

# Master ISL, projet tueuré, de recherche

## Répondre à des requêtes SQL en respectant la vie privée

J.-F. COUCHOT<sup>1</sup>, V. SONIGO<sup>1</sup>, and N. BENNANI<sup>2</sup>

<sup>1</sup>Institut FEMTO-ST, UMR 6174 CNRS, Université de Franche-Comté, Besançon, France.

`vsonigo,couchot@femto-st.fr`

<sup>2</sup>Laboratoire d'InfoRmatique en Image et Systèmes d'information, UMR 5205 CNRS, INSA Lyon, France., `nadia.bennani@insa-lyon.fr`

## 1 Contexte

L'Union Européenne souhaite responsabiliser les acteurs traitant des données en imposant que certaines décisions algorithmiques critiques ne soient prises par IA que si celle-ci est explicable, transparente et sûre [Com21]. Cette explicabilité dépend directement de la possibilité d'évaluer l'IA intensément et notamment parfois sur des données stockées dans des bases de données véritables. Cependant, les Systèmes de Gestion de Bases de Données (SGBD) peuvent manipuler des données personnelles et sensibles, et celles-ci ne doivent être révélées qu'aux personnes destinataires de ladite tâche (RGPD). A contrario, ces données ne doivent pas être mises à la connaissance des personnes qui développent des algorithmes (de traitement, d'évaluation d'IA) sur ces bases, qui les corrigent aussi par la suite. Chiffrer ces données est nécessaire, pseudonymiser les identifiants est indispensable, mais cet ensemble est nullement suffisant. Savoir anonymiser les données pour interdire toute ré-identification, toute acquisition de connaissances personnelles est un enjeu crucial à l'ère du BIG DATA. Différentes classes d'approches existent pour atteindre l'objectif d'anonymisation des résultats des requêtes : des méthodes syntaxiques (k-anonymat et dérivées), par ajout de bruit, ... La confidentialité différentielle (Differential Privacy en anglais, DP) qui fait partie de cette seconde catégorie est acceptée comme un standard puisqu'elle permet de mesurer en termes probabilistes la quantité d'information fuitée à chaque réponse de requête. Ce confort a deux prix : d'une part, l'expression probabiliste est moins compréhensible par la personne en charge qu'un k-anonymat. D'autre part, le bruit introduit se fait souvent au détriment de l'utilité si l'on n'y prend pas garde.

Des extensions d'anonymisation aux SGBD usuels existent aujourd'hui. Dynamic DataMasking<sup>1</sup> (MS SQL Server), Masking with maxscale<sup>2</sup> (MariaDB), PostgreSQL Anonymizer<sup>3</sup> proposent des méthodes à base de pseudonymisation, de k-anonymat et d'ajout de bruit pour nettoyer plusieurs tables d'un SGBD. [KTH<sup>+</sup>19, WZL<sup>+</sup>19] proposent quant-à eux des méthodes permettant de nettoyer par confidentialité différentielle globale les réponses aux requêtes de type Select-Project-Join-Aggregation (SPJA) en autorisant plus ou moins de contributions par entité sensible (utilisateur, par exemple). [WZL<sup>+</sup>19] est complètement outillée<sup>4</sup> dans le langage ZETA-SQL de Google (nommée DP-ZETA-SQL ensuite).

## 2 Objectifs

Les objectifs de ce projet de recherche sont multiples.

1. Bibliographique d'abord. Il s'agira de comprendre les différentes méthodes d'anonymisation de données.
2. Pratique ensuite. Il s'agira ensuite d'évaluer a minima PostgreSQL Anonymizer pour l'approche syntaxique essentiellement et DP-ZETA-SQL sur des benchmarks reconnus.
3. Théorique enfin. Il s'agira d'évaluer comment des mécanismes récents à base de DP locale [DJW13], de d-privacy [CABP13], de Renyi-DP [Mir17], pourraient être exploitées particulièrement lorsqu'il s'agit de publier des releases. Le nettoyage par DP locale de publications de données catégorielles pour construire des histogrammes (COUNT/GROUP BY) a récemment beaucoup progressé [ACABX21,

---

1. <https://docs.microsoft.com/en-us/sql/relational-databases/security/dynamic-data-masking>

2. <https://mariadb.com/kb/en/mariadb-enterprise/mariadb-maxscale-21-masking/>

3. <https://postgresql-anonymizer.readthedocs.io/en/stable/>

4. <https://github.com/google/differential-privacy/tree/main/examples/zetasql>

ACABX22] et l'étude de ces progrès pour leur intégration fait partie des objectifs théoriques de ce travail.

### 3 Suites possibles

Ce projet pourra être suivi par un stage de recherche en laboratoire, financé et d'un doctorat dans le cadre d'un projet de recherche, pour lequel une bourse est déjà acquise.

### Références

- [ACABX21] Héber H Arcolezi, Jean-François Couchot, Bechara Al Bouna, and Xiaokui Xiao. Random sampling plus fake data : Multidimensional frequency estimates with local differential privacy. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 47–57, 2021.
- [ACABX22] Héber H Arcolezi, Jean-François Couchot, Bechara Al Bouna, and Xiaokui Xiao. Improving the utility of locally differentially private protocols for longitudinal and multidimensional frequency estimates. *Digital Communications and Networks*, 2022.
- [CABP13] Konstantinos Chatzikokolakis, Miguel E Andrés, Nicolás Emilio Bordenabe, and Catuscia Palamidessi. Broadening the scope of differential privacy using metrics. In *International Symposium on Privacy Enhancing Technologies Symposium*, pages 82–102. Springer, 2013.
- [Com21] European Commission. Laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts, 2021.
- [DJW13] John C Duchi, Michael I Jordan, and Martin J Wainwright. Local privacy and statistical minimax rates. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 429–438. IEEE, 2013.
- [KTH<sup>+</sup>19] Ios Kotsogiannis, Yuchao Tao, Xi He, Maryam Fanaeepour, Ashwin Machanavajjhala, Michael Hay, and Gerome Miklau. Privatesql : a differentially private sql query engine. *Proceedings of the VLDB Endowment*, 12(11) :1371–1384, 2019.
- [Mir17] Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th computer security foundations symposium (CSF)*, pages 263–275. IEEE, 2017.
- [WZL<sup>+</sup>19] Royce J Wilson, Celia Yuxin Zhang, William Lam, Damien Desfontaines, Daniel Simmons-Marengo, and Bryant Gipson. Differentially private sql with bounded user contribution. *arXiv preprint arXiv :1909.01917*, 2019.