

A Language for the Composition of Privacy-Enforcement Techniques

Ronan-Alexandre Cherrueau, Rémi Douence and Mario Südholt

Ascola team, Mines Nantes, Inria, LINA

École des Mines de Nantes, Nantes, France

Email: `firstname.lastname@mines-nantes.fr`

Abstract—Today’s large-scale computations, *e.g.*, in the cloud, are subject to a multitude of risks concerning the divulging and ownership of private data. Privacy risks are mainly addressed using encryption-based techniques. These make data private, but costly to operate. Furthermore, today’s computations have to ensure privacy properties in the context of complex software compositions; however, no general support for the declarative definition and implementation of privacy-preserving applications has been put forward.

This article presents an approach to the correct composition of privacy-preserving applications in the cloud. Our approach provides language support for the composition of encryption- and fragmentation-based privacy-preserving algorithms. This language comes with a set of laws that allows us to verify privacy properties. Finally, we introduce implementation support in Scala that ensures privacy properties by construction using advanced features of Scala’s type system.

Keywords—*Language, Fragmentation, Encryption, Typing, Laws*

I. INTRODUCTION

The generalization of large-scale service-based computations executed over mutualized resources, notably in the context of cloud computing, has considerably increased the risk of losing control or even ownership of one’s personal data. In particular, the confidentiality and integrity of private data are at risk. Currently, privacy-preserving computations use frequently encryption techniques in order to preserve such properties of private data. However, these techniques have important drawbacks. They are costly to apply. They result in large amount of data being kept in one place (that are a target for attacks). And they do not allow to flexibly handle subsets of encrypted data.

In order to improve on these characteristics, alternative approaches have been explored. Data fragmentation [1], [2], in particular, consists in dividing data sets into parts and store them in different places. No unauthorized party can thus extract sensitive information from the pieces. Fragmentation allows to eliminate (most of) the computational overhead incurred by en/decryption. It allows for the distribution of data on a multitude of sites and supports the handling of subsets of data sets. Because fragmentation-based approaches frequently use encryption for parts of the data and computations, handling privacy properties results in complex compositions of privacy techniques. However, no comprehensive composition approach

for the construction of privacy-preserving computations has been put forward until now.

Figure 1 shows the need for such a comprehensive composition approach for the construction of privacy-preserving computations. Figure 1a illustrates a privacy-preserving query (of the number, per day, of meetings Alice had in her office last week) on a local application. The query is easy to formulate because Alice’s data is stored locally and is not subject to privacy problems. However, the same computation in the context of cloud computing requires to encrypt the database. Then, Alice has to decrypt the database on her computer to perform the query and encrypt the result once again, which is obviously not efficient. Illustration 1b shows how the composition of fragmentation with encryption and client-side computation can improve the query efficiency while preserving privacy. In Particular, the composition of these three techniques makes the query operates without any decryption overhead. However, the formulation of queries is far more difficult and error-prone in this case.

Currently, compositions of privacy-preserving strategies are programmed using traditional programming means (languages or frameworks) that do not provide any correctness guarantees. This article makes the following contributions, in this context:

- A motivation for the need of dedicated means for the correct composition of different privacy-preserving strategies, in particular, strategies based on encryption, fragmentation and client side computation. (Sec. II)
- A language for the declarative definition of a range of (composed) strategies for the enforcement of privacy-centric properties. The language is a SQL-like query language extended with abstractions for fragmentation and encryption. The language comes equipped with a set of laws that ensure the correct composition of privacy strategies and enables the transformation of a privacy-preserving query over a local application to a privacy-preserving query over a cloud application. (Sec. III)
- An implementation of the language in Scala that harnesses advanced typing properties in order to enforce the correct composition of privacy strategies. Specifically, the implementation prevents the compilation if the composition is not correct. (Sec. IV)

In addition, the paper discusses related work (Sec. V) and provides a conclusion and future works (Sec. VI).

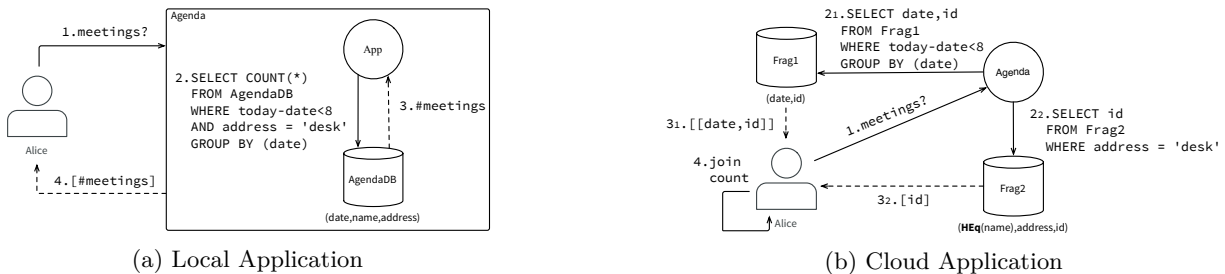


Fig. 1: Privacy-Preserving Agenda as Local (a) and cloud (b) Application

II. MOTIVATION

Consider an application developed for the cloud in order to improve, for instance, its availability, data replication properties or integration with other applications. Typically, the basic workflow consists in first the outsourcing of user data to a database service, and then applying applications hosted in the cloud to that data. For the case of the agenda application, one has first to subscribe to a cloud database service and outsource meeting data, and then install the agenda application on a cloud platform.

The privacy concern. In the age of cloud computing, cloud application programmers outsource without considering the personal data of individuals. Outsourcing reveals personal data to cloud providers that can share information with third parties. This is the conundrum of cloud computing: how to share data and keep personal data private.

To handle the privacy concern, two types of approaches have been introduced. The first approach focuses on supervised cloud. It includes systems relying on access control and policy enforcement [3], [4] to support accountability for violations [5]. The second approach focuses on an unsupervised cloud. It is a set of techniques to keep data private, such as encryption [6], fragmentation [2], differential privacy [7] and client-side computation [8]. Each technique has good privacy properties but, in general, do not offer as many guarantees as policy enforcement. For instance, data encryption fits for storage protection but not for computations and differential privacy is useful for statistical databases only. The main advantage is that they are applicable to real world problems.

This article handles the privacy concern in unsupervised clouds. We show that even if each technique taken by itself is limited, the composition of such techniques is much more effective and enables the development of expressive privacy-preserving cloud applications. In the remainder of this article, we will focus on the composition of three classes of techniques: encryption, fragmentation and client-side computations.

A. Encryption, Fragmentation and Client-Side Computations

Encryption. Encryption [6] is the process of encoding information before it is outsourced in such a way that only authorized parties can read it.

Historically, encryption is the first approach to the protection of private databases in the cloud. Using, for

instance, a symmetric encryption algorithm, a client encrypts its data before outsourcing it in the cloud. All queries can then simply be executed by first returning the necessary data from the database to the client in its encrypted state, be decrypted by her, and execute the query on her side. This approach is, however, far too expensive to be practical, since Alice has to bring back a large amount of data. One solution to this problem are recent encryption schemes, so called homomorphic ones, that execute query directly on encrypted data.

Theoretically, fully homomorphic schemes enable arbitrary operations to be performed on encrypted data. However, these schemes are prohibitively expensive to compute. Reasonably efficient homomorphic schemes are currently known only for a small set of operations. A deterministic encryption scheme [9], *e.g.*, allows to efficiently check the equality of values by comparing encrypted data. For this reason, query execution over encrypted data is often seen as practical only if corresponding efficient homomorphic encryption schemes are available [10].

Fragmentation. (Vertical) fragmentation [1], [2] is the process of separating information into non-linkable fragments in such a way that only authorized parties can recompose the original information. Fragmentation is applicable when associations of data are sensitive rather than individual data items themselves.

Fragmentation is tightly coupled to the notion of privacy constraints. A privacy constraint specifies which data are sensitive and should, therefore, be kept confidential. In our agenda application, for instance, meetings should satisfy two privacy constraints. An agenda meeting is the triplet $(date, name, address)$ that represents the meeting date, the name of the contact and the meeting location. The first constraint is $\{date, address\}$ so that an attacker cannot locate Alice by associating an address to a meeting date. The second constraint should be $\{name\}$ so that an attacker cannot infer the name of Alice's contacts. Here, fragmentation aims to make privacy associations such as $\{date, address\}$ safe by splitting the triplet $(date, name, address)$ in two.

Client-side computation. Client-side computation [8] designates the concept of letting clients perform sensitive computations on their computer and upload only the results. Especially, client-side computation stores private information on the client side. Thus, computations on private data are performed on the client side. Services that require strong guarantees on the truthfulness of the

result, such as billing service, can use, *e.g.*, zero-knowledge protocols [8] to ensure the integrity of the query results and privacy.

B. Composition of Privacy-Enforcement Techniques

The application. Figure 1b shows a common case in which a composition of techniques is required to obtain efficient and private application in the cloud. The application is the cloud version of the local one (Fig. 1a) and requires a composition of encryption, fragmentation and client-side computation.

To distribute this agenda application in a privacy-preserving manner, the programmer adopts the two privacy constraints $\{date, address\}$ and $\{name\}$ introduced in Sec. II-A. A programmer then has to choose a configuration for the application that satisfies the privacy constraints. First, she may, for instance, fragment the database in two. The first fragment contains the dates. The second contains the names and addresses. Now, the $\{date, address\}$ constraint is safe unless the two fragments are joined. Then, the programmer encrypts names with an homomorphic encryption that supports equality, thus ensuring that the $\{name\}$ constraint is satisfied.

Based on such a configuration, the query that computes the number of meetings Alice had per day at her desk last week (1 in Fig. 1b) is distributed on both fragments. It applies selection and grouping operations on the first fragment (2_1) to obtain the dates and identifiers of the meetings of the past week. At the same time, it applies the selection on the second fragment (2_2) to get identifiers of the office meetings. The application then fully takes advantages of the cloud, whether for storage or querying. However, the agenda has to join results of both fragments to count the number of people. Because that operation is not privacy preserving, finally, the rest of the query is executed on Alice’s side (3),(4).

Another useful query built on top of cloud configuration is the number of people Alice has met last week (Fig. 2). The query uses a *left-first* strategy of fragmentation to be efficient. It applies the selection on the left fragment to get identifiers of meetings the past week (2), (3). It then applies projection and grouping operations to the right fragment, and the programmer uses identifiers obtained from the left fragment to reduce the number of components (4). The grouping operation requires to compare names, but this is fine since names are encrypted with an homomorphic scheme that supports equality testing. Moreover, the overhead due to encryption is minimal because the selection largely reduces the number of comparisons.

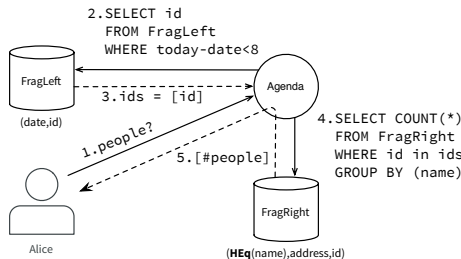


Fig. 2: Number of people Alice met last week

C. Language-Support to the Rescue

The above examples, as well as the large majority of real cloud applications, require several privacy techniques to be composed. In this case, the formulation of queries is far more difficult and errors easily slip in. For instance, the request in Fig. 2 is safe regarding privacy solely because of step (2) that drops the dates from the result. If step (2) does not drop dates, then the agenda application can join dates with addresses and violate the $\{date, address\}$ privacy constraint. Even without considering a specific request, composition makes it easy to write non-sense programs. For example, trying to encrypt twice the same column in a database does not make sense. Similarly, trying to sum values encrypted with a symmetric scheme is not reasonable.

To help programmers compose privacy techniques, we propose a functional language that focuses on privacy protection and query computation. As input, we have a privacy relation and apply privacy/query functions in turn to reduce the number of components in the relation. Together with the language, a set of algebraic laws specifies how to transform a local program to a privacy-preserving distributed cloud application.

We are convinced that this formalization constitutes an appropriate abstraction to help programmers reason on query and privacy techniques composition. The corresponding language-level and implementation support the declarative definition of privacy-preserving cloud applications and their efficient execution.

III. LANGUAGE-BASED COMPOSITION OF PRIVACY-ENFORCEMENT TECHNIQUES

This section formally describes our language for composition of encryption, fragmentation and client side computation techniques. First, we introduce its commands based on a SQL-like query language extended with abstractions for fragmentation and encryption. Second, we show how a query can be transformed by introducing privacy commands. Finally, we formalize and generalize our approach by providing composition laws.

A. Language Description

Our language is based on database queries obeying relational algebra properties.

A relation (*i.e.*, a table) is a set of tuples. For instance, an agenda stores a meeting as a triple $(date, name, address)$ representing the meeting date, the name of the contact and the meeting location. Our language offers four query functions:

- A selection σ filters tuples of a relation. For instance $\sigma_{(today-date < 8) \wedge (address = desk)}$ keeps meetings whose date is at most a week old and has desk as meeting venue.
- A projection π keeps a subset of the columns of a relation. For instance π_{date} keeps only the date of meetings.

- A grouping *group* definition groups together tuples of a relation. For instance, $group_{date}$ creates a group of tuples for each date.
- An aggregation *fold* computes a single value for groups. For instance, $count = fold (+) 0$ counts the number of meeting in a group.

Note that we omit a product operation of tables because it is not required by our examples but it could be easily introduced. These functions can be composed (\circ) in order to define complex queries according to the following grammar:

$$Q ::= Q \circ Q \mid \sigma \mid \pi \mid group \mid fold f k$$

Note that a sequence of function compositions is read from right to left. For instance, the number of meetings per day at desk last week is:

$$count \circ group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8) \wedge (address = desk)}$$

Our language also provides privacy-related pairs of functions:

- $crypt_{s,as}/decrypt_{s,as}$ encrypts/decrypts with a scheme s the components of tuples corresponding to π_{as} . For instance, $crypt_{heq,name}$ encrypts the contacts in the agenda, although they still can be compared because heq is an homomorphic encryption that supports equality.
- $frag_{\pi_{as}}/defrag_{\pi_{as}}$ vertically fragments/defragments a table of tuples into two tables of tuples, so that the tuples of the first table contain only the components corresponding to π_{as} and the tuples of the second table contain the remaining components (noted $\pi_{\bar{as}}$). For instance, $frag_{\pi_{date}}$ fragments the agenda table into a first table for the dates and a second table for the names and addresses. Note [1] that in both tables each tuple also contains an index in order to reconstruct the original tuples.
- $frag_{\sigma_p}/defrag_{\sigma_p}$ horizontally fragments/defragments a table of tuples into two tables. For instance, $frag_{\sigma_{today-date < 8}}$ fragments the agenda into two tables of meetings. The first table contains only the meetings of last week and the second table contains the older meetings.

These privacy functions can be composed to make queries privacy-aware as defined by the following grammar:

$$Q_p ::= Q_p \circ Q_p \mid Q \mid crypt \mid decrypt \mid frag \mid defrag$$

For instance, when the agenda is fragmented ($date$ on one host, $name$ and $address$ on another host) and the identity of contacts is encrypted, the query for the number of meetings per day at desk last week is:

$$\begin{aligned} & count \circ defrag_{\pi_{date}} \\ & \circ (group_{date} \circ \pi_{date} \circ \sigma_{today-date < 8}, \sigma_{address = desk}) \\ & \circ frag_{\pi_{date}} \circ crypt_{heq,name} \end{aligned}$$

It can be read as: encrypt the contacts' names, then fragment the table, select the meeting dates of last week on the first fragment and group them, select the meeting with desk as venue on the second fragment, defragment the results in order to get a group by day but only for

the desk venue, and finally count the number in groups. Note that there is no need to call the decrypt function here because $count$ does not require values, but $count$ must be executed on the client side because $defrag$ discards privacy protection. The next section shows how to transform the local query into the private one.

B. Making a Query Private with Transformations

The local query that computes the number of meeting per day at desk last week is:

$$Q_1 \equiv count \circ group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8) \wedge (address = desk)}$$

It can be decentralized by introducing fragmentation and encryption. In the following, the transformation delays the discarding of privacy protection (*i.e.*, “push” $defrag/decrypt$ to the left in a query); intuitively, this means that more computations are performed in the cloud rather than on the client side:

$$\begin{aligned} Q_1 & \equiv count \circ group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8) \wedge (address = desk)} \\ & \circ defrag_{\pi_{date}} \circ frag_{\pi_{date}} \end{aligned}$$

we add $frag$ then $defrag$ at the beginning (*i.e.*, right) of the query (this is correct since these functions are inverse)

$$\begin{aligned} & \equiv count \circ group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8) \wedge (address = desk)} \\ & \circ defrag_{\pi_{date}} \circ frag_{\pi_{date}} \\ & \circ decrypt_{heq,name} \circ crypt_{heq,name} \end{aligned}$$

similarly we introduce encryption

$$\begin{aligned} & \equiv count \circ group_{date} \circ \pi_{date} \circ defrag_{\pi_{date}} \\ & \circ (\sigma_{(today-date < 8)}, \sigma_{(address = desk)}) \\ & \circ frag_{\pi_{date}} \circ decrypt_{heq,name} \circ crypt_{heq,name} \end{aligned}$$

$decrypt$ and $defrag$ must be executed at the owner's place since they discard protection, so we delay them to the left. $defrag$ commutes with σ , but the selection must be applied only to the relevant fragment

$$\begin{aligned} & \equiv count \circ group_{date} \circ defrag_{\pi_{date}} \\ & \circ (\pi_{date} \circ \sigma_{(today-date < 8)}, \pi_{date} \circ \sigma_{(address = desk)}) \\ & \circ frag_{\pi_{date}} \circ decrypt_{heq,name} \circ crypt_{heq,name} \end{aligned}$$

$defrag$ and π commutes, but the projection must be applied to each fragment

$$\begin{aligned} & \equiv count \circ group_{date} \circ defrag_{\pi_{date}} \\ & \circ (\pi_{date} \circ \sigma_{(today-date < 8)}, \sigma_{(address = desk)}) \\ & \circ frag_{\pi_{date}} \circ decrypt_{heq,name} \circ crypt_{heq,name} \end{aligned}$$

the projection on components missing in the fragment can be simplified

$$\begin{aligned} & \equiv count \circ defrag_{\pi_{date}} \\ & \circ (group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8)}, \sigma_{(address = desk)}) \\ & \circ frag_{\pi_{date}} \circ decrypt_{heq,name} \circ crypt_{heq,name} \end{aligned}$$

$defrag$ and $group$ commutes, but grouping must be performed only in the relevant fragment

$$\begin{aligned} & \equiv count \circ defrag_{\pi_{date}} \\ & \circ (group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8)}, \\ & \quad \sigma_{(address = desk)} \circ decrypt_{heq,name}) \\ & \circ frag_{\pi_{date}} \circ crypt_{heq,name} \end{aligned}$$

$decrypt$ commutes with $frag$, must be applied to both fragments and it can be simplified when components are

Identity Laws:

$$id \equiv decrypt_{s,a,s} \circ crypt_{s,a,s} \quad (1)$$

$$id \equiv defrag_{\pi_{a,s}} \circ frag_{\pi_{a,s}} \quad (2)$$

$$id \equiv defrag_{\sigma_p} \circ frag_{\sigma_p} \quad (3)$$

Projection Laws:

$$\pi_a \circ decrypt_{s,a} \equiv decrypt_{s,a} \circ \pi_a \quad (4)$$

$$\pi_{a\bar{a}} \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (\pi_a, \pi_{\bar{a}}) \quad (5)$$

$$\pi_a \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (\pi_a, \pi_a) \quad (6)$$

$$\pi_{\bar{a}} \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (\pi_{\bar{a}}, \pi_{\bar{a}}) \quad (7)$$

$$\pi_a \circ defrag_{\sigma_p} \equiv defrag_{\sigma_p} \circ (\pi_a, \pi_a) \quad (8)$$

Grouping Laws:

$$group_a \circ decrypt_{s,b} \equiv decrypt_{s,b} \circ group_a \text{ if } a \notin \mathcal{P}(b) \quad (9)$$

$$group_a \circ decrypt_{s,b} \equiv decrypt_{s,b} \circ group_{s_a} \text{ if } a \in \mathcal{P}(b) \quad (10)$$

$$group_a \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (group_a, id) \quad (11)$$

$$group_{\bar{a}} \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (id, group_{\bar{a}}) \quad (12)$$

$$group_a \circ defrag_{\sigma_p} \equiv defrag_{\sigma_p} \circ (group_a, group_a) \quad (13)$$

Selection Laws:

$$\sigma_p \circ decrypt_{s,a} \equiv decrypt_{s,a} \circ \sigma_p \text{ if } dom(p) \notin \mathcal{P}(a) \quad (14)$$

$$\sigma_p \circ decrypt_{s,a} \equiv decrypt_{s,a} \circ \sigma_{s_p} \text{ if } dom(p) \in \mathcal{P}(a) \quad (15)$$

$$\sigma_{pa \wedge p\bar{a} \wedge pa\bar{a}} \circ defrag_{\pi_a} \equiv \sigma_{pa\bar{a}} \circ defrag_{\pi_a} \circ (\sigma_{pa}, \sigma_{p\bar{a}}) \quad (16)$$

$$\sigma_{pa \wedge t \wedge t} \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (\sigma_{pa}, id) \quad (17)$$

$$\sigma_{t \wedge p\bar{a} \wedge t} \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (id, \sigma_{p\bar{a}}) \quad (18)$$

$$\sigma_{p'} \circ defrag_{\sigma_p} \equiv defrag_{\sigma_p} \circ (\sigma_{p'}, \sigma_{p'}) \quad (19)$$

Aggregating Laws:

$$count \circ decrypt_{s,a,s} \equiv count \quad (20)$$

Protection Composition Laws:

$$f \circ id \equiv id \circ f \equiv f \quad (21)$$

$$(f1, f2) \circ (g1, g2) \equiv (f1 \circ f2, g1 \circ g2) \quad (22)$$

$$frag_{\pi_a} \circ decrypt_{s,a} \equiv (decrypt_{s,a}, id) \circ frag_{\pi_a} \quad (23)$$

$$frag_{\pi_a} \circ decrypt_{s,\bar{a}} \equiv (id, decrypt_{s,\bar{a}}) \circ frag_{\pi_a} \quad (24)$$

$$decrypt_{s,a} \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (decrypt_{s,a}, id) \quad (25)$$

$$decrypt_{s,\bar{a}} \circ defrag_{\pi_a} \equiv defrag_{\pi_a} \circ (id, decrypt_{s,\bar{a}}) \quad (26)$$

$$frag_{\sigma_p} \circ decrypt_{s,a} \equiv (decrypt_{s,a}, decrypt_{s,a}) \circ frag_{\sigma_p} \quad (27)$$

$$decrypt_{s,a} \circ defrag_{\sigma_p} \equiv defrag_{\sigma_p} \circ (decrypt_{s,a}, decrypt_{s,a}) \quad (28)$$

Fig. 3: Some Laws of the Algebra of Fragmentation and Encryption

missing in the fragment

$$\begin{aligned} &\equiv count \circ defrag_{\pi_{date}} \\ &\quad \circ (group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8)}, \\ &\quad \quad decrypt_{heq,name} \circ \sigma_{(address=desk)}) \\ &\quad \circ frag_{\pi_{date}} \circ crypt_{heq,name} \end{aligned}$$

$decrypt$ commutes with σ , the predicate of σ does not work on encrypted data, thus it does not require homomorphic encryption to compute selection

$$\begin{aligned} &\equiv count \circ decrypt_{heq,name} \circ defrag_{\pi_{date}} \\ &\quad \circ (group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8)}, \sigma_{(address=desk)}) \\ &\quad \circ frag_{\pi_{date}} \circ crypt_{heq,name} \end{aligned}$$

$decrypt$ commutes with $defrag$

$$\begin{aligned} &\equiv count \circ defrag_{\pi_{date}} \\ &\quad \circ (group_{date} \circ \pi_{date} \circ \sigma_{(today-date < 8)}, \sigma_{(address=desk)}) \\ &\quad \circ frag_{\pi_{date}} \circ crypt_{heq,name} \end{aligned}$$

finally $count$ computes the number of groups but it does not rely on the value of tuples, so they do not require to be decrypted ■

C. Laws for Composition

We now generalize our transformational approach to formally-defined composition laws following Backus's approach [11] and thus get an axiomatic semantics. Our language obeys the laws detailed in Figure 3. They specify how protection and query functions interact, in particular how they commute. We review them briefly.

Identity Laws specify that pairs of protection functions are inverse to each other. Applying a protection then discarding a protection results in no protection. When oriented from left to right, these three rules can be used to introduce privacy functions in a query.

All the other rules specify when privacy protection discarding and query functions commute. They can be

used to delay the discarding of protection and preform computations in the cloud rather than on the client side. The *Projection Laws*, (4) specifies that projection and $decrypt$ commute. When fragmentation is used, laws (5)-(8) specify a projection becomes a pair of projections (one per fragment).

Grouping Laws (9),(10) specify that $decrypt$ and $group$ commute. When groups are based on encrypted components, $group$ must take into account encryption (10). Vertical fragmentation and $group$ commute, when groups can be computed in a single fragment (11),(12). Horizontal fragmentation and $group$ commute and groups must be computed in both fragments (13).

Selection Laws are quite similar to *Grouping Laws*, but the different cases are based on predicate parts dealing with one fragment, the other or both. In particular, (16) requires to select tuples in both vertical fragments *after* the defragmentation (*i.e.*, at owner's place). In law (20), $count$ does not access values so that $decrypt$ can be discarded.

Finally, *Protection Composition Laws* commute and distribute functions in order to apply previous laws.

IV. IMPLEMENTATION

The functional language for the composition of privacy-aware query-based applications introduced above permits the definition of privacy-preserving cloud applications. This section presents our Scala-based [12], [13] framework¹, a prototype that makes it possible to program sophisticated privacy-aware cloud applications. It can be viewed as an operational semantics of our language. In particular, our implementation harnesses Scala's type system to prevent the compilation of a program if the composition of privacy techniques is not correct. This is achieved

¹The sources of our framework and examples are available on the Github platform at <https://github.com/rcherrueau/phand>

```

1 for {
2   _ <- configure[Date,Name,Addr]
3   q <- query (db => {
4     val r1 =  $\sigma$  (db) (lastweek  $\wedge$  atdesk)
5     val r2 =  $\pi$  (r1) (date)
6     val r3 = group (r2) (date)
7     val r4 = count (r3); r4
8   })
9 } yield q

```

(a) Local Application

```

1 for {
2   _ <- configure[Date,Name,Addr]
3   _ <- crypt (_2) (HEq(_))
4   _ <- fragV (_1) (Site1(_), Site2(_))
5   qL <- queryL (fragL => {
6     val r1 =  $\sigma$  (fragL) ( $\sigma_{lift}$  lastweek)
7     val r2 =  $\pi$  (r1) ( $\pi_{lift}$  date)
8     val r3 = group (r2) (date); r3
9   })
10  qR <- queryR (fragR => {
11    val r1 =  $\sigma$  (fragR) ( $\sigma_{lift}$  atdesk)
12    val r2 =  $\pi$  (r1) (id); r2
13  })
14 } yield count (gather (qL, qR))

```

(b) Cloud Application

Fig. 4: Privacy-Preserving Agenda as Local (a) and Cloud (b) Application

through the satisfaction of the laws introduced previously. We rely on property-based testing with ScalaCheck ² to argue their correctness.

A. From Theory to Practice

One of the main differences between the language and the actual implementation is that we distinguish functions that compute queries from functions that change the shape of the database. In the previous section, the language uses the pointwise application of one function to the result of another, which successively reduces the database until the result is obtained. This abstraction helps reasoning about the query design. But, for real world programming we do not want to reduce our database and lose components. In practice, two different levels of computation are more useful: one that uses database components to compute queries, and another one that modifies the database shape and applies privacy protection.

The programs in Fig. 4 query the number of meetings per day at the desk last week, for a local (a) and a cloud (b) application. The local application only uses functions to query the database (σ , π , **group** and **count**). In contrast, the cloud application also composes functions for privacy protection (**crypt** and **fragV**).

The Guardian Monad. In functional programming (FP), functions with side effects such as protection functions are performed in a *monad* [14], a pattern that makes it easy to chain function calls. In particular, FP deals with side effects using a *state monad* [15] that attaches state information to function calls. For this reason, our framework provides a state monad suitable for our purposes. Our state monad, called *guardian*, has a progressive state, and is defined based on the following requirements:

- The state is the database.
- Query functions only access the components of the database. The application of a query function returns some result that can be used as input for a second query, without modifying the database. For instance, the selection in figure 4a accesses the content of the database and returns the result in **r1** (line 4). Values in **r1** are then used for the projection (1.5) and so on, until

the **count** operation (1.7). At the end of the program (1.9), the **q** variable gets the result of the query. But, the database is unchanged from the input.

- Protection functions only modify the shape of a database. For instance, the **crypt** instruction in figure 4b modifies the database by encrypting the second column with a homomorphic scheme that supports equality testing (1.3). Similarly, the **fragV** instruction splits the database vertically on the first column (1.4). It distributes the left fragment on site number one and the right fragment on site number two. Henceforth, the querying should be done on both fragments (1.5-13).

The protection functions **crypt** and **fragV** are complemented by discarding functions **decrypt** and **defragV**. The discarding functions from the framework differ from the ones in the language because the framework ones discard protection on the database, whereas the functions of the language discard protections on the query result. In the language example of the cloud application (Sec. III-B), the call of *defrag π_{date}* joins at the client side the result of both fragments and hands it over to the *count* computation. In the corresponding implementation (Fig. 4b), calling a **defragV** will join fragments but not the result of the queries. For this reason, the framework provides the **gather** instruction, which is applied at the query level. It brings back data at the client side and discards protection. Given this, the program 4b gathers the result from both fragments in order to count the number of meetings (1.14). Note that in this example like in the example of the language, the **gather** does not have to decrypt names because the right query drops them.

Monads and types. In example 4b it does not make sense to encrypt the second column twice; encrypting once is enough. Likewise, it does not make sense to query the fragmented/encrypted database in the same manner as the local database; the new query request should be split and work on encrypted data. Furthermore, it does not make sense to group on encrypted data that do not support equality testing; grouping on encrypted data requires the use of an homomorphic scheme. To put it simply, composing privacy techniques can easily lead to runtime errors!

A main feature of monads is the use of types to exhibit what it means to execute chained function applications. Given this, the guardian monad exhibits useful information

²<http://scalacheck.org/>

```

1 for {
2   _ <- configure[Date,Name,Addr]
3   _ <- crypt (_2) (HEq(_))
4   // ill-typed, encryption of
5   // all-ready encrypted column:
6   _ <- crypt (_2) (HEq(_))
7   // ...
8 } yield ()

```

Fig. 5: Twice encryption does not type check

```

1 for {
2   _ <- configure[Date,Name,Addr]
3   _ <- crypt (_2) (HEq(_))
4   _ <- fragV (_1) (Site1(_), Site2(_))
5   // ill-typed, query on a non-local
6   // database.
7   q <- query (db => { /* ... */ })
8   // ...
9 } yield ()

```

Fig. 6: Fragmentation requires querying on fragments

at the type level to help the programmer write programs that are well compose. For instance, the type of the guardian at the end of program 4b is:

```

1 Guard[
2   Site0[DB[Raw[Date] |: Raw[Name] |: Raw[Addr]]],
3   (Site1[DB[Raw[Date] |: Id]],
4   Site2[DB[HEq[Name] |: Raw[Addr] |: Id]]),
5   Site0[List[Int]]]

```

with:

- 1.2 The shape of the database at the start of the computation. Here the guardian only accepts, as input, database that stores date, name and address components (*i.e.*, `DB[Date |: Name |: Addr]`) in plain form (*i.e.*, `Raw`). Components are uploaded at client side (*i.e.*, `Site0`).
- 1.3-4 The shape of the database at the end of the computation. Here the guardian transforms the database into two fragments. The first one stores dates in plain form at site one. The second one stores names and addresses at site two. The guardian also encrypts names with an encryption that supports equality (*i.e.*, `HEq`).
- 1.5 The type of the query result. Here, the type of the number of meetings per day at the desk last week (*i.e.*, `List[Int]`). The `Site0` annotation informs the programmer that a part of the query is computed at the client side.

The guardian monad gives useful type information and uses this type information to check, during the compilation of the program, that the composition of privacy technique goes well. In other cases, the program does not compile and does not produce an executable. Hence, trying to encrypt an already encrypted column does not compile (Fig. 5). Querying a fragment with a local approach does not compile (Fig. 6). Grouping/Filtering on encrypted data that do not support the equality test does not compile (Fig. 7). Generally speaking, the implementation satisfies the laws of figure 3. We rely on property-based testing with `ScalaCheck` to argue their correctness.

Finally, monad bindings enable the naming of query result such as in program 4b. Value `qL` contains the result of the left query. Value `qR` contains the result of the

```

1 for {
2   _ <- configure[Date,Name,Addr]
3   // Symmetric encryption of Name.
4   // Symmetric doesn't support equality
5   // testing.
6   _ <- crypt (_2) (Symmetric(_))
7   q <- query (db => {
8     // ill-typed, name doesn't
9     // support equality testing
10    group (db) (name)
11  })
12 // ...
13 } yield ()

```

Fig. 7: Grouping with encryption requires support equality

```

1 for {
2   _ <- configure[Date,Name,Addr]
3   _ <- crypt (_2) (HEq(_))
4   _ <- fragV (_1) (Site1(_), Site2(_))
5   // Queries on left fragment to get identifiers
6   // of meetings the past week:
7   ids <- queryL (fragL => {
8     val r1 =
9       σ (fragL) (σlift lastweek)
10    val r2 = π (r1) (id); r2
11  })
12   q <- queryR (fragR => {
13     // Reduces the number of
14     // elems with ids of left
15     val r1 = σ (fragR) {
16       case (_,_,id) =>
17         ids.exists(id)
18     }
19     val r2 = group (r1) (name)
20     val r3 = count (r2); r3
21  })
22 } yield q

```

Fig. 8: Number of people met last week – left-first strategy

right query. Naming is essential when the programmer wants to implement a profitable strategy like the *left-first* strategy seen in the motivation section (figure 2). Program 8 implements the left-first strategy. It simply consists of naming the result of the left fragment, and then use it in the right fragment.

B. Feedback on the Implementation with Scala

The Scala programming language is good for generalization. We harness an advanced use of Scala implicits and type members to perform type-level computations [16]. This enables the definition of arity-polymorphic databases, so that the guardian monad can be implemented once and for all. Without an arity-polymorphic database we would have to write as many guardian monads as a database could contain attributes, just like for tuples. In the previous examples (*e.g.*, Fig. 4b), the presence of integers prefixed by an underscore is the direct consequences of type-level computation. Integers with an underscore are Church encodings of the natural numbers at the type level. They make it possible to identify, at compile time, which column has to be encrypted, and the type of both fragments after fragmentation.

Scala unifies functional and object-oriented programming. Under the hood, data are objects and operations method calls. This object model is good for modularity and generalization, for instance, with subtyping. However,

it makes type inference less powerful than that of ML-like functional languages that use the Hindley-Milner algorithm. Because of that, the guardian monad sometimes requires the programmer to explicitly specifying the type to help the compiler infer type parameters. The code examples we presented above are a beautified version, omitting a few type annotations that are needed by the compiler, so that readers can more easily understand the intention of the guardian monad. The code with all necessary type annotation is available on the Github platform.

Finally, the guardian monad is a prototype which proves the validity of our approach. In the future we intend to bind the current implementation with popular Scala libraries such as Akka³ for the distribution and Slick³ for the mapping with real relational databases.

V. RELATED WORK

In the following, we compare our work to three sets of related work: approaches that focus on data fragmentation and encryption, related work providing language support for privacy properties and approaches, such as sticky policies and security-aware objects that enable privacy properties to be expressed and enforced directly.

Data fragmentation is a recent technique that strives to ensure strong confidentiality properties without the query overhead of encryption-based approaches. A recent overview [2] presents a wide range of fragmentation techniques and corresponding algorithms. However, none of the discussed approaches include, unlike ours, declarative means for the construction of fragmentation algorithms, their composition especially with encapsulation techniques and formal property guarantees over implementations.

Several language-based approaches have been proposed for other privacy properties. Tetali *et al.* [17], have proposed analysis techniques for compositions of different types of homomorphic encryption algorithms. Fournet *et al.* [8] define ZQL, a query language operating over annotated database schemas that use strong typing in order to ensure security properties of generated implementations in F# and C++. Reed and Pierce [7] propose a specialized type system to enforce privacy guarantees by means of differential privacy. However, none of these approaches provide language support for the composition of fragmentation and encapsulation techniques.

Another domain of related work consists in support for the expression of privacy properties in the form of policies and constraints over accesses to runtime objects. Sticky policies [18], represent a class of policies that enable the abstract definition of privacy properties and their enforcement through runtime annotations. Self-protecting software systems [19], such as self-defending objects [20], support the protection of privacy properties of runtime entities by strong encapsulation and access control of these entities. Again however, none of these approaches, as well as other policy-based and encapsulation-based privacy techniques, support properties involving the composition of fragmentation and encapsulation techniques.

VI. CONCLUSION

This paper addresses the problem of how to define and enforce privacy properties in the context of unsupervised cloud with different privacy-enforcing techniques. We have considered privacy properties that are formulated in terms of compositions of data fragmentation, encryption and client-side computation. We have provided programming language support for the definition of such composed privacy-enforcement strategies and an implementation on top of the Scala language that, using a specialized type system, ensures privacy properties by construction. We have also provided a set of laws that ensures privacy properties and that are satisfied by the language mechanisms we provide. As future work we intend to extend the set of fragmentation techniques and approaches to encryption in order to obtain a full-fledged composition theory for these two classes of privacy-enforcing techniques.

REFERENCES

- [1] G. Aggarwal, M. Bawa, P. Ganesan *et al.*, “Two can keep A secret: A distributed architecture for secure database services,” in *CIDR*, 2005.
- [2] S. D. C. di Vimercati, R. F. Erbacher, S. Foresti *et al.*, “Encryption and fragmentation for data confidentiality in the cloud,” in *FOSAD*, 2013.
- [3] S. Pearson and M. C. Mont, “Sticky policies: An approach for managing privacy across multiple parties,” *IEEE Computer*, vol. 44, no. 9, 2011.
- [4] Y.-Y. Chen, P. A. Jamkhedkar, and R. B. Lee, “A software-hardware architecture for self-protecting data,” in *CCS*, 2012.
- [5] D. J. Weitzner, H. Abelson, T. Berners-Lee *et al.*, “Information accountability,” *Communication of the ACM*, 2008.
- [6] D. Salomon, *Data privacy and security: encryption and information hiding*. Springer Science & Business Media, 2003.
- [7] J. Reed and B. C. Pierce, “Distance makes the types grow stronger: a calculus for differential privacy,” in *ICFP*, 2010.
- [8] C. Fournet, M. Kohlweiss *et al.*, “ZQL: A compiler for privacy-preserving data processing,” in *USENIX Security*, 2013.
- [9] B. Schneier, *Applied cryptography - protocols, algorithms, and source code in C (2. ed.)*. Wiley, 1996.
- [10] M. Naehrig, K. Lauter, and V. Vaikuntanathan, “Can homomorphic encryption be practical?” in *CCSW*, 2011.
- [11] J. W. Backus, “Can Programming Be Liberated From the von Neumann Style? A Functional Style and its Algebra of Programs,” *Commun. ACM*, vol. 21, no. 8, 1978.
- [12] M. Odersky, L. Spoon, and B. Venners, *Programming in Scala*. Artima Inc, 2008.
- [13] P. Chiusano and R. Bjarnason, *Functional Programming in Scala*. Manning Publications Co., 2014.
- [14] P. Wadler, “Comprehending monads,” *Mathematical Structures in Computer Science*, vol. 2, no. 4, 1992.
- [15] —, “The essence of functional programming,” in *Principles of Programming Languages*, 1992.
- [16] B. C. d. S. Oliveira, A. Moors, and M. Odersky, “Type classes as objects and implicits,” in *OOPSLA*, 2010.
- [17] S. D. Tetali, M. Lesani, R. Majumdar *et al.*, “MrCrypt: static analysis for secure cloud computations,” in *OOPSLA*, 2013.
- [18] G. Karjoth, M. Schunter, and M. Waidner, “Platform for enterprise privacy practices: Privacy-enabled management of customer data,” in *Privacy Enhancing Technologies*, 2002.
- [19] E. Yuan, N. Esfahani, and S. Malek, “A systematic survey of self-protecting software systems,” *TAAS*, vol. 8, no. 4, 2014.
- [20] J. W. Holford, W. J. Caelli, and A. W. Rhodes, “Using self-defending objects to develop security aware applications in java,” in *ACSC*, 2004.

³<http://akka.io/>; <http://slick.typesafe.com/>