

A formal study of collaborative access control in distributed datalog

Serge Abiteboul, Pierre Bourhis, Victor Vianu

► **To cite this version:**

Serge Abiteboul, Pierre Bourhis, Victor Vianu. A formal study of collaborative access control in distributed datalog. Wim Martens ; Thomas Zeume. ICDT 2016 - 19th International Conference on Database Theory , Mar 2016, Bordeaux, France. <<http://drops.dagstuhl.de/opus/portals/lipics/index.php?semnr=16002>>. <hal-01290497>

HAL Id: hal-01290497

<https://hal.inria.fr/hal-01290497>

Submitted on 18 Mar 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A formal study of collaborative access control in distributed datalog

Serge Abiteboul¹, Pierre Bourhis², and Victor Vianu³

- 1 INRIA-Saclay & ENS Cachan
fname.lname@inria.fr
- 2 CNRS, CRIStAL, UMR 9189
fname.lname@univ-lille1.fr
- 3 UCSD & INRIA-Saclay
lname@cs.ucsd.edu

Abstract

We formalize and study a declaratively specified collaborative access control mechanism for data dissemination in a distributed environment. Data dissemination is specified using distributed datalog. Access control is also defined by datalog-style rules, at the relation level for extensional relations, and at the tuple level for intensional ones, based on the derivation of tuples. The model also includes a mechanism for “declassifying” data, that allows circumventing overly restrictive access control. We consider the complexity of determining whether a peer is allowed to access a given fact, and address the problem of achieving the goal of disseminating certain information under some access control policy. We also investigate the problem of information leakage, which occurs when a peer is able to infer facts to which the peer is not allowed access by the policy. Finally, we consider access control extended to facts equipped with provenance information, motivated by the many applications where such information is required. We provide semantics for access control with provenance, and establish the complexity of determining whether a peer may access a given fact together with its provenance. This work is motivated by the access control of the Webdamlog system, whose core features it formalizes.

1 Introduction

The personal *data* and favorite *applications* of Web users are typically distributed across many heterogeneous devices and systems. In [20], a novel *collaborative access control mechanism* for a distributed setting is introduced in the context of the language Webdamlog, a datalog-style language designed for autonomous peers [3, 2]. The experimental results of [20] indicate that the proposed mechanism is practically feasible, and deserves in-depth investigation. In the present paper, we provide for the first time formal grounding for the mechanism of [20] and answer basic questions about the semantics, expressiveness, and computational cost of such a mechanism. In the formal development, we build upon *distributed datalog* [17, 21], which abstracts the core of Webdamlog, while ignoring certain features, such as updates and delegation.

In this investigation, as in Webdamlog, access control is *collaborative* in the following sense. The system provides the *means* to specify and infer access rights on disseminated information, thus *enabling* peers to collectively enforce access control. The system is agnostic as to how peers are motivated or coerced into conforming to the access control policy. This can be achieved in various ways, from economic incentives to legal means (see, e.g., [31]), possibly relying on techniques such as encryption or watermarking (see, e.g., [5]). We do not address these aspects here.

The access control of [20] that we formalize and study here works as follows. First, each peer specifies which other peers may access each of its extensional relations using *access-*



licensed under Creative Commons License CC-BY

Leibniz International Proceedings in Informatics

LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

control-list rules. This provides in a standard manner an initial coarse-grained (relation-at-a-time) access control, enforced locally by each peer. Next, facts can be derived among peers using *application rules*. Access control is extended to such facts based on their provenance: to see a propagated fact, a peer must have access to the extensional relations used by the various peers in producing the fact. This enables controlling access to data *disseminated* throughout the entire network, at a fine-grain (i.e., tuple) level. This capability is a main distinguishing feature of Webdamlog’s access control model. The access control also includes a *hide* mechanism that allows circumventing overly restrictive access control on some disseminated facts, thus achieving a controlled form of “declassification” for selected peers.

Access control in distributed datalog raises a variety of novel semantic, expressiveness and complexity issues. How complex is it to check whether a peer has the right to access a propagated fact? What are the appropriate complexity measures in this distributed setting? Does the access control mechanism prevent leakage of unauthorized information? What does it mean to extend access control to facts equipped with their *provenance*? Is there an additional cost? These are some of the fundamental questions we study, described in more detail next.

While the experimental results of [20] suggest that the computational cost of the proposed mechanism is modest, we show formally that its complexity is reasonable. Specifically, we prove that the data complexity of determining whether a peer can access a given fact is PTIME-complete (with and without *hide*).

We next consider the problem of information leakage, which occurs when a peer is able to infer some facts to which the peer is not allowed access by the policy. We show that, while undecidable in general, information leakage can be tested for certain restricted classes of policies and is guaranteed not to occur for more restricted classes.

One of the challenges of access control is the intrinsic tension between access restrictions and desired exchange of information. We consider the issue of achieving the goal of disseminating certain information under some access control policy. The goal is specified as a distributed datalog program. We show that it is undecidable whether a goal can be achieved without declassification (i.e., without *hide*). We study the issue of finding a policy without *hide* that achieves a maximum subset of the specified goal. While any goal can be achieved by extensive use of *hide*, we show, more interestingly, how this can be done with *minimal* declassification.

In many applications, it is important for inferred facts to come with *provenance* information, i.e., with traces of their derivation. We demonstrate that adding such a requirement has surprising negative effects on the complexity. For this, we introduce an intermediate measure between data and combined complexity, called *locally-bounded* combined complexity that allows making finer distinctions than the classical measures in our context. The intuition is that the peers are seen as part of the data and not of the schema, which is more in the spirit of a Web setting. We show that the locally bounded complexity of query answering increases from PTIME-complete to PSPACE-complete when it is required that the query answer carries *provenance* information.

The organization is as follows. Section 2 recalls the distributed datalog language [3]. In Section 3, we formalize the core aspects of the access control mechanism of [20], establish the complexity of answering queries under access control. Information leakage is studied in Section 4. The issue of achieving some dissemination goal under a particular access control policy is the topic of Section 5. Access control in the presence of provenance is investigated in Section 6. Finally, we discuss related work and conclude. Proofs are provided in an appendix.

2 Distributed Datalog

In this preliminary section, we formally define a variant of distributed datalog, which captures the core of Webdamlog [3].

The language. We assume infinite disjoint sets Ext of extensional relation symbols, Int of intensional relation symbols, \mathcal{P} of *peers* (e.g. p, q), \mathcal{D}_p of *pure data values* (e.g., a, b), and \mathcal{V} of *variables* (e.g., x, y, X, Y). For relations, we use symbols such as R, S, T . The set \mathcal{D} of *constants* is $\mathcal{P} \cup \mathcal{D}_p \cup Ext \cup Int$. A *schema* is a mapping σ whose domain $dom(\sigma)$ is a finite subset of \mathcal{P} , that associates to each p a finite set $\sigma(p)$ of relation symbols in $Int \cup Ext$, with associated arities. Let σ be a schema, $p \in dom(\sigma)$. A relation R in $\sigma(p)$ is denoted by $R@p$, and its arity by $arity(R@p)$. We denote $ext(p) = \sigma(p) \cap Ext$, $int(p) = \sigma(p) \cap Int$, $ext(\sigma) = \cup_{p \in dom(\sigma)} ext(p)$, and $int(\sigma) = \cup_{p \in dom(\sigma)} int(p)$. An *instance* I over σ is a mapping associating to each relation schema $R@p$ a finite relation over \mathcal{D} of the same arity. For a tuple \bar{a} in $I(R@p)$, the expression $R@p(\bar{a})$ is called a (p -)*fact* in $R@p$. An *extensional* instance is one that is empty on $int(\sigma)$. Observe that $R@p$ and $R@q$, for distinct p, q , are distinct relations with no a priori semantic connection, and possibly different arities. Note also that an expression $R@p(a_1, \dots, a_k)$ for R, p, a_1, \dots, a_k in \mathcal{D} is a *fact* for a schema σ if: p is a peer in $dom(\sigma)$, R is a relation schema in $\sigma(p)$, and $arity(R@p) = k$. Note that relations may contain pure data values, peers, as well as relation symbols. Finally, (U)CQ denotes (unions) of conjunctive queries (see [6]).

► **Definition 1** (distributed datalog). A *d-datalog* program P over schema σ is a finite set of rules of the form

- $Z_0@z(\bar{x}_0) :- R_1@p(\bar{x}_1), \dots, R_k@p(\bar{x}_k)$ where
- $p \in dom(\sigma)$, $k \geq 0$, and for every $i \geq 1$, R_i is in $\sigma(p)$ and \bar{x}_i is a vector of variables and constants in \mathcal{D} of the proper arity;
- $z \in dom(\sigma) \cup \mathcal{V}$, $Z_0 \in Int \cup \mathcal{V}$; and
- each variable occurring in the head appears in \bar{x}_i for some $i \geq 1$.

Note that the relation or peer names in the head may be variables. Note also that all the relations in the body of a rule come from the same peer. Although we define a global d-datalog program, one should think of each peer p as having its separate program consisting of all the rules whose bodies use relations at p .

► **Example 2.** Consider the rules:

- $Album@Alice(x) :- Album@Bob(x)$
- $Album@z(x) :- Album@Bob(x), Friend@Bob(z)$
- $Z@z(x) :- Album@Bob(x), FriendPhotos@Bob(Z, z)$

Bob uses the first rule to publish his photos in Alice's album, and the second to publish his photos in all of his friends' albums (peer variable z). In the last rule, different names can be used for the relations where the friends keep their photos (variable Z for a relation name).

A d-datalog program defines the meaning of intensional relations from given extensional relations. The semantics is in the spirit of the datalog semantics. More precisely:

► **Definition 3** (Semantics). Let P be a d-datalog program over some schema σ . The *immediate consequence operator* Γ_P on instances over σ is defined as follows. Let I be an instance over σ .

Consider a rule $Z_0@z(\bar{x}_0) :- R_1@p(\bar{x}_1), \dots, R_k@p(\bar{x}_k)$ of P . An *instantiation* of the rule in I is a mapping ν from its variables to the active domain (the set of values occurring in P , I , or $\text{dom}(\sigma)$), extended with the identity on constants, such that:

- for each $i \geq 1$, $R_i@p(\nu(\bar{x}_i)) \in I$; and
- $\nu(Z_0)@p(\nu(\bar{x}_0))$ is a fact for schema σ .

$\Gamma_P(I)$ is obtained by adding to I all facts $\nu(Z_0)@p(\nu(\bar{x}_0))$ where ν is an instantiation in I of some rule $Z_0@z(\bar{x}_0) :- R_1@p(\bar{x}_1), \dots, R_k@p(\bar{x}_k)$ of P . Note that Γ_P is monotonic. The *semantics* of P for an extensional instance I , denoted $P(I)$, is the mapping associating to each extensional instance I the *projection* on the intensional relations of P of the *least fixpoint* of Γ_P containing I .

Observe that a rule may “attempt” to derive an improper fact, for which $\nu(z)$ is not in $\text{dom}(\sigma)$, or $\nu(Z_0)$ is not a relation in $\sigma(\nu(z))$, or the arity is incorrect. In such cases, the fact is simply *not* derived.

► **Remark.** Consider a rule with variable peer or relation name. Suppose for instance that both are variables. A *head-instantiation* ν of that rule for a schema σ is a mapping over Z_0, z such that $\nu(z)$ is a peer of σ , $\nu(Z_0)$ an intensional relation of $\sigma(\nu(z))$, and $\text{arity}(\nu(Z_0)) = |\bar{x}_0|$. One can define similarly the notion of head-instantiation for a rule with only a variable peer or only a variable relation name. It is easy to see that the program obtained by replacing each rule by all its head-instantiations has the same semantics as the original. So if the set of peers is fixed (known in advance), one can assume that, for each rule, the name of the relation and the peer in the head are constants.

3 The access control model

In this section, we formalize the core aspects of the access control mechanism of [20]. The focus here is on the READ privilege; we will ignore the GRANT privilege (allowing a peer to define permissions on another peer’s relations) and the WRITE privilege (allowing a peer to push data to another peer’s relations), see [20]. We also provide in this section basic expressiveness and complexity results on access control.

The extensional relations at a given peer are owned by the peer. The peer can give READ privilege on these extensional relations to other peers. This is specified at each peer p using an intensional relation $acl@p$ (for *access control list*) of arity 2. A fact $acl@p(R, q)$ states that peer q is allowed to read the extensional relation $R@p$.

In the following, we assume that for each peer p , $acl \in \text{int}(p)$ and $\text{arity}(acl@p) = 2$. For instance, a rule “ $acl@p(R, z) :- Likes@p(z)$ ” can be used in a program to grant access to relation $R@p$ to all the peers z that are in relation $Likes@p$.

A d-datalog *program* P with access control (denoted d-datalog_{ac}) over some schema σ is a finite set of d-datalog rules $Z_0@z(\bar{x}_0) :- R_1@p(\bar{x}_1), \dots, R_k@p(\bar{x}_k)$, where R_1, \dots, R_k are not *acl* and the rules are of one of the following two kinds:

- **Application rule:** Z_0 is not *acl*; and
- **Access control rule:** The rule head is $acl@p(Z, z)$ for some terms Z, z .

Given a program P , the set of application rules forms the *application program* of P , denoted P_{app} , and the set of access control rules forms the (*access control*) *policy* of P , denoted P_{pol} . Facts of the form $acl@p(R, q)$ are called *access control facts*, and the others are called *application facts*. It should be noted that no such distinction is made in Webdamlog.

We distinguish here between access control and application rules to be able to formally compare access control policies.

The meaning of an access control policy P_{pol} for a given extensional instance I is clear in the absence of intensional relations: use the access rules to compute at each peer the set of peers allowed to read its extensional relations. This yields relation-at-a-time, coarse-grained access control to the extensional relations. For intensional relations, we use tuple-level fine-grained access control. Intuitively, an intensional fact can be read by a peer p if it can be derived by some application of a rule from tuples that p is already allowed to access. Then, for a d-datalog_{ac} program P , P_{app} and P_{pol} may interact recursively: the derivation of an intensional fact may yield some new permission for an extensional relation, which, in turn, may enable the derivation of a new intensional fact, and so on. The fine-grained access control at the tuple level is illustrated in an example.

► **Example 4.** Consider the program P :

$$\begin{array}{lll} P_{pol} & acl@Bob(Album, z) & :- friends@Bob(z); \\ & acl@Bob(Tagged, z) & :- friends@Bob(z); \\ P_{app} & Album@z(x) & :- Album@Bob(x), Tagged@Bob(x, z); \end{array}$$

The access control rules allow Bob's friends access to his *Album* and *Tagged* relations. The application rule transfers to a given person the photos in which he/she is tagged. Consider a photo α with tagging *Sue*, assuming she is a friend of Bob. Then the picture α belongs (intensionally) to Sue's album. A friend of Bob who will ask to see Sue's album will see the photo α .

With standard access control, peers are only be able to control access to their local data. With the proposed mechanism, they further control the *dissemination* of their data. In other words, they can control what *other* peers should do with their data. This is achieved by propagating, together with data, permissions via application rules, based on *provenance* information about derived facts. A tuple derived by some instantiation of an application rule is accessible by a peer if that peer has access to each tuple in the body of the rule.

The semantics. To define the semantics of programs, we associate with each peer p in $dom(\sigma)$ and each relation $R@p$, $R \neq acl$, a relation $\widehat{R}@p$ of arity $arity(R) + 1$. Intuitively, $\widehat{R}@p(\bar{x}, q)$ says that peer q is allowed access to the fact $R@p(\bar{x})$. The semantics is defined using a d-datalog program. We describe next the construction of that program.

► **Definition 5.** (\widehat{P} construction) The semantics of a d-datalog_{ac} program P over some schema σ for an extensional instance I over σ is defined using a d-datalog program \widehat{P} (without access control) defined as follows. Its schema consists of: (i) the extensional and intensional relations of σ ; and (ii) intensional relations $\{\widehat{R}@p \mid R@p \in \sigma(p), R \neq acl\}$.

The rules of \widehat{P} are as follows: for a tuple \bar{x} of distinct variables,

1. $\widehat{R}@p(\bar{x}, p) :- R@p(\bar{x})$ for each peer p in σ and each $R \in ext(p)$ (each peer can read its own extensional relations);
2. $\widehat{R}@p(\bar{x}, z) :- acl@p(R, z), R@p(\bar{x})$ for each peer p in σ and each $R \in ext(p)$ (each peer z entitled to read $R@p$ can read all of its tuples);
3. for each rule $acl@p(Z, z) :- R_1@p(\bar{x}_1), \dots, R_k@p(\bar{x}_k)$ in P_{pol} , a rule $acl@p(Z, z) :- \widehat{R}_1@p(\bar{x}_1, p), \dots, \widehat{R}_k@p(\bar{x}_k, p)$;

4. for each rule $Z_0@z(\bar{x}_0) :- R_1@p(\bar{x}_1), \dots, R_k@p(\bar{x}_k)$ in P_{app} and for each intensional relation $R_0 \neq acl$ occurring in σ , a rule¹

$$\widehat{R}_0@z(\bar{x}_0, y) :- Z_0 = R_0, \widehat{R}_1@p(\bar{x}_1, y), \dots, \widehat{R}_k@p(\bar{x}_k, y), \widehat{R}_1@p(\bar{x}_1, z), \dots, \widehat{R}_k@p(\bar{x}_k, z)$$
5. A rule $R@p(\bar{x}) :- \widehat{R}@p(\bar{x}, p)$ for each $p \in dom(\sigma)$ and $R \in int(p)$ ($\widehat{R}@p$ defines the local facts visible at p).

The fourth item requires that both z (the next reader) and y (potential future readers) may access the facts in the body of the rule, in order to be allowed to see the fact derived by the rule. The third item is the analog for *acl*. Note that (3.) is simpler than (4.) because the relation *acl* is only defined locally.

Clearly, the size of \widehat{P} is linear in P and the image of σ . Moreover, it is independent of the data, i.e. $dom(\sigma)$ and I . Using \widehat{P} , we define two semantics for P : state semantics, and visibility semantics.

State semantics. State semantics provides for each peer the *local* intensional facts inferred by taking into account the combined effect of the access control rules and the application rules. More precisely, the state semantics of a d-datalog_{ac} program P over schema σ is a mapping $[P]$ associating to each extensional instance I over σ the set of facts

$$[P](I) = \{R@p(\bar{a}) \in \widehat{P}(I) \mid p \in dom(\sigma), R \in int(p)\}.$$

One can easily verify by induction that $[P](I) \subseteq P(I)$. (Recall that $P(I)$ is the access-control-free semantics). The inclusion may be strict because the derivation of a fact at a peer p may be blocked because p does not have access to some data.

Visibility semantics. This semantics captures more broadly the facts at *all* peers that a given peer is allowed to see. Indeed, in addition to their local state provided by $[P]$, peers also have permission to see facts residing at *other* peers. The facts that they are allowed to see are specified by the relations $\widehat{R}@q(-, p)$ defined by the program \widehat{P} . We say that such a fact is *visible* by a peer p . For each p , we denote by $[P]_p^\vee$ the mapping associating to each instance I over $ext(\sigma)$ the set of facts $\{R@q(\bar{u}) \mid \widehat{R}@q(\bar{u}, p) \in \widehat{P}(I)\}$. We refer to $[P]_p^\vee$ as the *visibility semantics* for peer p . Clearly, for each p , $[P]_p^\vee(I)$ and $[P](I)$ agree on $int(p)$.

Intuitively, if a fact $R@q(\bar{a})$ is visible by p , then p can access it by querying the relation $R@q$. More precisely, let P' be the program obtained by adding to P a rule $temp@p(\bar{u}) :- R@q(\bar{u})$ for some new relation $temp@p$ and vector \bar{u} of distinct variables. Then $temp@p(\bar{a}) \in [P'](I)$ iff $R@q(\bar{a}) \in [P]_p^\vee(I)$, i.e. $R@q(\bar{a})$ is visible by p . Thus, visibility semantics can be reduced to state semantics by the addition of such rules.

In addition to state and visibility semantics, we consider in Section 4 the facts that a peer may *infer* from the visible ones, possibly circumventing the access control policy. We will refer to this as *implicit visibility*.

Hiding access restrictions. The above access control mechanism may be too constraining in some situations. We next consider means of relaxing it. To do so, we introduce a *hide* annotation that can be attached to atoms in rule bodies, e.g., $[hide R@q(\bar{x})]$. Intuitively, such an annotation lifts access restrictions on $R@q(\bar{x})$ by “hiding its provenance”.

We illustrate this feature with an example. For further illustration, Example 30 in the appendix shows how the hide mechanism can be used to simulate accessing a relation with binding patterns [26].

¹ Strictly speaking, equalities $Z = R_0$ are not allowed in d-datalog, but these can be easily simulated by substituting the variable by the constant everywhere in the rule.

► **Example 6.** Consider the two rules:

$$Album@z(x) :- Album@Bob(x), friend@Bob(z)$$

$$Album@z(x) :- Album@Bob(x), [hide friend@Bob(z)]$$

The first rule is used by Bob to publish his photos in all of his friends albums. Suppose Sue is a friend. Will the photos in $Album@Bob$ be transferred to $Album@Sue$? Yes, but only if Sue has read privileges on both $Album@Bob$ and $friends@Bob$. However, it may be the case that Bob wishes to keep his list of friends private, but still let his friends see his album pictures. He can do this by “hiding” the access restrictions on $friends@Bob$ as in the second rule. Intuitively, Bob is in effect reducing the protection level of the $friend$ relation, in some sense “declassifying” it.

In the example, Bob declassifies *his own* extensional relation. As we will see, “hide” also allows a peer to declassify information received from *other* peers, thus overriding their access control restrictions. In the actual Webdamlog system [20], doing so requires the peer to have GRANT privilege on that piece of information. As previously mentioned, for simplicity we do not consider explicitly the GRANT mechanism here.

Programs with *hide* are defined as follows.

► **Definition 7.** A d -datalog_{ac} program with *hide* (denoted h -d-datalog_{ac}) over some schema σ consists of: (i) a d -datalog_{ac} program $P = P_{app} \cup P_{pol}$; and (ii) a function h (called the *hide* function) whose domain h is the set P_{app} of rules², such that for each rule r , $h(r)$ is a strict subset of the atoms in the body of r . The pair (P_{pol}, h) forms the *policy* of the program.

As in Example 6, the function h is represented using annotations. More precisely, in each rule, the atoms in $h(r)$ are annotated with the keyword *hide*. For instance, the rule r that is $A :- B_1, \dots, B_5$ with $h(r) = \{B_2, B_4\}$ is denoted: $A :- B_1, [hide B_2], B_3, [hide B_4], B_5$.

We next consider how *hide* annotations modify the semantics of access control. The semantics for h -d-datalog_{ac} programs is obtained by replacing item (4) of Definition 5 with:

(4') for each application rule $Z_0@z(\bar{x}_0) :- R_1@p(\bar{x}_1), \dots, R_k@p(\bar{x}_k)$ of P_{app} , for each intensional relation $R_0 \neq acl$ occurring in σ , and some new variable y , the rule $\widehat{R}_0@z(\bar{x}_0, y) :- Z_0 = R_0, \widehat{R}_1@p(\bar{x}_1, y_1), \dots, \widehat{R}_k@p(\bar{x}_k, y_k), \widehat{R}_1@p(\bar{x}_1, q_1), \dots, \widehat{R}_k@p(\bar{x}_k, q_k)$ where for each i , if $R_i@p(\bar{x}_i)$ is not hidden in the rule, $y_i = y$ and $q_i = z$; and if it is hidden, $y_i = q_i = p$.

Note that this imposes that both y (a potential future reader) and z (the site that will host the fact) can read the facts in the body of the rule *that are not annotated by hide*, in order for the reader to be allowed to see the fact derived by the rule. For a h -d-datalog_{ac} program P , we denote by $[P]$ the state semantics of P as defined by the above program.

The next result, namely Proposition 9, shows that the use of *hide* extends the expressive power of d -datalog_{ac} relative to state semantics. (One can obtain a similar result for visibility semantics.) This is illustrated by the following example.

► **Example 8.** Consider a peer p that has a binary extensional relation $R@p$. Suppose we wish to specify that peer q sees from $R@p$ exactly the tuples of the form $(x, 0)$, and no other peer sees anything from $R@p$. As a first attempt, one might use an intensional relation R_{export} and the rule: $R_{export}@q(x, 0) :- R@p(x, 0)$.

² Because of the way we define access control rules, *hide* annotations would have no effect on them.

However, either $acl@p(R, q)$ holds, so $R@p$ is entirely visible to q ; or not, and $R_{export}@q$ is empty. Considering *hide*, assume the existence of some extensional fact $ok_q@p()$ that only q can read. Then there is a solution: $R_{export}@q(x, 0) :- ok_q@p(), [hide R@p(x, 0)]$.

► **Proposition 9.** There is a h -d-datalog_{ac} program P over schema σ for which there is no d-datalog_{ac} program \bar{P} such that, for every extensional instance I over σ , $[P](I) = [\bar{P}](I)$.

Thus, the *hide* construct strictly increases the expressivity of the language. In fact, we will show in Section 5 that h -d-datalog_{ac} is in some sense expressively complete.

The complexity of access control. We consider throughout the paper the complexity of various problems related to access control. Typically, three kinds of complexity are considered in databases: data, query, and combined complexity. In d-datalog_{ac}, the distinction between data and schema/program is less clear. For instance, the set of peers affects both the schema and the data. If there are many peers, the global program may be large, even if each peer has a small program. To capture this situation, we consider a measure assuming that the size of the program *at each peer* is bounded. This gives rise to a novel notion of complexity that we call *locally-bounded combined complexity*. More precisely, for a decision problem whose input is an extensional instance I and a d-datalog_{ac} program P over some schema σ :

- The *combined complexity* is computed as a function of $|I|$, $|P|$, and σ .
- The *data complexity* is computed as a function of $|I|$ only (σ and P are fixed).
- The *locally-bounded combined complexity* is computed as a function of $|I|$ and $|dom(\sigma)|$, assuming some fixed bound on the size of the program at each peer (so $|P|$ is linear in the number of peers).

We begin by establishing the complexity of checking the visibility of a fact.

► **Theorem 10.** *Let σ be a schema, I an extensional instance, and P a h -d-datalog_{ac} program over σ . Determining whether a fact is in $[P]_p^V(I)$ for some peer p has PTIME-complete data and locally-bounded combined complexity, and EXPTIME-complete combined complexity,*

While the data and the locally-bounded combined complexities are the same in this case, we will see later that the two differ in other settings, allowing to draw finer distinctions than the classical notions.

Static analysis of policies. To conclude this section, we briefly discuss the issue of *comparing* policies relative to a given application program, based on the visible facts they allow. This leads to the notion of a policy being *more relaxed than* another. By reduction from containment of datalog programs, one can show that this is undecidable for given policies and application program. As for datalog containment, one can consider restrictions for which the policy comparison can be performed, e.g., “frontier-guarded” rules [8]. As an alternative to comparing policies, one can consider applying syntactic transformations to a given policy in order to relax or tighten it. For example, augmenting the *hide* function of a program, or adding rules to P_{pol} , always results in a more relaxed policy. Due to space limitations, we do not further consider these issues here.

4 Implicit visibility

The purpose of access control is to analyse the ability of peers to see unauthorized information. As discussed in Section 3, a peer can access information by examining its own state or by querying relations of other peers. But can a peer infer more information beyond what

is allowed according to the policy? We capture this using the notion of *implicit visibility* (i-visibility) that we formalize next. For this, we use the auxiliary notion of “visibility instance”. For a program P over σ and a peer p , we say that an instance I_p over σ is a *visibility instance* of p if there is some instance J over $\text{ext}(\sigma)$ for which $I_p = [P]_p^V(J)$. Now we define:

► **Definition 11.** Let P be a d -datalog_{ac} program over some schema σ , p a peer and I_p a visibility instance for p . A fact $R@q(\bar{u})$ (for some q, R) is *i(mplicitly)-visible* at p given I_p , if for each instance J over $\text{ext}(\sigma)$ such that $[P]_p^V(J) = I_p$, $R@q(\bar{u}) \in J \cup [P](J)$.

It turns out that facts beyond $[P]_p^V(J)$ may be i-visible at peer p . To see how such information “leakage” can occur, suppose that we have a rule $\text{acl}@q(R, p) :- Q@q(p)$, where $Q@q$ is an extensional relation. If peer p sees some fact in $R@q$, it can infer that it has access to $R@q$, so that $Q@q(p)$ holds, although the policy may not allow p to see $Q@q$. This may in turn provide additional information on other relations. Before exploring this formally, we introduce some restrictions of policies.

► **Definition 12.** Let σ be a schema and $P = P_{pol} \cup P_{app}$ a d -datalog program.

- The policy of P is *static* iff for each rule of P_{pol} , its body is empty;
- The policy of P is *simple* iff for each rule of P_{pol} , the atoms in its body are extensional;
- The policy of P is *local* for P_{app} iff for each peer p and rule of P_{pol} at p , the atoms in its body are either extensional, or intensional but not depending on non-local relations.

We can show that with static policy, no leakage can occur.

► **Proposition 13.** Let P be a d -datalog_{ac} program over σ with static policy. For each peer p and instance I over $\text{ext}(\sigma)$, the set of i-visible facts at p is precisely $[P]_p^V(I)$.

In contrast to the above, when P_{pol} contains arbitrary rules, i-visibility provides additional information, and is in fact undecidable.

► **Theorem 14.** *It is undecidable, given a d -datalog_{ac} program P over σ , a visibility instance I_p for p , and a fact $R@q(\bar{u})$, whether $R@q(\bar{u})$ is i-visible at p given I_p . Moreover, undecidability holds even for programs with local access policies.*

The above undecidability result uses the fact that the *acl* relations are defined by datalog programs. We next show that i-visibility becomes decidable if recursion is disallowed in the definition of *acl* relations. The problem can be reduced to computing certain answers to datalog queries using exact UCQ views, which is known to be in co-NP [4]. However, using the fact that the views we use are particular UCQs, we can show that the complexity goes down to PTIME.

► **Theorem 15.** *The i-visibility problem for d -datalog_{ac} programs with simple policies is decidable in PTIME (data complexity).*

The i-visibility problem with hide. We now turn to the problem of i-visibility for d -datalog_{ac} programs with *hide*. The notions of visibility and i-visibility are adapted to this setting in the natural way. We first illustrate the fact that *hide* can lead to non-trivial i-visibility of facts, even when the *acl* policy is static.

► **Example 16.** Consider the following h - d -datalog_{ac} program P where P_{pol} consists of the rule $\text{acl}@q(Q, p) :-$ and P_{app} of the rules:

- $R_1@p(X) :- Q@q(), [\text{hide } R@q(X, Y)];$

■ $R_2@p(Y) :- Q@q(), [\text{hide } R@q(X, Y)].$

Consider the p -visibility instance $\{R_1@q(a), R_2@q(b)\}$. Note that p does not have access to $R@q$. However, it is clear that $R@q(a, b)$ is i -visible at p .

The following result shows that i -visibility is undecidable for h - d - datalog_{ac} programs even for static policies (when, by Proposition 13, no leakage occurs in the absence of *hide*). The proof is by reduction from finding certain answers to identity queries using exact datalog views, known to be undecidable [4].

► **Theorem 17.** *It is undecidable, given a h - d - datalog_{ac} program P over σ with static policy, in which *hide* is applied only to extensional relations, a peer p , a p -visibility instance I_p , and an extensional fact $R@q(\bar{a})$, whether $R@q(\bar{a})$ is i -visible at p given I_p .*

Testing information leakage. The previous result concerned i -visibility for a given instance. We finally consider the problem of testing whether a d - datalog_{ac} program has information leakage beyond that provided by the access control policy for *some* instance (the static analysis analog).

► **Definition 18.** A d - datalog_{ac} program P *leaks information* at p if for some p -visibility instance I_p there exists some fact $R@q(\bar{a}) \notin I_p$ that is i -visible at p given I_p .

We show that one cannot generally decide whether a program leaks information. However, one can do so for programs with *simple* policies. The undecidability is proved using a reduction from datalog program containment. The 2EXPTIME algorithm for simple policies is by reduction to an exponential set of inclusions of datalog programs into UCQs.

► **Theorem 19.** (1) *It is undecidable, given a d - datalog_{ac} program P and a peer p , whether P leaks information at p .* (2) *The problem is 2EXPTIME-complete if P has a simple acl policy.*

5 Achieving dissemination goals

We next consider the problem of achieving a specific data dissemination goal among peers, when a particular access control policy is imposed. The goal is specified by a d -datalog program. Clearly, a given goal may violate the policy, so it may be impossible to achieve it. We study the problem of determining whether achieving a goal is possible, and if not, how one might maximize what *can* be achieved. We then consider the issue of relaxing the access control policy in order to achieve the goal, using the *hide* mechanism. Not surprisingly, it is always possible to achieve a goal using *hide*. More interestingly, we will show how to do so while minimizing its use. But first, we consider what can be done without *hide*.

Strict adherence to the policy. Consider a policy P_{pol} and a goal d -datalog program P . We wish to know whether there is a d -datalog program P_{app} such that (i) P_{app} uses the relations of P and possibly additional intensional relations, and (ii) for each extensional instance I , $[(P_{pol} \cup P_{app})](I)$ and $P(I)$ agree on the intensional relations of P . In this case, we say that P_{app} *simulates* P under policy P_{pol} . We will see that it is generally impossible to find such a P_{app} without *hide*, and present restrictions on the policies that make it possible. When such a simulation does not exist, we will attempt to find a program that is as close as possible to the goal.

The next example illustrates how a policy may prevent achieving a goal even in the simplest setting. The example is more complicated than needed because we will also use it to illustrate finding a “maximum” simulation.

► **Example 20.** Consider the following policy and goal program:

$$\begin{array}{ll} P_{pol} & acl@p(R_1, r) :- ; & P & R@q(x) :- R_1@p(x); \\ & acl@p(R_2, r) :- ; & & R@q(x) :- R_2@p(x); \\ & acl@p(R_1, q) :- ; & & R@r(x) :- R@q(x) \end{array}$$

The d-datalog P does not simulate P under P_{pol} because q is not allowed to see the relation $R_2@p$ and therefore the relation $R@q$ does not hold tuples from $R_2@p$ under the policy P_{pol} . In such cases, we can try to find a program that is, in some sense, maximally achieves the goal. This is a nontrivial issue. In this example, a maximum application program is:

$$P_{app} : R@q(x) :- R_1@p(x); \mid R@r(x) :- R@q(x); \mid R@r(x) :- R_2@p(x).$$

Note that $[(P_{pol} \cup P_{app})] \subseteq P$ but $[(P_{pol} \cup P)] \subset [(P_{pol} \cup P_{app})]$ (as mappings).

The first result states that one cannot decide whether a program can be simulated under a particular policy.

► **Theorem 21.** *It is undecidable, given a policy P_{pol} and a goal d-datalog program P , whether there exists a d-datalog program P_{app} without `hide` such that P_{app} simulates P under P_{pol} . This holds even if P_{pol} is static.*

If such a simulation is not possible, can we find a “maximum simulation”? Let P be a d-datalog program over some schema σ and P_{pol} a policy program over σ . A d-datalog program P_{app} without `hide` is a *maximum simulation* of P under P_{pol} iff

1. $[(P_{pol} \cup P_{app})] \subseteq P$, and
2. for each P'_{app} such that $[(P_{pol} \cup P'_{app})] \subseteq P$, $[(P_{pol} \cup P_{app})] \subseteq [(P_{pol} \cup P'_{app})]$.

The question of whether a maximum simulation always exists remains open. Moreover, there does not exist an algorithm building a maximum simulation, *if such exists*.

► **Theorem 22.** *There is no algorithm that computes, given a d-datalog program P and a policy P_{pol} , a maximum simulation without `hide` P_{app} of P under P_{pol} , whenever such a maximum simulation exists. This holds even for local policies.*

While it is not known whether a maximum simulation always exists, we present informally a plausible candidate for a maximum simulation of P under P_{pol} and explore its potential. The program, denoted by $\text{MAC}(P_{pol}, P)$, is based on a simple idea: each peer collects all the extensional tuples that peer is allowed to see under P_{pol} , and then simulates P locally.

► **Definition 23.** Let P be a d-datalog program over some schema σ and P_{pol} a policy over the relations in σ . The program $P_{app} = \text{MAC}(P_{pol}, P)$ is constructed as follows:

1. For all peers $p, q, p \neq q$ and each (extensional or intensional) relation $R@q$, P_{app} has an intensional relation $R_q@p$ of the same arity as $R@q$. These relations allow p to perform a simulation of P with the data that p has access to.
2. For all peers $p, q, p \neq q$, and each extensional relation $R@q$, P_{app} has rules copying $R@q$ into $R_q@p$, if $acl@q(R, p)$ holds.
3. Finally, for each peer p , P_{app} has rules that simulate P locally with the data that p has access to.

Observe how $\text{MAC}(P_{pol}, P)$ interacts with P_{pol} . During the computation, some peer p may use rules in P_{pol} to derive a new fact $acl@p(R, q)$. This results in copying $R@p$ into $R_p@q$ which may lead to the derivations of more facts at p .

Note the connection between $\text{MAC}(P_{pol}, P)$ and P itself. By definition, $[(P_{pol} \cup P)] \subseteq [(P_{pol} \cup \text{MAC}(P, P_{pol}))]$. However, the inclusion may be strict. For instance, P may try to transfer a fact from p to q via a peer r that is not allowed to see this fact whereas it is possible to send this fact directly (with a different rule) without violating access rights.

It turns out, surprisingly, that $\text{MAC}(P_{pol}, P)$ is not always a maximum simulation of P under P_{pol} , and it is in fact undecidable whether $\text{MAC}(P_{pol}, P)$ is a maximum simulation for some given $(P_{pol}$ and P , even for local policies. However, $\text{MAC}(P_{pol}, P)$ is a maximum simulation if P_{pol} is static.

► **Theorem 24.** *Let P_{pol} be a local policy and P a d-datalog goal program over σ . (i) It is undecidable whether the program $\text{MAC}(P, P_{pol})$ is a maximum simulation of P under P_{pol} . (ii) If P_{pol} is static, then $\text{MAC}(P, P_{pol})$ is a maximum simulation of P under P_{pol} .*

Besides ensuring the existence of a maximal simulation, a simple policy is of interest for another reason: it guarantees that, if there exists some application program simulating P under P_{pol} , then P itself simulates P under that policy (details omitted).

Declassifying information. Let us now consider the issue of achieving a goal at the cost of declassifying information, in other words using the *hide* construct. There is an immediate solution that would consist in modifying every rule of the goal program P by hiding the entire body. The goal would be satisfied, but in a brutal way: each derived fact would be visible to all peers.

It is possible to realize the goal in a much more controlled way as illustrated by Example 8. In that example, special relations of the form $ok_q@p$ are used to limit as much as possible the visibility of data. The example suggests the following mild technical assumptions: (†) for all distinct peers $p, q \in \text{dom}(\sigma)$, (1.) σ contains a 0-ary extensional relation $ok_q@p$, and (2.) extensional instances of σ are assumed to contain the fact $ok_q@p()$.

We next show that (†) is sufficient to guarantee that the *hide* construct allows achieving any goal program by declassifying no more information than necessary.

► **Theorem 25.** *Let σ satisfy (†.1). For each policy P_{pol} and a d-datalog goal program P over σ , there exists an application P_{app} with *hide* over the same σ such that, for each extensional instance I satisfying (†.2), P_{app} simulates P under P_{pol} ; and on input I , a fact $R@p(u)$ is visible at $q \neq p$ for $(P_{pol} \cup P_{app})$ iff it is visible at q for $(P_{pol} \cup P)$.*

6 Accessing provenance

We considered so far the inference of individual facts using d-datalog_{ac} rules, subject to an access control policy. In many applications, it is essential for inferred facts to be accompanied by *provenance* information. In this section, we extend our approach to access control to cover provenance. We adopt a simple model of provenance of a fact, consisting of derivation trees tracing the application of the rules at different peers that participated in the inference of the fact. To simplify the presentation, we ignore *hide*. The definition of provenance can be easily adapted to the presence of *hide* (a *hide* annotation in a rule results in truncating the corresponding portion of the proof tree) and the complexity results continue to hold.

Consider a d-datalog_{ac} program P over schema σ . Let I be an extensional instance over σ , and $R@p(\bar{a})$ a fact in $P_{app}(I)$. A *provenance tree* for $R@p(\bar{a})$ is a derivation tree for $R@p(\bar{a})$ using P_{app} and I . Intuitively, we are interested in passing provenance information from peer to peer, so that a peer p not only knows that some fact $R@p(u)$ holds, but can also know how $R@p(u)$ has been derived.

► **Example 26.** Consider a schema σ with peers $\{p_0, p_1, p_2, p_3, p_4\}$, 0-ary extensional relations (propositions), $R@p_0, R@p_1$, and 0-ary intensional relations $S@p_2, S@p_3, S@p_4$. Let $I = \{R@p_0, R@p_1\}$. Consider the following application program:

$$P_{app} \quad S@p_2 :- R@p_0; \mid S@p_2 :- R@p_1; \mid S@p_3 :- S@p_2; \mid S@p_4 :- S@p_3.$$

Note that $S@p_4 \in P_{app}(I)$ and has two provenance trees (linear in this case):

$$S@p_4 \leftarrow S@p_3 \leftarrow S@p_2 \leftarrow R@p_1 \qquad S@p_4 \leftarrow S@p_3 \leftarrow S@p_2 \leftarrow R@p_0$$

Suppose we have the following access control rules in addition to P_{app} :

$$P_{pol} : \quad acl@p_0(R, p_2) :- ; \mid acl@p_0(R, p_4) :- ; \mid acl@p_1(R, p_3) :- ; \mid acl@p_1(R, p_4) :- .$$

Consider again the two provenance trees of $S@p_4 \in P_{app}(I)$. Neither satisfies the access control policy defined by P_{pol} . Indeed, the first tree violates the policy because p_2 does not have access to $R@p_1$. The second also violates the policy, because p_3 does not have access to $R@p_0$. If we add the access control rule: $acl@p_0(R, p_3) :-$ then the second provenance tree satisfies the access control policy.

Note the difference between visibility of a fact A by a peer p and visibility of its *provenance*. In order for A to be visible by p , it suffices for each fact involved in its derivation to be visible by the corresponding intermediate peer, based on its own access permissions, independently derived. In other words, peers may justify their permissions by derivations independent of each other and of the actual derivation of A . Visibility of provenance imposes a stronger condition, as it requires each intermediate peer to have access to the *entire history* of the partial derivation of p . As seen in the example, a fact A may itself be visible by p but not have any provenance tree visible by p . More formally we have:

► **Definition 27** (Provenance access control). Let P be a d-datalog_{ac} program over some schema σ and I an extensional instance over σ . A fact F has *visible provenance* if there exists a provenance tree T of F such that: For each internal node $R@p(\bar{a})$ in T and extensional fact $E@q(\bar{c})$ occurring in the subtree rooted at $R@p(\bar{a})$, we have that $acl@q(E, p) \in [P](I)$. For given P and I , $[P]^{prov}(I)$ denotes the set of facts that have visible provenance.

It is clear that visible provenance implies visibility. More precisely, one can show that for each P , σ , and each extensional instance I , $[P]^{prov}(I) \subseteq [P](I)$, but Example 26 shows the converse does not hold. We next show that, although the definition of provenance visibility is proof-theoretic, one can simulate it using a d-datalog program. However, unlike the program \hat{P} constructed earlier, the program simulating provenance visibility is exponential in the number of peers.

► **Proposition 28.** Let P be a d-datalog_{ac} program over some schema σ . There exists a d-datalog program (without access control) P^{prov} of size exponential in $dom(\sigma)$ (and polynomial in σ and P if $dom(\sigma)$ is fixed) with the same extensional relations as σ , such that for each extensional instance I , $[P]^{prov}(I)$ and $P^{prov}(I)$ agree on the intensional relations of σ .

The program P^{prov} (in the proof of the previous result) uses constants to denote sets of peers. An alternative would consist in using an extension of d-datalog with nesting, in the style of extensions of datalog with nesting [6]. (Such a nested datalog is used in the implementation in [20].)

The d-datalog program P^{prov} is exponential in the set $dom(\sigma)$ of peers. Is it possible to avoid the exponential blowup? The following complexity result implies a negative answer (subject to usual assumptions). Consider the problem of deciding, given an extensional

instance I and a program P , whether a fact is in $[P]^{prov}(I)$. Recall from Theorem 10 that the complexity of checking visibility of a fact has EXPTIME-complete combined complexity, and PTIME-complete data and locally-bounded combined complexity. Now we have:

► **Theorem 29.** *Let σ be a schema, I an extensional instance, and P a d -datalog_{ac} program over σ . Determining whether a fact is in $[P]^{prov}(I)$ has EXPTIME-complete combined complexity, PTIME-complete data complexity and PSPACE-complete locally-bounded combined complexity.*

Theorems 10 and 29 show that provenance visibility has the same combined and data complexity as the standard semantics, but different locally-bounded combined complexity. As a corollary, the exponential blowup in Proposition 28 cannot be avoided (unless PTIME = PSPACE). This highlights the usefulness of this complexity measure in making finer distinctions than the classical ones.

7 Related work

Database security and access control have been studied in depth (e.g., see [9]) since the earliest works on System R [28] and Ingres [30].

Controlling access to intensional facts in deductive languages is related to managing virtual views in SQL, which is handled differently among various database systems. When an authorized user accesses a view, it is usually evaluated with the privileges of the defining user (“definer’s rights”). Some systems (e.g. MySQL) allow the creator of a view to specify that later access to the view will be with respect to the privileges of the invoker of the view (“invoker’s rights”). This is similar in spirit to our approach.

The access control model we have described is fine-grained, unlike the SQL standard. Lefevre et al [19] propose a fine-grained access control model for implementing personal privacy policies in a relational database. They use query modification to enforce their policies, as we do, but their policy model and implementation are oriented towards a centralized database system. A commercial example of fine-grained access control is Oracle’s Virtual Private Database (VPD), which supports access control at the level of tuples or cells. VPD allows an administrator to associate an external function with a relation and automatically modifies queries to restrict access by tuple or cell. Alternative semantics for fine-grained access control have been investigated thoroughly [19, 27, 32]. Rizvi et al. [27] distinguish between Truman and Non-Truman models (the expression is motivated by the movie *The Truman Show* where the hero is unaware that he lives in an artificial environment). Query answers in our system follow the Truman paradigm: queries are not rejected because of lack of privilege but the user’s privileges limit the answers that are returned.

Fine-grained access control is also studied in [13], where predicate-based specification of authorization is supported. The inference of sensitive data from exposed data (that we study here under the name of *i*-visibility) is related to a notion studied in [33].

Our model of access control shares some features with the model of reflective database access control (RDBAC) in which access policies can be specified in terms of data contained in any part of the database. Olson et al. [22] formalize RDBAC using a version of datalog with updates [10] but their model does not include distribution, delegation, or the use of provenance. In Cassandra [18], access rights are specified using a language based on datalog with constraints. The language supports complex specifications based on “user roles”. On the other hand, fine-grained access control is not considered.

The use of provenance as a basis for access control was first noted in the context of provenance semirings [16, 7]. A security semiring can contain tuple-level security annotations

and define the rules by which they are propagated to query results. Another example of provenance-based access control is the work of Park et al. [24] in which access decisions are based on a transactional form of provenance.

The emergence of social networks and other Web 2.0 applications has led to new forms of access control. In online social networks, the distinguishing feature is that access control policy is expressed in terms of network relationships amongst members [11, 15], and this is one of the motivations of the model we presented. However, the model is intended to support the diverse requirements of access control in a variety of distributed applications.

The Webdamlog language was first described in [3] as a version of distributed datalog in which peers exchange not only facts, but also rules. Expressiveness and semantic issues were formally investigated, but access control was not considered. As already mentioned, we build here on the Webdamlog access control mechanism of [20]. Its main novelty is the specification of the access rights on an inferred tuple based on the access rights on the tuples used to derive it. The full access control mechanism of [20] is richer than the one described here, notably using also GRANT and WRITE privileges. They present an open-source implementation (with Bud [25] inside), and an experimental evaluation showing that the computational cost of access control is modest. In the Webdam project context, cryptographic techniques for enforcing access control in a distributed manner (and detecting security violations) have been considered in [5]. The techniques proposed there can be combined with those presented here.

Security in distributed systems has primarily focused on issues of remote authentication, authorization, and protection of data and distributed trust; such issues are outside the scope of our present work [1, 23].

8 Conclusion

We presented a first formal study of provenance-based access control in distributed datalog inspired by the collaborative access control mechanism of [20]. The results highlight the subtle interplay between declarative distributed computation, coarse-grained and fine-grained access control. Starting from coarse-grained access control on local extensional relations, distributed datalog computation yields fine-grained access control on derived facts based on their provenance. We also considered access control on tuples equipped with explicit provenance. We briefly studied the problem of information leakage, occurring when peers can infer unauthorized information from authorized data. We established the complexity of access control, as well as of various analysis tasks, such as detecting information leakage, comparing access policies, or the ability to achieve specified goals under a given policy. A challenging aspect of the framework is the fluid boundary of schema, data, and program, that has an impact on both semantics and complexity. For example, this led us to define a new complexity measure, locally-bounded combined complexity, that can make more subtle distinctions than classical data and query complexity.

In this first investigation, we have ignored some important aspects of the Webdamlog system presented in [20]. In Webdamlog, “nonlocal rules” allow dynamic deployment of rules from one peer to another. Most of the results presented here extend to non-local rules. We also ignored here the GRANT and WRITE privileges of Webdamlog. These raise new subtle issues, notably when access control updates are considered. Finally, delegation in Webdamlog allows peers to assign tasks to other peers. The access control of delegation is supported in Webdamlog by a mechanism called “sandboxing” that also raises interesting issues. These are left for future research.

References

- 1 M. Abadi, M. Burrows, B. Lampson, and G. Plotkin. A calculus for access control in distributed systems. In *ACM Trans. Program. Lang. Syst.*, 706-734, 1993.
- 2 S. Abiteboul, E. Antoine, G. Miklau, J. Stoyanovich, and J. Testard. [Demo] rule-based application development using WebdamLog. In *SIGMOD*, 2013.
- 3 S. Abiteboul, M. Bienvenu, A. Galland, and E. Antoine. A rule-based language for Web data management. In *PODS*, 2011.
- 4 S. Abiteboul and O. M. Duschka. Complexity of answering queries using materialized views. In *PODS*, pages 254–263. ACM, 1998.
- 5 S. Abiteboul, A. Galland, and N. Polyzotis. A model for web information management with access control. In *WebDB Workshop*, 2011.
- 6 S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- 7 Y. Amsterdamer, D. Deutch, and V. Tannen. Provenance for aggregate queries. In *PODS*, 2011.
- 8 V. Bárány, B. ten Cate, and M. Otto. Queries with guarded negation. *PVLDB*, 5(11):1328–1339, 2012.
- 9 E. Bertino and R. Sandhu. Database security-concepts, approaches, and challenges. *Dependable and Secure Computing, IEEE Transactions on*, 2(1):2–19, 2005.
- 10 A. Bonner. Transaction datalog: A compositional language for transaction programming. In *DBPL*. Springer, 1997.
- 11 B. Carminati, E. Ferrari, R. Heatherly, M. Kantarcioglu, and B. Thuraisingham. A semantic web based framework for social network access control. In *SACMAT*, pages 177–186, 2009.
- 12 A. K. Chandra, D. C. Kozen, and L. J. Stockmeyer. Alternation. *J. ACM*, 28(1):114–133, 1981.
- 13 S. Chaudhuri, T. Dutta, and S. Sudarshan. Fine grained authorization through predicated grants. In *ICDE*, pages 1174–1183. IEEE, 2007.
- 14 S. Chaudhuri and M. Y. Vardi. On the equivalence of recursive and nonrecursive datalog programs. *J. of Computer and System Sciences*, 54(1):61–78, 1997.
- 15 E. Ferrari. *Access Control in Data Management Systems*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2010.
- 16 T. J. Green, G. Karvounarakis, and V. Tannen. Provenance semirings. In *PODS*, pages 31–40, 2007.
- 17 G. Hulin. Parallel processing of recursive queries in distributed architectures. In *VLDB*, pages 87–96, San Francisco, CA, USA, 1989. Morgan Kaufmann Publishers Inc.
- 18 A. Lakshman and P. Malik. Cassandra: A decentralized structured storage system. *SIGOPS Oper. Syst. Rev.*, 44(2):35–40, 2010.
- 19 K. LeFevre, R. Agrawal, V. Ercegovac, R. Ramakrishnan, Y. Xu, and D. DeWitt. Limiting disclosure in hippocratic databases. In *VLDB*, pages 108–119. VLDB Endowment, 2004.
- 20 V. Z. Moffit, J. Stoyanovich, S. Abiteboul, and G. Miklau. Collaborative access control in WebdamLog. In *SIGMOD*, 2015.
- 21 W. Nejdl, S. Ceri, and G. Wiederhold. Evaluating recursive queries in distributed databases. *Knowledge and Data Engineering, IEEE Transactions on*, 5(1):104–121, 1993.
- 22 L. E. Olson, C. A. Gunter, and P. Madhusudan. A formal framework for reflective database access control policies. In *CCS '08: Proceedings of the 15th ACM conference on Computer and communications security*, pages 289–298, New York, NY, USA, 2008. ACM.
- 23 M. T. Özsu and P. Valduriez. *Principles of Distributed Database Systems, Third Edition*. Springer, 2011.
- 24 J. Park, D. Nguyen, and R. Sandhu. A provenance-based access control model. In *International Conference on Privacy, Security and Trust*, pages 137–144, 2012.
- 25 B. O. O. M. project. Bloom programming language. <http://www.bloom-lang.net/>.

- 26 A. Rajaraman, Y. Sagiv, and J. D. Ullman. Answering queries using templates with binding patterns. In *PODS*, pages 105–112. ACM, 1995.
- 27 S. Rizvi, A. Mendelzon, S. Sudarshan, and P. Roy. Extending query rewriting techniques for fine-grained access control. In *SIGMOD*, pages 551–562, New York, NY, USA, 2004. ACM Press.
- 28 P. G. Selinger, M. M. Astrahan, D. D. Chamberlin, R. A. Lorie, and T. G. Price. Access path selection in a relational database management system. In *SIGMOD Conference*, pages 23–34, 1979.
- 29 M. Sipser. *Introduction to the Theory of Computation*. International Thomson Publishing, 1996.
- 30 M. Stonebraker, G. Held, E. Wong, and P. Kreps. The design and implementation of INGRES. *ACM Trans. Database Syst.*, 1(3):189–222, Sept. 1976.
- 31 P. Upadhyaya, M. Balazinska, and D. Suciu. Automatic enforcement of data use policies with datalawyer. In *SIGMOD*, pages 213–225, 2015.
- 32 Q. Wang, T. Yu, N. Li, J. Lobo, E. Bertino, K. Irwin, and J.-W. Byun. On the correctness criteria of fine-grained access control in relational databases. In *VLDB*, pages 555–566, 2007.
- 33 H. Zhu, J. Shi, Y. Wang, and Y. Feng. Controlling information leakage of fine-grained access model in dbms. In *WAIM*, pages 583–590. IEEE, 2008.

Appendix for Section 3

► **Example 30.** This example illustrates how the hide mechanism can be used to simulate accessing a relation with binding patterns [26]. Suppose that peer p wishes to export an extensional binary relation R with binding pattern bf . The intuition is that one cannot obtain the entire relation, but if one provides bindings for the first column, peer p will provide the corresponding values in the second column. This is done as follows:

- $Seed@p(x) :- S@q(x)$
- $Q@q(x, y) :- Seed@p(x), [\text{hide } R@p(x, y)]$

Suppose the access control policy is such that p has read privilege on $S@q$, but q has no read privilege on $R@p$. Observe that $Seed@p$ is a copy of $S@q$, and $Q@q$ is the join of $Seed@p$ and $R@p$. Peer q cannot see $R@p$. But if q provides some values for the first column of $R@p$ (in relation $S@q$), then q will obtain in $Q@q$ the corresponding values for the second column of $R@p$.

Proof of Proposition 9 Consider the program P in Example 30. Suppose some program \bar{P} without hide simulates P . By definition of simulation, $[\bar{P}](I)$ does not contain $acl@p(R, q)$ because $[P](I)$ does not. Then consider the input $I = \{S@q(0), R@p(0, 1)\}$. Clearly, $[P](I)$ contains $Q@q(0, 1)$. On the other hand, using Lemma 31, it is seen that $[\bar{P}](I)$ does not contain $Q@q(0, 1)$, a contradiction. \square

Proof of Theorem 10 The combined and data complexities follow from the same complexities for datalog [6]. The PTIME upper bound for the locally-bounded combined complexity follows from the fact that the $d\text{-datalog}_{ac}$ program in Definition 5 uses rules of size linear in those of P (so bounded by a constant), and the number of rules is polynomial in P and $dom(\sigma)$. The lower bound follows from the fact that data complexity is already PTIME-hard. \square

Appendix for Section 4

Proof of Proposition 13 It is easy to see that for each instance I over $ext(\sigma)$, and peer p , $[P]_p^\forall(I) = [P]_p^\forall(I_p)$, where I_p coincides with I on extensional relations visible to p , and is empty everywhere else. The statement immediately follows. \square

Proof of Theorem 14 The proof is by reduction from datalog containment, known to be undecidable (see [6]). Consider two peers p and q . Let P_1, P_2 be datalog programs local at q , with disjoint intensional relations, defining relations $Q_1@q$ and $Q_2@q$ of some arity k . The schema of q consists of the relations of P_1 and P_2 together with extensional unary relations $R_1@q, R_2@q, R_3@q$, one k -ary relation $R_4@q$ and one 0-ary relation $secret@q$. The program P contains the rules of P_1 and P_2 and the following:

$$\begin{aligned}
 acl@q(R_1, p) & :- secret@q() \\
 acl@q(R_1, p) & :- Q_1@q(\bar{u}), R_4@q(\bar{u}) \\
 acl@q(R_2, p) & :- Q_1@q(\bar{u}), Q_2@q(\bar{u}), R_4@q(\bar{u}) \\
 acl@q(R_3, p) & :- R_2@q(u)
 \end{aligned}$$

where \bar{u} consists of k distinct variables. Let $I_p = \{R_1@q(1), R_3@q(1)\}$. Note that I_p is a visibility instance, since $I_p = [P]_p^\forall(I)$ for $I = \{R_1@q(1), R_2@q(1), R_3@q(1), secret@q()\}$. We claim that $secret@q()$ is i -visible at p given I_p iff $P_1(J) \subseteq P_2(J)$ for every instance J

over the extensional relations of P_1 and P_2 . Consider the *if* part. Suppose $P_1 \subseteq P_2$ and let J be an extensional instance such that $[P]_p^V(J) = \{R_1@q(1), R_3@q(1)\}$. We show that $secret@q()$ is in J . Since $R_3@q(1)$ is visible at p , $acl@q(R_3, p)$ must hold so $J(R_2@q) \neq \emptyset$. Since additionally no fact of $R_2@q$ is visible at p , it follows that $acl@q(R_2, p)$ is false so $J(Q_1@q) \cap J(Q_2@q) \cap J(R_4@q) = \emptyset$. Since $P_1(J) \subseteq P_2(J)$, it follows that $J(Q_1@q) \subseteq J(Q_2@q)$ so $J(Q_1@q) \cap J(R_4@q) = J(Q_1@q) \cap J(Q_2@q) \cap J(R_4@q) = \emptyset$. Thus, the body of the second rule is false. However, $R_1@q(1)$ is visible at p , so $acl@q(R_1, p)$ must be true. Thus, $secret@q()$ must hold in J . Since this holds for *every* such J , $secret@q()$ is i-visible at p .

Now consider the *only-if* part. The proof is by contradiction. Suppose that $P_1 \not\subseteq P_2$ and let J_0 be an instance over the extensional relations of P_1 and P_2 such that $P_1(J_0) - P_2(J_0) \neq \emptyset$. Let J be the instance extending J_0 with $\{R_1@q(1), R_2@q(1), R_3@q(1)\}$, and $J(R_4@q) = P_1(J_0) - P_2(J_0)$. Note that $[P]_p^V(J) = \{R_1@q(1), R_3@q(1)\}$. However, $secret@q()$ is false in J . It follows that $secret@q()$ is not i-visible at p , a contradiction. \square

Proof of Theorem 15 We reduce the problem to answering queries using views. Consider a d-datalog_{ac} program P over σ , p a peer, and I_p a visibility instance for p . Let $Acl(P)$ consist of all heads of rules in P_{pol} (we assume w.l.o.g. that they use no variables). Consider the view \mathbf{V} on $ext(\sigma)$ providing the truth values of $Acl(P)$ (obtained by evaluating P_{pol}), all relations in $ext(p)$, and the queries $acl@q(R, p) \wedge R@q(\bar{x})$ for $q \neq p$ (where \bar{x} are distinct variables). Thus, we can regard a view instance $\mathbf{V}(J)$ as a pair (θ, E) where θ is a truth assignment to $Acl(P)$ and E is an extensional instance for which $E|_{ext(p)} = J|_{ext(p)}$ and $E(R@q)$ is the answer to the query $acl@q(R, p) \wedge R@q(\bar{x})$ for $q \neq p$.

Clearly, $\mathbf{V}(J) = (\theta, E)$ provides enough information to compute $[P]_p^V(J)$. Indeed, consider the program \bar{P} obtained from P by replacing P_{pol} with the acl rules $acl@q(R, s) :-$ for all $acl@q(R, s)$ that are true by θ . It is easily checked that

$$[P]_p^V(J) = [\bar{P}]_p^V(E).$$

We say that a view instance (θ, E) is *compatible* with I_p if $[\bar{P}]_p^V(E) = I_p$. Note that in particular, this implies that $E = I_p|_{ext(\sigma)}$. Thus, all views compatible with I_p share the same $E = I_p|_{ext(\sigma)}$ and differ *only* in the θ component. Let $V_\theta = (I_p|_{ext(\sigma)}, \theta)$. As observed, compatibility of V_θ with I_p can be checked in PTIME. Let

$$\Theta(I_p) = \{\theta \mid V_\theta \text{ is compatible with } I_p\}.$$

For given θ , let $cert(\theta) = \cap\{[P](J) \mid \mathbf{V}(J) = V_\theta\}$. It is clear that the set of i-visible facts at p is $\bigcap_{\theta \in \Theta(I_p)} cert(\theta)$. Since the size of $\Theta(I_p)$ is constant, it is enough to show how to compute in PTIME each set $cert(V_\theta)$ for $\theta \in \Theta(I_p)$. Because P is a program with simple-acl policy, each fact $acl@q(R, s)$ is defined by a Boolean UCQ over $ext(q)$. Thus, the view \mathbf{V} is defined by UCQs, and $[P](J)$ can be simulated by a datalog program.

The idea of the proof is as follows. Suppose $\mathbf{V}(J) = V_\theta$. Then J must include all non-empty extensional relations of $I_p|_{ext(\sigma)}$. In addition, the UCQs defining the acl facts true by θ must be satisfied, so J must contain witnesses to the variables in the bodies of corresponding P_{pol} rules, of which there is a constant number. Suppose $|adom(E)| = n$. It follows that there is some constant m and extensional instances J_1, \dots, J_k , such that $|adom(J_i)| \leq n + m$ and k is polynomial in $(n + m)$, such that every J for which $\mathbf{V}(J) = V_\theta$ must contain J_i for some i . Let \mathcal{I} consist of the i for which $P_{pol}(J_i)$ agrees with θ (note that this includes falsifying the rules corresponding to the acl facts false by θ). It follows that $cert(\theta) = \cap\{[P](J_i) \mid i \in \mathcal{I}\}$, which can be computed in PTIME with respect to n . \square

Proof of Theorem 17 We sketch a reduction from the problem of finding certain answers to identity queries using exact datalog views, known to be undecidable [4]. Let $R@q$ be an extensional relation and P_V a datalog program at q defining a relation $V@q$. One cannot decide given a view of an instance over $R@q$ whether a tuple \bar{u} is in that instance.

Now consider a program P using peers q and p , where P_{pol} consists of the single rule $acl@q(S, p) :- S@q$ 0-ary. Let \bar{P}_V be obtained by replacing each occurrence of $R@q(\bar{u})$ in a rule body of P_V by $S@q(), [hide\ R@q(\bar{u})]$. In addition P has the rule $V@p(\bar{x}) :- V@q(\bar{x})$. Note that p does not have access to $R@q$.

Consider a visibility instance for p consisting of an instance over $V@p$ together with $S@q()$. A fact $R@q(\bar{u})$ i-visible by p iff \bar{u} is in the instance of $R@q$ for the datalog view $V@q$ defined by P_V . From [4], it follows that these are not computable. \square

Proof of Theorem 19 First consider (1). The proof is by reduction from equivalence of Boolean datalog programs, known to be undecidable (see [6]). Let P_1 and P_2 be Boolean datalog programs over some extensional relation $R@q$, using disjoint intensional relations and defining 0-ary relations $S_1@q$, resp. $S_2@q$. Let $T@q$ and $Q@q$ be 0-ary extensional relations. Let P_{app} consist of the rules for P_1 and P_2 , and P_{pol} of the rules:

$$\begin{aligned} acl@q(T, p) & :- S_1@q(), T@q(), Q@q() \\ acl@q(Q, p) & :- S_2@q(), T@q(), Q@q() \\ acl@q(S, p) & :- \end{aligned}$$

for every relation $S@q$ used by P_1 or P_2 . We claim that P leaks information at p iff P_1 and P_2 are not equivalent. Consider the “if” part. Suppose $P_1 \not\equiv P_2$. Let I be such that $I(T@q())$ and $I(Q@q())$ hold, $P_1(I(R@q))$ holds and $P_2(I(R@q))$ does not. Then $T@q()$ is visible at p . Therefore, $Q@q()$ is i-visible at p . However, $Q@q()$ is not visible at p because the second acl rule is false. Thus, P leaks information at p . The case when $P_2 \not\equiv P_1$ is symmetric. Conversely, suppose P_1 and P_2 are equivalent. Let J be an instance over $ext(\sigma)$. Suppose first that $S_1()$ and $S_2()$ hold in $[P]_p^y(J)$. There are two cases: (i) $T@q() \wedge Q@q()$ holds, and (ii) $T@q() \wedge Q@q()$ is false. If (i) holds, then everything is visible at p so there is no information leakage. If (ii) holds then p has no access to $T@q$ and $Q@q$ but sees that $S_1() \wedge S_2()$ hold. It can therefore infer that $T@q() \wedge Q@q()$ holds but this does not render any fact i-visible. Finally, suppose $S_1@q()$ and $S_2@q()$ are both false in $[P]_p^y(J)$. Then p can infer nothing about $T@q$ or $Q@q$ so again there is no leakage at p .

Now consider (2). We first prove the upper bound. The proof is by reduction to an exponential set of inclusions of datalog programs into UCQs, each known to be decidable in 2EXPTIME [14]. Let P be a d-datalog_{ac} program over σ with simple-acl policy, and p a peer. Thus, each acl fact $a = acl@q(R, p)$ is defined by a UCQ ψ_a . Recall the proof of Theorem 15, providing an algorithm for checking whether a given fact is i-visible at a peer p given a set of visible extensional relations. Let σ_0 be a subset of $ext(\sigma)$ including $ext(p)$. We use the characterization in the proof of Theorem 15 of p -visible facts for a given instance of σ_0 to construct, for each relation $Q@q \in \sigma - \sigma_0$, a datalog program $D(\sigma_0, Q@q)$ and UCQ $\varphi(\sigma_0, Q@q)$, both with 0-ary answer, such that $D(\sigma_0, Q@q) \not\subseteq \varphi(\sigma_0, Q@q)$ iff there exists an extensional instance J such that:

- (i) J extends $J|_{\sigma_0}$ with witnesses to variables of each CQ in $\psi_{acl@q(R, p)}$ for which $R@q \in \sigma_0$;
- (ii) $\sigma_0 = ext(p) \cup \{R@q \in ext(\sigma) \mid J \models \psi_{acl@q(R, p)}, J|R@q \neq \emptyset\}$
- (iii) there exists a fact $Q@q(\bar{a})$ in $[P](J)$ such that \bar{a} uses only constants in P or in $J|_{\sigma_0}$.

It can be checked that P leaks information at p iff there exist $\sigma_0, Q@q$, and J satisfying (i) – (iii). This builds upon the characterization of i-visible facts provided in the proof of

Theorem 15 for *fixed* $J|\sigma_0$. Intuitively, the construction of $D(\sigma_0, Q@q)$ and $\varphi(\sigma_0, Q@q)$ is done as follows. The datalog program $D(\sigma_0, Q@q)$ ensures that any extensional instance J with non-empty answer satisfies (i), (iii) and $J \models \psi_{acl@q(R,p)}$, $J|R@q \neq \emptyset$ for every $R@q \in \sigma_0$. The UCQ $\varphi(\sigma_0, Q@q)$ is true iff $J \models \psi_{acl@q(R,p)}$ and $J|R@q \neq \emptyset$ for some $R@q \notin \sigma_0$. Thus an instance J witnesses $D(\sigma_0, Q@q) \not\subseteq \varphi(\sigma_0, Q@q)$ iff it satisfies $D(\sigma_0, Q@q)$ and violates $\varphi(\sigma_0, Q@q)$, i.e. if it satisfies (i) – (iii). Each such test can be done in 2EXPTIME and we need a number of tests exponential in $ext(\sigma)$ and linear in $dom(\sigma)$ and $int(\sigma)$. This yields the overall 2EXPTIME bound.

We now prove the lower bound of this result. We reduce the problem of containment of a datalog program into a union of conjunctive queries to the problem of the not leaking. The problem of containment of a datalog program into a union of conjunctive queries is 2 EXPTIME-hard [14].

Let σ be a schema. Let P_1 and Q be a datalog program and a conjunctive query. We reduce the problem of containment of P_1 in Q to the problem of not leaking for a d-datalog P at the peer p . There exists in particular an intensional relation G of arity 0 which is the goal relation of the datalog program.

We explain how to construct P . We consider that there exists three peers p , q and u . The schema σ' is built as follows: for each relation name R in σ , there exists a relation $R@q$. There are also the external relations $secret@q$, $o@q$ and the intensional relation $T@u$. All these new relations have an arity equal to 0.

The d-datalog P is constituted of the following rules:

- for each rule of P_1 of the form, $R(\bar{x}) :- A_1(\bar{y}_1), \dots, A_k(\bar{y}_k)$ there is a rule in P ,

$$R@q(\bar{x}) :- A_1@q(\bar{y}_1), \dots, A_k@q(\bar{y}_k)$$

- $T@u() :- G@q(), o@q()$

We denote by $Q@q$ be the query obtained from Q by changing any atom $A(\bar{x})$ of Q by $A@q(\bar{x})$. The access policy P_{pol} is constituted of the following rules :

- for each relation name R in σ , $acl@q(R, p) :-$ and $acl@q(R, u) :-$
- $acl@q(o, p) :-$
- $acl@q(o, u) :- Q@q()$
- $acl@q(o, u) :- secret@q()$

The only relation which is always hidden at p is $secret@q$. The relations $R@q$ are always visible at p . The relation $o@u$ is visible at p . Thus, there is a leak at p iff $secret@q$ is i -visible at p .

Let assume that P_1 is not included in Q . Let I_1 be an instance such that $I_1 \models P_1$ and $I_1 \not\models Q$. From I_1 , we build an instance I such that for any fact $R(\bar{a})$ in I_1 , there exists a fact $R@q(\bar{a})$. There also exist the facts $o@q()$ and $secret@q()$ in I . Clearly, p sees all the facts except $secret@q()$. Because I_1 satisfies P_1 then $G@q()$ is derived. Moreover, because u can see $o@q$, then $T@u()$ is derived. Due to the fact that p can see all the extensional facts except $secret@q()$, p can see also $T@u()$. Therefore, p can deduce that u can see $o@q()$. Because p can check from its visible tuples that $Q@q()$ is not satisfied, it implies that $secret@q$ holds in I . Therefore, there is a leak at p .

Let assume that there is a leak at p . As noticed before, p sees all the relations except $secret@q()$. Therefore, the only possible leak is $secret@q()$. Due to the leak, there exists an instance I such that p can deduce from $[P]_p^V(I)$ that $secret@q$ holds in I . First, we prove by contradiction, that $T@u()$ holds in $[P, P_{pol}](I)$. Let assume that $T@u()$ does not holds in

$[P, P_{pol}](I)$. Let J be the instance obtained by removing $secret@q()$ from I . Then $[P]_p^V(J)$ is equal to $[P]_p^V(I)$. Therefore, there is no leak, whence a contradiction. Thus, $T@u()$ holds in $[P, P_{pol}](I)$. It implies that $G@q()$ holds in $[P, P_{pol}](I)$ and $o@u$ is visible at u . It implies that (i) $Q@q()$ is true or (ii) $secret@q()$ is in I . We can prove by contradiction that $Q@q()$ is not satisfied by I . Therefore, we can build from I , an instance I_1 such that P_1 is satisfied and Q is not satisfied. \square

Appendix for Section 5

We first show the following lemma, used in some of the proofs of this section.

► **Lemma 31.** *Consider a d -datalog_{ac} program $(P_{pol} \cup P_{app})$ where the policy is static and P_{app} has no hide. Let I be an extensional instance, p a peer, and I_p the restriction of I to the extensional relations visible by p according P_{pol} . Then $[(P_{pol} \cup P_{app})](I)$ agrees with $[(P_{pol} \cup P_{app})](I_p)$ on $int(p)$.*

Proof As the policy is static, satisfaction of the *acl* rules does not depend on the instance. Therefore, the extensional tuples that can be seen by p are exactly those in I_p and the tuples in $int(p)$ are derived from these visible tuples. Thus, $[(P_{pol} \cup P_{app})](I)$ agrees with $[(P_{pol} \cup P_{app})](I_p)$ on $int(p)$. \square

The following illustrates why the restriction to static policies is needed in the lemma.

► **Example 32.** Consider

- the policy $P_{pol}: acl@q(T, p) :- secret@q()$;
- the program $P_{app}: R@p() :- T@q()$;
- and the instance $I = \{T@q(); secret@q()\}$.

Observe that $R@p()$ holds in $[(P_{pol} \cup P_{app})](I)$. However, the instance I_p as defined in the lemma equals $\{T@q()\}$, and $R@p()$ does not hold in $[(P_{pol} \cup P_{app})](I_p)$. Thus, Lemma 31 does not generally hold for policies that are not static.

Proof of Theorem 21 This is by reduction from datalog containment. Let P_1, P_2 be two datalog programs over the same extensional relations and distinct intensional relations computing respectively, R_1, R_2 . Recall that one cannot decide whether $P_1 \subseteq P_2$, i.e., whether for each extensional I , (*) each tuple in $P_1(I)(R_1)$ is also in $P_2(I)(R_2)$.

Let P be the program consisting of:

$$\begin{aligned} A_1@p(u_1) &:- B_1@p(v_1), \dots, B_n@p(v_n) && \text{if } A_1(u_1) :- B_1(v_1), \dots, B_n(v_n) \text{ in } P_1 \\ A_1@p(u_1) &:- B_1@p(v_1), \dots, B_n@p(v_n) && \text{if } A_1(u_1) :- B_1(v_1), \dots, B_n(v_n) \text{ in } P_2 \\ R@q(u) &:- R_1@p(u), secret@p() \\ R@q(u) &:- R_2@p(u) \end{aligned}$$

Suppose P_{pol} specifies that all extensional relations are visible by q with the exception of $secret@p$. There exists P_{app} such that (P_{pol}, P_{app}) simulates P iff (*).

First suppose (*) hold. Then the rule with *secret* in the body has no effect. Thus (P_{pol}, P) simulates P .

Now suppose (*) does not hold. Suppose that (P_{pol}, P_{app}) simulates P for some P_{app} . Observe that we are in the conditions of Lemma 31 since the *acl* relations don't depend on intensional relations. Let I_0 be an instance such that a fact $R_1@p(u)$ is derived but $R_2@p(u)$

is not. Consider I_1 obtained by extending I_0 with the fact $secret@p()$. Since (P_{pol}, P_{app}) simulates P_{app} , it derives $R_1@p(u)$ with I_1 , so $R@q(u)$, but does not derive $R@q$ with I_0 . But by Lemma 31, (P_{pol}, P_{app}) yields the same q -tuples for I_0 and I_1 , a contradiction. Thus, such a P_{app} does not exist. \square

Proof of Theorem 22 The proof is by contradiction. Suppose that such an algorithm exists. We use a reduction from datalog containment. Let $R_1@q(), R_2@q()$ be derived by datalog programs P_1, P_2 running at peer q with the same extensional relations but distinct intensional ones. Let us give to the algorithm the input:

- the policy $P_{pol}: acl@q(S, p) :- R_1@q;$
- the target program $P: T@p() :- R_2@q(), S@q.$

Consider also

- the application program $P_{app}: T@p() :- S@q().$

First observe that, for this input, there is always a maximum application. For two cases may occur:

1. $P_1 \subseteq P_2$. But then P_{app} is a maximum.
2. $P_1 \not\subseteq P_2$. But then one can show that one cannot do better than P , i.e., P itself is a maximum.

Let Q be the application returned by the algorithm for this input. We show that:

$$(*) P_1 \not\subseteq P_2 \text{ iff for each } I, [(P_{pol}, Q)](I)(T@p) \text{ is empty.}$$

For suppose $(*)$ holds. Then we have reduced $P_1 \subseteq P_2$ (that is undecidable) to the emptiness problem for datalog (that is decidable). Thus there is no such algorithm.

To prove $(*)$, first suppose that $P_1 \subseteq P_2$. Then consider an instance I for which R_1 holds for $P_1(I)$ (we can ignore w.l.g. the case where P_1 is not satisfiable). Then R_2 holds in $P_2(I)$. Thus the program P_{app} returns $T@p$, i.e. $[(P_{pol}, P_{app})](I)(T@p)$ is nonempty. Hence $[(P_{pol}, Q)](I)(T@p)$ is nonempty since Q is maximum.

Now suppose $P_1 \not\subseteq P_2$ and suppose that $(+)$ there is an instance I such that $[(P_{pol}, P_{app})](I)(T@p)$ is nonempty. In I , a proof by Q of $T@p$ can use only $S@q$ (all the other facts are not visible by p). Now consider an instance J that derives $R_1@q$ and not $R_2@q$ and where $S@q$ holds. The proof used in I also holds. Thus $T@p$ is derived by Q for J . But it should not according to P , a contradiction. Hence there is no such I , which concludes the proof. \square

Proof of Theorem 24 We prove (i) by reduction from datalog containment. Let σ be a schema. Let P_1 and P_2 be two datalog programs over σ . Let G_1 and G_2 be the goal relations of P_1 and P_2 . Let p and q be two peers. We build the schema σ' as follows: for each relation R in σ , there exists a relation $R@q$ in σ' . We add an extensional relation $T@q$ and two intentional relations $O_1@p$ and $O_T@p$. These three relations have an arity equal to 0. We build P as follows:

- for each rule in P_1 or in P_2 , $B(\bar{x}) :- A_1(\bar{x}_1), \dots, A_k(\bar{x}_k)$, rule

$$B@q(\bar{x}) :- A_1@q(\bar{x}_1), \dots, A_k@q(\bar{x}_k)$$

- the rule $O_1@p() :- G_1@q()$ is in P
- the rule $O_T@p() :- T@q()$ is in P

P_{pol} consists of the following two rules $acl(T, p) :- G_1 @ q()$ and $acl(T, p) :- G_2 @ q()$.

Let P_{app} be the d-datalog program equal to $MAC(P_{pol}, P)$. We prove that P_{app} is a maximum simulation P over P_{pol} iff P_2 is not contained in P_1 .

Let assume that P_2 is contained in P_1 . Then a maximum simulation is P' consisting of the rules of P and the rule $O_1 @ p() :- T @ q()$. Note that if p can see $T @ q$, it implies that $G_1 @ q$ is true or $G_2 @ q$ is true. Because P_2 is included in P_1 , if $G_2 @ q()$ holds then $G_1 @ q()$ holds, which shows the correctness of the previous rule.

Let assume that P_2 is not contained in P_1 . We prove by contradiction that P_{app} is not a maximum application. The two only tuples that could not be derived by P_{app} are $O_1 @ p()$ and $O_T @ p()$. In both cases, due to the fact that P_1 is not included in P_2 , we can find to instances I and J not deriving the same tuples but having the same access control facts. Therefore, there is a contradiction.

Now consider (ii). Let P'_{app} be another program such that $[P_{pol}, P'_{app}] \subseteq P$. Let I be an input and $R @ p(u)$ a fact over σ derived by (P_{pol}, P'_{app}) on input I . We show that $R @ p(u)$ is also derived by (P_{pol}, P_{app}) on input I . Since (P_{pol}, P'_{app}) does not violate the policy, and the policy does not depend on I , p has all the facts needed to derive $R @ p(u)$. And since (P_{pol}, P_{app}) simulates P , it must also derive $R @ p(u)$. Thus (P_{pol}, P_{app}) is maximum. \square

Proof of Theorem 25 The program P' is defined as follows. P'_{pol} consists of P_{pol} together with the rules $acl @ p(ok_q, q) :-$ for all peers q . P'_{app} is obtained by modifying P_{app} as follows. First, for every rule having $R @ p$ or $R' @ p$ in the head, add to its body the atom $ok_p @ p()$. Second, add the rules

$$R @ p(\bar{u}) :- ok_q @ p(), [\text{hide } R @ p(\bar{u})], [\text{hide } R' @ p(\bar{u}, q)]$$

for each $q \neq p$. It is clear that P' has the desired properties. \square

Appendix for Section 6

Proof of Proposition 28 Observe first that a given intensional fact $R @ p(\bar{a})$ may have provenance trees whose extensional facts are visible by distinct (even disjoint) sets of peers. The program P^{prov} must remember, for each such $R @ p(\bar{a})$, all sets F of peers such that $R @ p(\bar{a})$ has some provenance tree whose extensional facts are visible by the peers in F . Since a d-datalog program cannot create sets, and since the set of peers is fixed, P^{prov} uses distinct constants associated to each subset F of $dom(\sigma)$. By slight abuse of notation, we also denote by F the constant associated to the set F of peers. The extensional relations are those of σ . The intensional relations of P^{prov} are:

- the intensional relations of \widehat{P} .
- for each $p \in dom(\sigma)$, a relation $\overline{acl} @ p$ of arity 2;
- for each $R \in (\sigma(p) \cap Int) - \{acl\}$, an intensional relation $\overline{R} @ p$ of arity $arity(R @ p) + 1$

The rules are the following:

1. rules (1-4) of \widehat{P}
2. $\overline{acl} @ p(R, \{f\}) :- acl @ p(R, f)$
for each $p \in dom(\sigma)$, and each extensional relation R in $\sigma(p)$.
3. $\overline{acl} @ p(R, F) :- \overline{acl} @ p(R, F_1), \overline{acl} @ p(R, F_2)$,
for each $F_1, F_2 \subseteq dom(\sigma)$, $F = F_1 \cup F_2$.

4. $\overline{R}@q(\bar{x}_0, F) :- L_1, \dots, L_k$
for each rule $R_0@q(\bar{x}_0) :- R_1@p(\bar{x}_1), \dots, R_k@p(\bar{x}_k)$ of P_{app} ,
where $L_i = [R_i@p(\bar{x}_i), \overline{acl}@p(R_i, F_i)]$ if R_i is extensional, and $L_i = \overline{R}_i@p(\bar{x}_i, F_i)$ if R_i is intensional, for $F_i \subseteq dom(\sigma)$ and $F = \cap_i F_i$ that includes q .
5. $R@p(\bar{x}) :- \overline{R}@p(\bar{x}, v)$,
for each $R \in (\sigma(p) \cap Int) - \{acl\}$, where v is a new variable.

Intuitively, the relation $\overline{acl}@p$ provides the *sets* of peers allowed to see each extensional relation of P . Note that this may lead to using the exponentially many constants corresponding to all the subsets of the set of peers. Note that the $\overline{acl}@p$ is defined using the $acl@p$ relation that is computed according to the usual access control policy for individual facts, defined by \hat{P} . A fact $\overline{R}@p(\bar{a}, F)$ provides a *set* F of peers for which there exists some provenance tree of $R@p(\bar{a})$ such that all of its extensional facts are visible by all the peers in F .

One can verify by induction on the depth of the provenance tree that $P^{prov}(I) = [P]^{prov}(I)$, intuitively the set F carries a set of peers that includes all the peers to which the leaves of the provenance tree belongs. \square

Proof of Theorem 29 The EXPTIME and PTIME upper bounds for combined and data complexity follow from Proposition 28 and the same upper bounds for datalog. The lower bounds also follow from the EXPTIME-hardness and PTIME-hardness of combined and data complexity for datalog. Consider the locally-bounded combined complexity. For the upper bound, we exhibit an APTIME algorithm to determine whether a given fact is visible by a peer, and use the fact that APTIME = PSPACE [12]. Consider the set of instantiations of rules in P in the active domain of I . Since P is linear in $dom(\sigma)$ and all rules have bounded size, this is polynomial in $I \cup dom(\sigma)$. The APTIME algorithm works as follows. Suppose we wish to determine whether $R@p(\bar{a})$ has visible provenance. First, the acl relations are computed using the program \hat{P} . By Theorem 10, this has PTIME locally-bounded combined complexity. Next, we need to check the existence of a provenance tree for $R@p(\bar{a})$, satisfying the conditions of Definition 27. The algorithm non-deterministically produces pairs $\langle Q@q(\bar{c}), F \rangle$, where $Q@q(\bar{c})$ is a ground fact and F a set of peers. Intuitively, F consists of the set of peers so far used along the path from $R@p(\bar{a})$ to $Q@q(\bar{c})$ in a top-down derivation of $R@p(\bar{a})$. If $Q@q$ is extensional, this indicates that access permission to $Q@q$ is needed for all peers in $F - \{q\}$ (recall that q always has access to $Q@q$).

In more detail, $\langle Q@q(\bar{c}), F \rangle$ is initialized to $\langle R@p(\bar{a}), \{p\} \rangle$. Existential and universal moves occur as follows. Given a current $\langle Q@q(\bar{c}), F \rangle$, where $Q@p$ is intensional, an existential move chooses an instantiated rule of P ,

$$Q@q(\bar{c}) :- R_1@f(\bar{a}_1) \dots R_n@f(\bar{a}_n).$$

Then, a universal move transitions nondeterministically to each $\langle R_i@f(\bar{a}_i), F \cup \{f\} \rangle$. The computation terminates and accepts iff $Q@q$ is extensional and $F - \{q\} \subseteq \{f \mid acl@q(Q, f) \in \hat{P}(I)\}$. The length of each computation branch is polynomial because only a polynomial number of sets F of peers is generated along each branch, the number of facts of the form $Q@q(\bar{c})$ is polynomial, and the computation can be pruned (and fails) as soon as a pair $\langle Q@q(\bar{c}), F \rangle$ is generated twice. Clearly, $R@p(\bar{a})$ is accepted by the above APTIME algorithm iff $R@p(\bar{a}) \in [P]^{prov}(I)$.

For the lower bound, we use a reduction from Quantified Boolean Formula (QBF), known to be PSPACE-complete [29]. Let $\psi = Q_1x_1 \dots Q_nx_n\alpha(x_1, \dots, x_n)$ be a QBF ($Q_i \in \{\exists, \forall\}$). We can assume w.l.o.g. that all negations in α are pushed to the leaves (so they apply only

to variables). We construct $dom(\sigma), \sigma, P, I$ of size polynomial in ψ , and a fact $R@p$ for some $p \in dom(\sigma)$ and proposition R , such that ψ is true iff R is visible by p under global-proof semantics.

We use the set of peers $dom(\sigma) = \{p, e, a_i, t_i, f_i, l_i \mid 1 \leq i \leq n\}$. The schema σ is defined as follows:

peer	schema
<i>start</i>	$R : 0$
a_i	$T_i, F_i : 0$
t_i	$T_i : 0$
f_i	$F_i : 0$
l_i	$L : 0, L_{\neg} : 0, S : 1, S_{\neg} : 1$
e	$E : 1, F_{\wedge} : 3, F_{\vee} : 3, W : 1$

The extensional relations are $F_{\wedge}@e, F_{\vee}@e$, which represent the formula α , and $W@e, L@l_i, L_{\neg}@l_i, S@l_i, S_{\neg}@l_i, 1 \leq i \leq n$. More precisely, the following facts are in I :

- for each subformula φ of α , $F_{\wedge}(\varphi, \varphi_1, \varphi_2)$ if $\varphi = \varphi_1 \wedge \varphi_2$, $F_{\vee}(\varphi, \varphi_1, \varphi_2)$ if $\varphi = \varphi_1 \vee \varphi_2$
- $L@l_i(), L_{\neg}@l_i(), 1 \leq i \leq n$

Also,

- $W@e$ contains all peers
- $S@l_i$ contains all peers other than f_i and $S_{\neg}@l_i$ contains all peers other than $t_i, 1 \leq i \leq n$.

Intuitively, there are two sorts of rules. In a first stage, we have rules that simulate choices of truth values for $x_1 \dots, x_n$. If x_i is existentially quantified, at least one choice must succeed, i.e. lead to satisfaction of α . If x_i is universally quantified, both choices must succeed. In a second stage, we have rules that evaluate α . The access control policy P_{pol} connects the truth values chosen for the x_i to the truth of the formula, with the desired quantification for each variable.

The rules in P_{app} are the following:

- $R@start :- T_1@t_1 | F_1@f_1$ if x_1 is existentially quantified
- $R@start :- T_1@a_1, F_1@a_1$
 $T_1@a_1 :- T_1@t_1$
 $F_1@a_1 :- F_1@f_1$
 if x_1 is universally quantified

For each $i, 1 \leq i < n$,

- $T_i@t_i :- T_{i+1}@t_{i+1} | F_{i+1}@f_{i+1}$ if x_{i+1} is existentially quantified
- $T_i@t_i :- T_{i+1}@a_{i+1}, F_{i+1}@a_{i+1}$
 $T_{i+1}@a_{i+1} :- T_{i+1}@t_{i+1}$
 $F_{i+1}@a_{i+1} :- F_{i+1}@f_{i+1}$
 if x_{i+1} is universally quantified
- $T_n@t_n :- E@e(\alpha)$ and $F_n@f_n :- E@e(\alpha)$
- $E@e(x) :- F_{\wedge}@e(x, y, z), E@e(y), E@e(z)$
- $E@e(x) :- F_{\vee}@e(x, y, z), E@e(y)$
- $E@e(x) :- F_{\vee}@e(x, y, z), E@e(z)$
- $E@e(x_i) :- L@l_i$
- $E@e(\neg x_i) :- L_{\neg}@l_i$

The access control policy P_{pol} allows full access by all peers to relations $F_{\wedge}@e$ and $F_{\vee}@e$, using the rules

$$\begin{aligned}acl@e(F_{\wedge}, x) &:- W@e(x) \\acl@e(F_{\vee}, x) &:- W@e(x)\end{aligned}$$

For $L@l_i$ and $L_{\neg}@l_i$, permissions are defined by the rules:

$$\begin{aligned}acl@l_i(L, x) &:- S@l_i(x) \\acl@l_i(L_{\neg}, x) &:- S_{\neg}@l_i(x)\end{aligned}$$

Intuitively, the last two permission rules say that, in order for x_i to evaluate to true, its truth value may not be *false* (so peer f_i is excluded as an ancestor of $L@l_i()$ in the proof tree of $R@p$) and in order for $\neg x_i$ to be true, the value chosen for x_i may not be *true* (so peer t_i is excluded as an ancestor of $L_{\neg}@l_i()$ in the proof tree). It can be checked that $R@start$ is visible by $start$ under global-proof semantics iff ψ is true. Note that σ and I are polynomial in ψ . Moreover, the size and number of rules of P at each peer are bounded by a constant. \square