

Privacy through the Lens of Data Flows

Shujun LI (李树钧)

Director, [Institute of Cyber Security for Society \(iCSS\)](#)

Professor of Cyber Security, School of Computing
University of Kent, UK

<http://www.hooklee.com/>

 @hooklee

Privacy in the cyber-physical world



- Information about your machine / operating system
 - IP addresses, time zone, language, screen size, fonts installed, ...
- Sensor data made accessible via web browsers
 - GPS coordinates, camera data, microphone data, ...
- Information stored within or by web browsers
 - Web browser type and version, web cookies, web browsing histories and caches, search keywords on search engines, auto completion data, bookmarks, account details (e.g., Google account for Google Chrome), web browser extensions installed, ...
- Data stored in web browser extensions
 - May or may not be in the web browser (can be stored online)
- **User-generated content (UGC)**
 - Search queries, online profiles, online posts, ...
- ...

What online services?

- Online social networks (OSN) or social media (SM)

- Microblogging systems

- Online chat rooms

- Instant messaging systems

- Web forums

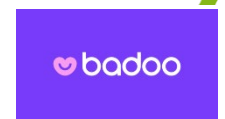
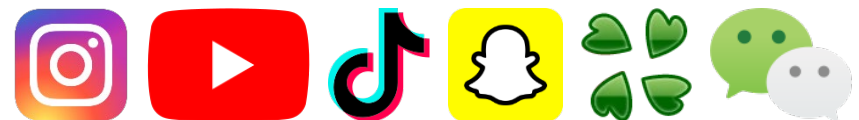
- Image and video sharing portals

- Search engines

- **Any** websites or web-based applications and (mobile) apps supporting UGCs

- App marketplaces, online shopping, online booking, news outlets, online dating, e-petition, online maps, ...

- ... **Booking.com** **JustGiving™**



What UGC?

- Profile data
 - Mostly personal data: name(s), facial image (profile image), addresses, geo-location, gender, age, profession, personal interests, financial status, information of family members, etc.
- Online posts, comments, reviews, ratings, ...
 - They may contain (a lot of!) personal data.
 - And they are normally linked to profile data.
- ...
- Personal and sensitive data of **yourself** and **other people**
 - Neither you nor other people may have noticed you did it.
 - If you disclosed other people's personal data and they saw it, they may not feel comfortable to (publicly) object to your disclosure.
 - You can't control what others will share about you!

Your data may not be just yours!

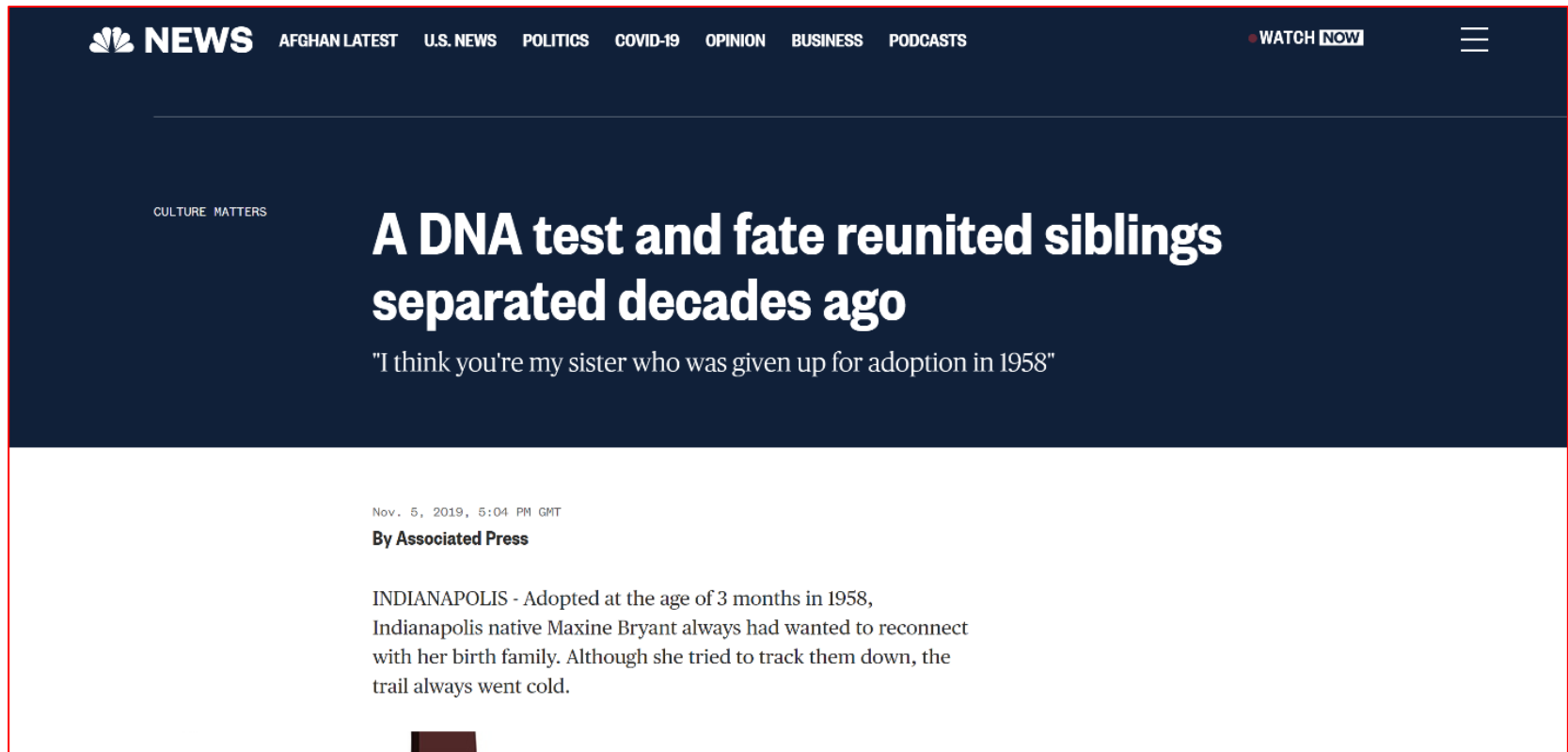
- Example: Your (sur)name(s) and your DNA data are your whole **extended** family's, ...



Family Tree

Your data may not be just yours!

- Example: Your DNA data is your whole **extended** family's, including "unknown" relatives.

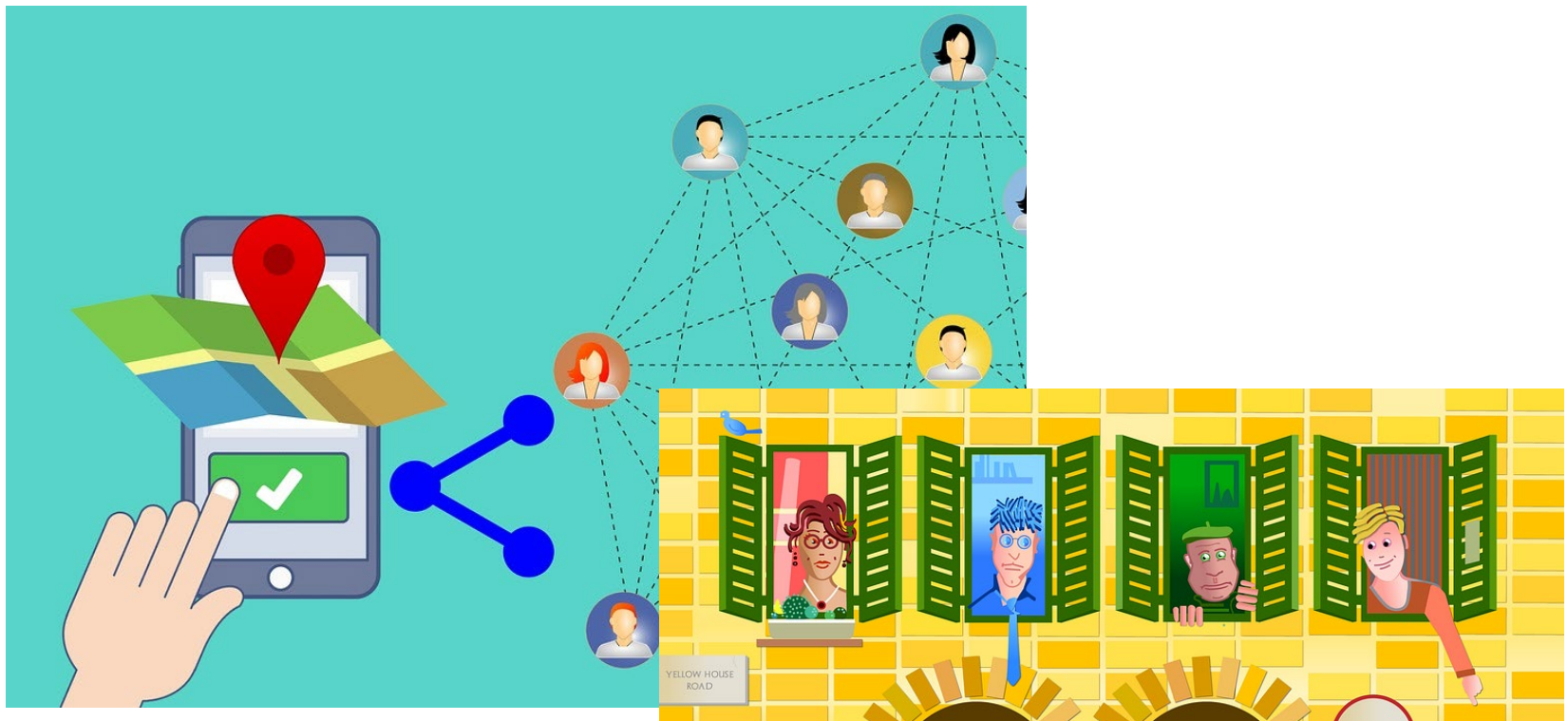


The screenshot shows the top navigation bar of the NBC News website with categories like 'AFGHAN LATEST', 'U.S. NEWS', 'POLITICS', 'COVID-19', 'OPINION', 'BUSINESS', and 'PODCASTS'. The main article title is 'A DNA test and fate reunited siblings separated decades ago' with a sub-headline 'I think you're my sister who was given up for adoption in 1958'. The byline is 'By Associated Press' and the date is 'Nov. 5, 2019, 5:04 PM GMT'. The lead paragraph reads: 'INDIANAPOLIS - Adopted at the age of 3 months in 1958, Indianapolis native Maxine Bryant always had wanted to reconnect with her birth family. Although she tried to track them down, the trail always went cold.'

Source: <https://www.nbcnews.com/news/nbcblk/dna-test-fate-reunited-siblings-separated-decades-ago-n1076006>

Your data may not be just yours!

- Example: Your geo-location, life patterns, religion, cultural background also belong to your neighbours, friends and those living nearby!



Where are the data?

- On your devices or on the Internet/Web?
- Where exactly?



Example: CCTV and smart cities

- CCTV and cameras everywhere?



Example: Privacy in public spaces



A project of **Global Voices**, we are a global anti-censorship network of bloggers and activists dedicated to protecting freedom of expression online.

f t r Newsletter About Us Search

Activism Privacy Netizen Report About Our Community Archives

The Russian Art of Meta-Stalking

Posted 7 April 2016 21:38 GMT



Photo: Shane Huang / Flickr

The next time you ride



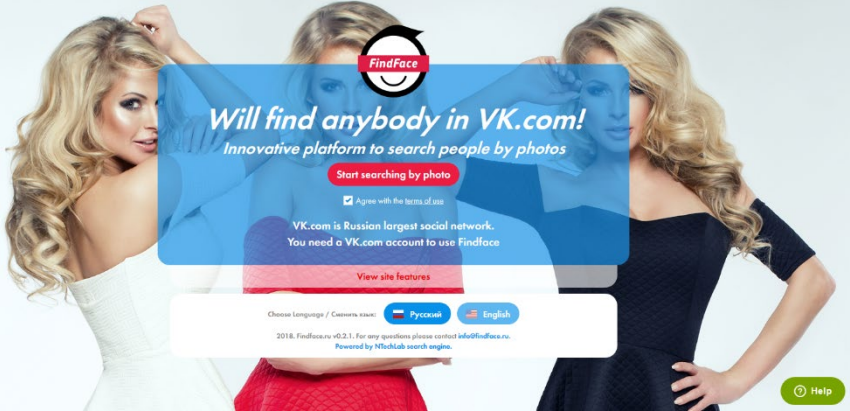
Written by Kevin Rothrock

Read this post in [Ελληνικά](#), [Español](#), [English](#)

Harassment, Media, Tech & Tools

Russia

Eastern & Central Europe



Who are data consumers and why?



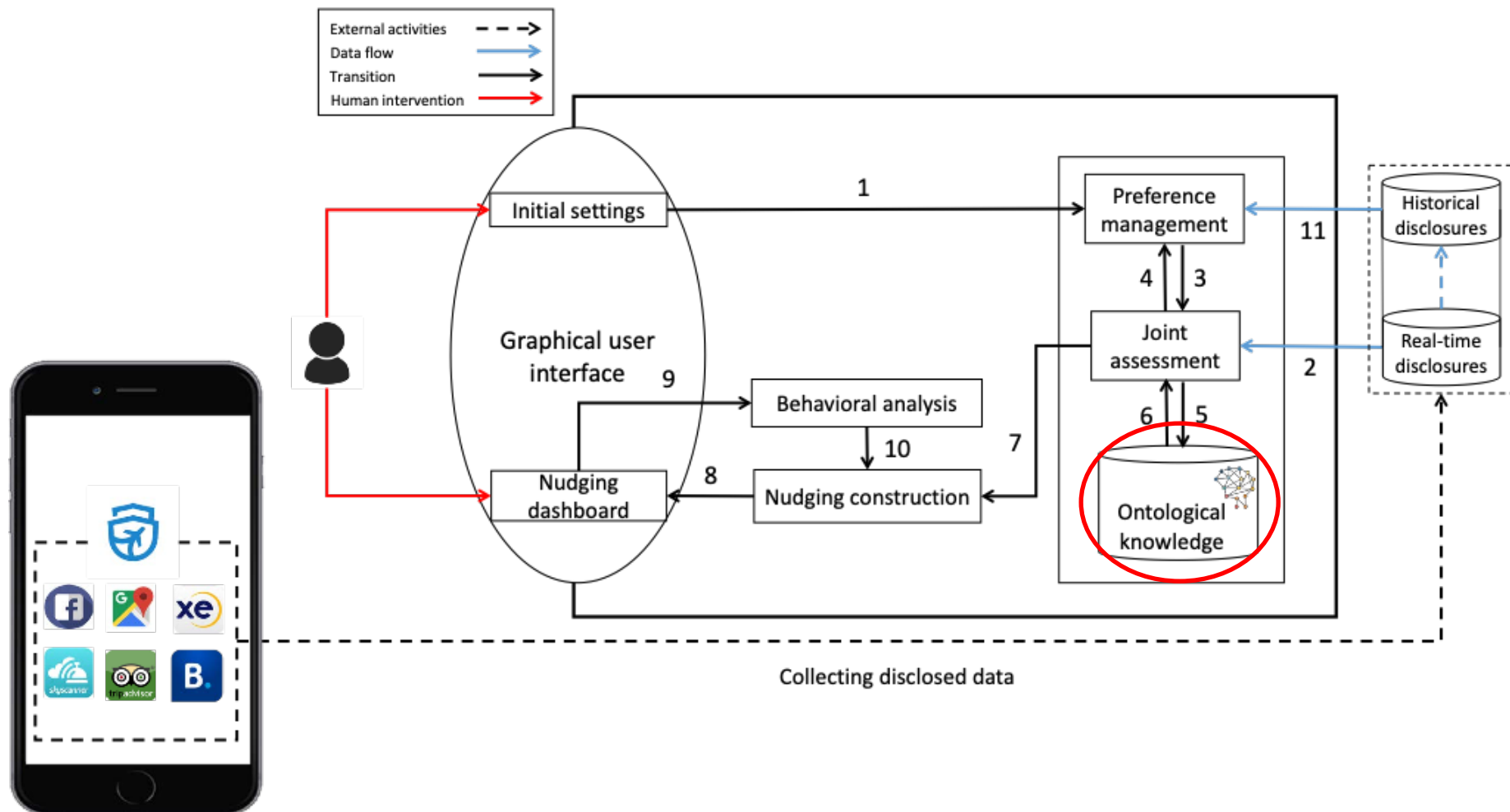
The **PRIVELT** Project

- Title: PRiVacy-aware personal data management and Value Enhancement for Leisure Travellers (**PRiVELT**)
- Funder:  **Engineering and Physical Sciences Research Council**
- Call: Trust, Identity, Privacy and Security in the Digital Economy 2.0 (2018)
- Budget: £~1.4m
- Duration: 10/2018 – 06/2023 (57 months)
- Website: <https://privelt.ac.uk/>

(Part of) The (former) project team

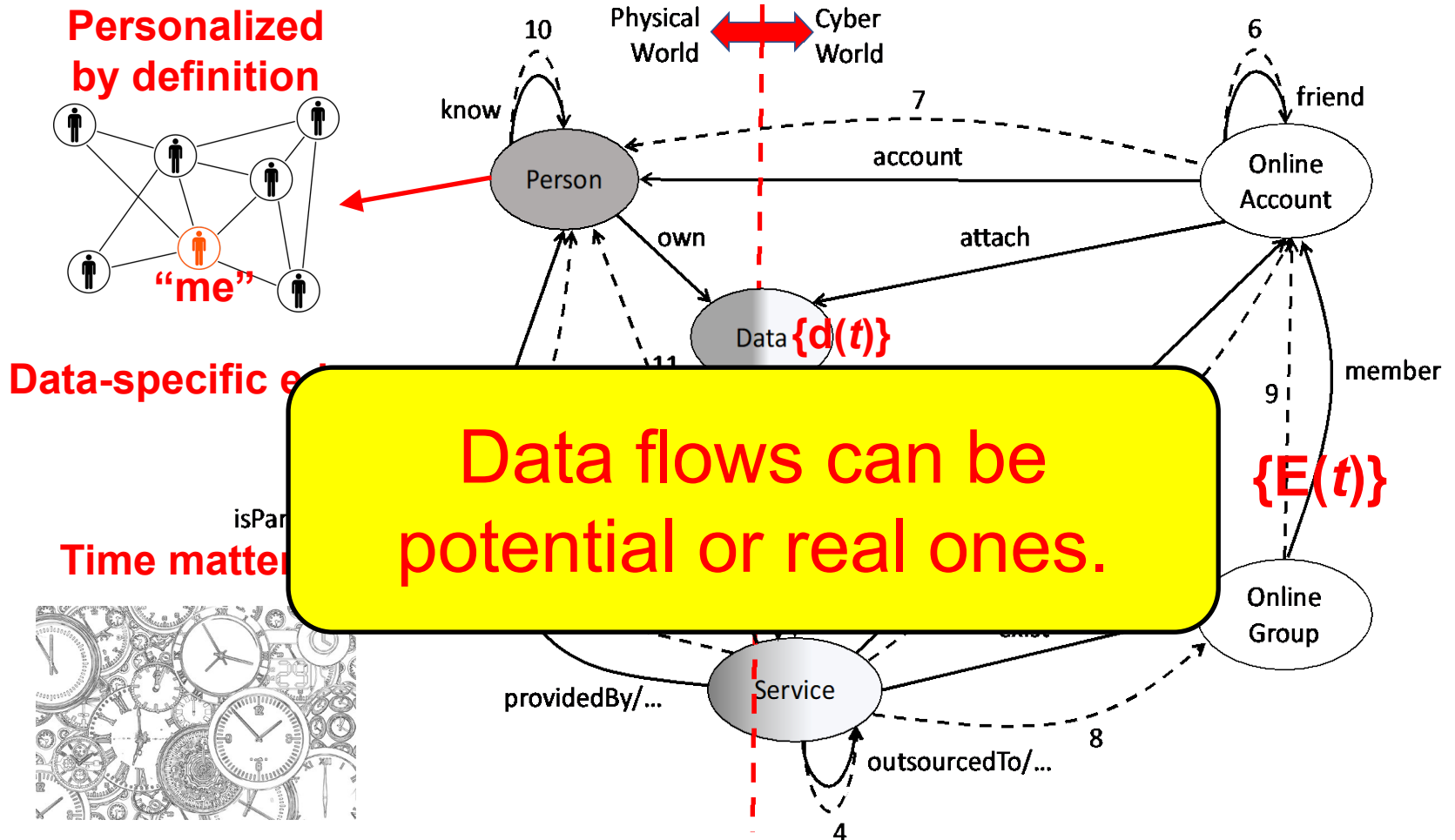


The vision: user-centric, server-less, from privacy awareness to nudging



Yang Lu, Shujun Li, Athina Ioannou and Iis Tussyadiah (2019) [From Data Disclosure to Privacy Nudges: A Privacy-aware and User-centric Personal Data Management Framework](#). In *Proc. DependSys 2019*, Springer. doi:10.1007/978-981-15-1304-6_21

Data sharing (flow) ontology



Yang Lu and Shujun Li (2020) [From Data Flows to Privacy Issues: A User-Centric Semantic Model for Representing and Discovering Privacy Issues](#). In *Proc. HICSS 2020*, University of Hawai'i at Mānoa. doi: 10.24251/HICSS.2021.651

Is a data flow graph complex?

- Number of nodes: **large**
 - “Me”: the “center” / owner of the graph
 - All data item and data packages about “me”
 - All people your data can flow to (could be **anyone**)
 - All physical and online services you data can flow to
 - All organizations your fata can flow to
 - *Exercise: Check your password manager!*
- Number of edges: **huge**
 - Relationships between different types of nodes
 - Often more than one edge between any two nodes
 - *Exercise: Check the data dashboard of your Google Account or your account with any other online service!*

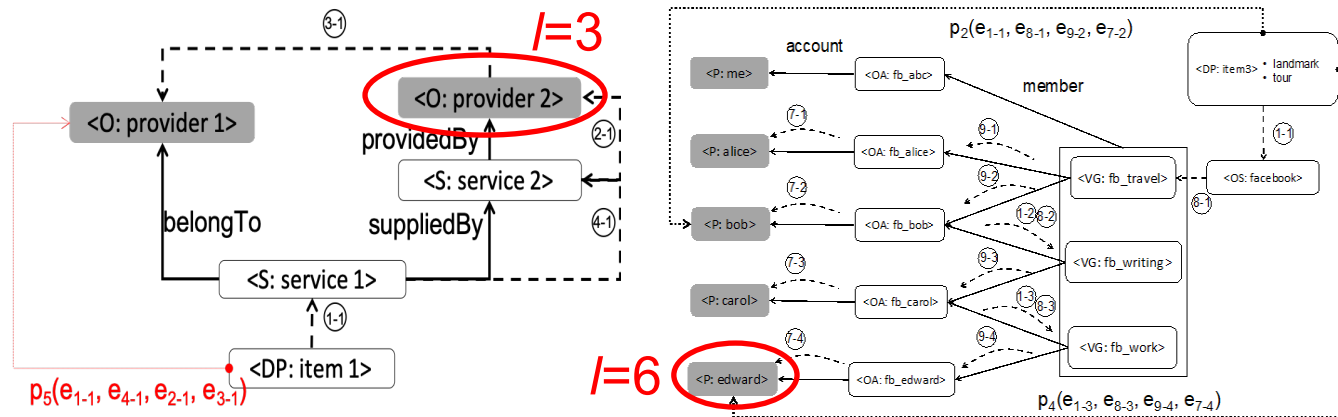
- Out-degree (of “me” node), given a time window
 - The amount of data shared
- Average nodal degree (of a data consumer node)
 - The average amount of “my” data disclosed to that data consumer
- Node / Link connectivity (of the whole graph)
 - The number of “essential” data consumers / data sharing activities
- Centrality metrics (of data consumer nodes)
 - For identifying major (potentially “hidden”) data consumers
- The longest path(s) originating from “me”
 - For identifying the most “hidden” data consumer(s)
- Network type (of the whole graph)
 - Small-world network, scale-free network or something else?
- ...

“Topological” privacy issues

- A specific privacy issue with a specific data item or a data package corresponds to a **data flow path**.
- A specific privacy issue with more than one data items and/or data packages corresponds to **a set of data flow paths**.
- A specific **type** of privacy issues of one or more data item / package type(s) is **a set of sets of data flow paths**.
- All of them can be **described** and **potentially detected** via their **common topological features**.

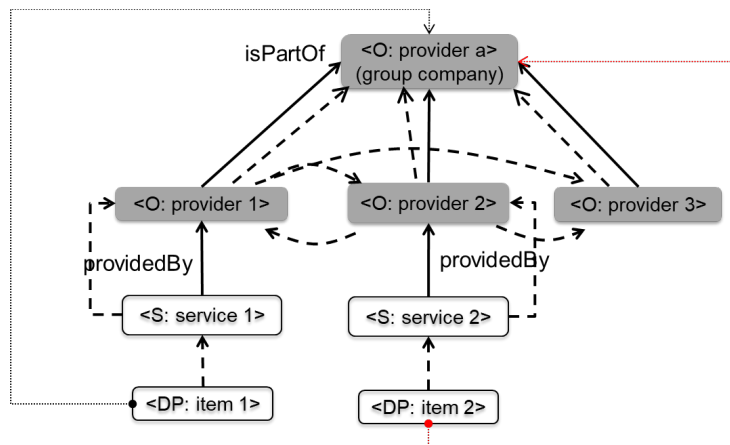
“Topological” privacy issue #1

- Data shared with (potentially) **unknown consumers**
 - **Hypothesis**: the longer the data flow path length between “me” and a data consumer node is, the more likely the user is unaware of the data consumer
 - **Risk assessment**: $r=f(l)$, where l is the path length
 - **Visualization**: show a ranked list of all potential unknown data consumers with decreasing values of r
 - **Detection (naïve method)**: $r > r_t \Rightarrow$ issue an alert



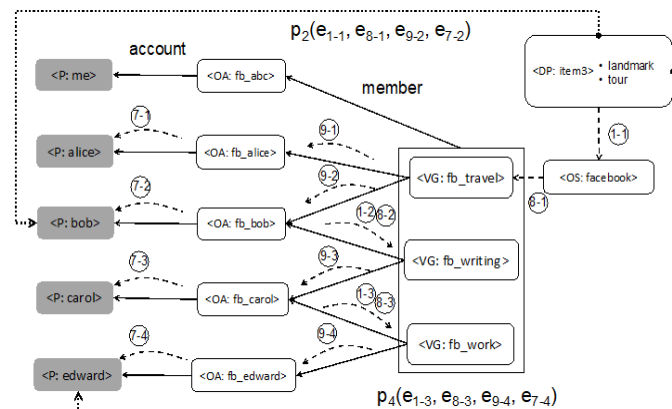
“Topological” privacy issue #2

- Indirect (= potentially unknown) **data aggregator**
 - **Hypothesis**: given a tree whose root node is a data consumer, the taller the tree is, the more likely the root node is an unknown (super) data aggregator
 - **Risk assessment**: $r=f(h)$, where h is the tree's height
 - **Visualization**: show a ranked list of all potential data aggregators with decreasing values of r
 - **Detection (naïve method)**: $r > r_t \Rightarrow$ issue an alert



“Topological” privacy issue #3

- Data shared with **too many consumers**
 - **Hypothesis**: given a tree whose root is a data node, the bigger the tree is, the more likely the data has been over-shared too much
 - **Risk assessment**: $r=f(n)$, where n is the total number of nodes in the tree minus 1 (the root node)
 - **Visualization**: show the whole tree
 - **Detection (naïve method)**: $r > r_t \Rightarrow$ issue an alert



Automatic reasoning is possible!

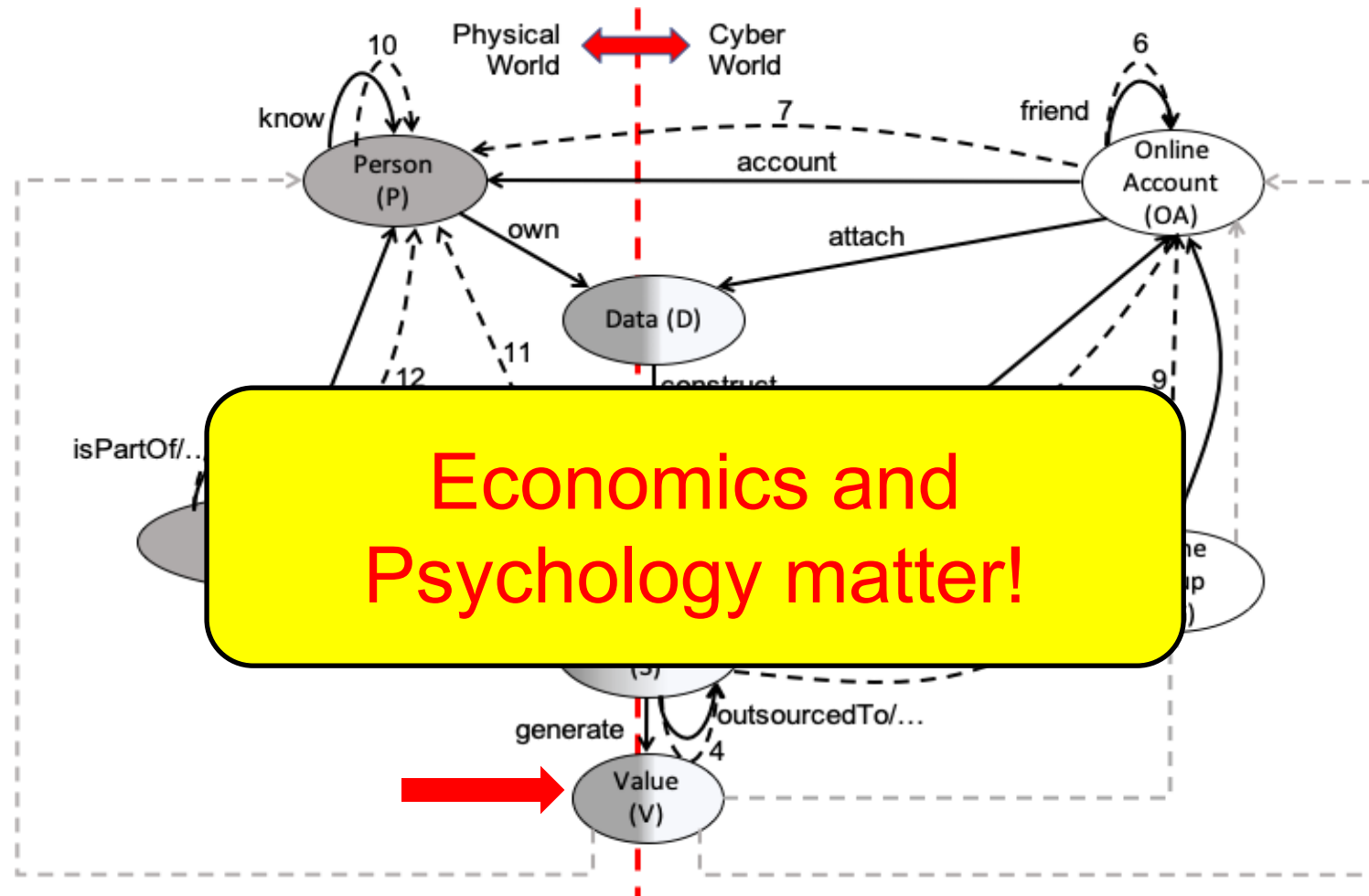
The image displays three overlapping screenshots of a DL query interface, demonstrating automatic reasoning. Each screenshot shows a query input field, execution buttons, and a results panel.

Top-left screenshot:
DL query: Query (class expression)
Service_Provider that access some (Data that has some Sensitive)
Buttons: Execute, Add to ontology
Query results: Direct superclasses (1 of 1)
● Service_Provider
Instances (11 of 11)
List of instances: Agoda, Booking.com, GoToGate, Kayak, OpenTable, Princline.com, Rentalcars.com, TravelJigsaw, flygresor.se, mytrip.com, supersaver. Red boxes highlight the top five and bottom five items.

Top-right screenshot:
DL query: Query (class expression)
Person that access some (Data_Package that has some Location) and access some (Data_Package that has some Event)
Buttons: Execute, Add to ontology
Query results: Instances (1 of 1)
◆ edward

Bottom-right screenshot:
DL query: Query (class expression)
Data_Package that has some Entertainment that (flowTo some Work)
Buttons: Execute, Add to ontology
Query results: Instances (1 of 1)
◆ item3

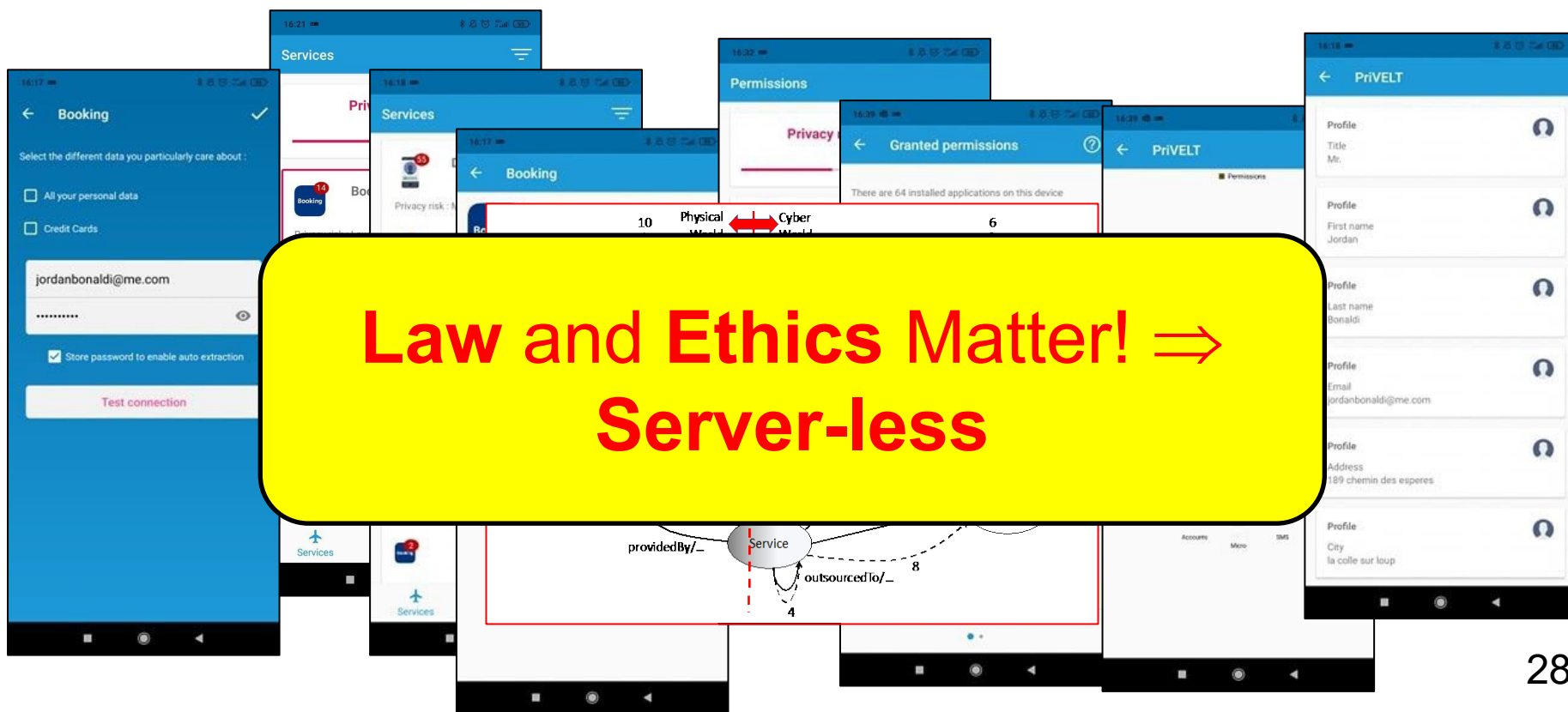
Adding returned values



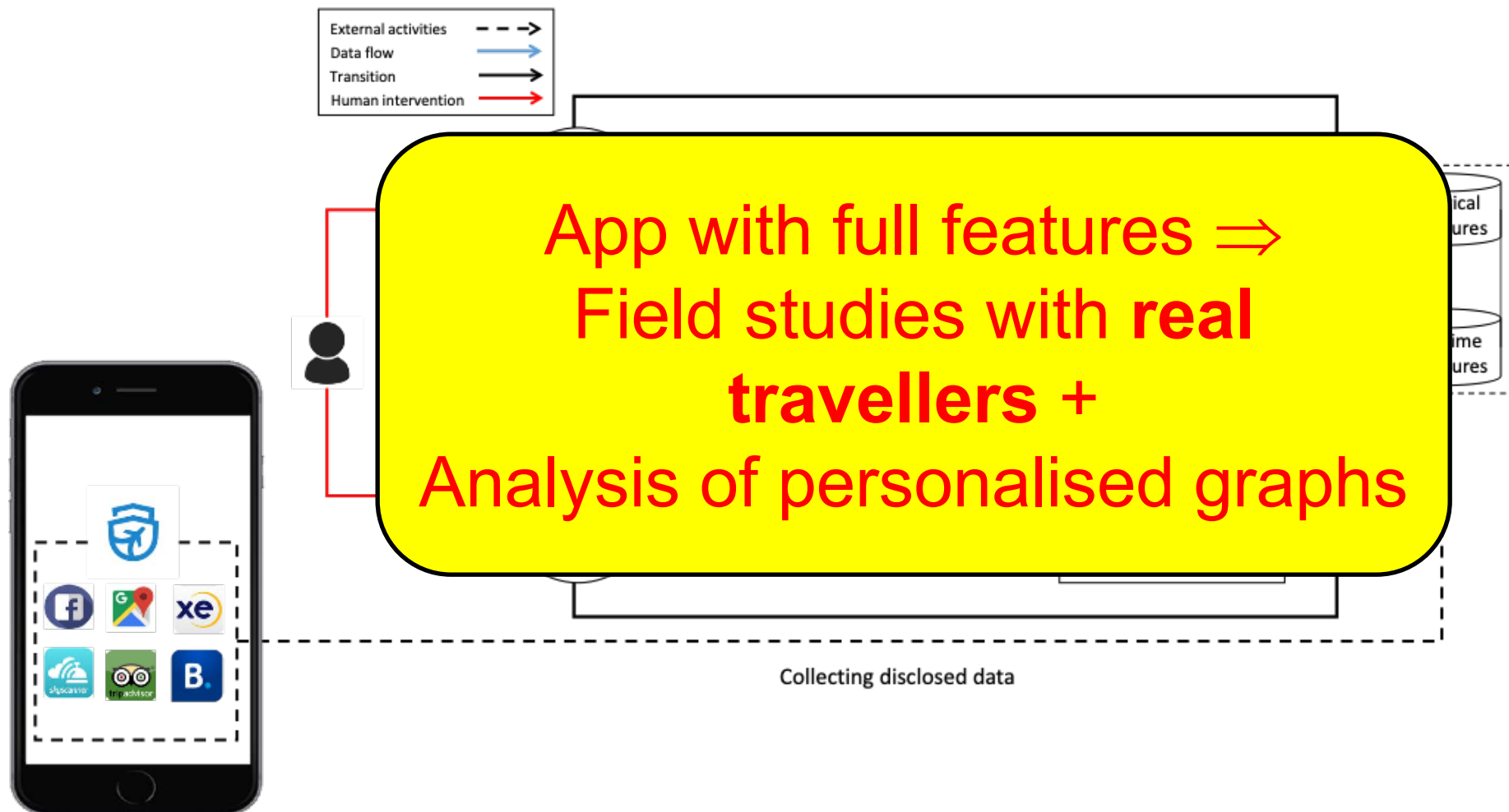
Yang Lu and Shujun Li (2022) [From Data Flows to Privacy-Benefit Trade-offs: A User-Centric Semantic Model](#). *Security and Privacy*, 5(4):e225, 24 pages, [John Wiley & Sons, Inc.](#) doi: 10.1002/spy2.225

Personalized data flow graphs

- User-centric and service-independent tools are needed to build “my” data flow graph.
- ⇒ We are developing an Android app for this.

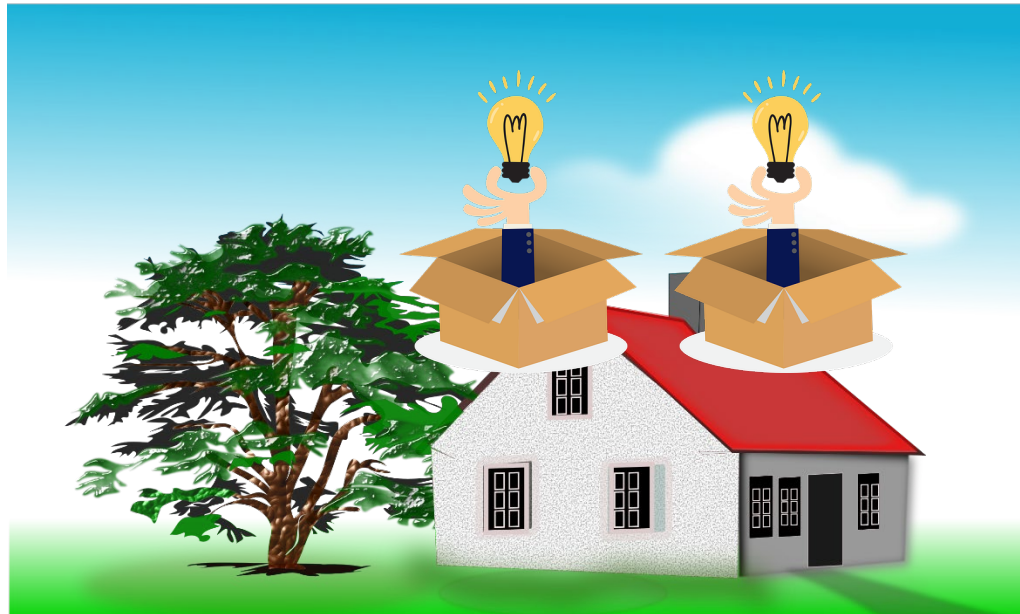


The vision: user-centric, server-less, from privacy awareness to nudging



Yang Lu, Shujun Li, Athina Ioannou and Iis Tussyadiah (2019) [From Data Disclosure to Privacy Nudges: A Privacy-aware and User-centric Personal Data Management Framework](#). In *Proc. DependSys 2019*, Springer. doi:10.1007/978-981-15-1304-6_21

Take-Home Messages



Anything useful from the talk?

- Data privacy problems **in all domains** can be studied using **data flow graphs**.
- Such data flow graphs are **personalised** around a node called “me”.
 - **User-centric** tools are needed to engage users.
- There are essential research questions across **multiple disciplines**.
 - There are **legal**, **ethical** and **economic** implications!
- Although PriVELT as a project has completed, our work is still **ongoing**.
 - ⇒ We welcome different types of **collaborations!**

Privacy through lens of data flows



Shujun LI (李树钧)

<http://www.hooklee.com/>