

# Balancing expressiveness and inexpressiveness in view design: full version

Michael Benedikt and Pierre Bourhis and Louis Jachiet and Efthymia Tsamoura

## Abstract

We study the design of data publishing mechanisms that allow a collection of autonomous distributed data-sources to collaborate to support queries. A common mechanism for data publishing is via *views*: functions that expose derived data to users, usually specified as declarative queries. Our autonomy assumption is that the views must be on individual sources, but with the intention of supporting integrated queries. In deciding what data to expose to users, two considerations must be balanced. The views must be sufficiently expressive to support queries that users want to ask – the *utility* of the publishing mechanism. But there may also be some expressiveness restriction. Here we consider two restrictions, a *minimal information* requirement, saying that the views should reveal as little as possible while supporting the utility query, and a *non-disclosure* requirement, formalizing the need to prevent external users from computing information that data owners do not want revealed. We investigate the problem of designing views that satisfy both an expressiveness and an inexpressiveness requirement, for views in a restricted declarative language (conjunctive queries), and for arbitrary views.

## 1 Introduction

The value of data is increased when data owners make their data available through publicly-accessible interfaces. The value is magnified even further when multiple data owners publish information from related datasets; this allows users to answer queries that require linking information across datasources.

But the benefits of data publishing come with a corresponding risk of revealing too much. For example, here may be information that a data owner wishes to protect, and a user may be able to infer this information either from the published data in isolation, or from the data published by all parties as a whole. There is thus a need to provide data publishing mechanisms that are simultaneously expressive enough to be useful – they enable

users to answer appropriate queries – while satisfying some expressiveness restriction.

In data publishing much of the focus has been on disclosure via the familiar mechanism of *views* – declarative queries whose output is made available to users as a table. In this context the competing requirements on a publishing mechanism have been primarily studied in isolation. There is extensive work on analysis of the utility of a set of views: namely given a query, can it be answered using the views, see, e.g. (Halevy 2001; Calvanese et al. 2012). There has also been research concerning analysis of whether a given set of views obeys some expressiveness *restriction*. The negation of answerability is clearly too weak a restriction, since it just guarantees that on some instance the query answer can not be computed from the view images. A relevant query-based notion of expressiveness restriction is *data-independent privacy*: given the views and a set of “secret” queries, check that the secret query answers can not be computed from the views on any source data. Variants of this problem have been studied in (Nash and Deutsch 2007; Benedikt et al. 2016; Benedikt, Cuenca Grau, and Kostylev 2018; Benedikt et al. 2019). Expressiveness restrictions with a similar flavor have also been studied in the context of ontologies (Bonatti and Sauro 2013). But the question of whether there is a query-independent expressiveness restriction appropriate for views based on traditional database queries, as well as the question of how one obtains views that satisfy both expressiveness and inexpressiveness requirements, has not been considered in the context of traditional queries and views, to the best of our knowledge.

In the multi-agent systems community there is work concerned with protocols allowing a set of agents to learn a distribution on source data without revealing certain local information (e.g. (Fernández-Duque and Goranko 2016)). The notion of security is “exact and information-theoretic”, as in this work: they deal with getting the exact answer to a query using a deterministic algorithm, but without restricting the computational

power used in answering. But the setting, unlike in this work, is propositional.

A larger body of work comes from privacy research, considering the design of mechanisms that achieve a mix of expressiveness and inexpressiveness goals, usually denoted as “utility” and “privacy”. However, the focus there has been on mechanisms based on probabilistic transformations or, more generally, probabilistic protocols (see for example (Chaum, Crépeau, and Damgård 1988; Dwork 2006; Dwork and Roth 2014)).

In contrast, we consider the question of designing views that use *traditional database queries*, with no randomization, so that conflicting requirements of expressiveness and inexpressiveness are satisfied. Both our expressiveness and inexpressiveness requirements will be given in terms of *exact information-theoretic criteria*: they will be defined in terms of what queries can be answered exactly (as opposed to probabilistically) by a party with unlimited computation power. Due both to the difference in the mechanisms we consider and the requirements we impose, our contribution has a very different flavor from prior lines of work.

**Example 1.** *A health study hosted by a government agency holds information about certain treatments, with an abstraction of the data being a database with schema  $\text{Trtmnt}(\text{pid}, \text{tinfo})$ , where  $\text{pid}$  might be a national insurance number.*

*Demographic information about patients is stored by another agency, in a table  $\text{Patient}(\text{pid}, \text{age}, \text{address})$ . The agencies are completely autonomous, perhaps even in distinct administrative regions. But they want to co-operate to support certain queries over the data that legitimate researchers might wish; for example about the relationships between treatment and age:*

$$Q = \exists \text{pid}, \text{address}, \text{tdate}. \text{Trtmnt}(\text{pid}, \text{tinfo}, \text{tdate}) \wedge \text{Patient}(\text{pid}, \text{age}, \text{address})$$

*Of course, the two parties could agree to an encryption schema on the patient identifiers, and then expose encrypted versions of their local schemas. But this would require both strong co-operation of the parties, and the use of views beyond traditional database queries.*

*But assuming that the parties are restricted to using traditional queries, what is the most restrictive thing that they can do while supporting the ability to answer  $Q$ ? Intuitively, the most restrictive views would correspond to one party revealing the projection of the  $\text{Trtmnt}$  table on  $\text{pid}$  and  $\text{tinfo}$  and the other revealing the projection of the  $\text{Patient}$  table.*

*If this intuition is correct, it would imply that nothing the parties can do with traditional queries can avoid revealing which patients had particular treatments. That is, they have no choice but to allow an external party to learn the answer to query*

$$p = \exists \text{tdate}. \text{Trtmnt}(\text{pid}, \text{tinfo}, \text{tdate})$$

*Our results will validate this obvious answer – in this example, the projections described above are the  $CQ$  views that reveal minimal information, and one can not support the disclosure of the join query while protecting a query on these join attributes. Further we will show that even using encryption – indeed, even using any deterministic function – the parties can not not reveal less information while supporting exact answering of the query  $Q$ . In contrast, we will also show that in some cases the parties can obtain combinations of expressiveness and inexpressiveness requirements by using counter-intuitive view combinations.*

Our goal here is to look at the problem of designing independent views over multiple relational data-sources that satisfy both expressiveness and inexpressiveness requirements. Our expressiveness requirement (usefulness) will be phrased in terms of the ability to answer a relational query, where answering is the traditional deterministic notion used in database theory and KR. For expressiveness limitations we require the views to be useful but to *minimize information* within a class of views. We also consider an expressiveness limitation specified by *non-disclosure* of a set of *secret queries*. Our contributions include formalizing these notions, characterizing when minimal information views exist and what form they take, and determining when views exist that satisfy utility and expressiveness limitations. We look at these problems both for views given in the standard view language of conjunctive queries, and for arbitrary views. We also consider the impact of background knowledge on these problems.

**Organization.** Section 2 gives database and logic preliminaries, and then goes on to formalize our expressiveness and inexpressiveness restrictions. Section 3 deals with the variant of the problem where the only views considered are conjunctive query views, while Section 4 shows how the situation changes when arbitrary views can be utilized. Section 5 gives some extensions to the situation when background knowledge about the sources is present. We close with conclusions in Section 6. In this extended version we include details of most proofs in the appendix.

## 2 Preliminaries

### 2.1 Basic definitions

The bulk of this subsection reviews standard definitions from databases and KR. But it includes two notions, DCQs and distributed schemas, that are less standard.

**Databases and queries.** A *schema* consists of a finite set of relations, each with an associated arity (a non-negative number). An *instance* of a schema is an assignment of each relation in the schema of arity  $n$  to a collection (finite or infinite) of  $n$ -tuples of elements. Given an instance  $\mathcal{I}$  and a relation  $R$ , we let  $\mathcal{I}(R)$  be the  $n$ -tuples assigned to  $R$  in  $\mathcal{I}$ . The *active domain* of

an instance  $\mathcal{I}$ , denoted  $\text{adom}(\mathcal{I})$ , is the set of elements occurring in some tuple of some  $\mathcal{I}(R)$ . A *fact* consists of a relation  $R$  of arity  $n$  and an  $n$ -tuple  $\mathbf{t}$ . We write such a fact as  $R(\mathbf{t})$ . An instance can equivalently be thought of as a set of facts.

An  $n$ -ary *query* is a function from instances of a fixed schema  $\mathcal{S}$  to some set. We refer to  $\mathcal{S}$  as the *input schema* of the query. A *Conjunctive Query* (CQ) is a formula of the form  $\exists \mathbf{y} \bigwedge_i A_i$  where  $A_i$  are atoms from the schema. A Boolean CQ (BCQ) is a CQ with no free variables. Following the notation used in several places in the literature (e.g. (Arenas, Barceló, and Reuter 2011)) we define a *union of CQs* (UCQ) to be a disjunction of CQs satisfying the *safety condition* where the free variables of each disjunct are the same. Disjunctions of CQs where the safety condition is dropped will play a key role in this paper. The terminology for such queries is less standardized, but we refer to them as *disjunctions of CQs* (DCQs). Relational algebra is the standard algebraic presentation of first-order relational queries: queries are build up from relations by union, difference, selection, and product. (Abiteboul, Hull, and Vianu 1995).

For a logical formula  $\rho$  with free variables  $x_1, \dots, x_n$  and an instance  $\mathcal{I}$ , a *variable binding*  $\sigma$  for  $\rho$  in  $\mathcal{I}$  is a mapping taking each  $x_i$  to an element of  $\text{adom}(\mathcal{I})$ . We can apply  $\sigma$  to  $\rho$  to get a new formula  $\sigma(\rho)$  where each  $x$  is replaced by  $\sigma(x)$ . Assuming an ordering of the free variables as  $x_1, \dots, x_n$  we may identify a  $k$ -tuple  $\mathbf{t}$ , writing  $\rho(\mathbf{t})$  to mean that  $t_i$  is substituted for  $x_i$  in  $\rho$ .

A *homomorphism* between instances  $\mathcal{I}_1$  and  $\mathcal{I}_2$  is a function  $f$  from  $\text{adom}(\mathcal{I}_1)$  to  $\text{adom}(\mathcal{I}_2)$  such that  $R(c_1, \dots, c_m) \in \mathcal{I}_1$  implies  $R(f(c_1), \dots, f(c_m)) \in \mathcal{I}_2$ . The notion of homomorphism from a CQ to an instance and from a CQ to a CQ is defined similarly. In the case of CQ-to-CQ homomorphisms we additionally require that the mapping be the identity on any free variables or constants. The *output* of a CQ  $Q$  on an instance  $\mathcal{I}$ , denoted  $Q(\mathcal{I})$  consists of the restrictions to free variables of  $Q$  of the homomorphisms of  $Q$  to  $\mathcal{I}$ . The output of a UCQ is defined similarly. We can choose an ordering of the free variables of  $Q$ , and can then say that the output of  $Q$  on  $\mathcal{I}$  consists of  $n$ -tuples. We write  $\mathcal{I}, \mathbf{t} \models Q$  for an  $n$ -tuple  $\mathbf{t}$ , if  $\mathbf{t} \in Q(\mathcal{I})$ . We analogously define  $\mathcal{I}, \sigma \models Q$  for a variable binding  $\sigma$ . We sometimes refer to a homomorphism of a BCQ into an instance as a *match*. For a logical formula  $\rho(\mathbf{x})$  and a tuple of elements  $\mathbf{t}$ ,  $\rho(\mathbf{t})$  denotes the formula where each  $x_i$  is substituted with  $t_i$ .

A CQ  $Q_0$  is a *subquery* of a CQ  $Q$  if the atoms of  $Q_0$  are a subset of the atoms of  $Q$  and a variable of  $Q_0$  is free in  $Q_0$  if and only if it is free in  $Q$ . A *strict subquery* of  $Q$  is a subquery of  $Q$  that is not  $Q$  itself;  $Q$  is *minimal* if there is no homomorphism from  $Q$  to a strict subquery of  $Q$ .

A *view* over a schema  $\mathcal{S}$  consists of an  $n$ -ary relation  $V$  and a corresponding  $n$ -ary query  $Q_V$  over relations from  $\mathcal{S}$ . Given a collection of views  $\mathcal{V}$  and instance  $\mathcal{I}$ , the *view-image* of  $\mathcal{I}$ , denoted  $\mathcal{V}(\mathcal{I})$  is the instance that interprets each  $V \in \mathcal{V}$  by  $Q_V(\mathcal{I})$ . We thus talk about *CQ views*, *UCQ views*, etc: views defined by formulas within a class.

**Distributed data and views.** A *distributed schema* (d-schema)  $\mathcal{S}$  consists of a finite set of *sources*  $\text{Srcs}$ , with each source  $s$  associated with a *local schema*  $\mathcal{S}_s$ . We assume that the relations in distinct local schemas are pairwise disjoint. In Example 1 our distributed schema consisted of two sources, one containing Trtmnt and the other containing Patient. A *distributed instance* (d-instance) is an instance of a distributed schema. For a source  $s$ , an  $s$ -instance is an instance of the local schema  $s$ . Given a d-instance  $\mathcal{D}$ , we denote by  $\mathcal{D}_s$  the restriction of  $\mathcal{D}$  to relations in  $s$ . If d-schema  $\mathcal{S}$  is the disjoint union of  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , and we have sources  $\mathcal{I}_1$  for  $\mathcal{S}_1$  and  $\mathcal{I}_2$  for  $\mathcal{S}_2$ , then we use  $(\mathcal{I}_1, \mathcal{I}_2)$  to denote the union of  $\mathcal{I}_1$  and  $\mathcal{I}_2$ , which is an instance of  $\mathcal{S}$ .

For a given d-schema a *distributed view* (d-view)  $\mathcal{V}$  is an assignment to each source  $s$  of a finite set  $\mathcal{V}_s$  of views over its local schema. Note that here is our “autonomy” assumption on the instances: we are free to design views on each local source, but views can not cross sources. We can similarly talk about CQ-based d-views, relational algebra-based d-views, etc.

**Existential Rules.** Semantic relationships between relations can be described using *existential rules*, logical sentences of the form  $\forall \mathbf{x}. \lambda \rightarrow \exists \mathbf{y}. \rho$ , where  $\lambda$  and  $\rho$  are conjunctions of relational atoms. The notion of a formula  $\rho$  *holding* in  $\mathcal{I}$  (or  $\mathcal{I}$  *satisfying*  $\rho$ , written  $\mathcal{I} \models \rho$ ) is given by first-order logic. Let  $\Sigma$  be a set of rules,  $\mathcal{I}$  be a finite instance, and  $Q$  be a BCQ. We write  $\mathcal{I} \wedge \Sigma \models Q$  to mean that every instance containing  $\mathcal{I}$  and satisfying  $\Sigma$  also satisfies  $Q$ . For two CQs  $Q$  and  $Q'$ , we similarly write  $Q \wedge \Sigma \models Q'$  to mean that every instance satisfying  $Q$  and  $\Sigma$  satisfies  $Q'$ .

## 2.2 Problem formalization

We now give the key definitions in the paper, capturing our expressiveness requirements (“usefulness”) and expressiveness limitations (“minimal usefulness” and “non-disclosure”).

**Definition 1.** Two d-instances  $\mathcal{D}$  and  $\mathcal{D}'$  are *indistinguishable by a d-view*  $\mathcal{V}$  (or just  *$\mathcal{V}$ -indistinguishable*) if  $V(\mathcal{D}) = V(\mathcal{D}')$  holds, for each view  $V \in \mathcal{V}$ .

Since each view  $V \in \mathcal{V}_s$  is defined over relations occurring only in  $\mathcal{V}_s$ , we can equivalently say that  $\mathcal{D}$  and  $\mathcal{D}'$  are  $\mathcal{V}$ -indistinguishable if  $V(\mathcal{D}_s) = V(\mathcal{D}'_s)$  holds, for each  $V \in \mathcal{V}_s$ .

**Definition 2.** A d-view  $\mathcal{V}$  determines a query  $Q$  at a d-instance  $\mathcal{D}$  if  $Q(\mathcal{D}') = Q(\mathcal{D})$  holds, for each  $\mathcal{D}'$  that is

$\mathcal{V}$ -indistinguishable from  $\mathcal{D}$ .

The  $d$ -view  $\mathcal{V}$  is useful for  $Q$  if  $\mathcal{V}$  determines  $Q$  on every  $d$ -instance (for short, just “ $\mathcal{V}$  determines  $Q$ ”).

Usefulness for a given query  $Q$  will be our expressiveness requirement on  $d$ -views. Our first inexpressiveness requirement captures the idea that we want to reveal as little as possible:

**Definition 3** (Minimally useful views). *Given a class of views  $\mathcal{C}$  and a query  $Q$ , we say that a  $d$ -view  $\mathcal{V}$  is minimally useful for  $Q$  within  $\mathcal{C}$  if  $\mathcal{V}$  is useful for  $Q$  and for any other  $d$ -view  $\mathcal{V}'$  using views in  $\mathcal{C}$ ,  $\mathcal{V}'$  determines the view definition of each view in  $\mathcal{V}$ .*

We look at another inexpressiveness requirement that requires an external party to not learn about another query.

**Definition 4** (Non-disclosure). *A non-disclosure function specifies, for each query  $p$ , the set of  $d$ -views  $\mathcal{V}$  that are said to disclose  $p$ . We require such a function  $F$  to be determinacy-compatible: if  $\mathcal{V}_2$  discloses  $p$  according to  $F$  and  $\mathcal{V}_1$  determines each view in  $\mathcal{V}_2$ , then  $\mathcal{V}_1$  also discloses  $p$  according to  $F$ . If  $\mathcal{V}$  does not disclose  $p$ , we say that  $\mathcal{V}$  is non-disclosing for  $p$  (relative to the given non-disclosure function).*

When we can find minimal useful views, this tells us something about non-disclosure, since it is easy to see that if we are looking to design views that are useful and non-disclosing, it suffices to consider minimal information views, assuming they exist:

**Proposition 1.** *Suppose  $d$ -view  $\mathcal{V}$  is minimally useful for  $Q$  within  $\mathcal{C}$ , and there is a  $d$ -view based on views in  $\mathcal{C}$  that is useful for  $Q$  and non-disclosing for  $p$  according to non-disclosure function  $F$ . Then  $\mathcal{V}$  is useful for  $Q$  and non-disclosing for  $p$  according to  $F$ .*

There are many disclosure functions that are determinacy-compatible. But in our examples, our complexity results, and in Subsection 5.3, we will focus on a specific non-disclosure function:

**Definition 5.** *A  $d$ -view  $\mathcal{V}$  is universal non-inference non-disclosing (UN non-disclosing) for a CQ  $p$  if for each instance  $\mathcal{I}$  and each tuple  $\mathbf{t}$  with  $\mathcal{I}, \mathbf{t} \models p$ ,  $\mathcal{V}$  does not determine  $p(\mathbf{t})$  at  $\mathcal{I}$ . Otherwise  $\mathcal{V}$  is said to be UN disclosing for  $p$ .*

This non-disclosure function is clearly determinacy-compatible, so Proposition 1 will apply to it. Thus we will be able to utilize UN non-disclosure as a means of showing that certain  $d$ -views are *not* minimally useful.

**Variations: negative disclosure, background knowledge, and finite instances.** All of these notions can be additionally parameterized by background knowledge  $\Sigma$ , consisting of integrity constraints in some logic. Given  $d$ -instance  $\mathcal{D}$  satisfying  $\Sigma$  and a  $d$ -view  $\mathcal{V}$ ,  $\mathcal{V}$  determines a query  $Q$  over the  $d$ -schema at  $\mathcal{D}$  relative to  $\Sigma$  if: for every  $\mathcal{D}'$  satisfying  $\Sigma$  that

is  $\mathcal{V}$ -indistinguishable from  $\mathcal{D}$ ,  $Q(\mathcal{D}') = Q(\mathcal{D})$ . We say that  $\mathcal{V}$  is *useful* for a query  $Q$  relative to  $\Sigma$  if it determines  $Q$  on every  $d$ -instance satisfying  $\Sigma$ . We say  $\mathcal{V}$  is UN non-disclosing for query  $p$  with respect to  $\Sigma$  if  $\mathcal{V}$  does not determine  $p$  on any  $d$ -instance satisfying  $\Sigma$ .

By default, when we say “every instance”, we mean all instances, finite or infinite. There are variations of this problem requiring the quantification in both non-disclosure and utility to be over finite instances.

**Main problem.** We focus on the problem of determining whether minimally useful  $d$ -views exist for a given query  $Q$  and class  $\mathcal{C}$ , and characterizing such views when they do exist. When minimal useful  $d$ -views do not exist, we consider the problem of obtaining a  $d$ -view that is useful for  $Q$  and which minimizes the set of secrets  $p$  that are UN disclosed for  $p$ . We refer to  $Q$  as the *utility query*, and  $p$  as the *secret query*.

**Discussion.** Our notion of utility of views is *information-theoretic and exact*: a view is useful if a party with access to the view can compute the exact output of the query (as opposed to the correct output with high probability), with no limit on how difficult the computation may be. The generality of this notion will make our negative results stronger. And it turns out the generality will not limit our positive results, since these will be realized by very simply views.

Our notion of minimal usefulness is likewise natural if one seeks an ordering on sets of views measuring the ability to support exact information-theoretic query answering. Our query-based inexpressiveness notion, non-disclosure, gives a way of seeing the impact of minimal useful views on protecting information, and it is also based on information-theoretic and exact notions. We exemplify our notion of non-disclosure function with UN non-disclosure, which has been studied in prior work under several different names (Benedikt et al. 2016; Benedikt, Cuenca Grau, and Kostylev 2018; Benedikt et al. 2019). We choose the name “non-disclosing” rather than “private” for all our query-based expressiveness restrictions, since they are clearly very different from more traditional probabilistic privacy guarantees (Dwork and Roth 2014). On the one hand the UN non-disclosure guarantee is weak in that  $p$  is considered safe for  $\mathcal{V}$  (UN non-disclosed) if an attacker can never infer that  $p$  is true with absolutely certainty. Given a distribution on source instances, the information in the views may still increase the likelihood that  $p$  holds. On the other hand, the notion is quite strong in that it must hold on *every* source instance. Thus, although we do not claim that this captures all intuitively desirable properties of privacy, we do feel that it allows us to explore the ability to create views that simultaneously support the strong ability to answer certain queries in data integration and the strong inability to answer other queries.

**Restrictions and simplifications.** Although the utility and non-disclosure definitions make sense for any

queries, in this paper we will assume that  $Q$  and  $p$  are BCQs without constants (abusing notation by dropping “without constants”). While we restrict to the case of a single utility query and secret query here, all of our results have easy analogs for a finite set of such queries.

### 2.3 Some tools

Throughout the paper we rely on two basic tools.

**Canonical views.** Recall that we are looking for views that are useful for answering a CQ  $Q$  over a d-schema. The “obvious” set of views to try are those obtained by partitioning the atoms of  $Q$  among sources, with the free variables of the views including the free variables of  $Q$  and the variables occurring in atoms from different sources.

Given a CQ  $Q$  over a d-schema, and a source  $s$ , we denote by  $SVars(s, Q)$  the variables of  $Q$  that appear in an atom from source  $s$ . We also denote by  $SJVars(s, Q)$ , the “source-join variables of  $s$ ” in  $Q$ : the variables that are in  $SVars(s, Q)$  and which also occur in an atom of another source.

**Definition 6.** The canonical view of  $Q$  for source  $s$ ,  $CanView^s(Q)$ , has a view definition formed by conjoining all  $s$ -source atoms in  $Q$  and then existentially quantifying all bound variables of  $Q$  in  $SVars(s, Q) \setminus SJVars(s, Q)$ . The canonical d-view of  $Q$  is formed by taking the canonical view for each source.

In Example 1, the canonical d-view is what we referred to as “the obvious view design”. It would mean that one source has a view exposing  $\exists tdate \text{ Trtmnt}(pid, tinfo)$  – since  $pid$  is a source-join variable while  $tinfo$  is a free variable of  $Q$ . The other source should expose a view revealing  $\exists address. \text{Patient}(pid, age, address)$ , since  $address$  is neither a free variable nor shared across sources.

**The critical instance.** In the definition of UN non-disclosure of a set of views, we required that on any instance of the sources, a user who has access to the views cannot reconstruct the answer to the query  $p$ . For a certain class of views, one need only consider a single “most problematic” instance.

**Definition 7.** The critical instance of a schema  $S$  is the instance whose active domain consists of a single element  $*$  and whose facts are  $R(*, \dots, *)$  for all relational names  $R$  in  $S$ .

Note that every BCQ over the relevant relations holds on the critical instance of the source. The critical instance is the hardest instance for UN non-disclosure in the following sense:

**Theorem 1.** (Benedikt et al. 2016; Benedikt, Cuenca Grau, and Kostylev 2018) Consider any CQ views  $\mathcal{V}$ , and any BCQ  $p$ . If  $p$  is determined by  $\mathcal{V}$  at the critical instance of the source schema then it is determined by  $\mathcal{V}$  at each instance.

### 3 CQ views

Returning to Example 1, recall the intuition that the canonical d-view of  $Q$  are the “least informative views” that support the ability to answer  $Q$ . We start our analysis by proving such a result, but with two restrictions:  $Q$  must be a minimal CQ, and we only consider views specified by CQs:

**Theorem 2.** [Minimally useful CQ views] If CQ-based d-view  $\mathcal{V}$  determines a minimal Boolean CQ  $Q$ , then  $\mathcal{V}$  determines each canonical view  $CanView^s(Q)$  of  $Q$ .

We only sketch the proof here. The first step is to show that the determinacy of a CQ  $Q$  by a CQ-based d-view  $\mathcal{V}$  leads to a certain homomorphism of  $Q$  to itself. This step follows using a characterization of determinacy of CQ queries by CQ views via the well-known chase procedure (See (Benedikt, ten Cate, and Tsamoura 2016) or the “Green-Red chase” in (Gogacz and Marcinkowski 2015)). The second step is to argue that if this homomorphism is a bijection, then it implies that the canonical d-view determines  $Q$ , while if the homomorphism is not bijective, we get a contradiction of minimality. This step relies on an analysis of the chase characterization of determinacy factorizes over a d-schema.

**Consequences.** Combining Theorem 2 and Proposition 1 gives a partial answer to the question of how to obtain useful and non-disclosing views:

**Corollary 1.** For any non-disclosure function  $F$ , BCQs  $Q$  and  $p$ , if there is a CQ based d-view that is useful for  $Q$  and non-disclosing for  $p$  according to  $F$ , then the canonical d-view of  $Q^{\min}$  is such a d-view, where  $Q^{\min}$  is any minimal CQ equivalent to  $Q$ .

If we consider the specific non-disclosure notion, UN non-disclosure, we can infer a complexity bound from combining these results with prior work on the complexity of checking non-disclosure (Benedikt, Cuenca Grau, and Kostylev 2018):

**Corollary 2.** There is a  $\Sigma_2^P$  algorithm taking as input BCQs  $Q$  and  $p$  and determines whether there is a CQ-based d-view that is useful for  $Q$  and UN non-disclosing for  $p$ . If  $Q$  is assumed minimal then the algorithm can be taken to be in CONP.

The following example shows that the requirement that  $Q$  is minimal in Theorem 2 is essential.

**Example 2.** Consider two sources. The first source comprises the relations  $R_1$ ,  $R_2$  and  $R_3$ , while the second source comprises a single relation  $T$ . Consider also the conjunctions of atoms  $C_1$  and  $C_2$  defined as:

$$\begin{aligned} C_1 &= R_1(x, y) \wedge R_1(y, z) \wedge R_1(z, x) \wedge \\ &\quad T(x) \wedge R_2(w) \wedge R_3(q) \\ C_2 &= R_1(x', y') \wedge R_1(y', z') \wedge R_1(z', x') \wedge \\ &\quad T(x') \wedge R_2(y') \wedge R_3(z') \end{aligned}$$

The conjunct  $C_1$  states that there is  $R_1$ -triangle with one vertex in  $T$ , and in addition  $R_2$  and  $R_3$  are non-empty. The second conjunct states that there is an  $R_1$ -triangle with one vertex in  $T$ , a second in  $R_2$ , and a third in  $R_3$ . Consider now the query  $Q$  defined as

$$\exists x, y, z, x', y', z', q, w. C_1 \wedge C_2$$

Note that the conjunction of atoms  $C_1$  in  $Q$  is redundant. Indeed, the query  $Q^{\min} = \exists x', y', z'. C_2$  is equivalent to  $Q$ .

The canonical view of  $Q$  for the  $R_i$ 's source is:

$$\begin{aligned} \exists y, z, y', z', w, q. & R_1(x, y) \wedge R_1(y, z) \wedge R_1(z, x) \wedge \\ & R_1(x', y') \wedge R_1(y', z') \wedge R_1(z', x') \wedge \\ & R_2(w) \wedge R_3(q) \wedge R_2(y') \wedge R_3(z') \end{aligned}$$

In contrast, the canonical view of  $Q^{\min}$  for this source is

$$\begin{aligned} \exists y', z'. & R_1(x', y') \wedge R_1(y', z') \wedge R_1(z', x') \wedge \\ & R_2(y') \wedge R_3(z') \end{aligned}$$

Theorem 2 tells us that the canonical d-view of  $Q^{\min}$  is minimally useful within the class of CQ views. We claim that the canonical d-view of  $Q$  is not minimally useful for this class, and in fact reveals significantly more than the canonical d-view of  $Q^{\min}$ . Consider the secret query  $p = R_1(t, t)$  stating that there is an  $R_1$  self-loop. The canonical d-view of  $Q$  is UN non-disclosing for  $p$ . Indeed, given any instance  $\mathcal{D}$ , consider the instance  $\mathcal{D}'$  in which the  $T$  source is identical to the one in  $\mathcal{D}$ , but the  $R$  source is replaced by one where each node  $e$  in the canonical view for  $\mathcal{D}$  is in a triangle with distinct elements for  $y', z'$ . Such a  $\mathcal{D}'$  does not satisfy  $p$ , and is indistinguishable from  $\mathcal{D}$  according to the canonical d-view of  $Q$ .

In contrast, the canonical d-view is UN disclosing for  $p$ , since on the critical instance the returned bindings have  $x$  as only  $*$ , and from this we can infer that the witness elements for  $y$  and  $z$  can only be  $*$ . By Proposition 1 the canonical d-view of  $Q$  can not be minimally useful within the class of CQ views.

## 4 Arbitrary views

In the previous section we showed that the canonical d-view is minimally useful within the class of CQ views, assuming that the utility query is minimized. We now turn to minimally informative useful views, not restricting to views given by CQ view definitions.

Our goal will be to arrive at a generalization of the notion of canonical d-view that give the minimal information over arbitrary d-views that are useful for a given BCQ  $Q$ . That is we want to arrive at an analog of Theorem 2 replacing “CQ views” by “arbitrary views” and “canonical view” by a generalization.

### 4.1 An equivalence class representation of minimal useful views

Recall that general views are defined by queries, where a query can be any function on instances. An *Equivalence Class Representation of a d-view* (ECR) consists of an equivalence relation  $\equiv_s$  for each source  $s$ . An ECR is just another way of looking at a d-view defined by a set of arbitrary functions on instances: given a function  $F$ , one can define an equivalence relation by identifying two local instances when the values of  $F$  are the same. Conversely, given an equivalence relation then one can define a function mapping each instance to its equivalence class. A d-instance  $\mathcal{D}$  is indistinguishable from a d-instance  $\mathcal{D}'$  by the d-view specified by ECR  $\langle \equiv_s : s \in \text{Srcs} \rangle$  exactly when  $\mathcal{D}_s \equiv_s \mathcal{D}'_s$  holds for each  $s$ . Determinacy of one d-view by another corresponds exactly to the refinement relation between the corresponding ECRs. We will thus abuse notation by talking about indistinguishability, usefulness, and minimal usefulness of an ECR, referring to the corresponding d-view.

Our first step will be to show that there is an easy-to-define ECR whose corresponding d-view is minimally useful. For a source  $s$ , an  $s$ -context is an instance for each source other than  $s$ . Given an  $s$ -context  $C$  and an  $s$ -instance  $\mathcal{I}$ , we use  $(\mathcal{I}, C)$  to denote the d-instance formed by interpreting the  $s$ -relations as in  $\mathcal{I}$  and the others as in  $C$ .

We say two  $s$ -instances  $\mathcal{I}, \mathcal{I}'$  are  $(s, Q)$ -equivalent if for any  $s$ -context  $C$ ,  $(\mathcal{I}, C) \models Q \Leftrightarrow (\mathcal{I}', C) \models Q$ . We say two d-instances  $\mathcal{D}$  and  $\mathcal{D}'$  are globally  $Q$ -equivalent if for each source  $s$ , the restrictions of  $\mathcal{D}$  and  $\mathcal{D}'$  over source  $s$ , are  $(s, Q)$ -equivalent.

Global  $Q$ -equivalence is clearly an ECR. Via “swapping one component at a time” we can see that the corresponding d-views is useful for  $Q$ . It is also not difficult to see that this d-view is minimally useful within the class of all d-views:

**Proposition 2.** *The d-view corresponding to global  $Q$ -equivalence is minimally useful for  $Q$  within the collection of all d-views.*

Note that the result can be seen as an analog of Theorem 2. From it we conclude an analog of Corollary 1:

**Proposition 3.** *If there is any d-view that is useful for BCQ  $Q$  and non-disclosing for BCQ  $p$ , then the d-view given by global  $Q$ -equivalence is useful for  $Q$  and non-disclosing for  $p$ .*

**From an ECR to a concrete d-view.** We now have a d-view that is minimally useful within the set of all d-views, but it is given only as the ECR global  $Q$ -equivalence, and it is not clear that there are any views in the usual sense – isomorphism-invariant functions mapping into relations of some fixed schema — that correspond to this ECR. Our next goal is to show that

global  $Q$ -equivalence is induced by a d-view defined using standard database queries.

A *shuffle* of a CQ is a mapping from its free variables to themselves (not necessarily injective). Given a CQ  $Q$  and a shuffle  $\mu$ , we denote by  $\mu(Q)$  the CQ that results after replacing each variable occurring in  $Q$  by its  $\mu$ -image. We call  $\mu(Q)$  a *shuffled query*. For example, consider the query  $\exists y. R(x_1, x_2, x_2, y) \wedge S(x_2, x_3, x_3, y)$ . Then, the query  $\exists y. R(x_2, x_1, x_1, y) \wedge S(x_1, x_1, x_1, y)$  is a shuffle of  $Q$ .

The *canonical context query for  $Q$  at source  $s$* ,  $\text{CanCtxt}^s(Q)$ , is the CQ whose atoms are all the atoms of  $Q$  that are *not* in source  $s$ , and whose free variables are  $\text{SJVars}(s, Q)$ .

**Definition 8.** For a source  $s$ , a BCQ  $Q$  and a variable binding  $\sigma$  for  $Q$ , a shuffle  $\mu$  of  $\text{CanView}^s(Q)$  is invariant relative to  $\langle \sigma, \text{CanCtxt}^s(Q) \rangle$  if for any d-instance  $\mathcal{I}'$  where  $\mathcal{I}', \sigma \models \text{CanCtxt}^s(Q)$ , we have  $\mathcal{I}', \sigma \models \mu(\text{CanCtxt}^s(Q))$ .

Note that we can verify this invariance by finding a homomorphism from  $\sigma(\mu(\text{CanCtxt}^s(Q)))$  to  $\sigma(\text{CanCtxt}^s(Q))$ .

Invariance talks about every binding  $\sigma$ . We would like to abstract to bindings satisfying a set of equalities. A *type* for  $\text{SJVars}(s, Q)$  is a set of equalities between variables in  $\text{SJVars}(s, Q)$ . The notion of a variable binding satisfying a type is the standard one. For a type  $\tau$ , we can talk about a mapping  $\mu$  being *invariant relative to  $\langle \tau, \text{CanCtxt}^s(Q) \rangle$* : the invariance condition holds for all bindings  $\sigma$  satisfying  $\tau$ .

For a source  $s$  and a CQ  $Q$ , let  $\tau_1, \dots, \tau_n$  be all the equality types over the variables in  $\text{SJVars}(s)$  and  $\mathbf{x}$  be the variables in  $\text{SJVars}(s)$ .

**Definition 9.** The invariant shuffle views of  $Q$  for source  $s$  is the set of views  $V_{\tau_1}, \dots, V_{\tau_n}$  where each  $V_{\tau_i}$  is defined as  $\tau_i(\mathbf{x}) \wedge \bigvee_{\mu} \mu(\text{CanView}^s(Q))$ , where  $\mathbf{x}$  are the source-join variables of  $Q$  for source  $s$ , and where the disjunction is over shuffles invariant relative to  $\tau_i$ .

Note that since the domain of  $\mu$  is finite, there are only finitely many mappings on them, and thus there are finitely many disjuncts in each view up to equivalence.

We can show that global  $Q$  equivalence corresponds to agreement on these views:

**Proposition 4.** For any BCQ  $Q$  and any source  $s$ , two  $s$ -instances are  $(s, Q)$ -equivalent if and only if they agree on each invariant shuffle view of  $Q$  for  $s$ .

Putting together Proposition 3 and 4, we obtain:

**Theorem 3.** [Views that are minimally useful among all views] The invariant shuffle views are minimally useful for  $Q$  within the class of all views.

This yields a corollary for non-disclosure analogous to Corollary 1:

**Corollary 3.** If an arbitrary d-view  $\mathcal{V}$  is useful for BCQ  $Q$  and non-disclosing for BCQ  $p$ , then the d-view containing, for each source  $s$ , the invariant shuffle views of  $Q$  for  $s$ , is useful and non-disclosing. In particular, some DCQ is useful for  $Q$  and non-disclosing for  $p$ .

Notice that the invariant shuffle views can be *unsafe*: the variables that are contained in a relational atom within the distinct disjuncts may be distinct. Of course, they can be implemented easily by using a wildcard to represent elements outside the active domain. However, we can convert unsafe views to an “information-equivalent” set of safe relational algebra views.

**Proposition 5.** For every view defined by a DCQ (possibly unsafe), there is a finite set of relational algebra-based views  $\mathcal{V}'$  that induces the same ECR. Applying this to the invariant shuffle views for a CQ  $Q$ , we can find a relation algebra-based d-view that is minimally useful for  $Q$  within the class of all views.

The intuition behind the proposition is to construct separate views for different subsets of the variables that occur as a CQ disjunct. A view with a given set of variables  $S$  will assert that some CQ disjunct with variables  $S$  holds and that no disjunct corresponding to a subset of  $S$  holds.

**Example 3.** Consider a d-schema with two sources, one containing a ternary relation  $R$  and the other containing a unary relation  $S$ . Consider the utility query  $Q$ :

$$\exists x_1, x_2, y. R(x_1, x_2, y) \wedge R(y, x_2, x_1) \wedge S(x_1) \wedge S(x_2)$$

The canonical view for the  $R$  source  $\text{CanView}^R(Q)$  is  $\exists y. R(x_1, x_2, y) \wedge R(y, x_2, x_1)$ .

Observe that  $Q$  is a minimal CQ, and hence by Theorem 2 the canonical d-view of  $Q$  is minimally useful among the CQ-based d-views. We will argue that this d-view is not minimally useful for  $Q$  among all d-views, by arguing that it discloses more secrets than the shuffle views disclose.

Consider the secret query  $p = \exists x. R(x, x, x)$ . We show that canonical view is UN disclosing for  $p$ . Consider the critical instance of the  $R$ -source. An external party will know that the instance contains  $\{R(*, *, y_0), R(y_0, *, *)\}$  for some  $y_0$ . On the other hand, if  $y_0 \neq *$ , then the d-view would reveal  $x_1 = y_0, x_2 = *$ . So  $y_0$  must be  $*$ , and therefore  $p$  is disclosed. By Corollary 1, we know that no CQ-based d-views can be UN non-disclosing for  $p$  and useful for  $Q$ .

The shuffle views of  $Q$  are always useful for  $Q$ . We will show that they are UN non-disclosing for  $p$ . Let us start by deriving the invariant shuffle views for the  $R$  source. There are two types,  $\tau_1$  in which  $x_1 = x_2$ , and  $\tau_2$  in which the variables are not identified.

For a binding satisfying  $\tau_1$ , the canonical view of  $Q$  for the  $R$  source is equivalent to

$$\exists y. R(x_1, x_1, y) \wedge R(y, x_1, x_1)$$

Since there is only one free variable in it, there is only one invariant shuffle, the identity. Thus

$$V_{\tau_1} = \exists y.R(x_1, x_1, y) \wedge R(y, x_1, x_1)$$

For bindings satisfying  $\tau_2$  there are several shuffles invariant for  $\text{CanCtxt}^R(Q) = S(x_1) \wedge S(x_2)$ : the identity, the shuffle which swaps  $x_1$  and  $x_2$ , the shuffle in which  $x_1$  and  $x_2$  both go to  $x_1$ , and the shuffle in which both  $x_1$  and  $x_2$  go to  $x_2$ . Thus we get the view  $V_{\tau_2}$  defined as  $x_1 \neq x_2$  conjoined with:

$$\begin{aligned} & \exists y.R(x_1, x_2, y) \wedge R(y, x_2, x_1) \vee \\ & \exists y.R(x_1, x_1, y) \wedge R(y, x_1, x_1) \vee \\ & \exists y.R(x_2, x_2, y) \wedge R(y, x_2, x_2) \vee \\ & \exists y.R(x_2, x_1, y) \wedge R(y, x_1, x_2) \end{aligned}$$

This last view is unsafe, but via Proposition 5 we can convert it into a safe relational algebra view that yields the same ECR,  $V_{\tau_2}^{\text{safe}}$  defined as  $x_1 \neq x_2$  conjoined with:

$$\begin{aligned} & \neg(\exists y.R(x_1, x_1, y) \wedge R(y, x_1, x_1)) \wedge \\ & \neg(\exists y.R(x_2, x_2, y) \wedge R(y, x_2, x_2)) \wedge \\ & (\exists y.R(x_1, x_2, y) \wedge R(y, x_2, x_1) \vee \\ & \exists y.R(x_2, x_1, y) \wedge R(y, x_1, x_2)) \end{aligned}$$

We now argue that the shuffle views are UN non-disclosing for  $p$ . This is because in any  $d$ -instance we can replace each fact  $R(x_0, x_0, x_0)$  by facts  $R(x_0, x_0, c)$  and  $R(c, x_0, x_0)$  for a fresh  $c$ , obtaining an indistinguishable instance where  $p$  does not hold. Hence by Proposition 1, the canonical views of  $Q$  can not be minimally useful within the class of all views, or even within the class of relational algebra views.

## 5 Background knowledge

We now look at the impact of a background knowledge on the sources. We start with the case of a background theory  $\Sigma$  in which each sentence is *local*, referencing relations on a single source.

### 5.1 Extension of results on CQ views to local background knowledge

It is easy to show that we can not generalize the prior results for CQ views to arbitrary local background knowledge. Intuitively using such knowledge we can encode design problems for arbitrary views using CQ views; see the appendix for details. Thus we restrict to *local constraints*  $\Sigma$  that are *existential rules*. We show that the results on CQ views extend to this setting. We must now consider utility queries  $Q$  that are *minimal with respect to*  $\Sigma$ , meaning that there is no strict subquery equivalent to  $Q$  under  $\Sigma$ . By modifying the chase-based approach used to prove Theorem 2, we show:

**Theorem 4.** [Minimally useful CQ views w.r.t. local rules] Let  $\Sigma$  be a set of local existential rules,  $Q$  a CQ

minimal with respect to  $\Sigma$ . Then the canonical  $d$ -view of  $Q$  is minimally useful within the class of CQ views relative to  $\Sigma$ .

The analog of Corollary 1 follows from the theorem:

**Corollary 4.** If any CQ based  $d$ -view is useful for  $\Sigma$ -minimal  $Q$  and non-disclosing for  $p$  relative to  $\Sigma$ , then the canonical  $d$ -views of  $Q$  is useful for  $Q$  and non-disclosing for  $p$ .

**Consequences for decidability.** Theorem 4 shows that even in the presence of arbitrary local existential rules  $\Sigma$  it suffices to minimize the utility under  $\Sigma$  and check the canonical  $d$ -view for non-disclosure under  $\Sigma$ . For arbitrary existential rules, even CQ minimization is undecidable. But for well-behaved classes of rules (e.g. those with terminating chase (Cuenca Grau et al. 2013; Baget et al. 2014), or frontier-guarded rules (Baget et al. 2011)) we can perform both minimization and UN non-disclosure checking effectively (Benedikt et al. 2019).

### 5.2 Extension of results on arbitrary views to the presence of local constraints

We can also revise the results on arbitrary  $d$ -views to account for local existential rules  $\Sigma$ . The notion of a shuffle being  $\Sigma$ -invariant is defined in the obvious way, restricting to instances that satisfy the constraints in  $\Sigma$ . The  $\Sigma$ -invariant shuffle views are also defined analogously; they are DCQ views, but can be replaced by the appropriate relational algebra views. Following the prior template, we can show:

**Theorem 5.** [Minimally useful views w.r.t. local existential rules] For any set of local existential rules  $\Sigma$ , the  $\Sigma$ -invariant shuffle views of  $Q$  are minimally useful within the class of all views, relative to  $\Sigma$ .

The result has effective consequences for “tame” rules (e.g. with terminating chase) with no non-trivial invariant shuffles. In such cases, the  $\Sigma$ -invariant shuffle views degenerate to the canonical  $d$ -view, and we can check whether the canonical  $d$ -view is non-disclosing effectively (Benedikt et al. 2016; 2019).

### 5.3 Non-local background knowledge

The simplest kind of non-local constraint is the replication of a table between sources. Unlike local constraints, these require some communication among the sources to enforce. Thus we can consider a replication constraint to be a restricted form of source-to-source communication.

We will see that several new phenomena arise in the presence of these constraints. Recall that with only local constraints, we have useful  $d$ -views with minimal information. We can not guarantee the existence of such views in the presence of replication:

**Proposition 6.** There is a schema with replication constraint  $\Sigma$  and a BCQ  $Q$  where no  $d$ -view is minimally useful for  $Q$  within the class of all views w.r.t.  $\Sigma$ .



Our proof of Proposition 6 uses a schema with unary relations  $R, S, T$ . There are two sources:  $R$  and  $T$  are in different sources, and  $S$  is replicated between the two sources. Let  $Q$  be equal to  $\exists x. R(x) \wedge S(x) \wedge T(x)$ .

Identify the active domain of our instances with integers, and consider the functions  $F_1(x) = x + 1$  and  $F_2(x) = 2 - x$ . Let  $\text{Str}_1$  be the function that applies  $F_1$  to a relation element-wise, and similarly define  $\text{Str}_2$  using  $F_2$ . Notice that  $F_1$  maps 0 to 1 and 1 to 2, while  $F_2$  maps 0 to 2 and 1 to 1.

We define an ECR  $\mathcal{V}_1$  on each source by relating two instances  $\mathcal{I}$  and  $\mathcal{I}'$  of the source exactly when  $\mathcal{I}'$  can be obtained from applying  $\text{Str}_1$  on  $\mathcal{I}$  some number of times (applying it to both relations of the source) or vice versa. That is  $\mathcal{V}_1$  is the equivalence closure of  $\text{Str}_1$ . Let  $\mathcal{V}_2$  be defined analogously using  $\text{Str}_2$ . We argue that both  $\mathcal{V}_1$  and  $\mathcal{V}_2$  are useful. We focus on  $\mathcal{V}_1$  since the proof for  $\mathcal{V}_2$  is analogous. Usefulness follows from the following claim. If we have two valid d-instances  $(\mathcal{I}_1, \mathcal{I}_2)$  and  $(\mathcal{I}'_1, \mathcal{I}'_2)$  with all local instances non-empty then 1. We cannot have  $\mathcal{I}'_1 = \text{Str}_1^i(\mathcal{I}_1)$ ,  $\mathcal{I}'_2 = \text{Str}_1^j(\mathcal{I}_2)$  with  $i \neq j$ ; and similarly for  $\text{Str}_2$ . 2. We cannot have  $\mathcal{I}'_1 = \text{Str}_1^i(\mathcal{I}_1)$  and  $\mathcal{I}_2 = \text{Str}_2^j(\mathcal{I}'_2)$ ; and similarly for  $\text{Str}_2$ .

Now suppose  $\mathcal{V}$  were a set of views minimally useful for  $Q$ . We must have  $\mathcal{V}_1$  and  $\mathcal{V}_2$  determine  $\mathcal{V}$ . Thus in particular if we have two local instances that agree on either  $\mathcal{V}_1$  or  $\mathcal{V}_2$ , then they agree on  $\mathcal{V}$ .

Consider a d-instance  $\mathcal{D}$  with  $R = \{0\}$ ,  $S = \{0, 1\}$ ,  $T = \{0\}$ , and a d-instance  $\mathcal{D}'$  with  $R = \{1\}$ ,  $S = \{1, 2\}$ ,  $T = \{2\}$ . These are both valid instances (i.e. the replication constraint is respected). But  $\mathcal{D}'$  is obtained from  $\mathcal{D}$  by applying  $\text{Str}_1$  on the source with  $R$ , and by applying  $\text{Str}_2$  on the other source.

Thus  $\mathcal{D}'$  and  $\mathcal{D}$  are indistinguishable by  $\mathcal{V}$ , but  $Q$  has a match in  $\mathcal{D}$  but not in  $\mathcal{D}'$ . This contradicts the assumption that  $\mathcal{V}$  is useful for  $Q$ .

Given Proposition 6 we will focus on obtaining useful views that minimize the queries that are UN non-disclosed for the remainder of the subsection. We can use the technique in the proof of Proposition 6 to show a more promising new phenomenon: there may be d-views that are useful for CQ  $Q$  and UN non-disclosing for CQ  $p$ , but they are much more intricate than any query related to the canonical d-view of  $Q$ . In fact we can show that with a fully-replicated relation in the utility query we can get useful and UN non-disclosing views whenever this is not ruled out for trivial reasons:

**Theorem 6.** *If BCQ  $Q$  contains a relation of non-zero arity replicated across all sources then there is a d-view that is useful for  $Q$  and UN non-disclosing for BCQ  $p$  if and only if there is no homomorphism of  $p$  to  $Q$ .*

Further we can use the same d-view for every such  $p$  without a homomorphism into a given  $Q$ ! Thus even though we do not have minimally useful d-views, we have d-views that are optimal from the perspective of

UN non-disclosure and utility for a fixed  $Q$ .

We sketch the idea of the proof of Theorem 6. Given utility query  $Q$  and local instance  $\mathcal{I}$  we can form the “product instance” of  $\mathcal{I}$  and  $Q$ . The elements of a product instance will be pairs  $(x, c)$  where  $x$  is a variable of  $Q$  and  $c$  an element of  $\mathcal{I}$ , and there will be atom  $R((x_1, c_1), \dots, (x_n, c_n))$  in the product exactly when there are corresponding atoms in  $\text{CanView}^s(Q)$  and  $\mathcal{I}$ . Thus we will have homomorphisms from this instance to both  $\mathcal{I}$  and  $\text{CanView}^s(Q)$ . We use ECRs that make an  $s$ -instance  $\mathcal{I}$  equivalent to all instances formed by iterating this product construction of  $\mathcal{I}$  with  $\text{CanView}^s(Q)$ . The d-views corresponding to these ECRs will be UN non-disclosing because in the product we will have a fresh copy of any partial match of  $Q$ , and so the only way for the secret query  $p$  to hold in the product will be if it has a homomorphism into  $Q$ , which is forbidden by hypothesis. The replication constraint ensures that if we have two d-instances that are equivalent, where the interpretation of the replicated relation is non-empty, the number of iterations of the product construction is the same on each source. Using this fact we can ensure that the d-views are useful.

The views used in Theorem 6 are not isomorphism-invariant: like the views from Proposition 6, the product construction can be seen as applying some value transformation on the elements of each instance. We can show — in sharp contrast to the situation with only local constraints — that with replication it may be essential to use d-views based on queries that are not isomorphism-invariant.

**Proposition 7.** *There is a d-schema with a replication constraint, along with BCQs  $Q$  and  $p$  such that there is a d-view useful for  $Q$  and UN non-disclosing for  $p$ , but there is no such d-view based on queries returning values in the active domain and commuting with isomorphisms.*

Proposition 7 shows that in the presence of replication, we can do much more with arbitrary d-views than we can even with relational algebra-based d-views.

## 6 Discussion and outlook

We have studied the ability to design views that satisfy diverse goals: expressiveness requirements in terms of full disclosure of a specified set of queries in the context of data integration, and inexpressiveness restrictions in terms of either minimal utility or minimizing disclosure of queries. Our main results characterize information-theoretically minimal views that support the querying of a given CQ.

We consider only a limited setting; e.g. CQs for the utility query. Our hope is that the work can serve as a basis for further exploration of the trade-offs in using query-based mechanisms in a variety of settings.

## References

- Abiteboul, S.; Hull, R.; and Vianu, V. 1995. *Foundations of Databases*. Addison-Wesley.
- Arenas, M.; Barceló, P.; and Reutter, J. L. 2011. Query languages for data exchange: Beyond unions of conjunctive queries. *Theory Comput. Syst.* 49(2):489–564.
- Baget, J.-F.; Mugnier, M.-L.; Rudolph, S.; and Thomazo, M. 2011. Walking the complexity lines for generalized guarded existential rules. In *IJCAI*.
- Baget, J.; Garreau, F.; Mugnier, M.; and Rocher, S. 2014. Extending acyclicity notions for existential rules. In *ECAI*.
- Benedikt, M.; Bourhis, P.; ten Cate, B.; and Puppis, G. 2016. Querying visible and invisible information. In *LICS*.
- Benedikt, M.; Bourhis, P.; Jachiet, L.; and Thomazo, M. 2019. Reasoning about disclosure in data integration in the presence of source constraints. In *IJCAI*.
- Benedikt, M.; Cuenca Grau, B.; and Kostylev, E. V. 2018. Logical foundations of information disclosure in ontology-based data integration. *Artif. Intell.* 262:52–95.
- Benedikt, M.; ten Cate, B.; and Tsamoura, E. 2016. Generating plans from proofs. In *TODS*.
- Bonatti, P. A., and Sauro, L. 2013. A confidentiality model for ontologies. In *ISWC*.
- Calvanese, D.; De Giacomo, G.; Lenzerini, M.; and Rosati, R. 2012. View-based Query Answering in Description Logics: Semantics and Complexity. *J. Comput. Syst. Sci.* 78(1):26–46.
- Chaum, D.; Crépeau, C.; and Damgård, I. 1988. Multi-party unconditionally secure protocols. In *STOC*.
- Cuenca Grau, B.; Horrocks, I.; Krötzsch, M.; Kupke, C.; Magka, D.; Motik, B.; and Wang, Z. 2013. Acyclicity notions for existential rules and their application to query answering in ontologies. *JAIR* 47:741–808.
- Dwork, C., and Roth, A. 2014. The algorithmic foundations of differential privacy. *Found. & Trends in Th. Comp. Sci.* 9(3&4):211–407.
- Dwork, C. 2006. Differential privacy. In *ICALP*.
- Fernández-Duque, D., and Goranko, V. 2016. Secure aggregation of distributed information. *Discrete Appl. Math.* 198(C):118–135.
- Gogacz, T., and Marcinkowski, J. 2015. The hunt for a red spider: Conjunctive query determinacy is undecidable. In *LICS*.
- Halevy, A. Y. 2001. Answering queries using views: A survey. *VLDB J.* 10(4):270–294.
- Nash, A., and Deutsch, A. 2007. Privacy in GLAV information integration. In *ICDT*.
- Nash, A.; Segoufin, L.; and Vianu, V. 2010. Views and queries: Determinacy and rewriting. *TODS* 35(3).

## APPENDIX

### Proofs for Section 2: basic properties of the definitions

#### Proof of Proposition 1

Recall the statement:

Suppose d-view  $\mathcal{V}$  is minimally useful for  $Q$  within  $\mathcal{C}$ , and there is a d-view based on queries in  $\mathcal{C}$  that is useful for  $Q$  and non-disclosing for  $p$  according to non-disclosure function  $F$ . Then  $\mathcal{V}$  is useful for  $Q$  and non-disclosing for  $p$  according to  $F$ .

*Proof.* By assumption  $\mathcal{V}$  is useful for  $Q$ . Suppose  $\mathcal{V}'$  based on views in  $\mathcal{C}$  is useful for  $Q$  and non-disclosing for  $p$  according to  $F$ . Since  $\mathcal{V}$  is minimally useful in  $\mathcal{C}$ ,  $\mathcal{V}'$  determines the view definition of each view in  $\mathcal{V}$ . If  $\mathcal{V}$  disclosed  $p$  according to  $F$ , then  $\mathcal{V}'$  would disclose  $p$  according to  $F$ , since  $F$  is a non-disclosure function, hence is determinacy-compatible.  $\square$

## Proofs for Section 3: balancing expressiveness and inexpressiveness inside the class of views defined by conjunctive queries

### Proof of Theorem 2

Recall the statement:

If CQ views  $\mathcal{V}$  determine a minimal Boolean CQ  $Q$ , then  $\mathcal{V}$  determines each canonical view  $\text{CanView}^s(Q)$  of  $Q$ .

We will make use of the following property of minimal queries, which is easy to verify:

**Lemma 1.** *If  $Q$  is minimal, then there does not exist a homomorphism  $h$  from  $Q$  into itself that maps two different variables occurring in  $Q$  to the same variable.*

We will begin the proof by showing that the determinacy of a CQ  $Q$  by a set of CQ views  $\mathcal{V}$  leads to the existence of a certain homomorphism of  $Q$  to itself.

Let  $Q$  be a Boolean conjunctive query and  $\mathcal{V}$  be an arbitrary set of conjunctive query views. We will need to review an algorithm for checking determinacy. We fix a signature for our queries and views, which we refer to as the *original signature*. From it we derive a *primed signature*, containing a relation  $R'$  for each  $R$  in the original signature. Given a formula  $\varphi$  in the original signature, we let  $\varphi'$  be formed by replacing every relation  $R$  in  $\varphi$  by  $R'$ . We use a similar notation for a set of facts  $S$  in the original signature. In particular, for a conjunctive query  $Q$  in the original signature,  $Q'$  refers to the conjunctive query obtained by replacing every relation symbol by its primed counterpart.

Given a view  $V(\mathbf{x})$  defined by conjunctive query  $\exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$ , the corresponding *forward view definition* (for  $V$ ) is the rule:

$$\varphi(\mathbf{x}, \mathbf{y}) \rightarrow V(\mathbf{x})$$

while the *inverse view definition* for  $V$  is the rule:

$$V(\mathbf{x}) \rightarrow \exists \mathbf{y} . \varphi(\mathbf{x}, \mathbf{y})$$

We let  $\text{ForView}(\mathcal{V})$  and  $\text{BackView}(\mathcal{V})$ , denote the forward and inverse view definitions for views in  $\mathcal{V}$ , and  $\text{ForView}'(\mathcal{V})$  and  $\text{BackView}'(\mathcal{V})$  the same axioms but for primed copies of the base predicates.

**Example 4.** *Suppose we have a set of views  $\mathcal{V}$  consisting only of view  $V(x)$  given by CQ  $\exists y, z. R(x, y) \wedge S(x, z)$ . Then  $\text{ForView}(\mathcal{V})$  contains the rule:*

$$R(x, y) \wedge S(x, z) \rightarrow V(x)$$

*while  $\text{BackView}'(\mathcal{V})$  contains the rule*

$$V(x) \rightarrow \exists y, z. R'(x, y) \wedge S'(x, z)$$

Letting  $\Sigma_{Q, \mathcal{V}}$  to be the axioms above, we can easily see that determinacy can be expressed as:

$$Q \wedge \Sigma_{Q, \mathcal{V}} \models Q'$$

This is a containment of CQs under existential rules, and Algorithm 1 checks this via the chase procedure, which is complete for such containments (see Section 2).

Intuitively, the algorithm iteratively chases with the axioms above in rounds, checking for a match of  $Q'$  after each round. The correctness is implicit in (Nash, Segoufin, and Vianu 2010); see also (Benedikt, ten Cate, and Tsamoura 2016; Gogacz and Marcinkowski 2015).

**Theorem 7.** *Algorithm 1 returns true if and only if  $\mathcal{V}$  determines  $Q$ .*

Given an arbitrary input query  $Q$ , we use a superscript to distinguish among the different sets computed during each iteration of Algorithm 1, e.g.,  $F_0^i(Q, \mathcal{V})$  denotes instance  $F_0$  during the  $i$ -th iteration of Algorithm 1 when run on  $Q$  with  $\mathcal{V}$ . For any  $Q, \mathcal{V}$  such that the algorithm returns **true**, and any  $j \in 0 \dots 5$ , let  $F_j^\infty(Q, \mathcal{V})$  denote the content of  $F_j(Q, \mathcal{V})$  at the point where the algorithm terminates; this will be  $F_j^l(Q, \mathcal{V})$  for some  $l$ . Note that in line 11 of Algorithm 1 we add to each  $F_0^l(Q, \mathcal{V})$  only the subset of ground atoms from  $F_4^l(Q, \mathcal{V})$  that belong to the original schema, since the other atoms will be re-derived in the next round.

For any fact  $F = R'(\mathbf{c})$  in the primed signature, let  $\text{UnPrime}(F)$  be the corresponding fact  $R(\mathbf{c})$ . For a set of facts  $S$ , we let

$$\text{UnPrime}(S) = \{\text{UnPrime}(F) : F \in S \text{ in the primed signature}\}$$

The following lemma states the easy fact that iterating the “chase and backchase” steps in lines 3 to 4 gives a homomorphic pre-image of what we started with:

Figure 1: Determinacy( $Q, \mathcal{V}$ )

```

1:  $F_0 := \text{canondb}(Q)$ 
2: while true do % Next 2 lines create view facts from  $F_0$  and then primed facts based on views
3:    $F_1 := \text{Chase}_{\text{ForView}(\mathcal{V})}(F_0)$ 
4:    $F_2 := \text{Chase}_{\text{BackView}(\mathcal{V})}(F_1)$ 
5:   if  $\exists h : Q' \rightarrow F_2$  mapping each free variable  $v$  of  $Q'$  into  $c_v \in \text{adom}(\text{canondb}(Q))$  then return true
6:   end if
7:    $F_3 := \text{Chase}_{\text{ForView}(\mathcal{V})}(F_2)$  % These 2 lines create view facts from primed facts and then original facts
   based on views
8:    $F_4 := \text{Chase}_{\text{BackView}(\mathcal{V})}(F_3)$ 
9:    $F_5 = \text{restrict } F_4 \text{ to original schema}$ 
10:  if  $F_0 \neq F_5$  then % no fixpoint, so repeat with expanded original facts
11:     $F_0 := F_0 \cup F_5,$ 
12:  else
13:    return false
14:  end if
15: end while

```

**Lemma 2.** For each  $l \geq 1$ , there is a homomorphism  $\mu_l : \text{UnPrime}(F_2^l(Q, \mathcal{V})) \rightarrow F_0^l(Q, \mathcal{V})$  that preserves all elements of facts in  $F_0^l(Q, \mathcal{V})$ .

*Proof.* Let  $\mathcal{V}_i \subseteq \mathcal{V}$  be the set of all views providing data from the  $i$ -th datasource and let  $V_{i,1} \dots V_{i,j_i}$  enumerate the view predicate names in  $\mathcal{V}_i$ , where  $V_{i,j}$  has definition  $\mathbf{B}_{i,j}$ . Consider the forward view definition, forward transfer as well as, for each  $V_{i,j}$ , the inverse view definition associated going from  $V_{i,j}$  to the primed relations. That is, the axioms:

$$\begin{aligned} \mathbf{B}_{i,j}(\mathbf{x}, \mathbf{y}) &\rightarrow V_{i,j}(\mathbf{x}) \\ V_{i,j}(\mathbf{x}) &\rightarrow \exists \mathbf{y} \mathbf{B}'_{i,j}(\mathbf{x}, \mathbf{y}) \end{aligned}$$

where  $\mathbf{B}_{i,j}(\mathbf{x})$  is some conjunction of atoms from the  $i$ -th datasource and  $\{\mathbf{y}\}$  is disjoint from  $\{\mathbf{x}\}$ .

Each null  $\lambda$  appearing in  $F_2^l(Q, \mathcal{V})$  is generated by a chase step associated with Line 4 of Algorithm 1, triggered by a unique fact  $V_{i,j}(\sigma)$ , where  $\sigma$  is a binding of  $\mathbf{x}$ , with  $\lambda$  corresponding to some variable  $y_\lambda$  in  $\mathbf{B}'_{i,j}(\mathbf{x}, \mathbf{y})$ . For each fact  $V_{i,j}(\sigma)$ , we know we can extend to a binding  $\sigma'$  for  $\mathbf{y}$  such that  $\mathbf{B}_{i,j}(\sigma')$  is in  $F_0^l(Q, \mathcal{V})$ . Our homomorphism will map  $\lambda$  to  $\sigma'(v_\lambda)$ .  $\square$

An analogous reasoning applies to lines 7 to 8:

**Lemma 3.** For any CQ  $Q$ , each  $l \geq 1$  there exists a homomorphism  $\mu_l : F_5^l(Q, \mathcal{V}) \rightarrow \text{UnPrime}(F_2^l(Q, \mathcal{V}))$  that preserves all elements occurring in  $\text{UnPrime}(F_2^l(Q, \mathcal{V}))$ .

Using the two lemmas above, we can prove by induction on  $l$ :

**Lemma 4.** For any CQ  $Q$ , CQ views  $\mathcal{V}$ , and number  $l$ , there exists a homomorphism  $\text{UnPrime}(F_2^l(Q, \mathcal{V})) \rightarrow \text{canondb}(Q)$  that is the identity on elements  $c_v$ . By priming each relation, we can get  $\mu : F_2^l(Q, \mathcal{V}) \rightarrow \text{canondb}(Q')$ .

The following simple proposition relates the sets produced by the algorithm on input  $Q$  with the sets produced by the same algorithm on one of the canonical views  $\text{CanView}^s(Q)$ .

**Lemma 5.** For any  $l \geq 1$  and any source  $s$ ,  $F_0^l(\text{CanView}^s(Q), \mathcal{V})$  is the same as the set of atoms in source  $s$  within  $F_0^l(Q, \mathcal{V})$ .

We are now ready for the proof of Theorem 2:

*Proof.* Assume that  $\mathcal{V}$  does not determine  $\text{CanView}^s(Q)$ , for some datasource  $s$ , and that  $\mathcal{V}$  determines  $Q$ . From the latter assumption and Theorem 7 we conclude that there exists a homomorphism  $h_1$  from  $Q'$  into  $F_2^\infty(Q, \mathcal{V})$ .

By composing  $h_1$  with the  $\nu$  produced by Lemma 4 we obtain a homomorphism  $h_Q$  from  $Q'$  into itself. By unpriming, we can change  $h_Q$  to a homomorphism from  $Q$  to itself, and we will sometimes abuse notation by considering  $h_Q$  in either way. Note that  $h_1$ , and also  $h_Q$ , maps each variable  $x$  in  $\text{SJVars}(s, Q)$  to some element of the form  $c_v$ ,

since each such  $x$  appears in atoms over distinct sources, and only elements of the form  $c_v$  appear in such atoms within  $F_2^\infty(Q, \mathcal{V})$ .

We consider two cases.

Case 1:  $h_Q$  is injective. Thus  $h_Q$  is both an injection on  $Q'$  and an injection from source-join variables to the constant corresponding to a source-join variable. Then  $(h_Q)^{-1}$  composed with  $h_1$  is an injective homomorphism from  $Q'$  into  $F_2^\infty(Q, \mathcal{V})$  that maps each  $x$  in  $\text{SJVars}(s, Q)$  to  $c_x$ . By Lemma 5 we can identify the source  $s$  atoms of  $F_2^\infty(Q, \mathcal{V})$  with the atoms of  $F_2^\infty(\text{CanView}^s(Q), \mathcal{V})$ . Applying this identification and abusing notation as described above we can see  $(h_Q)^{-1}$  as a homomorphism of  $Q$  into  $F_2^\infty(\text{CanView}^s(Q), \mathcal{V})$  that maps each  $v$  in  $\text{SJVars}(s, Q)$  to  $c_v$ . But then Algorithm 1 would have returned true for  $\text{CanView}^s(Q)$ . By Theorem 7, this is a contradiction of the fact that  $\mathcal{V}$  does not determine  $\text{CanView}^s(Q)$ .

Case 2:  $h_Q$  is not injective. But now, by Lemma 1, we have a contradiction of minimality.  $\square$

## Proof of Corollary ??, and further optimizations

Recall the statement of the corollary:

There is a  $\Sigma_2^P$  algorithm taking as input BCQs  $Q$  and  $p$  without constants, that determines whether there is a CQ based d-view that is useful for  $Q$  and UN non-disclosing for  $p$ . If  $Q$  is assumed minimal then the algorithm can be taken to be in CONP.

*Proof.* We know it suffices to find a minimal query  $Q^{\min}$  and check that the canonical views of  $Q^{\min}$  are UN non-disclosing for  $p$ . This means guessing a homomorphic image  $Q^{\min}$  of  $Q$  and checking that  $Q^{\min}$  is UN non-disclosing for  $p$ . (Benedikt, Cuenca Grau, and Kostylev 2018) showed that there is a CONP procedure that takes as input views and a query  $p$ , and checks whether or not the views are UN non-disclosing for  $p$ . Composing the guess with that check gives the  $\Sigma_2^P$  bounds. When  $Q$  is minimal, we can just use the CONP check of (Benedikt, Cuenca Grau, and Kostylev 2018).  $\square$

In fact, we can do slightly better. The canonical views are a way of decomposing a utility query. But it turns out that we can also apply the same construction underlying the canonical views as a way of decomposing a secret query. We start with the following observation, which states that in analyzing secrecy of a set of CQ views, we can break up the secret query into its canonical views and analyze them one at a time:

**Proposition 8.** *For any Boolean CQ  $p$  and CQ-based d-view  $\mathcal{V}$ ,  $\mathcal{V}$  is UN non-disclosing for  $p$  if and only if  $\mathcal{V}$  is UN non-disclosing for  $\text{CanView}^s(p)$  for some  $s$ .*

*Proof.* In one direction, suppose there is  $s$  such that  $\mathcal{V}$  is UN non-disclosing for  $\text{CanView}^s(p)$ . We show that  $\mathcal{V}$  is UN non-disclosing for  $p$  by showing that the critical instance  $\mathcal{D}_0$  for the d-schema is  $\mathcal{V}$ -indistinguishable from some other d-instance where  $p$  fails. We know that there is some  $s$ -instance  $\mathcal{I}_s$  that is  $\mathcal{V}^s$  indistinguishable from  $\mathcal{D}_s$  but in which  $\text{CanView}^s(p)(*)$  does not hold. Let  $\mathcal{D}'$  be formed from taking  $\mathcal{I}_s$  on source  $s$  and taking the critical instance on the other sources, choosing the other components arbitrarily. Clearly  $\mathcal{D}'$  is  $\mathcal{V}$ -indistinguishable from  $\mathcal{D}$ . If  $p$  held on  $\mathcal{D}'$ , the only possible witness would be the critical tuple, since this is the only binding. But the critical tuple fails the conjuncts on source  $s$ .

In the other direction, suppose  $\mathcal{V}$  is UN disclosing for  $\text{CanView}^s(p)$  for some  $s$ . We will show that  $\mathcal{V}$  is UN disclosing for  $p$ . We know that for each  $s$ , letting  $\mathcal{I}_s$  be the critical instance for source  $s$ , if we apply the disjunctive chase procedure from (Benedikt et al. 2019; 2016) to  $\mathcal{I}_s$  then  $\text{CanView}^s(* \dots *)$  holds. But when we apply the disjunctive chase procedure to the critical instance for the d-schema, this is the same as applying it to each component. Thus  $p$  holds with  $* \dots *$  as a witness, so  $\mathcal{V}$  is UN disclosing for  $p$ .  $\square$

The result above is about a fixed set of CQ views. But using Corollary 1 we can lift it to an observation about decomposing the secret query in searching for the existence of useful and UN non-disclosing views:

**Proposition 9.** *For any Boolean CQs  $Q$  and  $p$ , there is a CQ-based d-view that is useful for  $Q$  and UN non-disclosing for  $p$  if and only if for some  $s$ , the canonical views of  $Q$  are UN non-disclosing for  $\text{CanView}^s(p)$ .*

*Proof.* If there are CQ views that are useful for  $Q$  and UN non-disclosing for  $p$ , then the canonical views of  $Q$  are such a set of views, by Corollary 1. Thus by the previous proposition, they are UN non-disclosing for some  $\text{CanView}^s(p)$ .

In the other direction, if the canonical views of  $Q$  are UN non-disclosing for  $\text{CanView}^s(p)$  for some  $s$ , then by the proposition above they are UN non-disclosing for  $p$ , and thus these views serve as a witness.  $\square$

Proposition 9 implies that to test for useful and UN non-disclosing views, we need only take each source  $s$  and test whether  $\text{CanView}^s(Q)$  for source  $s$  is UN non-disclosing for the canonical view of  $\text{CanView}^s(p)$  for source  $s$ . A witness to failure of such a test requires first a deterministic computation that consists of chasing forward and backward with the view definitions, then a series of guesses of homomorphism of  $\text{CanView}^s(Q)$ , followed by the guess of a homomorphism of  $\text{CanView}^s(p)$ , thus CONP in the maximum cardinality over all  $s$  of  $\text{CanView}^s(p)$  and  $\text{CanView}^s(Q)$ .

## Proofs for Section 4: balancing expressiveness and inexpressiveness within the class of arbitrary views

### Proof of Proposition 2

Recall the statement:

The set of views corresponding to global  $Q$ -equivalence is minimally useful for  $Q$  within the collection of all views.

*Proof.* To see that  $Q$ -equivalence is useful, consider two d-instances  $\mathcal{D}$  and  $\mathcal{D}'$  that are globally  $Q$ -equivalent. Source-by-source, we can modify  $\mathcal{D}$  on a source  $s$  to be the same as  $\mathcal{D}'$  on source  $s$ , and the results of  $Q$  is not changed, by  $s, Q$ -equivalence.

Assume  $\equiv$  is useful and  $\mathcal{D} \equiv \mathcal{D}'$ . Fix a context  $C$  for source  $s$  such that  $(\mathcal{D}_s, C) \models Q$ . We will argue that  $(\mathcal{D}'_s, C) \models Q$ . Note that for every source  $s$ , identical instances for source  $s$  must be  $\equiv_s$ , since  $\equiv_s$  is an equivalence relation. Thus  $(\mathcal{D}_s, C) \equiv (\mathcal{D}'_s, C)$ . Since  $\equiv$  is useful for  $Q$ , this means  $(\mathcal{D}'_s, C) \models Q$  as required.  $\square$

### Proof of Proposition 4

Recall the statement:

For any BCQ  $Q$  and any source  $s$ , two  $s$ -instances are  $(s, Q)$ -equivalent if and only if they agree on each invariant shuffle view of  $Q$  for  $s$ .

The proof will go through an intermediate view, also given by an ECR.

**Definition 10.** Given two  $s$ -instances  $\mathcal{I}_1$  and  $\mathcal{I}_2$ , we say that  $\mathcal{I}_1$  and  $\mathcal{I}_2$  are invariant shuffle equivalent (relative to  $Q$ ) if: Whenever  $\mathcal{I}_1, \sigma$  satisfies  $\text{CanView}^s(Q)(\mathbf{x})$  then there is some shuffle  $\mu$  invariant relative to  $\langle \sigma, \text{CanCtxt}^s(Q) \rangle$ , such that  $\mathcal{I}_2, \sigma$  satisfies  $\mu(\text{CanView}^s(Q))(\mathbf{x})$  and similarly with the role of  $\mathcal{I}_1$  and  $\mathcal{I}_2$  reversed.

We show:

**Proposition 10.** For any source  $s$  invariant shuffle equivalence is identical to  $s, Q$ -equivalence.

*Proof.* First, suppose  $\mathcal{I}_1, \mathcal{I}_2$  are local instances for source  $s$  that are invariant shuffle equivalent, and suppose we have a match of  $Q$  in  $(\mathcal{I}_1, C)$  via  $h^{1,C}$ . We want to show that there is a match in  $(\mathcal{I}_2, C)$ .

We know that the variables in  $\text{SJVars}(s, Q)$  are mapped by  $h^{1,C}$  into  $\mathcal{I}_1$ . Let  $h_0$  be the restriction of  $h^{1,C}$  to the variables of  $\text{SJVars}(s, Q)$ . Then  $\mathcal{I}_1, h_0$  satisfies  $\text{CanView}^s(Q)$ . Thus by shuffle equivalence there is some shuffle  $\mu$  invariant relative to  $\langle h_0, \text{CanView}^s(Q) \rangle$ , such that  $\mathcal{I}_2, h_0 \models \mu(\text{CanView}^s(Q))$ , with witness  $h_2$  extending  $h_0$ . We also know that  $C, h_0$  satisfies  $\text{CanCtxt}^{S_0}(Q)$ , since  $h^{1,C}$  witnesses this as well. Applying the definition of shuffle invariance,  $C, h_0$  satisfies  $\mu(\text{CanCtxt}^{S_0}(Q))$ . Let  $h^{\mu,C}$  be a homomorphism witnessing this. Note that since  $h^{\mu,C}$  extends  $h_0$  and  $h_0$  restricts  $h^{1,C}$ ,  $h^{\mu,C}$  and  $h^{1,C}$  agree on their common variables. Define  $h^{2,C}$  by mapping the variables in  $\text{SVars}(s, Q)$  as in  $h_2$ , and those variables outside of  $\text{SJVars}(s, Q)$  as in  $h^{\mu,C}$ . Since these are two compatible homomorphisms,  $h^{2,C}$  witnesses that  $(\mathcal{I}_2, C) \models Q$ . This completes the argument that invariant shuffle equivalence implies global  $Q$ -equivalence.

We now show that global  $Q$ -equivalence implies shuffle equivalence. Suppose  $s$ -instances  $\mathcal{I}_1, \mathcal{I}_2$  are globally  $Q$ -equivalent, and  $\mathcal{I}_1, \sigma$  satisfies  $\text{CanView}^s(Q)(\mathbf{x})$ . We will show that there is a shuffle  $\mu$ , invariant relative to  $\langle \sigma, \text{CanView}^s(Q) \rangle$ , such that  $\mathcal{I}_2, \sigma \models \mu(\text{CanView}^s(Q))(\mathbf{x})$ .

Let  $C_1$  be the canonical database of  $\sigma(\text{CanCtxt}^s(Q))$ . That is, for each source  $s$  other than  $s$ , we have a fact for each atom of  $s$  atom of  $Q$ , where each variable  $x$  of  $\text{SJVars}(s, Q)$  is replaced by  $\sigma(x)$  and each variable  $x$  not in  $\text{SJVars}(s, Q)$  is replaced by a fresh element  $c_x$ .

$Q$  clearly holds in  $(\mathcal{I}_1, C_1)$ . So by global  $Q$ -equivalence,  $Q$  holds in  $(\mathcal{I}_2, C_1)$  via some homomorphism  $h$ . Note that in  $(\mathcal{I}_2, C_1)$  the only elements that are shared between  $\mathcal{I}_2$ -facts and  $C_1$ -facts lie in the range of  $\sigma$ . Thus  $h$  must map the variables in  $\text{SJVars}(s, Q)$  to the image of  $\sigma$ . The binding  $\sigma$  may not be injective, but we let  $\sigma^{-1}$  be “some inverse” that is, any function from the range of  $\sigma$  to variables such that for any  $c$  in the range of  $\sigma$   $\sigma(\sigma^{-1}(c)) = c$ . Let  $\mu$  map any variable  $x \in \text{SJVars}(s, Q)$  to  $\sigma^{-1}(h(x))$ . So  $\sigma(\mu(x)) = h(x)$ .

We first claim that  $\mu$  is invariant relative to  $\langle \sigma, \text{CanCtxt}^s(Q) \rangle$ . We show this by arguing that  $h$  is a homomorphism from  $\sigma(\mu(\text{CanCtxt}^s(Q)))$  to  $\sigma(\text{CanCtxt}^s(Q))$ . By definition of  $\mu$ , we have for each atom  $A(x_1 \dots x_m, y_1 \dots y_n)$  of  $\text{CanCtxt}^s(Q)(\mathbf{x})$ ,

$$\begin{aligned} A(h(x_1) \dots h(x_m), h(y_1) \dots h(y_n)) = \\ A(\sigma(\mu(x_1)) \dots \sigma(\mu(x_m)), h(y_1) \dots h(y_n)) \end{aligned}$$



$A(h(x_1) \dots h(x_m), h(y_1) \dots h(y_n))$  is in  $\sigma(\text{CanCtxt}^s(Q))$  since by assumption  $h$  is a homomorphism from  $Q$  to  $(\mathcal{I}_2, C_1)$ ,  $C_1$  is the canonical database of  $\sigma(\text{CanCtxt}^s(Q))$ , and  $\mathcal{I}_2$  is an  $s$ -instance, and hence cannot contain any facts over the relations in  $(\mu(\text{CanCtxt}^s(Q))) (\sigma)$ . Thus

$$A(\sigma(\mu(x_1)) \dots \sigma(\mu(x_m)), h(y_1) \dots h(y_n))$$

lies in  $\sigma(\text{CanCtxt}(Q))$  as required.

We next claim that  $\mathcal{I}_2, \sigma \models \mu(\text{CanView}^s(Q))$ . The witness will be the extension  $h'$  of  $\sigma$  that maps all variables in  $\text{SVars}(s, Q) - \text{SJVars}(s, Q)$  via  $h$ .

Consider an atomic formula  $A(x_1 \dots x_m, y_1 \dots y_n)$  of  $\text{CanView}^s(Q)$ , where  $\mathbf{x}$  are free variables of  $\text{CanView}^s(Q)$ . Therefore  $A(\mu(x_1) \dots \mu(x_m), y_1 \dots y_n)$  is a generic atom of  $\mu(\text{CanView}^s(Q))$ . To argue that  $h'$  is a homomorphism that witnesses  $\mathcal{I}_2, \sigma \models \mu(\text{CanView}^s(Q))$ , we need to argue that

$$A(\sigma(\mu(x_1)) \dots \sigma(\mu(x_m)), h(y_1) \dots h(y_n))$$

holds in  $\mathcal{I}_2$ .

But by the definition of  $\mu$ , this is equivalent to showing that

$$A(h(x_1) \dots h(x_m), h(y_1) \dots h(y_n))$$

holds in  $\mathcal{I}_2$ . But this follows since  $h$  is a homomorphism of  $Q$  into  $(\mathcal{I}_2, C_1)$ .  $\square$

To complete the proof of Proposition 4 we show:

**Proposition 11.** *For any CQ  $Q$  and source  $s$ , two  $s$ -instances are invariant shuffle equivalent if and only if they agree on each invariant shuffle view of  $Q$  for  $s$ .*

*Proof.* We first show that if  $\mathcal{I}_1$  and  $\mathcal{I}_2$  agree on the invariant shuffle views of  $Q$ , they are invariant shuffle equivalent. Suppose  $\mathcal{I}_1, \sigma \models \text{CanView}^s(Q)$ , and let  $\tau$  be the type of  $\sigma$ . Since the identity is invariant relative to  $\tau$ , we have  $\mathcal{I}_1, \sigma$  satisfies  $V_\tau$ , and thus  $\mathcal{I}_2, \sigma$  must satisfy it. Therefore there is  $\mu$  that is invariant relative to  $\tau$  such that  $\mathcal{I}_2, \sigma \models \mu(\text{CanView}^s(Q))$ . Since  $\tau$  is of type  $\sigma$ , we have  $\mu$  is invariant relative to  $\langle \sigma, \text{CanCtxt}^s(Q) \rangle$ . Arguing symmetrically for  $\mathcal{I}_2$ , we see that  $\mathcal{I}_1$  and  $\mathcal{I}_2$  are invariant shuffle equivalent.

In the other direction, suppose  $\mathcal{I}_1$  and  $\mathcal{I}_2$  are invariant shuffle equivalent. We will argue that they agree on each invariant shuffle view  $V_\tau$ .

Towards that end, suppose  $\mathcal{I}_1, \sigma \models V_\tau$ . That is,  $\mathcal{I}_1, \sigma \models \tau \wedge \mu(\text{CanView}^s(Q))$  for some  $\mu$  that is invariant relative to  $\tau$ . Let  $\sigma'$  be the pre-image of  $\sigma$  under  $\mu$ : that is, the variable binding defined by  $\sigma'(x) = \sigma(\mu(x))$ . Then  $\mathcal{I}_1, \sigma' \models \text{CanView}^s(Q)$  by definition. Thus by invariant shuffle equivalence, there is  $\mu'$  invariant for  $\sigma'$ ,  $\text{CanCtxt}^s(Q)$  such that  $\mathcal{I}_2, \sigma' \models \mu'(\text{CanView}^s(Q))$ .

Let  $\mu'' = \mu'(\mu)$ . We will show that  $\mu''$  witnesses that  $\mathcal{I}_2, \sigma \models V_\tau$ . We first verify invariance:

**Claim 1.**  $\mu''$  is invariant relative to  $\tau, \text{CanCtxt}^s(Q)$ .

*Proof.* Suppose  $\sigma_0$  satisfies  $\tau$ , and  $\mathcal{I}_0, \sigma_0 \models \text{CanCtxt}^s(Q)$ . Let  $\sigma'_0$  be the pre-image of  $\sigma'$  under  $\mu$ . Note that a shuffle that is invariant for  $\sigma$  must be invariant for  $\sigma_0$ , since  $\sigma_0$  satisfies all the equalities that  $\sigma$  does. Similarly a shuffle that is invariant for  $\sigma'$  must be invariant for  $\sigma'_0$ . We can see that the following chain of implications:

$$\begin{aligned} \mathcal{I}_0, \sigma_0 &\models \mu(\text{CanCtxt}^s(Q)) \text{ by invariance of } \mu \text{ for } \tau \\ \mathcal{I}_0, \sigma'_0 &\models \text{CanCtxt}^s(Q) \text{ by definition of } \sigma'_0 \\ \mathcal{I}_0, \sigma'_0 &\models \mu'(\text{CanCtxt}^s(Q)) \text{ by invariance of } \mu' \text{ for } \sigma'_0 \\ \mathcal{I}_0, \sigma_0 &\models \mu(\text{CanCtxt}^s(Q)) \text{ by definition of } \sigma'_0 \text{ again} \end{aligned}$$

$\square$

We now show that  $\mathcal{I}_2, \sigma$  satisfies the corresponding shuffled query.

**Claim 2.**  $\mathcal{I}_2, \sigma \models \mu''(\text{CanView}^s(Q))$ .

*Proof.* We make the following observation. For any instance  $\mathcal{I}$ , bindings  $\sigma_0$ , CQs  $R$ , and shuffles  $\mu_0$ , let  $\sigma_1$  be the pre-image of  $\sigma_0$  under  $\mu_0$ . Then  $\mathcal{I}, \sigma_0 \models \mu(R)$  if and only if  $\mathcal{I}, \sigma_1 \models R$ .

Let  $\sigma''$  be the pre-image of  $\sigma$  under  $\mu''$ . Note that  $\sigma''$  is also the pre-image of  $\sigma'$  under  $\mu'$ .

From the observation above, we see that the following are equivalent:

$$\begin{aligned}\mathcal{I}_2, \sigma &\models \mu(\text{CanView}^s(Q)) \\ \mathcal{I}_2, \sigma'' &\models \text{CanView}^s(Q) \\ \mathcal{I}_2, \sigma' &\models \mu'(\text{CanView}^s(Q))\end{aligned}$$

which gives the proof of the claim.  $\square$

Putting together the two claims, We conclude that  $\mathcal{I}_2, \sigma$  satisfies  $V_\tau$  as required, which completes the proof of Proposition 11.  $\square$

## Proof of Proposition 5

Recall the statement:

For every view defined by a DCQ (possibly unsafe), there is a set of relational algebra views  $\mathcal{V}'$  that induces the same ECR. Applying this to the invariant shuffle views for a CQ  $Q$ , we can find a relational algebra-based d-view that is minimally useful for  $Q$  within the class of all views.

Putting the conclusion “same ECR” another way: if  $V$  is the original DCQ view, then we obtain a finite set of views  $\mathcal{V}'$  with the definition of each view in relation algebra, such that  $\mathcal{V}'$  determines  $V$  and  $V$  determines  $\mathcal{V}'$ .

Clearly if we have this for a single DCQ view  $V$ , we obtain it for a finite set of views (and hence for a d-view) by applying the construction to each view in the set.

We consider a DCQ  $V(x_1 \dots x_n)$  defined by  $\bigvee_i \varphi_i$ . For each  $\varphi_i$  let  $\text{Vars}_i$  be the set of variables within them, and for each subset  $S$  of the vars let  $D_S$  be the set of  $i$  such that  $\varphi_i$  uses variables  $S$ .

We will partially order the disjuncts  $\varphi_i$  based on containment of the set of variables in them.

Given a set of variables  $S = x_{j_1} \dots x_{j_k}$  with  $D_S \neq \emptyset$ , create a view  $V_S(x_{j_1} \dots x_{j_k})$  defined by

$$\bigvee_{i \in D_S} \varphi_i(x_{j_1} \dots x_{j_k}) \wedge \neg \left( \bigvee_{S' \subsetneq D_S, \text{Vars}(\varphi_j)=S'} \varphi_j \right)$$

We claim that the set of views  $V_S$  determines  $V$  and vice versa.

In one direction suppose  $\mathcal{I}_1$  and  $\mathcal{I}_2$  agree on each  $V_S$ , and  $\mathcal{I}_1 \models V(\mathbf{t})$ . Choose  $i$  such that  $\mathcal{I}_1 \models \varphi_i(\mathbf{t})$  with  $S_i = \text{Vars}_i$  minimal. Let  $\mathbf{t}'$  be the subtuple of  $\mathbf{t}$  corresponding to the variables of  $\varphi_i$ . Then  $\mathcal{I}_1 \models V_{S_i}(\mathbf{t}')$  and hence  $\mathcal{I}_2 \models V_{S_i}(\mathbf{t}')$ . From this we see that  $\mathcal{I}_2 \models V_S(\mathbf{t})$ .

In the other direction, suppose  $\mathcal{I}_1$  and  $\mathcal{I}_2$  agree on  $V$ , and  $\mathcal{I}_1 \models V_S(\mathbf{t})$ . Fix  $\varphi_i$  with variables from  $S$  such that  $\mathcal{I}_1 \models \varphi_i(\mathbf{t})$ . We need to show  $\mathcal{I}_2 \models V_S(\mathbf{t})$ . We can assume by induction that  $V$  determines  $V_{S'}$  for each  $S'$  that is a proper subset of  $S$ . First consider the case where  $S$  consists of all variables. Then  $\mathcal{I}_1 \models V(\mathbf{t})$  hence  $\mathcal{I}_2 \models V(\mathbf{t})$ , and thus there is some  $j$  such that  $\mathcal{I}_2 \models \varphi_j(\mathbf{t})$ . If  $\varphi_j$  contains all the variables of  $\varphi_i$ . Using the induction hypothesis and  $\varphi_j$  we can conclude that  $\mathcal{I}_2 \models V_S(\mathbf{t})$  as required. If  $\varphi_j$  contains a proper subset  $S'$  of the variables in  $S$ , then we have  $\mathcal{I}_2 \models \varphi_k(\mathbf{t}')$  for  $\mathbf{t}'$  a proper subtuple of  $\mathbf{t}$ . Choose  $\varphi_k$  and  $\mathbf{t}'$  with this property such that the variables  $S'$  involved are minimized. The  $\mathcal{I}_2 \models V_{S'}(\mathbf{t}')$  so by the induction hypothesis  $\mathcal{I}_1 \models V_{S'}(\mathbf{t}')$ , which contradicts the facts that  $\mathcal{I}_1 \models V_S(\mathbf{t})$ .

Next consider the case where  $S$  is a proper subset of the variables. We extend  $\mathbf{t}$  to  $\mathbf{t}'$  choosing elements outside the active domain of both  $\mathcal{I}_1$  and  $\mathcal{I}_2$ .  $\mathcal{I}_1 \models V(\mathbf{t}')$ , and  $\mathcal{I}_2 \models V(\mathbf{t}')$ . Thus we have a proper subtuple  $\mathbf{t}''$  of  $\mathbf{t}'$  and a disjunct  $\varphi_k$  such  $\mathcal{I}_2 \models \varphi_k(\mathbf{t}'')$ . As above, we can choose  $\mathbf{t}''$  minimal. By our choice of the elements in  $\mathbf{t}' - \mathbf{t}$ , we must have  $\mathbf{t}''$  a subtuple of  $\mathbf{t}$ . If  $\mathbf{t}'' = \mathbf{t}$ , then we can conclude that  $\mathcal{I}_2 \models V_S(\mathbf{t})$  as required. If  $\mathbf{t}''$  is a proper subtuple, we argue by contradiction of the induction hypothesis as above.

## Proofs for Section 5: balancing requirements in the presence of constraints

### Minimally useful CQs beyond local existential rules

In the body of the paper we mentioned that we can not generalize results such as Corollary 1 from Section 3 to the presence of arbitrary local constraints. The following example exhibits the problem.

Consider the schema and utility query  $Q$  from Example 3. Let  $\Sigma$  consist of the view definitions for the invariant shuffle views in the example. We claim that the canonical views of  $Q$  are not minimally useful for  $Q$  within the set of CQ views, relative to  $\Sigma$ .

To see this, we consider the secret query  $p$  from Example 3.

Relative to  $\Sigma$  we have a distributed view based on CQ view definitions that is useful, and UN non-disclosing, namely the views those simply export each shuffle view. However, as observed in Example 3, the canonical views are-UN disclosing for  $p$ . Thus, applying Proposition 1, the canonical views can not be minimally useful among the class of all CQ views relative to  $\Sigma$ .

### Proof of Theorem 4

Recall the statement:

Let  $\Sigma$  be any set of existential rules that are local. Suppose that  $Q$  is minimal with respect to  $\Sigma$ . If CQ views  $\mathcal{V}$  determine CQ  $Q$  over all instances satisfying  $\Sigma$ , then  $\mathcal{V}$  determines each canonical view  $\text{CanView}^s(Q)$  of  $Q$  over all instances satisfying  $\Sigma$ .

We proceed as in the case of no constraints. We modify the determinacy algorithm by chasing with the local constraints in each round, giving us the Algorithm 2.

Figure 2:  $\text{Determinacy}(Q, \mathcal{V}, \Sigma)$

```

1:  $F_0 := \text{Chase}_\Sigma(\text{canondb}(Q))$ 
2: while true do
3:    $F_1 := \text{Chase}_{\text{ForView}(\mathcal{V})}(F_0)$ 
4:    $F_2 := \text{Chase}_{\text{BackView}'(\mathcal{V})}(F_1)$ 
5:    $G_2 := \text{Chase}_{\Sigma'}(F_2)$ 
6:   if  $\exists h : Q' \rightarrow G_2$  mapping each free variable  $v$  of  $Q'$  into  $c_v \in \text{adom}(\text{canondb}(Q))$  then
7:     return true
8:   end if
9:    $F_3 := \text{Chase}_{\text{ForView}'(\mathcal{V})}(G_2)$ 
10:   $F_4 := \text{Chase}_{\text{BackView}(\mathcal{V})}(F_3)$ 
11:   $G_4 := \text{Chase}_\Sigma(F_4)$ 
12:   $G_5 := \text{restrict } G_4 \text{ to the original signature}$ 
13:  if  $F_0 \neq G_5$  then
14:     $F_0 := F_0 \cup G_5$ 
15:  else
16:    return false
17:  end if
18: end while

```

As with Algorithm 1, it is a straightforward exercise to see that this algorithm correctly checks determinacy:

**Theorem 8.** (Benedikt, ten Cate, and Tsamoura 2016) *Algorithm 2 returns true if and only if  $\mathcal{V}$  determines  $Q$  relative to  $\Sigma$ .*

Note that the chase may not terminate, thus this is only a semi-decision procedure.

As before, we have a homomorphism from the output of the algorithm to its input:

**Lemma 6.** *For any CQ  $Q$ , views  $\mathcal{V}$  and number  $l$ , there exists a homomorphism  $\nu : \text{UnPrime}(F_2^l(Q, \mathcal{V})) \rightarrow \text{canondb}(Q)$  that is the identity on elements  $c_v$ . In particular, if Algorithm 2 returns **true**, there is a homomorphism from  $F_2^\infty(Q, \mathcal{V})$  to  $\text{canondb}(Q)$  that is the identity on elements  $c_v$ .*

We have the same relation of the output of the algorithm when run on the canonical view to the output when run on  $Q$ :

**Lemma 7.** For any  $l \geq 1$  the source  $s$  atoms in  $F_0^l(Q, \mathcal{V})$  are the same as the atoms of  $F_0^l(\text{CanView}^s(Q), \mathcal{V})$ .

We can now complete the proof of Theorem 4 as in the case without constraints. Locality of constraints ensures that the homomorphism  $h_1$  maps each variable  $x$  in  $\text{SJVars}(s)$  to some element of the form  $c_v$ .

## Proof of Theorem 5

Recall the statement:

For any set of local existential rules  $\Sigma$ , the  $\Sigma$ -invariant shuffle view of  $Q$  are minimally useful within the class of all views, relative to  $\Sigma$ .

We can generalize the ECR global  $Q$ -equivalence to global  $Q$ - $\Sigma$ -equivalence, looking only at contexts that satisfy  $\Sigma$ . Using the same argument we see that the views corresponding to this ECR is minimally useful for  $Q$  relative to  $\Sigma$ .

Let  $\sigma$  be a mapping of  $\text{SJVars}(s)$  into some instance  $\mathcal{I}$ . A shuffle  $\mu$  of  $\text{CanView}^s(Q)$  is  $\Sigma$ -invariant relative to  $\langle \sigma, \text{CanCtxt}^s(Q) \rangle$  if whenever  $\mathcal{I}', \sigma \models \text{CanCtxt}^s(Q)$  and  $\mathcal{I}'$  satisfies  $\Sigma$  then  $\mathcal{I}', \sigma \models \mu(\text{CanCtxt}^s(Q))$ . Note that since the rules are local, the notion of a context satisfying them is well-defined. Invariance is decidable whenever query containment for CQs under  $\Sigma$  is decidable; for example, this is the case when  $\Sigma$  is a set of dependencies with terminating chase.

Fixing  $\Sigma$  and two  $s$ -instances  $\mathcal{I}_1$  and  $\mathcal{I}_2$ , we say that  $\mathcal{I}_1$  and  $\mathcal{I}_2$  are  $\Sigma$ -invariant shuffle equivalent if whenever  $\mathcal{I}_1, \sigma$  satisfies  $\text{CanView}^s(Q)(\mathbf{x})$  then there is some shuffle  $\mu$  which is  $\Sigma$ -invariant relative to  $\langle \sigma, \text{CanView}^s(Q) \rangle$ , such that:  $\mathcal{I}_2, \sigma \models \mu(\text{CanView}^s(Q))(\mathbf{x})$  and vice versa.

We can now extend Proposition 10, following the same proof:

**Proposition 12.** Suppose  $\Sigma$  consists of local existential rules. Then for all instances satisfying the rules  $\Sigma$ -invariant shuffle equivalence is identical to global  $Q$ - $\Sigma$  equivalence.

*Proof.* First, suppose  $\mathcal{I}_1, \mathcal{I}_2$  satisfy all local rules and are  $\Sigma$ -invariant shuffle equivalent. Consider a context  $C$  that satisfies local rules  $\Sigma$  and suppose we have a match of  $Q$  in  $(\mathcal{I}_1, C)$  via  $h^{1,C}$ . We want to show that there is a match in  $(\mathcal{I}_2, C)$ . This will be exactly as in the case without constraints.

We know that the variables in  $\text{SJVars}(s, Q)$  are mapped by  $h^{1,C}$  into  $\mathcal{I}_1$ . Let  $h_0$  be the restriction of  $h^{1,C}$  to the variables of  $\text{SJVars}(s, Q)$ . Then  $\mathcal{I}_1, h_0$  satisfies  $\text{CanView}^s(Q)$ . Thus by  $\Sigma$ -invariant shuffle equivalence there is a shuffle  $\mu$  which is  $\Sigma$ -invariant relative to  $\langle h_0, \text{CanView}^s(Q) \rangle$ , such that  $\mathcal{I}_2, h_0 \models \mu(\text{CanView}^s(Q))$ , with witness  $h_2$  extending  $h_0$ . We also know that  $C, h_0$  satisfies  $\text{CanCtxt}^s(Q)$ , since  $h^{1,C}$  witnesses this as well. Applying the definition of  $\Sigma$ -invariance, and noting that  $C$  satisfies  $\Sigma$  by assumption, we infer that  $C, h_0$  satisfies  $\mu(\text{CanCtxt}^s(Q))$ . Let  $h^{\mu,C}$  be a homomorphism witnessing this. Note that since  $h^{\mu,C}$  extends  $h_0$  and  $h_0$  restricts  $h^{1,C}$ ,  $h^{\mu,C}$  and  $h_0$  agree on their common variables. Define  $h^{2,C}$  by mapping the variables in  $\text{SVars}(s, Q)$  as in  $h_2$ , and those variables outside of  $\text{SJVars}(s, Q)$  as in  $h^{\mu,C}$ . Since these are two compatible homomorphisms,  $h^{2,C}$  witnesses that  $(\mathcal{I}_2, C) \models Q$ . This completes the argument that  $\Sigma$ -invariant shuffle equivalence implies  $Q$ - $\Sigma$ -equivalence.

We now show that global  $Q$ - $\Sigma$ -equivalence implies  $\Sigma$ -invariant shuffle equivalence. Suppose  $s$ -instances  $\mathcal{I}_1, \mathcal{I}_2$  are  $Q$ - $\Sigma$ -equivalent, and  $\mathcal{I}_1, \sigma$  satisfies  $\text{CanView}^s(Q)(\mathbf{x})$ . We will show that there is a shuffle  $\mu$ ,  $\Sigma$ -invariant relative to  $\langle \sigma, \text{CanCtxt}^s(Q) \rangle$  such that  $\mathcal{I}_2, \sigma \models \mu(\text{CanView}^s(Q))(\mathbf{x})$ .

Let  $C_1$  be the context defined in two steps. We first proceed as in the case without constraints: for each source  $s$  other than  $s$ , we have a fact for each  $s$  atom of  $Q$ , where each variable  $x$  of  $\text{SJVars}(s, Q)$  is replaced by  $\sigma(x)$  and each variable  $x$  not in  $\text{SJVars}(s, Q)$  is replaced by a fresh element  $c_x$ . In the second step, we perform the chase construction with  $\Sigma$  to get an instance that satisfies the local constraints.

$Q$  clearly holds in  $(\mathcal{I}_1, C_1)$ . So by  $Q$ - $\Sigma$ -equivalence,  $Q$  holds in  $(\mathcal{I}_2, C_1)$  via some homomorphism  $h$ . As before, the only elements shared between  $\mathcal{I}_2$  and  $C_1$  lie in the range of  $\sigma$ . Thus  $h$  must map the variables in  $\text{SJVars}(s, Q)$  to the image of  $\sigma$ . For each  $c$  in the image of  $\sigma$ , choose a variable  $v_c$  such that  $\sigma$  maps  $v_c$  to  $c$ . Let  $\mu$  map any variable  $x \in \text{SJVars}(s, Q)$  to  $v_{h(x)}$ . Thus  $\sigma(\mu(x)) = h(x)$ .

We claim that  $\mu$  is  $\Sigma$ -invariant relative to  $\langle \sigma, \text{CanCtxt}^s(Q) \rangle$ . We show this by arguing that  $h$  is a homomorphism from  $\mu(\text{CanCtxt}^s(Q))(\sigma)$  to the chase under  $\Sigma$  of  $\text{CanCtxt}^s(Q)(\sigma)$ . By definition of  $\mu$ , we have for each atom  $A(x_1 \dots x_m, y_1 \dots y_n)$  of  $\text{CanCtxt}^s(Q)(\mathbf{x})$ ,

$$A(h(x_1) \dots h(x_m), h(y_1) \dots h(y_n)) = A(\sigma(\mu(x_1)) \dots \sigma(\mu(x_m)), h(y_1) \dots h(y_n))$$

$A(h(x_1) \dots h(x_m), h(y_1) \dots h(y_n))$  is in  $\text{Chase}_\Sigma(\text{CanCtxt}^s(Q)(\sigma))$  since  $h$  is a homomorphism into  $(\mathcal{I}_2, C_1)$  and thus for facts in  $\mu(\text{CanCtxt}^s(Q))(\sigma)$ , it must map into  $C_1$ .

Thus we conclude

$$A(\sigma(\mu(x_1)) \dots \sigma(\mu(x_m)), h(y_1) \dots h(y_n)) \in \text{Chase}_\Sigma(\text{CanCtxt}^s(Q)(\sigma))$$

This completes the proof that  $h$  is a homomorphism into the chase, and thus the proof that  $\mu$  is  $\Sigma$ -invariant relative to  $\langle \sigma, \text{CanCtxt}^s(Q) \rangle$ .

We next claim that  $\mathcal{I}_2, \sigma \models \mu(\text{CanView}^s(Q))$ . The witness will again be the extension of  $\sigma$  that maps all variables in  $\text{SVars}(s, Q) - \text{SJVars}(s, Q)$  via  $h$ .

Consider an atomic formula  $A(x_1 \dots x_m, y_1 \dots y_n)$  of  $\text{CanView}^s(Q)$  where  $\mathbf{x}$  are free variables of  $\text{CanView}^s(Q)$ . That is,

$$A(\mu(x_1) \dots \mu(x_m), \mathbf{y})$$

is a generic atom of  $\mu(\text{CanView}^s(Q))$ . We know that  $A(h(x_1) \dots h(x_m), h(y_1) \dots h(y_n))$  holds in  $\mathcal{I}_2$ , since  $h$  is a homomorphism into  $\mathcal{I}_2$ . Thus

$$A(v_{h(x_1)}, \dots, v_{h(x_m)}, h(y_1) \dots h(y_n))$$

holds of  $\sigma$  in  $\mathcal{I}_2$ , by definition of  $v_c$ . From this we see that

$$A(\mu(x_1) \dots \mu(x_m), h(y_1) \dots h(y_n))$$

holds of  $\sigma$  in  $\mathcal{I}_2$  as required.  $\square$

Recall that the  $\Sigma$ -invariant shuffle views of  $Q$  for  $s$  and types  $\tau$  are defined analogously to the case without local rules, as  $\tau(\mathbf{x}) \wedge \bigvee_\mu \mu(\text{CanView}^s(Q))$  where the disjunction is over  $\Sigma$ -invariant shuffles of  $\tau$ .

The following result is proven exactly as in the case without background knowledge:

**Proposition 13.** *For any Boolean CQ  $Q$ , and any source  $s$  two  $s$ -instances are  $\Sigma$ -invariant shuffle equivalent if and only if they agree on each  $\Sigma$ -invariant shuffle view of  $Q$  for  $s$ .*

Putting the prior results together gives us the extension of Theorem 5.

## The power of replication: proof of Theorem 6

Recall the statement:

If BCQ  $Q$  contains a relation of non-zero arity replicated across all sources then there are views that are useful for  $Q$  and UN non-disclosing for BCQ  $p$  if and only if there is no homomorphism of  $p$  to  $Q$ .

Note that the condition on  $p$  and  $Q$  can be restated as saying that  $Q$  does not logically entail  $p$ .

For notational simplicity, we keep the replication constraints implicit, assuming that the replicated predicates are named  $T$  in each source and  $Q$  refers to this “global”  $T$ .

One direction of the theorem is clear: if there is a homomorphism of  $p$  to  $Q$  and  $\mathcal{V}$  are useful for  $Q$ , then  $\mathcal{V}$  can not be UN non-disclosing for  $p$ , since on any instance where  $Q$  holds, the views will disclose  $p$ .

For the other direction, we show, as in the case without constraints, that there is a single d-view that is  $Q$ -universal: it works for any  $p$  such that there is no homomorphism from  $p$  to  $Q$ . We provide views that are not isomorphism-invariant, assuming that the active domain of instances is  $\text{Pair}(\mathcal{N})$  defined below.

$\text{Pair}(\mathcal{N})$  is the set that contains  $\mathcal{N}$  and is closed under pairing: when  $x$  and  $y$  are allowed, then so is  $(x, y)$ , the ordered pair consisting of  $x$  and  $y$ . Note that all elements in  $\text{Pair}(\mathcal{N})$  have a finite pairheight, where pairheight is defined as follows:  $\text{pairheight}(x) = 0$  for  $x$  an integer of  $\mathcal{N}$ , and  $\text{pairheight}((x, y)) = \max(\text{pairheight}(x), \text{pairheight}(y)) + 1$ .

Given two instances  $\mathcal{I}_1$  and  $\mathcal{I}_2$  for the same schema, the synchronous product of  $\mathcal{I}_1$  and  $\mathcal{I}_2$  is the instance defined as follows:

- elements of the instance are pairs  $(x, y)$  with  $x \in \mathcal{I}_1, y \in \mathcal{I}_2$ .
- for each relation  $R$ , we have  $R((x_1, y_1) \dots (x_n, y_n))$  holds exactly when  $R(x_1, \dots x_n)$  holds in  $\mathcal{I}_1$  and  $R(y_1, \dots y_n)$  holds in  $\mathcal{I}_2$

Note that the projection on the first component is a homomorphism of the product to instance  $\mathcal{I}_1$  and projection on the second component is a homomorphism to  $\mathcal{I}_2$ .

By a *position* we mean a relation  $S$  and a number between 1 and the arity of  $S$ . That is, a position describes an argument of a relation. We assume that each variable of  $Q$  is associated with a unique integer index.

We consider the transformation  $\text{Str}$  on a d-instance that maps each local instance to its product with  $\text{canondb}(Q)$ .

Note that  $\text{Str}(\mathcal{D})$  is over the same schema as  $\mathcal{D}$  and that  $\text{Str}(\mathcal{D})$  validates the replication constraint. Furthermore we suppose that the domain of  $\text{canondb}(Q)$  is included in  $\mathcal{N}$  which ensures that the minimal pairheight of elements in the relation  $T$  of  $\text{Str}(\mathcal{D})$  will be the minimal pairheight of elements in the relation  $T$  of  $\mathcal{D}$  plus one.

**Definition 11.** We define  $\equiv_s$  on  $s$ -instances as the reflexive transitive closure of  $\mathcal{R}$  where  $\mathcal{R}$  is defined as  $\mathcal{I} \mathcal{R} \mathcal{I}'$  when  $\mathcal{I} = \text{Str}(\mathcal{I}')$  for  $\mathcal{I}, \mathcal{I}'$  two  $s$ -instances.

**Definition 12.** We define  $\equiv_G$  on  $d$ -instances as the reflexive transitive closure of  $\mathcal{R}$  where  $\mathcal{R}$  is defined as  $\mathcal{D} \mathcal{R} \mathcal{D}'$  when  $\mathcal{D}' = \text{Str}(\mathcal{D})$  for  $\mathcal{D}, \mathcal{D}'$  two  $d$ -instances.

**Definition 13.** The ECR  $\equiv$  is defined as  $\mathcal{D} \equiv \mathcal{D}' \Leftrightarrow \bigwedge_{s \in \text{Srcs}} \mathcal{D}_s \equiv_s \mathcal{D}'_s$ .

**Proposition 14.** For two  $d$ -instances  $\mathcal{D}$  and  $\mathcal{D}'$ , when  $\mathcal{D} \models Q$  then  $\mathcal{D} \equiv \mathcal{D}'$  if and only if  $\mathcal{D} \equiv_G \mathcal{D}'$ .

*Proof.* Clearly  $\mathcal{D} \equiv_G \mathcal{D}'$  implies  $\mathcal{D} \equiv \mathcal{D}'$ . Now let us suppose that  $\mathcal{D} \equiv \mathcal{D}'$  with  $\mathcal{D} \models Q$  and let us show that  $\mathcal{D} \equiv_G \mathcal{D}'$ .

Since  $\mathcal{D} \equiv \mathcal{D}'$ , we have, for each source  $s$ , there exists  $i_s$  such that  $\mathcal{D}_s = \text{Str}^{i_s}(\mathcal{D}'_s)$  or  $\mathcal{D}'_s = \text{Str}^{i_s}(\mathcal{D}_s)$ , where the notation  $\text{Str}^i$  means iterating  $\text{Str}$   $i$  times. On a  $d$ -instance  $\mathcal{D}$  where the replicated relation  $T$  is not empty, the minimal pairheight of elements appearing in position  $T[1]$  needs to be equal for all sources  $s$ . Therefore we cannot have distinct sources  $s$  and  $s'$  where  $i_s \neq i_{s'}$ . And we cannot have  $s$  and  $s'$  such that  $\mathcal{D}_s = \text{Str}^{i_s}(\mathcal{D}'_s)$  and  $\mathcal{D}_{s'} = \text{Str}^{i_{s'}}(\mathcal{D}_{s'})$  with  $i_s, i_{s'} > 0$ . Therefore, when the replicated predicate is not empty,  $\mathcal{D} \equiv \mathcal{D}'$  implies  $\mathcal{D} \equiv_G \mathcal{D}'$ .

When  $\mathcal{D} \models Q$  we have that the replicated predicate  $T$  is not empty (since it appears in  $Q$ ) which proves that  $\mathcal{D} \equiv_G \mathcal{D}'$ .  $\square$

We now show that  $\equiv$  is the view that we want.

**Proposition 15.** The view corresponding to ECR  $\equiv$  is useful for  $Q$ .

*Proof.* Let  $\mathcal{D}$  and  $\mathcal{D}'$  be two  $d$ -instances and let us suppose  $\mathcal{D} \equiv \mathcal{D}'$  and  $\mathcal{D} \models Q$ .

By Proposition 14, we have that  $\mathcal{D} \equiv_G \mathcal{D}'$ . Recall that there is a homomorphism from  $\text{Str}(\mathcal{D})$  to  $\mathcal{D}$ . Thus it is clear that if  $\text{Str}(\mathcal{D}) \models Q$  it must be that  $\mathcal{D} \models Q$ . On the other hand if we have any match  $h$  of  $Q$  in  $\mathcal{D}$ , we can extend it to a match of  $Q$  in the product by taking any variable  $x$  of  $Q$  to  $(h(x), c_x)$  where  $c_x$  is the constant corresponding to  $x$  in  $\text{canondb}(Q)$ . From  $\mathcal{D} \equiv_G \mathcal{D}'$  either  $\mathcal{D}$  or  $\mathcal{D}'$  can be obtained from the other by applying  $\text{Str}$ , and so  $\mathcal{D} \models Q$  implies  $\mathcal{D}' \models Q$ .  $\square$

**Proposition 16.** The view corresponding to ECR  $\equiv$  is UN non-disclosing for  $p$ .

*Proof.* Given an instance  $\mathcal{D}$ , we know that  $\text{Str}(\mathcal{D})$  has a homomorphism into  $\text{canondb}(Q)$ . Therefore there is no homomorphism from  $p$  into  $\text{Str}(\mathcal{D})$  because if there were, we would have a homomorphism from  $p$  into  $\text{canondb}(Q)$ , a contradiction of the assumption that  $Q$  did not entail  $p$ .  $\square$

## Proof of Proposition 7

Recall the statement:

There is a  $d$ -schema with a replication constraint, along with Boolean CQs  $Q$  and  $p$  such that there is a  $d$ -view which is useful for  $Q$  and UN non-disclosing for  $p$ , but there is no  $d$ -view whose view definitions return only facts containing values in the active domain and which commute with isomorphisms.

For a query  $Q_V$  that always returns tuples containing only values in the active domain of the input, we will use the following restricted version of the isomorphism-invariance property:

If  $\mathcal{I}$  and  $\mathcal{I}'$  are source instances for  $Q_V$  with  $\mathcal{I}'$  formed from  $\mathcal{I}$  via an isomorphism that is the identity on  $\text{adom}(V(I))$ , then  $\mathcal{I}$  and  $h(\mathcal{I})$  agree on  $Q_V$ .

For example, any query in relational algebra has this property.

Consider the following Boolean CQs (existentially quantifiers omitted):

$$Q = P(x) \wedge P(y) \wedge S(y) \wedge S(z) \wedge T(z) \wedge T(x) \wedge R(w)$$

$$p = S(x) \wedge T(x) \wedge P(x)$$

It is easy to see that  $Q$  does not entail  $p$ . Because there is a replicated relation mentioned in the query, Theorem 6 implies that we can get a  $d$ -view that is useful for  $Q$  and UN non-disclosing for  $p$ .

Now, by way of contradiction, fix a  $d$ -view  $\mathcal{V}$  that is useful for  $Q$ , returns only facts containing only values in the active domain, and satisfies the isomorphism-invariance property.

**Proposition 17.** *Suppose d-view  $\mathcal{V}$  is useful for  $Q$ . Given two instances  $\mathcal{I}_1$  and  $\mathcal{I}_2$  for the source whose unreplicated relation is  $G$ , if  $\mathcal{I}_1$  and  $\mathcal{I}_2$  agree on the replicated relation and agree on each view in  $\mathcal{V}$  for this source, then they must agree on  $G$ .*

*Proof.* Let  $c \in G(\mathcal{I}_1)$  and form  $\mathcal{D}_1$  by choosing a context for the other unshared relations with each relation containing only  $c$ . Let  $\mathcal{D}_2$  be formed similarly from  $\mathcal{I}_2$ . Then  $\mathcal{D}_1$  and  $\mathcal{D}_2$  agree on all views in  $\mathcal{V}$ . There is a match of  $Q$  on  $\mathcal{D}_1$  so the same must be true of  $\mathcal{D}_2$ , since  $\mathcal{V}$  is useful for  $Q$ . Clearly the witness can only be  $c$ , thus  $c \in G(\mathcal{I}_2)$ .  $\square$

By way of contradiction, we fix a d-view  $\mathcal{V}$  that is useful for  $Q$  and UN non-disclosing for  $p$ , where each view in  $\mathcal{V}$  returns facts containing only elements of the active domain and satisfies the isomorphism-invariance property.

Let  $\mathcal{D}_0$  be the critical instance, recalling that this has only a single element  $*$  with every relation holding of it. We show that  $\mathcal{D}_0$  contradicts that  $\mathcal{V}$  is UN non-disclosing for  $p$ . Note that  $\text{adom}(\mathcal{V}_0) \subseteq \{*\}$  since we assume views return facts containing only elements in the active domain of the input instance.

**Proposition 18.** *Let  $\mathcal{D}$  agree with  $\mathcal{D}_0$  on  $\mathcal{V}$ . For any unreplicated relation  $G$ , we have  $\mathcal{D} \models \forall x.((G(x) \wedge x \neq *) \rightarrow R(x))$ .*

*Proof.* Let  $\mathcal{I}$  be the restriction of  $\mathcal{D}$  to the  $G$  source, and suppose  $G(v)$  holds in  $\mathcal{I}$  for  $v \neq *$ . Let  $c$  be a fresh value and  $\mathcal{I}'$  be the result of applying an isomorphism swapping  $v$  and  $c$ . Since  $v \neq *$  the isomorphism-invariance property implies that  $\mathcal{I}$  and  $\mathcal{I}'$  agree on the views. Since  $\mathcal{I}'$  and  $\mathcal{I}$  disagree on  $G$ , Proposition 17 implies they cannot agree on the replicated relation. Since  $c$  was fresh and no other change occurred outside of the swap of  $v$ , this implies that  $R(v)$  must have held.  $\square$

By usefulness of  $\mathcal{V}$  and the fact that  $\mathcal{D}_0$  has a match of  $Q$ , we know that  $\mathcal{D}$  has a match of  $Q$  with  $x_0, y_0, z_0, w_0$ . We first consider the case where two of  $x_0, y_0, z_0$  are the same. In this case  $p$  has a match.

Now suppose the match has all of  $x_0, y_0, z_0$  distinct. At most one of them can be  $*$ , so assume  $y_0$  and  $z_0$  are not  $*$ . Proposition 18 implies that they are both in  $R$ . Since  $y_0$  and  $z_0$  are not in  $\text{adom}(\mathcal{V}_0)$ , when we consider the result of swapping  $y_0$  and  $z_0$ , it does not impact the views, by the isomorphism-invariance property. Since this swap also does not change  $R$ , we can apply Proposition 17 to conclude that all local sources agree on  $y_0$  and  $z_0$ . In particular we have  $P(y_0) \wedge S(z_0)$ , thus we have  $P(z_0) \wedge S(y_0)$ . But then either  $y_0$  or  $z_0$  gives a match of  $p$ .