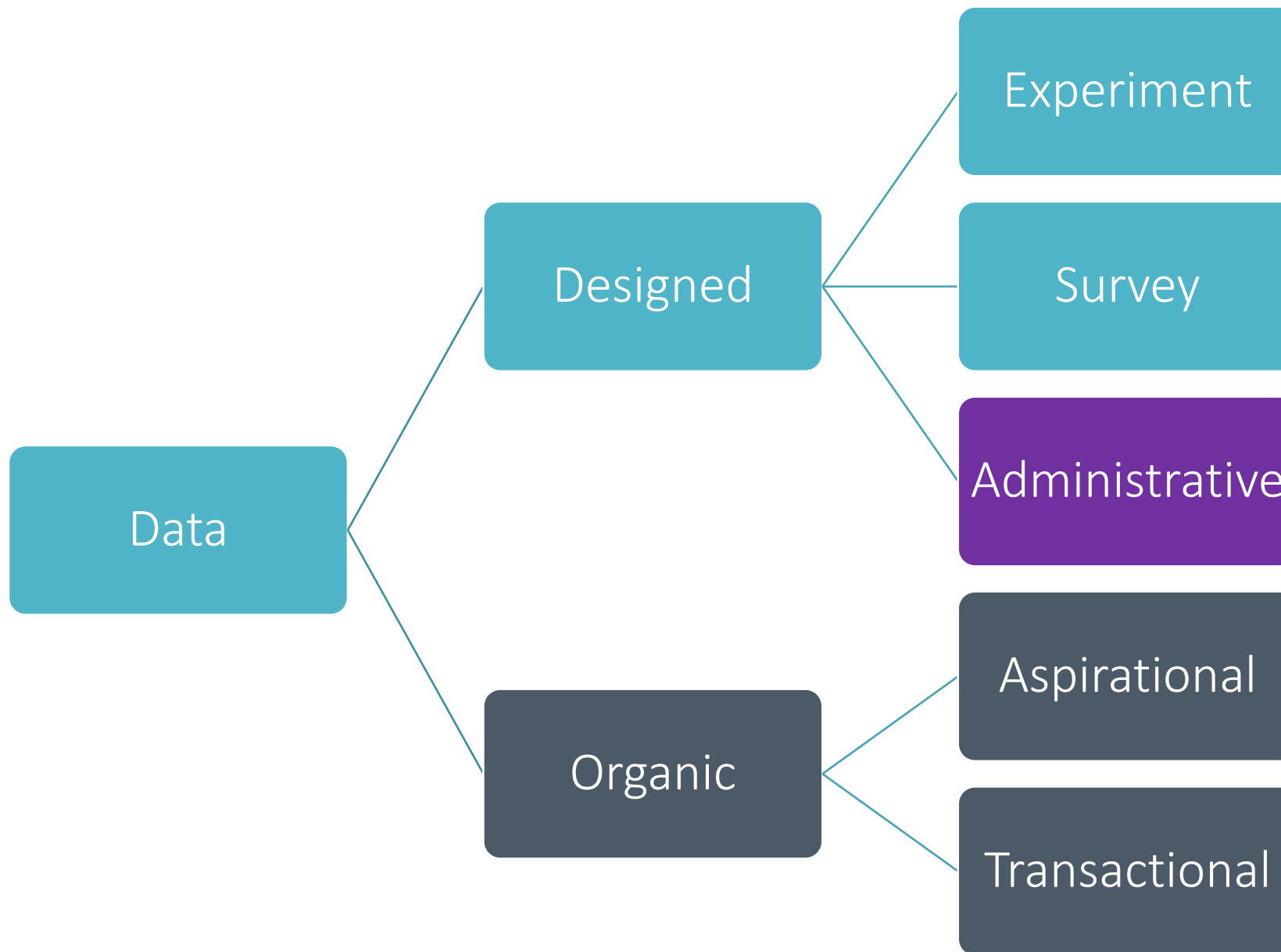
Communicate results and distribute data

Variety of analysis methods suited for different data types

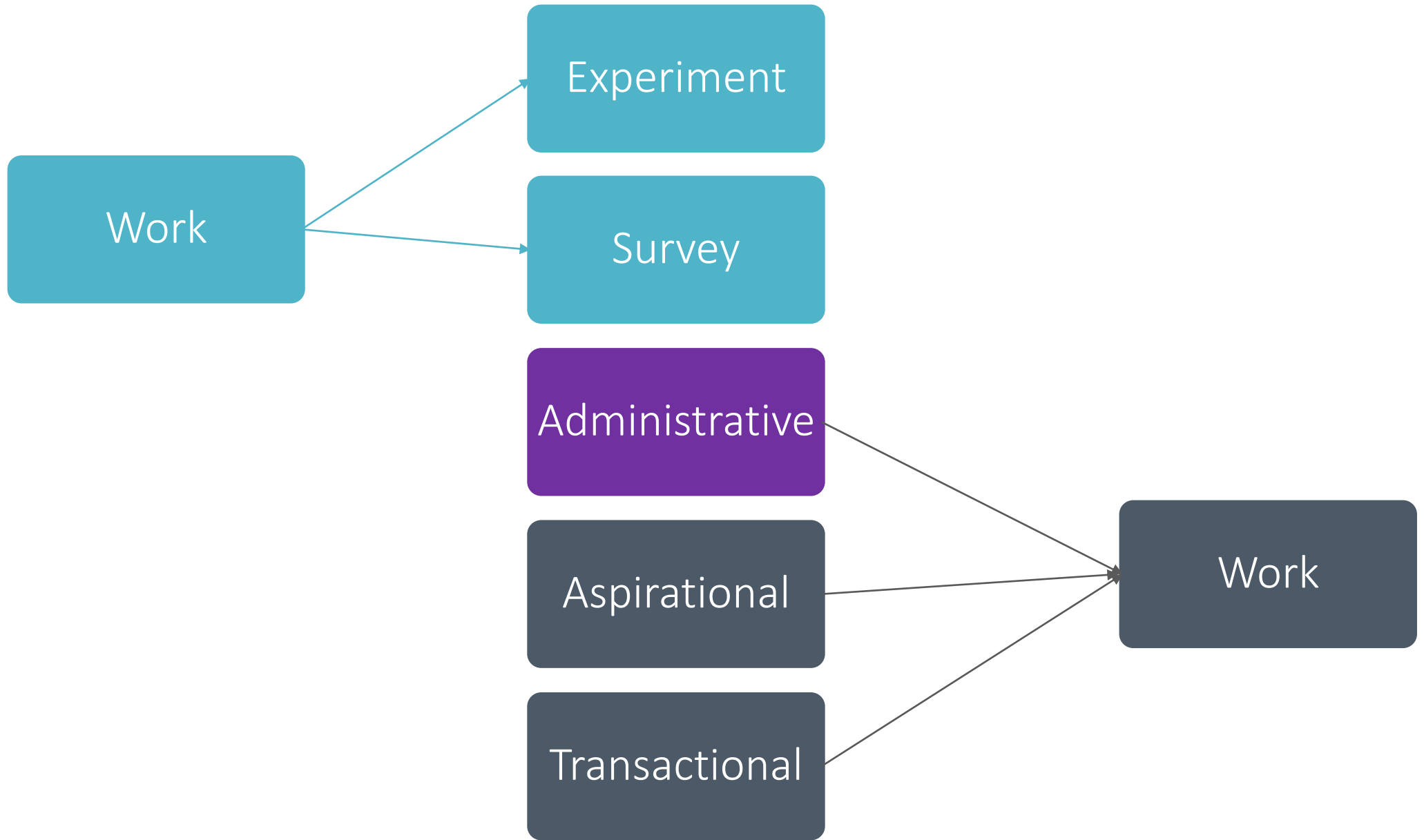
Curate and manage data

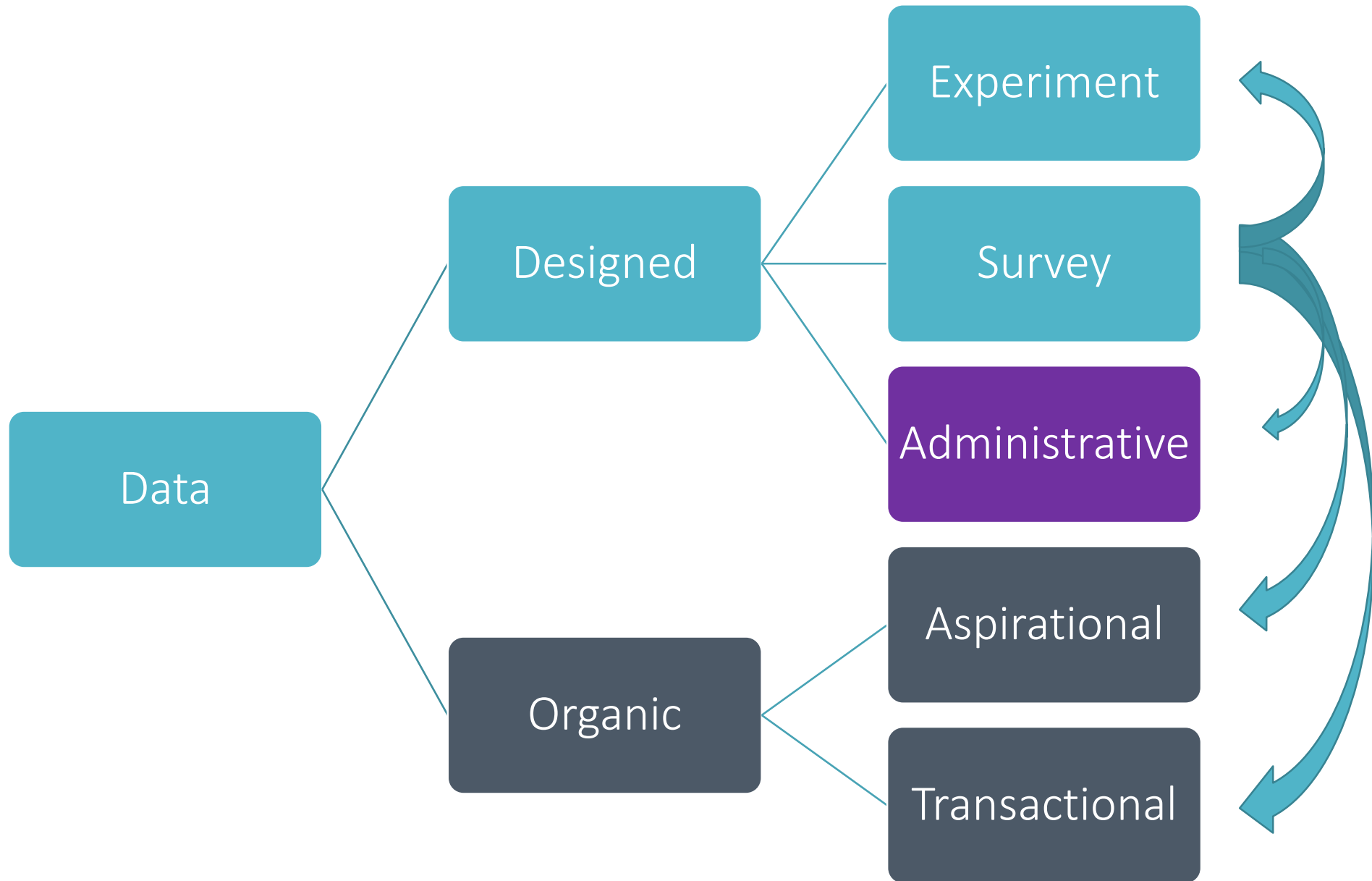
Collect data and understand how data are generated through administrative and other processes.

Formulate research goal and know which data are best suited to achieve this goal.



Source: Roberto Rigobon, [Discussion on Applications and Issues with Using Commercial Data in Research](#), BEA Expert Meeting on Exploiting Commercial Data for Official Economic Statistics November 19, 2015







DOMAIN EXPERT

User, analyst, or leaders with deep subject matter expertise related to the data, its appropriate use, and its limitations

SYS ADMIN

Team member responsible for defining and maintaining a computation infrastructure that enables large scale computation

METHODOLOGIST

Team member with experience applying formal research methods, including survey methodology and statistics

COMPUTER SCIENTIST

Technically skilled team member with education in computer programming and data processing technology

3. Some new vocabulary (skills)

Data Output/Access

Data Analysis

Data Curation/Storage

Data Generating Process

Research Questions



Example: map visualization / privacy / **GDPR**

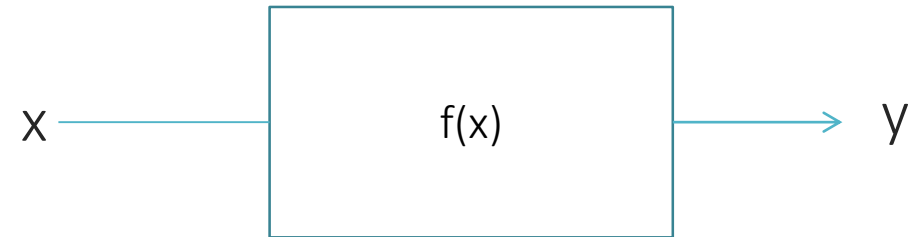
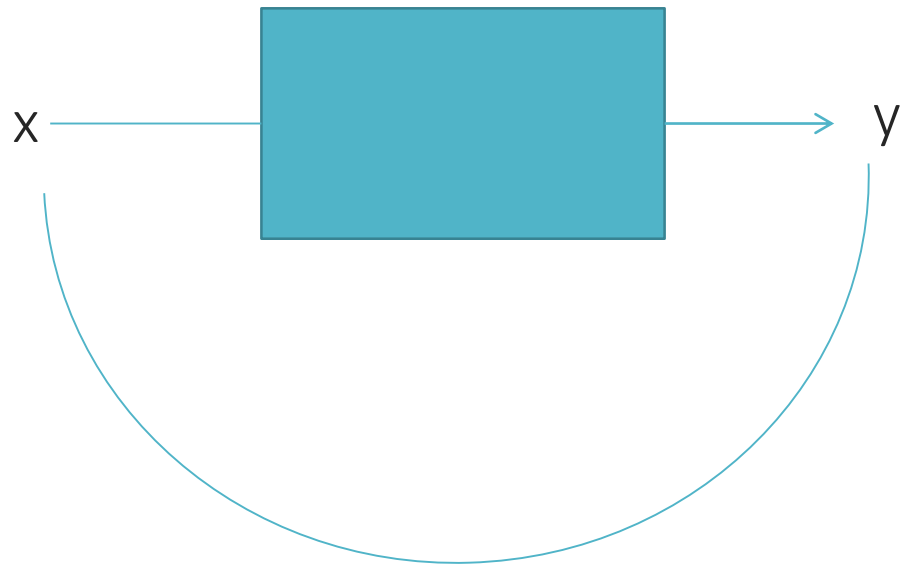
Example: **Hadoop MapReduce**;
High Frequency Data; **Machine Learning**

Example: Record Linkage; **Database**
Hadoop Distributed File System

Examples: geolocated social media + survey
+ administrative data

Examples: Behavior of interest
(political participation/job searches)

Machine Learning



Database Management

Text files and scripting language

- Your data is small
- Your analysis is simple
- You do not expect to repeat analyses over time

Statistical packages

- Your data is modest in size
- Your analysis maps well to your chosen statistical package

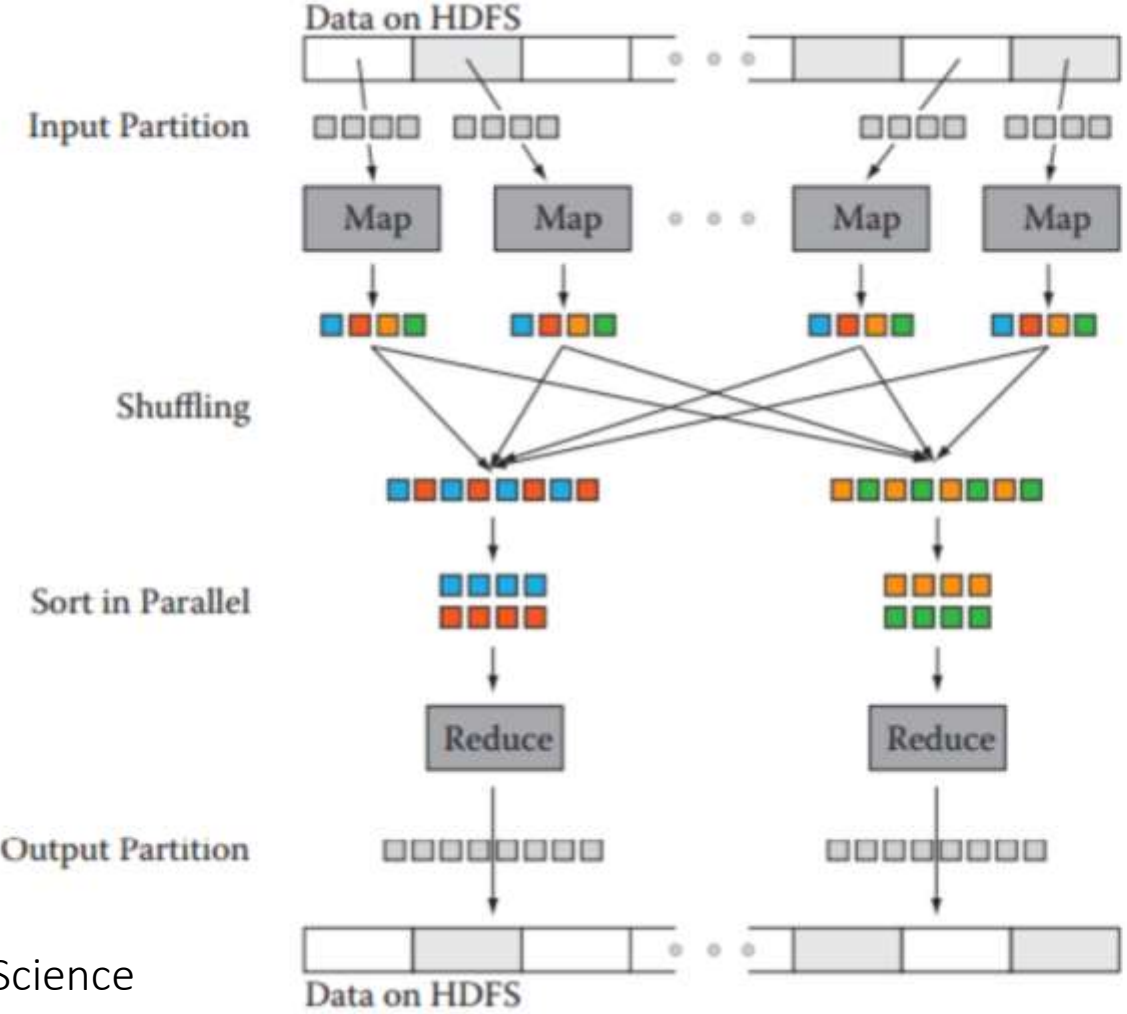
Relational database

- Your data is structured
- You will be analyzing data repeatedly over time

NoSQL database

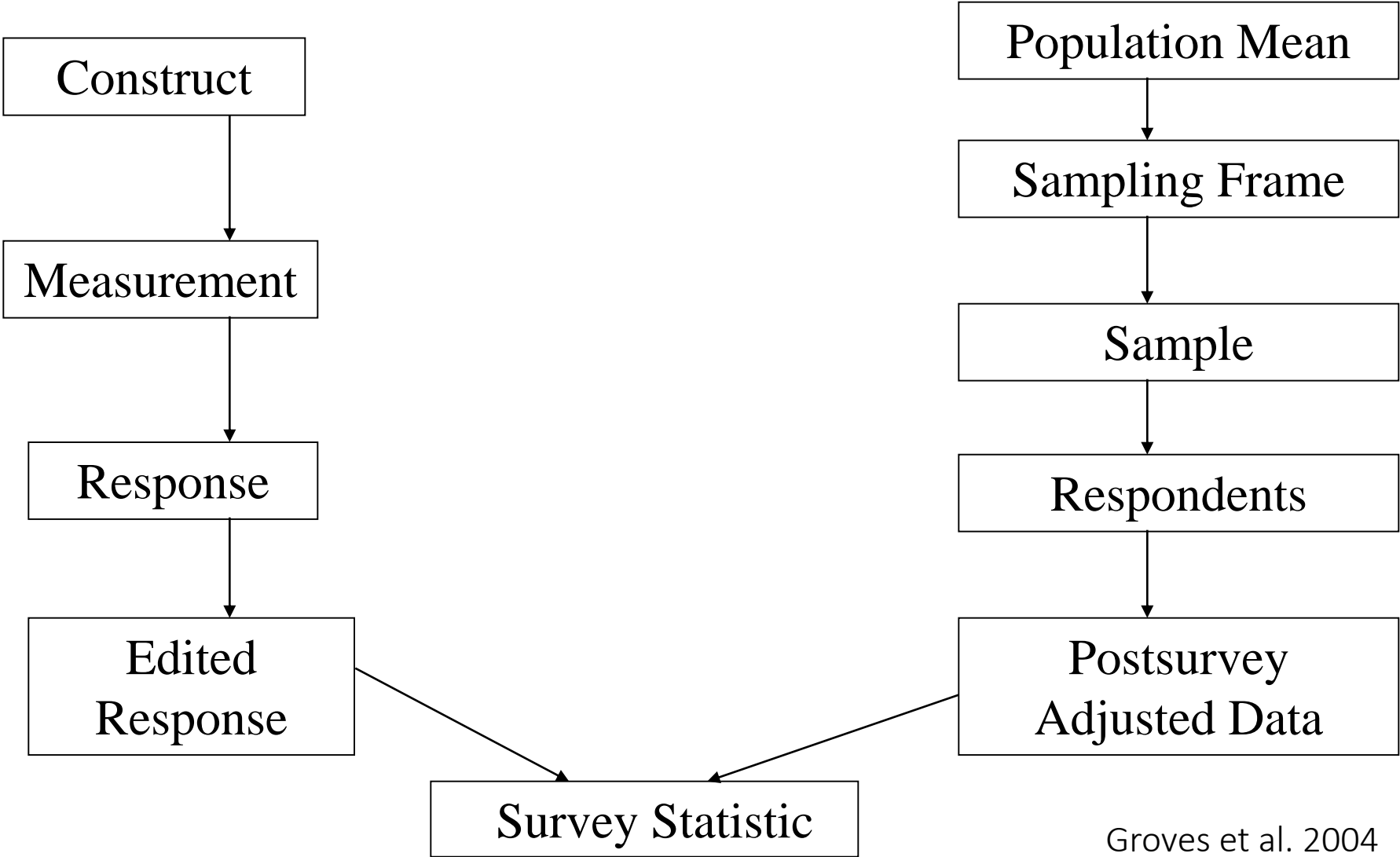
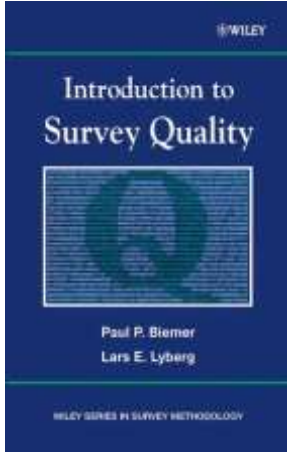
- Your data is unstructured
- Your data is extremely large

BD Programming – MapReduce

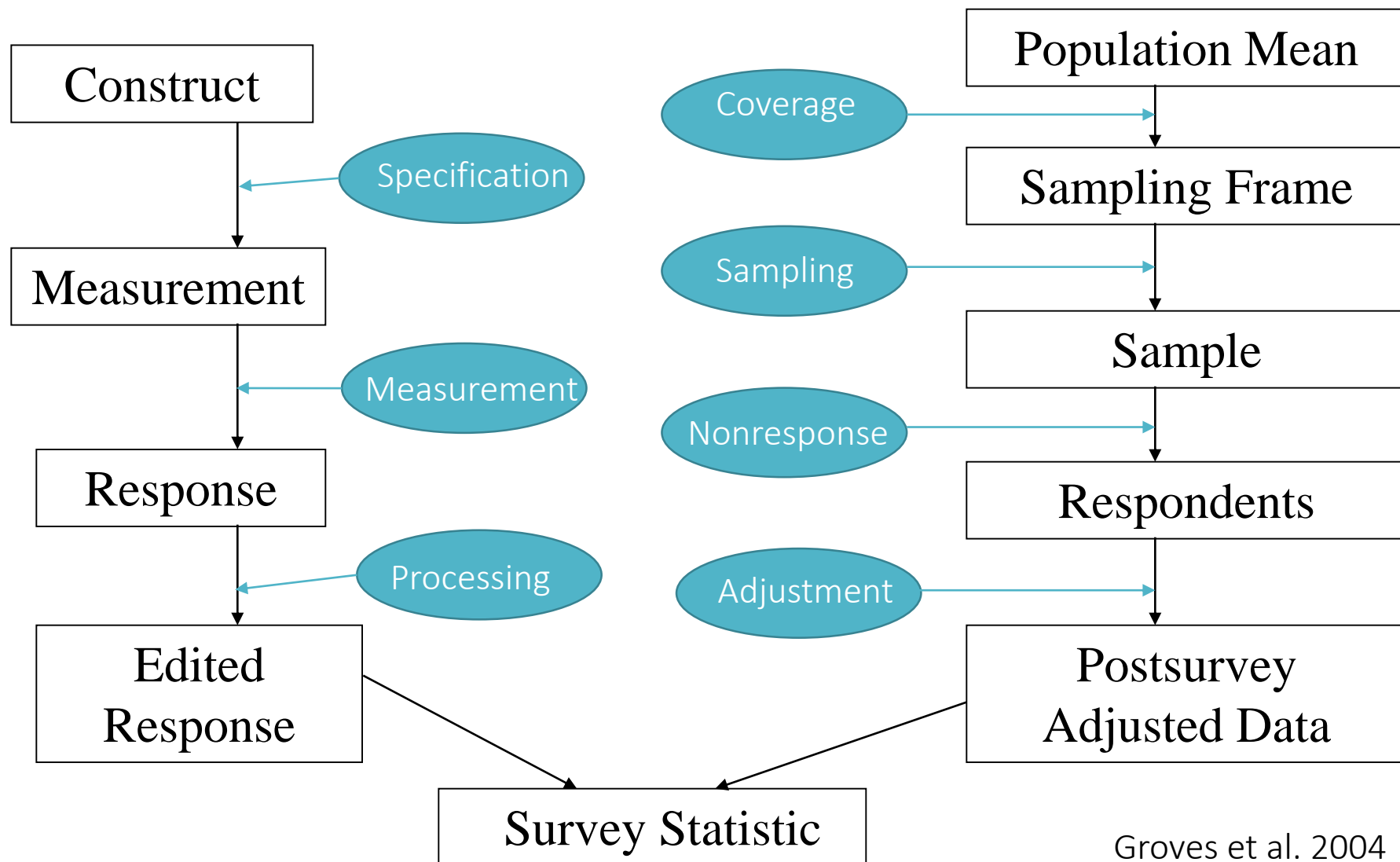
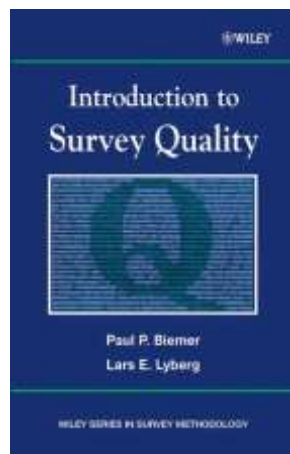
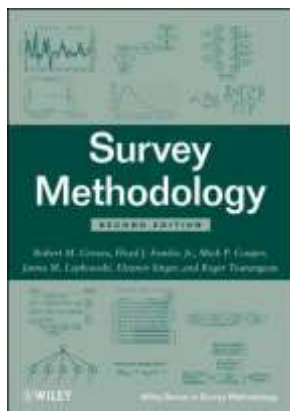


4. What we bring to the table

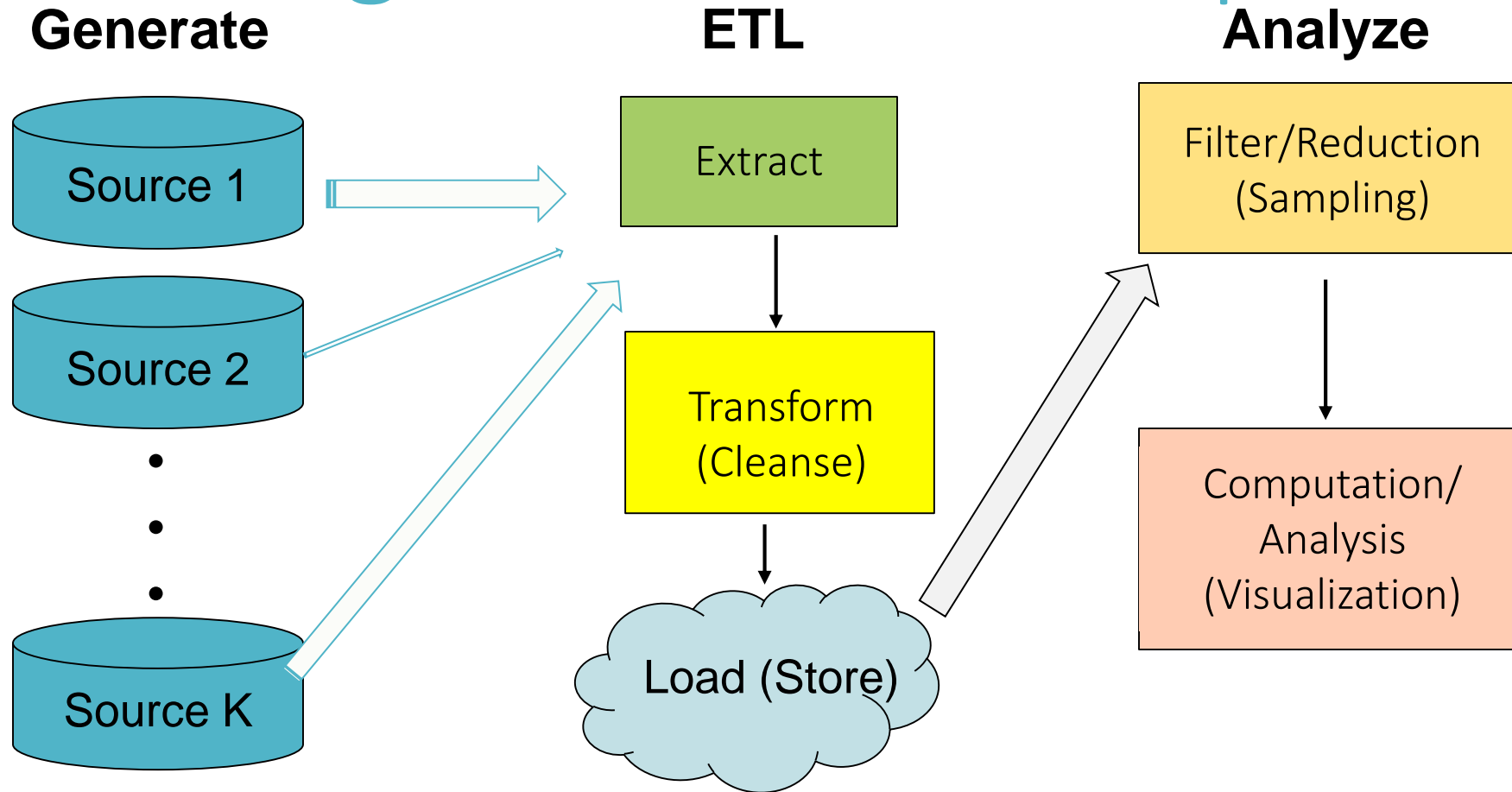
Data Generating Process



Data Generating Process

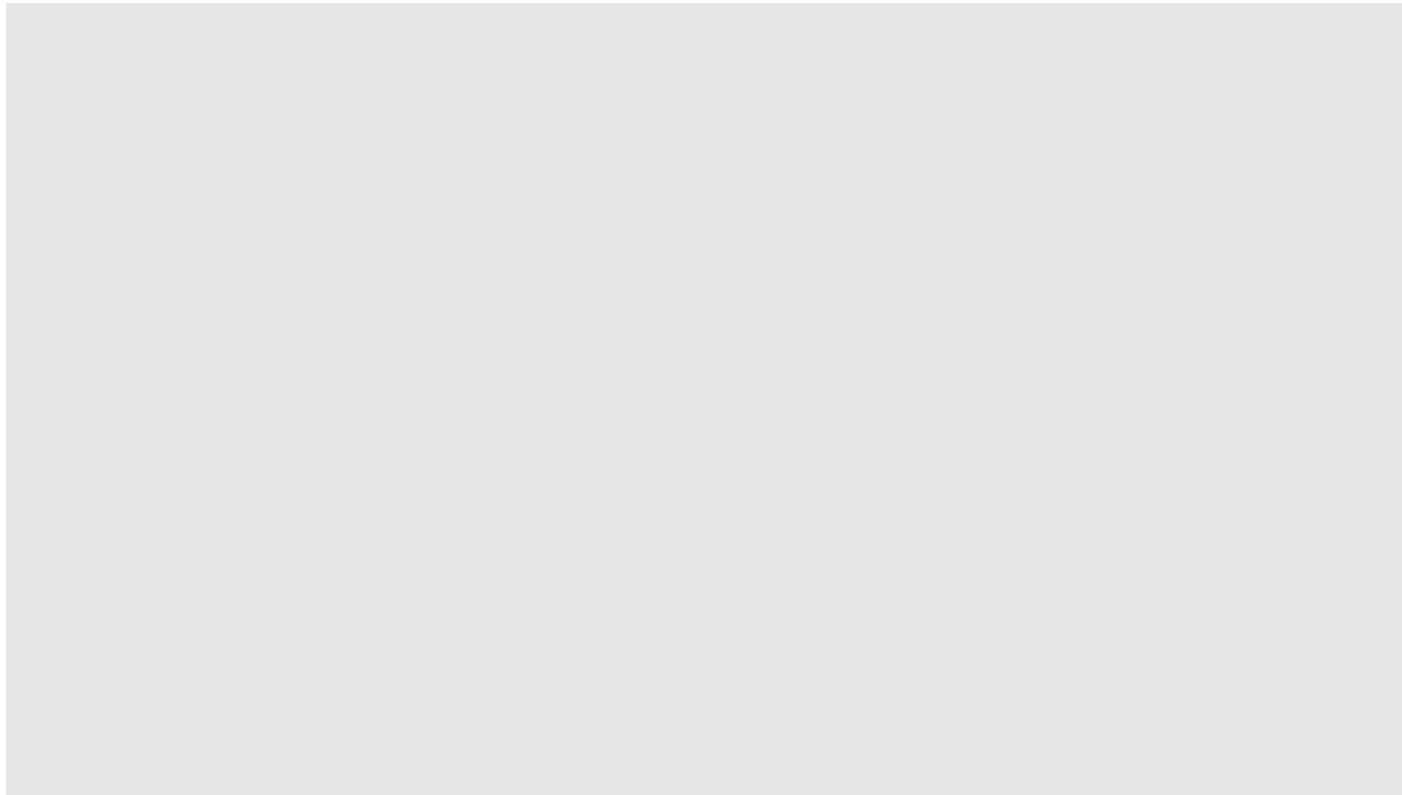


Big Data Process Map

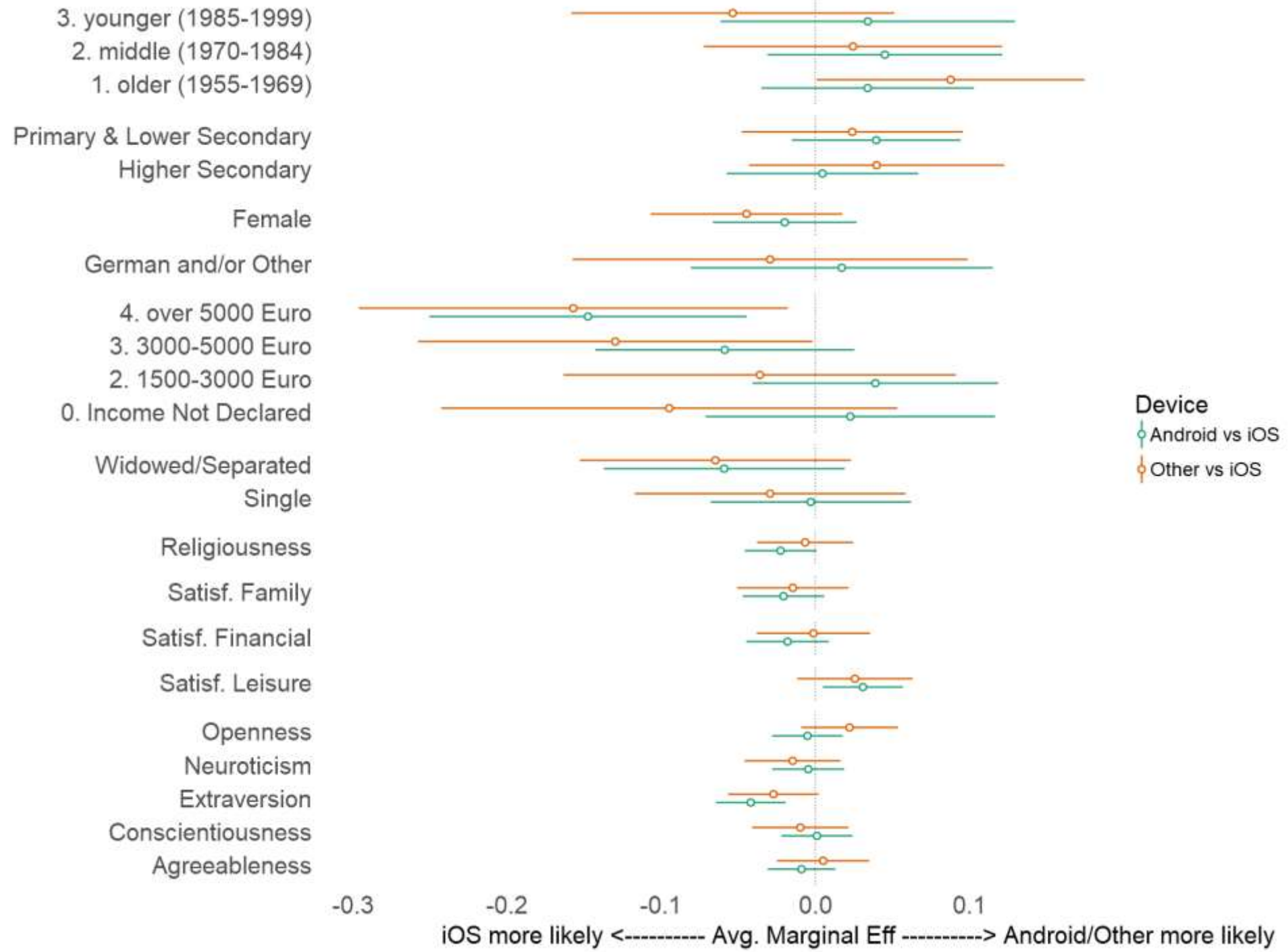


Source: Paul Biemer in Japac, Kreuter et al. 2015 – AAPOR Task Force Report

Boston Street Bumps



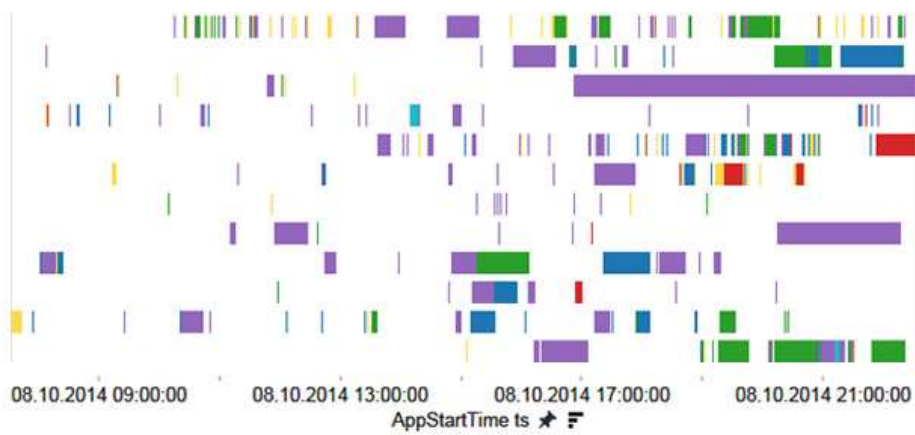
5. Back to the IAB - Example



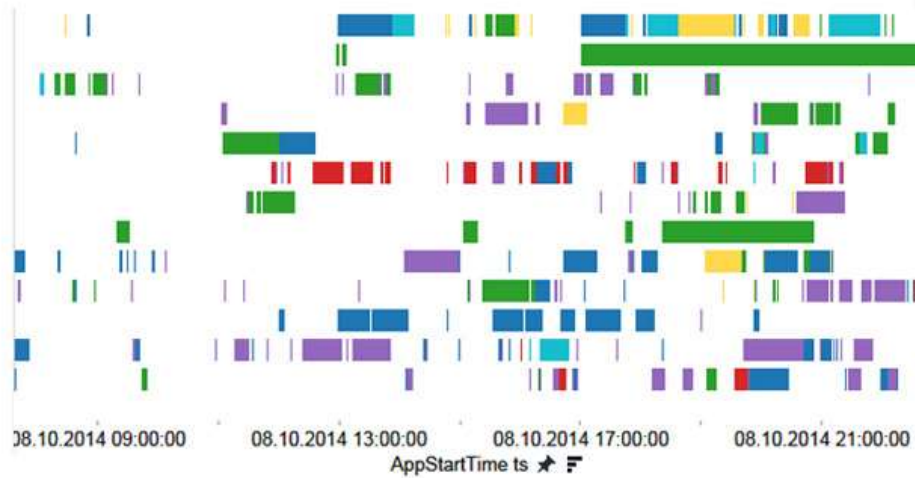
... typical attempts to find analysis help

DataFest

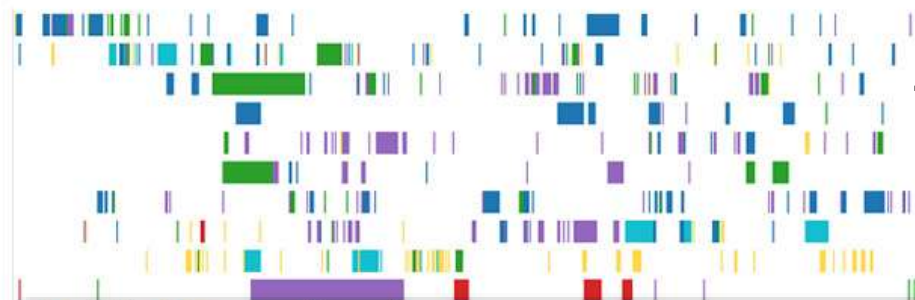




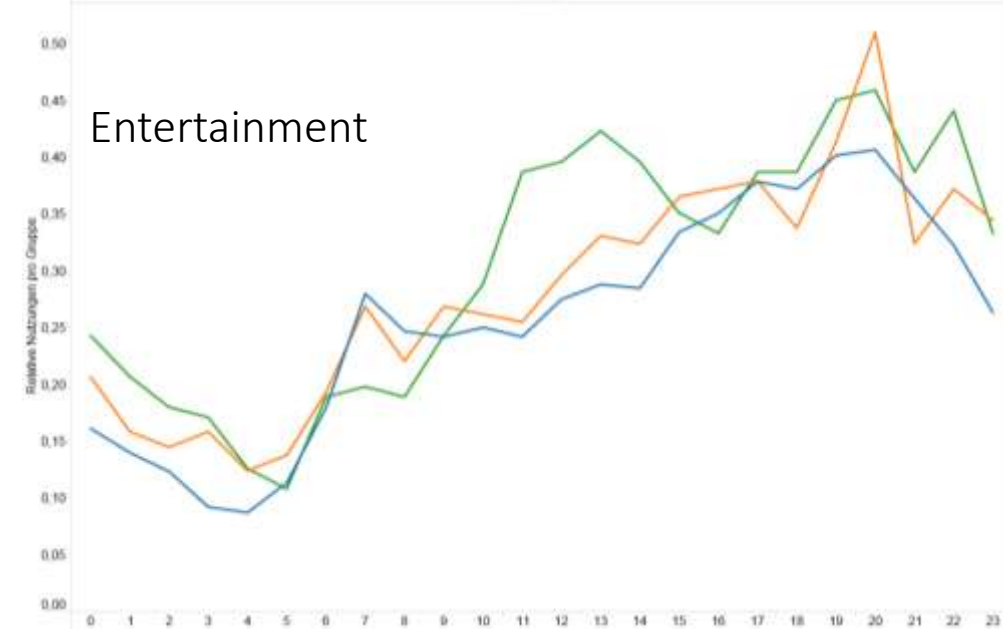
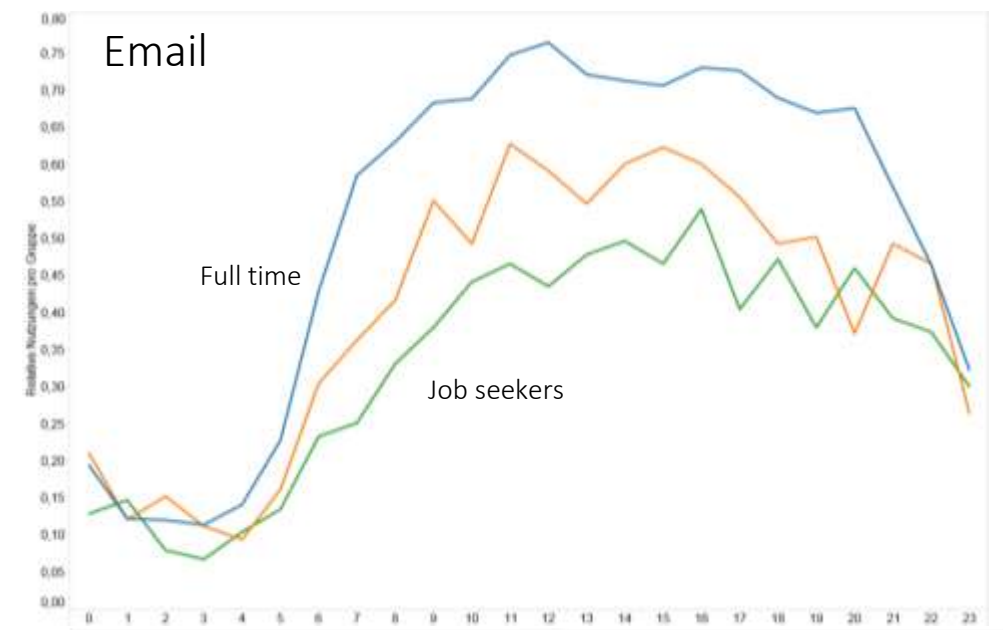
Full-time employed
→ App use past 5pm



Part-time employed
→ App use at noon



Job seekers
→ Continuous app use



Summary

1. Methodologists need to help teams find or create the right data for a specific purpose
2. It is easy to overlook what is missing. Privacy and confidentiality even more important
3. Methodologists should make use of all available data, and they can (with a few new skills)
4. Frameworks can help assess quality of single and combined data sources
5. Working in teams and with unlikely partners can accelerate

THANK YOU!

fkreuter@umd.edu

Twitter: @fraukolos