

Privacy in Process Mining: Motivation, Method and Research Challenges

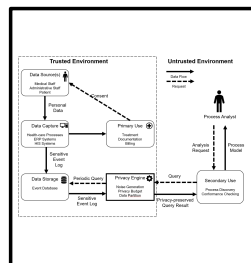
Agnes Koschmider

Group Process Analytics, Kiel University, Germany

Overview

CaseID	Date	Time	Activity	Role
1212				
1212				
1212	10-12-2018	11:00	A	Pera
1212	12-12-2018	09:00	B	Sara
1212	01-01-2019	13:12	C	Milo
1212	08-01-2019	11:38	D	Sara
1212	02-02-2019	14:34	A	Pera
1212	10-12-2018	11:00	A	Pera
1212			C	Milo
1212			B	Sara
1212			D	Sara
1212			A	Pera
1212			F	Milo
1212			C	Milo
1212	06-01-2019	09:18	D	Sara
1212	08-01-2019	11:38	H	Sara
1212	08-01-2019	11:00	B	Sara
1212	08-01-2019	11:43	C	Pera
1212	08-01-2019	09:51	D	Sara
1212	11-01-2019	10:45	G	Ellen

Motivation

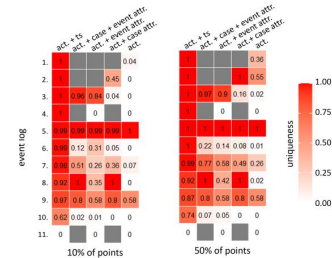
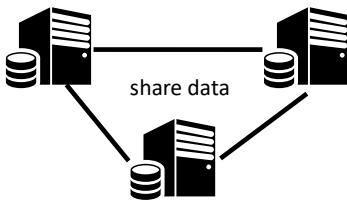


Method



Research Challenges

Privacy Risks



S. Nuñez von Voigt, S.A. Fahrenkrog-Petersen, D. Janssen, A. Koschmider, F. Tschorsch, F. Mannhardt, O. Landsiedel, M. Weidlich: Quantifying the Re-identification Risk of Event Logs for Process Mining - Empirical Evaluation Paper. CAISE 2020: 252-267

06.09.2021

Running Example

	Case attributes			Activity	Timestamp	Event attribute	
	Patient	Birth	Gender	Activity	Timestamp	Doctor	
Case	104	1935	Male	Blood Test	03/03/19 17:43	Dr. Scott	Events
	104	1935	Male	CT	03/05/19 18:15	Dr. Doe	
	104	1935	Male	Surgery	03/07/19 08:23	Dr. Doe	
	104	1935	Male	Rehab	03/10/19 09:36	John Brown	
Case	105	1968	Male	Blood Test	03/03/19 23:28	Dr. Fox	
	105	1968	Male	MRT	03/04/19 23:53	Dr. White	
Case	106	1990	Female	Session	03/03/19 12:34	Dr. Black	
	106	1990	Female	Abortion	03/08/19 16:23	Dr. Scott	
Case	107	1968	Male	Blood Test	03/02/19 18:25	Dr. Scott	
	107	1968	Male	MRT	03/06/19 11:32	Dr. Fox	

Transform Event Log

Case	Birth	Gender	Activity	Timestamp	Doctor
104	1935	Male	[Blood Test, CT, ...]	[03/03/19, 03/05/19, ...]	[Scott, Doe, ...]
105	1968	Male	[Blood Test, MRT, ...]	[03/03/19, 03/04/19, ...]	[Fox, White, ...]
106	1990	Female	[Session, Abortion]	[03/03/19, 03/08/19]	[Black, Scott]
107	1968	Male	[Blood Test, MRT]	[03/02/19, 03/06/19]	[Scott, Fox]

Quantify Uniqueness

Considering case attributes:
given case attribute **Gender**

Case	Birth	Gender	Activity	Timestamp	Doctor
104	1935	Male	[Blood Test, CT, ...]	[03/03/19, 03/05/19, ...]	[Scott, Doe, ...]
105	1968	Male	[Blood Test, MRT, ...]	[03/03/19, 03/04/19, ...]	[Fox, White, ...]
106	1990	Female	[Session, Abortion]	[03/03/19, 03/08/19]	[Black, Scott]
107	1968	Male	[Blood Test, MRT]	[03/02/19, 03/06/19]	[Scott, Fox]

Quantify Uniqueness

Considering case attributes:
given case attribute

Case	Birth	Gender	Activity	Timestamp	Doctor
104	1935	Male	[Blood Test, CT, ...]	[03/03/19, 03/05/19, ...]	[Scott, Doe, ...]
105	1968	Male	[Blood Test, MRT, ...]	[03/03/19, 03/04/19, ...]	[Fox, White, ...]
→ 106	1990	Female	[Session, Abortion]	[03/03/19, 03/08/19]	[Black, Scott]
107	1968	Male	[Blood Test, MRT]	[03/02/19, 03/06/19]	[Scott, Fox]

$1/4 = 0.25 = 25\%$ re-identification risk

7

Quantify Uniqueness

Considering events as points:
 $p_1 = (\text{Activity}_1, \text{Timestamp}_1, \text{Doctor}_1)$

Case	Birth	Gender	Activity	Timestamp	Doctor
104	1935	Male	[Blood Test, CT, ...]	[03/03/19, 03/05/19, ...]	[Scott, Doe, ...]
105	1968	Male	[Blood Test, MRT, ...]	[03/03/19, 03/04/19, ...]	[Fox, White, ...]
106	1990	Female	[Session, Abortion]	[03/03/19, 03/08/19]	[Black, Scott]
107	1968	Male	[Blood Test, MRT]	[03/02/19, 03/06/19]	[Scott, Fox]

7

Quantify Uniqueness

Considering events as points:
 $p_2 = (\text{Activity}_2, \text{Timestamp}_2, \text{Doctor}_2)$

Case	Birth	Gender	Activity	Timestamp	Doctor
104	1935	Male	[Blood Test, CT, ...]	[03/03/19, 03/05/19, ...]	[Scott, Doe, ...]
105	1968	Male	[Blood Test, MRT, ...]	[03/03/19, 03/04/19, ...]	[Fox, White, ...]
106	1990	Female	[Session, Abortion]	[03/03/19, 03/08/19]	[Black, Scott]
107	1968	Male	[Blood Test, MRT]	[03/02/19, 03/06/19]	[Scott, Fox]

7

Quantify Uniqueness

Considering events as points:
 $p_1 = (\text{Activity}_1)$

Case	Birth	Gender	Activity	Timestamp	Doctor
104	1935	Male	[Blood Test, CT, ...]	[03/03/19, 03/05/19, ...]	[Scott, Doe, ...]
105	1968	Male	[Blood Test, MRT, ...]	[03/03/19, 03/04/19, ...]	[Fox, White, ...]
→ 106	1990	Female	[Session, Abortion]	[03/03/19, 03/08/19]	[Black, Scott]
107	1968	Male	[Blood Test, MRT]	[03/02/19, 03/06/19]	[Scott, Fox]

$1/4 = 0.25 = 25\%$ re-identification risk

7

Quantify Uniqueness

Considering events as points:

$$p_1 = (\text{Activity}_1, \text{Timestamp}_1)$$

Case	Birth	Gender	Activity	Timestamp	Doctor
104	1935	Male	[Blood Test, CT, ...]	[03/03/19, 03/05/19, ...]	[Scott, Doe, ...]
105	1968	Male	[Blood Test, MRT, ...]	[03/03/19, 03/04/19, ...]	[Fox, White, ...]
→ 106	1990	Female	[Session, Abortion]	[03/03/19, 03/08/19]	[Black, Scott]
→ 107	1968	Male	[Blood Test, MRT]	[03/02/19, 03/06/19]	[Scott, Fox]

$$2/4 = 0.50 = 50\% \text{ re-identification risk}$$

Quantify Uniqueness

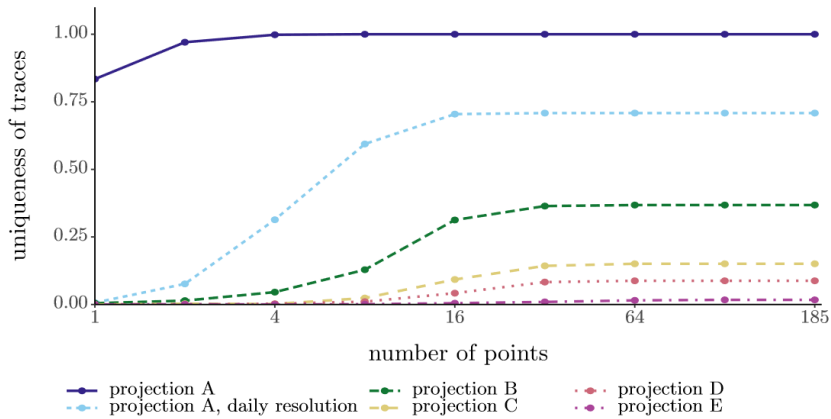
Considering events as points:

$$p_1 = (\text{Activity}_1, \text{Timestamp}_1, \text{Doctor}_1)$$

Case	Birth	Gender	Activity	Timestamp	Doctor
104	1935	Male	[Blood Test, CT, ...]	[03/03/19, 03/05/19, ...]	[Scott, Doe, ...]
105	1968	Male	[Blood Test, MRT, ...]	[03/03/19, 03/04/19, ...]	[Fox, White, ...]
→ 106	1990	Female	[Session, Abortion]	[03/03/19, 03/08/19]	[Black, Scott]
→ 107	1968	Male	[Blood Test, MRT]	[03/02/19, 03/06/19]	[Scott, Fox]

$$4/4 = 1.00 = 100\% \text{ re-identification risk}$$

Uniqueness for Cases of Sepsis Event log



Requirements for Privacy-Preserving Process Mining Techniques



Anonymity



Unlinkability



Transparency



Notice

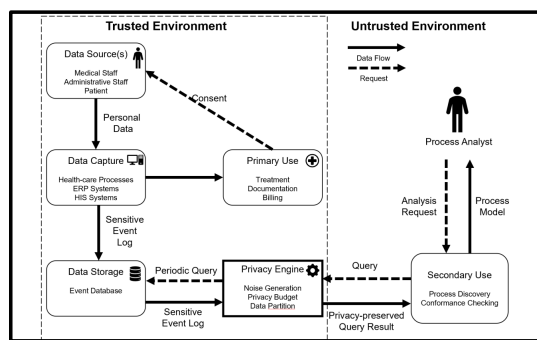


Accountability

	Anonymity	Unlinkability	Notice	Transparency	Accountability
TLKC	X	X			
PRETSA	X	X			
PPPM	X	X			
PRIPEL	X	X			
Multi-party computation	X				

- Beside the requirements of process mining techniques, also data, application and presentation are requirements

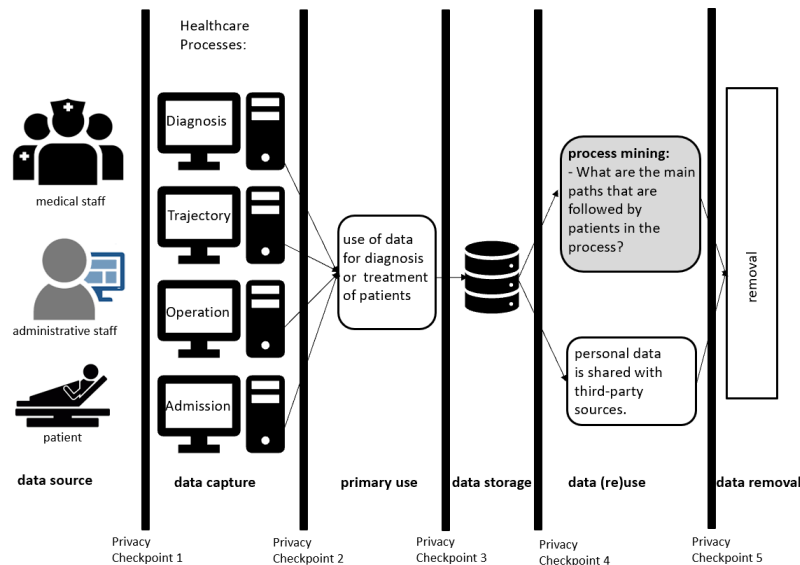
Privacy Preserving Process Mining



F. Mannhardt, A. Koschmider, N. Baracaldo, M. Weidlich, J. Michael: *Privacy-Preserving Process Mining: Differential Privacy for Event Logs*, Business & Information Systems Engineering 61(5), 2019

J. Michael, A. Koschmider, F. Mannhardt, N. Baracaldo, B. Rumpe: *User-Centered and Privacy-Driven Process Mining System Design for IoT*. CAiSE Forum 2019: 194-206, Springer

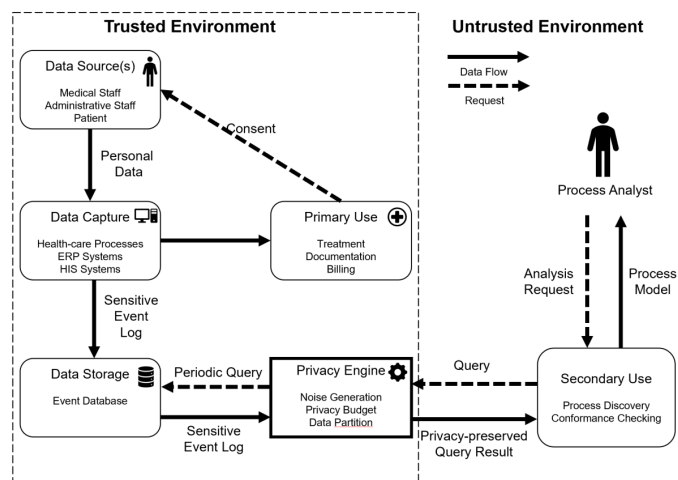
Identification of data passes and privacy checkpoints for hospital health processes



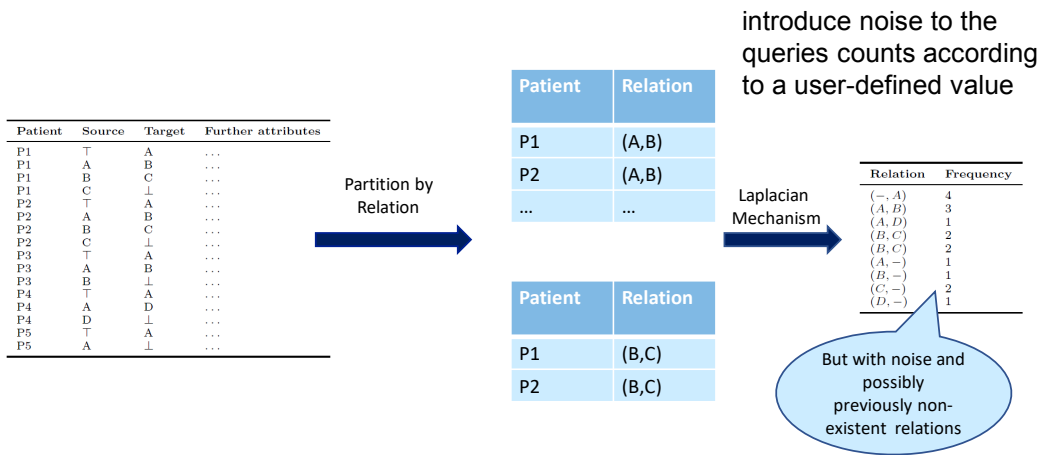
Laplacian mechanism is used to provide differential privacy for counting the number of records in a database

Privacy Model

- we assume a **centralized privacy approach**
- sensitive data is stored as an event log in protected data storage
- **privacy engine** acts as the **single point of access** for process mining algorithms and **introduces noise** to each query result
- no difference for data provider between the data used by the process mining algorithm regardless of whether his/her data is included or not



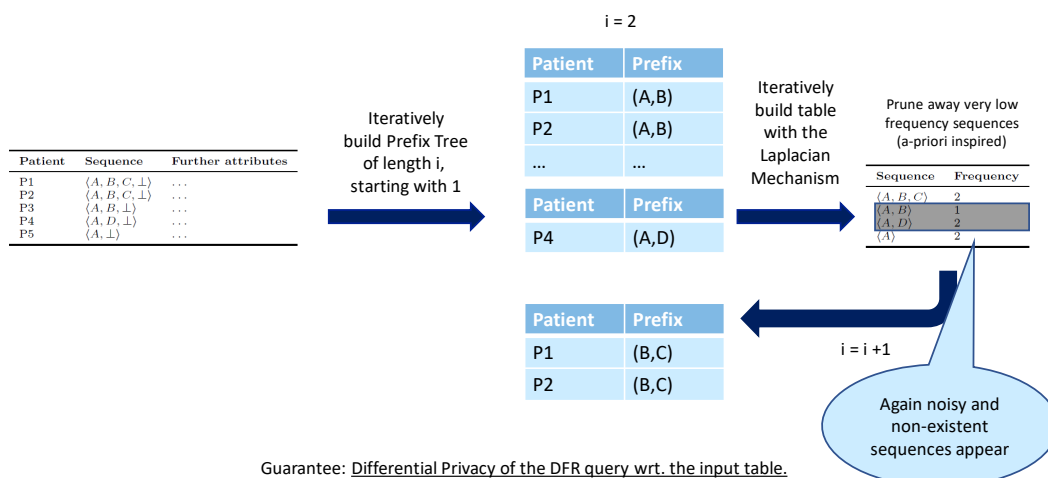
Our Initial Approach – Directly-Follows Relation (DFR)



Guarantee: Differential Privacy of the DFR query wrt. the input table.

if one would sequentially query information from the same data source, the privacy budget is reduced by the sum of the individual parameters

Our Initial Approach – Activity Sequences



Guarantee: Differential Privacy of the DFR query wrt. the input table.

we treat each trace as a sequence of identifiers

Current Steps

- Development of a log generator for synthetic, privacy-preserving event logs
 - use of Generative Adversarial Network (GAN)

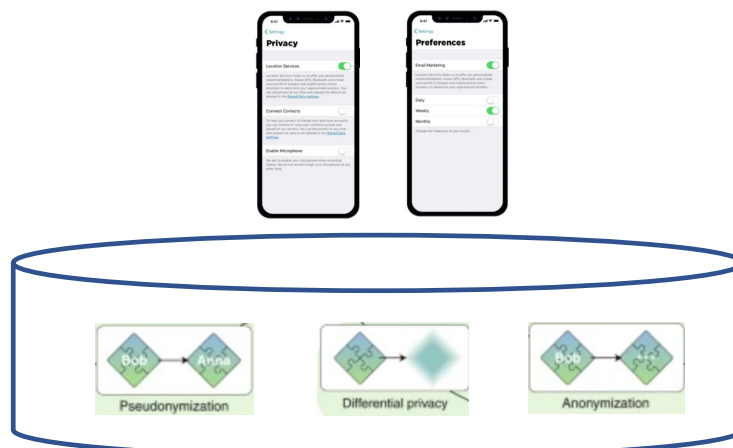


- Noise/Outlier quantification model



Challenges: Interpretable Quantification of Privacy Disclosure

- more reliable and interpretable metrics of privacy disclosure

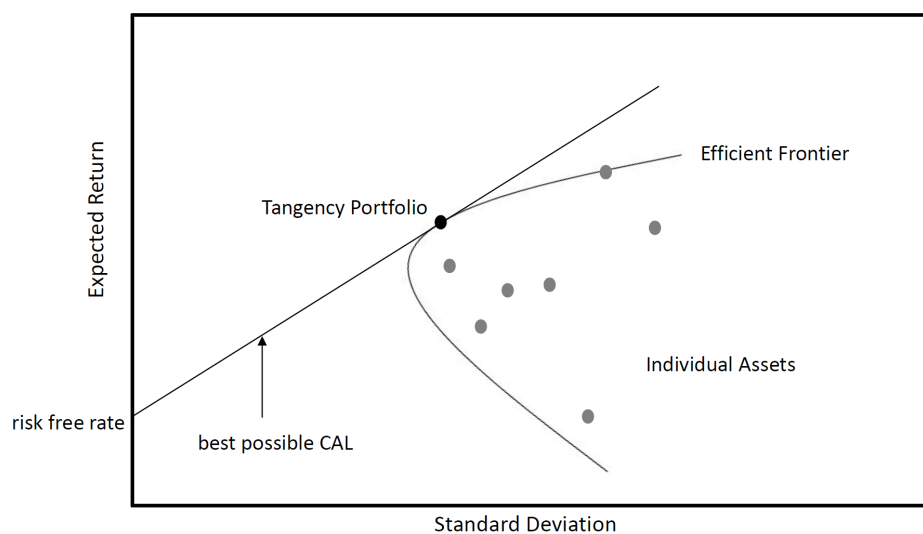


Challenges: Balancing Risk and Utility

- Trade-off between disclosure risk and utility

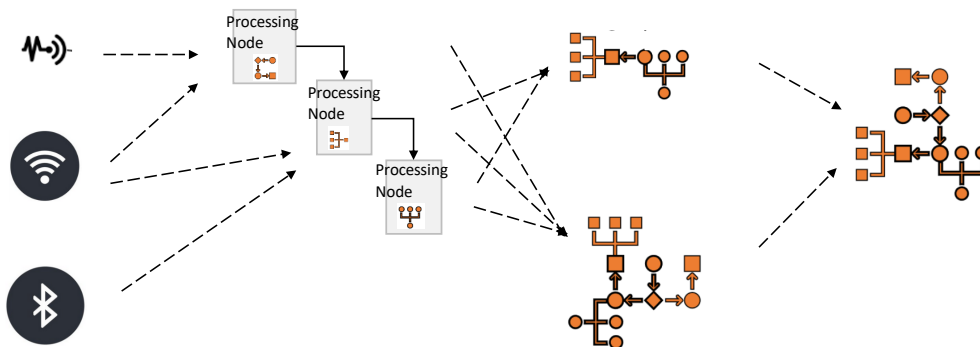
Elkoumy, G., Fahrenkrog-Petersen, S. A., Sani, M. F., Koschmider, A., Mannhardt, F., Voigt, S. N. V., Rafiei, M., & Waldthausen, L. V. Privacy and Confidentiality in Process Mining - Threats and Research Challenges. *ACM Transactions on Management Information Systems*, 2021, in press.

Challenges: Balancing Risk and Utility



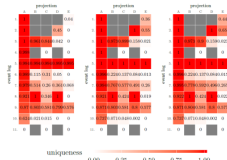
source: https://en.wikipedia.org/wiki/Efficient_frontier

- Distributed Privacy

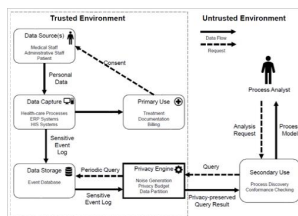


- Computational Challenges:
 - with increasing dimensions of attributes, it becomes more unpractical to achieve privacy-preserving process mining
- Traceability Challenge:
 - trace data-life cycle and ensure consent, right to be forgotten
- Transparency Challenge
 - notify who is using the data

Summary and Outlook



S. Nuñez von Voigt, S.A. Fahrenkrog-Petersen, D. Janssen, A. Koschmider, F. Tschorsch, F. Mannhardt, O. Landsiedel, M. Weidlich: Quantifying the Re-identification Risk of Event Logs for Process Mining. CAISE 2020: 252-267



F. Mannhardt, A. Koschmider, N. Baracaldo, M. Weidlich, J. Michael: Privacy-Preserving Process Mining - Differential Privacy for Event Logs. Business & Information Systems Engineering 61(5): 595-614 (2019)



G. Elkoumy, S.A. Fahrenkrog-Petersen, M. Fani Sani, A. Koschmider, F. Mannhardt, S. Nuñez von Voigt, M. Rafiei, L. von Waldhausen: Privacy and Confidentiality in Process Mining - Threats and Research Challenges, ACM Transactions of Management Information Systems, 2021