

# Introduction

M2 SIF - SED

Tristan Allard  
Univ. Rennes 1 / Irisa lab.  
`tristan.allard@irisa.fr`

Autumn 2023

# Organisation of the class

- ▶ **Teachings** : 7 CMs (lessons), 7 TPs (practical works).
- ▶ **Teachers** : Tristan Allard (main teacher, resp. class), Gaëtan Le Guelvout (head of the Trust & Security lab at the IRT BCOM).
- ▶ **Evaluation** : presentation (), final exam () .

# Themes of the class

## Privacy: focus on privacy-preserving data publishing techniques

1. Introduction  
*Privacy is not dead (especially anonymization).*
2. Lessons from the past: partition-based approaches  
*Not completely past, not completely learnt neither. . .*
3. Modern approaches: differential privacy  
*Today's "gold standard"?*
4. Attacks on privacy-preserving data publishing schemes  
*From re-identification attacks to membership inference attacks*

## Intellectual Property: securing multimedia content delivery

1. Digital watermarking  
*Some tattoo sheeps, others tattoo data.*
2. Secure distribution  
*"Keep an eye on me."*

# Privacy

*Focus: privacy-preserving data publishing*

# Progress of the Talk

Issues with Personal Data

Privacy: a Vague Notion

Privacy-Preserving Data Publishing: Roots and High-Level View

References

# Massive Data Collection

Which dimensions about us generate data today ?

- ▶ Geolocations (GPS-enabled devices, cellphone access points, online maps trajectories, public transport cards, ...)
- ▶ Electrical consumption (smart meters, smart plugs, ...)
- ▶ Health and physiological status (social security, wearables)
- ▶ Citizen's rights and duties (IDs, taxes, ...)<sup>1</sup>
- ▶ Social networks and messaging applications (friends and family, work, ...)
- ▶ Online activities (search queries, browsing history, ...)
- ▶ ...
- ▶ Students' connected beds at the CROUS<sup>2</sup> ?

Easier to answer to what does **not** generate data about us today ?

---

<sup>1</sup>See the Digital ID Wallet product: [https://youtu.be/Y5b0nLRte\\_A](https://youtu.be/Y5b0nLRte_A).

<sup>2</sup><https://www.ouest-france.fr/bretagne/rennes-35000/etudiants-rennes-dormez-vous-etes-surveilles-5227294/amp>

# The New Oil

“Personal data is the new oil of the internet and the new currency of the digital world.”

M. Kouneva, European Commissioner for Consumer Protection,  
March 2009



# Targeted Ads I

For ex : on social networks.



The image shows a Facebook advertisement banner. The top part is a dark blue navigation bar with the Facebook logo on the left, an 'Inscription' button, and login fields for 'Adresse électronique ou téléphone' and 'Mot de passe' with a 'Connexion' button. Below the navigation bar, the main banner features a background of overlapping fingerprint patterns in light blue and orange. On the left, there is a target icon with a red arrow hitting the bullseye. The main text reads 'Publicités Facebook' in large black font, followed by 'Plus d'un milliard de clients potentiels.' and 'Nous sommes là pour vous aider.' On the right side of the banner, there is a green 'Créer une publicité' button, followed by the text 'Bénéficiez une première fois de notre assistance technique gratuitement' and 'Appelez le +33 174180267 ou demandez que l'on vous rappelle'.

facebook [Inscription](#)

Adresse électronique ou téléphone  Mot de passe  [Connexion](#)

Garder ma session active [Mot de passe oublié ?](#)

 **Publicités Facebook**

Plus d'un milliard de clients potentiels.  
Nous sommes là pour vous aider.

[Créer une publicité](#)

Bénéficiez une première fois de notre assistance technique gratuitement

**Appelez le +33 174180267**  
ou demandez que l'on vous rappelle



# Targeted Ads II

From [2] (CACM '21).

The image shows a Facebook advertisement for misterbnb. The ad features a photo of two men in a city, with the text "misterbnb designed gay hotspots in the world" and "Stay Like a Gay Local". Below the photo, it says "Live like a gay local - feel welcome anywhere you go in over 130 countries. Check out our top cities New York, Paris, Barcelona, and Rome." and includes a "Book Now" button. To the right of the ad is a "Your ad preferences" sidebar with sections for "Your interests" and "Choose an interest to preview examples of ads you might see on Facebook or re".

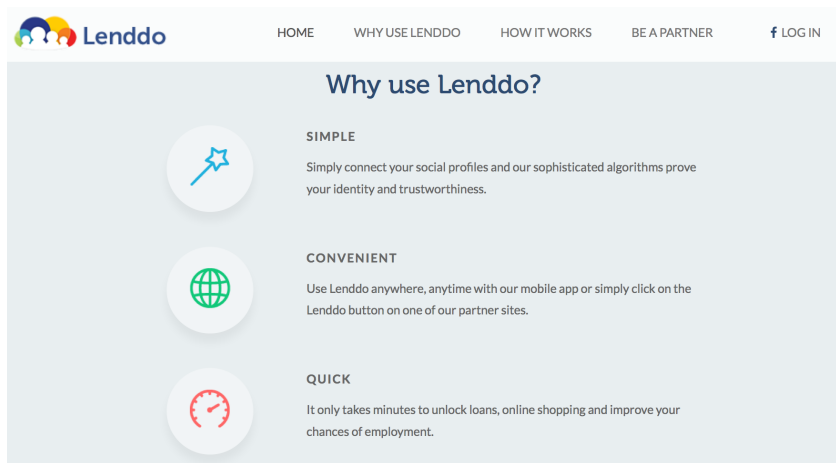
Figure: An example of a sensitive ad related to homosexuality

Code	Country	Homosexuality
AF	AFGHANISTAN	12.31
MR	MAURITANIA	0.99
QA	QATAR	2.35
SO	SOMALIA	1.44
PK	PAKISTAN	1.54
AE	UNITED ARAB EMIRATES	3.00
BN	BRUNEI	5.24
NG	NIGERIA	2.35
SA	SAUDI ARABIA	2.08
YE	YEMEN	1.08
IQ	IRAQ	3.20

Figure: Percentage of FB users (FFB) tagged with the interest "Homosexuality" in countries where being homosexual may lead to death penalty/

# Risk Optimization

For ex : credit scoring.



The screenshot shows the Lenddo website's navigation bar with the logo and links for HOME, WHY USE LENDDO, HOW IT WORKS, BE A PARTNER, and LOGIN. The main content area features a heading 'Why use Lenddo?' followed by three key benefits: SIMPLE, CONVENIENT, and QUICK, each with an icon and a brief description.

**Lenddo** HOME WHY USE LENDDO HOW IT WORKS BE A PARTNER **f** LOG IN

## Why use Lenddo?

**SIMPLE**  
Simply connect your social profiles and our sophisticated algorithms prove your identity and trustworthiness.

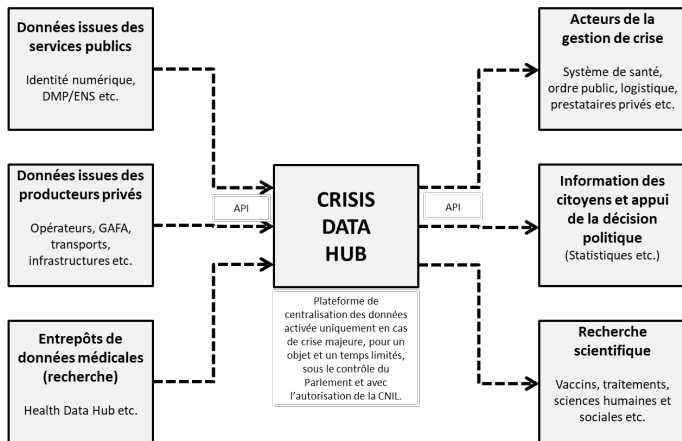
**CONVENIENT**  
Use Lenddo anywhere, anytime with our mobile app or simply click on the Lenddo button on one of our partner sites.

**QUICK**  
It only takes minutes to unlock loans, online shopping and improve your chances of employment.

<https://lenddo.com/>

# Crisis Mitigation?

For ex : the “crisis data hub”<sup>3</sup>.



<sup>3</sup>French Senate report num. 673 (3rd June 2021)

*Etc.*

(point/mass surveillance, impersonation, influencing voters, *etc.*)  
(health science, public transportation planning, energy transition,  
*etc.*)

# Progress of the Talk

Issues with Personal Data

Privacy: a Vague Notion

Privacy-Preserving Data Publishing: Roots and High-Level View

References

# Privacy and Computer Science I



**General goal:** Solve the eternal tradeoff between security/privacy and utility!

- ▶ Many meanings of “privacy”
- ▶ Many privacy-preserving techniques

# Privacy and Computer Science II

## Privacy-Preserving Techniques

- ▶ Any security technique applied to personal data can help (e.g., encryption schemes, access control mechanisms).
- ▶ Specificities:
  - ▶ **Specific security properties** required by the context (e.g., related to vote privacy)
  - ▶ **Specific tolerance to non-negligible leaks** because data must be **both used and protected** (e.g., disclosure of useful statistics by “anonymized” data).

# Privacy-Preserving Data Publishing

Privacy-Preserving Data Publishing (PPDP) :

- ▶ Publish *personal data* for analysis purposes (accurate aggregate queries) . . .
- ▶ . . . while preserving individuals' *privacy* (uncertain point queries)
- ▶ Also called *sanitization* (please do) and sometimes *anonymization* (please don't)



# Privacy-Preserving Data Publishing and the Law

Scope of the EU law :

- ▶ GDPR (General Data Protection Regulation) : the new European regulation about the protection of personal data.
- ▶ Protects personal data.
- ▶ It does not apply to personal data “made anonymous” .

⇒ Crucial to define “personal data” and “made anonymous” !

# An Encompassive Definition of Personal Data

**GDPR Article 4<sup>4</sup> (1)** : “*‘personal data’ means*

- ▶ **any information relating to an identified or identifiable natural person** (*‘data subject’*);
- ▶ *an identifiable natural person is* **one who can be identified, directly or indirectly,**
- ▶ *in particular by reference to an identifier such as*
  - ▶ **a name,**
  - ▶ **an identification number,**
  - ▶ **location data,**
  - ▶ **an online identifier**
  - ▶ *or to one or more factors specific to the* **physical, physiological, genetic, mental, economic, cultural or social identity** *of that natural person;”*

---

<sup>4</sup><https://gdpr-info.eu/art-4-gdpr/>

# But a Fuzzy Definition of Privacy-Preserving Data Publishing Techniques

## GDPR Recital 26 <sup>5</sup> :

- ▶ **“The principles of data protection should therefore not apply to anonymous information,**
- ▶ *namely information which does not relate to an identified or identifiable natural person*
- ▶ *or to **personal data rendered anonymous in such a manner that the data subject is not or no longer identifiable.***
- ▶ *This Regulation does not therefore concern the processing of such anonymous information, including for statistical or research purposes.”*

⇒ Crucial to design strong privacy-preserving data publishing techniques !

---

<sup>5</sup><https://gdpr-info.eu/recitals/no-26/>

# Progress of the Talk

Issues with Personal Data

Privacy: a Vague Notion

Privacy-Preserving Data Publishing: Roots and High-Level View

References

# Pseudonymization and Scandals

Historically, pseudonymization was considered as a valid anonymization method. But:

GIC 2002 [4] : Health data of the GIC X list of voters from the state of Massachusetts  $\Rightarrow$  Health folder of M. Weld (governor)

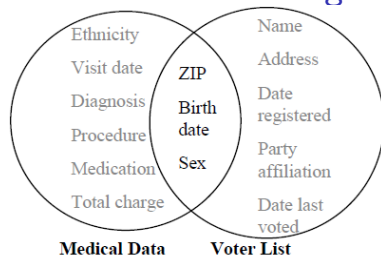
AOL 2006 [1] : search queries X ground investigation  $\Rightarrow$  Interests of Mrs T. Arnold (USA citizen)

Netflix 2006 [3] : private movie ratings on Netflix X a few public movie ratings on IMDB together with user names  $\Rightarrow$  Private movie ratings of two users (with high probability)

NYC 2013 : taxi trips X “dé-pseudonymisation” table or blogs of *stars*.

etc ...

# Pseudonymized Data: A Join is Enough



**Figure:** Gov. Weld's Case : Medical JOIN Voter ON (zip, DoB, sex)

## A straightforward disclosure

- ▶ Governor Weld's case lived in Cambridge and was part of the GIC dataset.
- ▶ In the voter list: 6 individuals had his birthdate, 3 of them were men, only one had Weld's zipcode.
- ▶ (zip, DoB, sex) : a *quasi-identifier*.

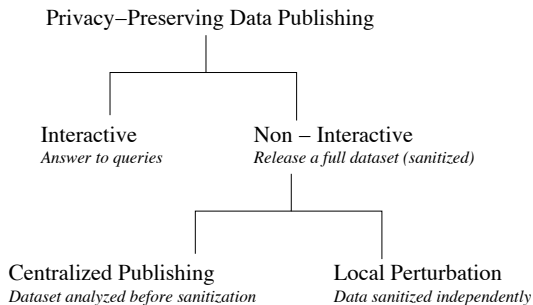
# Components of a Privacy-preserving Data Publishing Solution

Three components:

1. **Privacy model:** What does it mean for the data released to be privacy-preserving?
2. **Privacy mechanism:** How to produce the privacy-preserving data to be released?
3. **Utility metric:** How much useful is the released data?

Pseudonymity does not work. . . Which component(s) does it miss ?

# Variations on a Theme



- ▶ Privacy mechanism: in one of the following families: interactive, centralized publishing, local perturbation.
- ▶ Privacy model: essentially agnostic to these families.



**Let's go?**

# Progress of the Talk

Issues with Personal Data

Privacy: a Vague Notion

Privacy-Preserving Data Publishing: Roots and High-Level View

References

- [1] M. Arrington.  
AOL Proudly Releases Massive Amounts of Private Data.  
TechCrunch, 6th of August 2006.
- [2] J. G. Cabañas, A. Cuevas, A. Arrate, and R. Cuevas.  
Does facebook use sensitive data for advertising purposes?  
*Commun. ACM*, 64(1):62–69, Dec. 2020.
- [3] A. Narayanan and V. Shmatikov.  
Robust de-anonymization of large sparse datasets.  
In *2008 IEEE Symposium on Security and Privacy (S&P 2008)*,  
18-21 May 2008, Oakland, California, USA, pages 111–125,  
2008.
- [4] L. Sweeney.  
k-anonymity: a model for protecting privacy.  
*Int. J. Uncertain. Fuzziness Knowl.-Based Syst.*,  
10(5):557–570, 2002.