# Location Data: GPS, Wi-Fi, and Spatial Analytics

Class 2 of Digital Data Flows Masterclass:
*Emerging Technologies*

27 November, 2018 | Brussels

# DIGITAL DATA FLOWS MASTERCLASS: EMERGING TECHNOLOGIES

VUB · BRUSSELS PRIVACY HUB · FUTURE OF PRIVACY FORUM

## Curriculum

Session 1: Artificial Intelligence and Machine Learning - featuring Dr. Swati Gupta, Assistant Professor in the H. Milton Stewart School of Industrial and Systems Engineering at Georgia Tech; and Dr. Oliver Grau, Chair of ACM's Europe Technology Policy Committee, Intel Automated Driving Group, and University of Surrey

Session 2: Location Data: GPS, Wi-Fi, and Spatial Analytics

Session 3: Advertising Technologies: Online Data Flows, Behavioral Targeting, and Cross-Device Tracking

Session 4: Mobile Apps: Operating Systems, Software Development Kits (SDKs), and User Controls

Session 5: Transportation and Mobility: Video Analytics, Sensors, and Connected Infrastructure

Session 6: Biometric Data: Facial Recognition, Voice, and Digital Fingerprints

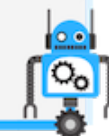Session 7: Tracking in Physical Spaces: Retail Technologies, Smart Homes, and the "Internet of Things"

Session 8: De-Identification: Multi-party Computing, Differential Privacy, and Homomorphic Encryption

## Date*

25 October, 2018 - side event, ICDPPC (Brussels) (with remote participation)

November 2018 - Brussels (with remote participation)

January 2019 - Brussels (with remote participation)

March 2019 - Virtual

April 2019 - Virtual

June 2019 - Virtual

July 2019 - Virtual

Sept. 2019 - Virtual

*dates may change.

All sessions are free and will support remote participation.
Priority registration may be held for government staff.
Enroll and receive updates on the full course at: www.fpf.org/classes
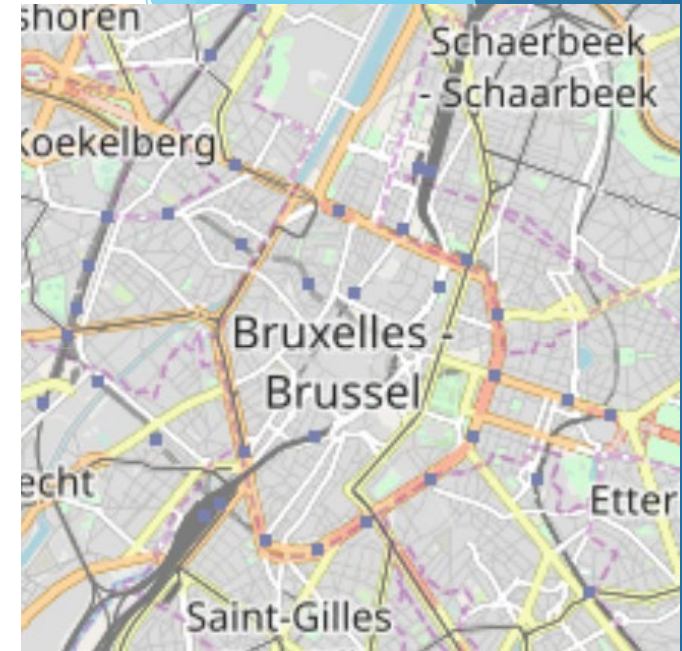
FUTURE OF PRIVACY FORUM

# AGENDA

I. Introduction to Geo-Location Data
II. Sources of Data: *Mobile Sensors, Wi-Fi Analytics*
III. Data Flows & Case Studies
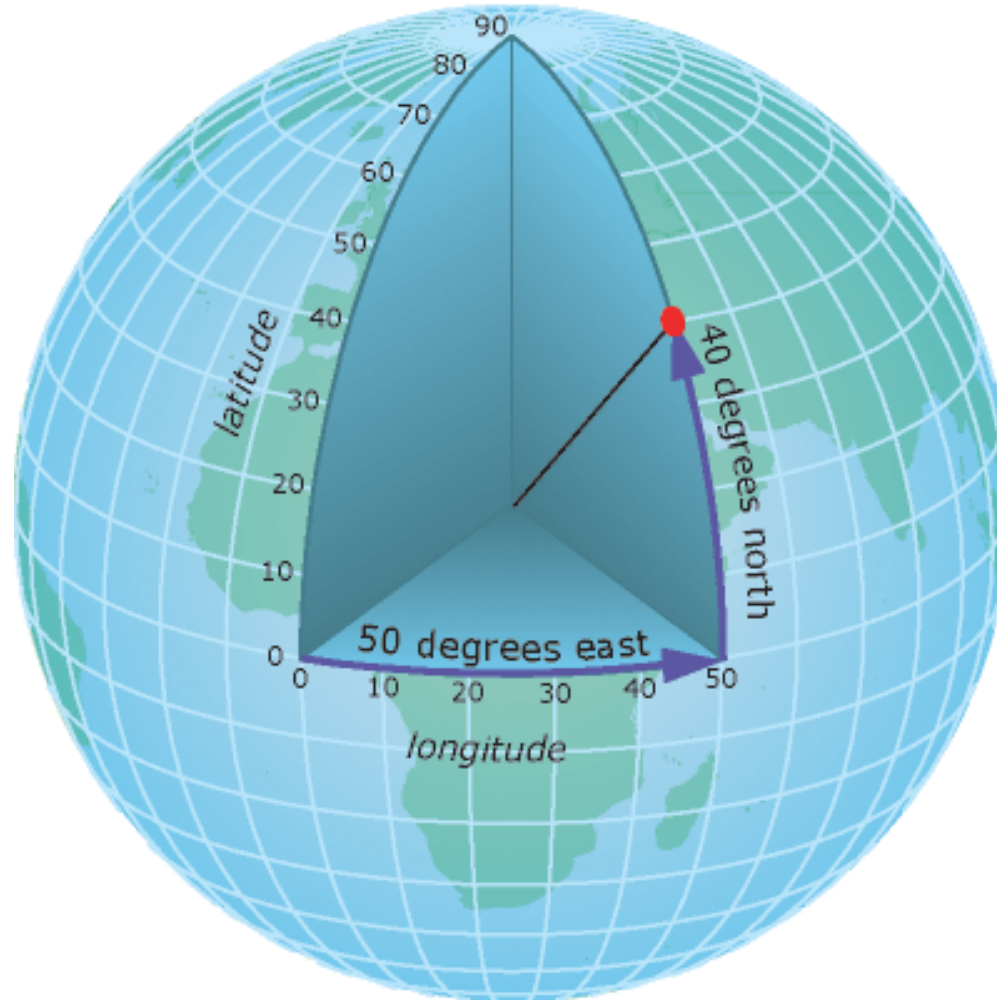IV. Current De-identification Methods

# I. Introduction

# Digital maps & Geographic Information Systems (GIS)

▶ A digital representation of the real world – a "**geobase**"

▶ Loaded with layers of additional information, static & dynamic

▶ Queried for e.g. *"what are geocoordinates of object XYZ?"* or *"at geocoordinates x,y what objects exist there?"*

▶ Visualized using a map projection

▶ Created and maintained using surveying, crowd sourcing and lots of computing & labor

# 6 dimensions of location data

- Latitude
- Longitude
- Altitude
- Time
- Frequency
- Precision

# Example: Uber ride in Berlin



Map data ©2018 GeoBasis-DE/BKG (©2009), Google

▶ Smartphone based location data

▶ Collection every 2 seconds

▶ Map matched to reduce inaccuracy

| Latitude | Longitude |
|---|---|
| 52.50333486 | 13.33955726 |
| 52.50333535 | 13.33955777 |
| 52.50333633 | 13.3395588 |
| 52.50333732 | 13.33955984 |
| 52.5033383 | 13.33956087 |
| 52.50333929 | 13.3395619 |
| 52.50333982 | 13.33956245 |
| 52.50334027 | 13.33956293 |
| 52.50334126 | 13.33956397 |
| 52.50334225 | 13.339565 |
| 52.50334323 | 13.33956603 |

# II. Sources of Data

# Sources of Location Data

▶ **"Location Services" and Platform Controls**

▶ **Hardware Sensors:**
  ▶ GNSS/GPS
  ▶ Nearby Cell Towers
  ▶ Nearby Wi-Fi Networks
  ▶ Beacons and Proximity
  ▶ Emerging Alternatives: LED, Audio

▶ **Connectivity Information,** e.g. CSLI

▶ **Wi-Fi Analytics (Tracking in Public Spaces)**

# Smartphone "Location Services"

▶ Operating system (e.g. iOS or Android) controls access to a device's geo-location

▶ Apps/websites must usually get user permission

▶ **Location Services** aggregates data from many different sources—including satellites, nearby cell towers, nearby Wi-Fi networks, and Bluetooth

# Hardware Sensors
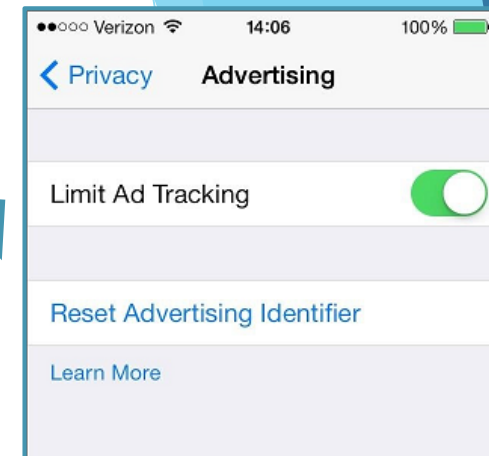


*iPhone 7*
*source: ifixit.com, Creative Commons*

**Accelerometer**
Ambient light
**Barometer**
**Bluetooth Radio**
Cameras
**Cellular Radio**
**Compass (Magnetometer)**
Face ID (iPhone)
**GPS Receiver**
**Gyroscope**
Microphones
Moisture sensor
Touch ID
**Wi-Fi Radio**

# Mobile sensors generate input for OS to provide standardized latitude-longitude

ID=NULL

Limit Ad Tracking

Reset Advertising Identifier

Learn More

*"Location Services"*

back to OS to improve services

*Accelerometer*
*Bluetooth*
*Cell Towers*
*Compass*
*GPS*
*Gyroscope*
*Wi-Fi Networks*

**38.9072˚ N, 77.0369˚ W**

**+ Identifier**
e.g. the device's mobile Ad ID:
6D92078A-8246-4BA4-AE5B-76104861E7DC

**Apps** + Third Parties using **SDKs** in Apps

# A closer look at…

*1. GPS*
*2. Cell Tower IDs*
*3. Wi-Fi Networks*
*4. Bluetooth Beacons*
*5. Alternatives – Audio and LED*

# 1. Global Navigation Satellite Systems

E.g. Global Positioning System (GPS) (U.S.)

*Allows devices to determine their location (latitude-longitude) using time signals transmitted by satellites.*

Challenges:
- Weather
- Buildings / urban environments
- Indoor positioning

# 2. Cell Tower IDs

Cell towers broadcast unique **Cell IDs**, which are compiled in both private and publicly available databases. Privately owned databases are often larger, some containing over 72 million unique cell towers.

Approximate location can be inferred by comparing detected Cell IDs and signal strengths with the known locations of cell towers.

| Cell Tower Database | Unique Cell Towers | Availability |
|---|---|---|
| OpenCellID | > 6 million | Public |
| Combain | > 72 million | Private |
| LocationAPI.org | > 72 million | Private |
| Mozilla | > 26 million | Public |
| Navizon | > 71 million | Private |
| Mylnikov GEO | > 15 million | Public |
| WiGLE | > 6 million | Private |

# 3. Nearby Wi-Fi Networks



Source: https://wigle.net/

# 3. Nearby Wi-Fi Networks



Source: https://wigle.net/

SSID: Free Hotel Wifi
BSSID: d4:62:4d:2c:c8:ec
Vendor: Ruckus Wireless

# 4. Bluetooth (Beacons)



Gimbal    Radius Networks    Estimote

Signal 360    GPShopper    Aruba

▶ Beacons are inexpensive radio transmitters that send one-way signals to devices equipped to receive them



Beacons placed in venue

Transmit one-way signal within close proximity to enabled devices

Opt-in process:

- App to interpret signal downloaded
- Bluetooth turned on
- Notifications for app enabled

Welcome to the shop

Device owner receives specific app-enabled notification

18

# 5. Proximity Alternatives – LED, Audio

Indoor Location-Based Services Using LED Lighting
**How it Works**

1. ByteLight-enabled GE LED fixtures "communicate" a unique light pattern using Visible Light Communication and Bluetooth Low Energy

2. Connected shoppers opt-in to "listen" with retailer's app on any smartphone and tablet with a camera and/or Bluetooth Smart

3. Camera detects unique light pattern and Bluetooth signal emitted by GE Lumination™ LED Luminaires; application notifies ByteLight platform of shopper's position and direction with sub-meter accuracy

4. Platform ties to retailer's digital marketing systems to deliver location-based services and personalized content to each shopper

ByteLight

shopkick

Let your device pick up our signal

Please allow Microphone access in order to get points at the store.

No | Get Signal

# Mobile Location Analytics (MLA)

*Devices with WiFi or Bluetooth capabilities broadcast their WiFi MAC Address and/or Bluetooth MAC address. Venues use MLA technology to detect how devices are moving within a space or to identify repeat visitors.*

- Looks like **68:A8:6D:E5:65:03.**

- Since different device manufacturers have been assigned groups of MAC addresses, your MAC indicates if your device is made by Apple, Samsung or another company.

- Most smartphones now randomize MAC addresses for privacy reasons.

THIS VENUE USES LOCATION INFORMATION FROM MOBILE DEVICES

For more information and your choices, visit
WWW.SMART-PLACES.ORG

FUTURE OF PRIVACY FORUM

VUB

BRUSSELS PRIVACY HUB

# III. Data Flows & Case Studies

# Location Data Ecosystem

**First Parties:**
▶ App or website that requests location
▶ Service providers (e.g. bikeshare company, mobile carrier's "cell site" location information)

**Third Parties:**
▶ Provider of an "SDK" (software development kit) integrated into an app to collect location information, e.g. for advertising or location analytics

# Location Data Ecosystem

**First Party Uses:**
- **Raw data** – *e.g. to analyze trends, user behavior, detect security threats, improve a geo-aware service*
- **Geo-fencing** – *e.g. to alert users of local promotions, events, or messages (e.g. Amber alerts)*

**Secondary Uses:**
- **Marketing profiles** across publishers or brands – e.g. coffee shop fan, frequent traveler
- **Measurement** of ad effectiveness (offline <-> online)
- **Data analysis:**
    - transportation analysis
    - city planning and Smart Cities

# Case Studies: Data Creation in Smart Communities Today

# When Selecting Data Sources, Spatial Precision is an Important Factor for Planners

# Transportation Behavior Is Changing – But Infrastructure and Budgets Have Not Kept Pace in U.S.

## Transportation Behavior is Changing

### Vehicle Ownership Trends: 2006 - 2012

**Vehicles per Person**

| | |
|---|---|
| 2006 | 0.79 |
| 2012 | 0.74 — **-6.3%** |

**Vehicles per Household**

| | |
|---|---|
| 2006 | 2.05 |
| 2012 | 1.93 — **-5.9%** |

Source: Michael Sivak, *Has motorization in the U.S. peaked?* UMTRI, January 2014

### By 2045, It Will Change Even More

- **32%** increase in urban population
- **30%** decrease in rural population
- **Up to 27%** more Vehicle Miles Traveled
- **44%** increase in trucks' freight volume

Source: US DOT, *Beyond Traffic* Final Report, January 2015

## Infrastructure Budgets Have Not Kept Up

### The Transportation Infrastructure Funding Gap: 2008 – 2028

| | | |
|---|---|---|
| Bridges | $12.8B per Year | $7.7B per Year |
| Roads | $91B per Year | $79B per Year |

0%   20%   40%   60%   80%   100%

Estimated Annual Funding   Funding Gap

Source: American Society of Civil Engineers, *2013 Report Card for America's Infrastructure*, March 2013

But according to the McKinsey Global Institute, 22% ($400B) per year could be saved globally by using data to optimize expenditures.

Source: McKinsey & Company, *Big Data vs. Congestion: Using Information to Improve Transport*, July 2015

FUTURE OF PRIVACY FORUM

VUB

BRUSSELS PRIVACY HUB

# Transportation is an Expensive, Dangerous Mystery

**A problem…**

8 Billion Hours Spent in Traffic costing over $101B in the US

**Transport is 27% US GHGs**

63,000+ structurally unsound bridges. Half of US roads under maintained. [1]

**An opportunity**

$130B/year US recommended transportation infrastructure spend.[1]

$3T/year global transport infrastructure spend expected.[2]

**22% infrastructure expenditures could be saved with data-driven techniques[2]** (and that doesn't include the externalities!)

# Example of Mobile Data – Fremont, California

**Location-Based Services Data Location**
*Circle radii vary: they accurately reflect the spatial precision of each unique data point*

**Navigation-GPS Data Location**
*Circle enlarged for visibility*

# Northern Virginia:
# Identifying and Prioritizing TDM Projects

**Transportation Demand Management**

## Scanning for Opportunities

**Need:** Evaluate and prioritize solutions to traffic when highway expansion is not an option due to widespread residential and commercial development

**Question to Answer:** Where are the highest volume of short trips between O-D pairs that could be converted to other modes?

**Challenge:** Northern Virginia had to scan hundreds of miles of roads to identify and prioritize the best TDM opportunities, which was not possible to do cost-effectively with conventional data sources

| TAZ ID | Avg Trip Duration (sec) | Avg Trip Speed (mph) | Sum under 1 mile | Sum under 3 mile |
|--------|-------------------------|----------------------|------------------|------------------|
| 851 | 1186 | 27 | 5% | 30% |
| 850 | 1433 | 27 | 6% | 25% |
| 849 | 1427 | 30 | 4% | 21% |
| 848 | 916 | 23 | 5% | 47% |
| 847 | 1420 | 27 | 9% | 39% |
| 846 | 1275 | 29 | 4% | 28% |
| 845 | 1180 | 23 | 6% | 38% |
| 844 | 1129 | 26 | 7% | 37% |
| 843 | 1504 | 27 | 5% | 25% |
| 842 | 1485 | 30 | 4% | 27% |
| 841 | 1460 | 26 | 7% | 31% |
| 840 | 1403 | 26 | 3% | 24% |
| 839 | 1177 | 25 | 4% | 37% |
| 838 | 1359 | 26 | 6% | 34% |
| 837 | 1272 | 28 | 3% | 30% |
| 836 | 1397 | 28 | 8% | 45% |
| 835 | 1732 | 33 | 6% | 36% |

# City of Lafayette, California:
# Pinpointing the Cause of Congestion Downtown

**Downtown Congestion Study**

## O-D for Select Link

**Need:** Evaluate and prioritize solutions to congestion in downtown corridor

**Question to Answer:** Understand what which type of trip causes congestion: School drop-offs, commuters to downtown, or "first/last mile" commuters to transit stop

**Challenge:** Studies were not providing satisfactory answers. The city had counts, but they didn't show origins and destinations, and surveys were inconclusive.



E-6 Northbridge

E-5 Walnut Creek

7%

18%

E-1 Orinda

39%

13%

E-2 MDB

9%

E-3 Moraga

E-4 St. Marys

Lafayette Reservoir

St. Mary's Ballfield

→ External Screenlines
→ Through Trip Movement
Lafayette Zone

# Charlotte, North Carolina:
# Calibrating a Travel Demand Model

**Hypothetical** Transport. Demand Modeling

### Origin-Destination for North Carolina MPO

**Need:** Accurate O/D for calibration or transportation demand model without expensive/time consuming survey for personal and medium/heavy duty commercial trips.

**Question:** How do travel patterns vary by demographic group and time of day?

**Challenge:** Planners need to understand how all groups travel, but MPO survey respondents were disproportionately higher income, making it difficult to determine the impact of plans on lower income travelers.

# IV. De-Identification: Current Methods

two or more objects can *not* be
at the *same* place at the *same* time

# "Identity" and "identification" according to Wikipedia

▶ Identity (philosophy), also called sameness, is whatever makes an entity definable and recognizable

▶ Identity (social science), individuality, personal identity, social identity, and cultural identity in psychology, sociology, and philosophy

▶ Identity (mathematics), an equality that holds regardless of the values of its variables

▶ Identification (information), the capability to find, retrieve, report, change, or delete specific data without ambiguity

# De-identifying location data adding ambiguity

Various methods exist, such as:

- **Replacing** identifiers with pseudo-identifiers, eg through hashing or lookup tables

- **Stripping identifiers**: (numeric) values that are relatable to individuals

- **Removing sections of data** that combined with other data could allow for identification e.g. begin/end of trip

- **Adding inaccuracy** in time and/or space

- **Aggregating** into "buckets" of time and space

1 month


1 day

Can location data be anonymous?

Yes. But it is very hard to achieve.

Taking an ongoing risk based approach is key.

Technical, organization and contractual measures can provide a "tripod" of assurance.

# Questions?