



Anonymisierung und Pseudonymisierung aus Sicht der Informatik

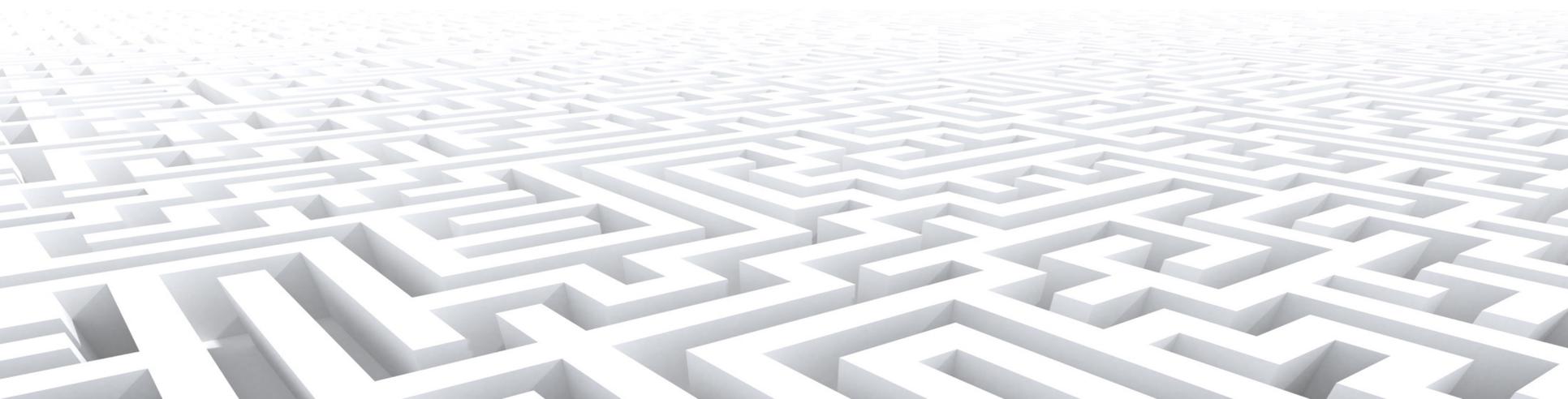
Prof. Dr. Hannes Federrath

Sicherheit in verteilten Systemen (SVS)

<http://svs.informatik.uni-hamburg.de>

Anonymisierung und Pseudonymisierung aus Sicht der Informatik

- Gliederung des Vortrags
 - Einordnung in den rechtlichen Kontext
 - Anonymisierung und Pseudonymisierung von Datensätzen
 - Grenzen der Leistungsfähigkeit von technischen Anonymisierungsverfahren
 - Schlussbemerkungen

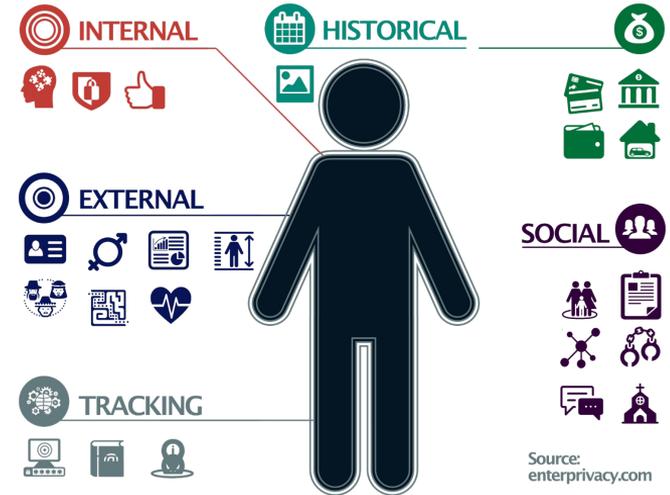


Personenbezogene Daten

Def. Personenbezogene Daten nach Art. 4 Nr. 1 DSGVO

»Personenbezogene Daten sind alle Informationen, die sich auf eine identifizierte oder identifizierbare natürliche Person (im Folgenden ›betroffene Person‹) beziehen.

Als identifizierbar wird eine natürliche Person angesehen, die direkt oder indirekt, insbesondere mittels Zuordnung zu einer Kennung wie einem Namen, zu einer Kennnummer, zu Standortdaten, zu einer Online-Kennung oder zu einem oder mehreren besonderen Merkmalen, die Ausdruck der physischen, physiologischen, genetischen, psychischen, wirtschaftlichen, kulturellen oder sozialen Identität dieser natürlichen Person sind, identifiziert werden kann.«



Name, Vorname	Telefonnummer		
Kontonummer	Geburtsdatum		
Kreditkarten-Nr.	E-Mail-Adresse	Religionszugehörigkeit	Geschlecht
Wohnort	Krankenversicherungs-Nr.	Steuernummer	Autokennzeichen

Indirekt identifizierbare natürliche Person

Daten sind auch dann personenbezogen, wenn die Zuordnung zu einer Person nur indirekt vorgenommen werden kann.

- entscheidend ist allein die (realistische) Möglichkeit einer Zuordnung
 - Zuordnung über Gruppenzugehörigkeit genügt bereits
 - Zuordnung auf Umwegen (Zusatzwissen, Big Data) ist ebenfalls eine realistische Möglichkeit
- Urteil des EuGH v. 19.10.2016, Az. C-582/14
- Es genügt auch, über rechtliche Mittel zu verfügen, die eine verantwortliche Stelle in die Lage versetzen, die Zuordnung zu einer Person vorzunehmen.
 - Hintergrund war: Personenbezug von dynamischen IP-Adressen für Webseitenbetreiber

Indirekt identifizierbare natürliche Person

■ Beispiel: Personenbezug von IP-Adressen

- entscheidend ist der Aufwand zur Herstellung des Personenbezugs
- wenn fehlender Personenbezug:
 - keine Anwendbarkeit des Datenschutzrechts
 - dürften unbeschränkt gespeichert und übermittelt werden
- Zulässigkeit der Speicherung: Zur Aufrechterhaltung des Dienstes
 - 14 Tage sind zu lang und damit nicht angemessen

■ Urteil des EuGH v. 19.10.2016, Az. C-582/14

- Es genügt auch, über rechtliche Mittel zu verfügen, die eine verantwortliche Stelle in die Lage versetzen, die Zuordnung zu einer Person vorzunehmen.
 - Hintergrund war: Personenbezug von dynamischen IP-Adressen für Webseitenbetreiber

- Auszug aus Erwägungsgrund 26 der DSGVO:

»Um festzustellen, ob eine natürliche Person identifizierbar ist, sollten alle Mittel berücksichtigt werden, die von dem Verantwortlichen oder einer anderen Person nach allgemeinem Ermessen wahrscheinlich genutzt werden, um die natürliche Person direkt oder indirekt zu identifizieren, wie beispielsweise das Aussondern.

Bei der Feststellung, ob Mittel nach allgemeinem Ermessen wahrscheinlich zur Identifizierung der natürlichen Person genutzt werden, sollten alle objektiven Faktoren, wie die Kosten der Identifizierung und der dafür erforderliche Zeitaufwand, herangezogen werden, wobei die zum Zeitpunkt der Verarbeitung verfügbare Technologie und technologische Entwicklungen zu berücksichtigen sind. ...«

Anonyme Daten

- **Auszug aus Erwägungsgrund 26 der DSGVO:**

»Die Grundsätze des Datenschutzes sollten [daher] nicht für anonyme Informationen gelten, d. h. für Informationen, die sich nicht auf eine identifizierte oder identifizierbare natürliche Person beziehen, oder personenbezogene Daten, die in einer Weise anonymisiert worden sind, dass die betroffene Person nicht oder nicht mehr identifiziert werden kann. Diese Verordnung betrifft somit nicht die Verarbeitung solcher anonymer Daten, auch für statistische oder für Forschungszwecke.«

- **Absolute Anonymität**

- Daten können unter keinen Umständen mehr zugeordnet werden

- **Faktische Anonymität**

- Einzelangaben können nur mit einem unverhältnismäßig großen Aufwand an Zeit, Kosten und Arbeitskraft zugeordnet werden

- Datenschutzrecht grundsätzlich nicht anwendbar
- Durch Leistungssteigerung von IT-Verfahren können faktisch anonyme Daten wieder personenbezogen werden!

Pseudonymisierte Daten

- Art. 4 Nr. 5 DSGVO Pseudonymisierung

»Verarbeitung personenbezogener Daten in einer Weise, dass die personenbezogenen Daten ohne Hinzuziehung zusätzlicher Informationen nicht mehr einer spezifischen betroffenen Person zugeordnet werden können, sofern diese zusätzlichen Informationen gesondert aufbewahrt werden und technischen und organisatorischen Maßnahmen unterliegen, die gewährleisten, dass die personenbezogenen Daten nicht einer identifizierten oder identifizierbaren natürlichen Person zugewiesen werden.«

Maßnahmen, die darauf abzielen, eine direkte Zuordnung (ohne Kenntnis einer Zuordnungsregel) zu einem Betroffenen zu unterbinden, ohne dabei in jedem Fall den Personenbezug dabei völlig aufzuheben

Anonymisierung und Pseudonymisierung von Datensätzen

Pseudonymisierte Daten

- »Einer Pseudonymisierung unterzogene personenbezogene Daten, die durch Heranziehung zusätzlicher Informationen einer natürlichen Person zugeordnet werden könnten, sollten als Informationen über eine identifizierbare natürliche Person betrachtet werden.« (aus: ErwG 26)

- Getrennte Datenhaltung
- Explizite Zuordnungsregel vorhanden

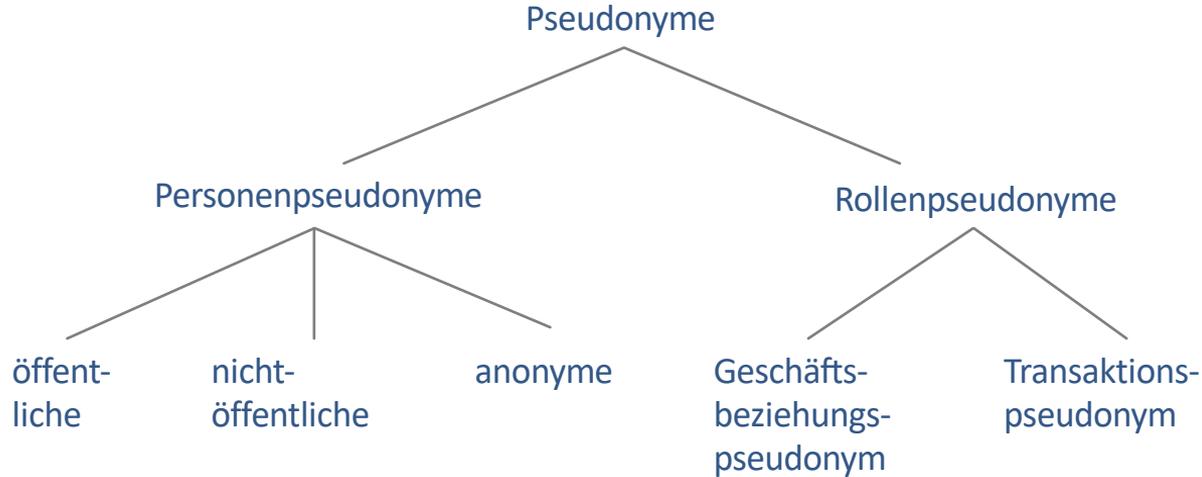
Name	P	P	Sex	Diagnosis
Alice	735	735	f	A
Bob	324	324	m	B
Carol	478	478	f	C
Dan	125	125	m	D

- Auch Merkmalskombinationen stellen ein Pseudonym dar, wenn sie Eindeutigkeit erzielen.

- Getrennte Datenhaltung
- Fehlen einer expliziten Zuordnungsregel
- Quasi-Identifiers (QIDs)

Name	Sex	Birth date	ZIP	Sex	Birth date	ZIP	Diagnosis
Alice	f	1953-06-11			1953-06-11	12345	A
Bob	m	1922-03-02			1922-03-02	98763	B
Carol	f	1973-05-20			1973-05-20	12390	C
Dan	m	1966-10-13			1966-10-13	98764	D

The diagram shows two overlapping circles. The left circle is labeled 'Name' and the right circle is labeled 'Diagnosis'. The intersection of the two circles is labeled 'Sex, Birth date, ZIP'.



Beispiele für Pseudonyme:

**Telefon-
nummer, E-
Mail-
Adresse**

**Konto-
nummer,
IP-Adresse**

**Biometrische Merkmale
(solange kein Register)**

**Künstlernername,
Nickname**

**Kennwort,
Zufallszahl**

Gute Skalierbarkeit bezüglich der Anonymität

Anonymisierte Daten

- **Vermeintlich einfache Umsetzbarkeit von Anonymisierung**
 - Entfernen der Datenfelder mit Personenbezug
- **Probleme**
 - fehlende Entscheidbarkeit, ob verbleibende Felder Quasi-Identifiers (QIDs)
 - Kontextwissen eines Angreifers ist zum Zeitpunkt der Anonymisierung unbekannt
- **Metriken zur Messung des Anonymisierungsgrades**
 - k-Anonymität (Sweeney, 2002)
 - Für alle Datensätze gilt, dass nach der Anonymisierung wenigstens k Datensätze nicht mehr voneinander unterscheidbar sind.
 - Differential Privacy (Dwork, 2006)
 - Für alle Datensätze, die sich in höchstens einem Eintrag unterscheiden, ist die Wahrscheinlichkeit kleiner als ein vorgegebener Wert e^ϵ , dass diese nach der Anonymisierung noch unterscheidbar sind.
 - ... (zahlreiche weitere)

Generalisieren
Aggregieren

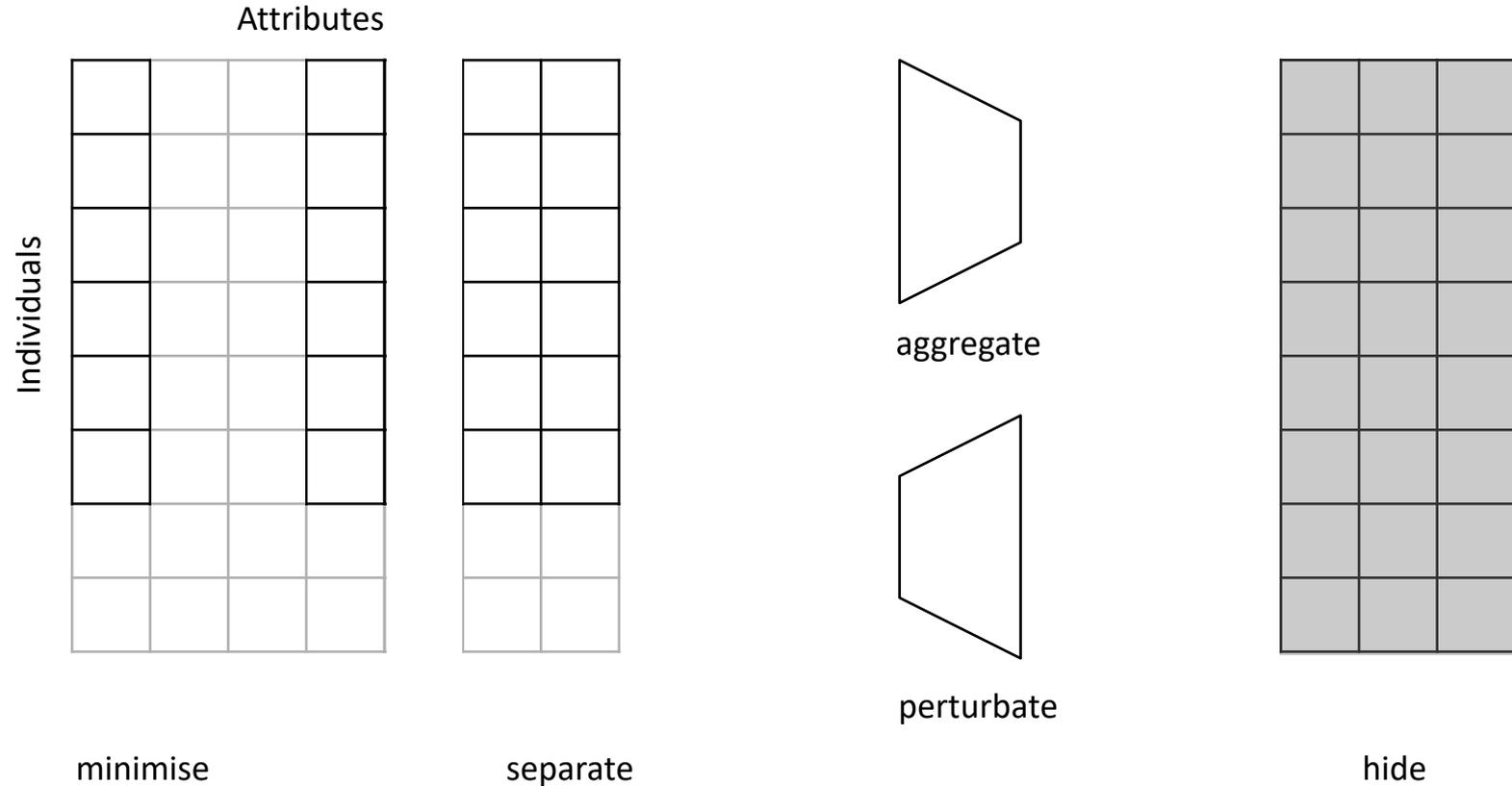
Perturbieren
Verrauschen

- Jeder Datensatz ist bezüglich einer vorgegebenen Kombination von QIDs ununterscheidbar von wenigstens $k-1$ anderen Datensätzen.

Name	Sex	ZIP	Diagnosis
Alice	f	12345	A
Bob	m	98763	B
Carol	f	12390	C
Dan	m	98764	D
Eve	f	98765	E

$k = 2$

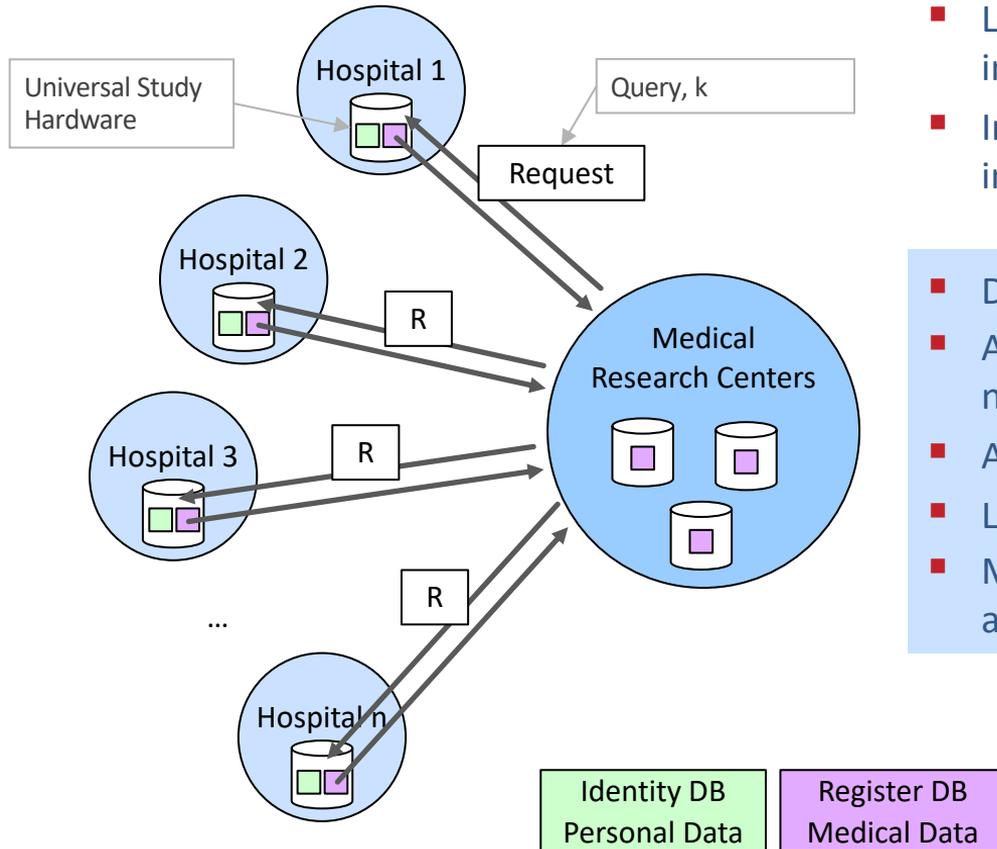
Name	Sex	ZIP	Diagnosis
	f	123**	A
	*	9876*	B
	f	123**	C
	*	9876*	D
	*	9876*	E



- Technisch
 - **Minimise**: Nur notwendige Daten speichern und verarbeiten
 - **Separate**: Daten verteilt verarbeiten und speichern
 - **Aggregate**: Daten auf das notwendige Maß zusammenfassen
 - **Perturbate**: Daten durch zufällige Störungen ungenau machen
 - **Hide**: Daten nicht in offener Form speichern

- Organisatorisch
 - **Enforce**: Durchsetzung einer Datenschutz-Policy (access control)
 - **Inform**: Betroffene über Datenverwendung informieren (P3P)
 - **Control**: Eingriffsmöglichkeit der Betroffenen (informed consent)
 - **Demonstrate**: Überprüfbarkeit (privacy management, logging)

Distributed calculation of anonymity level



- Large number of Centers (Hospitals) allows an increased number of cases into the Research Data
 - Informed Consent as a legal requirement for the inclusion of Study Participants (Patients)
- Data records remain in Centers
 - Approach to distributed gathering of research data needed for a concrete research question
 - Adaptive k-anonymous response from centers
 - Limited number of (similar) requests
 - Makes use of privacy-respecting federated learning and secure multi-party computation (SMPC)

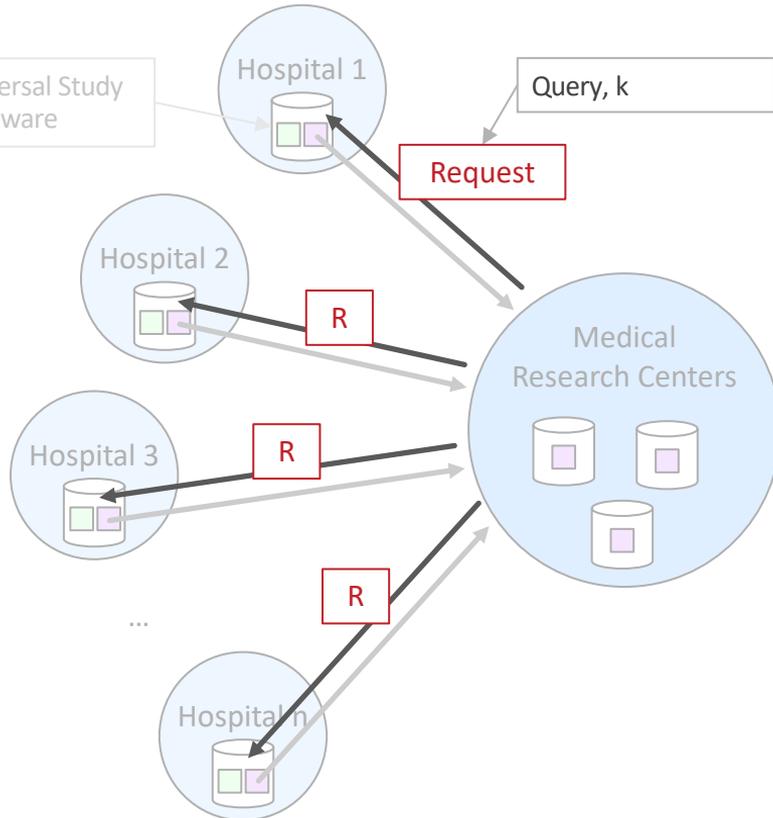
Agent k creator

- Central definition of requested variables, their range and their granularity
- Type of query (also repetitive query) with pre-defined k or patient number

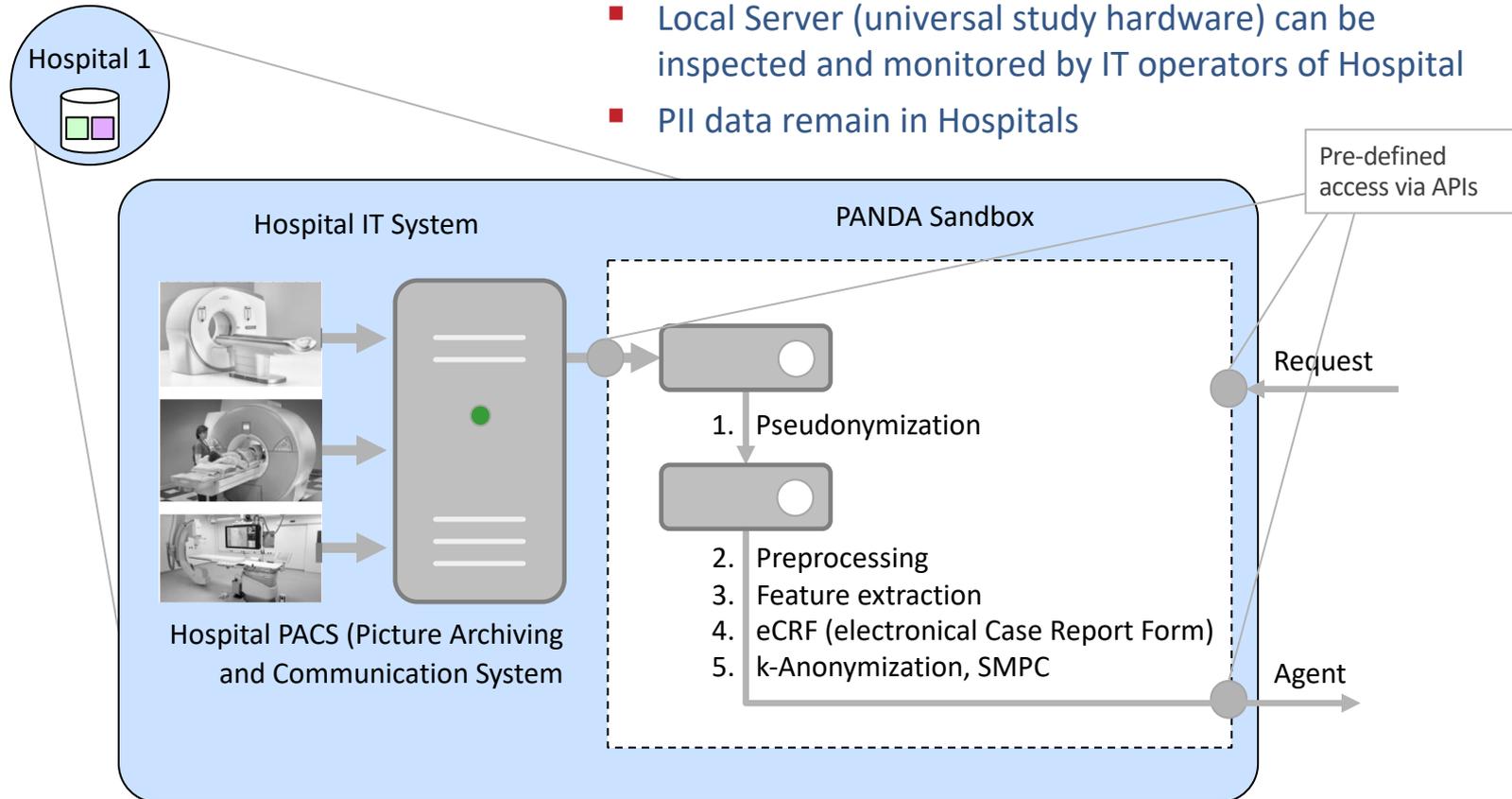
	Range	Granularity
Device	Solitaire Stent	4 x 40 mm only
Age	0 18 70	10
Symptoms (NIHSS-Score)	0 8 20	2
Infarct volume (ASPECT-Score)	7 10	2

K 3

Patients 297



Distributed calculation of anonymity level



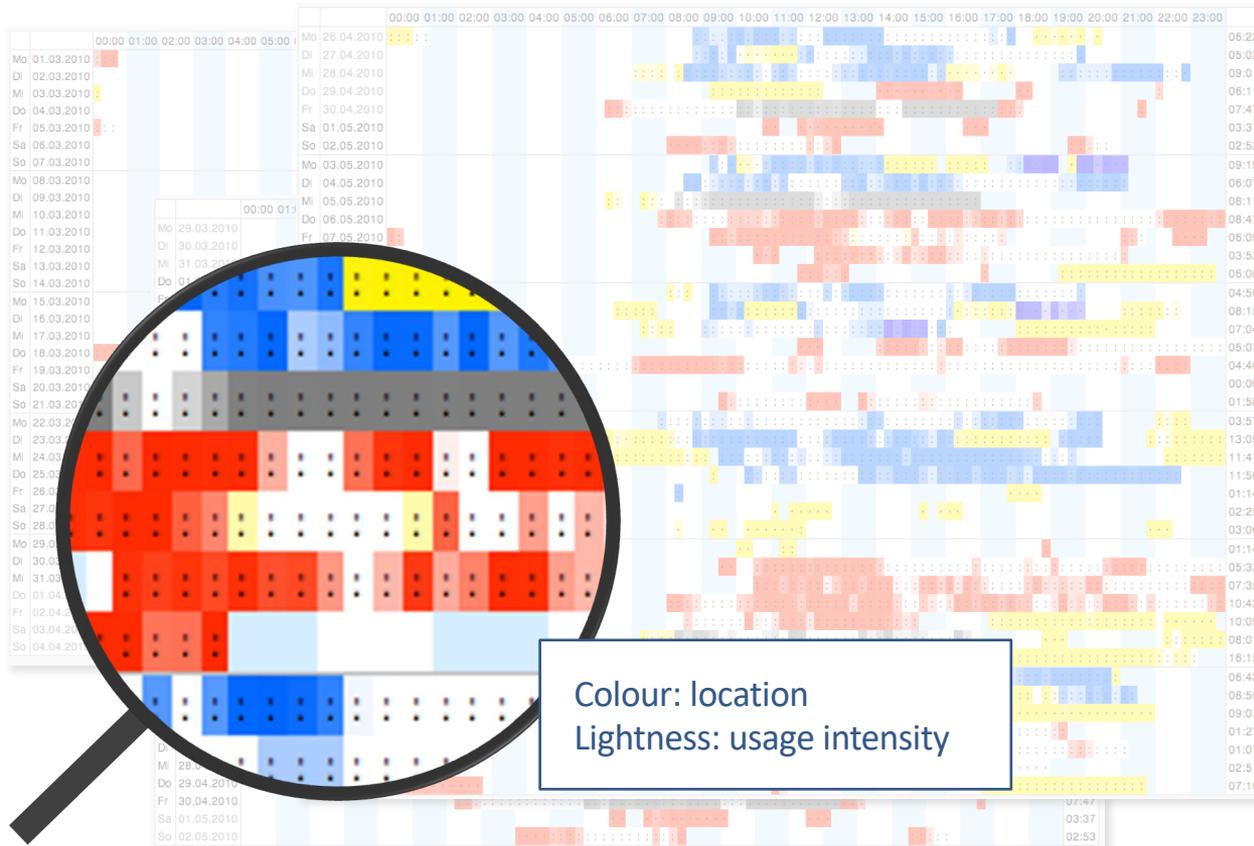
Grenzen des Schutzes durch Anonymisierung und Pseudonymisierung

Grenzen des Schutzes

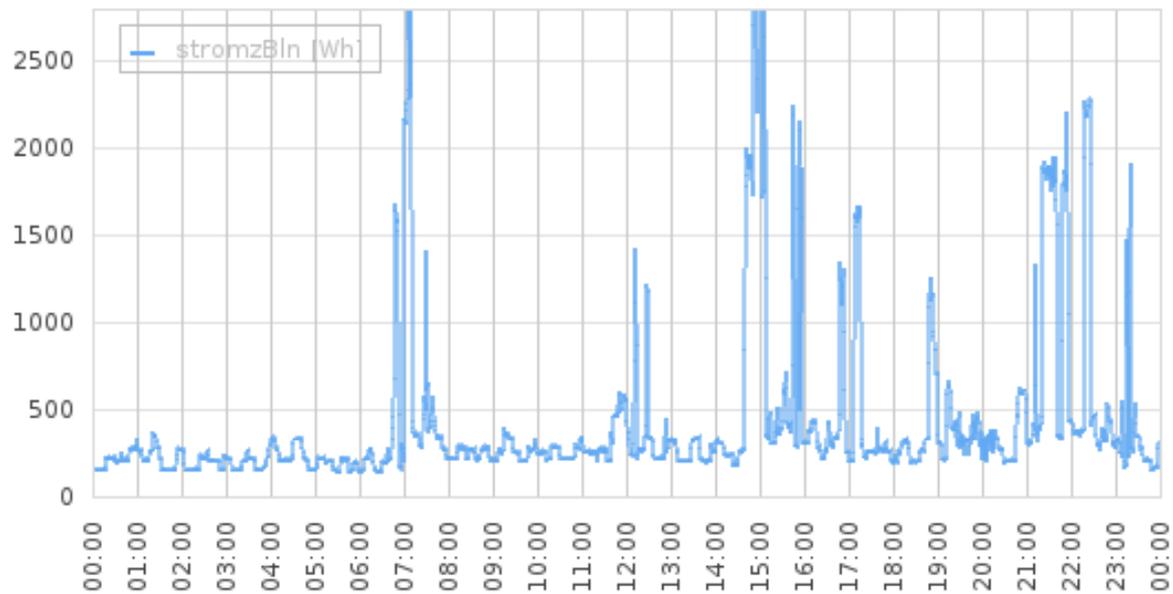
- Datenverarbeitung sollte mit dem klaren Bekenntnis betrieben werden, dass ein *perfekter Schutz niemals möglich ist*.
- Beispiele:
 - eigene Untersuchungen zur Kontrolle und Überwachung der Arbeitsleistung von Softwareentwicklern (nächtliches Arbeiten, Produktivitätseinbrüche, Zusammenarbeitsstrukturen im Team) und von Daten aus privaten Haushalten legen Lebensgewohnheiten offen
 - Der Fall »Strava Heatmap« (2018)
 - Auch Gruppeneigenschaften sind identifizierbare personenbezogene Daten
 - Analyse von 170 Mio. Open-Data-Datensätzen der Taxifahrten der NY Taxigesellschaft (2014)
 - Adressen der Interessenten von Nachtclubbesuchen wurden bekannt
 - KI-Systeme und Entfernung von identifizierenden Merkmalen

Realität hat gezeigt, dass es immer wieder möglich war, Personenbezug herzustellen (zu re-identifizieren), weil die Möglichkeiten der Anonymisierung und Pseudonymisierung zu blauäugig angewendet wurden

Usage timelines and ip-geo-tagging



Stromverbrauch verrät Infos über persönliche Lebensverhältnisse



Soziale Netze und Datenschutz – Der Fall »Strava Heatmap«

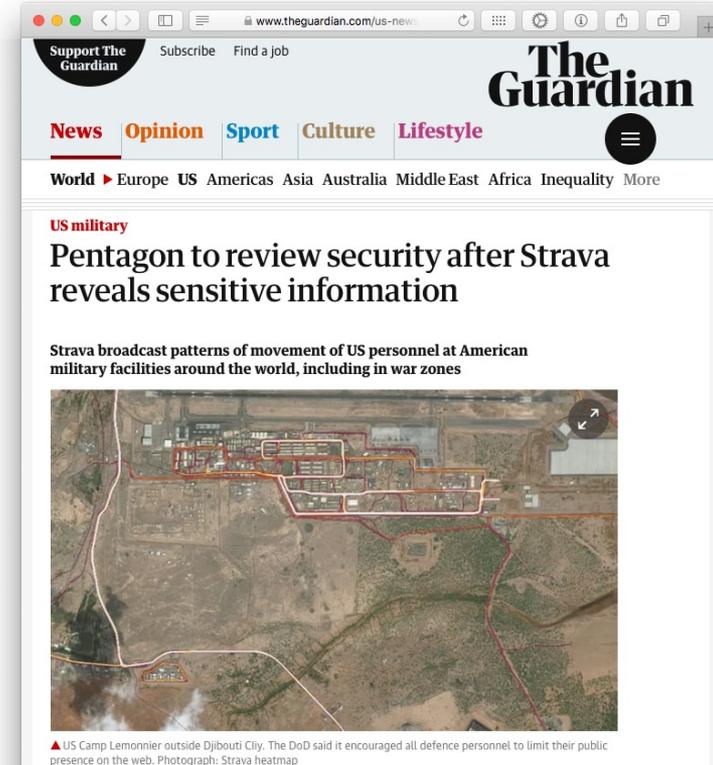
- Fitness-Tracker Website veröffentlicht beliebte Laufstrecken
- Soldaten des US-Militärs offiziell ausgestattet mit Fitness-Tracker
- Öffentliche »Heatmap« enthüllt ungewollt geheime US-Bases im Ausland



Sources:

<https://twitter.com/Nrg8000/status/957318498102865920>

<https://www.theguardian.com/us-news/2018/jan/29/pentagon-strava-fitness-security-us-military>



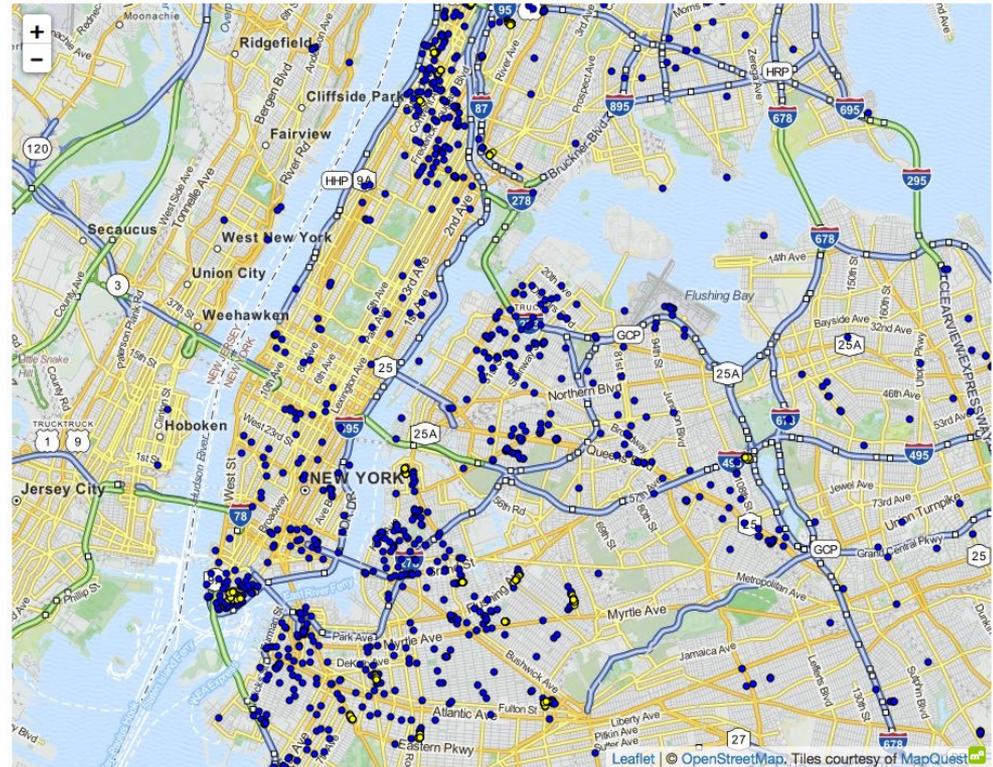
Pseudonymisierte Daten...

...können Persönlichkeitsrechte verletzen.

20 GByte of pseudonymisierter Daten von
170 Mio. Taxifahrten der New Yorker Taxi-
gesellschaft

Daten öffentlich abrufbar unter:

<http://www.andresmh.com/nyctaxitrips/>



Drop-off locations for trips starting at Larry Flynt's Hustler Club between
midnight and 6 am during 2013.

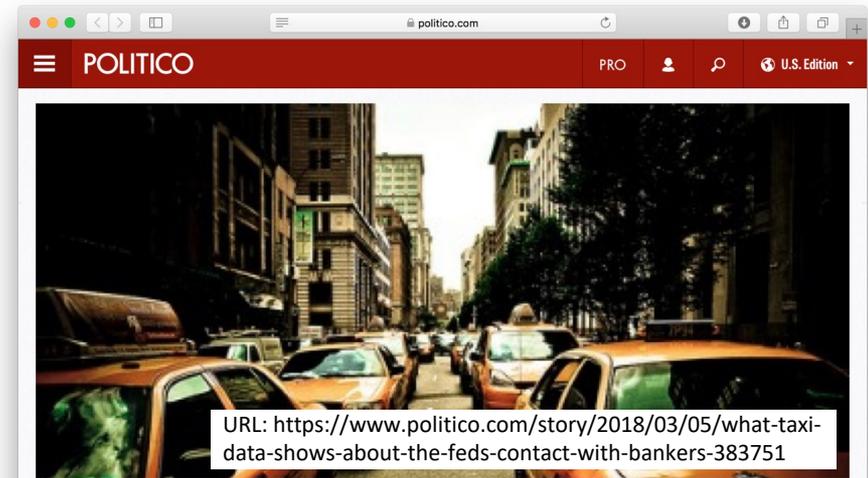
Source: <http://content.research.neustar.biz/blog/differential-privacy/stipRaw.html> (2014)

Anonymisierte Daten...

...können Geheimes verraten.



David Andrew Finer: What Insights Do Taxi Rides Offer into Federal Reserve Leakage? Working Paper, Booth School of Business, University of Chicago, March 2018. <https://research.chicagobooth.edu/-/media/research/stigler/pdfs/workingpapers/18whatinsightsdotaxiridesofferintofederalreserveleakage.pdf>



What taxi data shows about the Fed's contact with bankers

By VICTORIA GUIDA | 03/05/2018 12:59 PM EST

Share on Facebook Share on Twitter

Contact between the Federal Reserve Bank of New York and six of the largest U.S. banks seems to increase around the Fed's key interest-rate-setting meetings, according to a new academic study that raises questions about whether this could give top Wall Street institutions an unfair competitive edge.

The study, from the University of Chicago's Booth School of Business, examines granular data on cab rides released by New York's taxi regulator. It singles out lunchtime trips between the New York Fed and the major offices of six banks: Goldman Sachs, Citigroup, JPMorgan Chase, Morgan Stanley, Bank of New York Mellon and Bank of America.

It also includes potential off-site meetups by factoring in "coincidental drop-offs," where someone from the New York Fed and someone from a bank each appeared to be dropped off at the same location around the same time.

Lunchtime coincidental drop-offs happened about 50 percent more often between when an important meeting of the policy-making Federal Open Market Committee started through the following week, according to the study written by David Andrew Finer. That's an average of about 1.2 more taxi rides per meeting.

"I cannot conclusively demonstrate a link between rides and face-to-face meetings, but evidence that individuals are in very close proximity to each other more often around FOMC meetings would complement more indirect evidence of regular informal communication," the study says. It examines taxi data between 2009, when data was first made available, and 2014, before ride shares like Uber and Lyft became popular.

- **Deskriptive Big-Data-Analytik**
 - zur Auswertung, Sichtung und Aufbereitung von Daten; Beispiele:
 - Data Mining
 - Filterung, Klassifizierung und Priorisierung von Daten
- **Prädiktive Big-Data-Analytik**
 - Suche nach Indikatoren für einen möglichen Kausalzusammenhang
 - Einsichten in das Verhalten von Menschen
 - Trends und Verhaltensmuster zur Vorhersage künftigen Verhaltens
- **Präskriptive Big-Data-Analytik**
 - zur Erreichung bestimmter Ziele
 - personalisierte Selektion bei der Preisgestaltung
 - Beeinflussung öffentlicher Meinungsbildung
 - Einwirkung auf gesellschaftliche Entwicklungen

Prädiktive Big-Data-Analytik: Stecknadeln im Heuhaufen



- Personenbezogene Daten sind auch Daten, die als Ergebnis einer Big-Data-Analyse entstehen.
 - allgemein und ohne Herleitung aus Daten speziell der konkret betroffenen Person
 - Beispiele: Person wohnt in einem bestimmten Stadtteil; daraus Ableitung von Finanzkraft, Herkunft, sexueller Orientierung, Gesundheit
- Personenbezogene Daten sind auch Daten, deren Personenbezug durch Anonymisierung entfällt.
 - Möglichkeiten der Deanononymisierung und Ableitung von Eigenschaften dürften nicht unterschätzt werden
 - Beispiele: New York Taxi Data Analytics, Strava Heatmap
- ebenso kritisch pseudonymisierte, aggregierte, perturbierte, verschlüsselte Daten betrachten



Extraktion von Eigenschaften aus Machine-Learning-Modellen

- These: personenbezogene Daten in Machine-Learning-Modellen sind aufgrund von Aggregation und damit verbundenem Informationsverlust nicht mehr rückführbar auf einzelne Individuen
 - Beispiel: Daten für das Anlernen KI-Systemen von stammen aus verschiedenen Quellen mit einheitlichen Merkmalen
 - Model M_A aus Quelle A: Trainingsdaten bestehen überwiegend aus Männern (70%)
 - Model M_B aus Quelle B: Trainingsdaten enthalten überwiegend Frauen (70%)
 - Angriffsziel: Separieren von M_A und M_B
- Zur Erinnerung: Daten sind auch dann personenbezogen, wenn die Zuordnung zu einer Person nur indirekt vorgenommen werden kann.
 - Zuordnung über Gruppenzugehörigkeit genügt bereits
- Aktuelle Forschungsarbeiten zum sog. Property unlearning

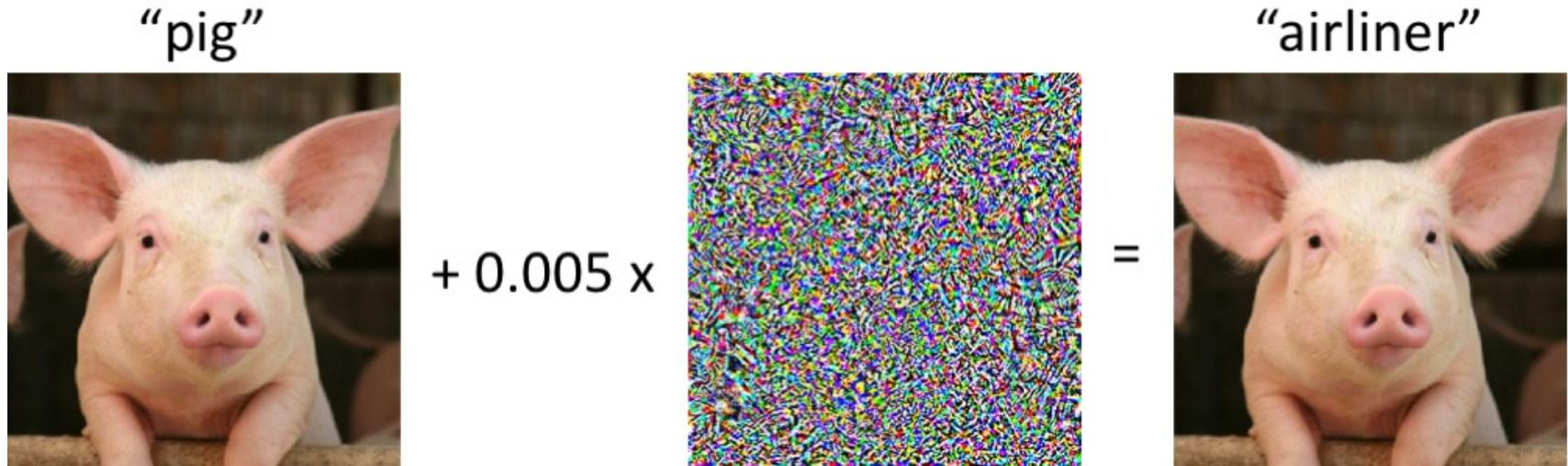


Shokri, Reza, et al. Membership inference attacks against machine learning models. IEEE Symposium on Security and Privacy (SP) 2017. IEEE, 2017.

Zugespitzte Frage: War ein bestimmtes Individuum in das Training Data Set eingeschlossen?

Verbergen von Eigenschaften beim Klassifizieren

- Adversarial learning



<https://www.designnews.com/electronics-test/yes-ai-can-be-tricked-and-its-serious-problem/161652909959780>

inf.uni-hamburg.de

 **Universität Hamburg**
DER FORSCHUNG | DER LEHRE | DER BILDUNG

DEPARTMENT OF INFORMATICS
SECURITY AND PRIVACY

[HOME](#) [COURSES](#) [THESES](#) [RESEARCH](#) [PEOPLE](#) [SERVICE](#) 



Foto: UHH/Denstorf

🏠 UHH → MIN-Fakultät → Fachbereich Informatik → Einrichtungen → Arbeitsbereiche → Security and Privacy → Home

WORKING GROUP ON «SECURITY AND PRIVACY»

Security and Privacy

Information systems become more and more important in critical infrastructures, while the Internet has evolved to a critical infrastructure itself. The secure operation of these infrastructures is vital and their failure can have severe impacts up to the loss of human lives.

Security refers to the fact that protection goals are achieved in the presence of malicious attacks and system failures. Typical security goals can be confidentiality, integrity, accountability, and availability. Security and privacy in information systems addresses both technical and organizational aspects, such as building and establishing security concepts and security infrastructures as well as risk analysis and risk management.

Privacy can be a conflicting goal to security, but they can also benefit from each other. Hence, it is necessary to balance both when developing secure information systems.

Prof. Dr. Hannes Federrath
Fachbereich Informatik
Universität Hamburg
Vogt-Kölln-Straße 30
D-22527 Hamburg

Telefon +49 40 42883 2358

hannes.federrath@uni-hamburg.de

<https://svs.informatik.uni-hamburg.de>