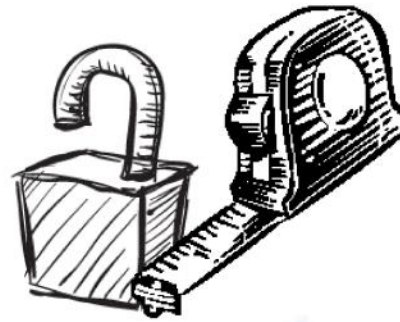


# ENGINEERING PRIVACY BY DESIGN



Carmela Troncoso

institute  
**iM**idea  
software

# PRIVACY BY DESIGN – LET'S HAVE IT!

INFORMATION AND PRIVACY COMMISSIONER OF ONTARIO



Privacy by Design

## Privacy by Design principles

1. Proactive not Reactive; Preventative not Remedial
2. Privacy as the Default Setting
3. **Privacy Embedded into Design**
4. Full Functionality: Positive-Sum, not Zero-Sum
5. End-to-End Security — Full Lifecycle Protection
6. Visibility and Transparency — Keep it Open
7. Respect for User Privacy — Keep it User-Centric

Cavoukian et al. (2010)

<https://www.ipc.on.ca/images/resources/7foundationalprinciples.pdf>

# PRIVACY BY DESIGN – LET'S HAVE IT!

INFORMATION AND PRIVACY COMMISSIONER OF ONTARIO



Privacy by Design

## Privacy by Design principles

1. Proactive not Reactive; Preventative not Remedial
2. Privacy as the Default Setting
3. **Privacy Embedded into Design**
4. Full Functionality: Positive-Sum, not Zero-Sum
5. End-to-End Security — Full Lifecycle Protection
6. Visibility and Transparency — Keep it Open
7. Respect for User Privacy — Keep it User-Centric

Cavoukian et al. (2010)

ARTICLE 25 EUROPEAN GENERAL DATA PROTECTION REGULATION



*“the controller shall [...] implement appropriate technical and organisational measures [...] which are designed to implement data-protection principles[...] in order to meet the requirements of this Regulation and protect the rights of data subjects.”*

<https://www.ipc.on.ca/images/resources/7foundationalprinciples.pdf>

<http://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32016R0679&from=EN>

# PRIVACY BY DESIGN – LET'S HAVE IT!

INFORMATION AND PRIVACY COMMISSIONER OF ONTARIO



Privacy by Design

## Privacy by Design principles

1. Proactive not Reactive; Preventative not Remedial
2. Privacy as the Default Setting
3. **Privacy Embedded into Design**
4. Full Functionality: Positive-Sum, not Zero-Sum
5. End-to-End Security — Full Lifecycle Protection
6. Visibility and Transparency — Keep it Open
7. Respect for User Privacy — Keep it User-Centric

Cavoukian et al. (2010)

ARTICLE 25 EUROPEAN GENERAL DATA PROTECTION REGULATION



*“the controller shall [...] implement appropriate technical and organisational measures [...] which are designed to implement data-protection principles[...] in order to meet the requirements of this Regulation and protect the rights of data subjects.”*



Actually... “Data Protection by design and by default”

<https://www.ipc.on.ca/images/resources/7foundationalprinciples.pdf>

<http://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32016R0679&from=EN>

# PRIVACY BY DESIGN – LET'S HAVE IT!

INFORMATION AND PRIVACY COMMISSIONER OF ONTARIO

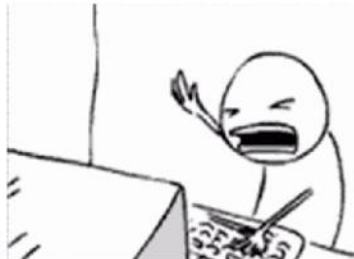


Privacy by Design

## Privacy by Design principles

1. Proactive not Reactive; Preventative not Remedial
2. Privacy as the Default Setting
3. **Privacy Embedded into Design**
4. Full Functionality: Positive-Sum, not Zero-Sum
5. End-to-End Security — Full Lifecycle Protection
6. Visibility and Transparency — Keep it Open
7. Respect for User Privacy — Keep it User-Centric

Cavoukian et al. (2010)



ARTICLE 25 EUROPEAN GENERAL DATA PROTECTION REGULATION



*“the controller shall [...] implement appropriate technical and organisational measures [...] which are designed to implement data-protection principles[...] in order to meet the requirements of this Regulation and protect the rights of data subjects.”*



Actually... “Data Protection by design and by default”

## BUT HOW ??????????????

<https://www.ipc.on.ca/images/resources/7foundationalprinciples.pdf>

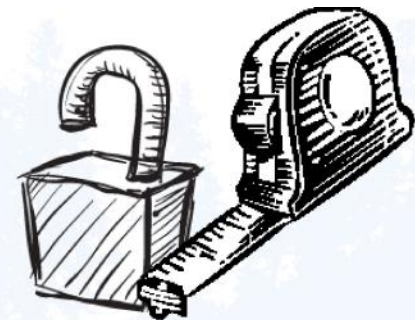
<http://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32016R0679&from=EN>

PART I:  
REASONING ABOUT PRIVACY WHEN  
DESIGNING SYSTEMS



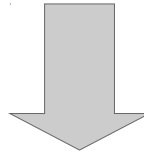
# THIS TALK: ENGINEERING PRIVACY BY DESIGN

PART II:  
EVALUATING PRIVACY IN PRIVACY-  
PRESERVING SYSTEMS



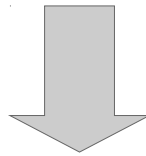
PART I:  
REASONING ABOUT PRIVACY WHEN  
DESIGNING SYSTEMS





# PRIVACY





# PRIVACY



WHY?? NOT ONLY MOTIVATION....  
IS PRIVACY ENGINEERING A CRAFT?



# ENGINEERING PRIVACY BY DESIGN 1.0

Two case studies:

- anonymous e-petitions: no identity attached to petitions
- privacy-preserving road tolling: no fine grained data sent to server

# ENGINEERING PRIVACY BY DESIGN 1.0

Two case studies:

- anonymous e-petitions: no identity attached to petitions
- privacy-preserving road tolling: no fine grained data sent to server

**THE KEY IS "DATA MINIMIZATION"**

# ENGINEERING PRIVACY BY DESIGN 1.0

Two case studies:

- anonymous e-petitions: no identity attached to petitions
- privacy-preserving road tolling: no fine grained data sent to server

**THE KEY IS “DATA MINIMIZATION”**

**BUT**, it’s not “data” that is minimized (in the system as a *whole*)

- kept in user devices
- sent encrypted to a server (only client has the key)
- distributed over multiple servers: only the user, or colluding servers, can recover the data

# ENGINEERING PRIVACY BY DESIGN 1.0

Two case studies:

- anonymous e-petitions: no identity attached to petitions
- privacy-preserving road tolling: no fine grained data sent to server

**THE KEY IS “DATA MINIMIZATION”**

**BUT**, it’s not “data” that is minimized (in the system as a *whole*)

- kept in user devices
- sent encrypted to a server (only client has the key)
- distributed over multiple servers: only the user, or colluding servers, can recover the data

**“DATA MINIMIZATION” IS A BAD METAPHOR!!!**

# UNPACKING “DATA MINIMIZATION”: PRIVACY BY DESIGN STRATEGIES

**OVERARCHING  
GOAL**

**MINIMIZING PRIVACY RISKS AND  
TRUST ASSUMPTIONS PLACED ON OTHER ENTITIES**

# UNPACKING "DATA MINIMIZATION": PRIVACY BY DESIGN STRATEGIES

OVERARCHING  
GOAL

MINIMIZING PRIVACY RISKS AND  
TRUST ASSUMPTIONS PLACED ON OTHER ENTITIES



THE ADVERSARY



Seda Gurses, Carmela Troncoso, Claudia Diaz. Engineering Privacy by Design Reloaded. Amsterdam Privacy Conference. 2015  
Seda Gurses and Claudia Diaz. "Two tales of privacy in online social networks." IEEE Security & Privacy Magazine. 2013

# UNPACKING "DATA MINIMIZATION": PRIVACY BY DESIGN STRATEGIES

OVERARCHING  
GOAL

MINIMIZING PRIVACY RISKS AND  
TRUST ASSUMPTIONS PLACED ON OTHER ENTITIES

Social  
Privacy



Other users  
3<sup>rd</sup> parties



THE ADVERSARY





# UNPACKING "DATA MINIMIZATION": PRIVACY BY DESIGN STRATEGIES

OVERARCHING  
GOAL

MINIMIZING PRIVACY RISKS AND  
TRUST ASSUMPTIONS PLACED ON OTHER ENTITIES

Social  
Privacy

Institutional  
Privacy  
(data protection)

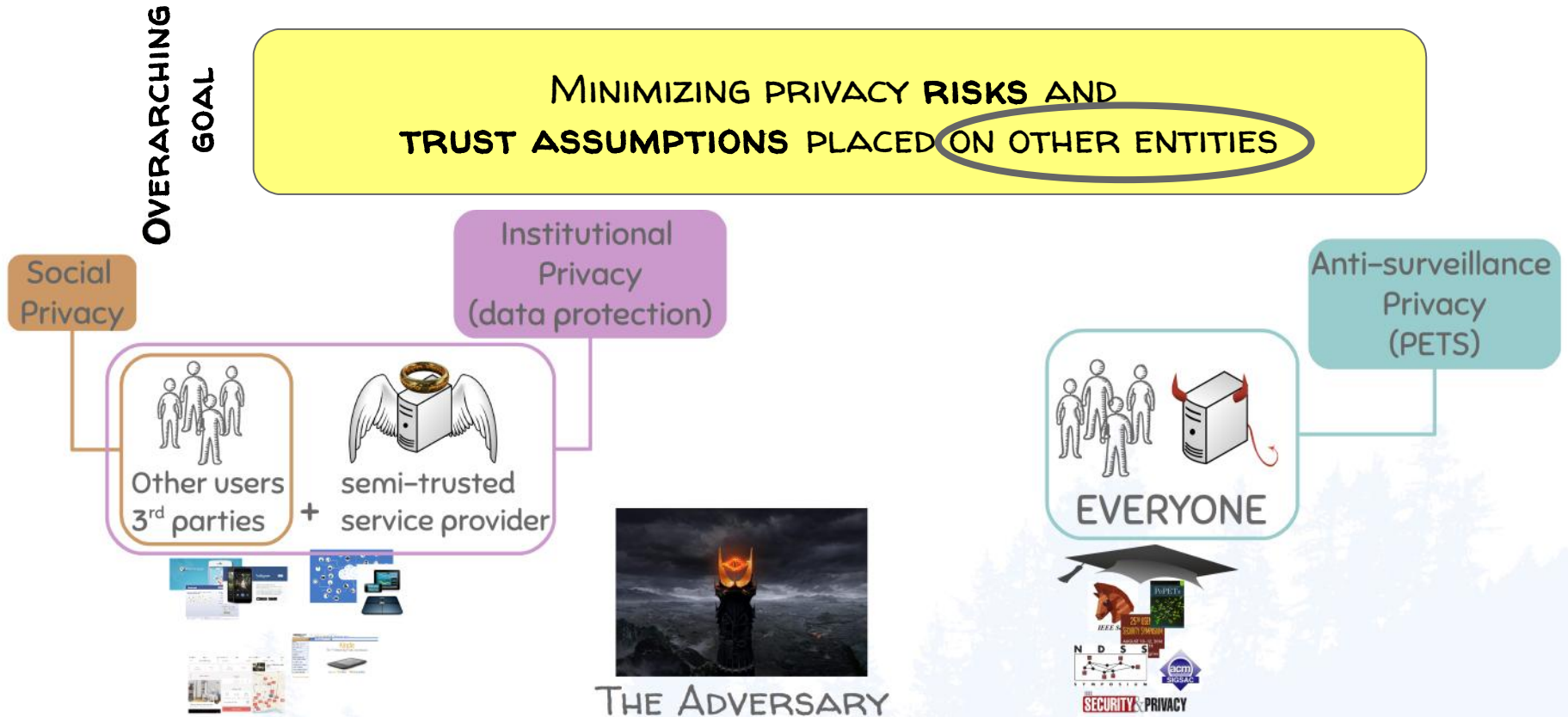
Other users  
3<sup>rd</sup> parties + semi-trusted  
service provider



THE ADVERSARY



# UNPACKING "DATA MINIMIZATION": PRIVACY BY DESIGN STRATEGIES



# UNPACKING “DATA MINIMIZATION”: PRIVACY BY DESIGN STRATEGIES

OVERARCHING  
GOAL

MINIMIZING PRIVACY RISKS AND  
TRUST ASSUMPTIONS PLACED ON OTHER ENTITIES

STRATEGIES

MINIMIZE  
COLLECTION

MINIMIZE  
DISCLOSURE

MINIMIZE  
LINKABILITY

MINIMIZE  
CENTRALIZATION

MINIMIZE  
REPLICATION

MINIMIZE  
RETENTION

# UNPACKING “DATA MINIMIZATION”: PRIVACY BY DESIGN STRATEGIES

OVERARCHING  
GOAL

MINIMIZING PRIVACY RISKS AND  
TRUST ASSUMPTIONS PLACED ON OTHER ENTITIES

STRATEGIES

MINIMIZE  
COLLECTION

MINIMIZE  
DISCLOSURE

MINIMIZE  
LINKABILITY

MINIMIZE  
CENTRALIZATION

MINIMIZE  
REPLICATION

MINIMIZE  
RETENTION

GREAT! BUT... HOW DO WE USE THESE STRATEGIES?

We make explicit the activities and reasoning in **PRIVACY ENGINEERING DESIGN** process

# CASE STUDY: ELECTRONIC TOLL PRICING

MOTIVATION: EUROPEAN ELECTRONIC TOLL SERVICE (EETS)

Toll collection on European Roads through On Board Equipment

Two approaches: Satellite Technology / DSRC

# CASE STUDY: ELECTRONIC TOLL PRICING

## MOTIVATION: EUROPEAN ELECTRONIC TOLL SERVICE (EETS)

Toll collection on European Roads through On Board Equipment

Two approaches: Satellite Technology / DSRC

## STARTING ASSUMPTIONS

- 1) Well defined functionality  
Charge depending on driving
  
- 2) Security, privacy & service integrity requirements  
Users location should be private  
No cheating clients
  
- 3) Initial reference system

# CASE STUDY: ELECTRONIC TOLL PRICING

## MOTIVATION: EUROPEAN ELECTRONIC TOLL SERVICE (EETS)

Toll collection on European Roads through On Board Equipment

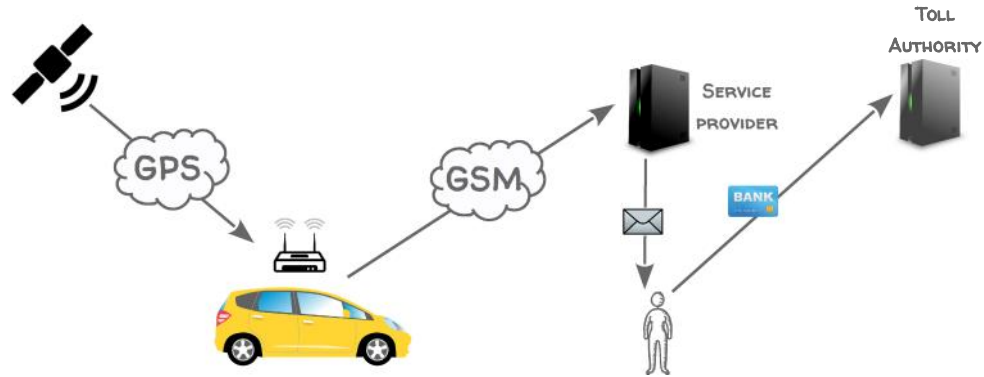
Two approaches: Satellite Technology / DSRC

## STARTING ASSUMPTIONS

- 1) Well defined functionality  
Charge depending on driving
- 2) Security, privacy & service integrity requirements  
Users location should be private  
No cheating clients
- 3) Initial reference system



# CASE STUDY: ELECTRONIC TOLL PRICING



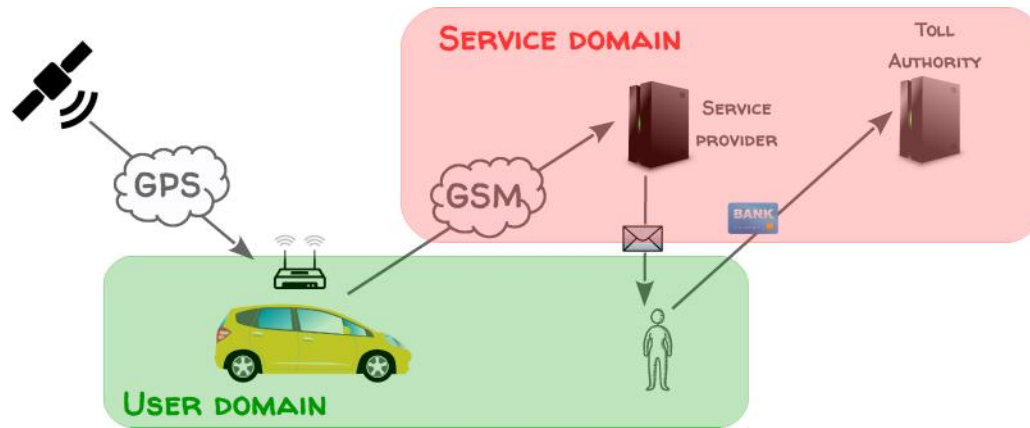
## ACTIVITY 1: CLASSIFY ENTITIES IN DOMAINS

**USER DOMAIN:** components under the control of the user, eg, user devices

**SERVICE DOMAIN:** components outside the control of the user, eg, backend system at provider



# CASE STUDY: ELECTRONIC TOLL PRICING

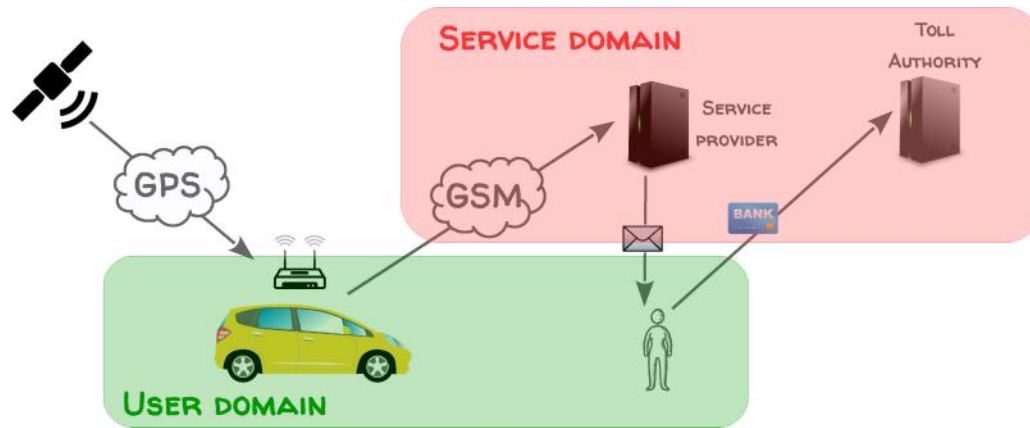


## ACTIVITY 1: CLASSIFY ENTITIES IN DOMAINS

**USER DOMAIN:** components under the control of the user, eg, user devices

**SERVICE DOMAIN:** components outside the control of the user, eg, backend system at provider

# CASE STUDY: ELECTRONIC TOLL PRICING



## ACTIVITY 1: CLASSIFY ENTITIES IN DOMAINS

**USER DOMAIN:** components under the control of the user, eg, user devices

**SERVICE DOMAIN:** components outside the control of the user, eg, backend system at provider

## ACTIVITY 2: IDENTIFY NECESSARY DATA FOR PROVIDING THE SERVICE

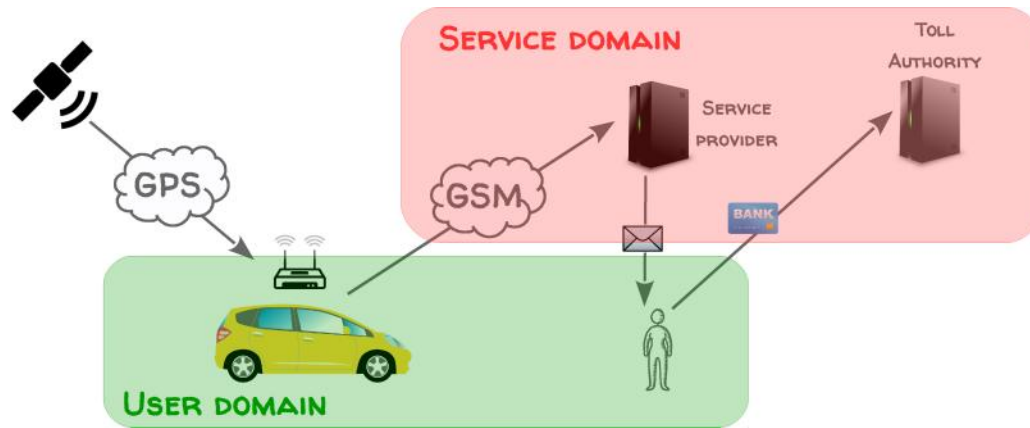
Location data – compute bill

Billing data – charge user

Personal data – send bill

Payment data – perform payment

# CASE STUDY: ELECTRONIC TOLL PRICING



## ACTIVITY 1: CLASSIFY ENTITIES IN DOMAINS

**USER DOMAIN:** components under the control of the user, eg, user devices

**SERVICE DOMAIN:** components outside the control of the user, eg, backend system at provider

## ACTIVITY 2: IDENTIFY NECESSARY DATA FOR PROVIDING THE SERVICE

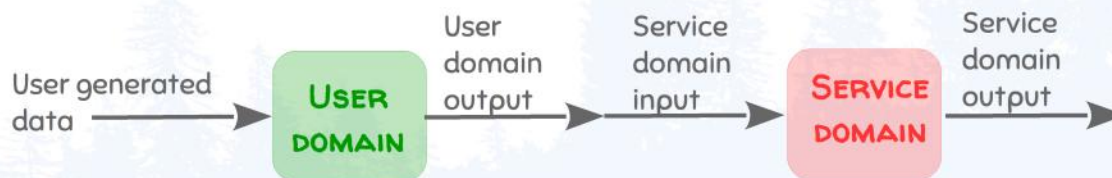
Location data – compute bill

Billing data – charge user

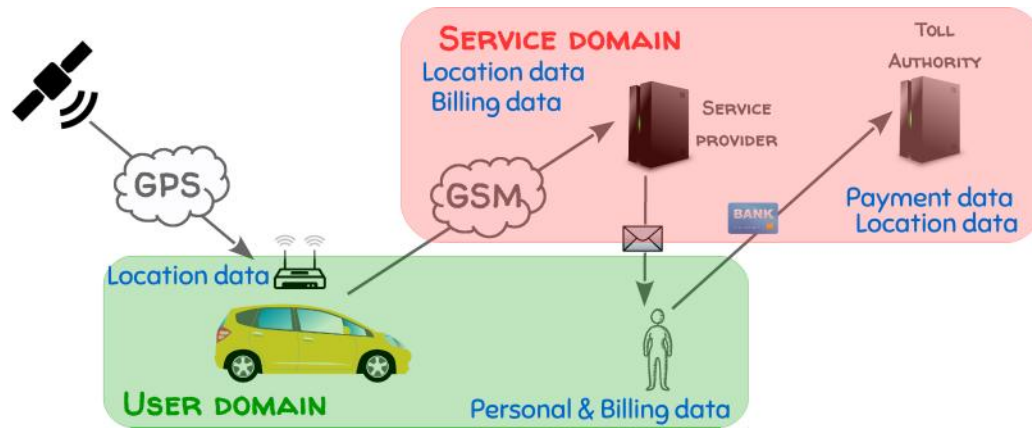
Personal data – send bill

Payment data – perform payment

## ACTIVITY 3: DISTRIBUTE DATA IN ARCHITECTURE



# CASE STUDY: ELECTRONIC TOLL PRICING



## ACTIVITY 1: CLASSIFY ENTITIES IN DOMAINS

**USER DOMAIN:** components under the control of the user, eg, user devices

**SERVICE DOMAIN:** components outside the control of the user, eg, backend system at provider

## ACTIVITY 2: IDENTIFY NECESSARY DATA FOR PROVIDING THE SERVICE

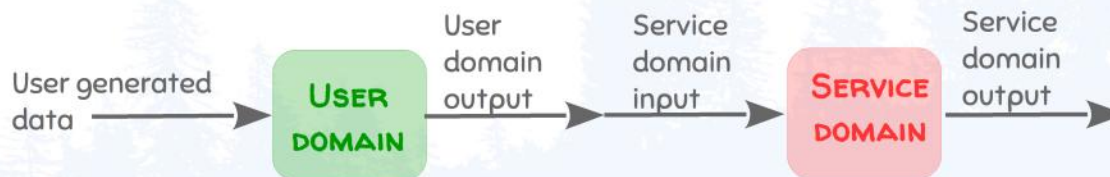
Location data – compute bill

Billing data – charge user

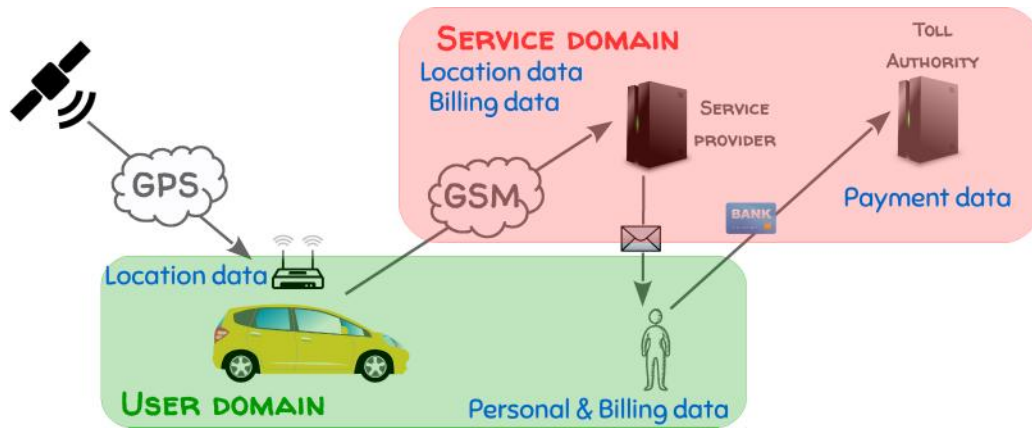
Personal data – send bill

Payment data – perform payment

## ACTIVITY 3: DISTRIBUTE DATA IN ARCHITECTURE



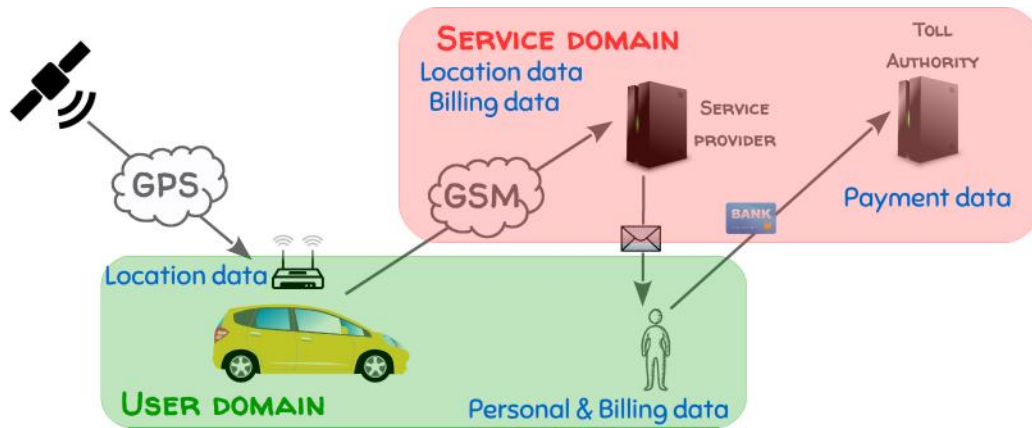
# CASE STUDY: ELECTRONIC TOLL PRICING



ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →



# CASE STUDY: ELECTRONIC TOLL PRICING

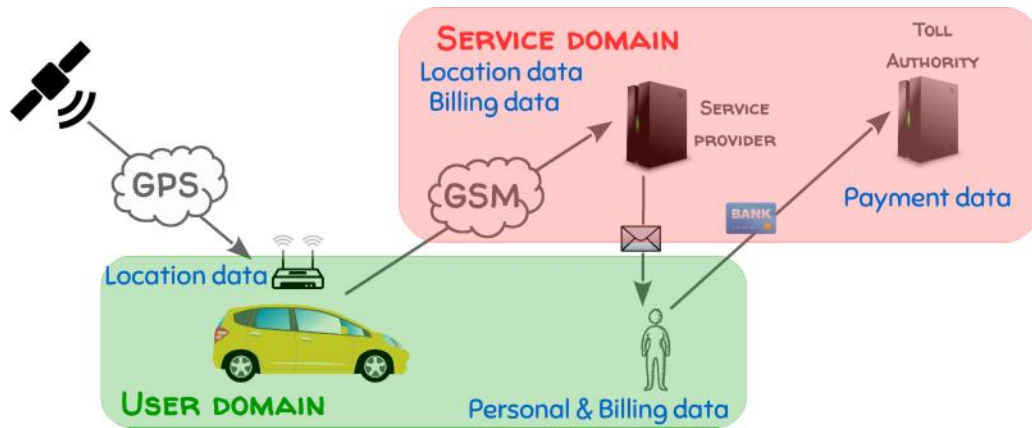


## ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →

not sending the data (local computations)  
encrypting the data  
advanced privacy-preserving protocols  
obfuscate the data  
anonymize the data



# CASE STUDY: ELECTRONIC TOLL PRICING



Trust Service to keep privacy of location data

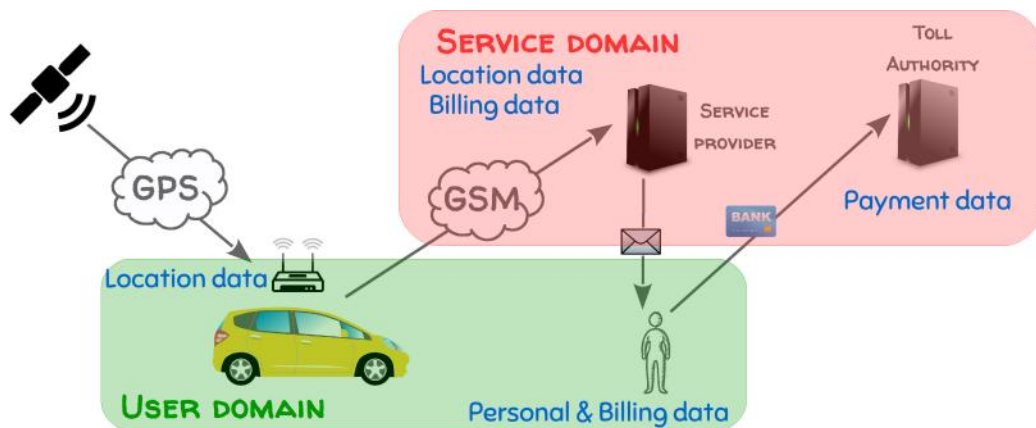
Risk of privacy breach

## ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →

not sending the data (local computations)  
encrypting the data  
advanced privacy-preserving protocols  
obfuscate the data  
anonymize the data



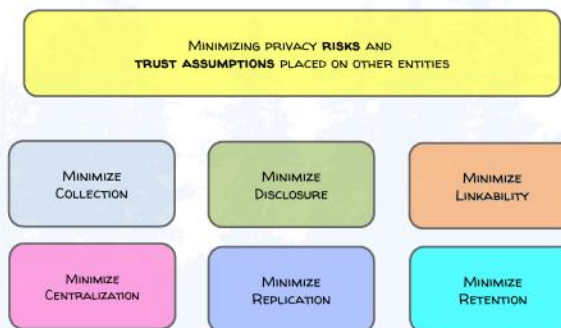
# CASE STUDY: ELECTRONIC TOLL PRICING



Location is not needed,  
only the amount to bill!

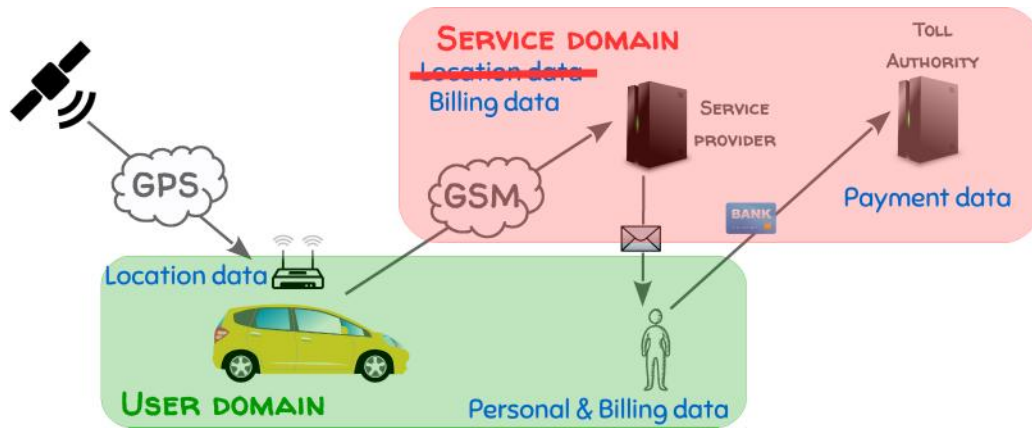
## ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →

not sending the data (local computations)  
encrypting the data  
advanced privacy-preserving protocols  
obfuscate the data  
anonymize the data





# CASE STUDY: ELECTRONIC TOLL PRICING



Location is not needed,  
only the amount to bill!

## ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →

not sending the data (local computations)

encrypting the data

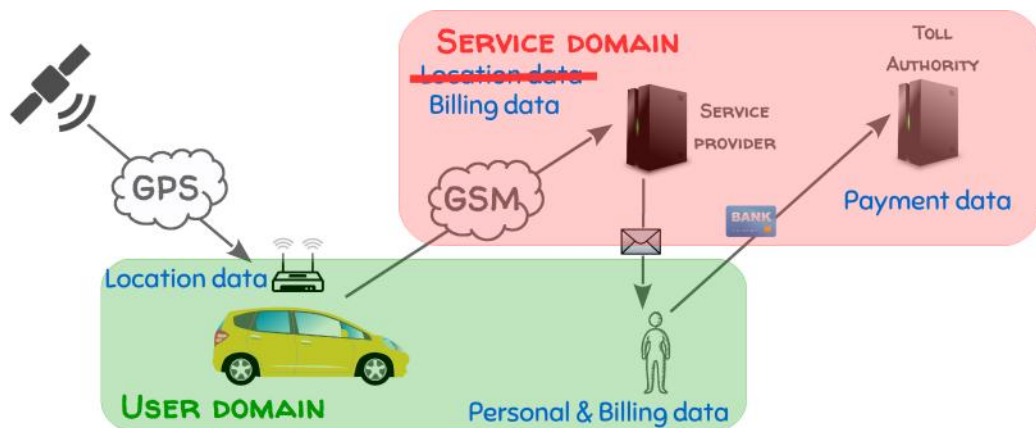
advanced privacy-preserving protocols

obfuscate the data

anonymize the data



# CASE STUDY: ELECTRONIC TOLL PRICING



Location is not needed,  
only the amount to bill!

Service integrity?

## ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →

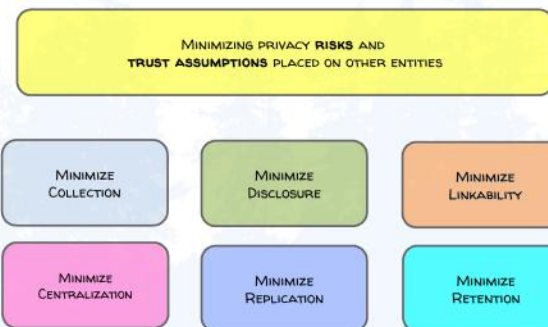
not sending the data (local computations)

encrypting the data

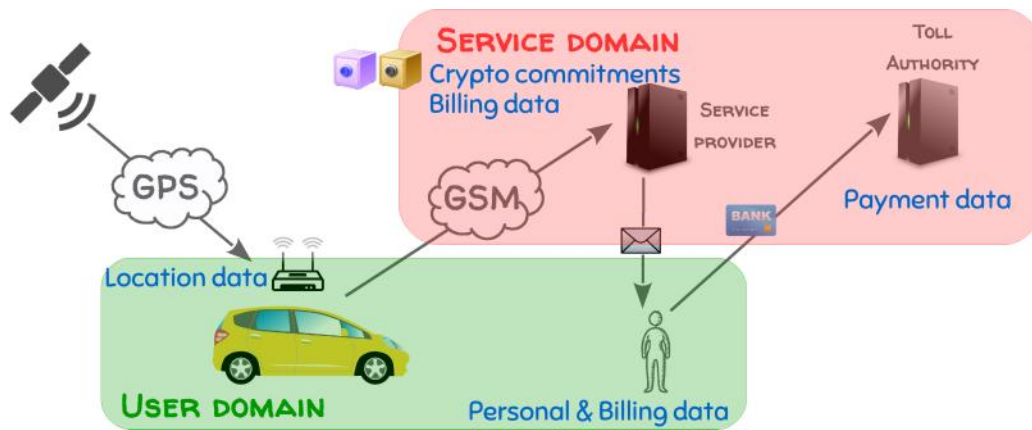
advanced privacy-preserving protocols

obfuscate the data

anonymize the data



# CASE STUDY: ELECTRONIC TOLL PRICING



Location is not needed,  
only the amount to bill!

Service integrity?

## ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →

not sending the data (local computations)

encrypting the data

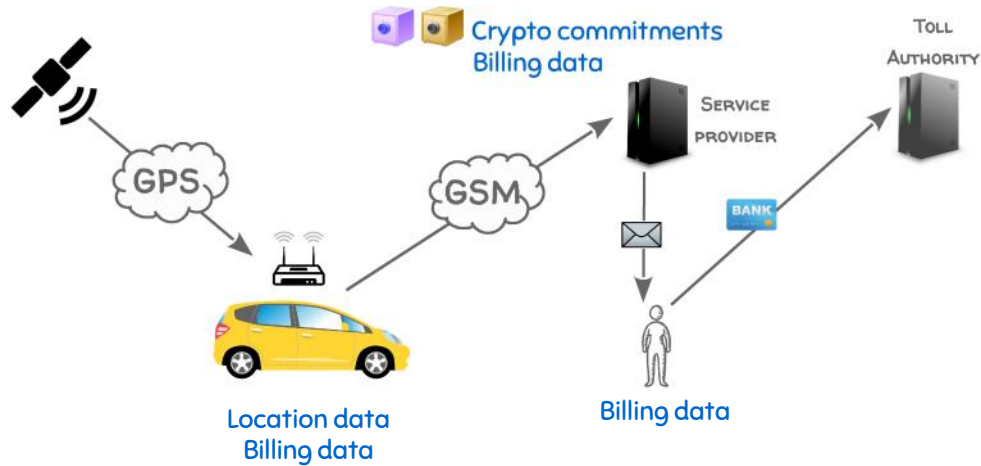
advanced privacy-preserving protocols

obfuscate the data

anonymize the data



# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING

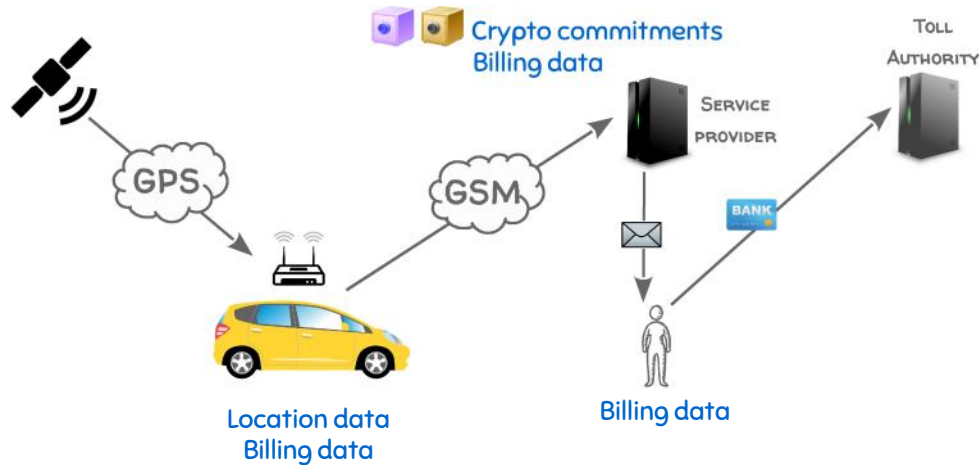


## Crypto Commitments to:

locations  

prices  

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



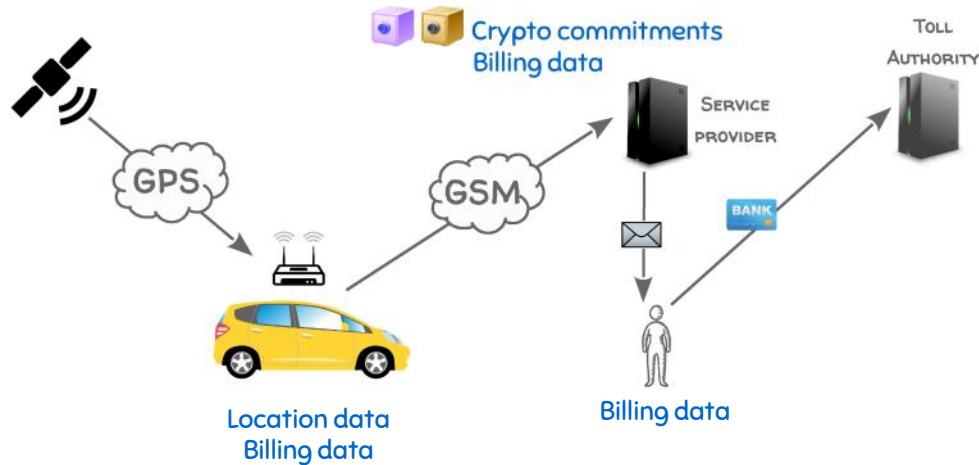
## Crypto Commitments to:

locations  

prices  

**Hiding – Given the commitment, no information about the locations/price can be gained**

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



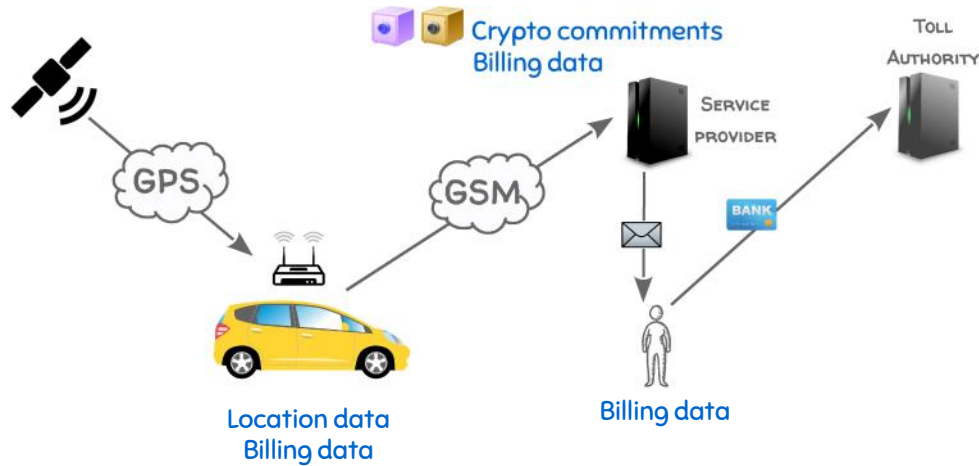
## Crypto Commitments to:

locations 

prices 

**Binding – Once committed to a set of locations/prices cannot be changed**

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING

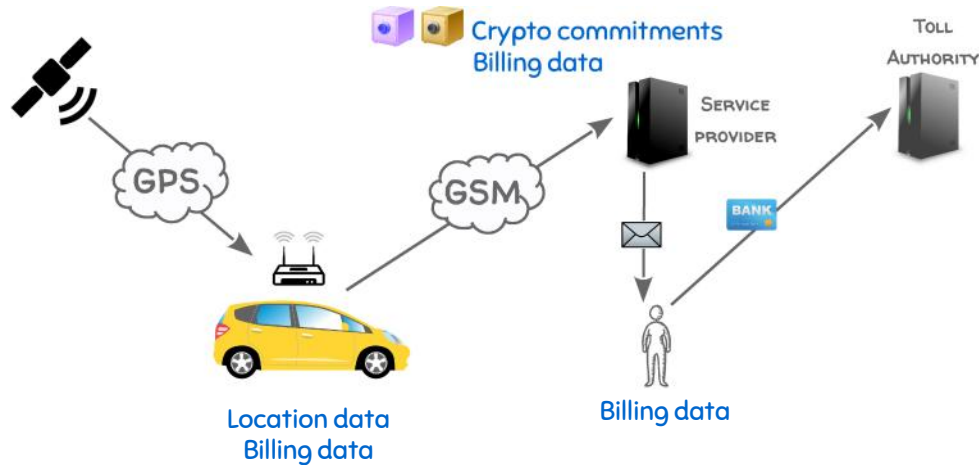


Crypto Commitments to: + ZK proofs that prices come from a correct policy

locations  A  B 

prices  

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



Homomorphic Commitments to: + ZK proofs that prices come from a correct policy



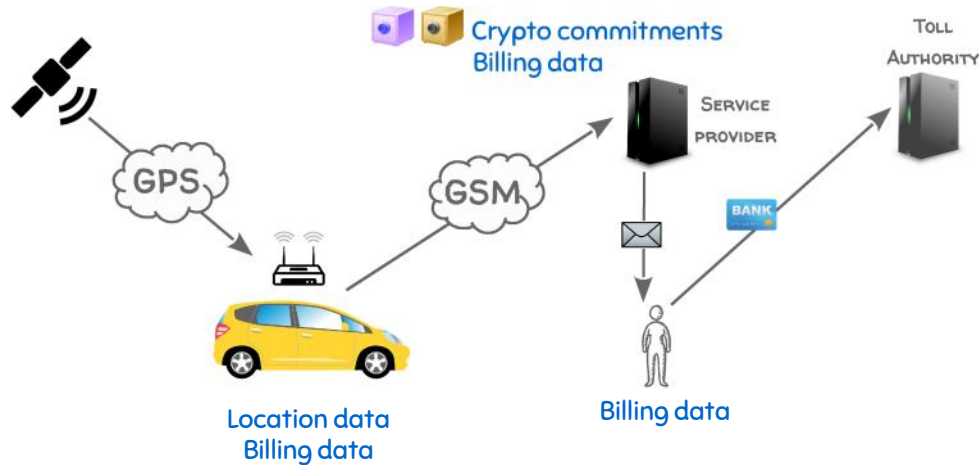
C. Troncoso, G. Danezis, E. Kosta, J. Balasch, B. Preneel. PriPAYD. Privacy-Friendly Pay-As-You-Drive Insurance. IEEE TDSC 2011

C. Troncoso, G. Danezis, E. Kosta, B. Preneel. PriPAYD. privacy friendly pay-as-you-drive insurance. WPES 2007

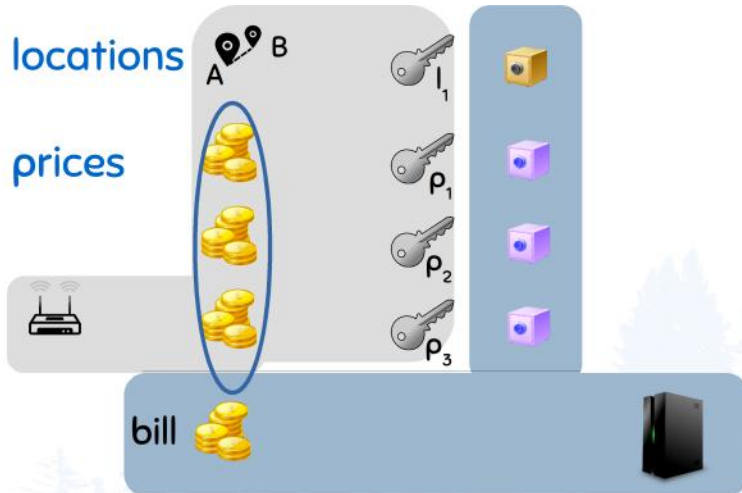
J. Balasch, A. Rial, C. Troncoso, B. Preneel, I. Verbauwhede, C. Geuens. PrETP. Privacy-Preserving Electronic Toll Pricing. USENIX Security 2010



# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



Homomorphic Commitments to: + ZK proofs that prices come from a correct policy

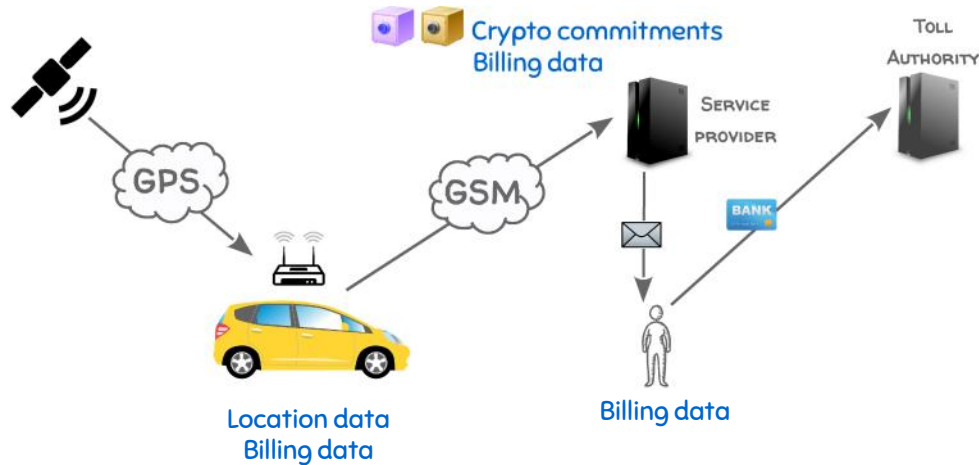


C. Troncoso, G. Danezis, E. Kosta, J. Balasch, B. Preneel. PriPAYD. Privacy-Friendly Pay-As-You-Drive Insurance. IEEE TDSC 2011

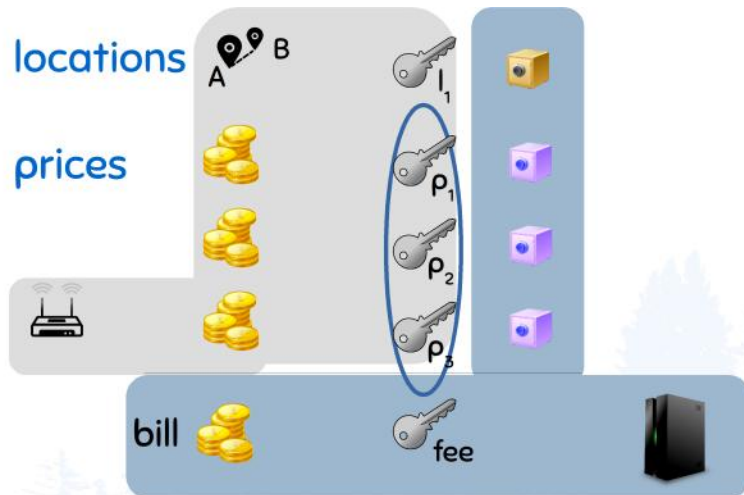
C. Troncoso, G. Danezis, E. Kosta, B. Preneel. PriPAYD. privacy friendly pay-as-you-drive insurance. WPES 2007

J. Balasch, A. Rial, C. Troncoso, B. Preneel, I. Verbauwhede, C. Geuens. PrETP. Privacy-Preserving Electronic Toll Pricing. USENIX Security 2010

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



Homomorphic Commitments to: + ZK proofs that prices come from a correct policy

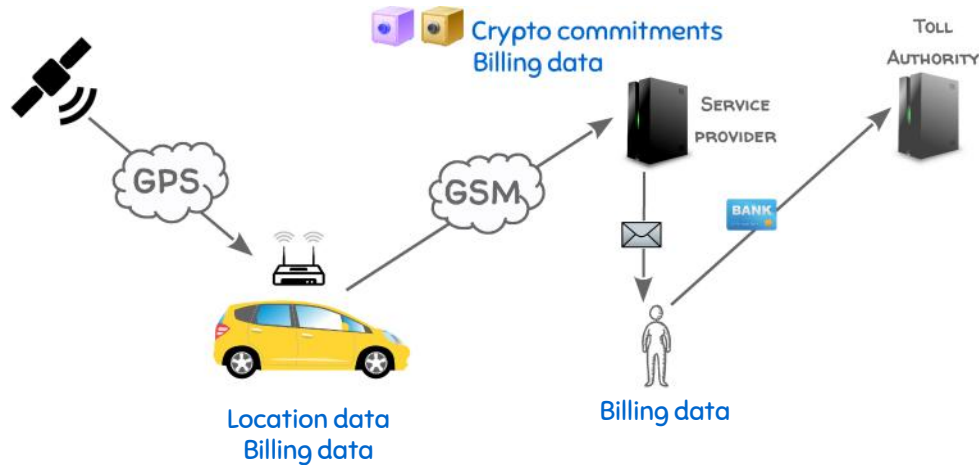


C. Troncoso, G. Danezis, E. Kosta, J. Balasch, B. Preneel. PriPAYD. Privacy-Friendly Pay-As-You-Drive Insurance. IEEE TDSC 2011

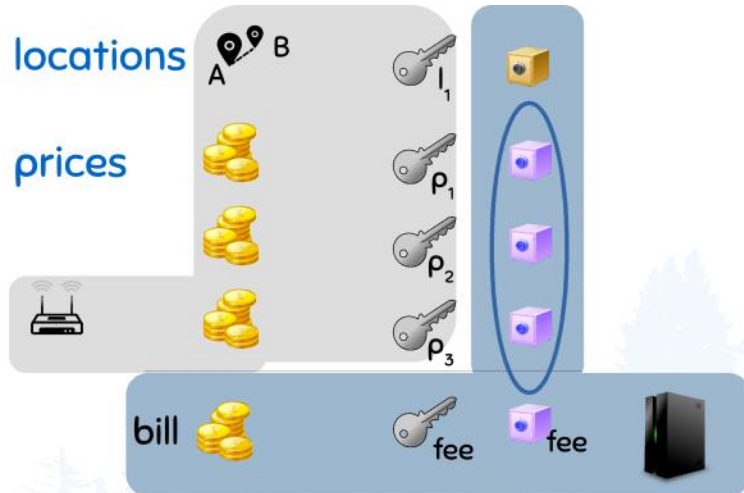
C. Troncoso, G. Danezis, E. Kosta, B. Preneel. PriPAYD. privacy friendly pay-as-you-drive insurance. WPES 2007

J. Balasch, A. Rial, C. Troncoso, B. Preneel, I. Verbauwhede, C. Geuens. PrETP. Privacy-Preserving Electronic Toll Pricing. USENIX Security 2010

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



Homomorphic Commitments to: + ZK proofs that prices come from a correct policy

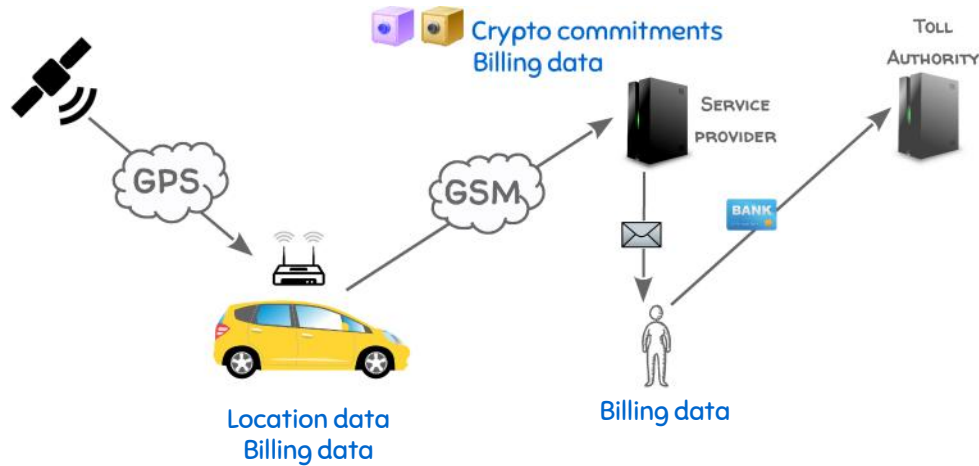


C. Troncoso, G. Danezis, E. Kosta, J. Balasch, B. Preneel. PriPAYD. Privacy-Friendly Pay-As-You-Drive Insurance. IEEE TDSC 2011

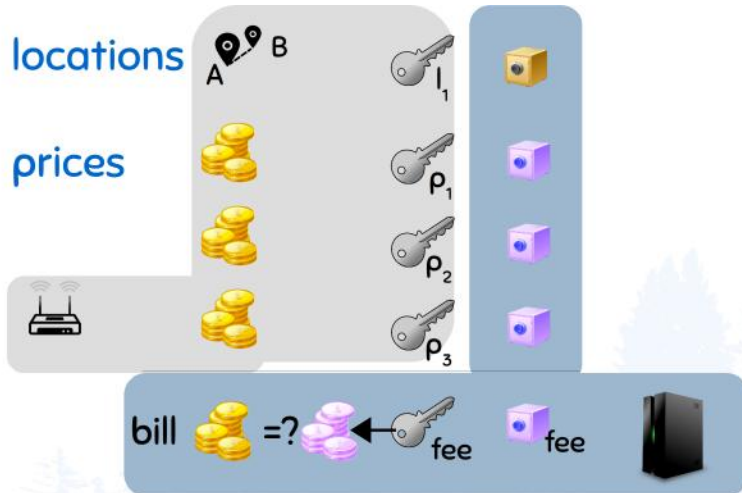
C. Troncoso, G. Danezis, E. Kosta, B. Preneel. PriPAYD. privacy friendly pay-as-you-drive insurance. WPES 2007

J. Balasch, A. Rial, C. Troncoso, B. Preneel, I. Verbauwhede, C. Geuens. PrETP. Privacy-Preserving Electronic Toll Pricing. USENIX Security 2010

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



Homomorphic Commitments to: + ZK proofs that prices come from a correct policy

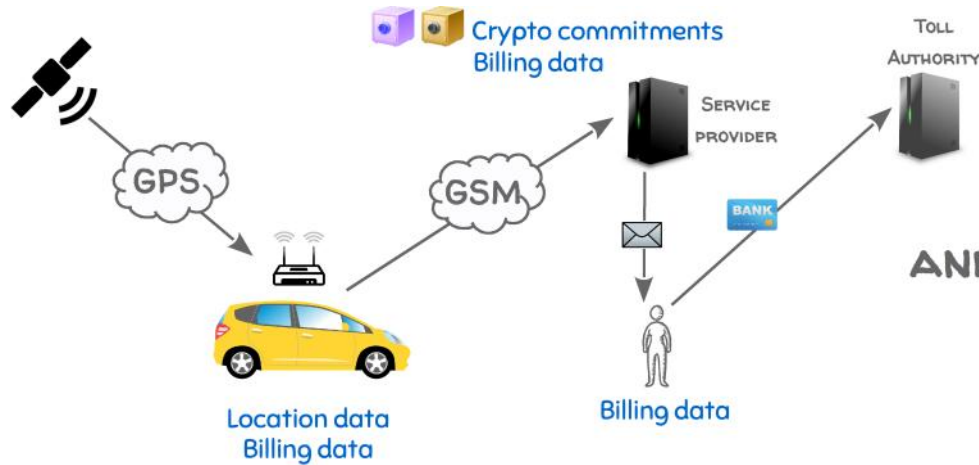


C. Troncoso, G. Danezis, E. Kosta, J. Balasch, B. Preneel. PriPAYD. Privacy-Friendly Pay-As-You-Drive Insurance. IEEE TDSC 2011

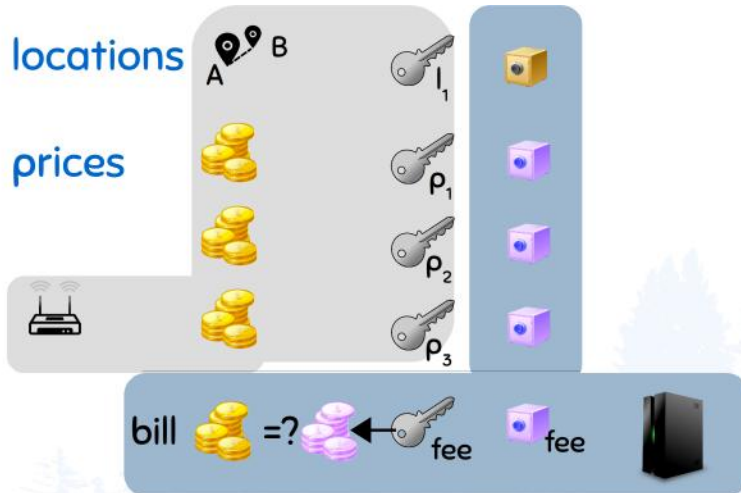
C. Troncoso, G. Danezis, E. Kosta, B. Preneel. PriPAYD. privacy friendly pay-as-you-drive insurance. WPES 2007

J. Balasch, A. Rial, C. Troncoso, B. Preneel, I. Verbauwhede, C. Geuens. PrETP. Privacy-Preserving Electronic Toll Pricing. USENIX Security 2010

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



Homomorphic Commitments to: + ZK proofs that prices come from a correct policy

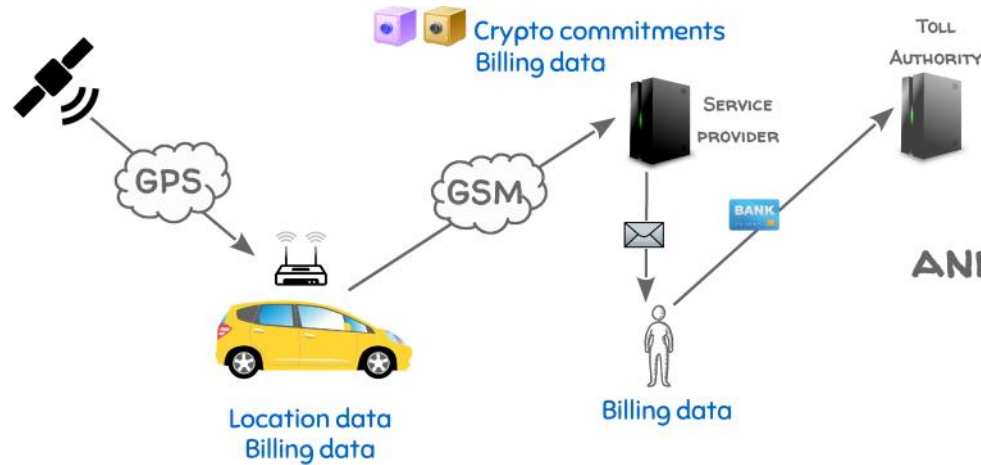


C. Troncoso, G. Danezis, E. Kosta, J. Balasch, B. Preneel. PriPAYD. Privacy-Friendly Pay-As-You-Drive Insurance. IEEE TDSC 2011

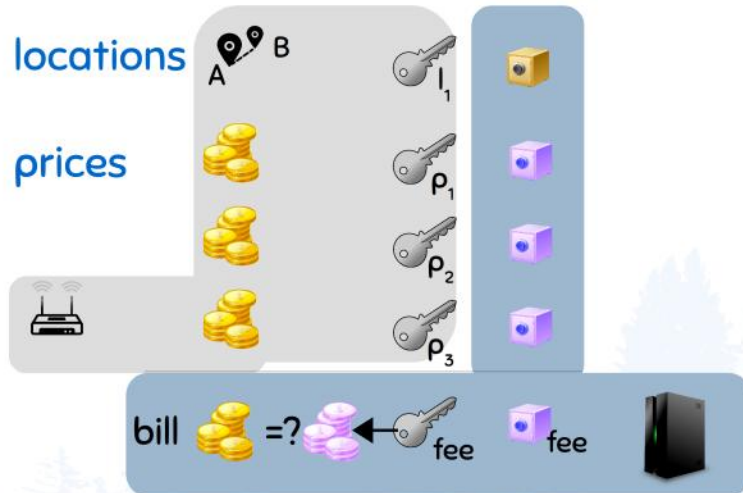
C. Troncoso, G. Danezis, E. Kosta, B. Preneel. PriPAYD. privacy friendly pay-as-you-drive insurance. WPES 2007

J. Balasch, A. Rial, C. Troncoso, B. Preneel, I. Verbauwhede, C. Geuens. PrETP. Privacy-Preserving Electronic Toll Pricing. USENIX Security 2010

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



Homomorphic Commitments to: + ZK proofs that prices come from a correct policy



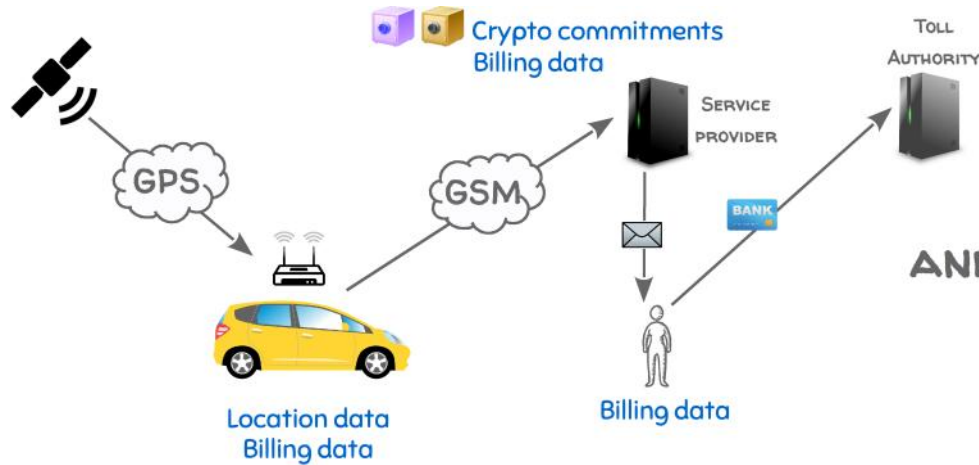
Random checks of location/price

C. Troncoso, G. Danezis, E. Kosta, J. Balasch, B. Preneel. PriPAYD. Privacy-Friendly Pay-As-You-Drive Insurance. IEEE TDSC 2011

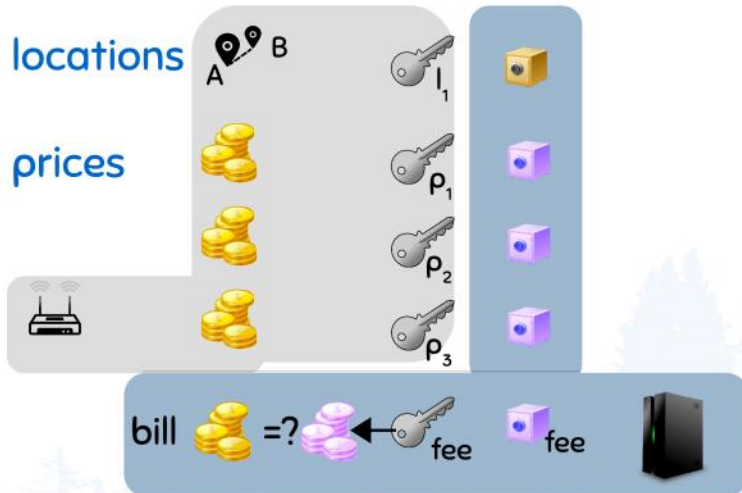
C. Troncoso, G. Danezis, E. Kosta, B. Preneel. PriPAYD. privacy friendly pay-as-you-drive insurance. WPES 2007

J. Balasch, A. Rial, C. Troncoso, B. Preneel, I. Verbauwhede, C. Geuens. PrETP. Privacy-Preserving Electronic Toll Pricing. USENIX Security 2010

# PRIVACY-PRESERVING ELECTRONIC TOLL PRICING



Homomorphic Commitments to: + ZK proofs that prices come from a correct policy



Random checks of location/price

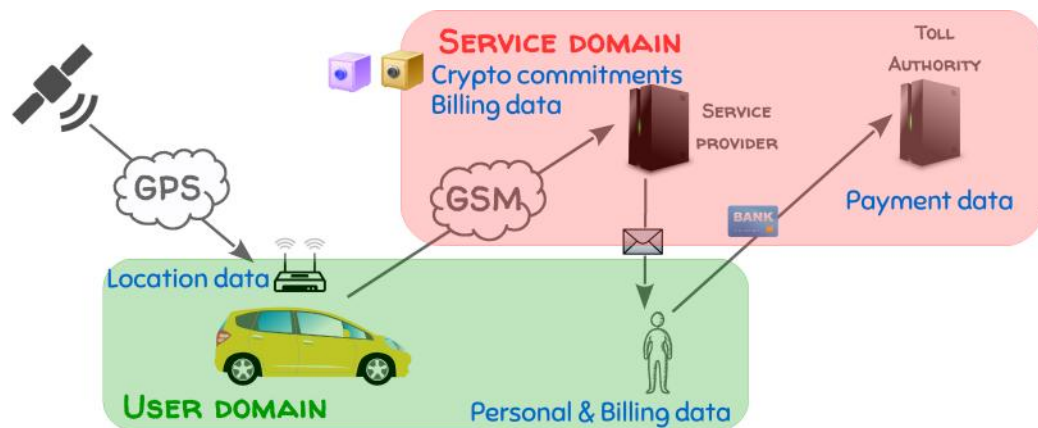
Complex Pricing Policies: Smart Metering ('11)

C. Troncoso, G. Danezis, E. Kosta, J. Balasch, B. Preneel. PriPAYD. Privacy-Friendly Pay-As-You-Drive Insurance. IEEE TDSC 2011

C. Troncoso, G. Danezis, E. Kosta, B. Preneel. PriPAYD. privacy friendly pay-as-you-drive insurance. WPES 2007

J. Balasch, A. Rial, C. Troncoso, B. Preneel, I. Verbauwhede, C. Geuens. PrETP. Privacy-Preserving Electronic Toll Pricing. USENIX Security 2010

# CASE STUDY: ELECTRONIC TOLL PRICING



Location is not needed,  
only the amount to bill!

Service integrity

## ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →

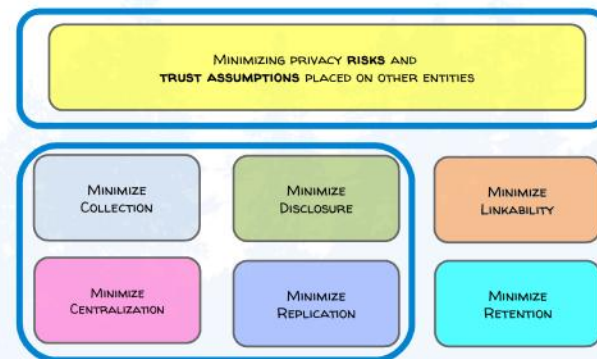
not sending the data (local computations)

encrypting the data

advanced privacy-preserving protocols

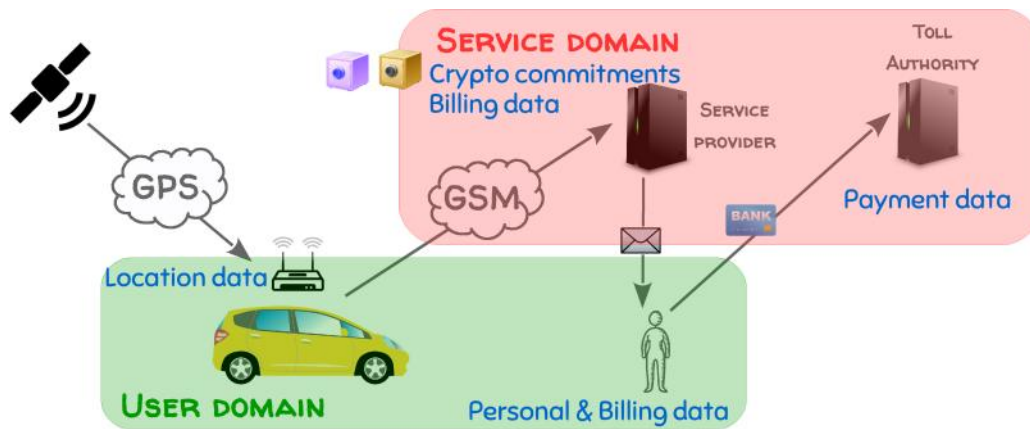
obfuscate the data

anonymize the data





# CASE STUDY: ELECTRONIC TOLL PRICING



Location is not needed,  
only the amount to bill!

Service integrity?

Requires deep knowledge of PETs  
Privacy ENABLING Technologies

## ACTIVITY 4: SELECT TECHNOLOGICAL SOLUTIONS FOLLOWING →

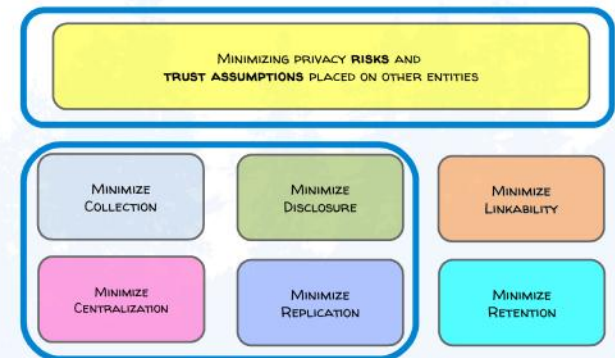
not sending the data (local computations)

encrypting the data

advanced privacy-preserving protocols

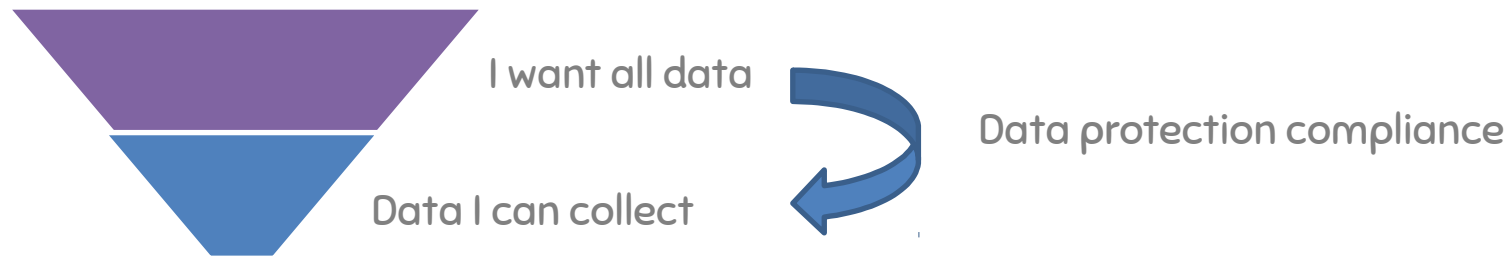
obfuscate the data

anonymize the data



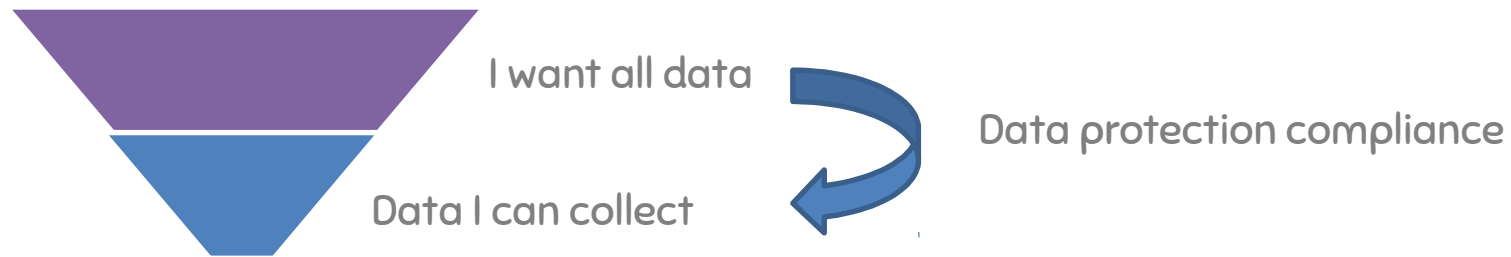
# A CHANGE IN OUR WAY OF THINKING....

## THE USUAL APPROACH

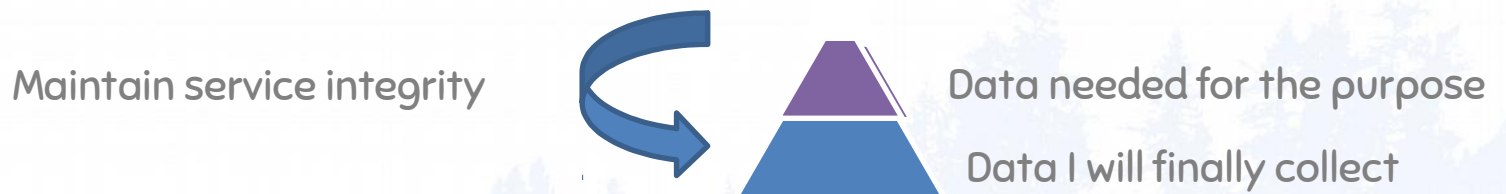


# A CHANGE IN OUR WAY OF THINKING....

## THE USUAL APPROACH



## THE PBD APPROACH



# A CHANGE IN OUR WAY OF THINKING....

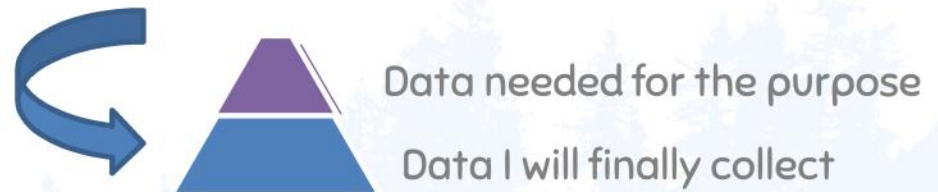
## THE USUAL APPROACH



Maintain security

**PETS**

## THE PBD APPROACH



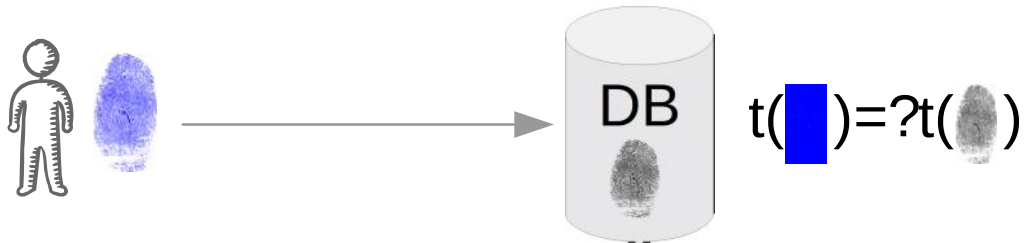
# OTHER CASE STUDIES: PRIVACY-PRESERVING BIOMETRICS

## THE USUAL APPROACH



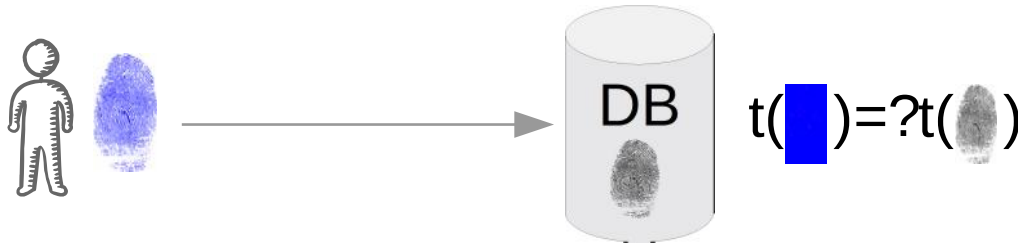
# OTHER CASE STUDIES: PRIVACY-PRESERVING BIOMETRICS

## THE USUAL APPROACH



# OTHER CASE STUDIES: PRIVACY-PRESERVING BIOMETRICS

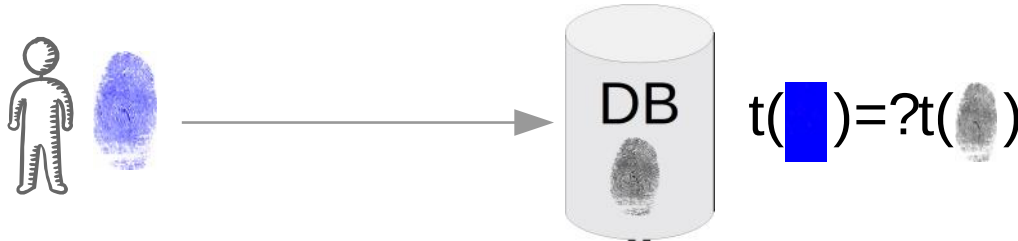
## THE USUAL APPROACH



Templates linkable across databases  
Reveal clear biometric  
Not revocable  
Many times not externalizable

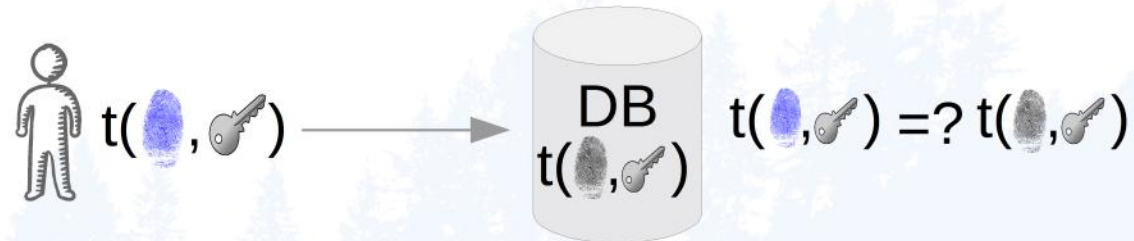
# OTHER CASE STUDIES: PRIVACY-PRESERVING BIOMETRICS

## THE USUAL APPROACH



Templates linkable across databases  
Reveal clear biometric  
Not revocable  
Many times not externalizable

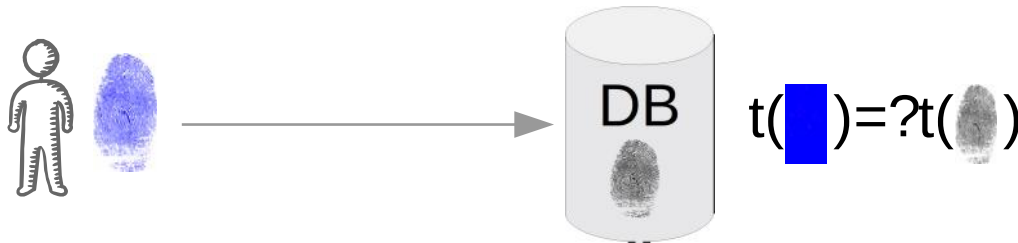
## THE PBD APPROACH





# OTHER CASE STUDIES: PRIVACY-PRESERVING BIOMETRICS

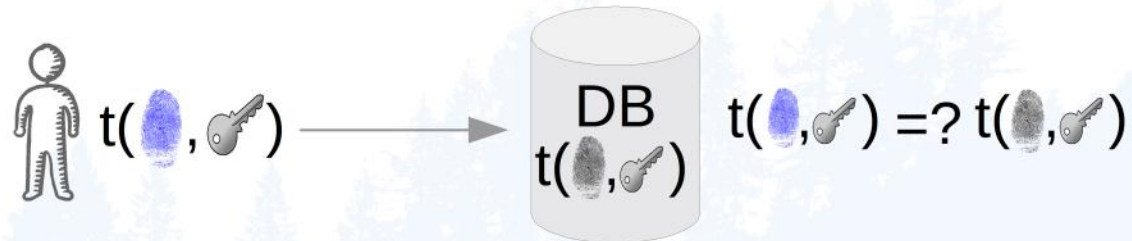
## THE USUAL APPROACH



Templates linkable across databases  
Reveal clear biometric  
Not revocable  
Many times not externalizable

Templates **NOT** linkable across databases  
Not Reveal clear biometric  
Revocable  
**EXTERNALIZABLE**

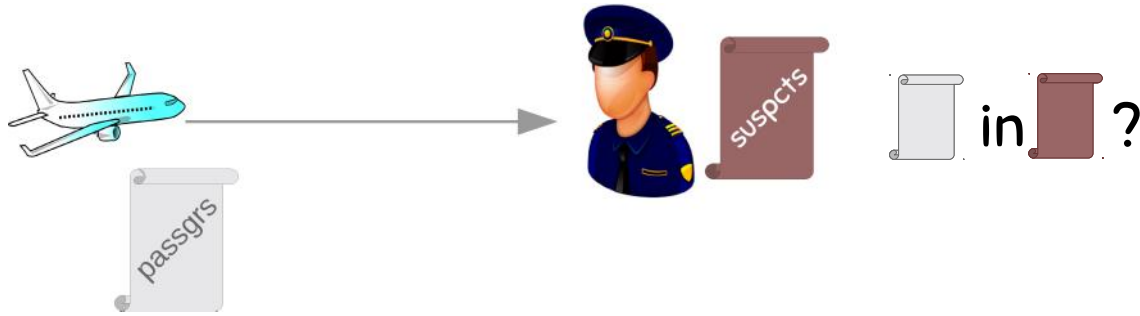
## THE PBD APPROACH



# OTHER CASE STUDIES:

## PRIVACY-PRESERVING PASSENGER REGISTRY

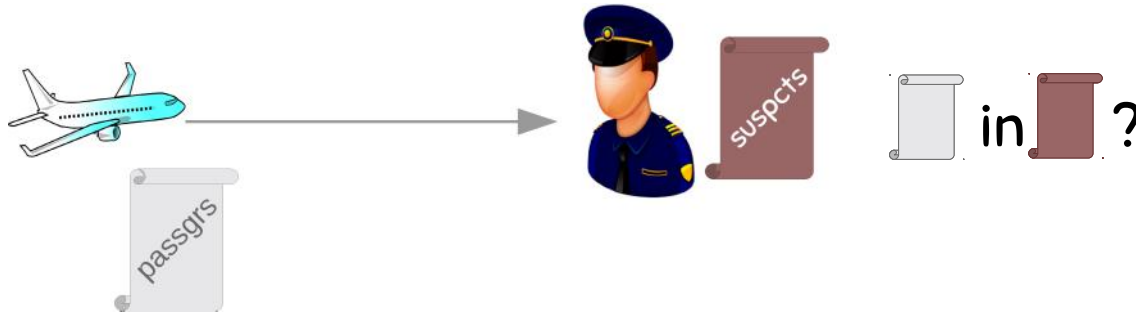
### THE USUAL APPROACH



# OTHER CASE STUDIES:

## PRIVACY-PRESERVING PASSENGER REGISTRY

### THE USUAL APPROACH

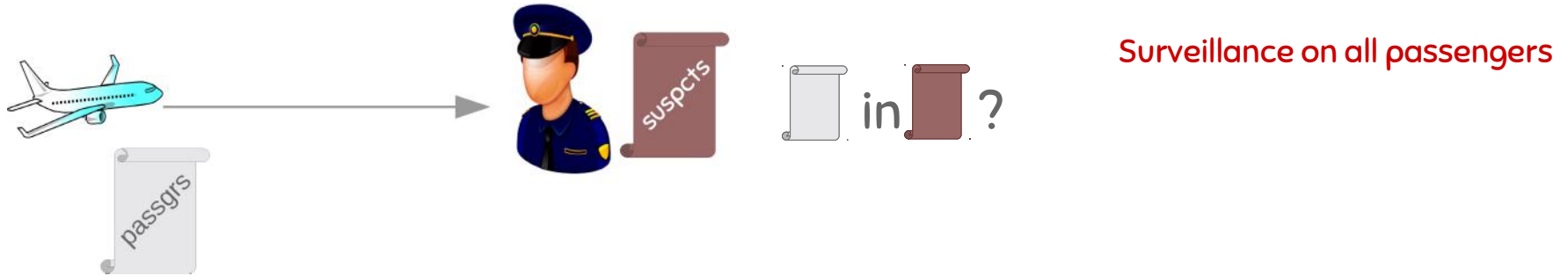


Surveillance on all passengers

# OTHER CASE STUDIES:

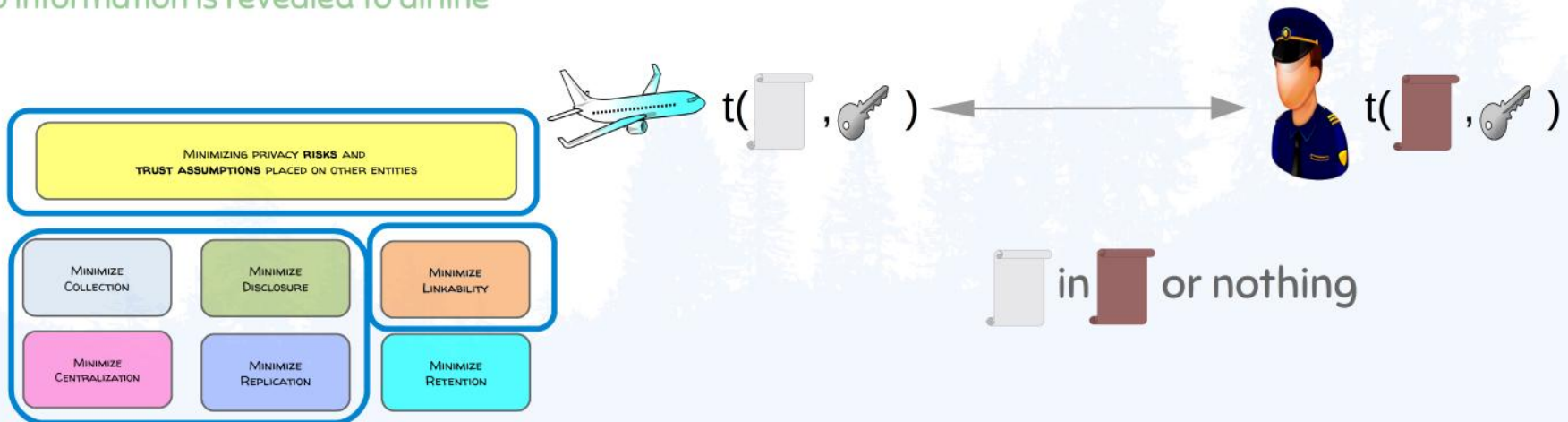
## PRIVACY-PRESERVING PASSENGER REGISTRY

### THE USUAL APPROACH

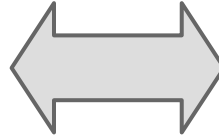


### THE PBD APPROACH

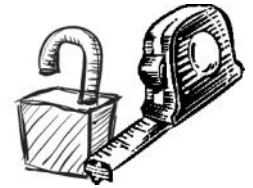
Only coincidences are revealed to law enforcement  
No information is revealed to airline



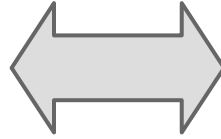
**PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS**



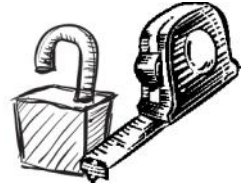
**PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS**



PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS

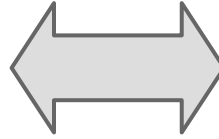


PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS

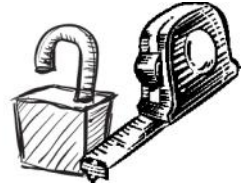


PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS



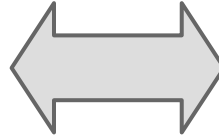
## PRIVACY-PRESERVING SOLUTIONS

**CRYPTO-BASED** VS ANONYMIZATION/OBFUSCATION

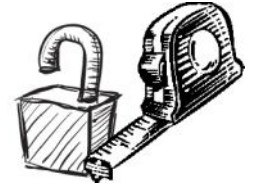
### WELL ESTABLISHED DESIGN AND EVALUATION METHODS

- Private searches
- Private billing
- Private comparison
- Private sharing
- Private statistics computation
- Private electronic cash
- Private genomic computations
- ...

PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS



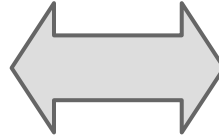
PRIVACY-PRESERVING SOLUTIONS

**CRYPTO-BASED** VS ANONYMIZATION/OBFUSCATION

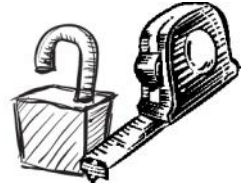
WELL ESTABLISHED DESIGN AND EVALUATION METHODS  
but expensive and require expertise



PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



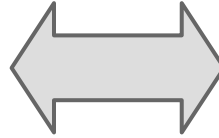
PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS



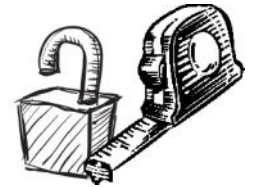
PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

cheap but...

PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



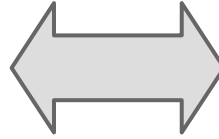
PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS



PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

cheap but...  
DIFFICULT TO DESIGN / EVALUATE

PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS

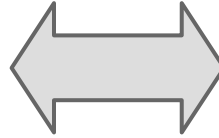


PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

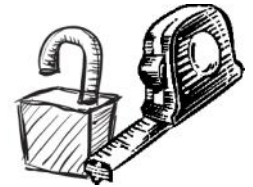
cheap but...  
DIFFICULT TO DESIGN / EVALUATE



PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS

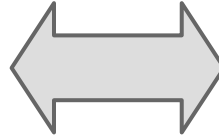


PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

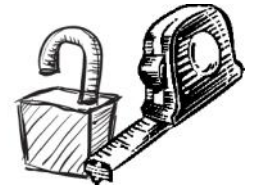
cheap but...  
DIFFICULT TO DESIGN / EVALUATE



PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS



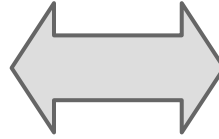
PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

cheap but...  
DIFFICULT TO DESIGN / EVALUATE



The adversary knows!

PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS



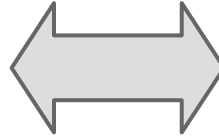
PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

cheap but...  
DIFFICULT TO DESIGN / EVALUATE

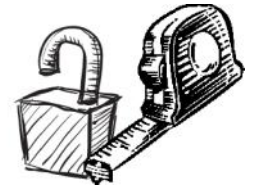


The adversary knows!

PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS



PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

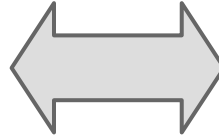
cheap but...  
DIFFICULT TO DESIGN / EVALUATE



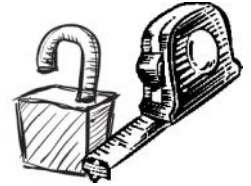
The adversary knows!



PART I:  
REASONING ABOUT  
PRIVACY WHEN DESIGNING  
SYSTEMS



PART II:  
EVALUATING PRIVACY IN  
PRIVACY-PRESERVING  
SYSTEMS



PRIVACY-PRESERVING SOLUTIONS  
CRYPTO-BASED VS ANONYMIZATION/OBFUSCATION

cheap but...  
DIFFICULT TO DESIGN / EVALUATE



The adversary knows!



KEY!! to design  systematically!



# WE NEED TECHNICAL OBJECTIVES – PRIVACY GOALS

**PSEUDONYMITY:** pseudonymous as ID (personal data!)

**ANONYMITY:** decoupling identity and action

**UNLINKABILITY:** hiding link between actions

**UNOBSERVABILITY:** hiding the very existence of actions

**PLAUSIBLE DENIABILITY:** not possible to prove a link between identity and action

**“OBFUSCATION”:** not possible to recover a real item from a noisy item

# WE NEED TECHNICAL OBJECTIVES – PRIVACY GOALS

**PSEUDONYMITY:** pseudonymous as ID (personal data!)

**ANONYMITY:** decoupling identity and action

**UNLINKABILITY:** hiding link between actions

**UNOBSERVABILITY:** hiding the very existence of actions

**PLAUSIBLE DENIABILITY:** not possible to prove a link between identity and action

**“OBFUSCATION”:** not possible to recover a real item from a noisy item

**WHY IS IT SO DIFFICULT TO EVALUATE THEM?**

# LET'S TAKE ONE EXAMPLE: ANONYMITY

Art. 29 WP's opinion on anonymization techniques:

3 criteria to decide a dataset is non-anonymous (pseudonymous):

- 1) is it still possible to single out an individual
- 2) is it still possible to link two records within a dataset (or between two datasets)
- 3) can information be inferred concerning an individual?

# LET'S TAKE ONE EXAMPLE: ANONYMITY

## 1) IS IT STILL POSSIBLE TO SINGLE OUT AN INDIVIDUAL

### On the Anonymity of Home/Work Location Pairs

Philippe Golle and Kurt Partridge

Palo Alto Research Center  
{pgolle, kurt}@parc.com

### Unique in the Crowd: The privacy bounds of human mobility

Yves-Alexandre de Montjoye<sup>1,2</sup>, César A. Hidalgo<sup>1,3,4</sup>, Michel Verleysen<sup>2</sup> & Vincent D. Blondel<sup>2,5</sup>

**Abstract.** Many applications benefit from user location data raises privacy concerns. Anonymization

location

<sup>1</sup>Massachusetts Institute of Technology, Media Lab, 20 Ames Street, Cambridge, MA 02139 USA, <sup>2</sup>Université catholique de Louvain, Institute for Information and Communication Technologies, Electronics and Applied Mathematics, Avenue Georges Lemaitre 4, B-1348 Louvain-la-Neuve, Belgium, <sup>3</sup>Harvard University, Center for International Development, 79 JFK Street, Cambridge, MA 02138, USA, <sup>4</sup>Instituto de Sistemas Complejos de Valparaíso, Paseo 21 de Mayo, Valparaíso, Chile, <sup>5</sup>Massachusetts Institute of Technology, Laboratory for Information and Decision Systems, 77 Massachusetts Avenue, Cambridge, MA 02139, USA.

We study fifteen months of human mobility data for one and a half million individuals and find that human mobility traces are highly unique. In fact, in a dataset where the location of an individual is specified hourly, and with a spatial resolution equal to that given by the carrier's antennas, four spatio-temporal points are enough to uniquely identify 95% of the individuals. We coarsen the data spatially and temporally to find a

“the median size of the individual's anonymity set in the U.S. working population is 1, 21 and 34,980, for locations known at the granularity of a census block, census tract and county respectively”



# LET'S TAKE ONE EXAMPLE: ANONYMITY

## 1) IS IT STILL POSSIBLE TO SINGLE OUT AN INDIVIDUAL

### On the Anonymity of Home/Work Location Pairs

Philippe Golle and Kurt Partridge

Palo Alto Research Center  
{golle, kurt}@parc.com

### Unique in the Crowd: The privacy bounds of human mobility

Yves-Alexandre de Montjoye<sup>1,2</sup>, César A. Hidalgo<sup>1,3,4</sup>, Michel Verleysen<sup>5</sup> & Vincent D. Blondel<sup>6</sup>

**Abstract.** Many applications benefit from using location data raises privacy concerns. Anonymization

location

<sup>1</sup>Massachusetts Institute of Technology, Media Lab, 20 Ames Street, Cambridge, MA 02139 USA, <sup>2</sup>Université catholique de Louvain, Institute for Information and Communication Technologies, Electronics and Applied Mathematics, Avenue Georges Lemaitre 4, B-1348 Louvain-la-Neuve, Belgium, <sup>3</sup>Harvard University, Center for International Development, 79 JFK Street, Cambridge, MA 02138, USA, <sup>4</sup>Instituto de Sistemas Complejos de Valparaíso, Paseo 21 de Mayo, Valparaíso, Chile, <sup>5</sup>Massachusetts Institute of Technology, Laboratory for Information and Decision Systems, 77 Massachusetts Avenue, Cambridge, MA 02139, USA.

We study fifteen months of human mobility data for one and a half million individuals and find that human mobility traces are highly unique. In fact, in a dataset where the location of an individual is specified hourly, and with a spatial resolution equal to that given by the carrier's antennas, four spatio-temporal points are enough to uniquely identify 95% of the individuals. We coarsen the data spatially and temporally to find a

“if the location of an individual is specified hourly, and with a spatial resolution equal to that given by the carrier's antennas, four spatio-temporal points are enough to uniquely identify 95% of the individuals.” [15 month, 1.5M people]

# LET'S TAKE ONE EXAMPLE: ANONYMITY

## 1) IS IT STILL POSSIBLE TO SINGLE OUT AN INDIVIDUAL

### On the Anonymity of Home/Work Location Pairs

Philippe Golle and Kurt Partridge

Palo Alto Research Center  
{pgolle, kurt}@parc.com

### Unique in the Crowd: The privacy bounds of human mobility

Yves-Alexandre de Montjoye<sup>1,2</sup>, César A. Hidalgo<sup>1,3,4</sup>, Michel Verleysen<sup>2</sup> & Vincent D. Blondel<sup>2,3</sup>

**Abstract.** Many applications benefit from user location data raises privacy concerns. Anonymization

<sup>1</sup>Massachusetts Institute of Technology, Media Lab, 20 Ames Street, Cambridge, MA 02139 USA, <sup>2</sup>Université catholique de Louvain, Institute for Information and Communication Technologies, Electronics and Applied Mathematics, Avenue Georges Lemaitre 4, B-1348 Louvain-la-Neuve, Belgium, <sup>3</sup>Harvard University, Center for International Development, 79 JFK Street, Cambridge, MA 02138, USA, <sup>4</sup>Instituto de Sistemas Complejos de Valparaíso, Paseo 21 de Mayo, Valparaíso, Chile, <sup>5</sup>Massachusetts Institute of Technology, Laboratory for Information and Decision Systems, 77 Massachusetts Avenue, Cambridge, MA 02139, USA.

location

We study fifteen months of human mobility data for one and a half million individuals and find that human mobility traces are highly unique. In fact, in a dataset where the location of an individual is specified hourly, and with a spatial resolution equal to that given by the carrier's antennas, four spatio-temporal points are enough to uniquely identify 95% of the individuals. We coarsen the data spatially and temporally to find a

### How Unique is Your Browser? *a report on the Panopticlick experiment*



Peter Eckersley  
Senior Staff Technologist  
Electronic Frontier Foundation  
pde@eff.org

83.6% had completely unique fingerprints  
(entropy: 18.1 bits, or more)

94.2% of "typical desktop browsers" were unique  
(entropy: 18.8 bits, or more)

web browser

# LET'S TAKE ONE EXAMPLE: ANONYMITY

## 1) IS IT STILL POSSIBLE TO SINGLE OUT AN INDIVIDUAL

### On the Anonymity of Home/Work Location Pairs

Philippe Golle and Kurt Partridge

Palo Alto Research Center  
{pgolle, kurt}@parc.com

### Unique in the Crowd: The privacy bounds of human mobility

Yves-Alexandre de Montjoye<sup>1,2</sup>, César A. Hidalgo<sup>1,3,4</sup>, Michel Verleysen<sup>2</sup> & Vincent D. Blondel<sup>5,6</sup>

**Abstract.** Many applications benefit from user location data raises privacy concerns. Anonymization

location

<sup>1</sup>Massachusetts Institute of Technology, Media Lab, 20 Ames Street, Cambridge, MA 02139 USA, <sup>2</sup>Université catholique de Louvain, Institute for Information and Communication Technologies, Electronics and Applied Mathematics, Avenue Georges Lemaitre 4, B-1348 Louvain-la-Neuve, Belgium, <sup>3</sup>Harvard University, Center for International Development, 79 JFK Street, Cambridge, MA 02138, USA, <sup>4</sup>Instituto de Sistemas Complejos de Valparaíso, Paseo 21 de Mayo, Valparaíso, Chile, <sup>5</sup>Massachusetts Institute of Technology, Laboratory for Information and Decision Systems, 77 Massachusetts Avenue, Cambridge, MA 02139, USA.

We study fifteen months of human mobility data for one and a half million individuals and find that human mobility traces are highly unique. In fact, in a dataset where the location of an individual is specified hourly, and with a spatial resolution equal to that given by the carrier's antennas, four spatio-temporal points are enough to uniquely identify 95% of the individuals. We coarsen the data spatially and temporally to find a

### How Unique is Your Browser? *a report on the Panopticlick experiment*



Peter Eckersley  
Senior Staff Technologist  
Electronic Frontier Foundation  
pde@eff.org

web browser

L. Sweeney, Simple Demographics Often Identify People Uniquely, Carnegie Mellon University, Data Privacy Working Paper 3, Pittsburgh 2000.

### Simple Demographics Often Identify People Uniquely

Latanya Sweeney  
Carnegie Mellon University  
latanya@andrew.cmu.edu

"It was found that 87% (216 million of 248 million) of the population in the United States had reported characteristics that likely made them unique based only on {5-digit ZIP, gender, date of birth}"

# LET'S TAKE ONE EXAMPLE: ANONYMITY

## 2) LINK TWO RECORDS WITHIN A DATASET (OR DATASETS)

### De-anonymizing Social Networks

Arvind Narayanan and Vitaly Shmatikov  
The University of Texas at Austin

#### Abstract

Operators of online social networks are increasingly sharing potentially sensitive information about users and their relationships with advertisers, application developers, and data-mining researchers. Privacy is typically protected by anonymization, i.e., removing names, addresses, etc. We present a framework for analyzing privacy and anonymity in social networks and develop a new re-identification algorithm targeting anonymized social-network graphs. To demonstrate its effectiveness on real-

associated with individual nodes are suppressed. Such suppression is often misinterpreted as removal of "personally identifiable information" (PII), even though PII may include much more than names and identifiers (see the discussion in Appendix B). For example, the EU privacy directive defines "personal data" as "any information relating to an identified or identifiable natural person [...]"; an identifiable person is one who can be identified, directly or indirectly, in particular by one or more factors relating to their physical, physi-

take two graphs representing social networks and map the nodes to each other based on the *graph structure alone*  
—no usernames, no nothing  
**NETFLIX PRIZE, KAGGLE CONTEST**

### An Automated Social Graph De-anonymization Technique

Kumar Sharad  
University of Cambridge, UK  
kumar.sharad@cl.cam.ac.uk

George Danezis  
University College London, UK  
g.danezis@ucl.ac.uk

#### ABSTRACT

We present a generic and automated approach to re-identifying nodes in anonymized social networks which enables novel anonymization techniques to be quickly evaluated. It uses machine learning (decision forests) to matching pairs of nodes in disparate anonymized sub-graphs. The technique uncovers artifacts and in-

Social network graphs in particular are high dimensional and feature rich data sets, and it is extremely hard to preserve their anonymity. Thus, any anonymization scheme has to be evaluated in detail, including those with a sound theoretical basis [11]. Techniques have been proposed to resist de-anonymization [8, 17, 22], however, Dwork and Naor have shown [7] that preserving privacy of

social graphs





# LET'S TAKE ONE EXAMPLE: ANONYMITY

## De-anonymizing Social Networks

Arvind Narayanan and Vitaly Shmatikov  
The University of Texas at Austin

### Abstract

Operators of online social networks are increasingly sharing potentially sensitive information about users and their relationships with advertisers, application developers, and data-mining researchers. Privacy is typically protected by anonymization, i.e., removing names, addresses, etc. We present a framework for analyzing privacy and anonymity in social networks and develop a new re-identification algorithm targeting anonymized social-network graphs. To demonstrate its effectiveness on real-

associated with individual nodes are suppressed. Such suppression is often misinterpreted as removal of "personally identifiable information" (PII), even though PII may include much more than names and identifiers (see the discussion in Appendix B). For example, the EU privacy directive defines "personal data" as "any information relating to an identified or identifiable natural person [...]"; an identifiable person is one who can be identified, directly or indirectly, in particular by one or more factors such as physical, physiological, genetic, mental, economic,

take two graphs representing social networks and map the nodes to each other based on the *graph structure alone*  
—no usernames, no nothing  
**NETFLIX PRIZE, KAGGLE CONTEST**

## An Automated Social Graph De-anonymization Technique

Kumar Sharad  
University of Cambridge, UK  
kumar.sharad@cl.cam.ac.uk

George Danezis  
University College London, UK  
g.danezis@ucl.ac.uk

social graphs

### ABSTRACT

We present a generic and automated approach to re-identifying nodes in anonymized social networks which enables novel anonymization techniques to be quickly evaluated. It uses machine learning techniques for recognizing pairs of nodes in disparate anonymized sub-graphs. The technique involves analyzing and in-

Social network graphs in particular are high dimensional and feature rich data sets, and it is extremely hard to preserve their anonymity. Thus, any anonymization scheme has to be evaluated in detail, including those with a sound theoretical basis [11]. Techniques have been proposed to resist de-anonymization [8, 17, 22], however, Dwork and Naor have shown [7] that preserving privacy of

Technique to automate graph de-anonymization based on machine learning.  
Does not need to know the algorithm!

# LET'S TAKE ONE EXAMPLE: ANONYMITY

## 2) LINK TWO RECORDS WITHIN A DATASET (OR DATASETS)

### De-anonymizing Social Networks

Arvind Narayanan and Vitya Shmatikov  
The University of Texas at Austin

#### Abstract

Operators of online social networks are increasingly sharing potentially sensitive information about users and their relationships with advertisers, application developers, and data-mining researchers. Privacy is typically protected by anonymization, i.e., removing names, addresses, etc.

We present a framework for analyzing privacy and anonymity in social networks and develop a new re-identification algorithm targeting anonymized social-network graphs. To demonstrate its effectiveness on real-

associated with individual nodes are suppressed. Such suppression is often misinterpreted as removal of "personally identifiable information" (PII), even though PII may include much more than names and identifiers (see the discussion in Appendix B). For example, the EU privacy directive defines "personal data" as "any information relating to an identified or identifiable natural person [...] an identifiable person is one who can be identified, directly or indirectly, in particular by one or more fac-

mental, econom

### An Automated Social Graph De-anonymization Technique

Kumar Sharad  
University of Cambridge, UK  
kumar.sharad@cl.cam.ac.uk

Georga Danezis  
University College London, UK  
g.danezis@ucl.ac.uk

## social graphs

#### ABSTRACT

We present a generic and automated approach to re-identifying nodes in anonymized social networks which enables novel anonymization techniques to be quickly evaluated. It uses machine learning (decision forests) to matching pairs of nodes in disparate anonymized sub-graphs. The technique uncovers artifacts and in-

Social network graphs in particular are high dimensional and feature rich data sets, and it is extremely hard to preserve their anonymity. Thus, any anonymization scheme has to be evaluated in detail, including those with a sound theoretical basis [1]. Techniques have been proposed to resist de-anonymization [8, 17, 22], however, Dwork and Naor have shown [7] that preserving privacy of

Link messages from same person  
with different pseudonyms

## Doppelgänger Finder: Taking Stylometry To The Underground

Sadia Afroz<sup>\*</sup>, Aylin Caliskan-Islami<sup>†</sup>, Ariel Stolerman<sup>‡</sup>, Rachel Greenstadt<sup>†</sup> and Damon McCoy<sup>†</sup>  
<sup>\*</sup>University of California, Berkeley <sup>†</sup>Drexel University <sup>‡</sup>George Mason University

**Abstract**—Stylometry is a method for identifying anonymous authors of anonymous texts by analyzing their writing style. While stylometric methods have produced impressive results in previous experiments, we wanted to explore their performance on a challenging dataset of particular interest to the security research community. Analysis of underground forums can provide key information about who controls a given bot network or other malicious actors and the flow and source of the information

Other information gleaned from underground forums is providing security researchers, law enforcement, and policy makers valuable information on how the market is segmented and specialized, the social dynamics of the community, and potential bottlenecks that are vulnerable to interventions. These advances have been accomplished primarily through

DE GRUYTER OPEN

Proceedings on Privacy Enhancing Technologies · 2016 (1):155–171

Rebekah Overdorf<sup>\*</sup> and Rachel Greenstadt

## Blogs, Twitter Feeds, and Reddit Comments: Cross-domain Authorship Attribution

**Abstract**—Stylometry is a form of authorship attribution that relies on the linguistic information to attribute

curity by serving as a verification or identification tool for digital text across the Internet.

As social media and micro-blogging sites increase in popularity, so does the need to identify the authors of these types of text. The accuracy with which stylometry can identify anonymous and pseudonymous authors has direct security implications. It can be used for verification of a person's claimed identity, or to identify the author of an anonymous threat should a suspect not be

## stylometry

# LET'S TAKE ONE EXAMPLE: ANONYMITY

## 2) LINK TWO RECORDS WITHIN A DATASET (OR DATASETS)

### De-anonymizing Social Networks

Arvind Narayanan and Vitya Shmatikov  
The University of Texas at Austin

#### Abstract

Operators of online social networks are increasingly sharing potentially sensitive information about users and their relationships with advertisers, application developers, and data-mining researchers. Privacy is typically protected by anonymization, i.e., removing names, addresses, etc. We present a framework for analyzing privacy and anonymity in social networks and develop a new re-identification algorithm targeting anonymized social-network graphs. It demonstrates its effectiveness on real-

associated with individual nodes are suppressed. Such suppression is often misinterpreted as removal of "personally identifiable information" (PII), even though PII may include much more than names and identifiers (see the discussion in Appendix B). For example, the EU privacy directive defines "personal data" as "any information relating to an identified or identifiable natural person [...] an identifiable person is one who can be identified, directly or indirectly, in particular by one or more fac-

### An Automated Social Graph De-anonymization Technique

Kumar Sharad  
University of Cambridge, UK  
kumar.sharad@cl.cam.ac.uk

Georgia Danezis  
University College London, UK  
g.danezis@ucl.ac.uk

#### ABSTRACT

We present a generic and automated approach to re-identifying nodes in anonymized social networks which enables novel anonymization techniques to be quickly evaluated. It uses machine learning (decision forests) to matching pairs of nodes in disparate anonymized sub-graphs. The technique uncovers artifacts and in-

Social network graphs in particular are high dimensional and feature rich data sets, and it is extremely hard to preserve their anonymity. Thus, any anonymization scheme has to be evaluated in detail, including those with a sound theoretical basis [1]. Techniques have been proposed to resist de-anonymization [8, 17, 22], however, Dwork and Naor have shown [7] that preserving privacy of

social graphs

Authorship attribution also works across domains!!

DE GRUYTER OPEN

Proceedings on Privacy Enhancing Technologies · 2016 (1): 155–171

Rebekah Overdorf\* and Rachel Greenstadt

### Blogs, Twitter Feeds, and Reddit Comments: Cross-domain Authorship Attribution

**Abstract.** Stylometry is a form of authorship attribution that relies on the linguistic information to attribute

curity by serving as a verification or identification tool for digital text across the Internet.

As social media and micro-blogging sites increase in popularity, so does the need to identify the authors of these types of text. The accuracy with which stylometry can identify anonymous and pseudonymous authors has direct security implications. It can be used for verification of a person's claimed identity, or to identify the author of an anonymous threat should a suspect not be

### Doppelgänger Finder: Taking Stylometry To The Underground

Sadia Afroz\*, Aylin Caliskan-Islami<sup>1</sup>, Ariel Stolerman<sup>1</sup>, Rachel Greenstadt<sup>1</sup> and Damon McCoy<sup>2</sup>  
<sup>1</sup>University of California, Berkeley <sup>2</sup>Drexel University <sup>3</sup>George Mason University

**Abstract.** Stylometry is a method for identifying anonymous authors of anonymous texts by analyzing their writing style. While stylometric methods have produced impressive results in previous experiments, we wanted to explore their performance on a challenging dataset of particular interest to the security research community. Analysis of underground forums can provide key information about who controls a given bot network or other malicious actors and the flow and source of the information

Other information gleaned from underground forums is providing security researchers, law enforcement, and policy makers valuable information on how the market is segmented and specialized, the social dynamics of the community, and potential bottlenecks that are vulnerable to interventions. These advances have been accomplished primarily through

Link messages from same person with different pseudonyms

stylometry

# “ANTI-SURVEILLANCE PETS” TECHNICAL GOALS

## PRIVACY PROPERTIES: ANONYMITY

### 3) INFER INFORMATION ABOUT AN INDIVIDUAL

#### Inference Attacks on Location Tracks

John Krumm

Microsoft Research  
One Microsoft Way  
Redmond, WA, USA  
jckrumm@microsoft.com

**Abstract.** Although the privacy threats and countermeasures associated with location data are well known, there has not been a thorough experiment to assess the effectiveness of either. We examine location data gathered from volunteer subjects to quantify how well four different algorithms can identify

“Based on GPS tracks from, we identify the latitude and longitude of their homes. From these locations, we used a free Web service to do a reverse “white pages” lookup, which takes a latitude and longitude coordinate as input and gives an address and name. [172 individuals]”

# LET'S TAKE ONE EXAMPLE: ANONYMITY

## 3) INFER INFORMATION ABOUT AN INDIVIDUAL

### Inference Attacks on Location Tracks

John Krumm

Microsoft Research  
One Microsoft Way  
Redmond, WA, USA  
jckrumm@microsoft.com

**Abstract.** Although the privacy threats and countermeasures associated with location data are well known, there has not been a thorough experiment to assess the effectiveness of either. We examine location data gathered from volunteer subjects to quantify how well four different algorithms can identify

“We investigate the subtle cues to user identity that may be exploited in attacks on the privacy of users in web search query logs. We study the application of simple classifiers to map a sequence of queries into the gender, age, and location of the user issuing the queries.”

### “I Know What You Did Last Summer” — Query Logs and User Privacy

Rosie Jones    Ravi Kumar    Bo Pang    Andrew Tomkins  
Yahoo! Research, 701 First Ave, Sunnyvale, CA 94089.  
{jonesr,ravikumar,bopang,atomkins}@yahoo-inc.com

#### ABSTRACT

We investigate the subtle cues to user identity that may be exploited in attacks on the privacy of users in web search query logs. We study the application of simple classifiers to map a sequence of queries into the gender, age, and location of the user issuing the queries. We then show how these classifiers may be carefully combined at multiple granularities to map a sequence of queries into a

ilities; this is the goal of this paper. We initiate the study of subtle cues to user identity that exist as vulnerabilities in web search query logs, which may be exploited in attacks on the privacy of users.

**Privacy attack models.** We begin with a characterization of two key forms of attack against which a query log privacy scheme must be resilient. The first is a *trace attack*, in which an attacker studies a privacy-enhanced version of a sequence of searches (*trace*) made

# LET'S TAKE ONE EXAMPLE: ANONYMITY

**MAGICAL THINKING!**  
THIS CANNOT HAPPEN IN GENERAL!



DATA ANONYMIZATION IS A WEAK PRIVACY MECHANISM

ONLY TO BE USED WHEN OTHER PROTECTIONS ARE ALSO APPLIED.

(CONTRACTUAL, ORGANIZATIONAL)

# LET'S TAKE ONE EXAMPLE: ANONYMITY

**MAGICAL THINKING!**  
THIS CANNOT HAPPEN IN GENERAL!



DATA ANONYMIZATION IS A WEAK PRIVACY MECHANISM

ONLY TO BE USED WHEN OTHER PROTECTIONS ARE ALSO APPLIED.

(CONTRACTUAL, ORGANIZATIONAL)

IMPOSSIBLE TO SANITIZE WITHOUT SEVERELY DAMAGING USEFULNESS

REMOVING PII IS NOT ENOUGH! – ANY ASPECT COULD LEAD TO RE-IDENTIFICATION

# LET'S TAKE ONE EXAMPLE: ANONYMITY

MAGICAL THINKING!  
THIS CANNOT HAPPEN IN GENERAL!



DATA ANONYMIZATION IS A WEAK PRIVACY MECHANISM

ONLY TO BE USED WHEN OTHER PROTECTIONS ARE ALSO APPLIED.  
(CONTRACTUAL, ORGANIZATIONAL)

IMPOSSIBLE TO SANITIZE WITHOUT SEVERELY DAMAGING USEFULNESS

REMOVING PII IS NOT ENOUGH! – ANY ASPECT COULD LEAD TO RE-IDENTIFICATION

RISK OF DE-ANONYMIZATION?  PROBABILISTIC ANALYSIS

$\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$



# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



1) MODEL THE PRIVACY-PRESERVING MECHANISM AS A PROBABILISTIC TRANSFORMATION

# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



1) MODEL THE PRIVACY-PRESERVING MECHANISM AS A PROBABILISTIC TRANSFORMATION

2) DETERMINE WHAT THE ADVERSARY WILL SEE

# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



1) MODEL THE PRIVACY-PRESERVING MECHANISM AS A PROBABILISTIC TRANSFORMATION

2) DETERMINE WHAT THE ADVERSARY WILL SEE

data

metadata

...

# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



1) MODEL THE PRIVACY-PRESERVING MECHANISM AS A PROBABILISTIC TRANSFORMATION

2) DETERMINE WHAT THE ADVERSARY WILL SEE

3) “INVERT” THE MECHANISM AS THE ADVERSARY WOULD DO

# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



1) MODEL THE PRIVACY-PRESERVING MECHANISM AS A PROBABILISTIC TRANSFORMATION

2) DETERMINE WHAT THE ADVERSARY WILL SEE

3) “INVERT” THE MECHANISM AS THE ADVERSARY WOULD DO  
**THE ADVERSARY KNOWS!!!**

# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



1) MODEL THE PRIVACY-PRESERVING MECHANISM AS A PROBABILISTIC TRANSFORMATION

**IF IT IS NOT PROBABILISTIC, IT IS NOT SECURE**

2) DETERMINE WHAT THE ADVERSARY WILL SEE

3) “INVERT” THE MECHANISM AS THE ADVERSARY WOULD DO  
**THE ADVERSARY KNOWS!!!**

# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



1) MODEL THE PRIVACY-PRESERVING MECHANISM AS A PROBABILISTIC TRANSFORMATION

**IF IT IS NOT PROBABILISTIC, IT IS NOT SECURE**

2) DETERMINE WHAT THE ADVERSARY WILL SEE

3) “INVERT” THE MECHANISM AS THE ADVERSARY WOULD DO  
**THE ADVERSARY KNOWS!!!**

4) COMPUTE PROBABILITY AFTER “INVERSION”



# PRIVACY EVALUATION IS A PROBABILISTIC ANALYSIS

## SYSTEMATIC REASONING TO EVALUATE A MECHANISM

Anonymity –  $\Pr[\text{identity} \rightarrow \text{action} \mid \text{observation}]$

Unlinkability –  $\Pr[\text{action A} \leftrightarrow \text{action B} \mid \text{observation}]$

Obfuscation –  $\Pr[\text{real action} \mid \text{observed noisy action}]$



1) MODEL THE PRIVACY-PRESERVING MECHANISM AS A PROBABILISTIC TRANSFORMATION

**IF IT IS NOT PROBABILISTIC, IT IS NOT SECURE**

2) DETERMINE WHAT THE ADVERSARY WILL SEE

3) “INVERT” THE MECHANISM AS THE ADVERSARY WOULD DO  
**THE ADVERSARY KNOWS!!!**

4) COMPUTE PROBABILITY AFTER “INVERSION”

5) MEASURE... MEAN ERROR, ENTROPY (ANY FLAVOUR), DIFF. PRIVACY

# “INVERSION”? WHAT DO YOU MEAN?

## 1) ANALYTICAL MECHANISM INVERSION

GIVEN THE DESCRIPTION OF THE SYSTEM, DEVELOP THE MATHEMATICAL EXPRESSIONS THAT EFFECTIVELY INVERT THE SYSTEM:

$$PR[OBS \mid REAL\ DATA, PET] \rightarrow PR[REAL\ DATA \mid OBS, PET]$$



# “INVERSION”? WHAT DO YOU MEAN?

## 1) ANALYTICAL MECHANISM INVERSION

GIVEN THE DESCRIPTION OF THE SYSTEM, DEVELOP THE MATHEMATICAL EXPRESSIONS THAT EFFECTIVELY INVERT THE SYSTEM:

$$PR[OBS \mid REAL\ DATA, PET] \rightarrow PR[REAL\ DATA \mid OBS, PET]$$

**NOT ALWAYS POSSIBLE – MAY REQUIRE APROX. OR SAMPLING**



# “INVERSION”? WHAT DO YOU MEAN?

## 1) ANALYTICAL MECHANISM INVERSION

GIVEN THE DESCRIPTION OF THE SYSTEM, DEVELOP THE MATHEMATICAL EXPRESSIONS THAT EFFECTIVELY INVERT THE SYSTEM:

$$PR[OBS | REAL DATA, PET] \rightarrow PR[REAL DATA | OBS, PET]$$

**NOT ALWAYS POSSIBLE – MAY REQUIRE APROX. OR SAMPLING**

## 2) MACHINE LEARNING (DATA DRIVEN)

TRAIN A CLASSIFIER TO BREAK THE MECHANISMS!



# “INVERSION”? WHAT DO YOU MEAN?

## 1) ANALYTICAL MECHANISM INVERSION

GIVEN THE DESCRIPTION OF THE SYSTEM, DEVELOP THE MATHEMATICAL EXPRESSIONS THAT EFFECTIVELY INVERT THE SYSTEM:

$$PR[OBS | REAL DATA, PET] \rightarrow PR[REAL DATA | OBS, PET]$$

**NOT ALWAYS POSSIBLE – MAY REQUIRE APROX. OR SAMPLING**

## 2) MACHINE LEARNING (DATA DRIVEN)

TRAIN A CLASSIFIER TO BREAK THE MECHANISMS!

**ONLY POSSIBLE IF ENOUGH DATA (THOUGH DATA CAN BE CREATED)**



# “INVERSION”? WHAT DO YOU MEAN?

## 1) ANALYTICAL MECHANISM INVERSION

GIVEN THE DESCRIPTION OF THE SYSTEM, DEVELOP THE MATHEMATICAL EXPRESSIONS THAT EFFECTIVELY INVERT THE SYSTEM:

$$PR[OBS | REAL DATA, PET] \rightarrow PR[REAL DATA | OBS, PET]$$

**NOT ALWAYS POSSIBLE – MAY REQUIRE APROX. OR SAMPLING**

## 2) MACHINE LEARNING (DATA DRIVEN)

TRAIN A CLASSIFIER TO BREAK THE MECHANISMS!

**ONLY POSSIBLE IF ENOUGH DATA (THOUGH DATA CAN BE CREATED)**



MUST TAKE INVERSION INTO ACCOUNT!! SYSTEMATIC DESIGN!!!

# “INVERSION”? WHAT DO YOU MEAN?

## 1) ANALYTICAL MECHANISM INVERSION

GIVEN THE DESCRIPTION OF THE SYSTEM, DEVELOP THE MATHEMATICAL EXPRESSIONS THAT EFFECTIVELY INVERT THE SYSTEM:

$$PR[OBS | REAL DATA, PET] \rightarrow PR[REAL DATA | OBS, PET]$$

NOT ALWAYS POSSIBLE – MAY REQUIRE APROX. OR SAMPLING

## 2) MACHINE LEARNING (DATA DRIVEN)

TRAIN A CLASSIFIER TO BREAK THE MECHANISMS!

ONLY POSSIBLE IF ENOUGH DATA (THOUGH DATA CAN BE CREATED)



MUST TAKE INVERSION INTO ACCOUNT!! SYSTEMATIC DESIGN!!!

THAT'S ANOTHER  
TALK.....

# TAKE AWAYS

PRIVACY BY DESIGN ROCKS!



BUT REALIZING IT IS NON-TRIVIAL



# TAKE AWAYS

PRIVACY BY DESIGN ROCKS!



BUT REALIZING IT IS NON-TRIVIAL

PART I:  
REASONING ABOUT PRIVACY WHEN  
DESIGNING SYSTEMS



Explicit privacy engineering activities

PART II:  
EVALUATING PRIVACY IN PRIVACY-  
PRESERVING SYSTEMS



Systematic reasoning for  
privacy evaluation

# TAKE AWAYS

PRIVACY BY DESIGN ROCKS!



BUT REALIZING IT IS NON-TRIVIAL

PART I:  
REASONING ABOUT PRIVACY WHEN  
DESIGNING SYSTEMS



Explicit privacy engineering activities



PART II:  
EVALUATING PRIVACY IN PRIVACY-  
PRESERVING SYSTEMS



Systematic reasoning for  
privacy evaluation

Fully fledged methodology?

Requirements? Evaluation?

# TAKE AWAYS

PRIVACY BY DESIGN ROCKS!



BUT REALIZING IT IS NON-TRIVIAL

PART I:  
REASONING ABOUT PRIVACY WHEN  
DESIGNING SYSTEMS



Explicit privacy engineering activities



PART II:  
EVALUATING PRIVACY IN PRIVACY-  
PRESERVING SYSTEMS



Systematic reasoning for  
privacy evaluation

Fully fledged methodology?

Requirements? Evaluation?

Accessible PETS

Understanding? Implementation?

# TAKE AWAYS

PRIVACY BY DESIGN ROCKS!



BUT REALIZING IT IS NON-TRIVIAL

PART I:  
REASONING ABOUT PRIVACY WHEN  
DESIGNING SYSTEMS



Explicit privacy engineering activities



Fully fledged methodology?

Requirements? Evaluation?

Accessible PETS

Understanding? Implementation?

PART II:  
EVALUATING PRIVACY IN PRIVACY-  
PRESERVING SYSTEMS



Systematic reasoning for  
privacy evaluation



Strong assumption's dependency

What does the adversary know?

Ad-hoc mechanisms (training!)

Lack of standard metrics

# THANKS!

## ANY QUESTIONS?

More about privacy:

<https://www.petsymposium.org/>

<http://www.degruyter.com/view/j/popets>

17TH PRIVACY ENHANCING TECHNOLOGIES SYMPOSIUM  
JULY 18–21, 2017

MINNEAPOLIS, MN, USA



2018 BARCELONA! DEADLINES: 31 AUG, 30 NOV, 28 FEB

[carmela.troncoso@imdea.org](mailto:carmela.troncoso@imdea.org)

<https://software.imdea.org/~carmela.troncoso/>

(these slides will be there soon)

Template: <http://www.brainybetty.com/>

Figures: [SlidesCarnival](#)

WHAT DO WE WANT THE DATA FOR...? STATISTICS!



## WHAT DO WE WANT THE DATA FOR...? STATISTICS!



“Wouldn't it be nice if I could send complex queries to a database to extract statistics, and it returned results that are informative, but leak very little information about any individual?”



# WHAT DO WE WANT THE DATA FOR...? STATISTICS!



“Wouldn't it be nice if I could send complex queries to a database to extract statistics, and it returned results that are informative, but leak very little information about any individual?”



MAGICAL THINKING?





## WHAT DO WE WANT THE DATA FOR...? STATISTICS!



“Wouldn't it be nice if I could send complex queries to a database to extract statistics, and it returned results that are informative, but leak very little information about any individual?”



**QUERY-BASED PRIVACY**  
**DIFFERENTIAL PRIVACY!**



# WHAT DO WE WANT THE DATA FOR...? STATISTICS!



“Wouldn't it be nice if I could send complex queries to a database to extract statistics, and it returned results that are informative, but leak very little information about any individual?”



**QUERY-BASED PRIVACY  
DIFFERENTIAL PRIVACY!**



Why is that possible (while anonymization was impossible):



## WHAT DO WE WANT THE DATA FOR...? STATISTICS!



“Wouldn't it be nice if I could send complex queries to a database to extract statistics, and it returned results that are informative, but leak very little information about any individual?”



**QUERY-BASED PRIVACY**  
**DIFFERENTIAL PRIVACY!**



Why is that possible (while anonymization was impossible):

The final result **DEPENDS ON MULTIPLE PERSONAL RECORDS**

However it **DOES NOT DEPEND MUCH ON ANY PARTICULAR ONE** (sensitivity)

## WHAT DO WE WANT THE DATA FOR...? STATISTICS!



“Wouldn't it be nice if I could send complex queries to a database to extract statistics, and it returned results that are informative, but leak very little information about any individual?”



**QUERY-BASED PRIVACY  
DIFFERENTIAL PRIVACY!**



Why is that possible (while anonymization was impossible):

The final result **DEPENDS ON MULTIPLE PERSONAL RECORDS**

However it **DOES NOT DEPEND MUCH ON ANY PARTICULAR ONE** (sensitivity)

Therefore adding a little bit of noise to the result, suffices to hide any record contribution

For full anonymization.... one would need to add a lot of noise to all the entries

## WHAT DO WE WANT THE DATA FOR...? STATISTICS!



“Wouldn't it be nice if I could send complex queries to a database to extract statistics, and it returned results that are informative, but leak very little information about any individual?”



**QUERY-BASED PRIVACY**  
**DIFFERENTIAL PRIVACY!**



Why is that possible (while anonymization was impossible):

The final result **DEPENDS ON MULTIPLE PERSONAL RECORDS**

However it **DOES NOT DEPEND MUCH ON ANY PARTICULAR ONE** (sensitivity)

Therefore adding a little bit of noise to the result, suffices to hide any record contribution

For full anonymization.... one would need to add a lot of noise to all the entries



**DIFFERENT ARCHITECTURE TO PROVIDE ROBUST PRIVACY!**

**A TTP HOLDS THE DATA!**

## WHAT DO WE WANT THE DATA FOR...? STATISTICS!



“Wouldn't it be nice if I could send complex queries to a database to extract statistics, and it returned results that are informative, but leak very little information about any individual?”



**QUERY-BASED PRIVACY**  
**DIFFERENTIAL PRIVACY!**



Why is that possible (while anonymization was impossible):

The final result DEPENDS ON MULTIPLE PERSONAL RECORDS

However it DOES NOT DEPEND MUCH ON ANY PARTICULAR ONE (sensitivity)

Therefore adding a little bit of noise to the result, suffices to hide any record contribution

For full anonymization.... one would need to add a lot of noise to all the entries



DIFFERENT ARCHITECTURE TO PROVIDE ROBUST PRIVACY!

A TTP HOLDS THE DATA!

ACTUALLY AFTER SOME USES... UTILITY DROPS

BETTER SUITED FOR ONE-TIME USE → DATA COLLECTION!