

As presented at NIKSUN WWSMC

July 25-27, 2011 | [www.niksun.com](http://www.niksun.com)

# A Theory of Privacy and Utility for Data Sources

Lalitha Sankar  
Princeton University

# Electronic Data Repositories

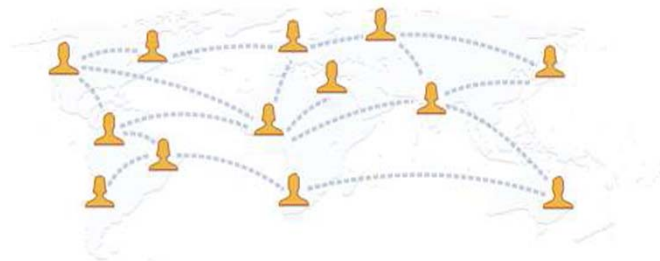
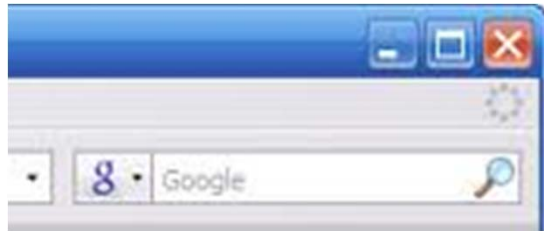
- Technological leaps in information processing, storage, and communications has led to the creation of vast electronic data repositories.



By simply clicking on a **blue button** icon, users will be able to download their medical (Medicare/Medicaid) data to their personal computers. – (PubMed Central)

# The Privacy Problem

- Many electronic information sources are **publicly accessible**
  - Google, Facebook, open governance, census, etc.



- **These electronic information sources can also leak private information!**

# Utility vs. Privacy

- **Utility** (benefit) of data repositories is in allowing legitimate users access to statistical/processed data.
  - e.g., census data
- However, individual information needs to be kept **private**
  - Private information (e.g., SSN, DoB) can be potentially inferred from revealed data.
- Private information is **application-specific**
  - DoB is private for medical but not DMV databases.
  - Census publications may not reveal name, SSN, DoB, address, tel. no. of any individual.
- **Need a framework that precisely quantifies the utility-privacy tradeoffs for any application.**



# Talk Outline

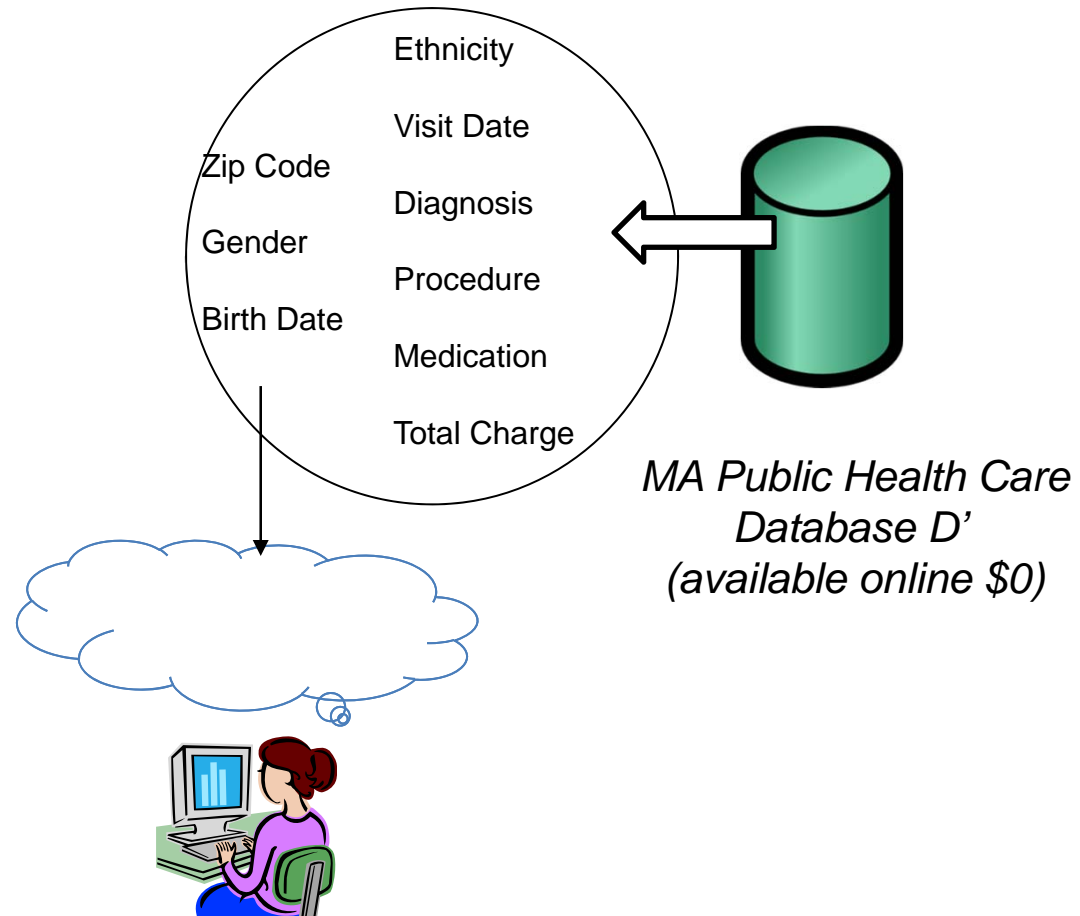
- Database privacy problem
- Smart grid privacy problems
- Summary and future work

# Talk Outline

- Database Privacy Problem
  - Source and Perturbation Model
  - Utility and Privacy Metrics
  - Examples
  - Related Results

# The Massachusetts Example

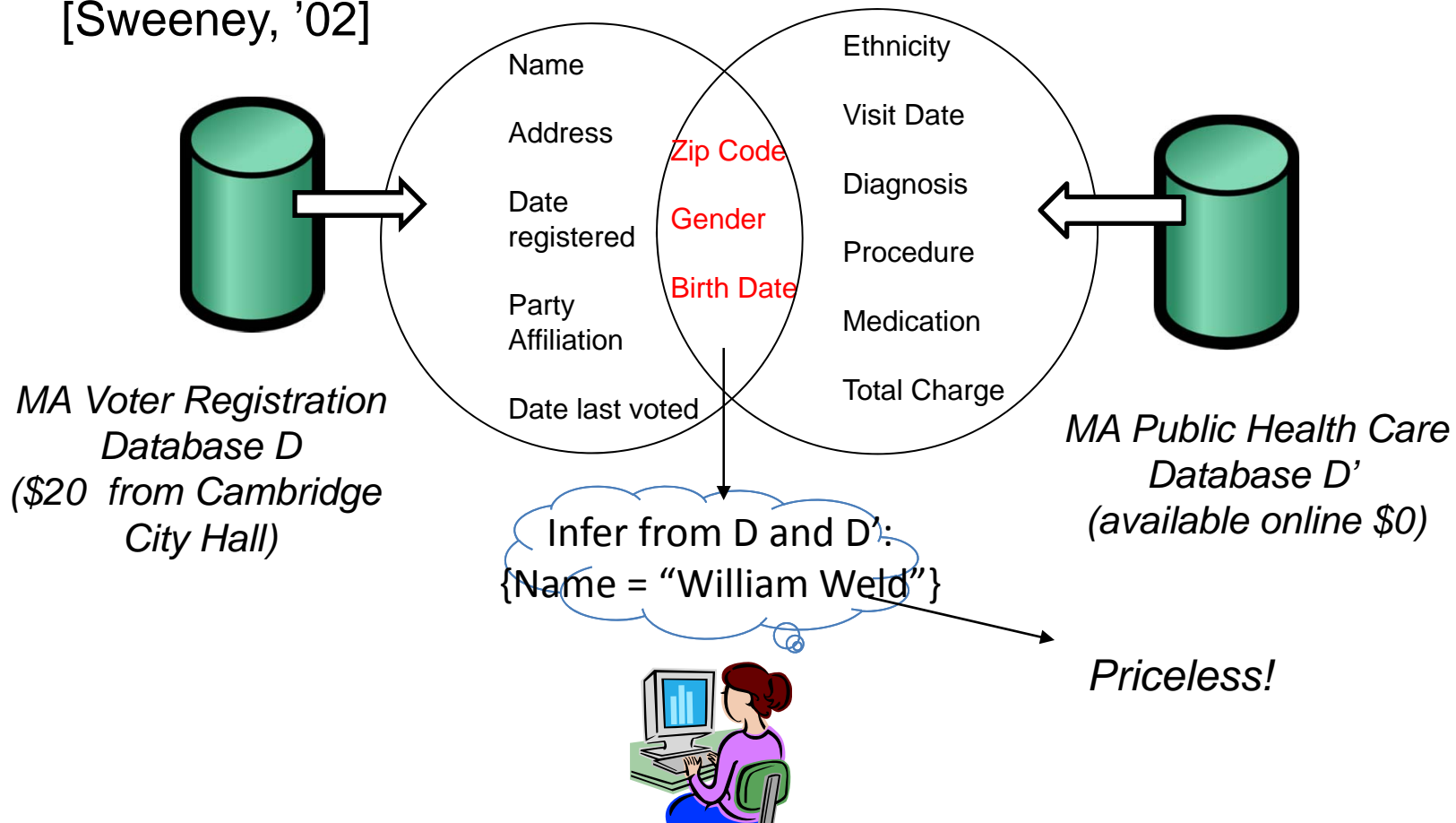
- Is it sufficient to hide personal information? [Sweeney, '02]



L. Sweeney, "k-anonymity: A model for protecting privacy," *Intl. J. Uncertainty, Fuzziness, and Knowledge-based Systems*, vol. 10, no. 5, pp. 557–570, 2002.

# The Massachusetts Example

- Unique identification via correlation from two public databases [Sweeney, '02]



L. Sweeney, "k-anonymity: A model for protecting privacy," *Intl. J. Uncertainty, Fuzziness, and Knowledge-based Systems*, vol. 10, no. 5, pp. 557–570, 2002.



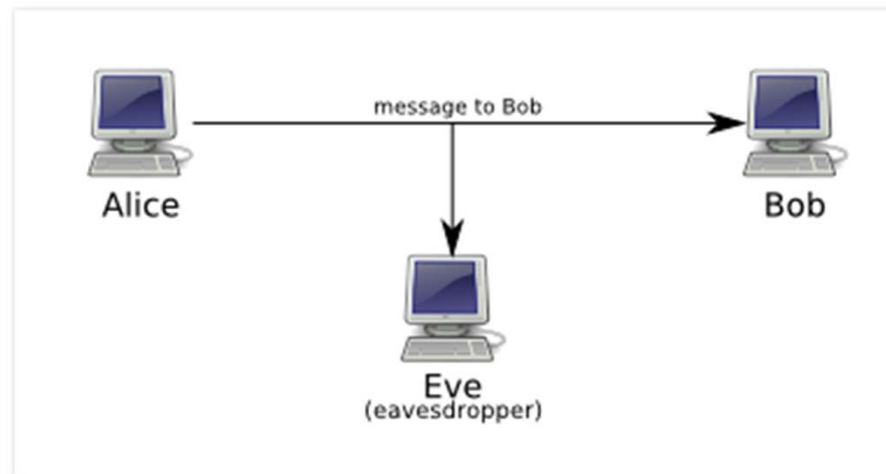
# The Privacy Problem is Pervasive

- Netflix, Reselling of medical data, Query logs, .....
- Sources leak information in unforeseeable ways
  - **Intra-source leaks**: hidden correlations between public and private information, e.g.: electronic health systems, census
  - **inter-source leaks**: correlation between sources [Sweeney, '02]
- But the electronic sources cannot be shut down
  - Tremendous utility provided.
  - Cannot shut down Google or Facebook
- **Can we disclose (utility) while guaranteeing privacy?**



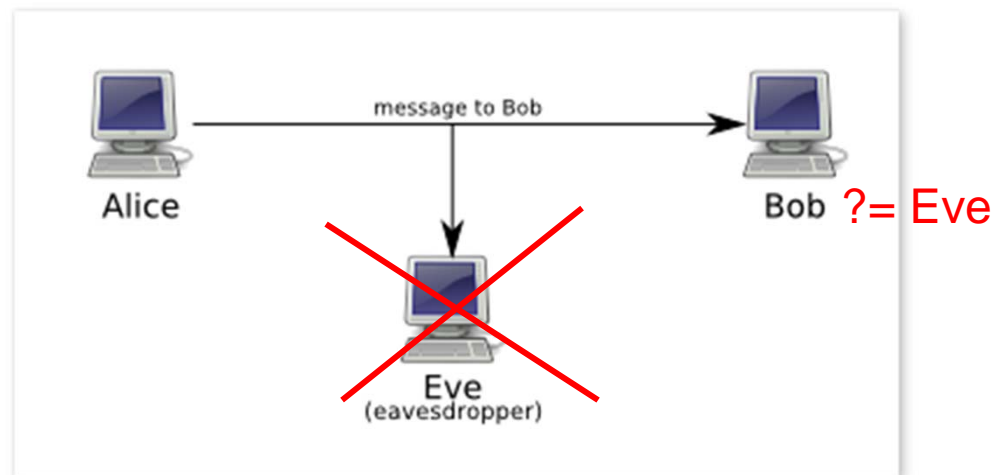
# Privacy vs. Secrecy!

- Privacy: the ability to prevent unwanted transfer of information (via inference or correlation) when legitimate transfers happen.
- **But privacy is not secrecy!**
- Secrecy Problem: Protocols and primitives clearly distinguish a malicious adversary vs. intended user and secret vs. non-secret data.
  - Encryption may be a solution.



# Privacy is not Secrecy!

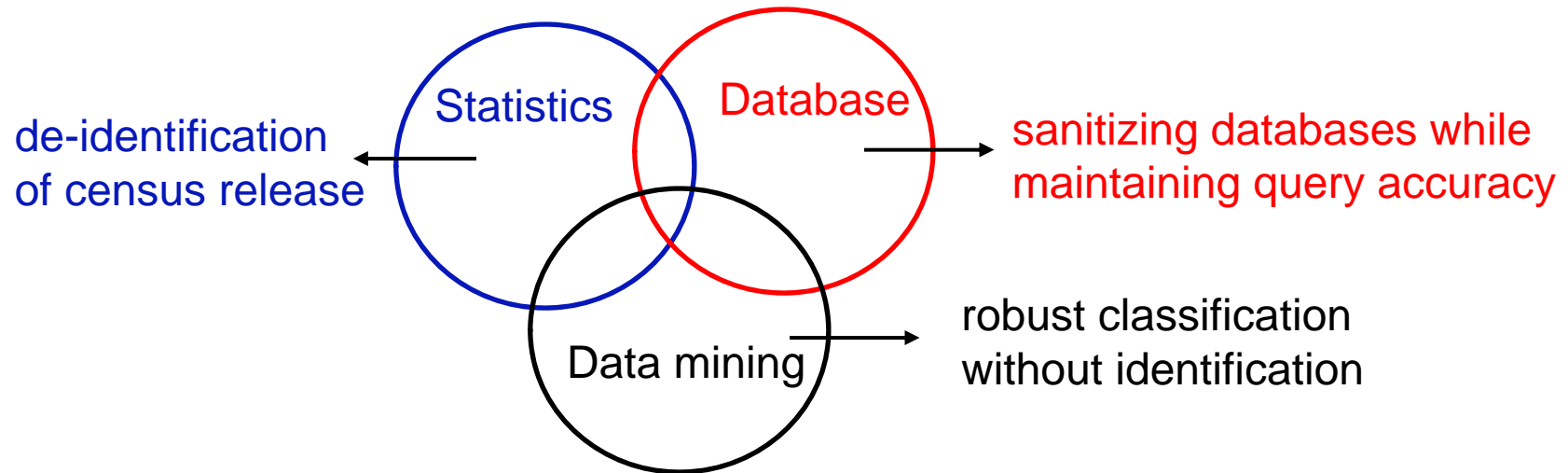
- Privacy: the ability to prevent unwanted transfer of information (via inference or correlation) when legitimate transfers happen.
- **But privacy is not secrecy!**
- Privacy problem: disclosing data provides informational utility while also enabling potential loss of privacy
  - Every user is potentially an adversary
  - Encryption is not a solution!





# Existing Approaches

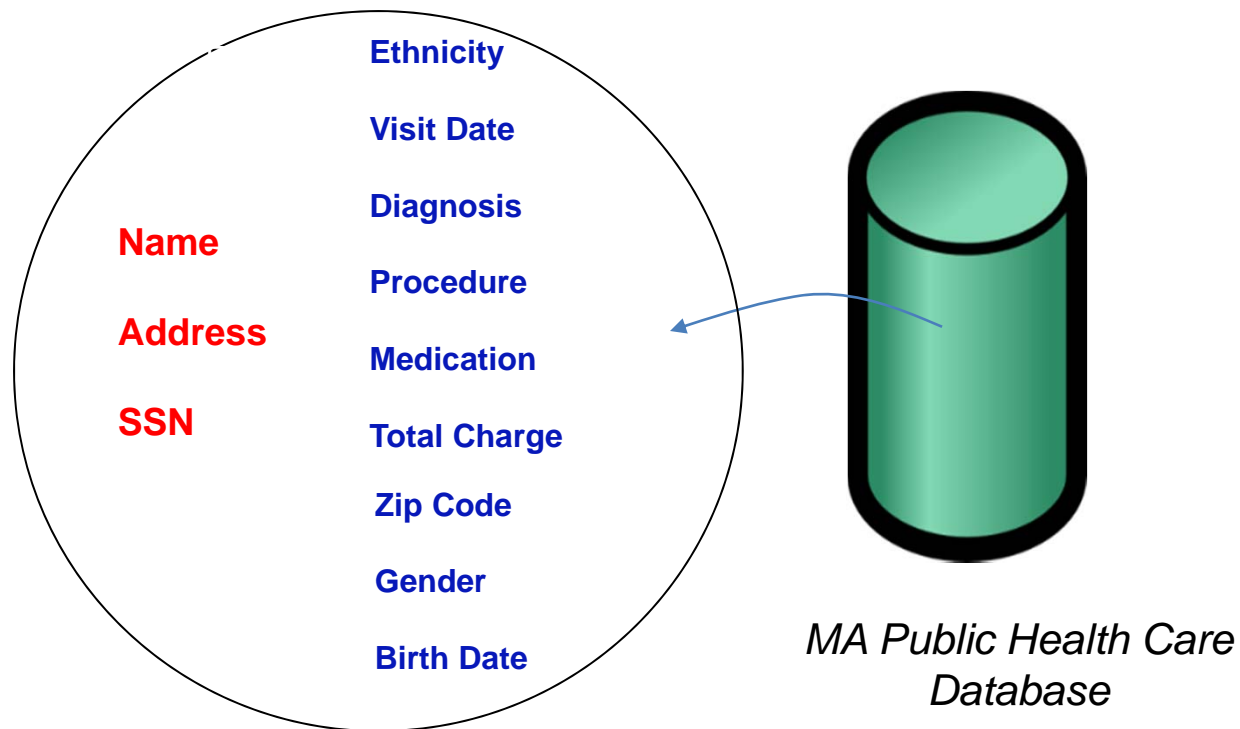
- Privacy problem lies at the intersection of multiple communities.



- **Application-specific** approaches without universal guarantees
- **CS Theory: differential privacy** – cryptography motivated definition
  - How to guarantee non-identification
  - Privacy paramount
- **Utility vs. privacy tradeoff remains unsolved.**

# Privacy Problem: A New Insight

- Any data source has **public** and **private** attributes

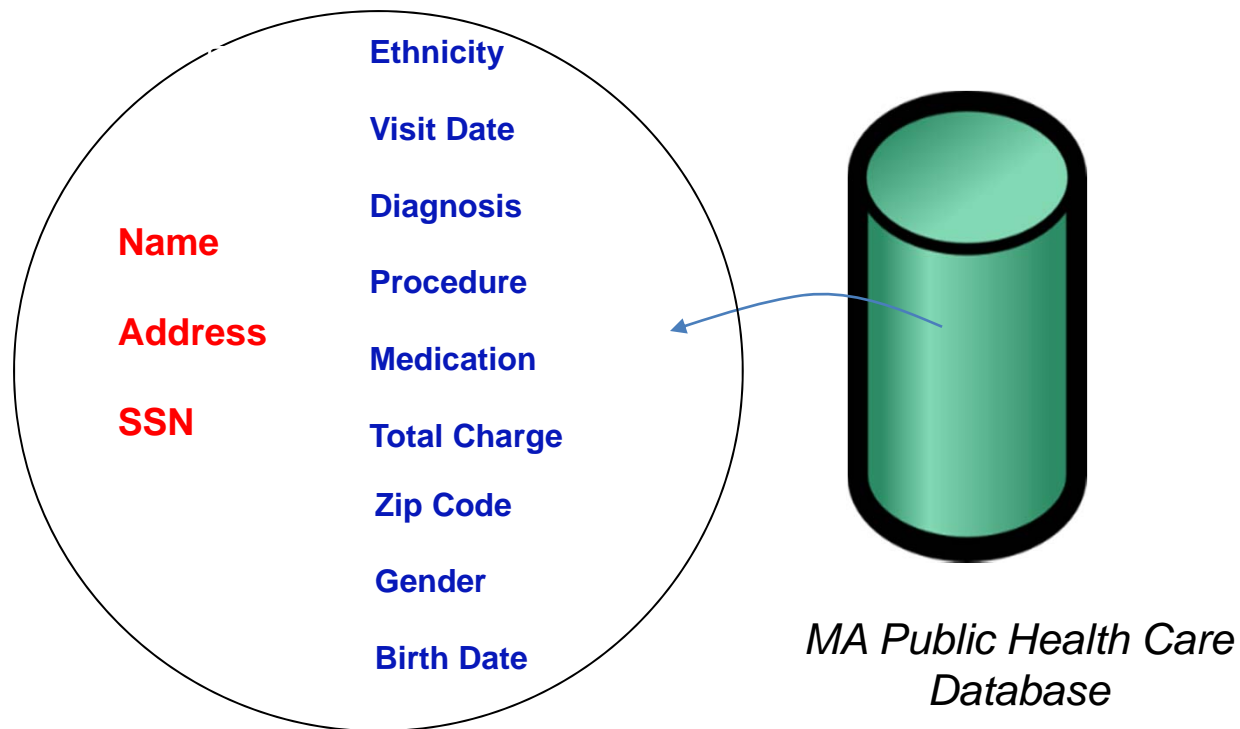


---

L. Sankar, S. R. Rajagopalan, and H. V. Poor. "Utility and privacy of data sources: Can Shannon help conceal and reveal information?," *ITA Workshop*, La Jolla, CA, Feb. 2010.

# Privacy Problem: A New Insight

- Any data source has **public** and **private** attributes
- Want to reveal public attributes maximally without revealing the private attributes



---

L. Sankar, S. R. Rajagopalan, and H. V. Poor. "Utility and privacy of data sources: Can Shannon help conceal and reveal information?," *ITA Workshop*, La Jolla, CA, Feb. 2010.

# Privacy Problem: A New Insight

- But... private and public attributes are correlated.
- Controlling privacy leakage amounts to controlling the correlation.
- Correlation can be controlled via perturbation of public attributes.
- Best U-P tradeoff: finding the minimal perturbation that achieves a desired correlation.
- Our contribution: a framework based on rate-distortion theory with universal metrics for utility and privacy.

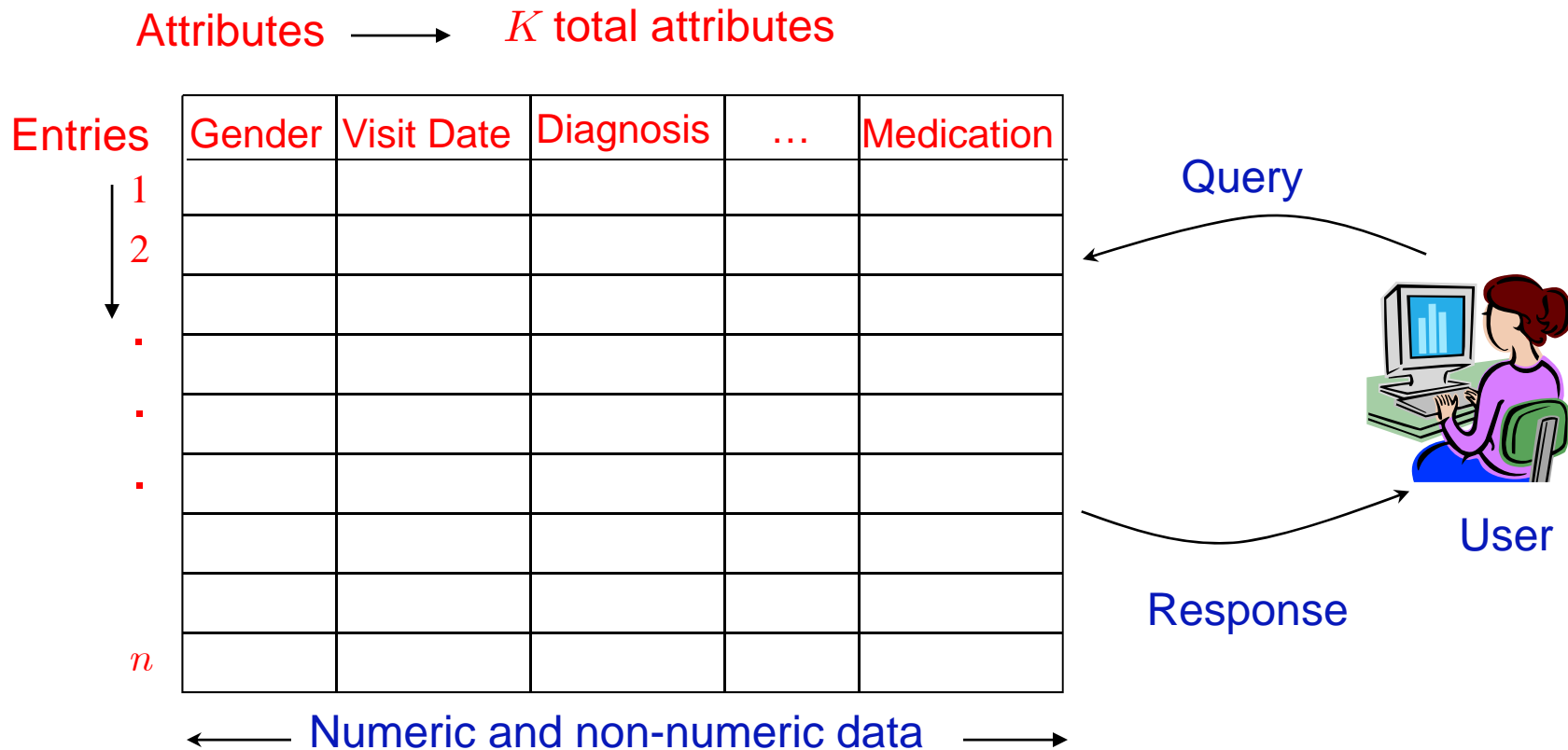
---

L. Sankar, S. R. Rajagopalan, and H. V. Poor. "Utility and privacy of data sources: Can Shannon help conceal and reveal information?," *ITA Workshop*, La Jolla, CA, Feb. 2010.



# The Database Privacy Problem

- A database is a table – rows: individual entries (total of  $n$ ); columns: attributes for each individual (total of  $K$ )



# Database: Source Model

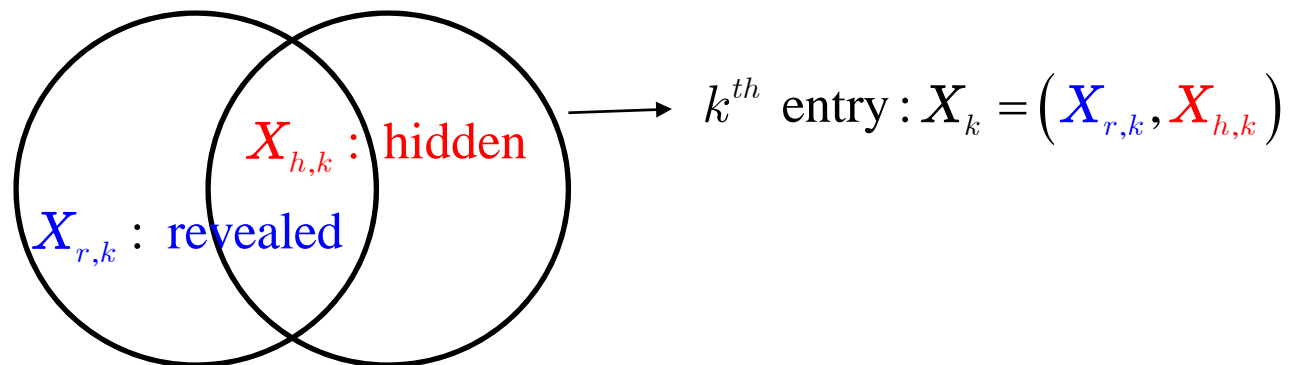
- A *real* database  $d$  is (typically) a table with  $n \gg 1$  rows (entries) and  $K$  columns (attributes)

Our model:

- Database  $d$  with  $n$  rows is a sequence of  $n$  i.i.d. observations of a vector random variable  $X = (X_1 X_2 \dots X_K)$  with the distribution

$$p_X(\mathbf{x}) = p_{X_1 X_2 \dots X_K}(x_1, x_2, \dots, x_K)$$

- Attributes divided into  $K_r$  public (revealed) and  $K_h$  private (hidden) variables, typically not disjoint



---

L. Sankar, S. R. Rajagopalan, and H. V. Poor. "A theory of utility and privacy of data sources," *Proc. of IEEE Intl. Symp. Inform. Theory*, Austin, TX, Jun. 13-18 2010.

# Database: Utility vs. Privacy

- **The Utility-Privacy Problem:**
  - How to reveal the **public** variables while hiding the **private** variables given that the two sets are correlated?

# Database: Utility vs. Privacy

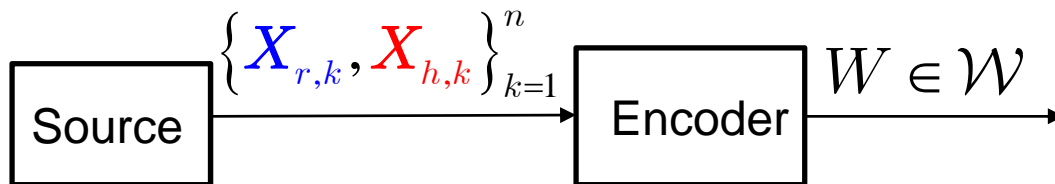
- The Utility-Privacy Problem: Rate distortion theory with privacy is a natural fit!

# Database: Utility vs. Privacy

- The Utility-Privacy Problem: Rate distortion theory with privacy is a natural fit!
- Encoder maps  $d(X^n)$  to a “sanitized” database (SDB)  $d'$

$$\text{Encoder} : X^n \rightarrow \mathcal{W} = \{SDB_1, SDB_2, \dots, SDB_M\}$$

- $M$ : number of revealed (“quantized”) databases



---

L. Sankar, S. R. Rajagopalan, and H. V. Poor. “A theory of utility and privacy of data sources,” *Proc. of IEEE Intl. Symp. Inform. Theory*, Austin, TX, Jun. 13-18 2010.

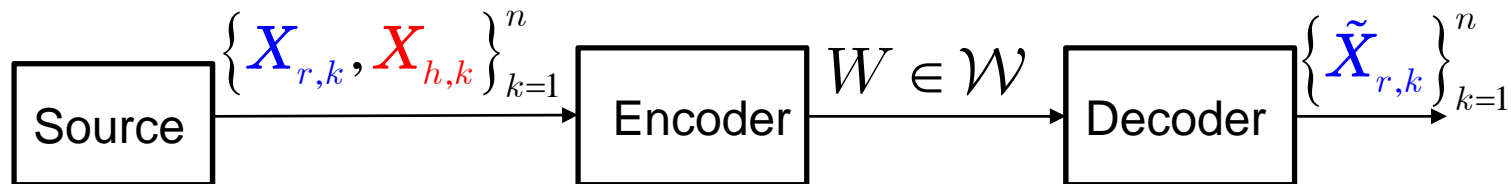
# Database: Utility vs. Privacy

- The Utility-Privacy Problem: **Rate distortion theory with privacy is a natural fit!**
- Encoder maps  $d(X^n)$  to a “sanitized” database (SDB)  $d'$

$$\text{Encoder: } X^n \rightarrow \mathcal{W} = \{SDB_1, SDB_2, \dots, SDB_M\}$$

- $M$ : number of revealed (“quantized”) databases
- Decoder: Uses  $d'$  to obtain a “reconstructed” database (for query processing)

$$\text{Decoder: } \mathcal{W} \rightarrow \tilde{X}_h^n$$




---

L. Sankar, S. R. Rajagopalan, and H. V. Poor. “A theory of utility and privacy of data sources,” *Proc. of IEEE Intl. Symp. Inform. Theory*, Austin, TX, Jun. 13-18 2010.

# Utility and Privacy Metrics

- Utility is a measure of closeness of  $d$  and  $d'$ .
- Map utility to fidelity (distortion)
  - Utility is affected by added noise, limited precision, suppression.
- Utility constraints  $\Delta_d \rightarrow$  bound on avg. distortion per entry (row)

$$\Delta_d \equiv \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \rho \left( X_{r,i}, \tilde{X}_{r,i} \right) \right] \leq D + \varepsilon$$

- $\rho$ : distance-based function (e.g.: Hamming, Euclidean, K-L)

# Privacy Metric

- Map privacy to equivocation
  - Privacy is a measure of ‘uncertainty’ about hidden data given revealed data.
- Privacy constraints  $\Delta_p \rightarrow$  equivocation on average per entry (row)

$$\Delta_p \equiv \frac{1}{n} H(\mathbf{X}_h^n | W) > E - \varepsilon$$

- $E$ : lower bound on the avg. privacy per entry



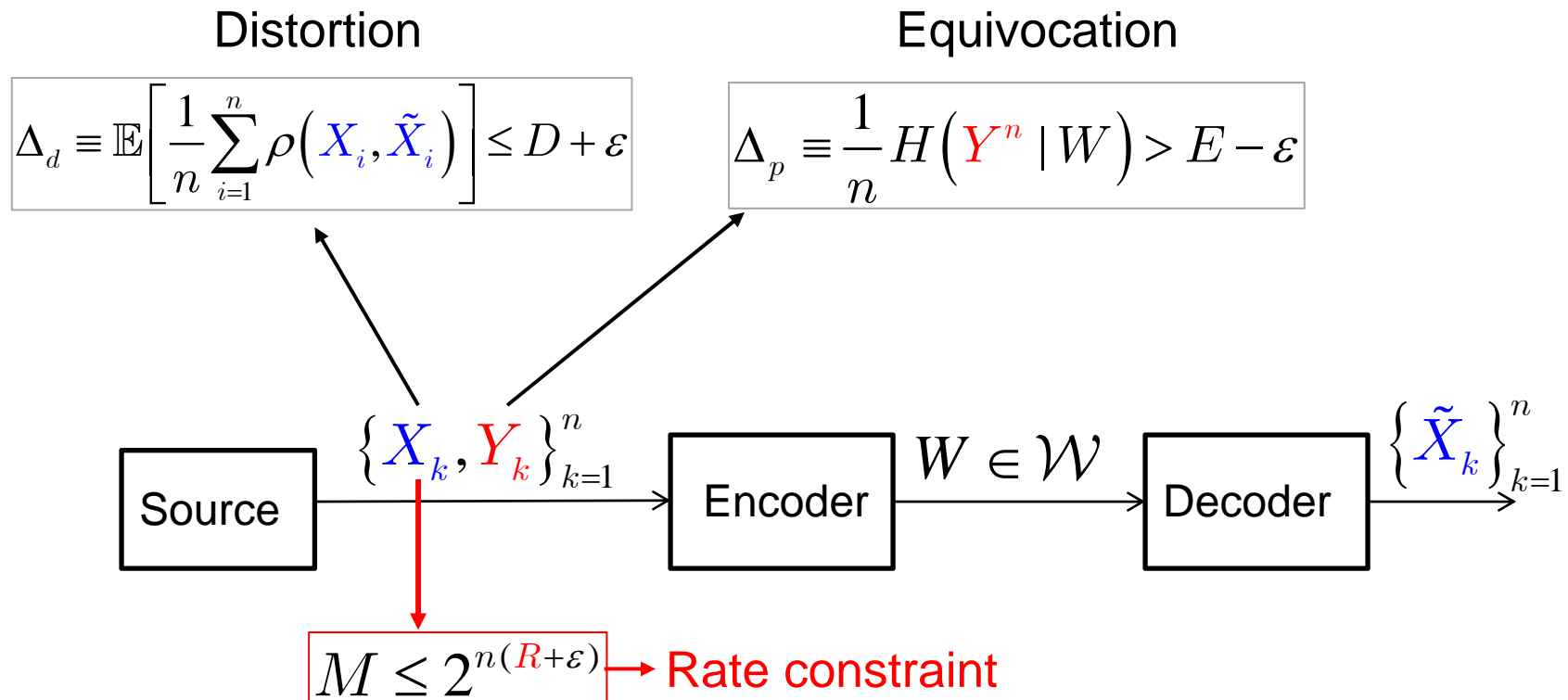
# The Utility-Privacy Tradeoff

- Utility-privacy tradeoff region ( $\mathcal{T}$ ) is

$$\mathcal{T} \equiv \{(D, E): (D, E) \text{ is feasible}\}$$

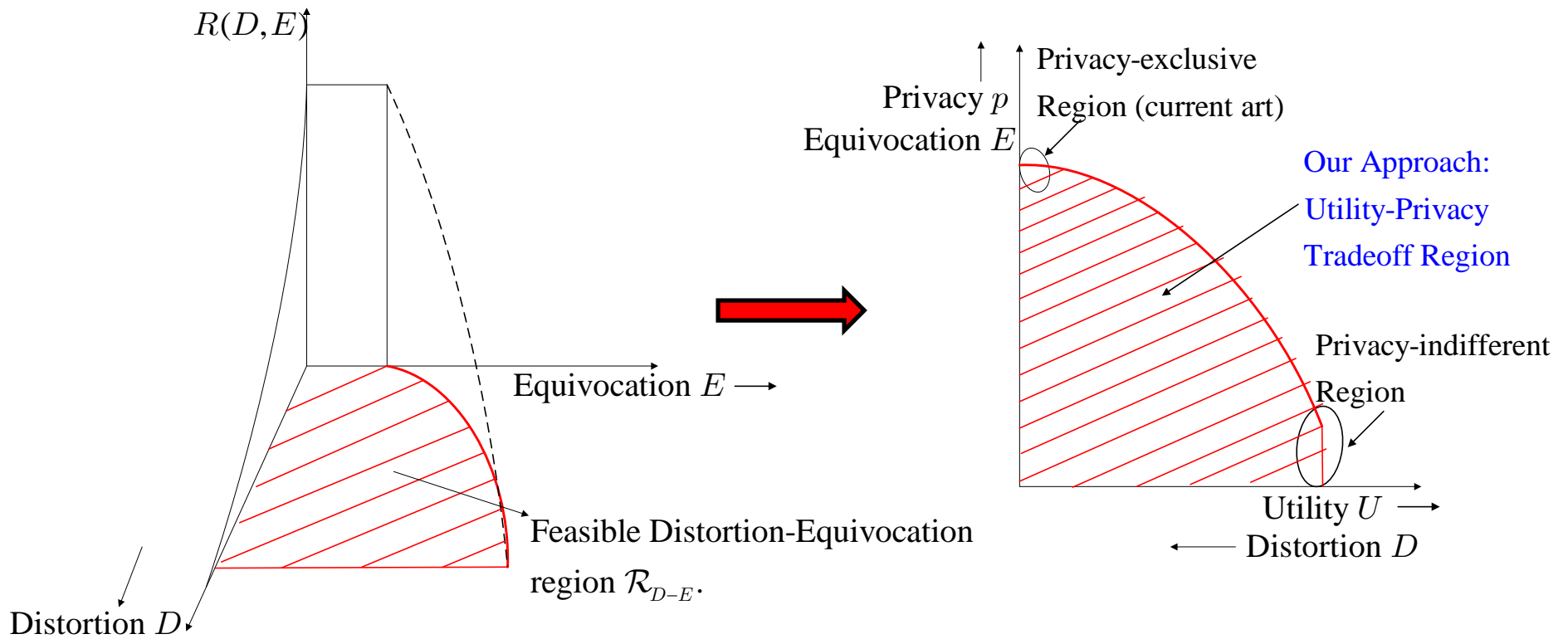
- How do we compute  $\mathcal{T}$ ?
- Add an additional rate constraint and map it to a rate-distortion-equivocation problem

# A Source Coding Problem with Privacy



- Simplified version of the database privacy problem with **additional rate constraint**
  - Rate constraint bounds the number of “quantized” sequences
  - For U-P tradeoff this seems **superfluous**

# Utility-Privacy/RDE Regions

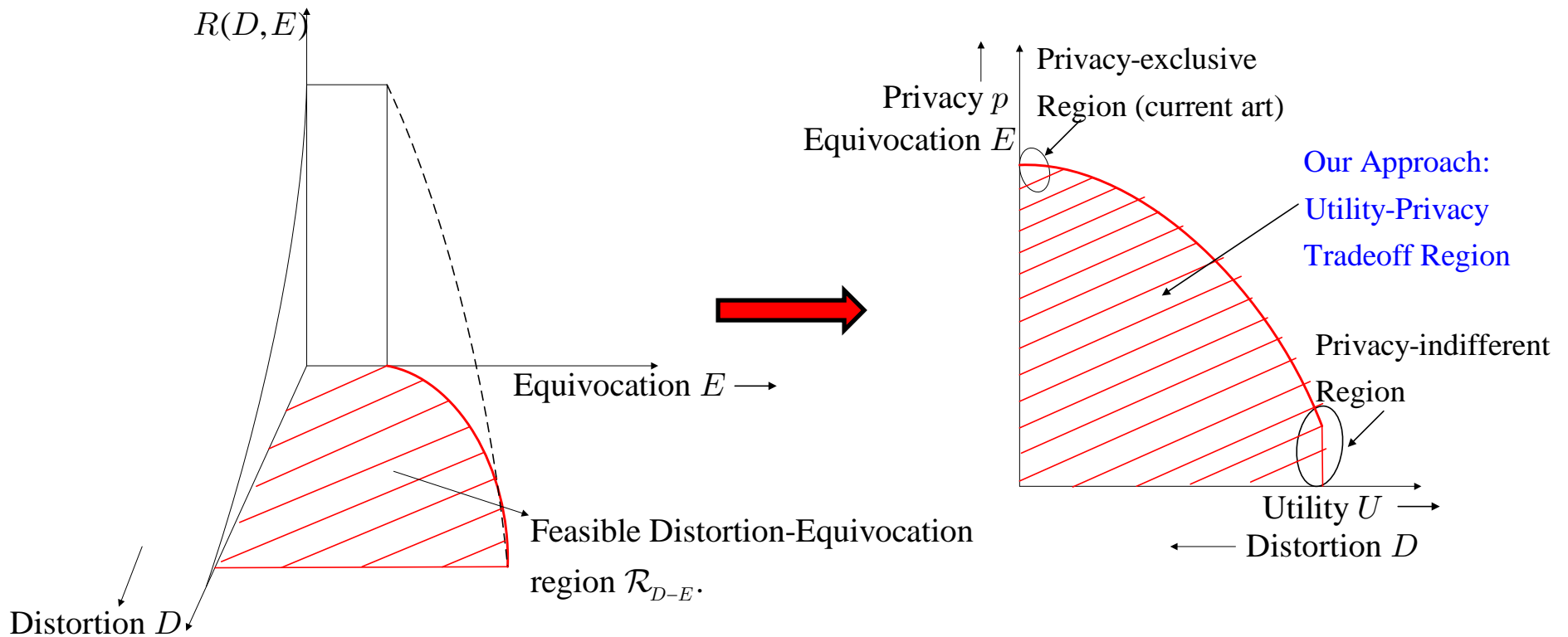


(a): Rate-Distortion-Equivocation Region

(b): Utility-Privacy Tradeoff Region

L. Sankar, S. Raj Rajagopalan, H. V. Poor, "A theory of privacy and utility in databases," submitted to the *IEEE Trans. Inform. Theory*, Feb. 2011.

# Utility-Privacy/RDE Regions



(a): Rate-Distortion-Equivocation Region

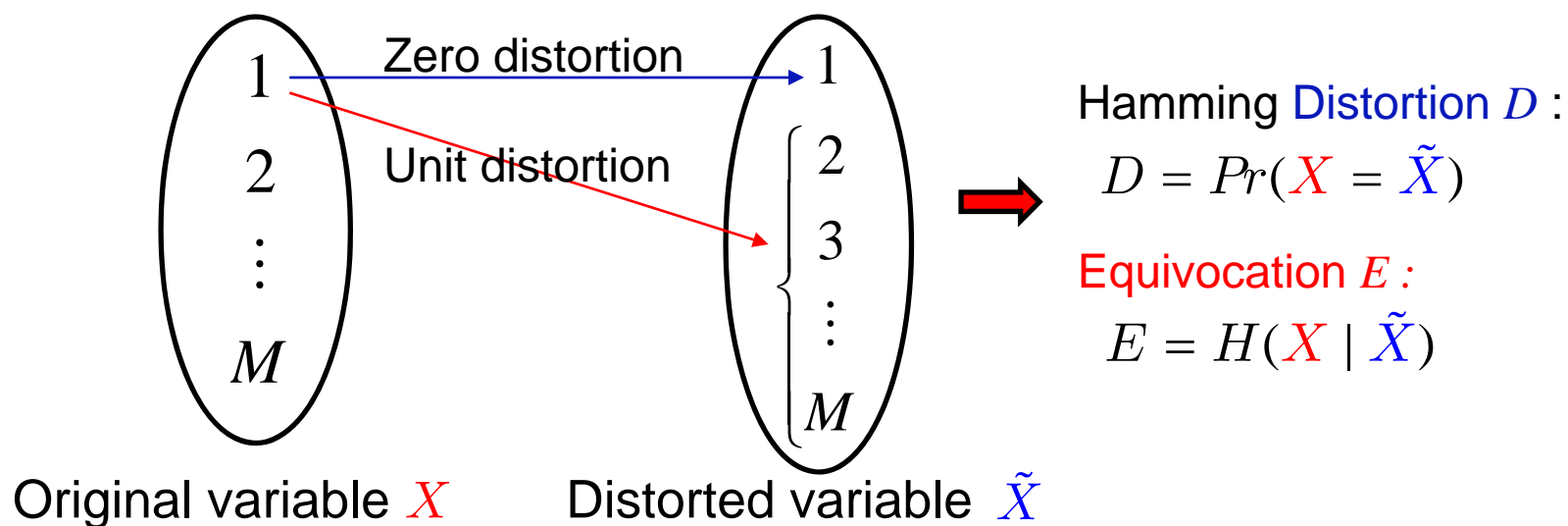
(b): Utility-Privacy Tradeoff Region

**For a database with utility and privacy constraints,  $\mathcal{T} = \mathcal{R}_{D-E}$ . [SRP, ISIT '10]**

L. Sankar, S. Raj Rajagopalan, H. V. Poor, "A theory of privacy and utility in databases," submitted to the *IEEE Trans. Inform. Theory*, Feb. 2011.

# Example 1: Categorical Database

- Categorical data: finite alphabet data with discrete distribution
  - e.g.: SSN, zipcode, etc.



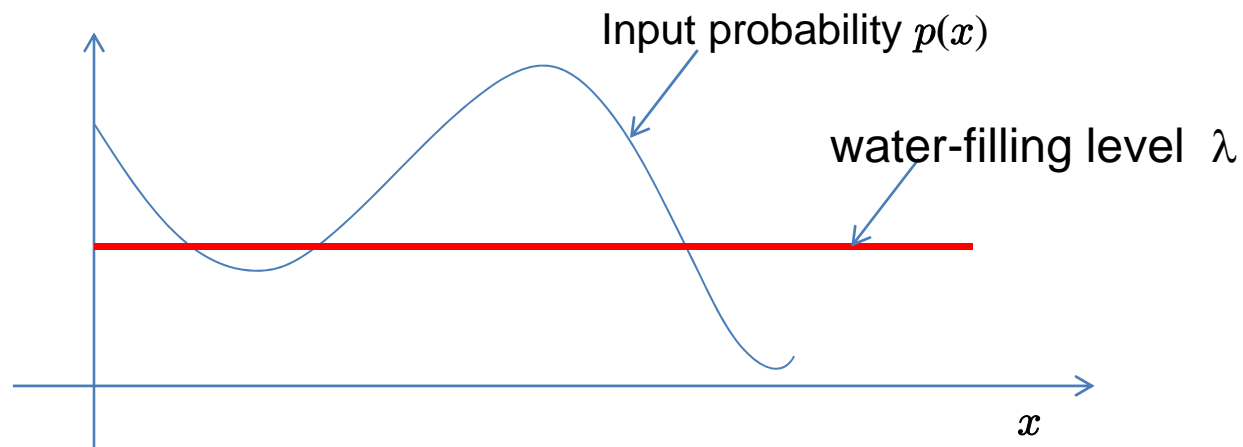
- The categorical database case has remained largely unaddressed in privacy research.

---

L. Sankar, S. R. Rajagopalan, and H. V. Poor. "An information-theoretic approach to privacy," *Proc. 48th Allerton Conf. Comm., Cntl., and Comp.*, Monticello, IL, Sep, 2010.

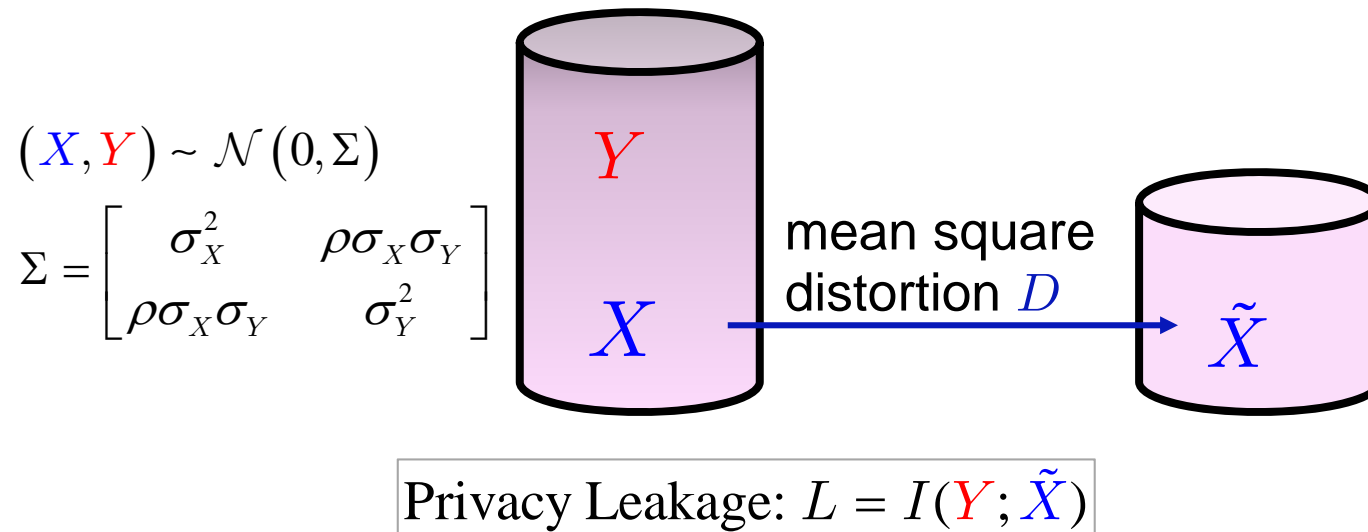
# Example 1: Categorical Database

- Optimal input to output mapping: reverse 'water-filling'
  - Only  $x$  with  $p(x) > \lambda$  revealed ( $\lambda$  : water-level).
- Eliminates samples with low probabilities (relative to water-level  $\lambda$ )
  - Equivalent to outlier aggregation/suppression (dominant statistical approaches)
  - Such samples reveal the most information
- As  $D \uparrow$ ,  $\lambda \uparrow$  (relative to distribution) to reveal fewer samples



# Example 2: Numerical Database

- Numerical data: finite/infinite alphabet real data
  - e.g.: results of medical tests, clinical trials, etc.
  - Medical research often assumes Gaussian distributed data



- Sanitized DB remains Gaussian distributed.
  - Gaussian  $\tilde{X}$  achieves minimal  $R(D, E)$  and maximal privacy  $\Gamma(D)$

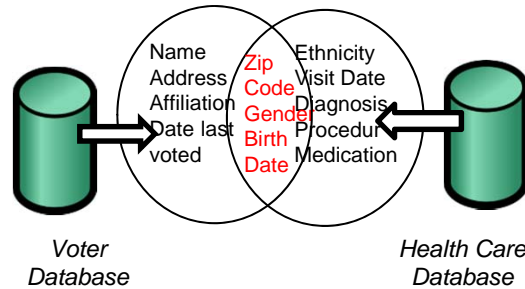
---

L. Sankar, S. R. Rajagopalan, and H. V. Poor. "An information-theoretic approach to privacy," *Proc. 48th Allerton Conf. Comm., Cntl., and Comp.*, Monticello, IL, Sep, 2010.

# Related and New Results

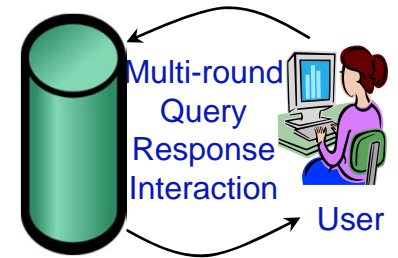
## The Side Information Problem

Model and U-P tradeoff for decoder side information



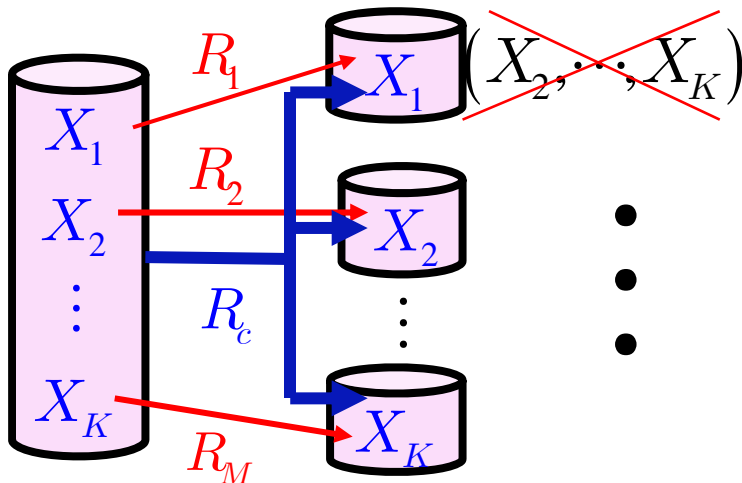
## The Successive Disclosure Problem

Conditions for no privacy leaks over successive queries relative to one-shot



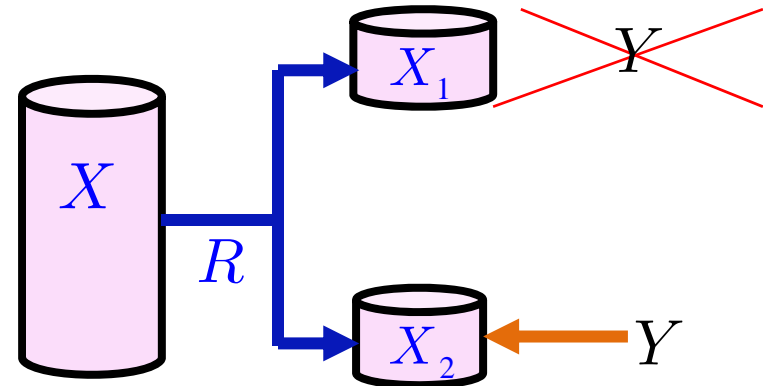
L. Sankar, S. Raj Rajagopalan, H. V. Poor, "A theory of privacy and utility in databases," submitted to the *IEEE Trans. Inform. Theory*, Feb. 2011.

## Multi-user Privacy



R. Tandon, L. Sankar, H. V. Poor, "Multiuser Privacy and Common Information", submitted to *ISIT 2011*.

## Discriminatory Coding and Privacy



R. Tandon, L. Sankar, H. V. Poor, "Discriminatory Lossy Source Coding", submitted to *Globecom 2011*.



# Talk Outline

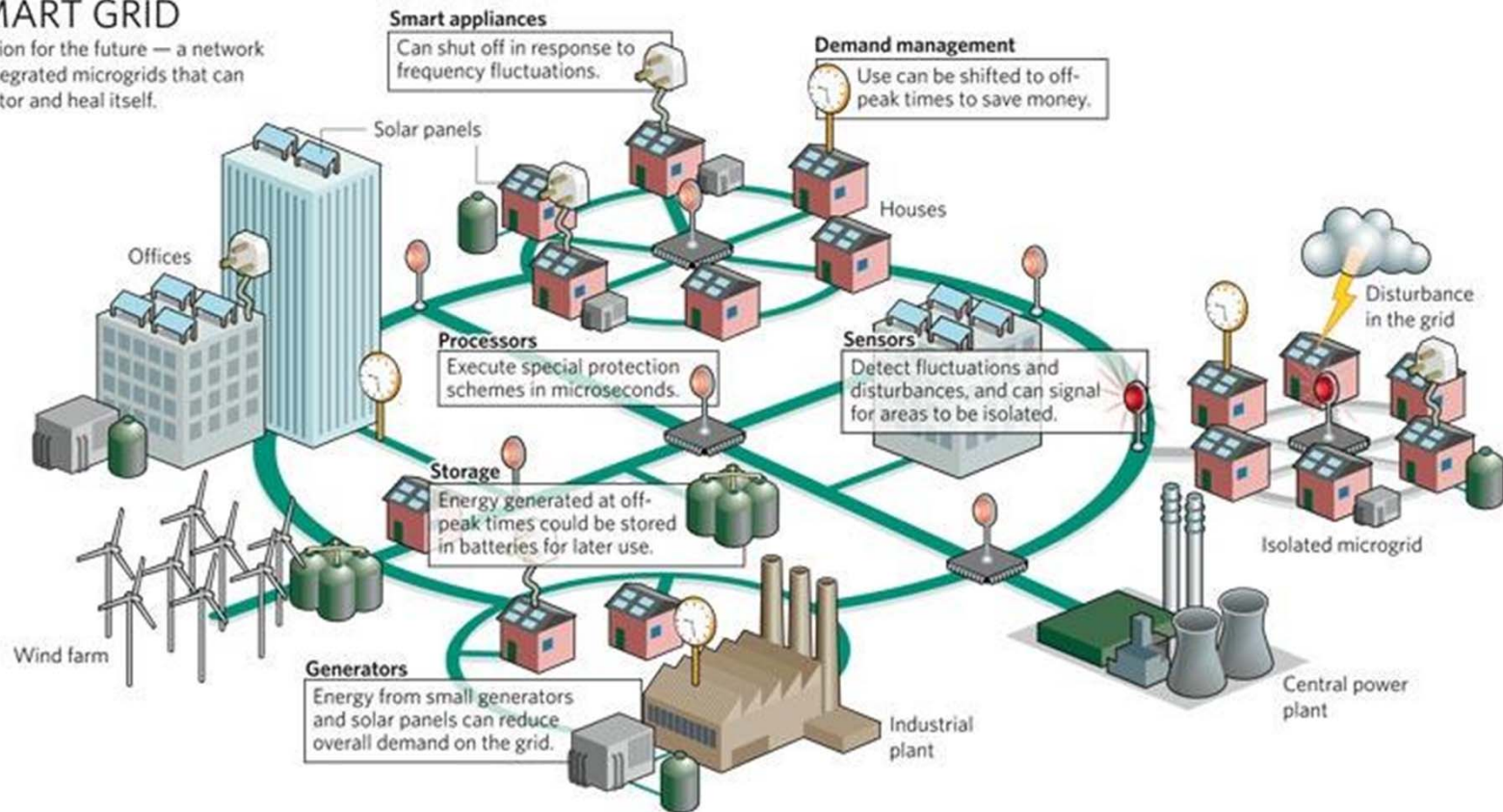
- Database privacy problems
- **Smart grid privacy problems**
- Summary and future work

# What is a Smart Grid?

- Smart Grid : Overlay electrical grid with sensors (phasor monitoring units - PMUs) and control systems (SCADA) to enable:
  - **reliable and secure** network monitoring, load balancing, energy efficiency via smart meters, and integration of new energy sources

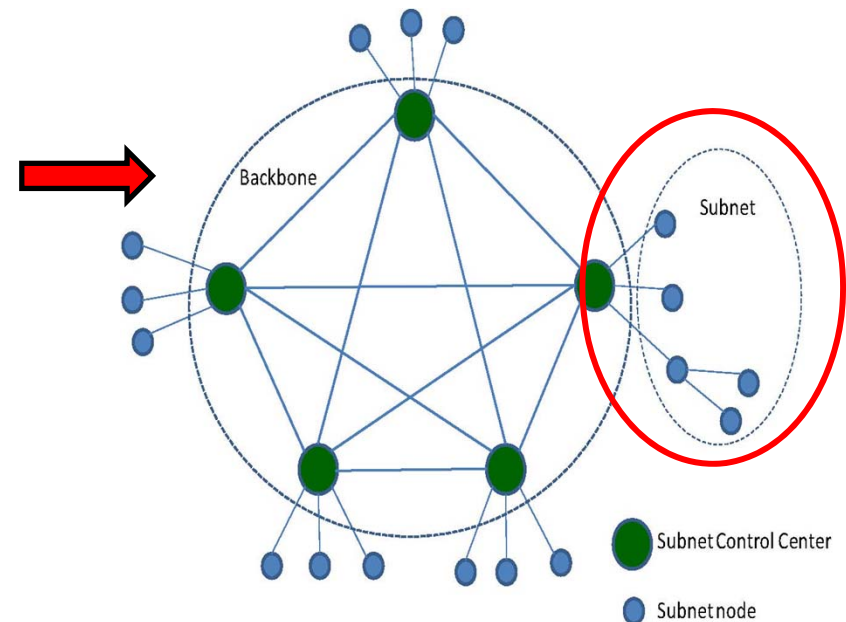
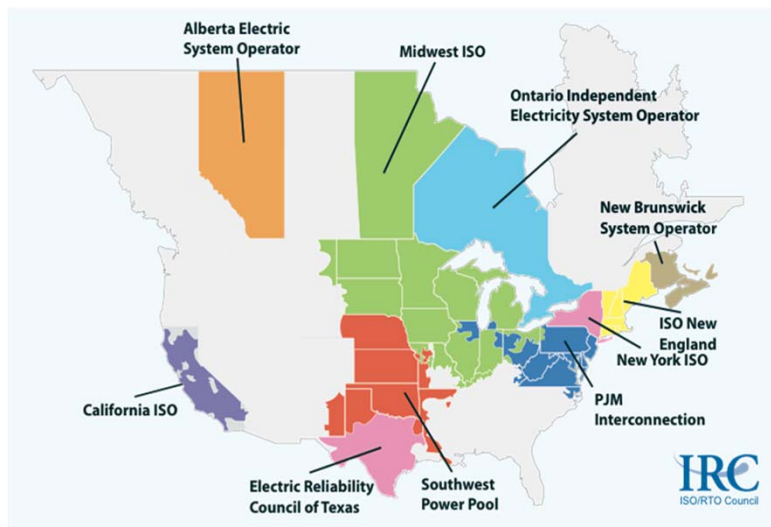
## SMART GRID

A vision for the future — a network of integrated microgrids that can monitor and heal itself.



# Smart Grid Privacy

- N.A. Grid: interconnected regional transmission organizations which:
  - need to share measurements on state estimation for **reliability** (utility)
  - wish to withhold information for economic **competitive** reasons (privacy)
- Leads to a new problem of **competitive privacy**



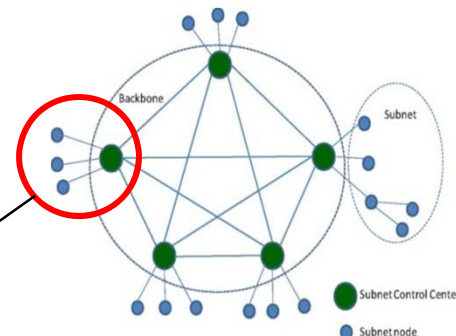
**L. Sankar**, S. Kar, R. Tandon, and H. V. Poor, "Competitive privacy in the smart grid: An information-theoretic approach," to appear, *IEEE SmartGridComm*, Oct. 2011.

# System Model

- Noisy measurements  $Y_k$  at RTO  $k$  with interference from other RTOs:

$$Y_k = \sum_{m=1}^M H_{k,m} X_m + Z_k, \quad k = 1, 2, \dots, M$$

$m^{\text{th}}$  system state



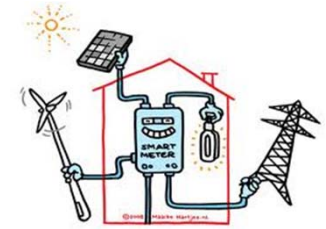
- Utility: Mean squared error state and its estimate
- Privacy: leakage of state from measurements and messages
- Cooperation leads to inevitable leakage of state information

[SKTP '11]: A one-shot Wyner-Ziv coding maximizes privacy for a desired utility at each RTO.

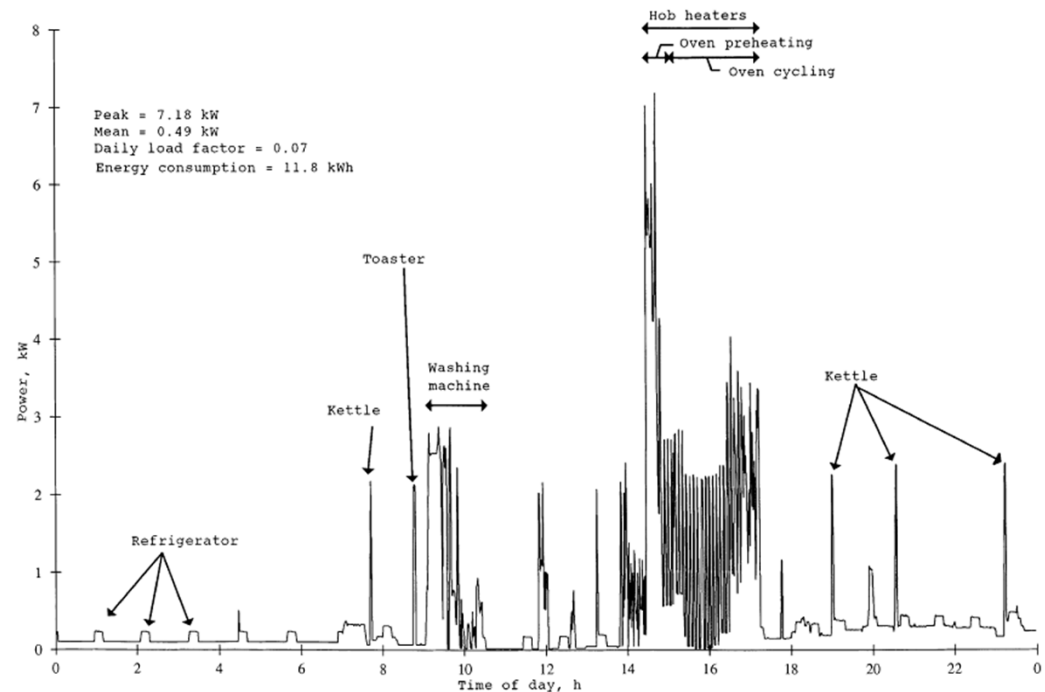
---

L. Sankar, S. Kar, R. Tandon, and H. V. Poor, "Competitive privacy in the smart grid: An information-theoretic approach," to appear, *IEEE SmartGridComm*, Oct. 2011.

# Smart Meter Privacy



- Smart meter is a critical enabler of the Smart Grid
- For consumers: Tariff- and load-aware appliance usage
- For electricity suppliers: Load balancing; **data mining** (analytics)
  - **Data mining: tremendous utility to supplier; huge consumer privacy risk**
- Utility-Privacy tradeoff via rate-distortion for sources with memory

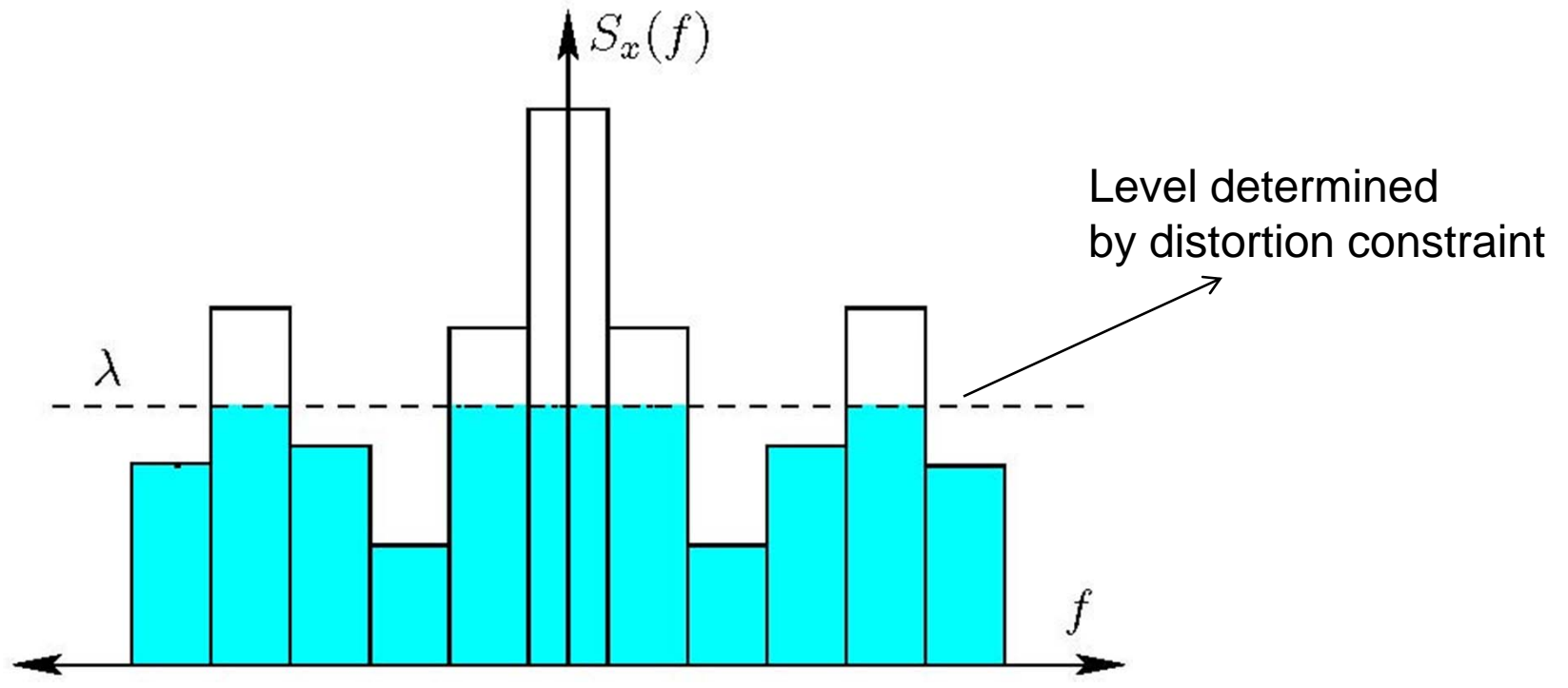


---

S. Rajagopalan, **L. Sankar**, S. Mohajer, and H. V. Poor, "Smart meter privacy: Utility-privacy tradeoff," to appear, *IEEE SmartGridComm*, Oct. 2011.

# Smart Meter Privacy: Results

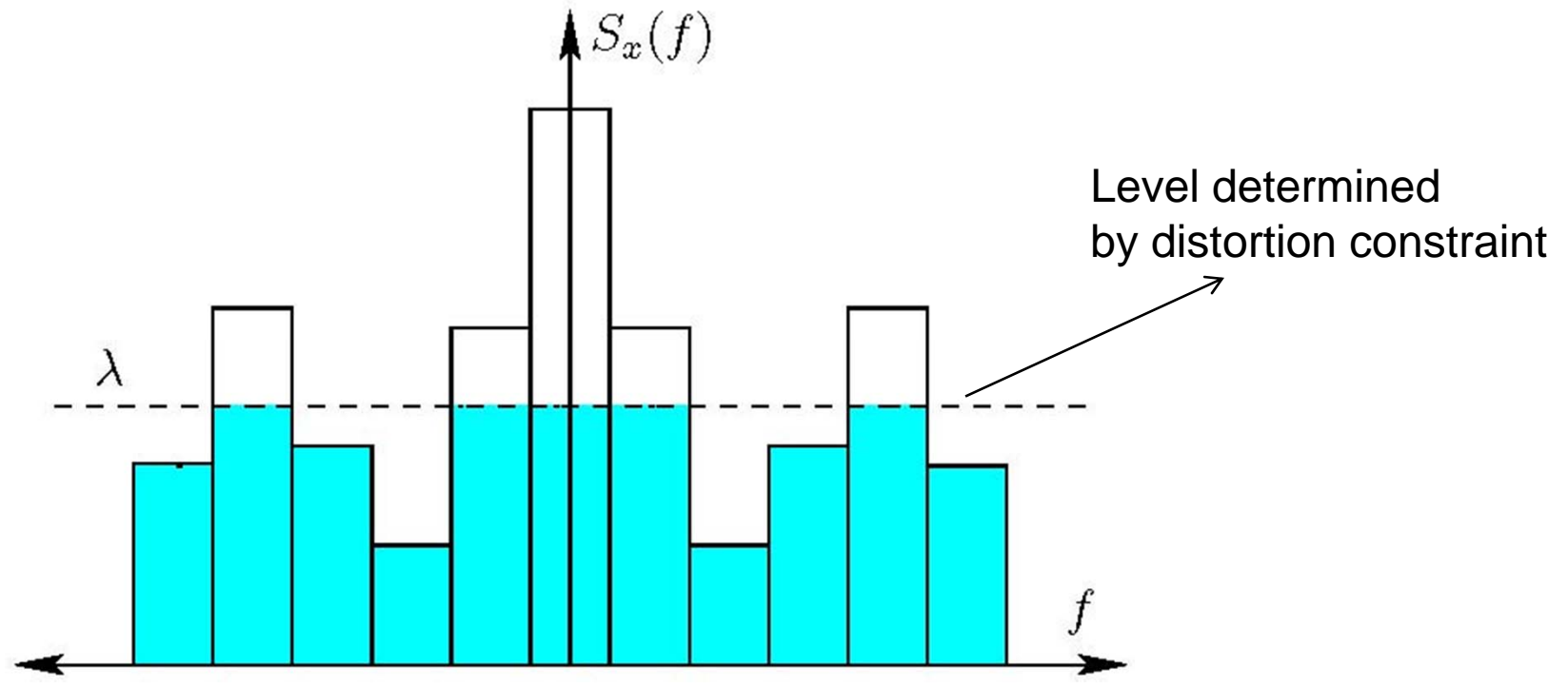
- Rate-distortion-leakage for sources with memory
  - ‘reverse water-filling’ in transform domain (without leakage or when hiding measurements itself)



S. Rajagopalan, L. Sankar, S. Mohajer, and H. V. Poor, “Smart meter privacy: An information-theoretic approach,” to appear, *IEEE SmartGridComm*, Oct. 2011.

# Smart Meter Privacy: Results

- Rate-distortion-leakage for sources with memory
  - ‘reverse water-filling’ in transform domain (without leakage or when hiding measurements itself)
  - With general inference model – generalized reverse water-filling solution



S. Rajagopalan, L. Sankar, S. Mohajer, and H. V. Poor, “Smart meter privacy: An information-theoretic approach,” submitted to *IEEE SmartGridComm*, Apr. 2011.

# Talk Outline

- Database privacy problem
- Smart grid privacy problems
- Summary and Future Work



# Summary

- The privacy problem is immediate and here to stay ... and multiply...
- One solution will not fit all applications...
- **But a framework provides the much needed abstraction**
- More needs to be done...



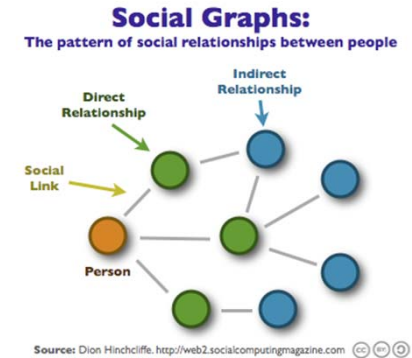
---

Trying to ward off regulators, the advertising industry has agreed on a standard icon — a little “i” — that it will add to most online ads that use demographics and behavioral data to tell consumers what is happening. – NY Times, Jan. 26, 2010.

# Future Work

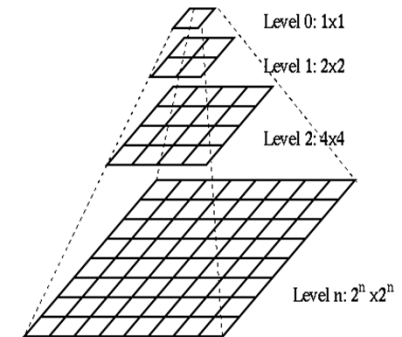
## Privacy in Social Networks:

- Quantifying privacy and utility in social networks
  - Information leakage due to social graph
  - How to quantify utility?



## Practical Privacy via Signal Processing:

- Compressive sensing, quantization, clustering, ...
- Universal lossy coding schemes



## Medical Database Privacy:

- De-identification and privacy
- Does synthetic data suffice?
- Need for re-identification?

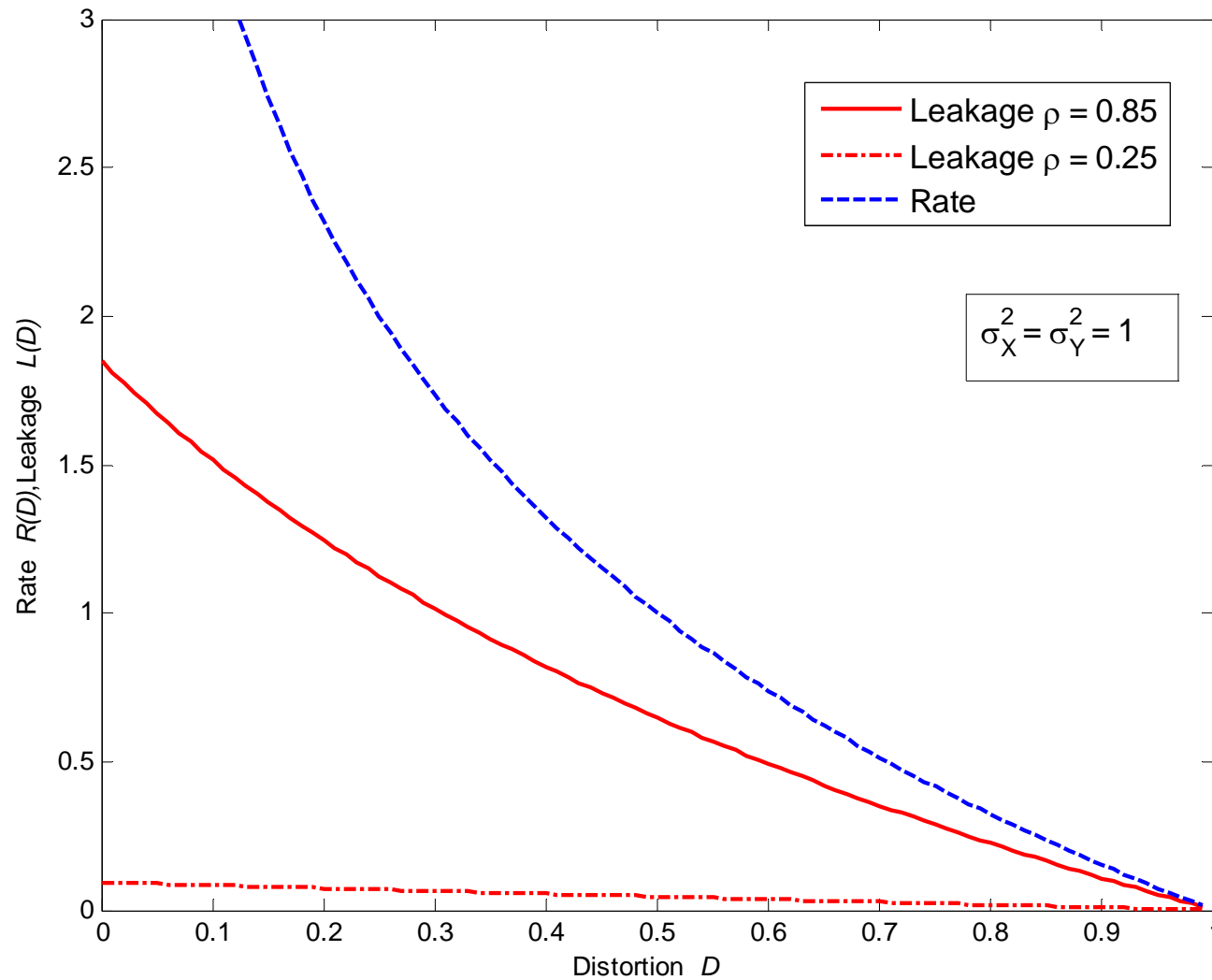


For more: ...

<http://www.princeton.edu/~lalitha>

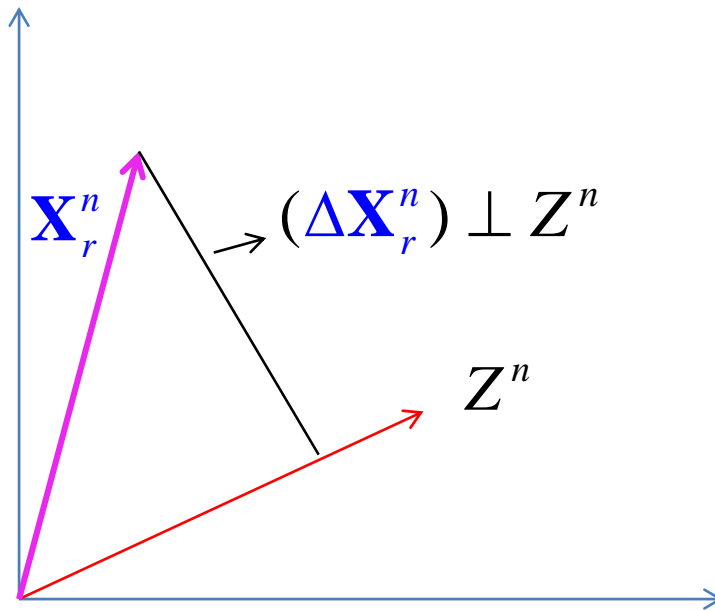
Thank you!

# Example 2: Numerical Database



# The Side Information Problem

- Optimality of Wyner-Ziv (Intuition):
  - Suffices to quantize the orthogonal variable  $\Delta X^n$
  - decoder uses  $\Delta X^n$  and  $Z^n$  to reconstruct the source
- Privacy achieved depends on the correlation between  $X^n$  and  $Z^n$



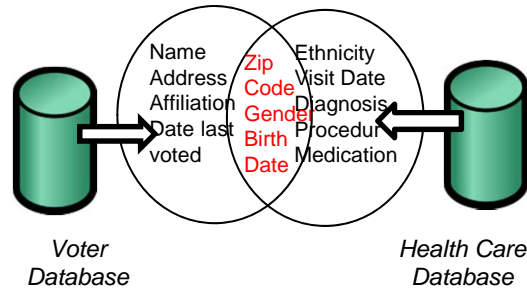
---

L. Sankar, S. Raj Rajagopalan, H. V. Poor, "A theory of privacy and utility in databases," submitted to the *IEEE Trans. Inform. Theory*, Feb. 2011.

# Related Results

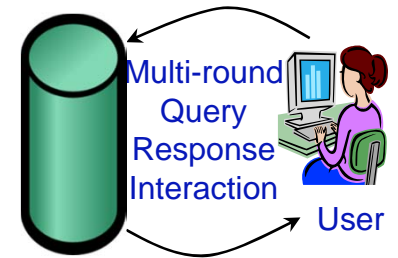
## The Side Information Problem

Model and U-P tradeoff for decoder side information



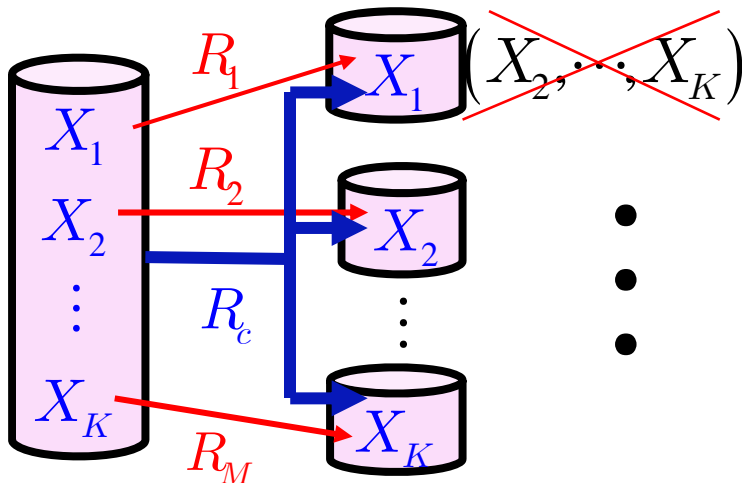
## The Successive Disclosure Problem

Conditions for no privacy leaks over successive queries relative to one-shot



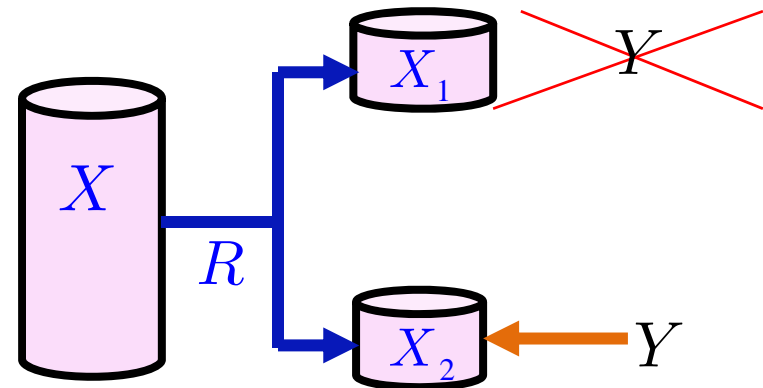
L. Sankar, S. Raj Rajagopalan, H. V. Poor, "A theory of privacy and utility in databases," submitted to the *IEEE Trans. Inform. Theory*, Feb. 2011.

## Multi-user Privacy



R. Tandon, L. Sankar, H. V. Poor, "Multiuser Privacy and Common Information," submitted to *ISIT 2011*.

## Discriminatory Coding and Privacy



R. Tandon, L. Sankar, H. V. Poor, "Discriminatory Lossy Source Coding," submitted to *Globecom 2011*.

# More Examples



- [The Netflix](#) competition [2006] to improve movie recommendations
  - Public training data set with movie preferences of 480,000 customers
  - Data was “de-identified” – stripped of specific personal details
- V. Shmatikov and A. Narayanan [ISSP, ‘08]
  - Compared film preferences of some anonymous customers with personal profiles on [imdb.com](#),
  - *Re-identification* using distinguishing information
- Netflix claimed
  - *“Anonymity of the study data is comparable to the strictest Federal standards for anonymizing personal health information.”*



---

A. Narayanan and V. Shmatikov, “Robust de-anonymization of large sparse datasets,” in *Proc. IEEE Intl. Symp. Security and Privacy*, Oakland, CA, May 2008, pp. 111–125.

# More Examples... Medical Data

- *New York Times* reports
  - Sale of clinical data is a huge and growing business.
  - *De-identified* information is “repackaged” and resold.
  - *New* regulations do NOT forbid sale of de-identified data.
- The opportunities for leakage are growing
  - Query logs, genetics, ...
- De-identification is NOT sufficient for safe disclosure of medical data!

