

DeCor: Decorrelating unsubscribed users' data from their identities

Anonymous Authors

Abstract

100-200 words that convince you to read this.

1 Introduction

1.1 Motivation

Web application companies face increasing legal requirements to protect users' data. These requirements pressure companies to properly delete and anonymize users' data when a user requests to *unsubscribe* from the service (i.e., revoke access to their personal data). For example, the GDPR requires that any user data remaining after a user unsubscribes is *decorrelated*, i.e., cannot be (directly or indirectly) used to identify the user [2].

Furthermore, keeping identifying and personal data when no longer strictly necessary increases companies' liability: the GDPR and other laws mandate that companies retain only user data that is relevant and necessary for their applications' purposes. To increase users' control over their data and decrease the amount of incriminating data stored in the application at any one point, users should be able to freely unsubscribe from the service to enter a privacy-preserving mode, and later resubscribe when they wish to use the service. [lyt: (this seems unconvincing, should come up with better motivation?)]

1.2 Goals

DeCor's goal is to provide the following properties while preserving an application's semantics:

- **Decorrelation:** Decorrelation ideally guarantees that it is impossible to distinguish between two records formerly associated with the same unsubscribed user and two records from different unsubscribed users.
- **Resubscription:** Users should be able to easily switch between a privacy-preserving unsubscribed mode and

an identity-revealing subscribed mode, without permanently losing their application data.

DeCor must implement these properties while ensuring (1) performance comparable to today's widely-used databases, and (2) easy adoption (decorrelation should be automated and not require modifications to application schemas or semantics).

1.3 Decorrelation Guarantees

We address applications in which application data consists of *data records* and computations (such as aggregations) that may be performed over these data records. Data records are considered sensitive, private data records when they contain *user identifiers*; for example, a row in a table containing a column of user IDs would be a user's private data record. Data records containing multiple, potentially different, user IDs are considered shared data records private to the identified users.

Perfect decorrelation is achieved when queries to the application reveal no information allowing an observer to determine if two user data records with different user IDs belong to the same user who has since unsubscribed (are *linked*), or belong to two different users. Observers gain no information that allows them to distinguish the two scenarios [lyt: Not sure how to define "distinguish" formally—equally probable? Perhaps noninterference can be cited here)].

In practical settings, however, perfect decorrelation is likely impossible: for example, a reposted screenshot may leak user IDs. Furthermore, an observer who can see application queries over time, or search web archives, can detect when a user unsubscribes and refer to prior snapshots in which user data records may have had the same user ID. Given these limitations, we seek to achieve the maximum decorrelation guarantees possible while assuming that the content of user data does not itself leak identifying information, and that observers cannot access past application database state. [lyt: Note: UIDs can also include things like email addrs, phone number, etc; so the notion of "identifiers" might be more broad than stated here].

Global Placeholder. A common strawman solution to decorrelation is to replace all unsubscribed user IDs with one global placeholder ID.

This results in the following privacy guarantees [lyt: how to formalize? probability stuff...]

- Extremely Unprivate: If only one user ever unsubscribes, all remnants can be linked back to that user’s identity.
- Extremely private: If all users unsubscribe, two remnants are equally as likely to belong to two individual users as they are to belong to one single user (assuming the number of users is unknown and follows some distribution... [lyt: I need to think about this a bit more, but I think the intuition is clear]).

The amount of linkable information [lyt: (need to define this)], and the resulting deviance from perfect decorrelation, decreases as more users unsubscribe. [lyt: not sure by how much though].

Furthermore, a global placeholder makes resubscription challenging: users can no identify which unsubscribed data records belong to them if all user-specific data has been erased from the system.

Data-Record Ghost IDs. To support stronger decorrelation guarantees while support resubscription, DeCor generates a unique ghost user for each data remnant. These ghosts, unlike a global placeholder, are indistinguishable from a real user in the system, ensuring that queries cannot correlate two ghosts with a single real (albeit unsubscribed) user. DeCor relies on coarse-grained schema annotations to establish which associations to decorrelate, and builds a dataflow computation resulting in materialized views that answer application queries. Use of dataflow automatically propagates the correct updates to materialized views.

However, some linkable information is still leaked. For example, ghost users may be randomly generated in a pattern identifiable by an observer (e.g., if all ghosts have usernames which are random numbers, or arbitrary animals, but a real user may have more human-friendly usernames).

[lyt: How to tie in differential privacy? Where there is some privacy budget that is lost on each query...? Maybe an observer needs to know the distribution of usernames in order to determine that there is a pattern that is machine-generated?]

[lyt: Or maybe there is some notion of noninterference we could use here? How do DeCor’s responses to queries deviate from what an ideally decorrelated system would produce?]

Unlike a global placeholder, DeCor allows users to reactivate their account and undo the decorrelation: user IDs can be linked back to a set of unique ghost IDs. This gives users the ability to freely unsubscribe to protect their privacy without worrying about losing their accounts. DeCor resubscribes users by transparently propagating updates to materialized views to expose real user identifiers in place of ghost identifiers.

2 Background and Related Work

Companies have developed frameworks to avoid deletion bugs (e.g., DELF [1] at Facebook), but many applications decorrelate only coarsely by associating remnants of data with a global placeholder for all deleted users, or do not decorrelate at all (e.g., replacing usernames with a pseudonym).

[lyt: (cite Reddit, Lobsters, others?)]

3 Design

4 Implementation

5 Evaluation

6 Discussion

Acknowledgments

Availability

[lyt: USENIX program committees give extra points to submissions that are backed by artifacts that are publicly available. If you made your code or data available, it’s worth mentioning this fact in a dedicated section.]

References

- [1] Katriel Cohn-Gordon, Georgios Damaskinos, Divino Neto, Joshi Cordova, Benoît Reitz, Benjamin Strahs, Daniel Obenshain, Paul Pearce, and Ioannis Papagiannis. DELF: Safeguarding deletion correctness in online social networks. In *29th USENIX Security Symposium (USENIX Security 20)*. USENIX Association, August 2020.
- [2] E Parliament. Regulation (eu) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing directive 95/46/ec (general data protection regulation). *Official Journal of the European Union*, L119(May 2016):1–88, 2016. <http://eur-lex.europa.eu/legal-content/EN/TXT/?uri=OJ:L:2016:119:TOC>.