

# gesis

Leibniz-Institut  
für Sozialwissenschaften



## More data, more problems?

Chancen und Herausforderungen der  
datafizierten Gesellschaft für die  
sozialwissenschaftliche Forschung

*Johannes Breuer, Sebastian Stier, Pascal Siegers*

Mitglied der

*Leibniz*  
Leibniz-Gemeinschaft

# Hintergrund

- Nutzer\_innen digitaler Medien produzieren immer mehr Daten
- solche digitalen (Verhaltens-) Spurdaten sind auch für die sozialwissenschaftliche Forschung äußerst interessant
  - ▶ Webtracking-Daten
  - ▶ Social-Media-Daten
  - ▶ ...

# Selbstauskünfte & digitale Spurdaten

- Selbstauskünfte zu Mediennutzung oft nicht verlässlich (siehe z.B. [Scharkow, 2016](#); [Araujo et al., 2017](#))
  - ▶ soziale Erwünschtheit
  - ▶ Schwierigkeiten bei Erinnerung oder mentaler Berechnung (z.B. Fragen zu spezifischen/seltenen Nutzungsweisen oder längeren Zeiträumen)
- digitale Spurdaten alleine oft problematisch für sozialwissenschaftliche Fragestellungen
  - ▶ fehlende Informationen zu Personen
  - ▶ fehlende relevante abhängige Variable (z.B. Einstellungen oder Verhalten)
  - ▶ u.U. forschungsethisch schwierig: z.B. kein Informed Consent

# Data Linking

- Kombination von Befragungs- und digitalen Spurdaten als Antwort auf die Limitationen der beiden Datentypen ([Stier et al., 2019](#))
- verschiedene Wege der Verknüpfung
- spezifische methodische und ethische Herausforderungen für solche verbundenen Daten

# Linking-Typologie

Ex Ante	Ex Post
<p><b>Aggregatebene</b></p> <ul style="list-style-type: none"> <li>Analyse von Audience Overlap (z.B. <a href="#">Mukerjee et al., 2018</a>; <a href="#">Nelson &amp; Webster, 2017</a>)</li> <li>Analyse aggregierter Publikums-/Nutzer_innenstatistiken (z.B. <a href="#">Nelson &amp; Webster, 2017</a>)</li> </ul>	<p><b>Aggregatebene</b></p> <ul style="list-style-type: none"> <li>Verknüpfung auf von Ebene von...           <ul style="list-style-type: none"> <li>Zeit (z.B. <a href="#">Mellon, 2014</a>; <a href="#">O'Connor et al., 2010</a>; <a href="#">Stier et al., 2018</a>)</li> <li>Themen (z.B. <a href="#">Pasek et al., 2019</a>)</li> <li>Geographie (z.B. <a href="#">Beauchamp, 2017</a>)</li> </ul> </li> </ul>
	<p><b>Öffentliche Akteure</b></p> <ul style="list-style-type: none"> <li>z.B. Politiker_innen oder Organisationen (siehe z.B. <a href="#">Karlsen &amp; Enjolras, 2016</a>; <a href="#">Quinlan et al., 2017</a>)</li> </ul>
<p><b>Individualebene</b></p> <ul style="list-style-type: none"> <li>Informed Consent in Befragung für Erhebung der digitalen Spurdaten:           <ul style="list-style-type: none"> <li>Webseitenbesuche (z.B. <a href="#">Guess, 2015</a>; <a href="#">Jürgens et al., 2019</a>; <a href="#">Möller et al., 2019</a>; <a href="#">Vraga &amp; Tully, 2018</a>)</li> <li>Smartphone-Daten (z.B. <a href="#">Boase &amp; Ling, 2013</a>; <a href="#">Jürgens et al., 2019</a>; <a href="#">Kreuter et al., 2019</a>)</li> <li>Sensordaten (z.B. <a href="#">Génois, Zens, Lechner, Rammstedt, &amp; Strohmaier, 2019</a>)</li> </ul> </li> </ul>	<p><b>Individualebene</b></p> <ul style="list-style-type: none"> <li>Informed Consent in Befragung für retrospektive Erhebung:           <ul style="list-style-type: none"> <li>über Social-Media-APIs (z.B. <a href="#">Al Baghal et al., 2019</a>; <a href="#">Haenschen, 2019</a>; <a href="#">Hofstra, Corten, van Tubergen, &amp; Ellison, 2017</a>; <a href="#">Hopp, Vargo, Dixon, &amp; Thain, 2018</a>; <a href="#">Vaccari et al., 2015</a>; <a href="#">Wells &amp; Thorson, 2017</a>)</li> <li>über Datenspende; z.B. persönliches Twitter-Archiv oder Export eigener Facebook-Daten (siehe z.B. <a href="#">Thorson et al., 2019</a>)</li> </ul> </li> </ul>

# Beispiele spezieller Herausforderungen

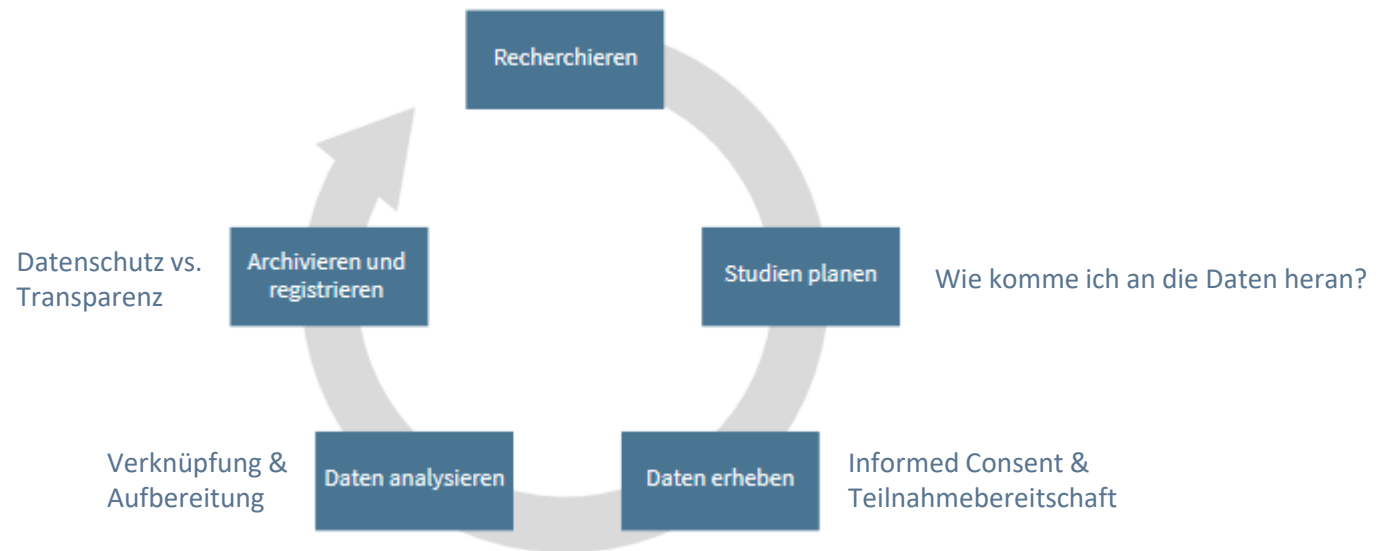
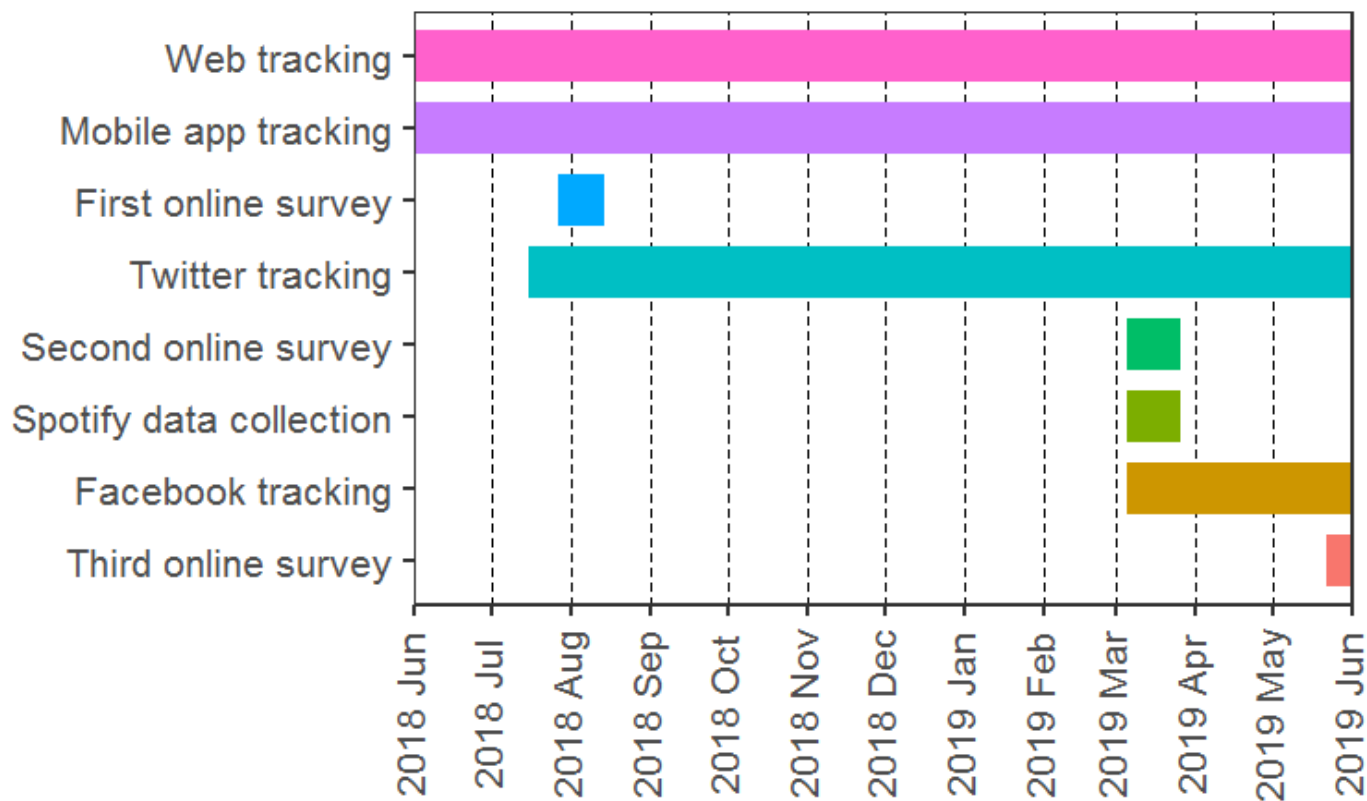


Abbildung: <https://www.gesis.org/angebot>

# Unser Projekt

- *respondi*-Web-Tracking-Panel
  - ▶ Daten für 1 Jahr (Juni 2018 – Mai 2019)
  - ▶ ~ N = 2000 Teilnehmer\_innen pro Monat
  - ▶ insg. ca. 94 Mio. Datenpunkte
  - ▶ für Teil der Teilnehmer\_innen: zusätzlich Tracking der mobilen App-Nutzung
- 3 Online-Befragungen der Panelist\_innen
  - ▶ thematischer Fokus: Mediennutzung & Politik
- Social-Media-Daten
  - ▶ Twitter
  - ▶ Facebook
  - ▶ Spotify

# Im Projekt erhobene Datentypen

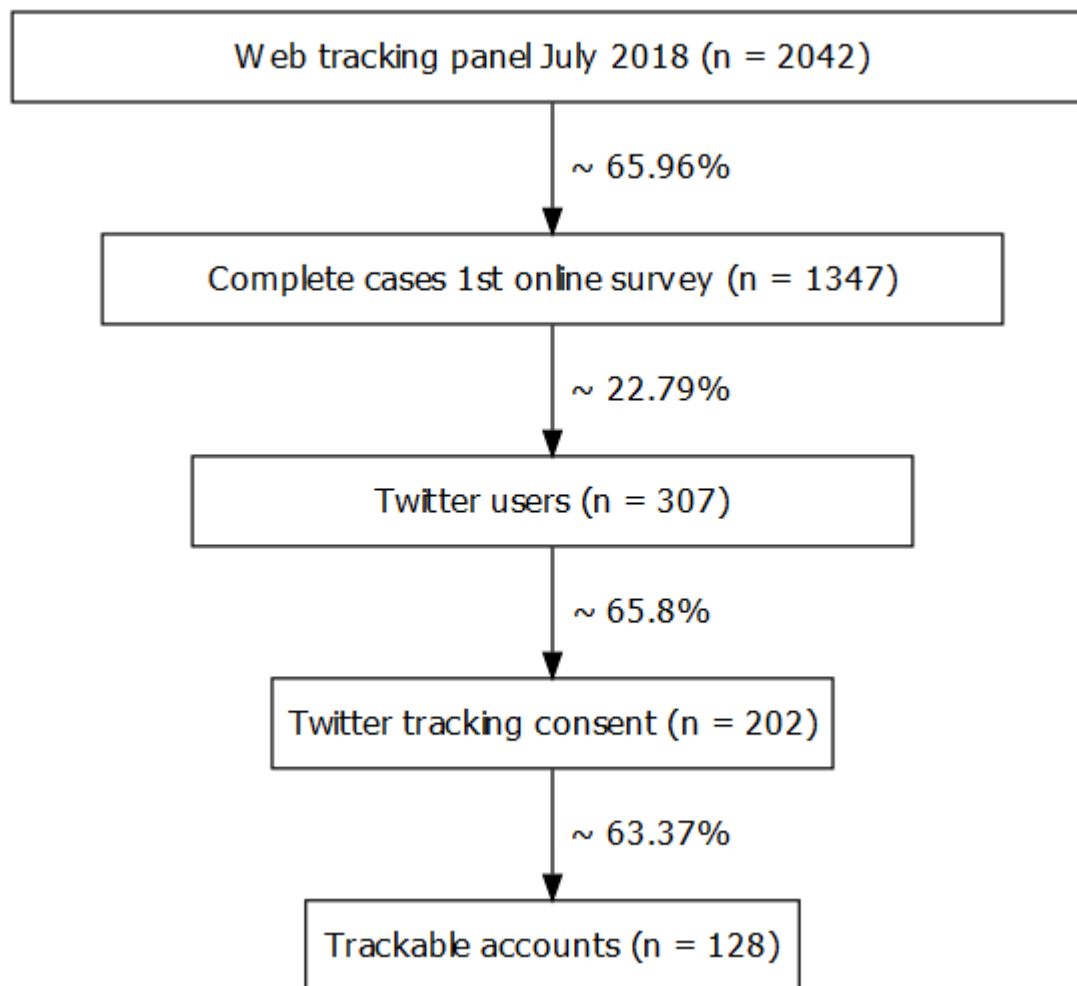




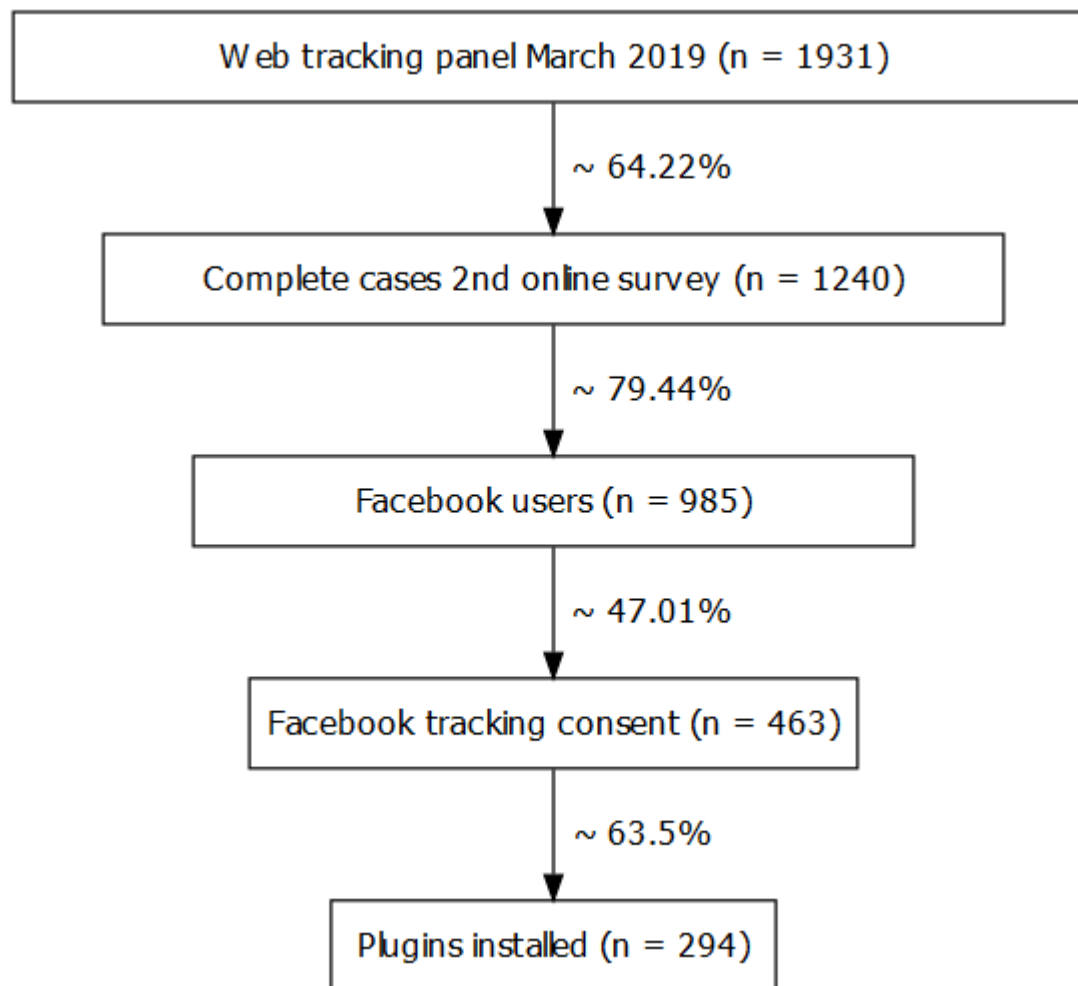
# Social-Media-Daten

- Twitter
  - ▶ kontinuierliches Tracking via Streaming API
- Facebook
  - ▶ Browser-Plugin ([Haim & Nienierza, 2019](#))
    - für Firefox & Chrome
    - sammelt öffentliche Posts (+ einige Metadaten zu diesen) aus dem Feed der Nutzer\_innen
- Spotify
  - ▶ Web-App von der KU Leuven
    - sammelt die 50 zuletzt abgespielten Songs, Playlists & Präferenzen
- Für alle Daten:
  - ▶ kurzer Informed Consent in Online-Befragung
  - ▶ ausführlichere Informationen (zu Datenschutz etc.) auf Webseite (im Fragebogen verlinkt)
  - ▶ Informed Consent & Privacy Information von [Al Bahgal et al. \(2019\)](#) übernommen, übersetzt und angepasst

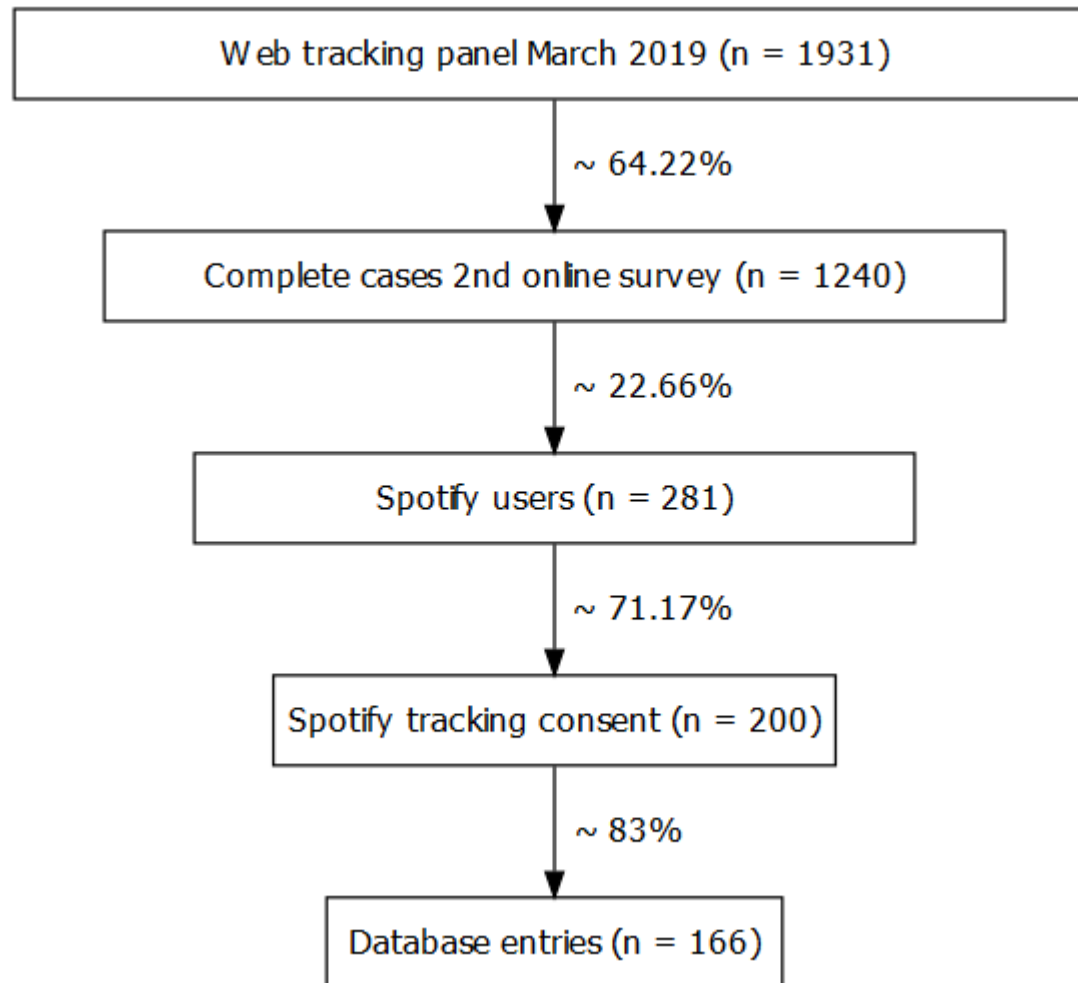
# Twitter-Daten



# Facebook-Daten



# Spotify-Daten



## Nächste Schritte & offene Fragen

- systematischer Bias in den Dropout-Stufen?
- Qualität/Nutzbarkeit der Social-Media-Daten?
  - ▶ Nutzungsaktivität quantitativ -> Datenpunkte pro Person
  - ▶ Nutzungsaktivität qualitativ -> z.B. passive vs. aktive Nutzung bei Twitter; Facebook als Nachrichtenquelle?
- geplante Lösung für Data Sharing
  - ▶ frei zugängliche aggregierte Daten über [GESIS datorium](#)
  - ▶ Individualdaten über [GESIS Secure Data Center](#)

# Vor- und Nachteile unseres Ansatzes

Vorteile	Nachteile
<ul style="list-style-type: none"><li>• Daten auf individueller Ebene</li><li>• große &amp; heterogene Stichprobe</li><li>• große Bandbreite an Daten</li><li>• Informierte Einwilligung der Teilnehmer_innen</li><li>• unaufwändiger Zugang zu Webtracking-Daten</li><li>• Facebook-Daten nicht über API</li></ul>	<ul style="list-style-type: none"><li>• Bias(es) in der Stichprobe</li><li>• Kosten für Webtracking-Daten</li><li>• Twitter &amp; Spotify über API<ul style="list-style-type: none"><li>• Risiko der Veränderung oder Schließung von APIs (<a href="#">Bruns, 2019</a>; <a href="#">Freelon, 2018</a>; <a href="#">Puschmann, 2019</a>)</li></ul></li><li>• Veränderungen in der Struktur des Facebook-Feeds u.U. problematisch für Browser-Plugin</li></ul>

# Alternative Datenzugänge

Datenart	Methode	Vorteile	Nachteile
<b>Web- und mobile App-Nutzung</b>	Eigene Software (+ eigenes Panel)	<ul style="list-style-type: none"> <li>Mehr Kontrolle über Rekrutierung und Erhebung</li> </ul>	<ul style="list-style-type: none"> <li>Hoher Aufwand (Erstellung der Software/App + Support; Rekrutierung der Teilnehmer_innen)</li> </ul>
<b>Social-Media-Daten</b>	„Datenspende“ (Nutzer_innen exportieren eigene Daten und teilen diese mit den Forscher_innen; siehe <a href="#">Halavais, 2019</a> ; <a href="#">Thorson et al., 2019</a> )	<ul style="list-style-type: none"> <li>Direkte Einbindung der Teilnehmer_innen</li> <li>Transparenz</li> </ul>	<ul style="list-style-type: none"> <li>Aufwand für Teilnehmer_innen (Export der Daten nicht selbsterklärend und i.d.R. nicht direkt möglich)</li> <li>Lösung zum Teilen der Daten (inkl. Anonymisierung) nötig</li> </ul>

Zusätzlich Variation der Rekrutierungswege/-reihenfolge möglich:

- z.B. Befragung (-> Webtracking) -> Social Media
- Reihenfolge hat Einfluss auf Zusammensetzung der Stichprobe & Biases

# Zusammenfassung

- Befragungs- und digitale Spurdaten haben jeweils eigene Limitationen
- Verknüpfung dieser beiden Datentypen ist vielversprechend für die sozialwissenschaftliche Forschung
- Verknüpfung bringt spezifische Herausforderungen für alle Phasen des Forschungsdatenzklus mit sich
- es gibt verschiedene Möglichkeiten, die Daten zu erheben und sie zu verknüpfen
- richtige Auswahl der Methoden hängt von der Fragestellung sowie den verfügbaren Ressourcen eines Projekts ab