

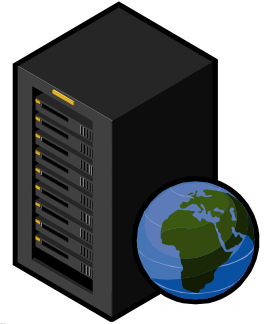


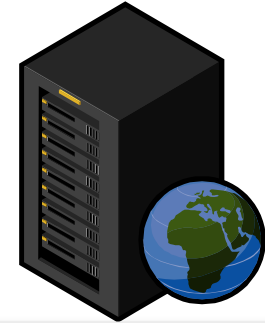
# Ethical Considerations and SocNet Research

Thorsten Strufe

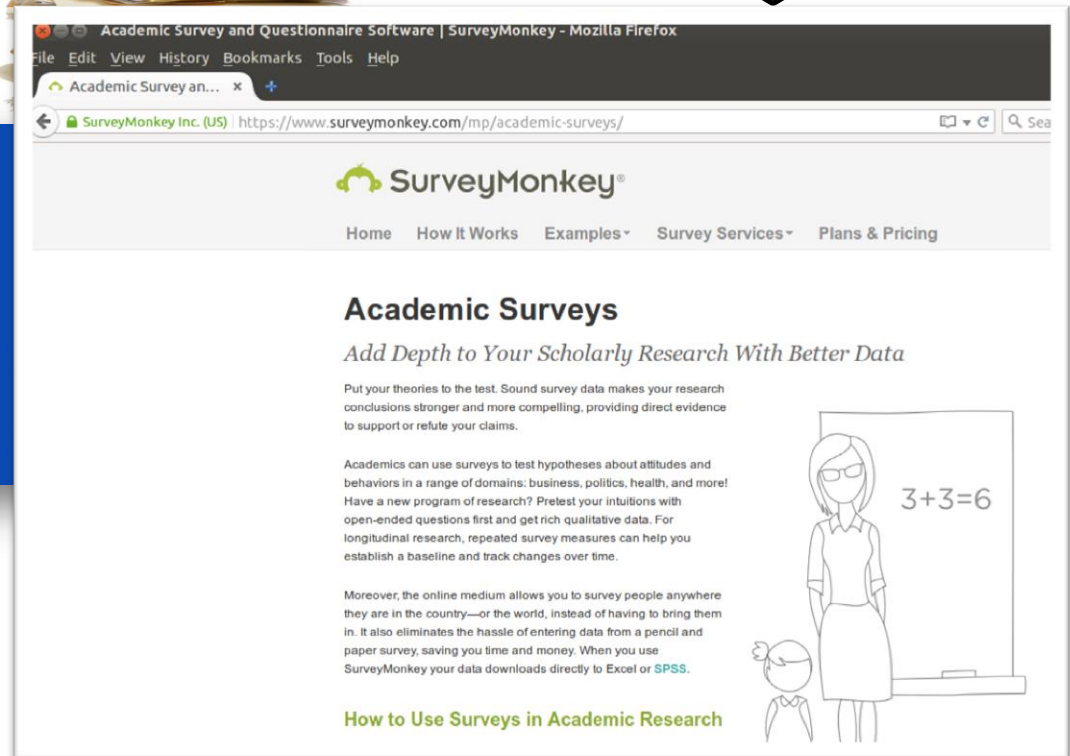
(Anne Lauber-Rönsberg, Stefan Köpsell, Joachim Scharloth)

Padova, 06.09.2016





OneDrive



Academic Survey and Questionnaire Software | SurveyMonkey - Mozilla Firefox

file Edit View History Bookmarks Tools Help

Academic Survey an... x +

SurveyMonkey Inc. (US) | <https://www.surveymonkey.com/mp/academic-surveys/>

SurveyMonkey®

Home How It Works Examples Survey Services Plans & Pricing

## Academic Surveys


*Add Depth to Your Scholarly Research With Better Data*

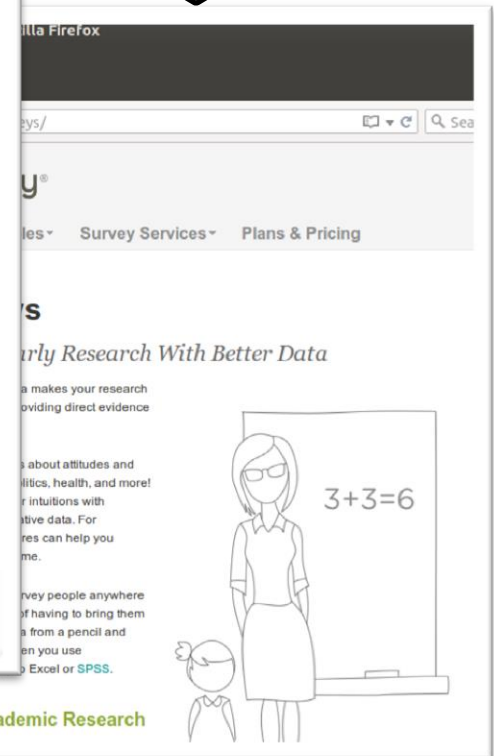
Put your theories to the test. Sound survey data makes your research conclusions stronger and more compelling, providing direct evidence to support or refute your claims.

Academics can use surveys to test hypotheses about attitudes and behaviors in a range of domains: business, politics, health, and more! Have a new program of research? Pretest your intuitions with open-ended questions first and get rich qualitative data. For longitudinal research, repeated survey measures can help you establish a baseline and track changes over time.

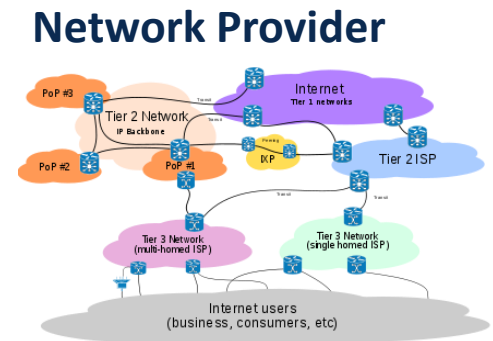
Moreover, the online medium allows you to survey people anywhere they are in the country—or the world, instead of having to bring them in. It also eliminates the hassle of entering data from a pencil and paper survey, saving you time and money. When you use SurveyMonkey your data downloads directly to Excel or SPSS.

[How to Use Surveys in Academic Research](#)





How to Use Surveys in Academic Research



- Subject
  - Online services
  - ...and their users
- Types of studies
  - Purely observational
  - Interactive (questionnaires, discussions, behavior)
- Who are (additional) adversaries?

# ***Ethics***

*...and why bother about it?*

- Institutional Revision Board (IRB), Ethics Commission | Committee
  - Committee to approve, monitor, review research involving humans
- Frequently three aspects:
  - (Informed) consent
  - Benefits (for society)
  - Management of risks (absence of damage to the subject...)
    - Commonly
      - Risk to their bodily well-being
      - Risk to their dignity and reputation -> „anonymize“



- Institutional Revision Board (IRB), Ethics Commission | Committee
  - Committee to approve, monitor, review research involving humans
- Frequently three aspects:
  - (Informed) consent
  - Benefits (for society)
  - Management of risks (absence of damage to the subject...)
    - Commonly
      - Risk to their bodily well-being
      - Risk to their dignity and reputation -> „anonymize“

# Where do you go, before your study?

- Institutional Revision Board (IRB), Ethics Commission | Committee
  - Committee to approve, monitor, review research involving humans
- Frequently
  - (Inform)
  - Benefit
  - Management
    - Commonly
      - Risk to their bodily well-being
      - Risk to their dignity and reputation -> „anonymize“

*But this breaks all my cool research!*  
*We're not evil, why bother anyways!?*



US prosecutor details illegal experiments, Nuremberg, Germany, Between October 1946 and August 1947

© United States Holocaust Memorial Museum, Washington, DC

The following slides cf.: Scharloth, “Research Ethics”, 2015

1. Required is the voluntary, well-informed, understanding consent of the human subject in a full legal capacity.
2. The experiment should aim at positive results for society that cannot be procured in some other way.
3. It should be based on previous knowledge (like, an expectation derived from animal experiments) that justifies the experiment.
4. The experiment should be set up in a way that avoids unnecessary physical and mental suffering and injuries.
5. It should not be conducted when there is any reason to believe that it implies a risk of death or disabling injury.
6. The risks of the experiment should be in proportion to (that is, not exceed) the expected humanitarian benefits.
7. Preparations and facilities must be provided that adequately protect the subjects against the experiment's risks.
8. The staff who conduct or take part in the experiment must be fully trained and scientifically qualified.
9. The human subjects must be free to immediately quit the experiment at any point when they feel physically or mentally unable to go on.
10. Likewise, the medical staff must stop the experiment at any point when they observe that continuation would be dangerous.

1. Required is the **voluntary, well-informed, understanding consent** of the human subject in a full legal capacity.
2. The experiment should **aim at positive results for society** that cannot be procured in some other way.
3. It should be **based on previous knowledge** (like, an expectation derived from animal experiments) that justifies the experiment.
4. The experiment should be set up in a way that **avoids unnecessary physical and mental suffering and injuries**.
5. It **should not be conducted** when there is any reason to believe that it implies a risk of death or disabling injury.
6. The **risks of the experiment** should be in proportion to (that is, not exceed) the expected **humanitarian benefits**.
10. Likewise, the medical staff must stop the experiment at any point when they observe that continuation would be dangerous.

- Respect for persons
  - Individuals should be treated as autonomous agents.
  - Persons with diminished autonomy are entitled to protection.
  - Informed consent
- **Beneficence**
  - Human subjects should not be harmed.
  - Research should maximize possible benefits and minimize possible harms.
- Justice
  - The benefits and risks of research must be distributed fairly.
- **Informed consent**
  - Subjects, to the degree that they are capable, must be given the opportunity to choose what shall or shall not happen to them.
  - The consent process must include three elements: *information*, *comprehension*, and *voluntariness*.
- Assessment of **risks** and benefits
  - The nature and scope of risks and benefits must be assessed in a systematic manner.
- Selection of subjects
  - There must be fair procedures and outcomes in the selection of research subjects.

But we surely can get around this, right?

---

- ***No.***

- ***At least in Europe, we can't:***

- Art. 13 EU-Charter:
  - **Freedom of the arts and sciences**
  - The arts and scientific research shall be free of constraint.
  
- Art. 8 (1) EU-Charter:
  - **Protection of personal data**
  - Everyone has the right to the protection of personal data concerning him or her.

The following slides cf: Lauber-Rönsberg, “Research Ethics”, 2015



- **Personal Data** – relating to an identified or identifiable person
  - Not applicable to *anonymized* data.
  - Still applicable to *pseudonymized* data.
- Data economy (data minimization): data processing limited to minimum amount of data
  - e.g. Art. 83 (1) GDPR as proposed by Commission
- Data Processing only if either **“informed consent”** or permitted by the law
- Data processing **only in accordance with the specified purpose** (further processing for scientific research *may* be permitted)
- **Right to opt out at any time**
- Plus:
  - **Collection of data directly from the data subject**
  - **right to be informed about collected personal data**
  - **right to correct the data**

And, btw: §1 and §12 of  
the Human Rights

- Belmont Report
  - §6 „Support the privacy of the patients identity, their motivation to join or refuse the experiment.”
- Conclusion 1: science must be ethically aware
- Conclusion 2: you need ***informed consent***
- Conclusion 3: Privacy of subjects has to be preserved

# *Privacy*

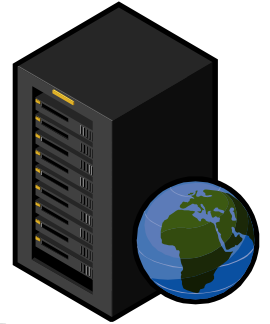
So what is this thing, anyways?

- *Which disclosures are people concerned about? („Study“ from '10)*

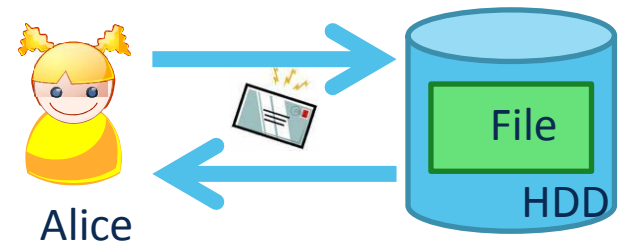
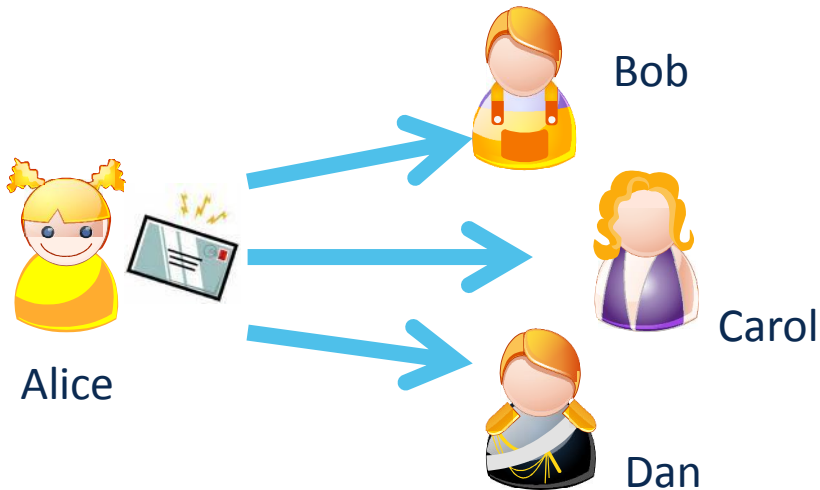
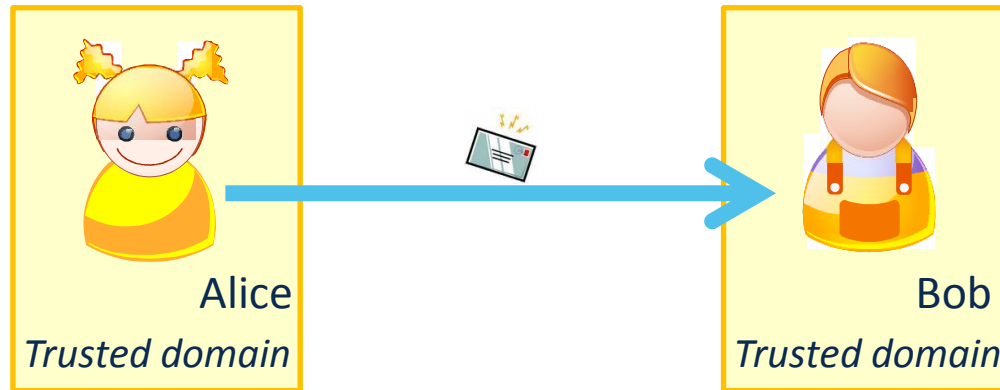


# *Privacy*

So what is this thing, anyways?

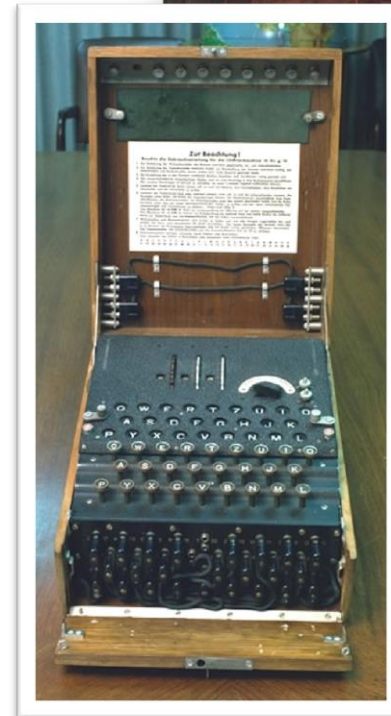


# The Classical Security View





- Data loss
  - Data accessible to unintended parties
- Manipulation and forgery
  - Tampered, spoofed data





- **Confidentiality**

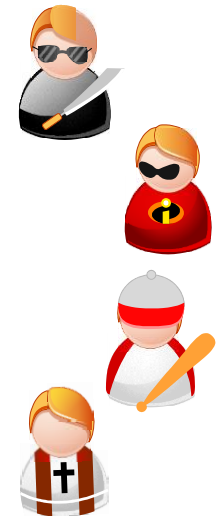
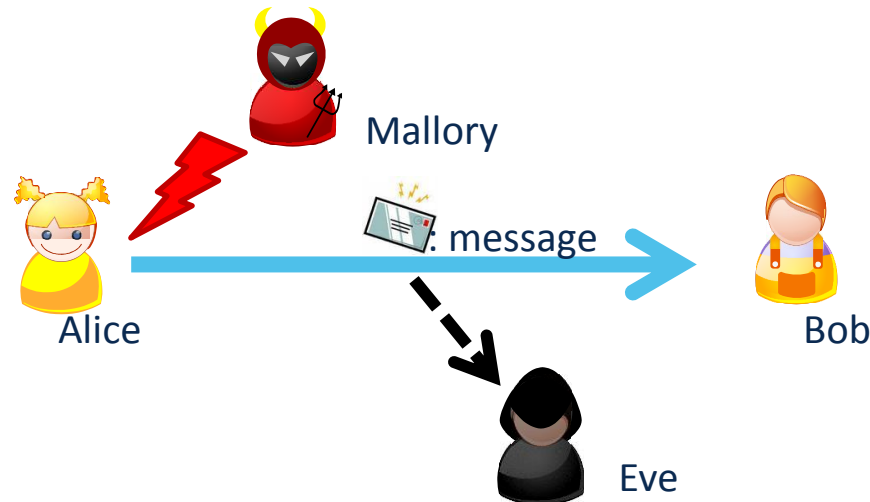
- Data transmitted or stored should only be revealed to the intended audience

- **Integrity**

- Modification of data is detected (identify source, first!)

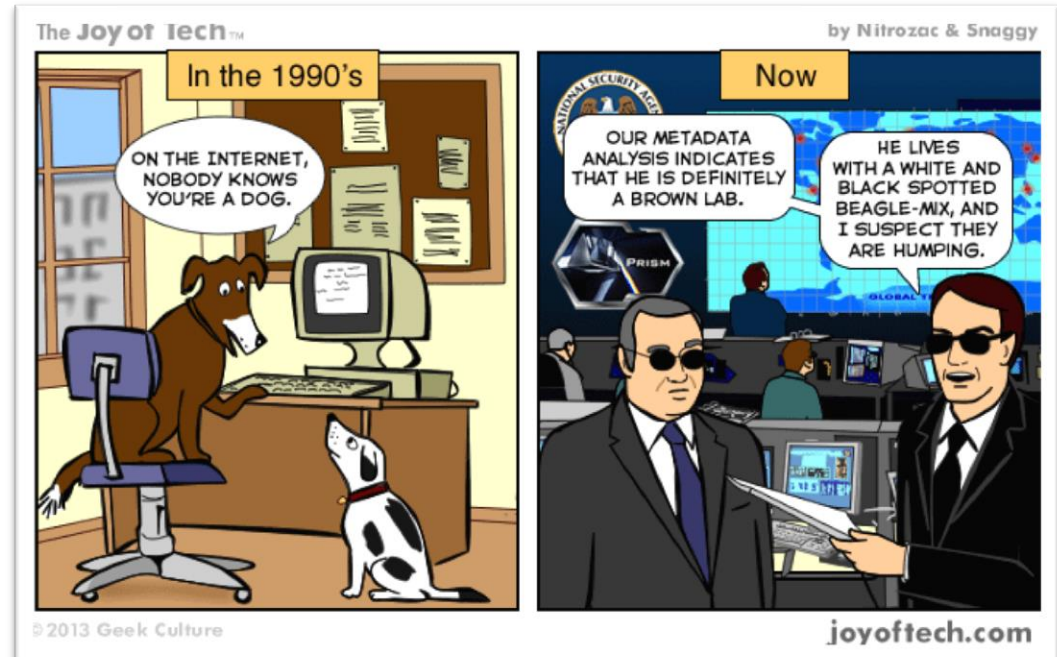
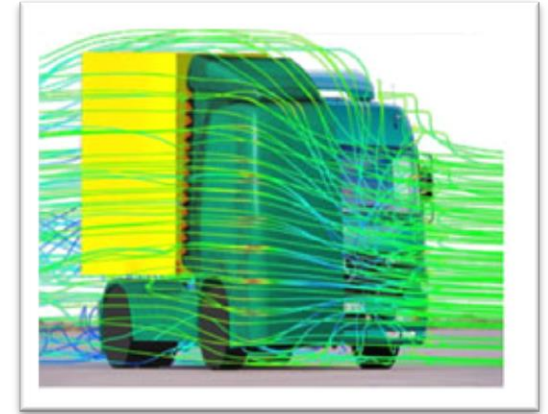
- **Availability**

- Services should function correctly upon request



- Protect data?
  - Rather: Protect integrity of individuals
  - Hence: Protect individuals FROM data
- 
- ***Hang on! What's all this „data“ about?***

- Data without any *relation* to *individuals*
  - Simulation data
  - Measurements from experiments
  
- Data *with relation to individuals*
  - Types
    - Content
    - Meta data
  - Revelation
    - Consciously
    - Unconsciously



- ***Metadata privacy***

- In controlled (opt-in!) study [1], participants
- *Called their family,...*
- *... adult establishments,*
- *... firearms dealer,*
- *... headshop, hydroponics- and hardware store,*
- *...different groups of medical specialists,*
- *...family and planned parenthood offices*

- ***Inference attacks***

- single-term lecture (students without any prior knowledge)
- Information (ab)used:
  - Partial profiles
  - Homophily
- Inferred (with high accuracy):
  - Gender
  - Age
  - Education level
  - Sexual preferences
  - Identity of anonymous profile
  - Expected tenure at employer

[1] <https://cyberlaw.stanford.edu/blog/2013/11/what%27s-in-your-metadata>

- Legally: **Personally Identifiable Information: PII**
  - **US:** Name, address (Phone, Email), national identifiers (tax, passports), IP address, driving (vehicle registration, drivers licence), biometrics (face, fingerprints), credit card numbers, date/place of birth (age, login name(s), gender, "race", grades, salary, criminal records)
  - **EU:** 'personal data' shall mean any information relating to an identified or identifiable natural person ('data subject'); an identifiable person is one who can be identified, directly or indirectly, by reference to an identification number or to one or more factors specific to his physical, physiological, mental, moral, economic, cultural or social identity [EU directive 95/46/EC]



- Samuel Warren, Louis Brandeis: **“The Right to Privacy”**, Harvard Law Review, Vol. IV, No. 5, 15<sup>th</sup> December **1890**
- **Reason:** “snapshot photography” (recent innovation at that time)
  - allowed newspapers to publish photographs of individuals without obtaining their consent.
  - private individuals were being continually injured
  - this practice weakened the “moral standards of society as a whole”
- **Consideration:**
  - basic principle of common law: individual shall have full protection in person and in property
  - “it has been found necessary from time to time to define anew the exact nature and extent of such protection”
  - “Political, social, and economic changes entail the recognition of new rights”
- **Conclusion:**
  - “right to be let alone”

- Principles
  - collect and process personal data **fairly and lawfully**
  - **purpose binding**
    - keep it only for one or more specified, explicit and lawful purposes
    - use and disclose it only in ways compatible with these purposes
  - **data minimization**
    - adequate, relevant and not excessive wrt. the purpose
    - retained no longer than necessary
  - **transparency**
    - inform who collects which data for which purposes
    - inform how the data is processed, stored, forwarded etc.
  - **user rights**
    - access to the data, correction, deletion
  - **keep the data safe and secure**

- Helen Nissenbaum: *Privacy as Contextual Integrity*, Washington Law Review, 2004
- close relation to data protection principles:
  - purpose binding
- Idea:
  - privacy violation, if:
    - violation of **Appropriateness**
      - the context „defines“ if revealing a given information is appropriate
      - **violation:** usage of information disclosed in one context in another context (even if first context is “public”)
    - violation of **Distribution**
      - the context „defines“ which information flows are appropriated
      - **violation:** inappropriate information flows



- Hang on... This only applies to data with relation to individuals!
- Again:
  - Data collection
    - requires ***informed consent*** (unless the *benefit* outweighs the *risk by far*)
    - (and provenance, make sure your subjects can act on their rights)
  - Processing of data with relation to individuals
    - Requires ***informed consent***
    - Is purpose-bound
- So may be we can remove the relation to individuals?



- Now, what does *this* mean, again?

- an-onymos <greek> (without calling the name, unnamed)
- **Attention:**
  - pseud|o|nymous <greek> (with pretense name)  
*Pseudonymized data falls under data protection laws*

- *Ref. Anne Lauber-Rönsberg:*
  - *Pseudonymized data falls under data protection laws*
  - *Anonymized data doesn't*
  
- an-onymos <greek> (without calling the name, unnamed)
  
- pseud|o|nymous <greek> (with pretense name)  
*since we already heard about it anyways*

- **Anonymity:**

- is the state of being not identifiable within a set of subjects, the ***anonymity set***.
- is the stronger, the larger the respective anonymity set is and the more evenly distributed the sending or receiving, respectively, of the subjects within that set is.

⇒ ***Quantity of Anonymity* within a particular setting depends on the number of users**

- **Unlinkability:**

- of two or more items of interest (IOIs, e.g., subjects, messages, actions, ...) from an attacker's perspective means that within the system, the attacker cannot sufficiently distinguish whether these IOIs are related or not.

⇒ ***Anonymity* in terms of *Unlinkability*:**

**Unlinkability between an identity (subject) and the IOI in question (message, data record etc.)**

*What can be disclosed?*

- Disclosure of **attributes**
  - Infer a (hidden) attribute of an individual



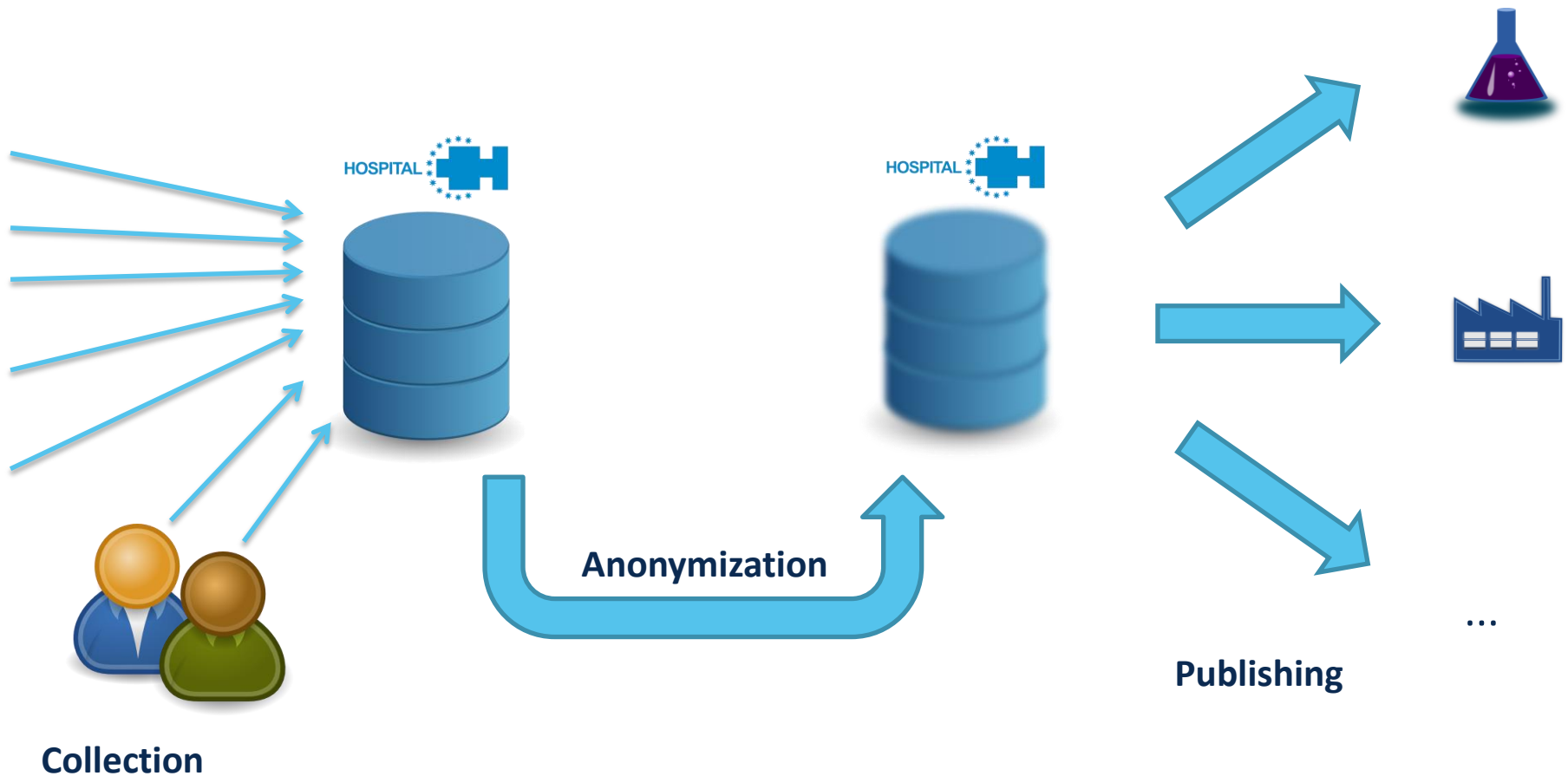
- Disclosure of **identity**
  - Identify an individual in a dataset



***Both must be prevented!***

- Goal: Obfuscate attributes
  - Identifying data (Identifier/name, PII)
  - Attributes (quasi identifiers, hidden attributes)
- Approaches
  - Encryption (hide for statistical evaluation)
    - Encrypt deterministically and delete key (frequency attacks!)
  - Perturbation (introduce error)
    - Add noise (Types: Gaussian noise, permute records, ...)
  - Generalization/Aggregation (decrease detail)
    - Suppressoin, binning
  - Modeling
    - Create model and generate synthetic data
  - (Anatomization: separate (quasi) identifiers from remaining data)



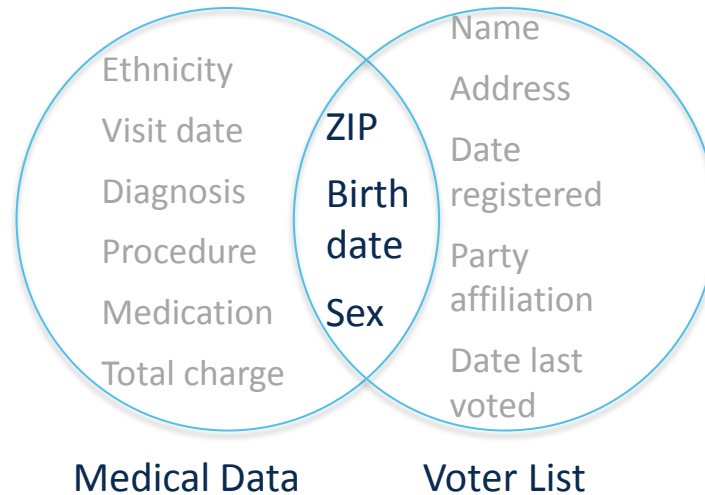


Explicit ID		Quasi ID			Sensitive		Non-sensitive	
SSN	Name	ZIP	Age	Sex	Disease	Salary	Q1	Q2
309-10-2346	Bob	47677	43	Male	Heart	3.000	a1	13
306-30-2349	Alice	47602	22	Female	Flu	5.000	a5	4
306-31-6548	Carol	47678	45	Female	Hepatitis	6.000	a4	22
309-80-2988	Dave	47905	31	Male	HIV	4.000	a1	12
316-11-9832	Marvin	47909	36	Male	Flu	10.000	a2	8

- Explicit identifiers must be removed
- Link between **Quasi-IDs** and sensitive attributes needs to be obfuscated

	Quasi ID			Sensitive		Non-sensitive	
	ZIP	Age	Sex	Disease	Salary	Q1	Q2
	47677	43	Male	Heart	3.000	a1	13
	47602	22	Female	Flu	5.000	a5	4
	47678	45	Female	Hepatitis	6.000	a4	22
	47905	31	Male	HIV	4.000	a1	12
	47909	36	Male	Flu	10.000	a2	8

- Explicit identifiers must be removed
- Link between **Quasi-IDs** and sensitive attributes needs to be obfuscated



- Re-identification through directly linking shared attributes
- 87% of US population show characteristics to be uniquely identifiable through {ZIP, Date of birth, Sex} (Census 1990)

	Quasi ID			Sensitive		Non-sensitive	
	ZIP	Age	Sex	Disease	Salary	Q1	Q2
	47677	43	Male	Heart	3.000	a1	13
	47602	22	Female	Flu	5.000	a5	4
	47678	45	Female	Hepatitis	6.000	a4	22
	47905	31	Male	HIV	4.000	a1	12
	47909	36	Male	Flu	10.000	a2	8

- Explicit identifiers must be removed
- Link between Quasi-IDs and sensitive attributes needs to be obfuscated
  - Generalization & Suppression
  - Anatomization & Permutation
  - Perturbation

	ZIP Code	Age	Disease
1	47677	29	Heart Disease
2	47602	22	Heart Disease
3	47678	27	Heart Disease
4	47905	43	Flu
5	47909	52	Heart Disease
6	47906	47	Cancer

$k=3$

	ZIP Code	Age	Disease
1	476**	2*	Heart Disease
2	476**	2*	Heart Disease
3	476**	2*	Heart Disease
4	4790*	≥40	Flu
5	4790*	≥40	Heart Disease
6	4790*	≥40	Cancer

- Groups of  $k$  records → resulting in  $k$ -anonymous table
- Probability  $1/k$  to link correct entry to known quasi-identifier
- Tradeoff between privacy and utility
  - larger groups normally result in less accurate data
- **Problem: Homogeneity in sensitive attributes**
  - Solution:  **$l$ -diversity** → at least  $l$  different values for each sensitive attribute in each equivalence class
  - **Problem:** meaning of “different”: different kinds of cancer → cancer
    - Solution:  **$t$ -closeness** (etc, etc)

- hide communication „meta-data“ (circumstances):
  - who is communicating with whom
  - how long, how often, how much data etc.
  - location
- terms:
  - sender anonymity, recipient anonymity
- existing solutions:
  - AN.ON ([anon.inf.tu-dresden.de](http://anon.inf.tu-dresden.de))
  - Tor ([tor.eff.org](http://tor.eff.org))

- OK, so let's anonymize!

Anonymization: remove all identifying attributes.

⇒ Linguists, anybody? Remove **ALL** content... (sorry!)

⇒ Social scientists? Sorry, same holds for you...

- ***So then, can I use the meta data on anonymized data?***



- AOL 2006
  - In Aug 2006, AOL published anonymized collection of search queries
  - Goal: help scientists understand what's happening
  - On Aug 9, 62 year old Thelma Arnold was identified
- Netflix 2006
  - In Oct 2006, anonymized data set of movie preferences released
  - Goal: challenge for better recommender systems!
  - On Oct 18, 2006, Narayan & Shmatikov published de-anonymization
- Location privacy
  - Location traces are very rare (Orange D4D, Nokia mobile data challenge)
  - Even more so since 2013:
  - Vincent Blondel et al.: 4 points (time & gps) identify 95% of individuals in 15 month dataset of 1.5 Mio people

*“Companies do not make money by giving researchers access to data. They do it to promote and encourage research in the field. Based on the AOL and Netflix incidents, I suspect that we will see a major chill hit the industry.*

*No high-tech company with large amounts of user data will ever again risk making it available to researchers without **first requiring them to sign a lengthy contract**. The risk of the data being de-anonymized (and the resulting public relations and legal trouble) is simply not worth it.”*

*-- Chris Soghoian, C|Net, 2007*

- Doing research on social media
  - Requires ethical considerations
  - **Informed consent** for data acquisition (and processing)
  - Anonymization before further processing/sharing!
- Make sure:
  - Clarify which data **really** is involved
  - What can be derived from this data
  - Take concrete measures to avoid or at least remove **as much as possible**
  - Reassess which effect these means **really** have
  - (Ask your local privacy/ethics expert)
  - Ponder if you share your data (you may have to? Do you have consent?)
- And btw, yes, the concept(s) of privacy are a little bit hard to grasp...

- Warren, Samuel, Brandeis, Louis. "The Right to Privacy", Harvard Law Review, Vol. IV, No. 5, 1890
- Cutillo, Leucio Antonio, et al. "Security and privacy in online social networks." Social Network Technologies and Applications. Springer US, 2010.
- Lauber-Rönsberg, Anne: "Research Ethics and Data Protection Laws". Online
- Narayanan, Arvind, and Vitaly Shmatikov. "How to break anonymity of the netflix prize dataset." arXiv preprint cs/0610105 (2006).
- Narayanan, Arvind, and Vitaly Shmatikov. "Robust de-anonymization of large sparse datasets." 2008 IEEE Symposium on Security and Privacy (sp 2008). IEEE, 2008.
- Nissenbaum, Helen. "*Privacy as Contextual Integrity*", Washington Law Review, 2004
- Pfitzmann, Andreas, and Hansen, Marit: "A terminology for talking about privacy by data minimization." Online: [https://dud.inf.tu-dresden.de/literatur/Anon\\_Terminology\\_v0.34.pdf](https://dud.inf.tu-dresden.de/literatur/Anon_Terminology_v0.34.pdf)
- Scharloth, Joachim. "Research Ethics: Principles and New Challenges". Online: [http://scharloth.com/slides/research\\_ethics/folie\\_19.html](http://scharloth.com/slides/research_ethics/folie_19.html)
- *All pictures credit wikimedia, unless stated differently*