



# Meaningful Data Erasure in the Presence of Dependencies

Vishal Chakraborty  
University of California (UC), Irvine  
vchakrab@uci.edu

Youri Kaminsky  
Hasso Plattner Institute,  
University of Potsdam  
youry.kaminsky@hpi.de

Sharad Mehrotra  
UC, Irvine  
sharad@ics.uci.edu

Felix Naumann  
Hasso Plattner Institute,  
University of Potsdam  
felixnaumann@hpi.de

Faisal Nawab  
UC, Irvine  
nawabf@uci.edu

Primal Pappachan  
Portland State University  
primal@pdx.edu

Mohammad Sadoghi  
UC, Davis  
msadoghi@ucdavis.edu

Nalini Venkatasubramanian  
UC, Irvine  
nalini@uci.edu

## ABSTRACT

Data regulations like GDPR require systems to support data erasure but leave the definition of “erasure” open to interpretation. This ambiguity makes compliance challenging, especially in databases where data dependencies can lead to erased data being inferred from remaining data. We formally define a precise notion of data erasure that ensures any inference about deleted data, through dependencies, remains bounded to what could have been inferred before its insertion. We design erasure mechanisms that enforce this guarantee at minimal cost. Additionally, we explore strategies to balance cost and throughput, batch multiple erasures, and proactively compute data retention times when possible. We demonstrate the practicality and scalability of our algorithms using both real and synthetic datasets.

### PVLDB Reference Format:

Vishal Chakraborty, Youri Kaminsky, Sharad Mehrotra, Felix Naumann, Faisal Nawab, Primal Pappachan, Mohammad Sadoghi, and Nalini Venkatasubramanian. Meaningful Data Erasure in the Presence of Dependencies. PVLDB, 18(10): 3435 - 3448, 2025.

doi:10.14778/3748191.3748206

### PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/HPI-Information-Systems/P2E2-Erasure>.

## 1 INTRODUCTION

Recently enacted data regulations [12, 13, 20, 22, 57] have codified the right to erasure and established principles of data minimization and integrity. This has driven academic [4, 47, 49, 50] and industry [18, 58] efforts to address the challenges of compliant data erasure. In semantically rich databases, simply deleting the user-specified data is insufficient, as it may be inferred from remaining

data. This problem arises in domains like social media, business, etc.—wherever applications store semantically linked user data.

Databases already support deletion of data beyond that specified by the user, particularly when required to preserve consistency [29, 41]. For example, deleting a record in a parent table can trigger cascading deletions [27] in child tables via foreign key constraints. Users can also define application-specific deletion logic through triggers [41]. Cascading deletions—motivated by regulatory compliance—have been studied in [46, 50], including “shallow” and “deep” deletions in graph databases [18] and annotation-driven deletions in [2]. However, these solutions are often bespoke to specific settings, and ad hoc—lacking formal guarantees. Critically, they adopt an operational view of deletion without grounding it in a principled notion of the user’s right to erasure, especially when semantic dependencies can enable inferences about deleted data.

Our goals in this paper are threefold: (a) to define a principled notion of deletion for databases that prevents deleted data from being inferred (by exploiting semantically dependent data) after deletion, (b) to develop effective and efficient ways to implement the developed deletion notion that minimizes additional data to be deleted, and (c) to understand (through a detailed experimental study across several domains and data sets) the implication and practical viability of such an approach.

**Pre-insertion Post Erasure Equivalence.** We formally define *Pre-insertion Post-Erasure Equivalence* (P2E2), a deletion principle that restricts inferences about deleted data using dependencies in the database to only those that were possible at the time of its insertion. P2E2 is defined at the cell level, with each cell assigned an expiration time by which it must be deleted. While some data regulations are vague on inference of deleted data using dependencies [15], some regulations [20] require to prevent “reconstruction (of deleted data) in an intelligible form.” Our formalization of data deletion in presence of data dependencies aligns with regulatory expectations for safe and effective deletion, selective retention, automated erasure and principles of data minimization and purpose limitation as outlined in CCPA [12] §1798.105(d), PIPEDA [13] PRINCIPLE 4.5.5, GDPR [22] ART. 5(1) & 17, LGPD [11] ART. 15, HIPAA [56] 45 CFR §164.310(d)(2)(i), and PDPA [52] Sec. 25.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 18, No. 10 ISSN 2150-8097.

doi:10.14778/3748191.3748206

*Example 1.1.* Consider a relation  $R(\text{Name}, \text{City}, \text{AreaCode})$  with the constraint: *if City = Irvine, then AreaCode  $\in \{714, 949\}$*  (implemented in PostgreSQL as a check constraint<sup>1</sup>). Starting with the tuple (Alice, NULL, NULL), first insert AreaCode = 949, resulting in the tuple (Alice, NULL, 949), followed by City = Irvine results in (Alice, Irvine, 949). Deleting AreaCode yields (Alice, Irvine, NULL), from which one can now infer that AreaCode was in {714, 949}—information that was *not inferable* at the time of AreaCode’s insertion. To satisfy P2E2, we must also delete City, ensuring no new inferences arise post-deletion. In contrast, if City = Irvine is inserted *before* AreaCode = 949, the same deletion of AreaCode does not require deleting City, since the constraint was already active and the inference about AreaCode was possible prior to its insertion. This highlights the subtlety of P2E2: deletion must prevent *new* inferences based on post-insertion context, which depends on the sequence of actions.

**Relational Dependency Rules.** To formally define P2E2, we introduce *relational dependency rules* (RDRs) that provide a simple, yet general, framework to express a variety of dependencies in data. RDRs use a SQL-like language and can express traditional dependencies including classes of both hard and soft constraints<sup>2</sup>, such as functional & inclusion dependencies, denial constraints [6, 38] and correlation constraints, respectively. RDRs, in addition, allow constraints to be limited to only a subset of data that satisfies the SQL query enabling specifications of conditional constraints such as (conditional) functional dependencies [8, 9, 25, 33], and similarity inclusion dependencies [34]. Frameworks, similar to RDRs, to express semantic constraints in database have previously been proposed in a variety of data processing contexts ranging from database design and cleaning [19, 44], consistent query answering [3], to database repair [6]. We adopt RDRs as they are based on SQL and hence expressive, and intuitive. Additionally, RDRs can allow specification of aggregation constraints (often found in organizational databases) which existing framework for specifying constraints such as [29] do not consider. Furthermore, RDRs can express regulatory-motivated semantic annotations [2, 18], dependencies discovered from data [1], and fine-grained erasure logic—e.g., Meta’s rule to delete user comments but retain messages [40], or selective retention under GDPR Art. 17(3) [22].

**Minimizing overhead.** Given a set of RDRs<sup>3</sup>, and a data item to be deleted, additional data may need to be deleted to ensure P2E2. Such additional deletions, as illustrated earlier, depend on the database state, both at the time of deletion, as well as the time of insertion of the data item being deleted. In Example 1.1, the database state at the deletion time of AreaCode was (Alice, Irvine, 465). However, depending upon the state at the time of *insertion* (viz. (Alice, NULL, NULL) versus (Alice, Irvine, NULL)) the set of additional deletions required to implement P2E2 differed. Realizing P2E2 requires tracking what changes were made to the database between the time the data was inserted to the time it is to be deleted, and to develop a logic to reason with RDRs to determine if such changes can cause leakage beyond what was already possible prior

to insertion. In addition, to minimize overhead, we would further like the set of additional deletions required (to implement P2E2) to be minimal. We note that such a logic to realize P2E2 cannot be easily encoded as simple deletion rules in the form of triggers<sup>4</sup>. We thus develop algorithms to realize P2E2 with minimal deletions and realize them in a middleware layer on top of the database.

Our problem of computing a minimal set of deletions to satisfy P2E2, formalized as the OPT-P2E2 problem, is related to the well-studied problem of minimal repair [3, 17]. The minimal repair problem takes as input a database  $D$  that is inconsistent with respect to a set of constraints  $C$  (e.g., denial constraints [17, 38, 53]) and computes a minimal set of modifications that restore consistency.

In contrast, P2E2 starts with a database consistent with respect to its constraints and a data item  $c$  to be deleted. The goal is to find a minimal subset of dependent data whose deletion ensures that no new inferences about  $c$  can be made beyond what was inferable at the time of its insertion. While both problems aim to minimize deletions, their objectives—and hence their solutions—differ. Repair algorithms cannot be directly applied to enforce P2E2. Therefore, we leverage RDRs to model data dependencies and design erasure mechanisms tailored to ensure P2E2.

**Erasure mechanisms.** We develop mechanisms to enforce P2E2 for: (a) *demand-driven erasure*, where users request deletion at any time, and (b) *retention-driven erasure*, where data is deleted after a predefined retention period [4], as in WhatsApp’s disappearing messages feature [59]. In both cases, the goal is to identify a minimal set of dependent data to delete to satisfy P2E2. We propose multiple strategies to compute this minimal set, including ILP, bottom-up tree traversal, top-down traversal (trading precision for throughput), and batched erasure. These methods offer different trade-offs in terms of overhead (time and space) and deletion scope. To further reduce overhead, we tailor optimizations to each erasure type. For demand-driven erasure, we introduce a grace period to batch deletions. For retention-driven erasure, we leverage advance knowledge of deletion times to schedule erasures cost-effectively—especially when derived data (e.g., aggregates or materialized views) must be reconstructed after base data is deleted.

**Evaluation.** We evaluate our approaches on real and synthetic datasets under various workloads, analyzing the cost and performance impact of P2E2. We compare exact and approximate algorithms, studying trade-offs between computational overhead and cost. Additionally, we examine the effects of batching, varying grace periods, and pre-computing retention periods for derived data. Our evaluation on five data sets shows that, on average, the number of extra deletions to guarantee P2E2 for a given cell is low and depends on factors such as the number of cells in a tuple, dependencies, and the number of insertions and deletions.

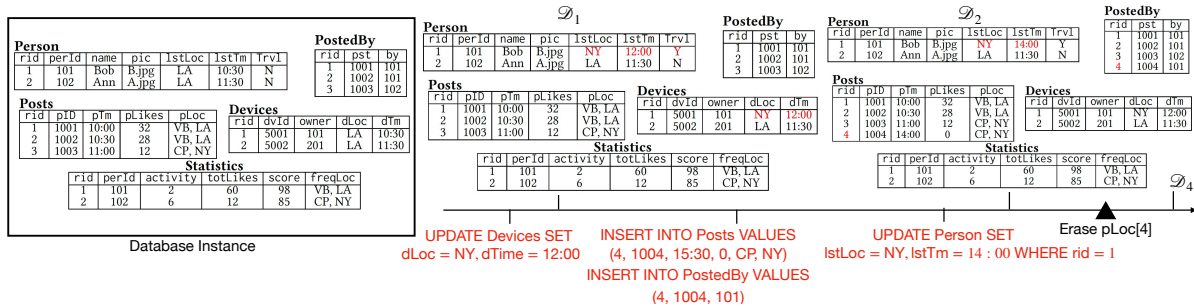
Section (Sec.) 2 formalizes data erasure and introduces its semantic guarantee, P2E2. Sec. 3 presents Relational Dependency Rules, and technical results for reasoning about cell-level erasure. Sec. 4 develops algorithms for demand-driven erasure, while Sec. 5 addresses retention-driven deletion. We evaluate our approaches and analyze P2E2’s overheads in Sec. 6. Sec. 7 reviews related work, and Sec. 8 concludes with directions for future research.

<sup>1</sup><https://www.postgresql.org/docs/current/ddl-constraints.html>

<sup>2</sup>Hard constraints are those that always hold while, soft constraints are likely to hold [17] on the database instance, though are not guaranteed to hold.

<sup>3</sup>We assume a list of RDRs capturing semantic constraints in the data as input.

<sup>4</sup>Trigger conditions only check the database state before and after the triggering event and, as such, are independent of the sequence of events that led to the state.



(a) Database Instance in a social media platform and timeline of operations for running example.

Dependence:  $\text{totLikes}(R) \not\sqsubseteq \text{pstLikes}(M)$   
Condition:  $\text{SELECT } S.\text{rid}, P.\text{rid AS } R, M$   
FROM Statistics S, Posts P, Person I, PostedBy B  
WHERE S.perID = I.perID AND B.pst = P.pID  
AND B.by = I.perID

(b) RDR R1

Dependence:  $\text{lastLoc}(U) \not\sqsubseteq \text{pLoc}(M)$   
Condition:  $\text{SELECT } I.\text{rid}, P.\text{rid AS } U, M$   
FROM Person I, Devices D, PostedBy B, Posts P (  
SELECT B.by AS by, MAX(P.pTm) AS max\_pTm FROM  
Postedby B, Posts P WHERE B.pst = P.pID  
GROUP BY B.by) AS last\_posts WHERE  
I.perID = D.owner AND I.perID = B.by AND  
B.by = last\_posts.by AND  
P.pTm = last\_posts.max\_pTm AND P.pTm > D.dTm

(c) RDR R2

Figure 1: Database instance, database states, and data dependencies (RDRs) for running example.

## 2 PRELIMINARIES

In this section, we introduce the formal semantics of data erasure and adopt standard notation and concepts from [28, 39]. But first we introduce a running example, which we use throughout the paper to illustrate notations and explain the intuition behind P2E2. **Running Example.** Consider a social media platform where users post content, upload photos, and use location-based services. The platform maintains the following database tables (Fig. 1a):

- Person: Stores user details, including name, profile picture, and travel status (Trvl), and last known location (lstLoc).
- Posts: Records user posts with timestamps, locations, and like counts (pLikes) and PostedBy: Tracks the author of each post.
- Device: Stores the last known location (dLoc) of a user’s device.

Attribute lstLoc in Person updates when the user’s device location (dLoc) changes or makes a new post with location (pLoc). - The platform also maintains a Statistics table for analytics, tracking average activity (activity), most frequented location

| Symbol                           | Meaning   |
|----------------------------------|---|
| $S = (\mathcal{R}, \mathcal{A})$ | Schema with sets $\mathcal{R}$ of relations & $\mathcal{A}$ of attributes |
| $R_i, A_j$                       | A relation $R_i \in \mathcal{R}$ and attribute $A_j \in \mathcal{A}$      |
| $rid_k$                          | Record-ID uniquely identifying records                                    |
| $\mathcal{D}_t$                  | Database state, i.e., set of cells in database at time $t$                |
| $\text{Cells}(\cdot)$            | Operator returning all cells in the given argument                        |
| $A_j(rid_k)$                     | Cell containing the value of attribute $A_j$ of $rid_k$                   |
| $\kappa(c)$                      | Creation (insertion) timestamp of cell $c$                                |
| $\eta(c)$                        | Expiration timestamp of cell $c$ (erasure time)                           |
| $Val(A(x), \mathcal{D}_t)$       | Value of attribute function $A(x)$ in $\mathcal{D}_t$                     |
| $dep(A(x)   \mathcal{D}_t)$      | Set of dependencies on $A(x)$ in $\mathcal{D}_t$                          |
| $\delta^-$                       | Instantiated relational dependency rule (RDR)                             |
| $\Delta^-(\mathcal{D}_t)$        | Set of all instantiated RDRs in $\mathcal{D}_t$                           |
| $Head(\delta^-)$                 | Head of instantiated RDR $\delta^-$                                       |
| $Tail(\delta^-)$                 | Tail of instantiated RDR $\delta^-$                                       |

Table 1: Notation Table

(freqLoc), and total likes (totLikes). The system enforces: (1) Updates to pLikes propagate to totLikes for the user; (2) Trvl in Person is set to Y if freqLoc  $\neq$  lstLoc; (3) Other Statistics attributes (e.g., activity) update periodically (e.g., weekly).

In the following, we present some notation and discuss our data and retention model required to formalize the problem setting.

**Functional Data Model.** We extend the functional data model [51] and its recent adaptation [39]. A database  $\mathcal{D}$  consists of relations  $\mathcal{R} = \{R_1, \dots, R_m\}$ , each with attributes  $R_i^{attr} = \{A_1, \dots, A_{n_j}\}$ . Attributes are referred to as  $A_j$  when the relation is clear from context. In Fig. 1a, Person is a relation, attributes perID, name, profPic, and lastLoc. Each relation  $R_i$  contains records  $r_k^i$ , uniquely identified by  $rid_k$ , its record-ID (rid). Records consist of cells, each corresponding to an attribute. We use the operator  $\text{Cells}(\cdot)$  which returns all the cells in the argument. To an attribute  $A_j \in R_i^{attr}$ , we associate the attribute function  $A_j : R_i(rid) \rightarrow \text{Cells}(\mathcal{D})$  that takes in as input a rid of a relation and returns the cell corresponding to attribute  $A_j$ . A cell is denoted as  $A_j(R_i(rid_k))$ , or  $A_j(rid_k)$  when the relation is clear. In the database instance in Fig. 1a, name(1) refers to the cell containing “Bob”. Each cell has an associated deletion cost given by the function  $Cost : \mathcal{D} \rightarrow \mathbb{R}$ . A relational functional schema is  $S = (\mathcal{R}, \mathcal{A})$  where  $\mathcal{A} = \bigcup_{i=1}^m R_i^{attr}$ .

**Database State.** At any time  $t$ , we write  $\mathcal{D}_t$  to denote the *database state*, i.e., the set of all records in the database at time  $t$  and their cell values. The value of a cell  $c$  in  $\mathcal{D}_t$  is given by  $Val(c, \mathcal{D}_t) \in \text{Dom}(A_j)$  where  $\text{Dom}(A_j)$ , the domain of  $A_j$ , is the set of all possible values  $A_j$  can take including NULL. As an example, in Fig. 1a,  $Val(\text{name}(1), \mathcal{D}_1) = \text{Bob}$ . With  $\mathcal{D}_t^-$  and  $\mathcal{D}_t^+$  we denote the states immediately before and after  $\mathcal{D}_t$ , respectively. A cell’s value becomes NULL (empty) when it is erased, and erasure of a record entails erasure of all cells within it which are not empty.

**Types of Data.** Relations are classified as *base* or *derived*. Base relations store personal data on which the user has direct control

(request deletion, rectification, etc.). Derived relations contain data which result from processing base and/or other derived relations. In the example, Posts is a base relation, while Statistics is derived.

*Base Relations.* A cell  $c$  in base relations has: (1) creation timestamp  $\kappa(c)$  – the time at which it is inserted, (2) erasure/expiration timestamp  $\eta(c)$  – the time at which it is deleted. The retention period of a cell is the interval between  $\kappa(c)$  and  $\eta(c)$ . When a record is inserted, for each cell  $c$  in the record,  $\kappa(c)$  is initialized with the time at which the record was inserted. The expiration time of each cell is set to a fixed or user-specified time at which it needs to be deleted from the database. For a base cell  $c$ , users can adjust  $\eta(c)$  to enable demand-driven erasure, which overrides the default retention period. Updating a cell is treated as an erasure followed by insertion, updating  $\kappa(c)$  to the time of the update but retaining the original  $\eta(c)$ . Base data erasure satisfies P2E2.

*Derived Relations.* Derived data are computed from base or other derived relations, with periodic recomputation. For a derived cell  $c$ , we denote with  $freq(c)$  the time period in which  $c$  has to be reconstructed at least once. For example,  $freq(c)$  for score is 30 days, and for totLikes, 7 days. When a derived cell  $c$  is recomputed, its creation timestamp  $\kappa(c)$  is updated to the recomputation time, while  $freq(c)$  remains unchanged. Derived data lack explicit erasure timestamps, as they may be reconstructed after base data erasure to prevent inferences. Derived data are deleted only when no longer required, ceasing further reconstruction. Particularly, users cannot directly ask derived data to be deleted.

### 3 P2E2

In a database, the value of a cell often depends on that of other cells. In this section, we will formally define RDRs as a means of expressing data dependencies. We discuss how RDRs can be instantiated, and how we can reason about inferences using RDRs.

#### 3.1 Specifying Background Knowledge by RDRs

RDRs express dependency (background knowledge) among cells without necessarily specifying how the cells are dependent. An RDR consists of two parts – a dependence statement (which itself has two parts: a head and tail) and a condition. The dependence statement specifies the dependency among cells, while the condition, is a query that identifies the cells on which a dependency holds and returns the rids corresponding to those cells.

*Definition 3.1 (Relational dependency Rules).* Let  $S = (\mathcal{R}, \mathcal{A})$  be a relational functional schema where  $A, A_1, \dots, A_p \in \mathcal{A}$  are attribute functions, and  $Q$  is a SQL query over the schema  $S$  that returns the rids  $X, X_1, \dots, X_p$  of records that satisfy the condition of the query. A *relational dependency rule (RDR)* is given by

$$\text{Dependence: } A(X) \not\perp A_1(X_1), \dots, A_p(X_p) \quad \text{Condition: } Q \quad (1)$$

In Eqn. 1, the *condition*  $Q$  identifies rids of a subset of records (from the set of relations in  $\mathcal{R}$ ) such that the records contain the attributes  $A, A_1, \dots, A_p$  which are dependent. The dependence part of the RDR  $A(X) \not\perp A_1(X_1), \dots, A_p(X_p)$  uses attribute function notation to express the dependency between attribute values amongst the selected records. In particular, the RDR expresses that value the attribute  $A(X)$  (*head* of the RDR) can take is **not independent** of the value of the attributes  $A_1(X_1), \dots, A_p(X_p)$  (*tail* of the RDR).

Thus, instantiated values of the attributes in the tail of the RDR can enable inference about the value of the head of the attribute.

We illustrate the RDR notation using a couple of semantic dependencies among the data in our example. In Fig. 1b, RDR R1 states the dependency between the number of likes (pLikes in Posts table) for a post authored by a person  $u$  and the total likes (totLikes in Statistics table) of  $u$ . The condition of R1 chooses rid pairs corresponding to the rids of the records in the Statistics and Post tables s.t. the two records are for the same person (due to join conditions in the WHERE clause).

As another illustration, we consider R2 (Fig. 1c) that states that the last location lStLoc of a person  $u$  and the location of the latest post (pLoc) made by  $u$  are dependent if the post was made before the time the location of the device of  $u$  was collected. We present a few other dependencies for our running example which will be used later. Let  $u$  be a person with a device  $d$  and makes a post  $p$ :

- **R3:** lStLoc  $\not\perp$  dLoc, i.e., the last known location (lStLoc) of  $u$  depends on the location of  $d$  (dLoc collected at dTm) if the last post made by  $u$ , pTm is earlier than dTm.
- **R4:** The location of  $p$  (pLoc) and  $u$ 's most frequented location (freqLoc) are dependent, i.e.,  $freqLoc(R) \not\perp pLoc(M)$ .
- **R5:** The location of  $d$  dLoc  $\not\perp$  pLoc of  $p$  if  $d$ 's location (dLoc) was collected within an hour of the time of  $p$ .
- **R6:** For  $u$ , Trv1  $\not\perp$  freqLoc, lStLoc.

RDRs may be discovered using existing work (such as [1, 2] which, motivated by data regulations, discovers data dependencies and annotations, and others like [8, 34]) or through data analysis. In our evaluations, we derived dependencies using existing work and manually. Moreover, rules used for data repairs, data cleaning, as well as provenance annotations can be easily expressed as RDRs. **Instantiated RDRs.** To use RDRs for specific cells, we need to define *instantiated RDRs*. Instantiations of a RDR, on a database state  $\mathcal{D}_t$ , are all possible dependencies (in  $\mathcal{D}_t$ ) between the relevant attribute functions (in  $\mathcal{D}_t$ ) which are specified in the dependence statement of the RDR and satisfy the condition of the RDR. When an attribute function is in the dependence as well as the condition, it is dropped from the latter. An instantiated RDR comprises just the dependence statement of the corresponding RDR, i.e., the head and the tail of the RDR, with all the variables  $X, X_1, \dots, X_p$  substituted with constants from the database state  $\mathcal{D}_t$  which are returned by the condition of the RDR. For example, the instantiations of RDR R1 in Fig. (1b) for Person(1) are:  $\delta_2^- : totLikes(1) \not\perp pstLikes(1)$ ,  $\delta_3^- : totLikes(1) \not\perp pstLikes(2)$ , and  $\delta_4^- : totLikes(1) \not\perp pstLikes(4)$  in  $\mathcal{D}_2$ .

*Definition 3.2 (Instantiations of RDRs).* Let  $\mathcal{D}_i$  be an instance over  $S = (\mathcal{R}, \mathcal{A})$ . Given an RDR as in Eqn. 1, for an instantiated attribute function  $A_k(x_k)$ , we associate the *instantiated RDR*, denoted  $\delta^-(\mathcal{D}_i, A_k(x_k))$ , obtained by assigning to the variables  $X = x, X_1 = x_1, \dots, X_p = x_p$  such that  $Val(A_1(x_1), \mathcal{D}_i), \dots, Val(A_p(x_p), \mathcal{D}_i)$  are returned by the query  $Q$  evaluated on  $\mathcal{D}_i$ . When clear from the context, we drop the database state  $\mathcal{D}_i$  and  $A_k(x_k)$  from  $\delta^-(\mathcal{D}_i, A_k(x_k))$  and write  $\delta^-$ . An instantiated RDR is of the following form where  $A_i(x_i)$  is the instantiated attribute function for  $A_i(X_i)$ .

$$\delta^- : A(x) \not\perp A_1(x_1), \dots, A_p(x_p) \quad (2)$$

We denote a set of RDRs with  $\Delta^-$ . Given a database state  $\mathcal{D}_t$ , we denote with  $\Delta^-(\mathcal{D}_t)$  the set of all instantiated RDRs on that state. With  $Head(\delta^-)$ , we refer to  $\{A(x)\}$ , the head of the instantiated RDR. The tail, denoted  $Tail(\delta^-)$  is given by the set  $\{A_1(x_1), \dots, A_p(x_p)\}$ . The condition is dropped. We denote with  $Cells(\delta^-)$  the set of all the attribute functions in  $\delta^-$ , i.e.,  $Head(\delta^-) \cup Tail(\delta^-)$ .

### 3.2 Dependency Sets

Given a set of instantiated RDRs, we define the notion of dependency sets. Dependency sets capture how an attribute function  $A(x)$  can be probabilistically influenced by, or can influence other data in the database. Instantiated RDRs can lead to the inference of an attribute function  $A(x)$  in two ways: (1) direct and (2) indirect.

**Direct inference.** Direct inference takes place through explicitly stated dependencies that involve  $A(x)$ .

*Example 3.3.* The instantiation  $\delta_5^- : \text{freqLoc}(1) \not\perp \text{pLoc}(4)$  of R4 and the instantiation  $\delta_6^- : \text{lastLoc}(1) \not\perp \text{pLoc}(4)$  of R2 lead to the direct inference of  $\text{pLoc}(4)$ .

**Indirect Inference.** An indirect inference on  $A(x)$  exists when a sequence of instantiated RDR can be used to infer  $A(x)$  through shared elements between each pair of instantiated RDRs. Continuing Eg. 3.3, let  $\delta_7^- : \text{lastLoc}(1) \not\perp \text{dLoc}(1)$  be another instantiated RDR. Observe that  $\text{pLoc}(1)$  does not directly depend on  $\text{devLoc}(1)$ . However,  $\text{pLoc}(1)$  depends on  $\text{lasLoc}(1)$  through  $\delta_6^-$ , which in turn depends on  $\text{dLoc}(1)$ .

*Definition 3.4.* Given a database state  $\mathcal{D}_t$  and an attribute function  $A(x) \in \mathcal{D}_t$ , and a set of  $\Delta^-(\mathcal{D}_t)$  of instantiated RDRs, we define the set of dependencies on  $A(x)$  in  $\mathcal{D}_t$ , denoted  $dep(A(x) \mid \mathcal{D}_t)$ , to contain  $\delta_i^- \in \Delta^-(\mathcal{D}_t)$  such that for all cells  $c \in Cells(\delta_i^-)$ , we have  $Val(c \mid \mathcal{D}_t) \neq NULL$  and one of the following holds: (1)  $A(x) \in Cells(\delta_i^-)$ ; (2) there exists  $\delta_1^-, \dots, \delta_\ell^-, \dots, \delta_K^-$  in  $\Delta^-(\mathcal{D}_t)$  such that, for  $1 < i \leq K$ , we have  $Tail(\delta_i^-) \cap Head(\delta_{\ell+1}^-) \neq \emptyset$  and  $A(x) \in Cells(\delta_1^-)$ .

### 3.3 Data Erasure Semantics

We formalize the semantic guarantees of data erasure in this section. We assume that at a given time, the database owner has access to the entire database at that time, and the dependencies. We do not consider adversarial scenarios wherein a malicious database owner maintains a copy of data secretly. Incorporating such a possibility goes well beyond the scope of erasure we are considering here.

Note that when a cell  $A(x)$  expires at  $t_e$ , its value is set to  $NULL$ . We write  $Val(A(x)) \leftarrow val$  to denote the assignment of the value  $val$  to the cell  $A(x)$ . In particular, with  $\mathcal{D}_{t_e^+} \cup \{Val(A(x)) \leftarrow val\}$  we denote the state that is identical to  $\mathcal{D}_{t_e^+}$  except that the cell  $A(x)$  which has the value  $val$ . We now formally define P2E2.

*Definition 3.5 (Pre-insertion Post Erasure Equivalence (P2E2)).* Given a set  $\Delta^-$  of RDRs, a cell  $A(x)$  with insertion time  $\kappa(A(x)) = t_b$ , expiration time  $\eta(A(x)) = t_e$  and  $Val(A(x), \mathcal{D}_{t_e}) = val$  that is not  $NULL$ , we say that *pre-insertion post erasure equivalence (P2E2)* holds for  $A(x)$  if:  $dep(A(x) \mid \mathcal{D}_{t_e^+} \cup \{Val(A(x)) \leftarrow val\}) \subseteq dep(A(x) \mid \mathcal{D}_{t_b})$ .

Informally, P2E2 states that the set of dependencies on an attribute function when it is inserted is the same as the set of dependencies on the attribute function after it has expired. Note that

$A(x)$  must be set to  $NULL$  after it has expired. In Fig. 1a, suppose Bob wants the location of their latest post (i.e., attribute  $\text{pLoc}$  in table  $\text{Posts}$ ,  $\text{rid}=4$  with  $\text{pID}=1004$ ) to be forgotten. We say P2E2 guarantee holds for  $\text{pLoc}(4)$  if the dependencies on  $\text{pLoc}(4)$  in state  $\mathcal{D}_1$  is the same as that in the state after the  $\blacktriangle$ , i.e., state  $\mathcal{D}_4$ .

We can now formally define the problem of minimal erasure.

*Definition 3.6 (OPT-P2E2).* Given a database state  $\mathcal{D}_t$ , the set of instantiated RDRs  $\Delta^-(\mathcal{D}_t)$ , and a cell  $A(x)$  in  $\mathcal{D}_t$ . Find a set  $\mathcal{T} = \{A(x)\} \cup \{A_i(x_i) \mid 1 \leq i \leq |\mathcal{T}|\}$  such that when the value of each  $A_i(x_i)$  is set to  $NULL$ , P2E2 holds in  $\mathcal{D}_{t^+}$  for  $A(x)$  and  $\sum_{A_i(x_i) \in \mathcal{T}} Cost(A_i(x_i))$  is minimized.

**THEOREM 3.7.** *The OPT-P2E2 problem is NP-HARD.*

### 3.4 Identifying Relevant Dependencies

Given a set of RDRs, and an attribute function  $A(x)$  for which P2E2 has to be guaranteed, where  $\kappa(A(x)) = t_b$  and  $\eta(\kappa)(x) = t_e$ , we need to instantiate RDRs to determine the sets  $dep(A(x) \mid \mathcal{D}_{t_b})$  and  $dep(A(x) \mid \mathcal{D}_{t_e})$  to compute  $\Delta^-(P2E2, A(x)) = dep(A(x) \mid \mathcal{D}_{t_e}) \setminus dep(A(x) \mid \mathcal{D}_{t_b})$ . Then we address each dependency in  $\Delta^-(P2E2, A(x))$  to ensure that P2E2 holds.

**Constructing  $dep(A(x) \mid \mathcal{D}_t)$ .** Recall that inference using RDRs can be direct or indirect. Direct inference occurs when an instantiated RDR  $\delta^-$  contains  $A(x) \in Cells(\delta^-)$ . Therefore we need to instantiate all such RDRs where  $A(x)$  is in the head or the tail of the dependence. To prevent indirect inference on  $A(x)$ , we need to consider all RDRs that lead to the direct inference of some attribute function  $A'(x)$  such that  $A'(x)$  leads to direct inference of  $A(x)$ .

With both sets of dependencies constructed as above, the difference  $\Delta^-(P2E2, A(x)) = dep(A(x) \mid \mathcal{D}_{t_e}) \setminus dep(A(x) \mid \mathcal{D}_{t_b})$  can be readily identified. However, constructing sets of dependencies is computationally intensive as it entails recursively instantiating RDRs to account for both direct and indirect inference.

**Constructing  $\Delta^-(P2E2, A(x))$ .** We note that instead of constructing the two sets of dependencies, namely,  $dep(A(x) \mid \mathcal{D}_{t_e})$  and  $dep(A(x) \mid \mathcal{D}_{t_b})$ , and then computing their difference, we can directly construct  $\Delta^-(P2E2, A(x))$ . To that end, let  $E$  be the set of cells erased and  $N$  be the set of cells inserted after state  $\mathcal{D}_{t_b}$ . Therefore,  $(\mathcal{D}_{t_b} - E) \cup N = \mathcal{D}_{t_e}$ . We first observe that the dependencies on  $A(x)$  in  $\mathcal{D}_{t_b} \setminus E$  are necessarily contained in the set of dependencies on  $A(x)$  in  $\mathcal{D}_{t_b}$ . This is because since some cells are erased, no new dependencies are introduced. Therefore, we only need to focus on the set of dependencies on  $A(x)$  introduced by the cells in  $N$ .

Observe that not all dependencies in the set  $dep(A(x) \mid \mathcal{D}_{t_e})$  are in the desired set  $\Delta^-(P2E2, A(x))$  of dependencies that violate P2E2. More specifically,  $dep(A(x) \mid \mathcal{D}_{t_e})$ , also contains dependencies that existed in the state  $\mathcal{D}_{t_b} - E$ . Observe that if all the cells of an instantiated RDR  $\delta^- \in dep(A(x) \mid \mathcal{D}_{t_e})$  were inserted before  $A(x)$ , it must have been in  $dep(A(x) \mid \mathcal{D}_{t_b} - E)$ , i.e., it does not violate P2E2. Consider the instantiated RDR in our running example  $\delta_5^-$  (E.g. 3.3), although this is in  $dep(\text{pLoc}(4) \mid \mathcal{D}_4)$ , it does not violate P2E2 as this dependency on  $\text{pLoc}(4)$  existed in the state  $\mathcal{D}_1$ . We show that for an instantiated RDR to be in  $\Delta^-(P2E2, A(x))$  it must contain at least one attribute function that was inserted after  $A(x)$ .

**THEOREM 3.8.** *Let  $A(x)$  be a cell with  $\kappa(A(x)) = t_b$  and  $\eta(A(x)) = t_e$ . Let  $E$  be the set of erased cells between the states  $\mathcal{D}_{t_b}$  and  $\mathcal{D}_{t_e}$ .*

The following holds: (1)  $dep(A(x) \mid \mathcal{D}_{t_b} - E) \subseteq dep(A(x) \mid \mathcal{D}_{t_b})$ ; (2)  $dep(A(x) \mid \mathcal{D}_{t_e}) - dep(A(x) \mid \mathcal{D}_{t_b} - E) = \{\delta_i^- \mid \delta_i^- \in dep(A(x) \mid \mathcal{D}_{t_e}) \wedge \exists c \in Cells(\delta_i^-) \text{ such that } \kappa(c) > t_b\}$ .

Consider the timeline in Fig. 1a, which shows the insertion of Bob's new post and other changes to the database. We want to guarantee P2E2 for the location of Bob's latest post (pLoc(4)). Not all the given dependencies hold for pLoc(4). RDR R3 is not applicable in this case. Since the difference between the time of the post (pTm(4)) and the time (dTm(1)) at which the location of Bob's device (dLoc(1)) was updated, pLoc(4) and dLoc(1) are not dependent. RDRs R2, R4, and R6 are applicable.

**Resolving Dependencies.** What remains is to resolve the dependencies in  $\Delta^-(P2E2, A(x))$  to ensure that the P2E2 condition for  $A(x)$  is satisfied. We want to identify a set  $\mathcal{T}$  of cells in  $\mathcal{D}_{t_e}$  such that when they are deleted (set to *NULL*), P2E2 guarantee holds for  $A(x)$ . We show that for each instantiated RDR that violates P2E2, we have to delete a cell from its head or tail.

## 4 SUPPORTING DEMAND-DRIVEN ERASURE

In this section, we focus on demand-driven erasure. Given a set of RDRs and a cell  $c_d$  for which P2E2 must hold, the first step is to instantiate the relevant RDRs (Sec. 4.1). Our first approach reduces OPT-P2E2 (Defn. 3.5) to ILP (Sec. 4.2) which can be solved using readily available ILP solvers. Optimal answers obtained from solvers often have high overheads. We develop a hypergraph-based approach (Sec. 4.3) that provides the optimal answer when RDRs are acyclic. Finally, we extend this to a heuristic approach (Sec. 4.4) that guarantees P2E2, but not at the least cost.

### 4.1 Instantiating RDRs

Given a database state  $\mathcal{D}_t$ , and an instantiated attribute function  $A(x)$ , we denote with  $\Delta^-(\mathcal{D}_t, A(x))$  the set of all instantiated relational dependency rules (RDRs) with  $A(x)$  in the head or tail. All RDRs with  $A(x)$  must be instantiated. Observe that to prevent *direct* and *indirect* inferences on  $A(x)$ , we need to consider instantiated RDRs in  $\Delta^-(\mathcal{D}_t, A(x))$  as well as all the instantiated attribute functions in  $\bigcup_{\delta^- \in \Delta^-(\mathcal{D}_t, A(x))} Cells(\delta^-)$  besides  $A(x)$ , and recursively so on.

---

#### Algorithm 1 Instantiating RDRs for P2E2

---

```

1: procedure DEPINST( $\mathcal{D}_t, \Delta^-, A(x)$ )
2:    $Q \leftarrow \{A(x)\}$  ▷ Queue of cells
3:    $\mathcal{V} \leftarrow \emptyset, \mathcal{I} \leftarrow \emptyset$  ▷ Lists of instantiated attributes & RDRs
4:   while  $Q \neq \emptyset$  do
5:      $attf \leftarrow Q.pop$ 
6:     if  $attf \notin \mathcal{V}$  then
7:        $\mathcal{V} \leftarrow \mathcal{V} \cup \{attf\}, Rules \leftarrow \Delta^-(attf)$ 
8:       for rule  $\in Rules$  do
9:          $\delta^- \leftarrow eval(rule), \mathcal{I} \leftarrow \mathcal{I} \cup \{\delta^-\}$ 
10:        for tail  $\in Tail(\delta^-)$  do
11:           $Q.push(tail)$ 
12:   return  $\mathcal{V}, \mathcal{I}$ 

```

---

Given  $A(x)$  and a database  $\mathcal{D}_t$ , Alg. 1 generates the set  $\Delta^-(P2E2, A(x))$  of instantiated RDRs. The algorithm iteratively instantiates RDRs

corresponding to direct or indirect inference of  $A(x)$ . The *eval()* function (in line 9) evaluates the condition of the rule, determines if the rule could have been used for inference in the state  $\mathcal{D}_{\kappa(A(x))}$  (Thm. 3.8) and returns the instantiated RDRs.

Alg. 1 performs a breadth-first traversal over the attribute function dependency graph, starting from a target cell  $A(x)$ . Let  $u$  be the number of unique attribute functions on which  $A(x)$  is dependent and  $r = |\Delta^-|$  the number of RDR templates. In the worst case, each of the  $u$  attributes may match up to  $r$  rules, each involving up to  $a$  attribute functions. Thus, the total number of instantiated rules is  $O(ura)$ , and the number of enqueued attributes is  $O(ura)$ . Since each attribute is visited at most once and each rule is evaluated once per attribute, the overall time complexity is  $O(ura)$ , and the space complexity is  $O(u + ur)$ . In practice, the algorithm is efficient when the dependency graph is sparse and rule arity is small.

### 4.2 ILP-Approach

We present a reduction from the OPT-P2E2 problem to integer linear programming (ILP). To provide an intuition of the reduction, we introduce some concepts and notation.

*Definition 4.1 (Induced Bipartite Graph).* For a cell  $A(x)$  in a state  $\mathcal{D}_t$ , given the set  $\Delta^-(P2E2, A(x)) = \{\delta_1^-, \dots, \delta_n^-\}$  of instantiated RDRs, we define the induced bipartite graph  $\mathcal{B}(\Delta^-(\mathcal{D}_t, P2E2(A(x)))) = (V = V_L \cup V_R, E = E_H \cup E_T)$  where  $V_L = \{\delta_i^- \mid 1 \leq i \leq n\}$  and  $V_R = \{c_j \mid c_j \in Cells(\delta_i^-)\}$  are the bipartition of the vertex set  $V$ . The set  $E = E_H \cup E_T$  of edges contains, for every  $\delta_i^- = c_{i_1} \dashv\dashv c_{i_2} \dots c_{i_{n_i}}$ , an edge  $(\delta_i, c_{i_1}) \in E_H$  and for,  $2 \leq j \leq n_i$ , edges  $(\delta_i, c_{i_j}) \in E_T$ .

**Reduction.** Consider, for a database state  $\mathcal{D}_t$ , a cell  $A(x)$ , denoted with  $c_d$ , for which P2E2 must hold, and for the set  $\Delta^-(\mathcal{D}_t, A(x))$  of instantiated RDRs, the induced bipartite graph  $\mathcal{B}(\Delta^-(\mathcal{D}_t, A(x))) = (V_L \cup V_R, E_H \cup E_T)$ . We introduce the following variables: for  $\delta_i \in V_L$  a binary variable  $b_i$ ; for each  $c_j \in V_R$ , a binary variable  $a_j$ ; for each edge  $(\delta_i, c_j) \in E_H$ , a binary variable  $h_i^j$ ; for each edge  $(\delta_i, c_j) \in E_T$  a binary variable  $t_i^j$ . We specify the constraints in the following.

- (1) For P2E2 to hold for unit  $c_d$ , it must be erased. So,  $a_d = 1$ .
- (2) For each unit  $c_j$ , where  $1 \leq j \leq m$ , that needs to be erased, all instantiated RDRs where  $c_j$  is in the head can be used to infer it. To prevent this, we need to address the instantiated RDR. This is stated, for all  $i$ , using the constraints  $a_j = h_i^j$ .
- (3) For each instantiated RDR  $\delta_i$ , if the unit in its head is hidden, then we set  $b_i = h_i^i$ .
- (4) For an instantiated RDR  $\delta_i$ , where  $1 \leq i \leq n$ , to prevent the inference of the unit in its head, a unit from the tail has to be erased. So we introduce the constraint  $\sum_{j \in \{i_1, \dots, i_{n_i}\}} t_i^j \geq b_i$ .
- (5) For each unit  $c_j$ , for  $1 \leq j \leq m$ , if the unit is erased, then all of the tail edges incident on it must indicate so which is ensured by, for all  $1 \leq i \leq n$ , we have  $a_j = t_i^j$ .
- (6)  $W = \min \sum_{j=1}^m a_j$  minimizes the number of units erased.

Let the ILP system above be  $O$ . It has  $O(nm)$  binary variables and  $O(mn)$  constraints. For a cost model where, for  $1 \leq j \leq m$ , the cost of erasing cell  $c_j$  is  $Cost(j) \in \mathbb{N}$ , we use  $W = \min \sum_{j=1}^m c_j Cost(j)$  to minimize the total cost of P2E2-guarantee.

### 4.3 Dependence Hypergraph

Here, we present a new approach to solve OPT-P2E2. Given a set of instantiated RDRs  $\Delta^-(P2E2, A(\mathbf{x}))$  we construct a *dependence hypergraph* where the vertices are cells and a hyperedge connects, for an instantiated RDR  $\delta^-$ , its head  $Head(\delta^-)$  to its tail  $Tail(\delta^-)$ .

*Definition 4.2 (Dependence Hypergraph).* For a cell  $A(\mathbf{x})$  in a database state  $\mathcal{D}_t$ , given the set  $\Delta^-(P2E2, A(\mathbf{x}))$  of instantiated RDRs, we define the *dependence hypergraph*  $\mathcal{H}(\Delta^-(P2E2, A(\mathbf{x}))) = (V, E)$  where  $V = \bigcup_{\delta_i^- \in \Delta^-(\mathcal{D}_t)} Cells(\delta_i^-)$  is the set of vertices and  $E = \{(Head(\delta_i^-), Tail(\delta_i^-)) \mid \delta_i^- \in \Delta^-(\mathcal{D}_t)\}$  is the set of edges.

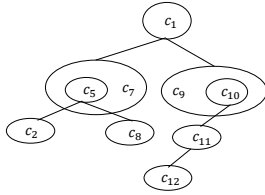


Figure 2: Example of a dependence hypergraph.

For a dependence hypergraph  $\mathcal{H}(\Delta^-(P2E2, A(\mathbf{x}))) = (V, E)$ , a vertex  $v \in V$  is called a *root* if  $v$  is not in the tail of any instantiated RDR, i.e., for all  $\delta_i^- \in \Delta^-(\mathcal{D}_t)$  we have  $v \notin Tail(\delta_i^-)$ . Similarly, a vertex  $v \in V$  is called a *leaf* if  $v$  is not the head of any instantiated RDR, i.e., for all  $\delta_i^- \in \Delta^-(\mathcal{D}_t)$  we have  $v \notin Head(\delta_i^-)$ . Fig. 2 shows the dependence hypergraph for the set  $\{\delta_1^- : c_1 \not\perp c_5, c_7; \delta_2^- : c_5 \perp c_2; \delta_3^- : c_5 \not\perp c_8; \delta_4^- : c_1 \not\perp c_9, c_{10}; \delta_5^- : c_{10} \perp c_{11}; \delta_6^- : c_{11} \perp c_{12}\}$ . The vertex  $c_1$  is a root and  $c_9$  and  $c_8$  are leaves.

Next, we define paths and complete paths in a dependence hypergraph to characterize the P2E2-guarantee.

*Definition 4.3.* For a dependence hypergraph  $\mathcal{H}(\Delta^-(P2E2, A(\mathbf{x})))$ , and a sequence  $P : v_1, v_2, \dots, v_n$  of vertices, we define the following.

- We say that the sequence  $P$  is a *path* if the following hold:
  - (1) There exists  $\delta_i^-$  s.t.  $v_1 \in Head(\delta_i^-)$ .
  - (2) For  $1 \leq i < n$ , there exists  $\delta_i^-$  s.t.  $v_i \in Head(\delta_i^-)$  and  $v_{i+1} \in Tail(\delta_i^-)$
- We say that a sequence  $P' : v_b, \dots, v_e$  where  $1 \leq b \leq e \leq n$ , is a *subpath* of  $P$  if  $P$  and  $P'$  are both paths. A path is, trivially, a subpath of itself.
- We say that the sequence  $P$  is a *complete (sub-)path* if  $P$  is a (sub-)path such that  $v_n$  is a leaf.

In Fig. 2, the vertex  $c_1$  is a root, vertices  $c_7$  and  $c_8$  are leaves; the sequence  $P_1 : c_1, c_{10}, c_{11}$  is a path; the sequence  $P_2 : c_1, c_5, c_2$  is a complete path;  $P_3 : c_5, c_2$  is a subpath of  $P_2$ .

We can now use the definition of path above to characterize when P2E2 holds. In particular, if there exists a path in  $\mathcal{H}(\Delta^-(P2E2, A(\mathbf{x})))$  such that for all vertices on the path, there is a subpath. We show that this is both a necessary and a sufficient condition for P2E2. Note that simply defining a cover [35] would have been sufficient but not necessary for P2E2 to hold. In the dependence hypergraph in Figure 2, it suffices when the units  $c_1, c_7$ , and  $c_9$  are set to NULL for P2E2 to hold (but they do not cover the graph).

**Optimization.** Given a set  $\Delta^-$  of RDRs, a database state  $\mathcal{D}_t$ , an instantiated attribute function  $A(\mathbf{x})$ , Algorithm 1 can be easily adapted

### Algorithm 2 Optimizing the Dependence Hypergraph

```

1: procedure OPTPATH( $\mathcal{V}, \mathcal{I}$ ) ▷ Output of Algorithm 1
2:    $Q \leftarrow \{leafs(\mathcal{V})\}, S \leftarrow \emptyset$  ▷ Queue of cells and set of seen cells
3:   while  $Q \neq \emptyset$  do
4:      $attf \leftarrow Q.pop$ 
5:     if  $attf \notin S$  then
6:        $S \leftarrow S \cup \{attf\}$ 
7:        $Cost(attf) \leftarrow 1$  ▷ Initialize node cost
8:        $Rules \leftarrow \mathcal{I}(attf)$  ▷ All RDRs that contain attf
9:       for  $\delta^- \in Rules$  do
10:        if  $attf \in Head(\delta^-)$  then
11:           $Sum(Cost(attf), \min Cost(Tail(\delta^-)))$  ▷
12:             $Sum(Cost(attf), \min Cost(Tail(\delta^-)))$  Add cheapest node of tail to node's cost
13:        else
14:           $Q.push(Head(\delta^-))$  ▷ Add RDR head to queue
15:    $\mathcal{T} \leftarrow \{A(\mathbf{x})\}, Q \leftarrow \{A(\mathbf{x})\}, S \leftarrow \emptyset$  ▷ Set of cells to delete
16:   while  $Q \neq \emptyset$  do
17:      $attf \leftarrow Q.pop$ 
18:     if  $attf \notin S$  then
19:        $S \leftarrow S \cup \{attf\}, Rules \leftarrow \mathcal{I}(attf)$ 
20:       for  $\delta^- \in Rules$  do
21:        if  $attf \in Head(\delta^-)$  then
22:           $child \leftarrow \arg \min Cost(Tail(\delta^-))$ 
23:           $\mathcal{T} \leftarrow \mathcal{T} \cup child$ 
24:           $Q.push(child)$  ▷ Delete cheapest node in tail
25:            and continue traversal
26:   return  $\mathcal{T}$ 

```

to construct a dependence hypergraph. Given a dependence hypergraph, Algorithm 2 produces an optimal solution. We assume that the graph is cycle-free. This assumption is necessary only for proving optimality<sup>5</sup>.

In practice, we implement the following procedure. The cost of every node is set to 1, as a base cost to erase a single attribute function (Line 7). The *leafs* cannot incur a higher cost, as they do not cause additional erasures. Next, we traverse the tree upwards and compute the cost of every *inner node* as the sum of the nodes' cost and the minimal cost of every attached hyper edge (Line 11). When reaching the root node, we construct the complete path by traversing the tree to the bottom and always choosing the node with minimal cost (Line 21). Thus, the algorithm guarantees P2E2 at minimal cost but has to traverse the tree twice.

**THEOREM 4.4.** *When Algorithm 2 terminates, the set  $\mathcal{T}$  contains cells which, when set to NULL, guarantees P2E2 for the input  $A(\mathbf{x})$ . Moreover,  $\sum_{A_i(\mathbf{x}_i) \in \mathcal{T}} Cost(A_i(\mathbf{x}_i))$  is minimized.*

The algorithm consists of two phases: a bottom-up cost propagation to assign minimal deletion costs, followed by a top-down traversal to extract the optimal deletion set. Let  $n = |\mathcal{J}|$  be the number of instantiated RDRs,  $a$  the maximum rule arity, and  $C_{max}$  the maximum cell cost. The bottom-up phase runs in  $O(n \cdot a \cdot \log C_{max})$  time,

<sup>5</sup>If cycles are present, our approach produces a correct solution but not necessarily an optimal one. To ensure that cycles are not present in the hypergraph, for all pairwise instantiated dependencies  $\delta_1^-$  and  $\delta_2^-$ , if  $Tail(\delta_1^-) \cap Tail(\delta_2^-) \neq \emptyset$ , discard the instantiated RDR with the larger number of attribute functions in its tail. The solution produced by Algorithm 2 is optimal with respect to the thus obtained set of RDRs.

and the top-down phase adds  $O(C_{\max} \cdot a \cdot \log C_{\max})$  for path extraction. Hence, the total time complexity is  $O((n + C_{\max}) \cdot a \cdot \log C_{\max})$ , with space complexity  $O(n \cdot a)$ . In practice, the algorithm is efficient due to typically low arity and shallow dependency paths.

#### 4.4 Approximate Algorithm

In this section, we present an approximation variant of Algorithm 2, which, given a cell, determines a (not necessarily the smallest) set of dependent cells to delete for P2E2 to hold.

We adapt the algorithm to determine the minimum cost of guaranteeing P2E2 (Alg. 2) such that instead of constructing the entire dependence hypergraph and then traversing it bottom-up, a partial top-down construction of the dependence hypergraph is sufficient: whenever there is an instantiated RDR that has more than two attribute functions, we only instantiate the next *level* in the tree and choose the one that has the lowest cost to erase. The other attribute functions are not instantiated fully, thereby saving time.

The algorithm greedily constructs a partial dependence hypergraph top-down and avoids the bottom up traversal as in Alg. 2. Since each of the  $n$  RDRs (each with maximum arity  $a$ ) is instantiated at most once for each of its  $a$  cells, the total number of instantiated cells is  $O(an)$ . As the hypergraph is traversed only once, the time complexity is  $O(an)$  and the space complexity is  $O(an)$ .

Observe that greedily selecting the attribute function with the lowest deletion cost can lead to suboptimal outcomes—for example, forcing the deletion of a high-cost function later. Let  $\mathcal{T}$  denote the deletion set from the greedy algorithm and  $\mathcal{T}^*$  the optimal set. In general, the cost ratio  $\frac{\text{Cost}(\mathcal{T})}{\text{Cost}(\mathcal{T}^*)}$  can be unbounded. However, when minimizing the cardinality of the deletion set, if the number of cells per instantiated RDR is bounded by arity  $a = \max_{\delta^- \in \Delta^-(P2E2, A(x))} |\text{Cells}(\delta^-)| \geq 2$ , if a cell is in at most  $d$  instantiated RDRs, i.e.,  $d = \max_{c' \in \cup_{\delta^- \in \Delta^-(P2E2, A(x))} \text{Cells}(\delta^-)} |\{\delta^- | c' \in \text{Cells}(\delta^-)\}|$ , and the number of cells in the acyclic graph is  $n$ , then we can bound the approximation ratio as follows.

**THEOREM 4.5.** *Given an acyclic dependence hypergraph  $\mathcal{H}(\Delta^-(P2E2, A(x))) = (V, E)$ , with  $|V| = n$  and  $|\Delta^-(P2E2, A(x))| = r$ , let  $\mathcal{T}^*$  be the minimal set of cells that need to be deleted to guarantee P2E2 for  $A(x)$  and  $\mathcal{T}$  be the set of cells to be deleted to guarantee P2E2 with the approximation algorithm (in Sec. 4.4). The following holds:  $\frac{|\mathcal{T}|}{|\mathcal{T}^*|} = \min\left\{\frac{d}{r} \cdot \left(1 + \log_2\left(\frac{a \cdot r}{d}\right)\right), \frac{ad}{n-1} (1 + \log_2(an))\right\}$ .*

#### 4.5 Batching Erasures

The approaches discussed until now focus on the erasure of one cell. Since regulatory data erasure requirements allow for a reasonable delay between the time at which the data is requested to be erased and the actual erasure of the data (referred to as grace period and denoted with  $\Gamma$ ), it is possible to batch data erasures. The grace period can be used to batch multiple data erasure requests and instead of constructing and solving an individual optimization model for each cell, we attempt to construct models that allow for multiple cells to be erased such that the P2E2 holds for each of them.

Intuitively, we instantiate RDRs for the cells to be erased, which maximizes the possibility that the corresponding dependence hypergraphs have shared vertices. In practice, over a  $\Gamma$  period of time, we collect all cells that have to be erased such that P2E2 holds for

them. Let this set of cells be  $S$ . We instantiate RDRs for each cell  $s \in S$  at a time. Whenever an instantiated RDR corresponding to  $s$  contains a dependent cell  $s'$  also in  $S$ , we mark it to be set to *NULL* and only instantiate the RDRs for  $s'$ . This not only reduces the number of RDRs instantiated and, thus, the number of leafs in the tree, but also the time taken to traverse the tree. Moreover, fewer models (ILP or hypergraphs) need to be constructed and optimized.

### 5 RETENTION-DRIVEN ERASURE

So far, we have considered demand-driven erasures (a user wants to erase a cell  $c$  before its expiration time  $\eta(c)$ ) and batching such erasures. Now we turn to retention-driven erasure (where cell  $c$  is erased at its preset  $\eta(c)$ ). For such erasures, we investigate how to minimize the overheads of P2E2 on derived data. For base data, we adopt the batching approach discussed in Sec. 4.5.

Guaranteeing P2E2 for cells, often requires additional and potentially undesirable update and reconstruction of derived data. For example, suppose, for a derived cell  $c$ , with the parameter  $\text{freq}(c) = 1hr$ , depends on cells  $c_1, c_2$ , and  $c_3$ . It is reconstructed at 1pm, 2pm, 3pm, and so on. The cells  $c_1, c_2$ , and  $c_3$  expire at 1:30pm, 3:00pm, and 4:30pm, respectively. To guarantee P2E2 for the dependent cells, cell  $c$  needs to be reconstructed at 1pm, 1:30pm, 2:30pm, 3pm, 4pm, and at 4:30pm thus incurring additional overheads.

Retention-driven erasures offer an opportunity to reduce additional reconstructions due to P2E2 by exploiting the already known expiration times. We present an algorithm that, given a derived cell  $c$  and its dependencies  $c_1, \dots, c_n$  with corresponding erasure time intervals<sup>6</sup>  $(\eta_1^b, \eta_1^e), (\eta_2^b, \eta_2^e), \dots, (\eta_n^b, \eta_n^e)$ , determines an erasure schedule  $Sch(c)$  that takes into account when derived data has to be erased while maintaining the invariant that  $c$  is reconstructed at least once every  $\text{freq}(c)$ . Intuitively, we progressively build the reconstruction schedule  $Sch(c)$  by determining the maximum overlap between the retention periods of the dependent data to minimize the number of extra reconstructions due to P2E2.

#### Algorithm 3 Reconstruction Scheduler

---

```

1: procedure CREATESCHEDULE( $c, \text{freq}(c), \text{depSet}$ )
2:    $\kappa(c) \leftarrow \text{time.now}()$ 
3:    $Sch[0] \leftarrow \text{MAXOVERLAP}(\text{depSet}, \kappa(c) + \text{freq}(c))$ 
4:    $\text{depSet} \leftarrow \text{depSet} \setminus \{c_i \mid Sch[0] \in (\eta_i^b, \eta_i^e)\}$ 
5:    $i \leftarrow 1$ 
6:   while  $\text{depSet} \neq \emptyset$  do
7:      $Sch[i] \leftarrow \text{MAXOVERLAP}(\text{depSet}, Sch[i-1] + \text{freq}(c))$ 
8:      $\text{depSet} \leftarrow \text{depSet} \setminus \{c_i \mid Sch[i] \in (\eta_i^b, \eta_i^e)\}$ 
9:      $i \leftarrow i + 1$ 
10: procedure UPDATESCHEDULE( $c, \text{freq}(c), Sch, \text{depSet}$ )
11:   for  $(c_i, \eta_i^b, \eta_i^e) \in \text{depSet}$  do
12:     if  $\eta_i^b > Sch[0] \wedge \eta_i^e < Sch[1]$  then
13:       CREATESCHEDULE( $c, \text{freq}(c), \text{depSet}$ )

```

---

Our algorithm (Algorithm 3) is in two parts: Step 1 (Lines 1-11) creates the reconstruction schedule  $Sch(c)$  of the cell  $c$ , and Step 2

<sup>6</sup>Erasure time interval refers to the time interval in which a cell has to be erased. Usually  $\eta_i^b + \Gamma = \eta_i^e$ . However, here we allow for cells to have different time intervals in which they must be erased.

(Lines 12-18) updates the schedule when required to ensure that newly inserted dependencies are accounted for. At any given time,  $depSet$  for a derived cell  $c$  denotes the set  $\{(c_i, \eta_i^b, \eta_i^e) \mid 1 \leq i \leq n\}$  of all its dependencies, their insertion time, and their erasure time, respectively. The `MAXOVERLAP` function is a standard algorithm to find the maximum overlap given a set of time intervals.

Step 1: The first step finds a reconstruction time  $\rho$  that maximizes for  $1 \leq i \leq n$  the overlap between the erasure time intervals  $(\eta_i^b, \eta_i^e)$  of the dependent cells, and the time interval  $(\kappa(c), \kappa(c) + freq(c))$  in which  $c$  must be reconstructed at least once. Therefore, P2E2 holds for any cell  $c_i$  where  $\rho \in (\eta_i^e, \eta_i^b)$ . The algorithm iteratively finds the maximum overlap and creates a list  $Sch(c) : \rho_1, \dots, \rho_m$  of reconstruction times for  $c$  such that P2E2 holds for its dependencies in the database at the time of the construction of the schedule.

Step 2: The update procedure is called when a derived cell  $c$  is reconstructed. It checks whether there exists a dependent cell  $c_i$  that has to be erased after the current reconstruction but before the next scheduled reconstruction. If it does, then a new reconstruction schedule is created using the first step described above. Observe that if the erasure time of any dependent cell  $c_i$  is updated, this step (Step 2) ensures that P2E2 holds.

Let  $n$  be the number of dependent base cells for a derived cell  $c$ . Since the algorithm utilizes a greedy interval scheduling problem, the time complexity is  $O(n \log n)$  due to sorting, and the space complexity is  $O(n)$  to store intervals and the resulting schedule.

## 6 EVALUATION

We evaluate our approaches for guaranteeing P2E2. We compare the ILP approach (Sec. 4.2) with the graph-based algorithm, HGR (Sec. 4.3) and its approximate version, APX (Sec. 4.4). We analyze the efficiency and effectiveness of all the algorithms when applied to individual demand-driven erasures as well as with a set of erasures using our batching method, BATCH (Sec. 4.5). Additionally, we investigate our approach for retention-driven erasures, SCHEDULER (Alg. 3).

### 6.1 Experimental Setup

All experiments were run on an Ubuntu-based (20.04 LTS) server (Intel Xeon E5-2650; RAM: 256 GB). All algorithms are single-threaded, running on Java 11 and the datasets are stored in a PostgreSQL (v12.20) database. The ILP approach uses the Gurobi (v11.03) solver. **Datasets.** We use the following five datasets in our evaluations (the first four columns of Table 3a summarizes them and shows the number of RDRs and where/how they were derived.)

(1) *Twitter* [21]. This dataset contains a subset of tweets posted on  $\mathbb{X}$  (formerly Twitter) that represents a real-world instance of our running example. The RDRs were designed manually and express dependencies between individuals and their content.

(2) *Tax* [9]. This dataset is a synthetic dataset created using the real-world distribution of values of American tax records. We discovered all present DCs using [43]. The top-10 DCs that are not entirely comprised of equality predicates are transformed into RDRs. The RDRs include the conditional functional dependencies used in the original publication for data cleaning.

(3) *SmartBench* [32]. This dataset is based on real data collected from sensors deployed throughout the campus at the University of California, Irvine. RDRs capture dependencies between multiple

physical sensors which are used to compute derived metrics, e.g., occupancy from Wi-Fi AP locations.

(4) *HotCRP* [36]. This is a dataset containing a sample of real-world conference data. It stores authors, papers, conferences, and their relationships. The RDRs are generated by the method of [1].

(5) *TPC-H* [54]. This well-known benchmark dataset stores transactions of commercial actors: customers and their orders, as well as suppliers that fulfill those. The RDRs capture the links between the tables, i.e., foreign keys. In TPC-H, both customers and suppliers may delete their data. We create two separate scenarios and combine them proportionally to the number of customers and suppliers.

The set of RDRs for each dataset is cycle free. Some RDRs in the Twitter and SmartBench dataset join on non-key columns. To speed up the instantiations, we index those columns separately.

**Metrics.** We measure the (i) total number of deletions, (ii) time taken, and (iii) space overheads to guarantee P2E2. For demand-driven erasure, we also measure the number of reconstructions.

**Workload.** Given the lack of suitable deletion benchmarks [48], we evaluate demand- and retention-driven erasures separately, as well as varying combinations of each.

### 6.2 Experiments

**Experiment 1.** To test the cost of demand-driven erasures, we erase 100 random cells for each base data attribute in the RDRs. We depict the average runtime and model size for a single erasure in Fig. 3b. The total runtime is divided into the following four steps: (i) RDR instantiation (Algorithm 1); (ii) model construction (graph construction in HGR and APX and defining the ILP instance); (iii) model optimization (traversing the graph and keeping track of the minimal cost erasure for HGR and solving the ILP; there is no optimization phase in APX, as it greedily chooses the next edge to process); (iv) update to *NULL* (which modifies the database to guarantee P2E2.) We measure the average number of instantiated and deleted cells beyond the initial erased cell and compare our three P2E2 algorithms against three baseline implementations inspired by cascading deletions [27, 29]. These baselines guarantee P2E2 by identifying and deleting dependent cells but, unlike our approaches, they cannot account for the state of the database at the time of insertion of the data being deleted. Each baseline scans the entire dataset but varies in how it selects which cells to delete: *INST* deletes all dependent (i.e., instantiated) cells, mimicking traditional cascading deletes; *OPR* improves this and deletes exactly one cell per instantiated RDR; finally, *MINSET* computes a minimal set of cells such that at least one cell is deleted for each RDR (see Table 3a).

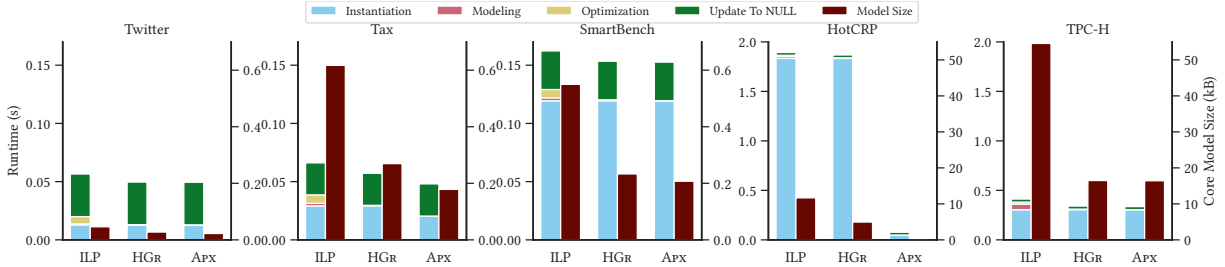
*Comparison with baselines:* As expected, all three baseline methods delete significantly more data than our P2E2 algorithms on every dataset except Tax. This is because, to guarantee P2E2 for a cell, we only need to consider data that was inserted after the cell and that is dependent on it. In high-volume datasets like Twitter and SmartBench, where a user generates many entries, the overhead of baselines is especially pronounced. In contrast, the Tax dataset does not incur additional deletions under *MINSET* because all user data is inserted simultaneously.

*Analysis of our P2E2 mechanisms:* We observe a stark difference between Twitter, Tax, SmartBench, and HotCRP and TPC-H. We

<sup>7</sup>[https://huggingface.co/datasets/enryu43/twitter100m\\_tweets](https://huggingface.co/datasets/enryu43/twitter100m_tweets)

| Dataset                  | # Cells    | # RDRs (Source) | # base, # derived | Instantiated cells |       |       | Deleted cells |       |       | Deleted cells for Baselines |        |        |
|--------------------------|------------|-----------------|-------------------|--------------------|-------|-------|---------------|-------|-------|-----------------------------|--------|--------|
|                          |            |                 |                   | ILP                | HGr   | Apx   | ILP           | HGr   | Apx   | INST                        | OPR    | MINSET |
| Σ (Twitter) <sup>7</sup> | 21 926 096 | 11*             | 7, 11             | 0.29               | 0.29  | 0.29  | 0.29          | 0.29  | 0.29  | 837.21                      | 517.8  | 321.1  |
| Tax [9]                  | 16 000 004 | 10 ([43])       | 15, 4             | 6.93               | 6.93  | 5.14  | 4.07          | 4.07  | 4.2   | 6.93                        | 6.6    | 4.07   |
| SmartBench [32]          | 94 424     | 5*              | 5, 1              | 5.43               | 5.43  | 5.43  | 5.43          | 5.43  | 5.43  | 1021.8                      | 1021.8 | 1021.8 |
| HotCRP [36]              | 122 697    | 3([1])          | 6, 3              | 134.4              | 134.4 | 4.81  | 3.78          | 3.78  | 3.78  | 815.0                       | 721.7  | 74.5   |
| TPC-H [54]               | 5 825 639  | 6 ([29])        | 8, 4              | 17.63              | 17.63 | 17.63 | 17.63         | 17.63 | 17.63 | 722.2                       | 722.2  | 722.2  |

(a) Summary of datasets and average number of instantiated & deleted cells. INST = all instantiated cells, OPR = one cell per rule, MINSET = minimal set, \*manually designed.



(b) Average runtime and space overhead

Figure 3: Evaluation of demand-driven erasures

highlight the similarity within these groups by using the same axis scaling. RDRs for the HotCRP and TPC-H datasets were created using data dependencies from schema constraints (or IND discovery). DBMSs already include a mechanism to delete data linked by foreign keys. Thus, there is no overhead to guarantee P2E2.

Interestingly, neither the number of rules, nor the dataset size determine the deletion complexity. P2E2 is more sensitive to the amount of related data, i.e., the number of instantiated cells. Consequently, the instantiation time takes the most amount of time for the SmartBench, HotCRP, and TPC-H dataset. In contrast, the instantiation time of the Tax dataset is lower although the number of instantiated cells is higher than in the SmartBench dataset. The RDRs in the Tax dataset specify connections only within one row of the data. For each deleted cell, we repeatedly query the same row, which is faster than scanning larger parts of the data as observed in SmartBench. In the Twitter dataset, the final update step dominates the cost, as the number of instantiated cells is low.

We observe that the ILP approach has the highest overheads (runtime and memory) for all three datasets. However, it is optimal in that it guarantees P2E2 using a minimum set of additionally deleted cells. Likewise, HGr always produces an optimal result, but consumes significantly less memory. Both approaches need to instantiate all available RDRs for all applicable cells, so their instantiation time is identical. However, the model construction and optimization overhead for HGr is negligible compared to the ILP approach. Apx instantiates fewer cells, so it outperforms the other approaches for all datasets. This behavior is especially noticeable for HotCRP. The optimal models have to instantiate a long chain of RDRs that turn out to be irrelevant to identify the cheapest deletion set. Due to its greedy nature, Apx avoids instantiating all those RDRs and significantly outperforms the rest of the algorithms. Since it does not keep track of an erasure cost, it also consumes less memory. However, it cannot guarantee optimality for its result set, as exemplified in the Tax dataset (see Table 3a).

**Experiment 2.** To investigate the influence of the degree of dependence (determined by the count of the number of cells that are part of instantiated RDRs) in the data on the number of deleted cells, we conducted the following experiment. For each dataset, we randomly sample 100 erasures. Each sampled erasure is processed by HGr given every subset of RDRs. In Fig. 4, we plot the average number of instantiated cells compared to the average number of deleted cells. The size of the data points signals the number of RDRs present in the processed subset. Intuitively, the more dependent cells a cell has, the more needs to be deleted. The number of rules does not have an immediate effect on the number of deleted cells, as larger and smaller points are mixed along the general trend line. However, the HotCRP dataset is an outlier. Due to the aforementioned long chain of instantiated RDRs, the number of instantiated cells is high, while the necessary deletions remain low.

**Experiment 3.** We evaluate the impact of BATCH, which exploits the grace period  $\Gamma$  to batch as many erasures as possible. First, we batch erasures based on a time interval. The number of erasures in such a batching strategy is workload-dependent. We further create batches based on fixed number of erasures to study the impact of batching across datasets in a workload-independent fashion.

*Batching based on time:* To unify the process for the Twitter and SmartBench datasets, we randomly sample 1000 erasures from a twelve-hour window in our data. For the Tax dataset, we assume

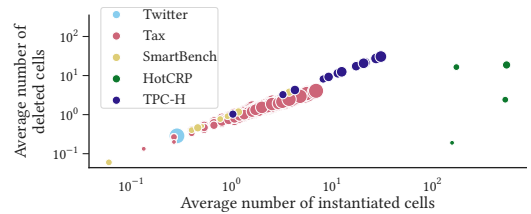


Figure 4: Impact of dependencies (log-axis) on P2E2 overhead; the size of data points signals number of RDRs

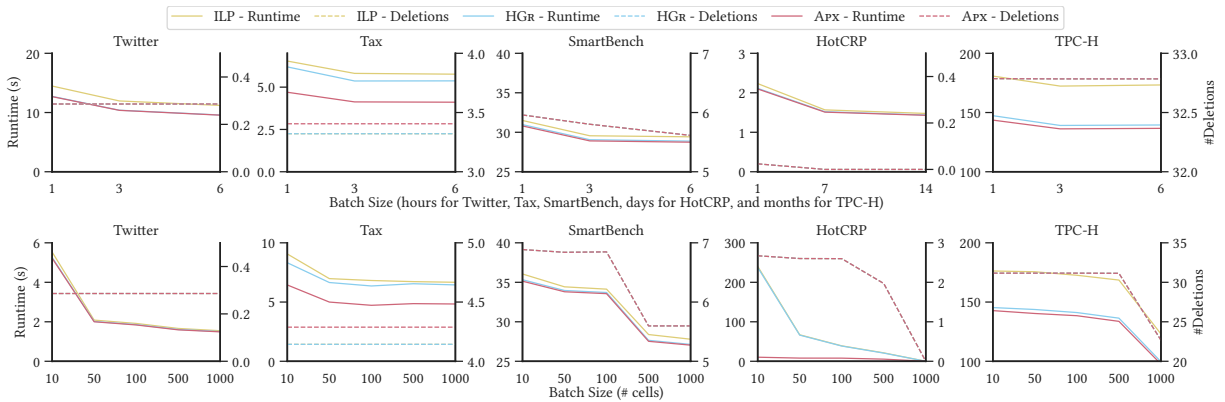


Figure 5: Evaluation of the batching

100 record updates per hour. After sampling the erasures, we use three different grace periods: one, three, and six hours. The HotCRP and TPC-H dataset operate on a different timescale. Therefore, we sample one month for the HotCRP dataset and one year for the TPC-H dataset. They use the batch sizes of one, seven, fourteen days, and one, three, and six months, respectively. The cumulated runtime and number of cells deleted for each batch size is depicted in Fig. 5. In general, larger batch sizes require fewer cells to hide and are processed faster. Initially, batching provides a larger benefit, as the impact of finding the first cells that are already instantiated is larger. In the Twitter and Tax dataset, the number of additionally deleted cells stays constant. In contrast, it scales similar to the runtime improvement in the HotCRP and TPC-H dataset. While in the SmartBench dataset, the number of additionally deleted cells scales linearly, the runtime exhibits a steeper slope between the batch sizes of one and three hours.

*Batching based on cardinality:* We randomly sampled 1000 erasures from the entire dataset, and grouped them in batches of size 10, 50, 100, 500, and 1000. We present the runtime and number of deleted cells for the entire set of erasures in Fig. 5. We observe a similar pattern as in the time-based method, larger batches perform better. In the Twitter, Tax, and HotCRP datasets, the scaling trend is similar to the previous experiment. For SmartBench, the biggest improvement occurs when enlarging the batch size from 100 to 500 cells. As the erasures are sampled uniformly from the entire dataset, they are less likely to overlap in the instantiated cells and improve runtime. Enlarging the batch size increases this likelihood. Similarly, in the larger TPC-H dataset, this improvement only happens for the largest batch size. In the HotCRP dataset, HGR and ILP are as performant as APX because there is no need to instantiate the long chain of RDRs, if parts of it are already in the batch.

**Experiment 4.** Here, we evaluate the effectiveness of SCHEDULER in reducing extra reconstructions of derived data for retention-driven erasures. As none of the datasets except SmartBench contained derived data, we manually added aggregate statistics and applicable RDRs to them. These represent real-world examples of typical derived data, e.g., the total likes of a Twitter profile, or the total sales of a supplier in TPC-H. On the one hand, they need to be updated regularly to reflect the underlying base data. On the other hand, it

is cost-prohibitive to compute them on the fly, so minimizing the number of necessary reconstructions is desirable.

We randomly sampled 100 cells of each derived data attribute and sequentially deleted all their associated base data. Based on the time-frame of the dataset, we varied the *freq* and the grace period,  $\Gamma$ . For both the  $\mathbb{X}$  and Tax dataset, we set a base frequency of one day and vary the grace period between 0–23 hours. The SmartBench dataset uses a *freq* of one hour and  $\Gamma$  between 0–55 minutes. For the HotCRP and TPC-H datasets, we choose a base frequency of one week and one month, and  $\Gamma$  between 0–6 and 0–30 days, respectively. In Fig. 6, we depict the average number of reconstructions that we save compared to the baseline. We observe the effectiveness of SCHEDULER in all cases, but it differs depending on the update characteristic of the dataset. There are three distinct patterns visible. First, HotCRP and SmartBench are insensitive to an increase of the grace period because only base data from the same time is aggregated. Thus, the maximal saving is reached as soon as we allow scheduling. The actual improvement (55.8% for HotCRP, 4.4% for SmartBench) differ based on the number of base data cells that are aggregated into a derived cell. This difference is also apparent between the Twitter and TPC-H dataset. Both these datasets experience “bursty” updates, so initially increasing the grace period reduces the number of reconstructions significantly, but the effect flattens off. The third update characteristic is exhibited in the Tax dataset. It is continuously updated, so there is a steady reduction of necessary reconstructions. Since no two updates happen at the same time, there is no benefit without a grace period. Given the large amount of base data cells per derived data cell, the overall improvement is the largest in this dataset.

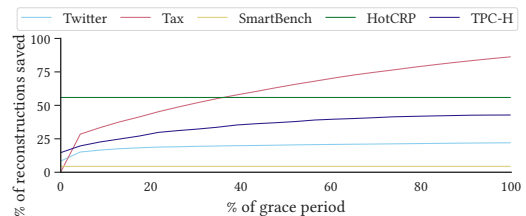
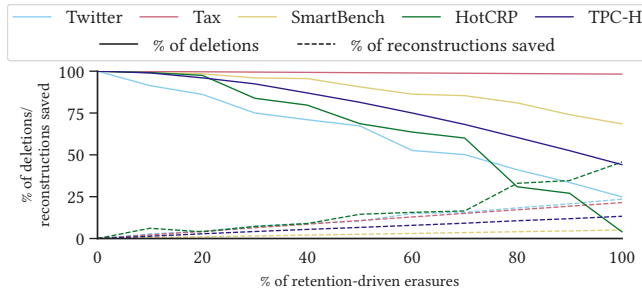


Figure 6: Number of saved reconstructions vs. grace period

**Experiment 5.** In this experiment, we combine both demand-driven and retention-driven erasures. To investigate the effect of different shares of retention-driven erasures, we employ a similar setting to Experiment 4. We vary the fraction of profiles that are erased using SCHEDULER (based on retention-time) vs. demand-driven on the fly. The demand-driven erasures are simulated by generating a random deletion time between the experiment start and the expiration time. For the entire experiment, we allow a grace period ( $\Gamma$ ) of one hour for Twitter and Tax, one minute for SmartBench, one day for HotCRP, and one week for TPC-H. In both erasure methods, we adopt a time-based batching.



**Figure 7: P2E2 overheads for fractions of demand- and retention-driven erasures**

We present the average number of deleted cells and the necessary reconstructions in Fig. 7. The figure shows that the number of deleted cells and the number of reconstructions reduce as the share of the retention-driven erasures increases. The improvement largely depends on the update characteristic, as explained for Experiment 4. Thus, the HotCRP dataset profits the most, while we cannot reduce the number of deleted cells for the Tax dataset. The variation in the amount of change depends on the number of base data cells per derived data cells. Some derived data cells that aggregate more data are more impactful, depending on the method used to delete them. When comparing 0% to 100% retention-driven erasures, we can save between 25% (SmartBench) and 90% (HotCRP) of deletions, and between 50% (HotCRP) and 20% (Twitter) of reconstructions.

## 7 RELATED WORK

**Deletion Semantics.** Regulation-driven deletions have been studied in several contexts: for timely deletion in LSM-tree [48], for cascading deletes in relationship graphs [18], for schema and annotations containing personal data [2], for deletion in blockchains [37], and in SQL context [46] that explores extensions to specify when data should be deleted. None of these consider the role of data inference in deletion. The need for formal specification of deletion semantics for regulatory compliance has been discussed in [4, 15, 45]. **Dependency rules.** RDRs draw inspiration from several lines of work in databases that explored dependency frameworks. These include: specification and reasoning frameworks for functional dependencies & denial constraints [44] in data cleaning [44] and consistent query answering [3, 5, 17, 24, 26], correlation constraint and causal constraints in causal databases, delta rules for generalized reasoning for a large class of dependencies (e.g., denial constraints, correlation, and causal constraints) for deletion-based data

repair [29], and provenance/lineage dependencies using semiring structures and annotations [16, 30, 31, 55] when available. RDRs extend these to include aggregate dependencies (like summation and max), and cell level dependencies which are essential to define fine grained data deletion. While reasoning frameworks based on above dependency specifications can be exploited to identify minimal deletion sets to make a database consistent to given constraints, they do not provide mechanisms to reason with relative changes in inferences across different database states (e.g., insertion and deletion states of a data). Discovering dependencies from data have also been explored in [1, 7, 8, 34, 42] which RDRs can express.

**Deletion in ML.** Recent regulatory efforts, such as the AI Act [23], raise deletion requirements in machine learning contexts. Naively deleting training data implies model retraining, which is often impractical [14]. Approaches such as SISA [10] propose partial retraining (where our retention-driven scheduling could be applied), while others apply differential privacy [60] to bound information leakage. Extending RDRs and P2E2 to cover dependencies between data and learned models remains an open and promising direction.

## 8 CONCLUSIONS AND FUTURE WORK

We formalize safe data erasure as Pre-insertion Post-Erasure Equivalence (P2E2), which resolves semantic ambiguity by providing strong guarantees on deleted data—filling a key gap in current systems [4]. We implement P2E2 using Relational Dependency Rules (RDRs), developing both exact and approximate algorithms.

While P2E2 provides a principled foundation for compliant data deletion, its adoption faces two key challenges. First, data dependencies may necessitate deleting more than what is requested to be erased. However, our evaluation across five domains—including highly dependent datasets—shows that our approaches effectively reduces this overhead. Future extensions may incorporate selectively applying P2E2 to data subsets, weighting dependencies, domain-specific inference, or relaxed variants of P2E2 enabling flexible trade-offs between deletion cost and retention obligations

Second, P2E2 relies on a specified set of dependencies through which deleted data can be inferred. While RDRs provide a structured way to encode known dependencies, discovering them may incur an overhead. In practice, such dependencies often exist elsewhere in the pipeline—e.g., business logic, analytics, consistency checks, or data cleaning—or can often be learned from data [1, 8, 42]. This requirement is consistent with the “reasonableness” standard in data regulations, which calls for reasonable, cost-effective conformative measures given current technology and implementation constraints. An interesting future direction is to extend P2E2 to protect against not only specified dependencies but any potential inferences, perhaps offering weaker, probabilistic guarantees.

Beyond refining P2E2 extensions and evaluating its applicability across domains, a key direction for future work is to extend deletion guarantees to data processing pipelines, where dependencies span multiple system components.

## ACKNOWLEDGMENTS

Chakraborty was supported by a fellowship from HPI@UCI. This work was supported by NSF Grants No. 2032525, 1545071, 1527536, 1952247, 2008993, 2133391, 2245372, and 2245373.

## REFERENCES

- [1] Archita Agarwal, Marilyn George, Aaron Jeyaraj, and Malte Schwarzkopf. 2022. Retrofitting GDPR Compliance onto Legacy Databases. *PVLDB* 15, 4 (2022), 958–970.
- [2] Kinan Dak Albab, Ishan Sharma, Justus Adam, Benjamin Kilimnik, Aaron Jeyaraj, Raj Paul, Artem Agvanyan, Leonhard Spiegelberg, and Malte Schwarzkopf. 2023. K9db: Privacy-Compliant Storage for Web Applications by Construction. In *Proceedings of the 17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23)*. USENIX Association, Boston, MA, USA, 99–116.
- [3] Marcelo Arenas, Leopoldo Bertossi, and Jan Chomicki. 1999. Consistent Query Answers in Inconsistent Databases. In *Proceedings of the Symposium on Principles of Database Systems (PODS)*. ACM, New York, NY, 68–79. <https://doi.org/10.1145/303976.303983>
- [4] Manos Athanassoulis, Subhadeep Sarkar, Zichen Zhu, and Dimitris Staratzis. 2022. Building Deletion-Compliant Data Systems. *IEEE Data Engineering Bulletin* 45, 1 (2022), 21–36. <http://sites.computer.org/debull/A22mar/p21.pdf>
- [5] Leopoldo Bertossi. 2006. Consistent Query Answering in Databases. *ACM SIGMOD Record* 35, 2 (2006), 68–76.
- [6] Leopoldo Bertossi. 2011. *Database Repairing and Consistent Query Answering*. Morgan & Claypool Publishers.
- [7] Tobias Bleifuß, Sebastian Kruse, and Felix Naumann. 2017. Efficient denial constraint discovery with Hydra. *PVLDB* 11, 3 (2017), 311–323.
- [8] Tobias Bleifuß, Thorsten Papenbrock, Thomas Bläsius, Martin Schirneck, and Felix Naumann. 2024. Discovering Functional Dependencies through Hitting Set Enumeration. *Proceedings of the International Conference on Management of Data (SIGMOD)* 2, 1 (2024), 1–24.
- [9] Philip Bohannon, Wenfei Fan, Floris Geerts, Xibei Jia, and Anastasios Kementsietsidis. 2007. Conditional Functional Dependencies for Data Cleaning. In *Proceedings of the IEEE International Conference on Data Engineering (ICDE)*. IEEE, 746–755.
- [10] Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. 2021. Machine Unlearning. In *Proceedings of the 2021 IEEE Symposium on Security and Privacy (S&P)*. IEEE, 141–159.
- [11] Brazil. 2018. Lei Geral de Proteção de Dados (LGPD) - Article 18(IV). <https://www.gov.br/cidadania/pt-br/acao-a-informacao/legpd> Brazilian Federal Law No. 13,709/2018. Last accessed on 2025-01-10.
- [12] California. 2018. Title 1.81.5. California Consumer Privacy Act of 2018 [1798.100 - 1798.199.100]. California Legislative Information. [https://leginfo.ca.gov/faces/codes\\_displayText.xhtml?division=3.&part=4.&lawCode=CIV&title=1.81.5](https://leginfo.ca.gov/faces/codes_displayText.xhtml?division=3.&part=4.&lawCode=CIV&title=1.81.5) California Civil Code. Last accessed on 2025-01-10.
- [13] Canada. 2000. Personal Information Protection and Electronic Documents Act (S.C. 2000, c. 5). <https://laws-lois.justice.gc.ca/ENG/ACTS/P-8.6/index.html> Justice Laws Website. Last accessed on 2025-01-10.
- [14] Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. 2022. Membership Inference Attacks from First Principles. In *Proceedings of the 2022 IEEE Symposium on Security and Privacy (S&P)*. IEEE, 1897–1914.
- [15] Vishal Chakraborty, Stacy Ann-Elvy, Sharad Mehrotra, Faisal Nawab, Mohammad Sadoghi, Shantanu Sharma, Nalini Venkatsubramanian, and Farhan Saeed. 2024. Data-CASE: Grounding Data Regulations for Compliant Data Processing Systems. In *Proceedings of the International Conference on Extending Database Technology (EDBT)*. OpenProceedings, Paestum, Italy, 108–115. [https://openproceedings.org/2024/conf/edbt/Camera\\_EDBT\\_Data\\_CASE\\_4.pdf](https://openproceedings.org/2024/conf/edbt/Camera_EDBT_Data_CASE_4.pdf)
- [16] James Cheney, Laura Chiticariu, and Wang-Chiew Tan. 2009. Provenance in Databases: Why, How, and Where. *Foundations and Trends in Databases* 1, 4 (April 2009), 379–474. <https://doi.org/10.1561/1900000006>
- [17] Jan Chomicki and Jerzy Marcinkowski. 2005. Minimal-Change Integrity Maintenance Using Tuple Deletions. *Information and Computation* 197, 1-2 (2005), 90–121.
- [18] Katriel Cohn-Gordon, Georgios Damaskinos, Divino Neto, Joshi Cordova, Benoît Reitz, Benjamin Strahs, Daniel Obenshain, Paul Pearce, and Ioannis Papagiannis. 2020. DELF: Safeguarding Deletion Correctness in Online Social Networks. In *Proceedings of the 29th USENIX Conference on Security Symposium (SEC)*. USENIX, USA, Article 60, 18 pages.
- [19] Michele Dallachiesa, Amr Ebaid, Ahmed Eldawy, Ahmed Elmagarmid, Ihab F. Ilyas, Mourad Ouzzani, and Nan Tang. 2013. NADEEF: A Commodity Data Cleaning System. In *Proceedings of the International Conference on Management of Data (SIGMOD)*. ACM, New York, NY, USA, 541–552.
- [20] Ghana Data Protection Commission. 2012. Data Protection Act. <https://nca.org.gh/wp-content/uploads/2020/09/Data-Protection-Act-2012.pdf> Accessed: 2025-04-02.
- [21] enryu43. 2023. Twitter 100 Million Tweets Dataset. [https://huggingface.co/datasets/enryu43/twitter100m\\_tweets](https://huggingface.co/datasets/enryu43/twitter100m_tweets) Accessed: 2025-01-25.
- [22] European Parliament and Council of the European Union. 2016. Regulation (EU) 2016/679 (General Data Protection Regulation). <https://eur-lex.europa.eu/eli/reg/2016/679/oj> Official Journal of the European Union. Last accessed on 2025-01-10.
- [23] European Parliament and Council of the European Union. 2024. Regulation (EU) 2024/1689: Artificial Intelligence Act. <https://eur-lex.europa.eu/eli/reg/2024/1689/oj/eng> Official Journal of the European Union. Last accessed on 2025-01-10.
- [24] Ronald Fagin, Benny Kimelfeld, Frederick Reiss, and Stijn Vansummeren. 2016. Declarative Cleaning of Inconsistencies in Information Extraction. *ACM Transactions on Database Systems (TODS)* 41, 1 (2016), 1–44.
- [25] Wenfei Fan, Floris Geerts, Xibei Jia, and Anastasios Kementsietsidis. 2008. Conditional functional dependencies for capturing data inconsistencies. *ACM Transactions on Database Systems (TODS)* 33, 2 (2008), 1–48.
- [26] Helena Galhardas, Daniela Florescu, Dennis Shasha, Eric Simon, and Cristian-Augustin Saita. 2001. Declarative Data Cleaning: Language, Model, and Algorithms. In *Proceedings of the International Conference on Very Large Databases (VLDB) (2013)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 371–380.
- [27] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. 2009. *Database Systems: The Complete Book (2nd ed.)*. Pearson Education.
- [28] Dan Geiger, Azaria Paz, and Judea Pearl. 1991. Axioms and algorithms for inferences involving probabilistic independence. *Information and Computation* 91, 1 (1991), 128–141.
- [29] Amir Gilad, Daniel Deutch, and Sudeepa Roy. 2020. On Multiple Semantics for Declarative Database Repairs. In *Proceedings of the International Conference on Management of Data (SIGMOD)*. ACM, New York, NY, USA, 817–831.
- [30] Todd J. Green, Grigoris Karvounarakis, Zachary G. Ives, and Val Tannen. 2007. Provenance Semirings. In *Proceedings of the Twenty-Sixth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS '07)* (Beijing, China). ACM, New York, NY, USA, 31–40.
- [31] Todd J. Green, Grigoris Karvounarakis, Zachary G. Ives, and Val Tannen. 2007. Update Exchange with Mappings and Provenance. In *Proceedings of the International Conference on Very Large Databases (VLDB)*. ACM, New York, NY, USA, 675–686.
- [32] Peeyush Gupta, Michael J Carey, Sharad Mehrotra, and oberto Yus. 2020. SmartBench: a benchmark for data management in smart spaces. *Proceedings of the VLDB Endowment* 13, 12 (2020), 1807–1820.
- [33] Ihab F. Ilyas, Volker Markl, Peter Haas, Paul Brown, and Ashraf Aboulnaga. 2004. CORDS: automatic discovery of correlations and soft functional dependencies. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD)*. ACM, New York, NY, USA, 647–658.
- [34] Youri Kaminsky, Eduardo HM Pena, and Felix Naumann. 2023. Discovering similarity inclusion dependencies. *Proceedings of the International Conference on Management of Data (SIGMOD)* 1, 1 (2023), 1–24.
- [35] Subhash Khot and Oded Regev. 2008. Vertex cover might be hard to approximate to within 2- $\epsilon$ . *J. Comput. System Sci.* 74, 3 (2008), 335–349.
- [36] Eddie Kohler. 2024. HotCRP: Conference Review System. <https://github.com/kohler/hotcrp> Accessed: 2025-01-25.
- [37] Michael Kuperberg. 2020. Towards Enabling Deletion in Append-Only Blockchains to Support Data Growth Management and GDPR Compliance. In *Proceedings of the IEEE International Conference on Blockchain (Blockchain)*. IEEE, Piscataway, NJ, USA, 393–400.
- [38] Andrei Lopatenko and Leopoldo Bertossi. 2007. Complexity of Consistent Query Answering in Databases under Cardinality-Based and Incremental Repair Semantics. In *Proceedings of the International Conference on Database Theory (ICDT)*. Springer, Berlin, Heidelberg, 179–193.
- [39] David Maier, Rachel Pottinger, AnHai Doan, Wang-Chiew Tan, Abdussalam Alawini, Hung Q Ngo, Babak Salimi, Harsh Parikh, Moe Kayali, Lise Getoor, Sudeepa Roy, and Dan Suciu. 2020. Causal Relational Learning. In *Proceedings of the International Conference on Management of Data (SIGMOD)*. ACM, New York, NY, USA, 241–256. <https://doi.org/10.1145/3318464.3389759>
- [40] Meta. 2017. Permanently Delete Your Facebook Account. <https://www.facebook.com/help/224562897555674>. Accessed: 2025-01-14.
- [41] MySQL. 2019. MySQL Triggers. <https://dev.mysql.com/doc/refman/9.0/en/trigger-syntax.html>. Accessed: 2025-01-10.
- [42] Eduardo HM Pena, Fabio Porto, and Felix Naumann. 2022. Fast Algorithms for Denial Constraint Discovery. *PVLDB* 16, 4 (2022), 684–696.
- [43] Eduardo H. M. Pena, Eduardo C. de Almeida, and Felix Naumann. 2019. Discovery of Approximate (and Exact) Denial Constraints. *PVLDB* 13, 3 (2019), 266–278.
- [44] Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. 2017. HoloClean: holistic data repairs with probabilistic inference. *PVLDB* 10, 11 (Aug. 2017), 1190–1201. <https://doi.org/10.14778/3137628.3137631>
- [45] Eduard Rupp, Emmanuel Symroutidis, and Jens Grossklags. 2022. Leave no data behind—empirical insights into data erasure from online services. *Proceedings on Privacy Enhancing Technologies* 3 (2022), 437–455.
- [46] Subhadeep Sarkar and Manos Athanassoulis. 2022. Query Language Support for Timely Data Deletion. In *Proceedings of the International Conference on Extending Database Technology (EDBT)*. OpenProceedings, Online, 418–429. <https://doi.org/10.48786/edbt.2022.35>

- [47] Subhadeep Sarkar, Jean-Pierre Banâtre, Louis Rilling, and Christine Morin. 2018. Towards Enforcement of the EU GDPR: Enabling Data Erasure. In *Proceedings of the 11th IEEE International Conference on Internet of Things (iThings 2018)*. IEEE, Halifax, Canada, 1–8.
- [48] Subhadeep Sarkar, Tarikul Islam Papon, Dimitris Staratzis, and Manos Athanassoulis. 2020. Lethe: A Tunable Delete-Aware LSM Engine. In *Proceedings of the International Conference on Management of Data (SIGMOD)*. ACM, New York, NY, USA, 893–908. <https://doi.org/10.1145/3318464.3389757>
- [49] Subhadeep Sarkar, Dimitris Staratzis, Ziehen Zhu, and Manos Athanassoulis. 2021. Constructing and analyzing the LSM compaction design space. *PVLDB* 14, 11 (jul 2021), 2216–2229. <https://doi.org/10.14778/3476249.3476274>
- [50] Supreeth Shastri, Vinay Banakar, Melissa Wasserman, Arun Kumar, and Vijay Chidambaram. 2020. Understanding and Benchmarking the Impact of GDPR on Database Systems. *Proceedings of the VLDB Endowment (PVLDB)* 13, 7 (mar 2020), 1064–1077.
- [51] David W Shipman. 1981. The functional data model and the data languages DAPLEX. *ACM Transactions on Database Systems (TODS)* 6, 1 (1981), 140–173.
- [52] Singapore. 2012. Personal Data Protection Act 2012 (PDPA) - Part IV: Retention Limitation Obligation. <https://sso.agc.gov.sg/Act/PDPA2012> Singapore Statutes Online. Last accessed on 2025-01-10.
- [53] Slawomir Staworko. 2007. *Declarative Inconsistency Handling in Relational and Semi-Structured Databases*. PhD Thesis. State University of New York at Buffalo. <https://cse.buffalo.edu/tech-reports/2007-11.pdf>
- [54] Transaction Processing Performance Council. 2021. *TPC-H Benchmark Specification, Version 2.17.3*. Technical Report. Transaction Processing Performance Council. <http://www.tpc.org/tpch/>
- [55] Benjamin E. Ujcich, Adam Bates, and William H. Sanders. 2018. A Provenance Model for the European Union General Data Protection Regulation. In *Proceedings of the International Provenance and Annotation Workshop (IPAW)*. Springer, 45–57.
- [56] United States. 2003. HIPAA Security Rule - 45 CFR §164.310(d)(2)(i). <https://www.ecfr.gov/current/title-45/subtitle-A/subchapter-C/part-164/subpart-C/section-164.310> U.S. Code of Federal Regulations. Last accessed on 2025-01-10.
- [57] Virginia. 2021. SB 1392: Consumer Data Protection Act (Virginia). <https://lis.virginia.gov/cgi-bin/legp604.exe?211+sum+SB1392> Virginia General Assembly. Last accessed on 2025-01-10.
- [58] Yang Wang, Pedro Giovanni Leon, Alessandro Acquisti, Lorrie Faith Cranor, Alain Forget, Norman Sadeh, Norman Sadeh, Lorrie Faith Cranor, Alain Forget, and Norman Sadeh. 2014. A Field Trial of Privacy Nudges for Facebook. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI)* (Toronto, Ontario, Canada) (CHI). ACM, New York, NY, USA, 2367–2376.
- [59] WhatsApp. 2020. About Disappearing Messages. [https://faq.whatsapp.com/673193694148537/?helpref=uf\\_share](https://faq.whatsapp.com/673193694148537/?helpref=uf_share) Last accessed on 2025-01-10.
- [60] Lefeng Zhang, Tianqing Zhu, Haibin Zhang, Ping Xiong, and Wanlei Zhou. 2023. FedRecovery: Differentially Private Machine Unlearning for Federated Learning Frameworks. *IEEE Transactions on Information Forensics and Security* 18 (2023), 4732–4746.