Privacy-preserving and Authenticated Data Cleaning on Outsourced Databases Thesis Defense

Boxiang Dong

THESIS COMMITTEE: Advisor: Prof. Wendy Hui Wang Prof. Yingying Chen Prof. David Naumann Prof. Antonio Nicolosi

> Department of Computer Science Stevens Institute of Technology

> > December 1, 2016

Dirty Data

Real-world datasets, particularly those from multiple sources, tend to be *dirty*.

Inaccuracy Multiple records that refer to the same entity

Inconsistency Violation of integrity constraints

Incompleteness Missing data values

Name	Street	City	Phone
John	Leonard	NY	518-457-5181
John	Lenard	NY	518-457-5181
Kevin		LA	213-974-3211
Mike	Main	Phil	518-457-5181

The ubiquitous dirty data: 40% of companies have suffered losses, problems, or costs due to data of poor quality [Eck02].

Data cleaning aims at detecting and removing errors, duplications, missing values, and inconsistencies to improve data quality.

- Data deduplication
- Data inconsistency repair
- Data imputation

Data cleaning is a labor-intensive and complex process. It can be NP-complete [BFFR05].

Data-Cleaning-as-a-Service

Outsourcing the data to a third-party data cleaning service provider provides a cost-effective way. E.g., Google's OpenRefine, Melissa Data.



Client with limited computational resources Server computationally powerful The third-party server is untrusted.

Result integrity The server may return incorrect data cleaning result.

- Software bugs
- Intention to save computational cost

Data privacy The outsourced data may include sensitive personal information.

- Medical information
- Financial record











Related Work

Data cleaning

- Data deduplication [GIJ⁺01, SAA10, YLKG07]
- Data inconsistency repair [PEM⁺15, BFG⁺07, BFFR05]

Privacy-preserving outsourced computation

- Encryption [SV10, PRZB12]
- Encoding [EAMY⁺13, CC04]
- Secure multiparty computation [TOEY11, LZL+15]
- Differential privacy [CMF⁺11, AHMP15]

Verifiable computing

- General-purpose verifiable computing [SVP+12, PHGR13]
- Function-specific verifiable computing [DLW13, LWM⁺12]

Outline

Introduction

- **2** Research Results
 - Authentication of Outsourced Data Deduplication
 - Verification of Similarity Search Approach (VS²)
 - Embedding-based Verification of Similarity Search Approach (*E*-*VS*²)
 - Experiments
 - Privacy-preserving Outsourced Data Deduplication
 - Privacy-preserving Outsourced Data Inconsistency Repair
- **3** Research beyond the Thesis
- Future Plan
- G Conclusion

Authentication of Outsourced Data Deduplication Boxiang Dong, Wendy Hui Wang. IEEE International Conference on Information Reuse and Integration (IRI), Pittsburgh, PA. July 2016. (Acceptance rate = 25%)

Data Deduplication

Data deduplication Eliminate near-duplicate copies.

• Record matching: Detect near-duplicate copies.



Data deduplication Eliminate near-duplicate copies.

• Record matching: Detect near-duplicate copies.

RID	Name	Street	City	Age
<i>r</i> ₁	John	Leonard	NY	45
<i>r</i> ₂	Kevin	Wicks	LA	31
r ₃	Mike	Main	Phil	22
			$\theta =$	= 2
			{ <i>r</i>	1}

 $s_q = (\text{John, Lenard, NY, 45})$

Outsourcing Framework

The client (data owner) outsources the record matching service to the untrusted server.



Assumption: The client is aware of the edit distance metric. We want to make sure that R^{S} is both sound and complete. Soundness $\forall s \in R^{S}$, $s \in D$ and $DST(s, s_q) \leq \theta$. Completeness $\forall s \in D$ s.t. $DST(s, s_q) \leq \theta$, $s \in R^{S}$.

Authentication

We aim at an authentication framework that satisfies the following objectives.



Preliminary - Merkle Tree

Merkle tree is a generalization of hash lists and hash chains.



- It allows efficient and secure verification of the contents of large data structures.
- Hash is computationally more efficient than edit distance calculation.

Preliminary - *B^{ed}*-Tree

 B^{ed} -Tree [ZHOS10] is a string indexing structure.



- Sort the strings in dictionary order.
- Store the longest common prefix (LCP) of the enclosed strings in every node.

Preliminary - B^{ed}-Tree

B^{ed}-Tree [ZHOS10] is a string indexing structure.



• $\forall N$, calculate $MIN_DST(s_q, N.LCP)$.

Preliminary - *B^{ed}*-Tree

B^{ed}-Tree [ZHOS10] is a string indexing structure.



- If $MIN_DST(s_q, N.LCP) > \theta$, then N is a MF-node.
- All strings covered by a MF-node must be dissimilar to sq.
- Avoid the edit distance calculation for NC-strings.
- Perform well with memory constraints.

Preliminary - Embedding

Embedding maps strings into Euclidean points in a similarity-preserving way.



- Euclidean distance calculation is much more efficient than edit distance computing, i.e., O(dst(p_i, p_j)) << O(DST(s_i, s_j)).
- SparseMap[HS] is a contractive embedding approach, i.e., dst(p_i, p_j) ≤ DST(s_i, s_j).
- The complexity is $O(cn^2)$, where c is a small constant, and n is the number of strings.

Solution in a Nutshell

We require the server to construct *verification object (VO)* to demonstrate the soundness and completeness of the result.



The client is able to efficiently detect any unsound or incomplete result returned by the server by checking the VO.

Outline

- Introduction
- **2** Research Results
 - Authentication of Outsourced Data Deduplication
 - Verification of Similarity Search Approach (*VS*²)
 - Embedding-based Verification of Similarity Search Approach (*E*-*VS*²)
 - Experiments
 - Privacy-preserving Outsourced Data Deduplication
 - Privacy-preserving Outsourced Data Inconsistency Repair
- **3** Research beyond the Thesis
- 4 Future Plan
- G Conclusion

VS^2 - Setup

We propose an authenticated string indexing structure, named MB-tree (Merkle B^{ed} -tree).



- The client signs the hash value in the root, and only keeps the signature of the *MB*-tree locally.
- The hash function is more efficient than edit distance calculation.

VS²-VO Construction

The server searches for the similar strings and constructs VO by traversing the MB-tree.



- Include all the C-strings and similar strings in VO.
- Substitute the large amount of NC-strings with the MF-nodes.

VS² - VO Verification

The client checks the soundness of completeness of R^{S} by verifying the *VO*.



VS² - VO Verification

The client checks the soundness and completeness of R^{S} by verifying the *VO*.



VS^2 - VO Verification

The client checks the soundness and completeness of R^{S} by verifying the *VO*.

		$\exists s \in \mathbb{R}^S$, but $s \notin D$	0	Compute $Sig(T)$ from VO
catches {		$\exists s \in R^S, \text{but} \; DST(s,s_q) > \theta$	ľ	$\forall s \in R^S$, check if $DST(s, s_q) \leq \theta$
	completeness violation	$ \exists s \in D \text{ s.t. } DST(s,s_q) \leq \theta \\ \text{but } s \not \in R^S $		$\label{eq:c-string} \begin{array}{l} \mbox{σ}, \mbox{ check if } DST(s,s_q) > \theta \\ \\ \mbox{π} \end{array}$

$ \begin{aligned} s_q &= \text{``Celestine ''} \\ \theta &= 4 \end{aligned} $	$R^{S} = \{s_{1}, s_{2}\}$ $VO = \{(((s_{1}, s_{2}, s_{3}), (s_{4}, s_{5}, s_{6})), ((s_{4}, s_{5}, s_{6}))))$	$(1, s_8, s_9), (LCP_{N_7}, h_{N_7})))\}$
for similar strings	$\left\{ \begin{array}{l} DST(s_1,s_q)=4\\ DST(s_2,s_q)=3<4 \end{array} \right.$	
for C-strings	$\begin{cases} DST(s_3, s_q) = 5 > 4\\ DST(s_4, s_q) = 9 > 4\\ DST(s_5, s_q) = 9 > 4\\ DST(s_6, s_q) = 8 > 4\\ DST(s_6, s_q) = 8 > 4\\ DST(s_8, s_q) = 8 > 4\\ DST(s_9, s_q) = 8 > 4 \end{cases}$	> 10 DST calculations
for MF-node	$MIN_DST(LCP_{N_7}, s_q) = 6 > 4$)

29/61

Outline

- Introduction
- **2** Research Results
 - Authentication of Outsourced Data Deduplication
 - Verification of Similarity Search Approach (VS²)
 - Embedding-based Verification of Similarity Search Approach (E-VS²)
 - Experiments
 - Privacy-preserving Outsourced Data Deduplication
 - Privacy-preserving Outsourced Data Inconsistency Repair
- **3** Research beyond the Thesis
- Future Plan
- G Conclusion

- The client constructs the MB-tree.
- The client applies *SparseMap* to embed strings into Euclidean points.



Key idea For any C-string s, if $dst(p, p_q) > \theta$, it must be true that $DST(s, s_q) > \theta$.

Distant Bounding Hyper-rectangle (DBH) A hyper-rectangle R in the Euclidean space is a DBH if $min_dst(p_q, R) > \theta$.

- **DBH-String** For any C-string *s*, if $dst(p, p_q) > \theta$, we call it a DBH-string.
- **FP-String** For any C-string *s*, if $dst(p, p_q) \le \theta$, we call it a FP-string.
 - Key idea
 To save the verification cost at the client side, the server should organize the set of DBH-strings into a small number of DBHs.
 - By only checking the Euclidean distance between the query point p_q and the DBHs, the client assures that all *DBH*-strings are dis-similar to s_q .

*E-VS*² - VO Construction



*E-VS*² - VO Construction



Theorem (NP-Completeness of DBH Construction)

Given a query string s_q , and a set of DBH-strings $\{s_1, \ldots, s_t\}$, let $\{p_1, \ldots, p_t\}$ be their Euclidean points. It is a NP-complete problem to construct a mimimum number of rectangles $\mathcal{R} = \{R_1, \ldots, R_k\}$ s.t. (1) $\forall i \neq j, R_i$ and R_j do not overlap; and (2) $\forall p_i$, there exists a R_i s.t. p_i is included in R_i .

- We design an efficient heuristic algorithm for the server to construct a small amount of *DBH*s.
- The complexity is cubic to the number of DBH-strings.

$E-VS^2$ - VO Construction

 p_q

 p_{12}

 p_{10}

The server includes the DBHs in the VO.



$$VO = \{(((s_1, s_2, (s_3, p_{R_1})), ((s_4, p_{R_2}), (s_5, p_{R_1}), (s_6, p_{R_1}))), \\ (((s_7, p_{R_2}), (s_8, p_{R_1}), s_9), (LCP_{N_7}, h_{N_7}))), \{R_1, R_2\}\}$$

*E-VS*² - VO Verification

The client checks the soundness and completeness of R^{S} by verifying the *VO*.



*E-VS*² - VO Verification

The client checks the soundness and completeness of R^{S} by verifying the *VO*.



Complexity Analysis

Phase	Measurement	VS^2	E-VS ²	
Sotup	Time	<i>O</i> (<i>n</i>)	O(cdn ²)	
Jetup	Space	<i>O</i> (<i>n</i>)	<i>O</i> (<i>n</i>)	
VO Construction	Time	<i>O</i> (<i>n</i>)	$O(n+n_{DS}^3)$	
VO COnstruction	VO Size	$(n_R + n_C)\sigma_S + n_{MF}\sigma_M$	$(n_R + n_C)\sigma_S + n_{MF}\sigma_M + n_{DBH}\sigma_D$	
VO Verification	Time $O((n_R + n_{MF} + n_C)C_{Ed})O((n_R + n_{MF} + n_{FP})C_{Ed} + n_{DBH}C_{Ed})$			
(<i>n</i> : # of strings in <i>D</i> ; <i>c</i> : a constant in [0, 1]; <i>d</i> : # of dimensions of Euclidean space; σ_{S} : the average length of the string; σ_{M} : Avg. size of a <i>MB</i> -tree node;				
σ_D : Avg. size of a DBH; n_R : # of strings in M^S ; n_C : # of C-strings;				
n _{FP} : # of FP-strings; n _{DS} : # of DBH-strings; n _{DBH} : # of DBHs;				
n_{MF} : # of MF nodes; C_{Ed} : the complexity of an edit distance computation;				
(C_{FI} : the comp	lexitv of Euclidean distan	ce calculation.)	

- *E-VS*² results in higher VO construction complexity at the server side.
- *E-VS*² dramatically saves the VO verification cost at the client side.

Outline

Introduction

- **2** Research Results
 - Authentication of Outsourced Data Deduplication
 - Verification of Similarity Search Approach (VS²)
 - Embedding-based Verification of Similarity Search Approach (*E-VS*²)
 - Experiments
 - Privacy-preserving Outsourced Data Deduplication
 - Privacy-preserving Outsourced Data Inconsistency Repair
- **3** Research beyond the Thesis
- Future Plan
- G Conclusion

Experiments - Setup

• Environment

Language C++ Testbed A Linux machine with 2.4 GHz CPU and 48 GB RAM

Datasets

Actors ¹ 260,000 lastnames Authors ² 1,000,000 full names

- Evaluation metric
 - VO construction time
 - VO verification time

¹http://www.imdb.com/interfaces
²http://dblp.uni-trier.de/xml/

Experiments - VO Construction Time

Time Performance of VO Construction



• $E-VS^2$ takes more time at the server side to construct VO, especially when θ is small.

Experiments - VO Verification Time

Time Performance of VO Verification



- VS² and E-VS² are significantly more efficient than the baseline approach in verification cost.
- The advantage of E- VS^2 is large when θ is small.

Outline

Introduction

- **2** Research Results
 - Authentication of Outsourced Data Deduplication
 - Verification of Similarity Search Approach (VS²)
 - Embedding-based Verification of Similarity Search Approach (*E*-*VS*²)
 - Experiments
 - Privacy-preserving Outsourced Data Deduplication
 - Privacy-preserving Outsourced Data Inconsistency Repair
- **3** Research beyond the Thesis
- Future Plan
- G Conclusion

α -Security against Frequency Analysis (FA) Attack ³

Define α -security to limit the success probability of frequency analysis attack.

Experiment $Exp_{A,\Pi}^{FA}()$ $p' \leftarrow A^{freq_{\epsilon}(e), freq(\mathcal{P})}$ Return 1 if p' = Decrypt(k, e)Return 0 otherwise

 α -security against FA attack if $Pr[Exp_{A,\Pi}^{FA}()=1] \leq \alpha$

Prada: Privacy-preserving Data-Deduplication-as-a-Service.

³Boxiang Dong, Ruilin Liu, Wendy Hui Wang.

International Conference on Information and Knowledge Management, 2014. (Acceptance rate=20%).

Privacy-preserving Outsourced Data Deduplication ⁴

We design two approaches to enable data deduplication and defend against the frequency analysis attack.

- Locality-sensitive Hashing Based Approach (LSHB)
- Embedding & Homomorphic Substitution Approach (EHS)



LSHB approach encodes strings into LSH values that

(1) preserve the string similarity; and

(2) are of the same frequency groupwise.

⁴Boxiang Dong, Ruilin Liu, Wendy Hui Wang. Prada: Privacy-preserving Data-Deduplication-as-a-Service. International Conference on Information and Knowledge Management, 2014. (Acceptance rate=20%).



EHS approach encodes strings into Euclidean points that

(1) preserve the string similarity; and

(2) are of uniform frequency.

Privacy-preserving Outsourced Data Deduplication

Experiment Results



Outline

Introduction

2 Research Results

- Authentication of Outsourced Data Deduplication
 - Verification of Similarity Search Approach (VS²)
 - Embedding-based Verification of Similarity Search Approach (*E*-*VS*²)
 - Experiments
- Privacy-preserving Outsourced Data Deduplication
- Privacy-preserving Outsourced Data Inconsistency Repair
- **3** Research beyond the Thesis
- 4 Future Plan
- G Conclusion

Functional dependency (FD) $X \to Y$ if $r_1[X] = r_2[X]$, then $r_1[Y] = r_2[Y]$.

FDs play a key role in identifying and fixing data inconsistency.

TID	Conference	Year	Country	Capital	City
<i>r</i> ₁	SIGMOD	2007	China	Beijing	Beijing
<i>r</i> ₂	ICDM	2014	China	Shanghai	Shenzhen
<i>r</i> ₃	KDD	2014	U.S.	Washington D.C.	New York City
<i>r</i> ₄	KDD	2015	Australia	Canberra	Sydney
<i>r</i> 5	ICDM	2015	U.S.	New York City	Atlantic City
		FD :	Country -	> Capital	

Indistinguishability against FD-preserving Chosen Plaintext Attack (IND-FCPA)

$$\begin{aligned} \textbf{Experiment} \ Exp_{A,\Pi}^{IND-FCPA}(\lambda) \\ & k \leftarrow KeyGen(\lambda) \\ (D_0, D_1) \leftarrow A^{O_{Encrypt(.)}}(k) \text{ s.t. } FD_0 = FD_1 \text{ and } |D_0| = |D_1| \\ & b \xleftarrow{\$} \{0, 1\} \\ & b' \leftarrow A^{O_{Encrypt(.)}}(k) \\ & \text{Return 1 if } b = b' \\ & \text{Return 0 otherwise} \end{aligned}$$

IND-FCPA if $\Pr[Exp_A^{IND-FCPA}(n)=1] \leq \frac{1}{2} + negl(n)$

Privacy-preserving Outsourced Data Inconsistency Repair

We consider two scenarios of the outsourced data inconsistency repair, and design two encryption/encoding approaches to provide robust privacy guarantee 5 .



⁵Boxiang Dong, Wendy Hui Wang, Jie Yang.

Secure Data Outsourcing with Adversarial Data Dependency Constraints.

International Conference on Big Data Security on Cloud, 2016. (Acceptance rate=23%).

Outline

Introduction

2 Research Results

- Authentication of Outsourced Data Deduplication
 - Verification of Similarity Search Approach (VS²)
 - Embedding-based Verification of Similarity Search Approach (*E*-*VS*²)
 - Experiments
- Privacy-preserving Outsourced Data Deduplication
- Privacy-preserving Outsourced Data Inconsistency Repair
- **3** Research beyond the Thesis
- 4 Future Plan
- G Conclusion

Research beyond the Thesis

- Authentication of outsourced data mining computations
 - Association rule mining [DBSec'13, ICDM'13, TSC'15]
 - Outlier mining (under review)
- Rank aggregation in the crowdsourcing setting (under review)
 - Rank inference
 - Task assignment with data privacy concern
- Data-as-a-commodity (under review)
 - Budget constraint
 - High quality (low inconsistency)

Outline

Introduction

2 Research Results

- Authentication of Outsourced Data Deduplication
 - Verification of Similarity Search Approach (VS²)
 - Embedding-based Verification of Similarity Search Approach (*E*-*VS*²)
 - Experiments
- Privacy-preserving Outsourced Data Deduplication
- Privacy-preserving Outsourced Data Inconsistency Repair
- **3** Research beyond the Thesis
- Future Plan
- G Conclusion

 Authenticated outsourced data inconsistency repair Challenge It is NP-complete to find a repair with the minimum cost.

Solution

- Convert the strings into Euclidean space.
- It is the *center of mass* that results in the smallest repair cost.
- Authenticated outsourced data imputation

Challenge It demands a similarity matrix between all values. Solution Create evidence imputation objects to verify the result in a probabilistic way. Privacy-preserving and authenticated data cleaning on outsourced databases.

- Define two security notions, namely α -security and *IND-FCPA*.
- Authentication of outsourced data deduplication.
- Privacy-preserving outsourced data deduplication.
- Privacy-preserving outsourced data inconsistency repair.
 - Privacy against FD attack.
 - Privacy against frequency analysis attack.

The suit of encryption, encoding, and authentication schemes address the security and privacy concerns in outsourced computing.

My Publications

IRI'16	Boxiang Dong, Hui (Wendy) Wang. $\overline{ARM:}$ Authenticated Approximate Record Matching for Outsourced Databases. IEEE International Conference on Information Reuse and Integration (IRI). Pittsburgh, PA. 2016. (Accentance rate = 25%).
BigDataSecurity'16	Boxiang Dong, Hui (Wendy) Wang, Jie Yang. Secure Data Outsourcing with Adversarial Data Dependency Constraints. IEEE International Conference on Big Data Security on Cloud (BigDataSecurity).
TSC'15	New York. 2016. (Acceptance rate = 23%). Boxiang Dong, Ruilin Liu, Hui (Wendy) Wang. Trust-but-Verify: Verifying Result Correctness of Outsourced Frequent Itemset
СІКМ'14	Mining. IEEE Transactions on Services Computing. 2015. Boxiang Dong, Ruilin Liu, Hui (Wendy) Wang. Prada: Privacy-preserving Data-Dedunlication-as-a-Service
ICDM'13	ACM International Conference on Information and Knowledge Management (CIKM). Shanghai, China. 2014. (Acceptance rate = 20%). Boxiang Dong, Ruilin Liu, Hui (Wendy) Wang. Interrity Verification of Outsourced Frequent Itemset Mining with Deterministic
DBSec'13	Guarantee. IEEE International Conference on Data Mining (ICDM). Dallas, Texas. 2013. (Acceptance rate = 19.7%). Boxiang Dong, Ruilin Liu, Hui (Wendy) Wang.
	Result Integrity Verification of Outsourced Frequent Itemset Mining. Annual IFIP WG 11.3 Conference on Data and Application Security and Privacy (DBSec). Newark, NJ. 2013.
IJIPM'10	Weifeng Sun, Juanyun Wang, <u>Boxiang Dong</u> , Mingchu Li, Zhenquan Qin. A Mediated RSA-based End Entity Certificates Revocation Mechanism with Secure Concern in Grid. International Journal of Information Processing and Machanet (JURN). 2010.
IIH-MSP'10	Weifeng Sun, <u>Boxiang Dong</u> , Zhenquan Qin, Juanyun Wang, Mingchu Li. A Low-Level Security Solving Method in Grid. International Conference on Intelligent Information Hiding and Multimedia Signal Processing (IIH-MSP). Darmstadt, Germany. 2010.

References I

[AHMP15] Tristan Allard, Georges Hébrail, Florent Masseglia, and Esther Pacitti. Chiaroscuro: Transparency and privacy for massive personal time-series clustering. In Proceedings of the ACM SIGMOD International Conference on Management of Data. pages 779-794, 2015. [BFFR05] Philip Bohannon, Wenfei Fan, Michael Flaster, and Raieev Rastogi, A cost-based model and effective heuristic for repairing constraints by value modification. In Proceedings of the ACM SIGMOD International Conference on Management of Data, pages 143-154, 2005. [BFG⁺07] Philip Bohannon, Wenfei Fan, Floris Geerts, Xibei Jia, and Anastasios Kementsietsidis. Conditional functional dependencies for data cleaning. In IEEE International Conference on Data Engineering, pages 746-755, 2007. [CC04] Tim Churches and Peter Christen. Some methods for blindfolded record linkage. BMC Medical Informatics and Decision Making, 4(1):9, 2004. [CMF⁺11] Rui Chen, Noman Mohammed, Benjamin CM Fung, Bipin C Desai, and Li Xiong. Publishing set-valued data via differential privacy. Proceedings of the VLDB Endowment, 4(11):1087-1098, 2011. [DLW13] Boxiang Dong, Ruilin Liu, and Hui Wendy Wang. Result integrity verification of outsourced frequent itemset mining. In Data and Applications Security and Privacy XXVII, pages 258–265, 2013. [EAMY⁺13] Durham E, Ashley, Kantarcioglu M., Xue Y., Kuzu M., and Malin Bradley. Composite bloom filters for secure record linkage. In IEEE Transactions on Knowledge and Data Engineering, 2013.

References II

[Eck02]	Wayne W Eckerson. Data quality and the bottom line. The Data Warehouse Institute Report, 2002.
[GIJ ⁺ 01]	Luis Gravano, Panagiotis G Ipeirotis, Hosagrahar Visvesvaraya Jagadish, Nick Koudas, Shanmugauelayut Muthukrishnan, Divesh Srivastava, et al. Approximate string joins in a database (almost) for free. In <i>Proceedings of the International Conference on Very Large Data Bases</i> , volume 1, pages 491–500, 2001.
[HS]	G Hjaltason and H Samet. Contractive embedding methods for similarity searching in metric spaces. Technical report, Computer Science Department, University of Maryland.
[LWM ⁺ 12]	Ruilin Liu, Hui Wendy Wang, Anna Monreale, Dino Pedreschi, Fosca Giannotti, and Wenge Guo. Audio: An integrity auditing framework of outlier-mining-as-a-service systems. In Machine Learning and Knowledge Discovery in Databases, pages 1–18. 2012.
[LZL ⁺ 15]	An Liu, Kai Zhengy, Lu Liz, Guanfeng Liu, Lei Zhao, and Xiaofang Zhou. Efficient secure similarity computation on encrypted trajectory data. In IEEE International Conference on Data Engineering, pages 66–77, 2015.
[PEM ⁺ 15]	Thorsten Papenbrock, Jens Ehrlich, Jannik Marten, Tommy Neubert, Jan-Peer Rudolph, Martin Schönberg, Jakob Zwiener, and Felix Naumann. Functional dependency discovery: An experimental evaluation of seven algorithms. <i>Proceedings of the VLDB Endowment</i> , 8(10):1082–1093, 2015.
[PHGR13]	Bryan Parno, Jon Howell, Craig Gentry, and Mariana Raykova. Pinocchio: Nearly practical verifiable computation. In IEEE Symposium on Security and Privacy (SP), pages 238–252, 2013.

References III

[PRZB12] Raluca Ada Popa, Catherine Redfield, Nickolai Zeldovich, and Hari Balakrishnan. Cryptdb: Processing queries on an encrypted database. Communications of the ACM, 55(9):103-111, 2012. [SAA10] Yasin N Silva, Walid G Aref, and Mohamed H Ali, The similarity join database operator. In IEEE International Conference on Data Engineering, volume 10, pages 892-903, 2010. [SV10] Nigel P Smart and Frederik Vercauteren. Fully homomorphic encryption with relatively small key and ciphertext sizes. In Public Key Cryptography-PKC, pages 420-443. 2010. [SVP⁺12] Srinath Setty, Victor Vu, Nikhil Panpalia, Benjamin Braun, Andrew J Blumberg, and Michael Walfish Taking proof-based verified computation a few steps closer to practicality. In The USENIX Security Symposium, pages 253-268, 2012. [TOEY11] Nilothpal Talukder, Mourad Ouzzani, Ahmed K Elmagarmid, and Mohamed Yakout. Detecting inconsistencies in private data with secure function evaluation. Technical report, Computer Science Department, Purdue University, 2011. [YLKG07] Su Yan, Dongwon Lee, Min-Yen Kan, and Lee C Giles. Adaptive sorted neighborhood methods for efficient record linkage. In Proceedings of the ACM/IEEE-CS Joint Conference on Digital Libraries, pages 185-194, 2007. [ZHOS10] Zhenjie Zhang, Marios Hadjieleftheriou, Beng Chin Ooi, and Divesh Srivastava. Bed-tree: an all-purpose index structure for string similarity search based on edit distance. In Proceedings of the International Conference on Management of Data, 2010.



Thank you!

Questions?