

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/181065/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Bian, Pengxin, Theodorakopoulos, Georgios , Pissis, Solon P. and Loukides, Grigorios 2025. Optimal string sanitization against strategic attackers. IEEE Transactions on Information Forensics and Security

Publishers page:

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



Supplemental Material: Optimal String Sanitization Against Strategic Attackers

Pengxin Bian, *Graduate Student Member, IEEE*, George Theodorakopoulos,
Solon P. Pissis, and Grigorios Loukides, *Senior Member, IEEE*

I. ADDITIONAL RELATED WORK

There are two main directions in privacy-preserving data publishing; data sanitization and anonymization. The former has been reviewed in Section II of the main paper. The latter aims to prevent the inference of information about individuals represented in a disseminated dataset [1]. Below we discuss data anonymization approaches as a supplement to Section II of the main paper.

Instead of protecting the privacy of sensitive patterns modeling confidential information as data sanitization does, data anonymization aims to protect information about individuals. For example, anonymization aims to prevent the disclosure of individuals' identities or attributes these individuals are not willing to be associated with. Existing string anonymization approaches propose algorithms based on two privacy principles: k -anonymity [2] and differential privacy [3].

To enforce k -anonymity, several algorithms [4]–[6] apply *condensation*, a methodology which splits a dataset into multiple groups of size at least k and then releases aggregate statistics about each group. These algorithms aim to create groups of similar entities, so that the quality loss incurred by releasing the group statistics is small. For example, [5], [6] are applied to a collection of strings and create groups of at least k similar strings. They represent each group of strings using summary statistics containing first and second order information about the distribution of the letters in the strings, while [6] also discusses the possibility of releasing third order information.

The algorithms based on differential privacy could be classified into those that are applied to a collection of strings and those that are applied to a single string. [7] proposes an algorithm to extract differentially-private variable-length n -grams with counts larger than a threshold and lengths larger than another threshold from a collection of strings. [8] releases a set of differentially-private top- k frequent strings from a collection of strings, where k denotes the number of frequent strings required, while [9] releases a set of differentially-private frequent subsequences from a

collection of strings under a maximum gap constraint. The approach of [10] proposes a sampling-based mechanism to enforce differential privacy on a single string and several algorithms, exact and heuristics, to achieve this while preserving data utility. The privacy goal of this approach is to prevent an attacker from inferring the presence in the disseminated string or absence from the disseminated string of any q -gram.

Since anonymization algorithms do not aim to hide sensitive patterns, they are not alternatives to our algorithms.

II. CONTEXTS CONSTRUCTION

Given a sensitive pattern $s \in S$, we find all its occurrences in the reference string $\mathcal{R} = \mathcal{R}[1] \dots \mathcal{R}[n]$ using the pattern matching algorithm of [11]. Then, for each occurrence (starting position i) of s , we assign $c_i^l = \mathcal{R}[i-l] \dots \mathcal{R}[i-1]$ and $c_i^r = \mathcal{R}[i+1] \dots \mathcal{R}[i+l]$. If any of the strings c_i^l or c_i^r overlaps with a sensitive pattern, we truncate it by removing the minimum number of letters from it so that it does not. This is to allow our algorithms to replace a pattern occurring right after c_i^l and right before c_i^r without replacing letters of c_i^l or c_i^r . An alternative would be to extend c_i^l or c_i^r , so that it includes the sensitive pattern(s) that overlap with $\mathcal{R}[i-l] \dots \mathcal{R}[i-1]$ or $\mathcal{R}[i+1] \dots \mathcal{R}[i+l]$. We will henceforth assume the first way, as the second one did not lead to a practical benefit, based on our experiments. Furthermore, we truncate c_i^l if $i-l < 1$ or c_i^r if $i+l > |\mathcal{R}|$, to ensure that the context occurs in \mathcal{R} .

III. BEST K NONSENSITIVE PATTERNS SELECTION

We use the pattern matching algorithm of [11] to find $\text{occ}_{\mathcal{R}}(c_i^l)$ and $\text{occ}_{\mathcal{R}}(c_i^r)$ (i.e., the positions of c_i^l and of c_i^r in \mathcal{R}). Then, for each position $i \in \text{occ}_{\mathcal{R}}(c_i^l)$, we find the closest position $j > i + |c_i^l| - 1$ such that $j \in \text{occ}_{\mathcal{R}}(c_i^r)$ (i.e., the closest starting position of c_i^r). Recall that U_{c_i} at this point contains all sensitive patterns that are associated with context c_i . If the string u' right after c_i^l and right before c_i^r is not sensitive (i.e., $u' = \mathcal{R}[i + |c_i^l|] \dots \mathcal{R}[j-1] \notin S$), then we add the pair $(\sum_{u \in U_{c_i}} d(u', u), u)$ into an initially empty priority queue, which is sorted in descending order of the first element of the pairs in it. As we consider the positions in $\text{occ}_{\mathcal{R}}(c_i^l)$, the priority queue will contain K elements at some point¹. Then, when we consider the current u' , we

¹If this never happens, we cannot have enough nonsensitive pattern replacements for c_i , so we report failure.

P. Bian and G. Loukides are with Department of Informatics, King's College London, 30 Aldwych, London WC2B 4BG, UK.

E-mail: {pengxin.bian, grigorios.loukides}@kcl.ac.uk

G. Theodorakopoulos is with Cardiff University, UK. Email: theodorakopoulosg@cardiff.ac.uk

S. P. Pissis is with CWI and the Vrije Universiteit, The Netherlands. Email: solon.pissis@cwi.nl

Manuscript received April 19, 2005; revised August 26, 2015.

compute $\sum_{u \in U_{c_l}} d(u', u)$ incrementally (i.e., term by term) and if at some point the current sum exceeds the sum of the top element in the priority queue, we skip this u' , as it cannot be one of the K best, and continue with the next one. Otherwise, the fully computed sum and this u form a pair which replaces the current top pair in the priority queue. Thus, the priority queue will have K elements, after all positions in $\text{OCC}_{\mathcal{R}}(c_l^\ell)$ are considered. Last, we remove each element from the priority queue and add its string into the set U'_{c_l} .

IV. LINEARIZATION OF THE PROGRAM FOR AQL-MILP

To prove that the linear program in Program 2 is equivalent to that in Program 1, we show below that, for any u' , $z_{u'}$ indeed takes the minimum value of T_{MIN} and $M \cdot \sum_u P(u' | u)$.

$$\text{Minimize } \sum_u \pi(u | c_l) \sum_{u'} P(u' | u) d_q(u', u) \quad (1)$$

subject to

$$\begin{aligned} & \sum_u \pi(u | c_l) P(u' | u) d_p(\hat{u}, u) \\ & \geq \min\{T_{\text{MIN}}, M \cdot \sum_u P(u' | u)\}, \forall \hat{u}, u' \end{aligned} \quad (2)$$

$$\sum_{u'} P(u' | u) = 1, \forall u \quad (3)$$

$$P(u' | u) \geq 0, \forall u, u' \quad (4)$$

Program 1: Average quality-loss-minimizing, privacy-constrained nonlinear program

$$\text{Minimize } \sum_u \pi(u | c_l) \sum_{u'} P(u' | u) d_q(u', u) \quad (5)$$

subject to

$$z_{u'} \leq \sum_u \pi(u | c_l) P(u' | u) d_p(\hat{u}, u), \forall \hat{u}, u' \quad (6)$$

$$z_{u'} \leq T_{\text{MIN}}, \forall u' \quad (7)$$

$$z_{u'} \geq T_{\text{MIN}} - m \cdot y_{u',1}, \forall u' \quad (8)$$

$$z_{u'} \leq M \cdot \sum_u P(u' | u), \forall u' \quad (9)$$

$$z_{u'} \geq M \cdot \sum_u P(u' | u) - m \cdot y_{u',2}, \forall u' \quad (10)$$

$$y_{u',1} + y_{u',2} \leq 1, \forall u' \quad (11)$$

$$y_{u',1}, y_{u',2} \in \{0, 1\}, \forall u' \quad (12)$$

$$\sum_{u'} P(u' | u) = 1, \forall u \quad (13)$$

$$P(u' | u) \geq 0, \forall u, u' \quad (14)$$

Program 2: Average quality-loss-minimizing, privacy-constrained linear program

Case 1: $y_{u',1} = 0, y_{u',2} = 0$.

In Case 1, Eq. 8 becomes

$$z_{u'} \geq T_{\text{MIN}} \quad (15)$$

and Eq. 10 becomes

$$z_{u'} \geq M \cdot \sum_u P(u' | u). \quad (16)$$

From Eq. 7 and Eq. 15, we infer that

$$z_{u'} = T_{\text{MIN}} \quad (17)$$

From Eq. 9 and Eq. 16, we infer that

$$z_{u'} = M \cdot \sum_u P(u' | u) \quad (18)$$

From Eq. 17 and Eq. 18, we get

$$T_{\text{MIN}} = M \cdot \sum_u P(u' | u) \quad (19)$$

Therefore, in Case 1 $z(u')$ could be either T_{MIN} or $M \cdot \sum_u P(u' | u)$ for a specific u' as they are equal.

Case 2: $y_{u',1} = 0, y_{u',2} = 1$.

In Case 2, Eq. 8 becomes

$$z_{u'} \geq T_{\text{MIN}} \quad (20)$$

and Eq. 10 becomes

$$z_{u'} \geq M \cdot \sum_u P(u' | u) - m. \quad (21)$$

From Eq. 7 and Eq. 20, we infer that

$$z_{u'} = T_{\text{MIN}} \quad (22)$$

Eq. 9 and Eq. 21 will always hold, because $m - M \cdot \sum_u P(u' | u) > 0$ for every u' always holds.

Therefore, in Case 2 $z_{u'} = T_{\text{MIN}}$.

Case 3: $y_{u',1} = 1, y_{u',2} = 0$.

In Case 3, Eq. 8 becomes

$$z_{u'} \geq T_{\text{MIN}} - m \quad (23)$$

and Eq. 10 becomes

$$z_{u'} \geq M \cdot \sum_u P(u' | u). \quad (24)$$

Eq. 7 and Eq. 23 will always hold because $T_{\text{MIN}} - m < 0$ always holds.

From Eq. 9 and Eq. 24, we infer that

$$z_{u'} = M \cdot \sum_u P(u' | u) \quad (25)$$

Therefore, in Case 3 $z_{u'} = M \cdot \sum_u P(u' | u), \forall u'$.

These are the only cases, due to Eq. 11. Thus, the linearization of the minimum operator in AQL-MILP is correct; $z_{u'}$ will always be the minimum one among $T_{\text{MIN}}, M \cdot \sum_u P(u' | u)$ for every $u' \in U'_{c_l}$.

V. ADDITIONAL EXPERIMENTAL RESULTS

A. Impact of main parameters

Figs. 1, 2, 3, and 4 present the results of APG and AQL, illustrating how each parameter (Q_{MAX} for APG-LP, T_{MIN} for AQL-MILP, ϵ for DP-LP, and R for the Baseline) affects its respective algorithm on the *msnbc*, *trucks*, *kasandr*, and *iot* datasets. They correspond to the experiments in the Impact of Main Parameters part of Section IV-D in the main paper.

Figs. 5, 6, 7, and 8 show the APG and AQL of each method for varying K on the *msnbc*, *trucks*, *kasandr*, and *iot* datasets. They correspond to the experiments in the Impact of Main Parameters part of Section IV-D in the main paper.

Figs. 9, 10, 11, 12, and 13 show the results of APG and AQL of each method for varying l on the *iot*, *ecoli*, *msnbc*, *trucks*, and *kasandr* datasets. The results are too data-dependent to draw general conclusions (e.g., they depend on how the values of the d_q matrix and $\pi(u | c_l)$ values change). Take *iot* dataset as an example. Fig. 9 shows that increasing l results in a lower AQL for all algorithms, as the values in their d_q matrices decrease. As l increases, the APG-LP algorithm has a higher privacy (APG increases); see Fig. 9a. This is because the differences between the probabilities $\pi(u | c_l)$ for the different strings $u \in U_{c_l}$ become smaller, while the values in the d_p matrix do not change (observe Eqs. 8 and 9 in the main paper). On the other hand, AQL-MILP has the same APG (equal to $T_{\text{MIN}} \cdot |U_{c_l}|$); see Fig. 9c. For DP-LP and Baseline, APG decreases (see Figs. 9e and 9g) because the $\pi(u | c_l)$ values change as explained above and the values in d_p stay the same (see Eqs. 5 and 6). On the other hand, for the *ecoli* dataset, increasing l results in a generally higher AQL for all algorithms (see Figs. 10b, 10d, 10f, and 10h), as the values in their d_q matrices increase, while APG stays the same (see Figs. 10a, 10c, 10e, and 10g) because in this case $\pi(u | c_l)$ does not vary with l , nor does the d_p matrix.

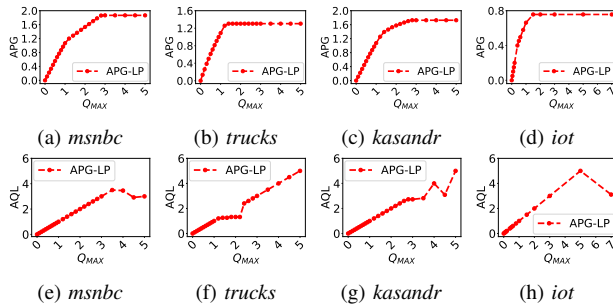


Fig. 1: (a)(b)(c)(d) APG and (e)(f)(g)(h) AQL for APG-LP vs. Q_{MAX} in the *msnbc*, *trucks*, *kasandr*, and *iot* datasets.

B. Hybrid Algorithms

Recall that the first hybrid algorithm APG-MILP maximizes APG while ensuring that AQL does not exceed Q_{MAX} and that each possible replacement u' gains at least T_{MIN}

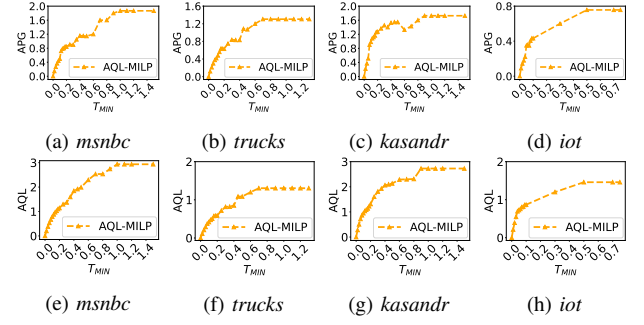


Fig. 2: (a)(b)(c)(d) APG and (e)(f)(g)(h) AQL for AQL-MILP vs. T_{MIN} in the *msnbc*, *trucks*, *kasandr*, and *iot* datasets.

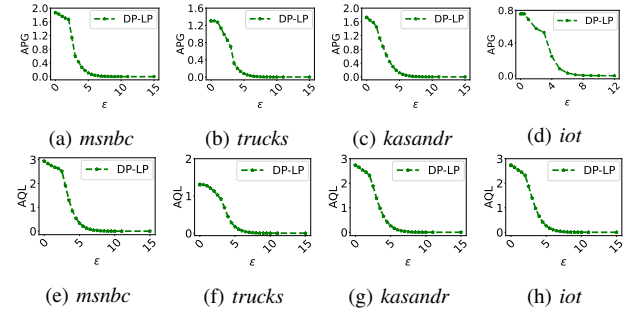


Fig. 3: (a)(b)(c)(d) APG and (e)(f)(g)(h) AQL for AQL-MILP vs. ϵ in the *msnbc*, *trucks*, *kasandr*, and *iot* datasets.

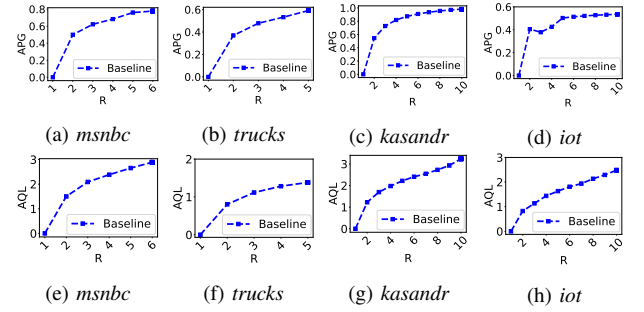


Fig. 4: (a)(b)(c)(d) APG and (e)(f)(g)(h) AQL for AQL-MILP vs. R in the *msnbc*, *trucks*, *kasandr*, and *iot* datasets.

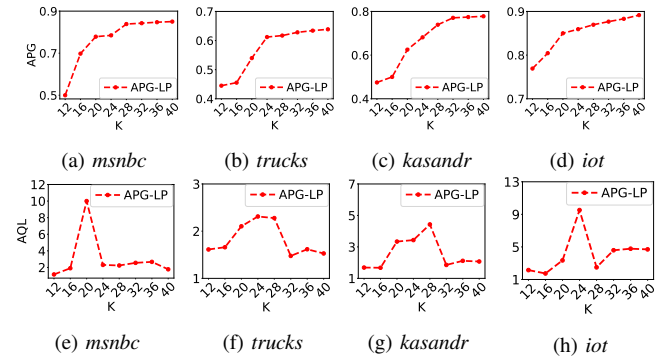


Fig. 5: APG and AQL for APG-LP vs. K in the *msnbc*, *trucks*, *kasandr*, and *iot* dataset.

privacy. Its mathematical programming formulation is shown in Program 3.

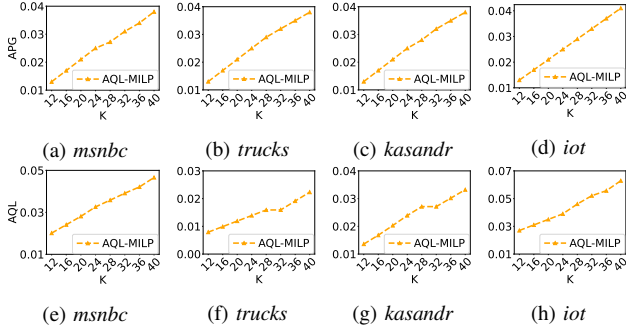


Fig. 6: APL and AQL for AQL-MILP vs. K in the *msnbc*, *trucks*, *kasandr*, and *iot* dataset.

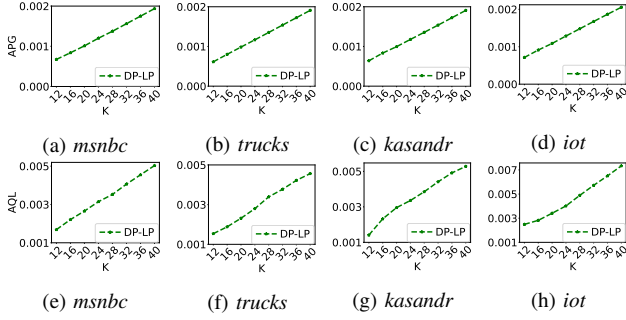


Fig. 7: APL and AQL for DP-LP vs. K in the *msnbc*, *trucks*, *kasandr*, and *iot* dataset.

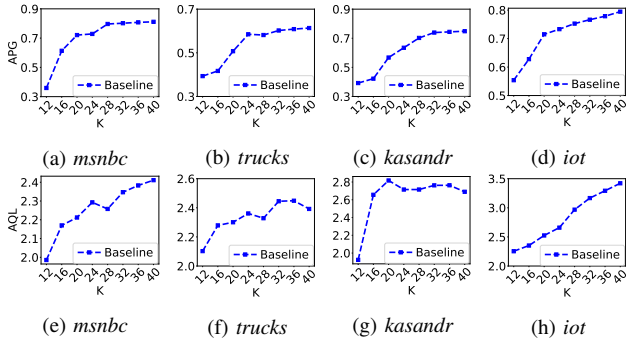


Fig. 8: APL and AQL for Baseline vs. K in the *msnbc*, *trucks*, *kasandr*, and *iot* dataset.

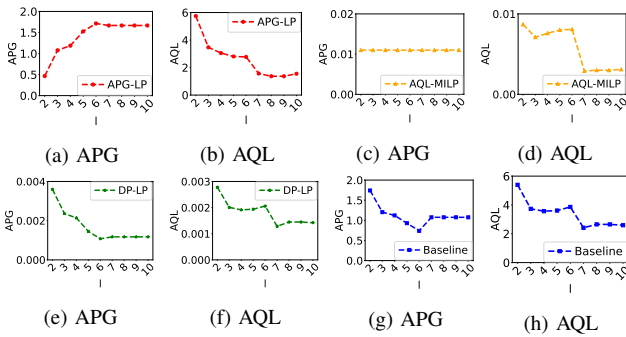


Fig. 9: (a)(c)(e)(g) APL and (b)(d)(f)(h) AQL vs. l for APG-LP, AQL-MILP, DP-LP, and Baseline. All results are for the *iot* dataset.

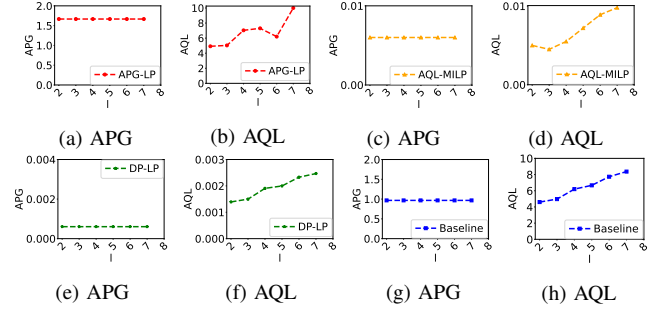


Fig. 10: (a)(c)(e)(g) APL and (b)(d)(f)(h) AQL vs. l for APG-LP, AQL-MILP, DP-LP, and Baseline. All results are for the *ecoli* dataset.

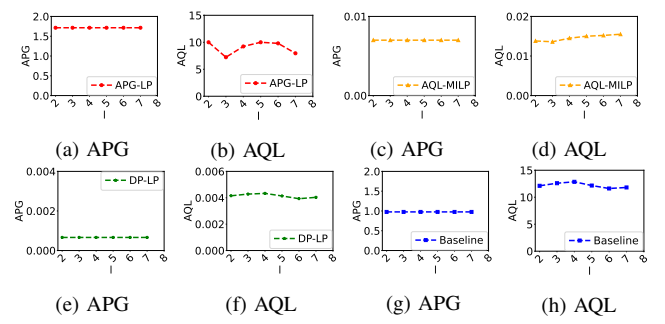


Fig. 11: (a)(c)(e)(g) APL and (b)(d)(f)(h) AQL vs. l for APG-LP, AQL-MILP, DP-LP, and Baseline. All results are for the *msnbc* dataset.

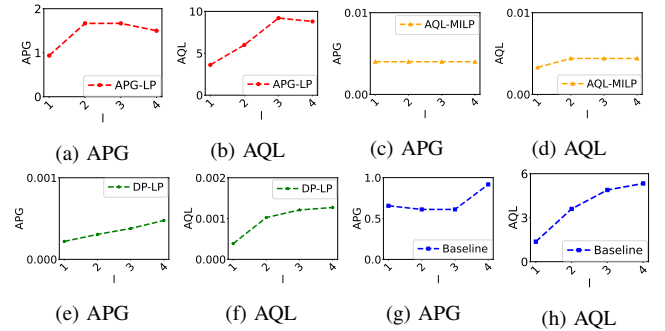


Fig. 12: (a)(c)(e)(g) APL and (b)(d)(f)(h) AQL vs. l for APG-LP, AQL-MILP, DP-LP, and Baseline. All results are for the *trucks* dataset.

and WAQL when varying Q_{MAX} and T_{MIN} respectively in *msnbc*, *trucks*, *ecoli*, and *iot* datasets, which correspond to the experiments in the Hybrid Algorithms part of Section IV-D in the main paper.

Recall that the second hybrid algorithm, DP-MILP, is derived by incorporating the differentially private constraint from DP-LP into AQL-MILP. The mathematical programming formulation for DP-MILP is shown in Program 4.

Similar to APG-MILP, we evaluate DP-MILP by varying T_{MIN} and ϵ separately. Fig. 16 and Fig. 17 show the results of WAPG and WAQL when varying T_{MIN} and ϵ respectively on the *msnbc*, *trucks*, *ecoli*, and *iot* datasets. These results supplement the experiments discussed in the Hybrid Algo-

Fig. 14 and Fig. 15 show the results of WAPG, WMPG

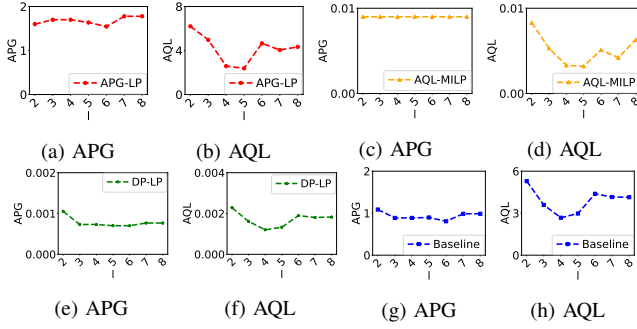


Fig. 13: (a)(c)(e)(g) APG and (b)(d)(f)(h) AQL vs. l for APG-LP, AQL-MILP, DP-LP, and Baseline. All results are for the *kasandr* dataset.

$$\text{Maximize } \sum_{u'} z_{u'} \quad (26)$$

subject to

$$\sum_u \pi(u | c_l) \sum_{u'} P(u' | u) d_q(u', u) \leq Q_{\text{MAX}} \quad (27)$$

$$z_{u'} \leq \sum_u \pi(u | c_l) P(u' | u) d_p(\hat{u}, u), \forall \hat{u}, u' \quad (28)$$

$$z_{u'} \leq T_{\text{MIN}}, \forall u' \quad (29)$$

$$z_{u'} \geq T_{\text{MIN}} - m \cdot y_{u',1}, \forall u' \quad (30)$$

$$z_{u'} \leq M \cdot \sum_u P(u' | u), \forall u' \quad (31)$$

$$z_{u'} \geq M \cdot \sum_u P(u' | u) - m \cdot y_{u',2}, \forall u' \quad (32)$$

$$y_{u',1} + y_{u',2} \leq 1, \forall u' \quad (33)$$

$$y_{u',1}, y_{u',2} \in \{0, 1\}, \forall u' \quad (34)$$

$$\sum_{u'} P(u' | u) = 1, \forall u \quad (35)$$

$$P(u' | u) \geq 0, \forall u, u' \quad (36)$$

Program 3: APG-MILP: Combination of APG-LP and AQL-MILP to maximize APG.

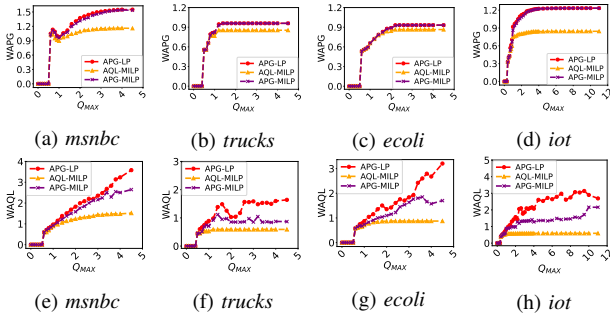


Fig. 14: WAPG and WAQL for APG-MILP vs. Q_{MAX} in the *msnbc*, *trucks*, *ecoli*, and *iot* datasets.

gorithms part of Section IV-D of the main paper.

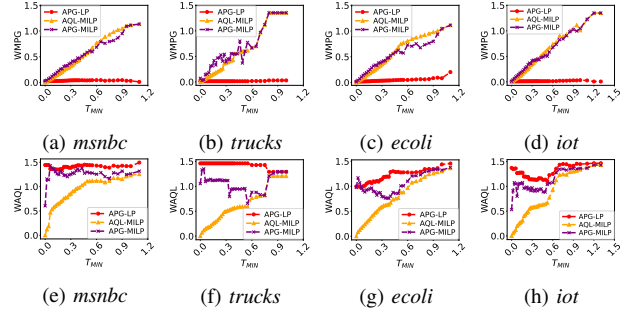


Fig. 15: WMPG and WAQL for APG-MILP vs. T_{MIN} in the *msnbc*, *trucks*, *ecoli*, and *iot* datasets.

$$\text{Minimize } \sum_u \pi(u | c_l) \sum_{u'} P(u' | u) d_q(u', u) \quad (37)$$

subject to

$$z_{u'} \leq \sum_u \pi(u | c_l) P(u' | u) d_p(\hat{u}, u), \forall \hat{u}, u' \quad (38)$$

$$z_{u'} \leq T_{\text{MIN}}, \forall u' \quad (39)$$

$$z_{u'} \geq T_{\text{MIN}} - m \cdot y_{u',1}, \forall u' \quad (40)$$

$$z_{u'} \leq M \cdot \sum_u P(u' | u), \forall u' \quad (41)$$

$$z_{u'} \geq M \cdot \sum_u P(u' | u) - m \cdot y_{u',2}, \forall u' \quad (42)$$

$$y_{u',1} + y_{u',2} \leq 1, \forall u' \quad (43)$$

$$P(u' | u_i) \leq P(u' | u_j) \cdot e^\epsilon, \forall u' \in U'_{c_l}, u_i, u_j \in U_{c_l} \quad (44)$$

$$y_{u',1}, y_{u',2} \in \{0, 1\}, \forall u' \quad (45)$$

$$\sum_{u'} P(u' | u) = 1, \forall u \quad (46)$$

$$P(u' | u) \geq 0, \forall u, u' \quad (47)$$

Program 4: DP-MILP: Combination of AQL-MILP and DP-LP to minimize the AQL

REFERENCES

- [1] B. C. M. Fung, K. Wang, R. Chen, and P. S. Yu, "Privacy-preserving data publishing: A survey of recent developments," *ACM Comput. Surv.*, vol. 42, no. 4, pp. 14:1–14:53, 2010.
- [2] P. Samarati and L. Sweeney, "Generalizing data to provide anonymity when disclosing information," in *Proc. ACM SIGMOD-SIGACT-SIGART Symp. Principles Database Syst. (PODS)*, vol. 98, no. 188, 1998, pp. 10–1145.
- [3] C. Dwork, F. McSherry, K. Nissim, and A. D. Smith, "Calibrating noise to sensitivity in private data analysis," in *Proc. Theory Cryptogr. Conf. (TCC)*, vol. 3876, 2006, pp. 265–284.
- [4] C. C. Aggarwal and P. S. Yu, "A condensation approach to privacy preserving data mining," in *Proc. Int. Conf. Extending Database Technol. (EDBT)*, vol. 2992, 2004, pp. 183–199.
- [5] —, "On anonymization of string data," in *Proc. SIAM Int. Conf. Data Mining (SDM)*, 2007, pp. 419–424.
- [6] —, "A framework for condensation-based anonymization of string data," *Data Min. Knowl. Discov.*, vol. 16, no. 3, pp. 251–275, 2008.
- [7] R. Chen, G. Ács, and C. Castelluccia, "Differentially private sequential data publication via variable-length n-grams," in *Proc. ACM Conf. Comput. Commun. Secur. (CCS)*, 2012, pp. 638–649.
- [8] L. Bonomi and L. Xiong, "A two-phase algorithm for mining sequential patterns with differential privacy," in *Proc. ACM Int. Conf. Inf. Knowl. Manag. (CIKM)*, 2013, pp. 269–278.

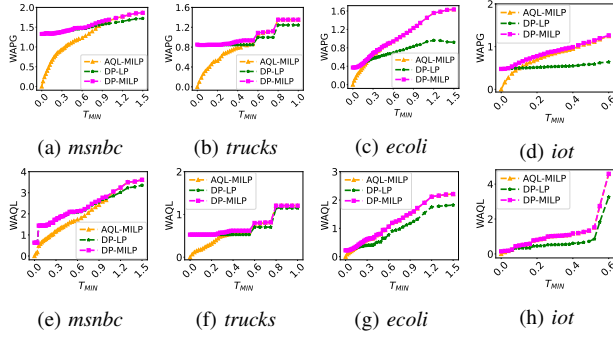


Fig. 16: WAPG and WAQL for DP-MILP vs. T_{\min} in the *msnbc*, *trucks*, *ecoli*, and *iot* datasets.

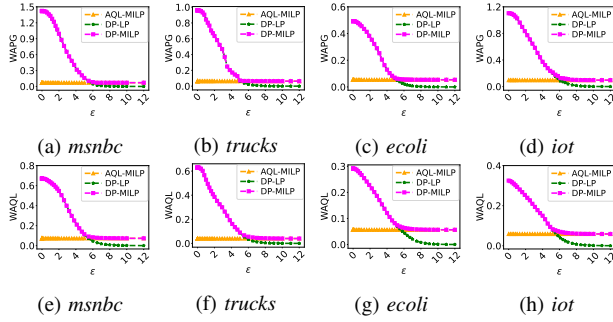


Fig. 17: WAPG and WAQL for DP-MILP vs. ϵ in the *msnbc*, *trucks*, *ecoli*, and *iot* datasets.

- [9] S. Xu, X. Cheng, S. Su, K. Xiao, and L. Xiong, “Differentially private frequent sequence mining,” *IEEE Trans. Knowl. Data Eng.*, vol. 28, no. 11, pp. 2910–2926, 2016.
- [10] H. Chen, C. Dong, L. Fan, G. Loukides, S. P. Pissis, and L. Stougie, “Differentially private string sanitization for frequency-based mining tasks,” in *Proc. IEEE Int. Conf. Data Mining (ICDM)*, 2021, pp. 41–50.
- [11] M. Crochemore, C. Hancart, and T. Lecroq, *Algorithms on strings*. Cambridge University Press, 2007.