

Theo yêu cầu của khách hàng, trong một năm qua, chúng tôi đã dịch qua 16 môn học, 34 cuốn sách, 43 bài báo, 5 sổ tay (chưa tính các tài liệu từ năm 2010 trở về trước) Xem ở đây

**DỊCH VỤ  
DỊCH  
TIẾNG  
ANH  
CHUYÊN  
NGÀNH  
NHANH  
NHẤT VÀ  
CHÍNH  
XÁC  
NHẤT**

Chỉ sau một lần liên lạc, việc dịch được tiến hành

Giá cả: có thể giảm đến 10 nghìn/1 trang

Chất lượng: Tao dựng niềm tin cho khách hàng bằng công nghệ 1. Bạn thấy được toàn bộ bản dịch; 2. Bạn đánh giá chất lượng. 3. Bạn quyết định thanh toán.

Tài liệu này được dịch sang tiếng việt bởi:

**[www.mientayvn.com](http://www.mientayvn.com)**

Từ bản gốc:

<https://drive.google.com/folderview?id=0B4rAPqlxIMRDcGpnN2JzSG1CZDO&usp=sharing>

Liên hệ để mua:

[thanhlam1910\\_2006@yahoo.com](mailto:thanhlam1910_2006@yahoo.com) hoặc [frbwrthes@gmail.com](mailto:frbwrthes@gmail.com) hoặc số 0168 8557 403 (gặp Lâm)

Giá tiền: 1 nghìn /trang đơn (trang không chia cột); 500 VND/trang song ngữ

Dịch tài liệu của bạn: [http://www.mientayvn.com/dich\\_tiang\\_anh\\_chuyen\\_nganh.html](http://www.mientayvn.com/dich_tiang_anh_chuyen_nganh.html)

## Datalog and Its Extensions for Semantic Web Databases 4 h 16 18/7

Abstract. Since the early 70s, data management played a central role in organizations and represented a challenging area of research. A number of languages have been proposed to model, query, and manipulate data, as well as for expressing very general classes of integrity constraints, inference procedures, and ontological knowledge. Such languages are nowadays crucial for many applications such as semantic data publishing and integration, decision support, and knowledge management. In this tutorial we first introduce Datalog, a powerful rule-based language originally intended for expressing complex queries over relational data, and that today is at the basis of languages for the specification of optimization and constraint satisfaction problems as well as of ontological constraints in data and knowledge bases. We then discuss the limitations of Datalog for the semantic web, in particular for ontological modeling and reasoning, and we present several extensions that allow to capture some of the ontology languages of the OWL family, the standard language for semantic data modeling on the semantic web.

### 1 Introduction

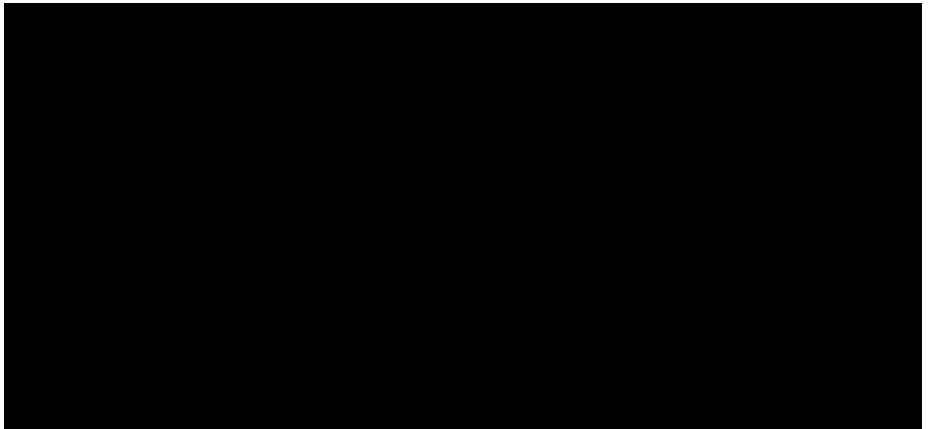
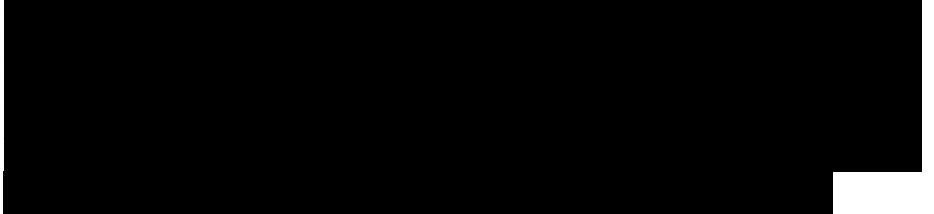
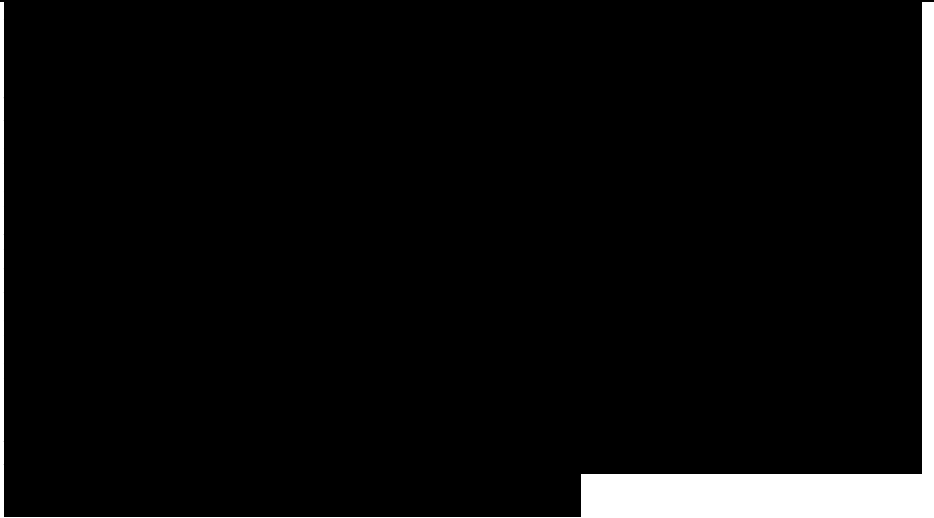
Data and knowledge base systems played a dominant role in computer science since the 70s when organizations massively adopted them to support their business operations and decision making activities. Yet, in the last decade,

Datalog và các phần mở rộng của nó cho cơ sở dữ liệu web ngữ nghĩa

Tóm tắt. Kể từ đầu thập niên 70, quản lý dữ liệu đóng vai trò trung tâm trong các tổ chức và là một lĩnh vực nghiên cứu đầy thách thức. Các nhà nghiên cứu đã đưa ra một số ngôn ngữ để mô hình hóa, truy vấn và thao tác dữ liệu, cũng như biểu diễn các loại ràng buộc toàn vẹn rất tổng quát, các thủ tục suy diễn, và tri thức ontology. Hiện nay, những ngôn ngữ như thế đóng vai trò quan trọng trong nhiều ứng dụng chẳng hạn như xuất bản và tích hợp dữ liệu ngữ nghĩa, hỗ trợ quyết định, và quản trị tri thức. Trong bài hướng dẫn này, trước hết chúng tôi giới thiệu Datalog, một ngôn ngữ mạnh dựa trên quy tắc lúc đầu được thiết kế để biểu diễn các truy vấn phức tạp trên dữ liệu quan hệ, và hiện nay là cơ sở của các ngôn ngữ đặc tả tối ưu hóa và các bài toán thỏa mãn ràng buộc cũng như các ràng buộc ontology trong cơ sở dữ liệu và cơ sở tri thức. Sau đó, chúng ta sẽ đề cập đến những mặt hạn chế của datalog trong web ngữ nghĩa, đặc biệt trong mô hình hóa và suy luận ontology, và chúng tôi cũng trình bày một số phần mở rộng để chúng ta có thể hiểu được một số ngôn ngữ ontology thuộc họ OWL, ngôn ngữ tiêu chuẩn để mô hình hóa dữ liệu ngữ nghĩa trên web ngữ nghĩa.

such systems became even more popular as data and knowledge turn out to be an intrinsic part of every individual and collective activity in our society. In this setting, a major problem is to represent information in such a way that software programs can access it and act as if they really understand its semantics. On the one hand, initiatives such as the semantic web defined languages like RDF(S), FLORA and OWL to support the creation of semantically annotated data, enabling ontological querying and reasoning. On the other hand, the Linked Open Data (LOD) community produced very large amounts of semantically enriched data that enabled a multitude of data-driven semantic web applications. The need for efficient processing of semantic data stimulated several research initiatives addressing data management problems such as representation, storage and querying. This tutorial is intended for people familiar with the basics of database and semantic web technologies who want to explore more in depth the connection between modeling languages used in these fields, and their practical adoption for knowledge representation and data management purposes.

**A Bit of History.** The data management problems we are facing today are not completely new. As an example, in the late 70s, Datalog [14] emerged as a prominent language from logic programming [30]. The term Datalog was coined by David Maier and reflects the intention of devising a counterpart of Prolog - the most prominent rule-



based formalism in computer science - for data processing. While Prolog is undecidable in general, if we consider the program as fixed, Datalog enjoys tractable reasoning complexity w.r.t. the size of the input database. Datalog's original aim was to be used as an expressive language for querying relational data; in fact, it adds recursion to the relational algebra, and therefore goes beyond the expressive power of select-project-join queries. Recursion is still important today for reasoning over complex paths in graphlike data which is abundant, for example, in the context of social networks and of the semantic web.

Applications of Datalog include data integration, reasoning about semistructured data, routing, security policy management, enterprise decision automation and many others. As a consequence, Datalog has evolved into a first-class formalism with efficient implementations such as DLV [21] and Clingo [25]. On the other hand, since Datalog **rules (luật, quy tắc)** are a representation of clauses in the function-free Horn fragment of first-order logic (FOL), Datalog revealed itself relevant also for semantic web applications such as ontological modeling and reasoning.

Example 1. Consider, as an example, the following Datalog rules expressing the knowledge that every female and every male is a person.  
 $person(X) \wedge female(X)$     $person(X) \wedge male(X)$ .

Intuitively, to construct the set of all persons, we need to take into account the union of all females and males.

$person(X) \leftarrow female(X)$     $person(X) \leftarrow male(X)$ .

In other words, some objects can be inferred to be persons, even without stating this fact explicitly. Datalog also provides a natural solution to some fundamental reasoning problems, such as the computation of the transitive closure of a binary relation, and is thus adequate for reasoning about graph reachability or connectedness. The rules

$ancestor(X, Y) \wedge parent(X, Y)$   
 $ancestor(X, Z) \wedge parent(X, Y), ancestor(Y, Z)$

express the ancestor relation between persons. Notice that in the last rule the predicate ancestor occurs in both sides of  $\wedge$ , which is an example of a recursive definition that is not possible in relational algebra. ■

From the other side, ontology languages for the semantic web were also designed to be able to express ontological information as the one above. In this setting, ontologies written in ontology languages are intended to describe and structure complex web resources, making them readily available for manipulation by automated agents [27,33]. The different ontology languages of the semantic web are based on the family of description logics (DLs) [5] which, in turn, are decidable fragments of FOL equipped with a convenient syntax.

DLs model a domain of interest in terms of concepts and roles, where concepts are interpreted as sets of individuals (i.e., constants), and roles as binary relations over them. A DL knowledge base  $K = (T, A)$  consists of a TBox  $T$  and an ABox  $A$ . The TBox  $T$  consists of axioms, where the most common axioms are

$ancestor(X, Y) \leftarrow parent(X, Y)$   
 $ancestor(X, Z) \leftarrow parent(X, Y), ancestor(Y, Z)$

statements of inclusion (a C 3) between pairs of concepts or roles. The ABox A is a set of facts about the participation of individuals in concepts and roles. In database terms, an ABox can be seen as a (possibly incomplete) relational database with unary and binary relations only, while a TBox is a set of expressive integrity constraints over the data.

As an example, consider the DL called SHIF, which has most of the features of the OWL languages. SHIF concepts are built by applying concept constructors on roles and other concepts. We use letters A, B, C for concept names and R, S for role names. A SHIF ABox A consists of ground atomic formulae of the form  $R(a, b)$  and  $A(a)$ , where R is a role and A is a concept. A SHIF TBox T consists of axioms of one of the following forms:

(1)  $A \sqcap B \sqsubseteq C$  is a concept inclusion. The complex concept  $A \sqcap B$  stands for conjunction of A and B, and the whole axiom states that each object that is both A and B is also C. This axiom translates into the FOL formula  $\forall X (PA(X) \wedge pB(X) \rightarrow pC(X))$ , where the predicates PA, pB, and pc represent the FOL predicates corresponding to the concept (resp., role) names. As an example, the axiom  $\text{Parent} \sqcap \text{Male} \sqsubseteq \text{Father}$  can express the knowledge that male parents are fathers.

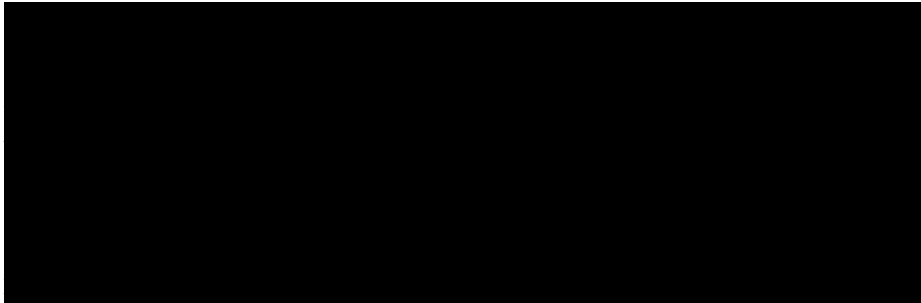
(2)  $A \sqsubseteq B \sqcup C$  is another form of concept inclusion, stating that each object that is A must also be B or C. The axiom can be written as the FOL formula  $\forall X (PA(X) \rightarrow (PB(X) \vee pc(X)))$

(X). As an example, the axiom Person C Female U Male states that every person must be female or male.

(3) A C yR.B is an inclusion employing on the right hand side a complex concept yR.B, denoting the set of all the objects such that all R “neighbors” are B, and is captured by  $\forall X \forall Y (pR(X,Y) \wedge pB(Y))$ . The whole axiom translates into the first-order logic formula  $\forall X \forall Y (PA(X) \wedge APR(X, Y) \wedge pB(Y))$ . An example of this axiom could be MetalDevice C yhasPart.Metal stating that all parts of a metal device are made of metal.

(4) A C 3R.B is an inclusion employing yet another kind of complex concept 3R.B on the right. 3R.B denotes the set of objects that have an R “neighbor” that is B, and is captured by  $\forall X \exists Y (pR(X, Y) \wedge pB(Y))$ . The full axiom translates into the formula  $\forall X (PA(X) \wedge \exists Y (pR(X,Y) \wedge pB(Y)))$ . For example, we can use Student C Battends.Course to express the requirement that each student must attend at least one course.

(5) AC  $\wedge^1$  R.B restricts the number of R “neighbors” B can have, and can make parts of R functional (functionality constraints are common integrity constraints in databases). The axiom says that each object that is A can be related via the role R to at most one object that is B. This axiom translates into the formula  $\forall X \forall Y_1 \forall Y_2 (PA(X) \wedge APR(XY) \wedge Aps(Y1) \wedge APRXY) \wedge Aps(Y2) \wedge Y1 = Y2$ . An example of this axiom could be PersonC  $\wedge^1$  hadldSocialSecNum, stating that a person can have at most one social security number.



(6) An axiom  $A \text{ disj } B$  states disjointness between concepts  $A$  and  $B$ , e.g.,  $\text{Student disj Professor}$  states that students and professors are disjoint sets. The axiom translates into the formula  $\forall X (PA(X) \wedge \neg PB(X))$ .

(7)  $R \subseteq S$  expresses the inclusion of  $R$  in  $S$ , e.g.,  $\text{brotherOf} \subseteq \text{relativeOf}$  captures the knowledge that brothers are relatives. This axiom translates into the formula  $\forall X \forall Y (pR(X, Y) \rightarrow pS(X, Y))$  of first-order logic.

(8) An axiom  $R \text{ inv } S$  allows to define inverse roles, and is translated into first-order logic as  $\forall X \forall Y (pR(X, Y) \rightarrow pS(Y, X))$ . For example,  $a$  is child of  $b$  iff  $b$  is a parent of  $a$ . This can be expressed as  $\text{parentOf} \subseteq \text{inv childOf}$ .

(9) Finally, the axiom  $\text{trans}(R)$  expresses transitivity of the role  $R$ , and is translated into the formula  $\forall X \forall Y \forall Z (pR(X, Y) \wedge pR(Y, Z) \rightarrow pR(X, Z))$ .

Example 2. Consider the Datalog rules of Example 1, the same semantics can be (rather succinctly) expressed in DL syntax as

$\text{Female} \subseteq \text{Person}$   
 $\text{Male} \subseteq \text{Person}$   
 $\text{parentOf} \subseteq \text{ancestorOf}$   
 $\text{trans}(\text{ancestorOf})$

where (1,2) are concept inclusions, (3) is a role inclusion, and (4) states that the role  $\text{ancestor}$  is transitive.

Datalog and DLs have several commonalities but also significant differences which need to be reconciled in order to make them interoperable in semantic web applications as also noticed in [35].



- Female  $\subseteq$  Person (1)
- Male  $\subseteq$  Person (2)
- parentOf  $\subseteq$  ancestorOf (3)
- trans(ancestorOf) (4)





Differences include disjunction and existential quantification in DL ontologies, as well as the different assumptions underlying the semantics of the languages (e.g., open vs closed-world assumption). Table 1 shows a partial translation of SHIF into Datalog. As it can be seen, some expressions have no Datalog counterpart.

Table 1. From the DL SHIF to Datalog

(1) An axiom  $A \sqcup B \sqcup C$  expresses disjunctive information and thus cannot be directly expressed in plain Datalog. One possibility is to employ, bearing the computational cost, disjunctive Datalog, which does support rules of the form  $PB(X) \vee PC(X) \wedge PA(X)$  [20].

(2) Notice that the translation of a functional constraint  $AC \leq 1$  R.B into first-order logic involves the equality predicate. Due to the slight semantic differences between DLs and Datalog, equality is treated differently in the two settings. For instance, in contrast to Datalog, DLs do not employ the so-called unique name assumption (UNA). Different constants are treated as different domain objects by the Datalog semantics, yet a pair of constants may denote the same object in the standard semantics of first-order logic (and therefore DLs).

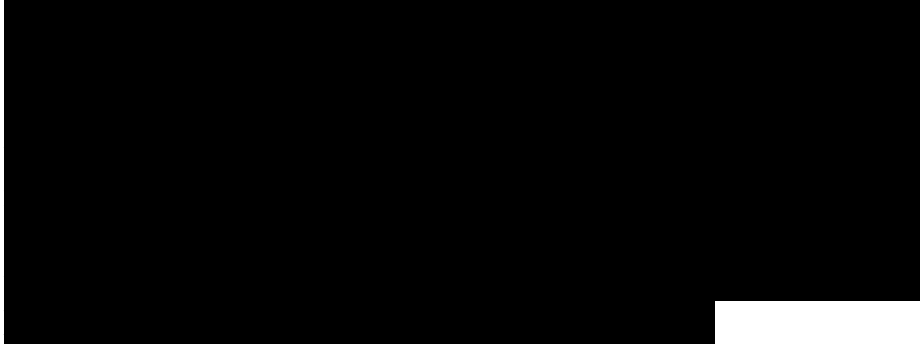
(3) A disjointness constraint  $A \text{ disj } B$  can cause inconsistency of an ontology. E.g., the ontology  $K = (A, T)$ , where  $A = \{\text{Blue}(ba/), \text{Red}(ba/)\}$  and  $T = \{\text{Red disj Blue}\}$ , is inconsistent, i.e., the first-order theory underlying  $K$  is unsatisfiable. Unfortunately, plain Datalog does

not have inconsistent programs and thus a satisfiability-preserving translation for disjointness constraints is not possible, in general.

(4) The axiom  $A \sqsubset B \sqsubset B$ , whose translation into first-order logic involves existential quantification, exposes a crucial difference between Datalog and DLs. Datalog was designed and intended for reasoning over finite databases, under the assumption that only the objects explicitly mentioned in the database exist. In contrast, DL-based ontologies support existential quantification, and are thus able to refer to objects that are not explicitly named in the ontology.

Challenge. We note that existential quantification, disjointness and functional constraints play an important role in knowledge representation. They are in fact necessary to represent even simple constructs of common ontology languages, and there are many important DLs (such as the ones underlying the Lite and DL profiles of OWL) in which ontologies may have infinite and only infinite models that can not be captured by the models of a Datalog program.

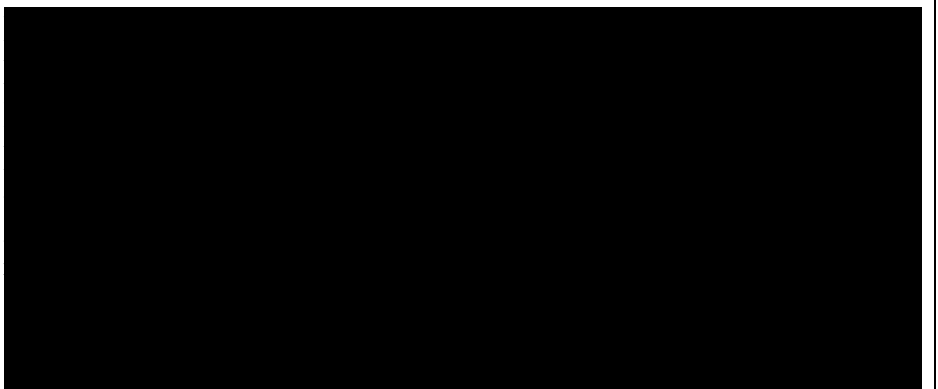
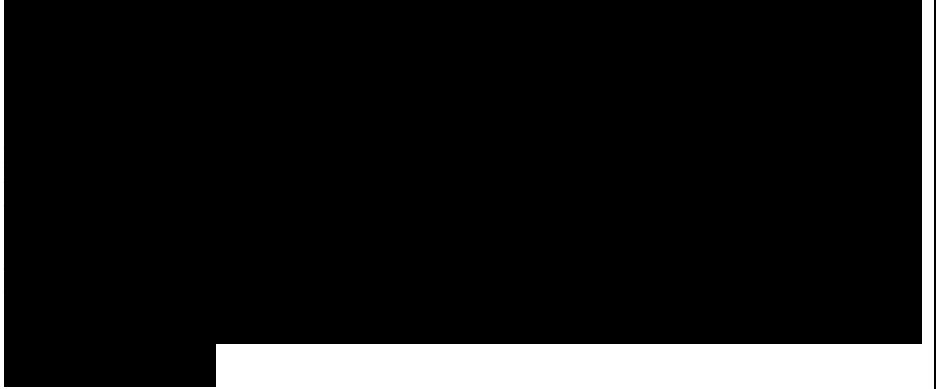
Adding existential quantification to Datalog is the most interesting and relevant extension. Unfortunately, a naive extension of Datalog with existential quantified variables in the head is undecidable [3], and thus our task is nontrivial. We will try to identify meaningful fragments of Datalog with existential variables in the head that have sufficient expressiveness for the above applications, but still retain the



decidability of reasoning.

Due to the existential quantification, correct and **terminating (hữu hạn, có kết thúc)** reasoning algorithms require to develop methods for reasoning about infinite structures without explicitly building them. In the area of DLs and modal logic, decidability, algorithms and complexity results have been shown for various logics that do not enjoy finite models. Most of these fragments allow for only a limited number of variables (often two), and impose some form of guardedness, which can be roughly understood as a syntactically restricted form of quantification that only allows to talk about relations between objects that are close to each other in a structure, and results in regular models that are conveniently similar to trees. While guardedness is of course a limitation, it is often claimed to be a robust reason for decidability [26,39]. Furthermore, there is wide evidence suggesting that guardedness is not overly restrictive for many knowledge representation problems, and it is implicitly or explicitly present in many of the popular languages. Inspired by the above considerations, we will present guarded extensions of Datalog as a solution to our problems.

Roadmap. In the next section we provide some preliminary notions about relational databases and queries. Syntax and semantics of Datalog is presented in Section 3, while Section 4 shows how Datalog can be provided with additional features to capture some of the most common DLs of the semantic web.



In particular, we overview Datalog<sup>±</sup> [9], a family of extensions to Datalog that can be used to specify ontological constraints over databases while preserving tractability of reasoning w.r.t. the size of the data. Finally, Section 5 draws some concluding remarks and collects several useful pointers to further reading material.

## 2 Preliminaries

In this section, we briefly recall some basics on relational databases, and (Boolean) conjunctive queries.

We define the following pairwise disjoint (infinite) sets of symbols: a set  $r$  of constants, that constitute the “normal” domain of a database, a set  $rN$  of labeled **nulls (giá trị không rõng)**, used as placeholders for unknown values, and thus can be also seen as (globally) existentially-quantified variables, and a set  $rV$  of **(regular) (chính quy)** variables, used in queries and dependencies. Different constants represent different values (unique name assumption), while different nulls may represent the same value. A lexicographic order is defined on  $r \cup rN$ , such that every value in  $rN$  follows all those in  $r$ . We denote by  $X$  sequences (or sets, with a slight abuse of notation) of variables or constants  $X_1, \dots, X_k$ , with  $k \geq 0$ . Throughout, let  $[n] = \{1, \dots, n\}$ , for any integer  $n \geq 1$ .

A relational schema  $R$  (or simply schema) is a set of relational symbols (or predicates), each with its associated arity. A position  $r[i]$ , in a schema  $R$ , is identified by a predicate  $G \in R$  and its  $i$ -th argument. A term  $t$  is a constant, null, or variable. An

atomic formula (or simply **atom** (**nguyên tử, nguyên tử**)) has the form  $r(t_1, \dots, t_n)$ , where  $r$  is an  $n$ -ngôi predicate, and  $t_1, \dots, t_n$  are terms. For an atom  $a$ , we denote  $\text{dom}(a)$  and  $\text{var}(a)$  the set of its constants and variables, respectively; these notations are naturally extended to sets of atoms. Conjunctions of atoms are often identified with the sets of their atoms. A relational instance (or simply instance)  $I$  for a schema  $R$  is a (possibly infinite) set of atoms of the form  $r(t)$ , where  $r$  is an  $n$ -ary predicate of  $R$ , and  $t \in (r \cup rN)^n$ . We denote as  $r(I)$  the set  $\{t \mid r(t) \in I\}$ . A database  $D$  for a schema  $R$  is a finite relational instance for  $R$  such that  $\text{dom}(D) \subseteq r$ .

A conjunctive query (CQ)  $q$  of arity  $n$  over a schema  $R$ , written as  $q/n$ , is an assertion the form  $p(X) \wedge \forall y \in (X, Y)$ , where  $y \in (X, Y)$ , called the body of  $q$  and denoted  $\text{body}(q)$ , is a conjunction of atoms over  $R$  having as arguments constants of  $r$  and variables of  $rV$ , and  $p$  is an  $n$ -ary predicate not occurring in  $R$ . A Boolean conjunctive query (BCQ) is a CQ of arity zero. Answers to queries are defined via homomorphisms. A homomorphism  $h$  from a set of atoms  $A_1$  to a set of atoms  $A_2$  is a mapping  $h : r \cup rN \cup rV \rightarrow r \cup rN \cup rV$  such that:  $t \in r$  implies  $h(t) = t$ , and  $r(t_1, \dots, t_n) \in A_1$  implies  $r(h(t_1), \dots, h(t_n)) \in A_2$ . The answer to a CQ  $q/n$  over an instance  $I$ , denoted  $q(I)$ , is the set of all  $n$ -tuples  $t \in r^n$  for which there exists a homomorphism  $h$  such that  $h(\text{body}(q)) \subseteq I$  and  $h(X) = t$ . A BCQ has only the empty tuple  $\{\}$  as possible answer, in which case

we say it has a positive answer. Formally, a BCQ has positive answer over  $I$ , denoted  $I \models q$ , iff  $\exists \theta q(I, \theta)$ , or, equivalently,  $q(I) \neq \emptyset$ .

### 3 Datalog

Datalog (see, e.g., [1,14]) has been used as a paradigmatic database programming and query language for over three decades. While it is rarely used directly as a query language in corporate application contexts, Datalog has influenced the development of popular query languages such as SQL, whose newer versions allow one to express recursive queries.

#### 3.1 Syntax

A Datalog rule  $p$  is an expression of the form  $head(p) \leftarrow body(p)$  where  $n \geq 0$ ,  $a_0, \dots, a_n$  are atoms over a relational schema which contain constants of  $r$  and variables of  $IV$ , and each variable occurring in  $a_0$  must appear in at least one of  $a_1, \dots, a_n$ . The atom  $a_0$  is called the head, denoted  $head(p)$ , while the set of atoms  $\{a_1, \dots, a_n\}$  is called the body, denoted  $body(p)$ . In other terms, a Datalog rule is a function-free Horn clause. A Datalog program  $P$  is a finite set of Datalog rules. The set of constants occurring in  $P$  is denoted  $dom(P)$ . An extensional predicate is a predicate that does not occur in the head of any rule of  $P$ , while an intensional predicate is a predicate that occurs in the head of some rule of  $P$ . The extensional (database) schema of  $P$ , denoted  $EDB(P)$ , consists of all the extensional predicates of  $P$ , whose values are given via an input database, while the intensional schema of  $P$ , denoted  $IDB(P)$ ,

consists of all the intensional predicates of P, whose values are computed by the program. The schema of P, written SCH(P), is the set of predicates EDB(P) U IDB(P). As we shall see, the semantics of a Datalog program is a mapping from databases for EDB(P) to databases for IDB(P).

Example 3. Consider the Datalog program Pgraph constituted by  
sp.reachable(X) sp-node(X)  
sp-reachable (Y) edge{X,Y), sp-reachable(X),

which takes as input EDB a directed graph given by a binary edge relation (hệ thức), plus a set of special nodes of this graph given by a unary relation sp-node. Clearly,

EDB{Pgraph) = {edge, sp-node} and IDB(Pgraph) = {sp-reachable}.

The above (recursive) program computes the set sp-reachable of all vertices in the graph which are reachable from special nodes. ■

### 3.2 Semantics

An interesting and elegant property of Datalog is the fact the there are three different but equivalent approaches to defining its semantics: a model-theoretic approach where the Datalog rules are considered as logical sentences asserting a property of the desired result, a fixpoint (điểm cố định, dấu chấm cố định) approach where the semantics are defined as a particular solution of a fixpoint equation, and a proof-theoretic approach which is based on obtaining proofs of facts. In the sequel, we discuss the model-theoretic and fixpoint semantics; for details on the proof-theoretic approach we refer the interested

$sp\_reachable(X) \leftarrow sp\_node(X)$   
 $sp\_reachable(Y) \leftarrow edge(X, Y), sp\_reachable(X),$

$EDB(P_{graph}) = \{edge, sp\_node\}$  and  $IDB(P_{graph}) = \{sp\_reachable\}$ .

reader to [1, Chapter 12, Section 12.4].

**Model-Theoretic Semantics.** The idea underlying this approach is to consider the given Datalog program as a set of first-order sentences (also called a first-order theory) which describes the desired outcome of the program. More precisely, to a rule  $p : a_0 \leftarrow a_1; \dots, a_n$  we associate the first-order sentence

$$\forall X_1 \dots \forall X_m (a_0 \wedge \dots \wedge a_n)$$

where  $X_1, \dots, X_m$  are the variables occurring in  $p$ . A database  $D$  satisfies  $p$ , denoted  $D \models p$ , if for each homomorphism  $h, \{h(a_1), \dots, h(a_n)\} \subseteq D$  implies  $h(a_0) \in D$ . The conjunction of the first-order sentences associated to the rules of a program  $P$  is denoted  $SP$ . A model of  $P$  is a database for  $SCH(P)$  that satisfies  $SP$ . The semantics of  $P$  on an input database  $D$ , denoted  $P(D)$ , is the unique C-minimal model of  $P$  containing  $D$ . Given an  $n$ -ary CQ  $q : p(X) \wedge p(X, Y)$  over  $SCH(P)$ , the answer to  $q$  w.r.t.  $P$  and  $D$  is the set of  $n$ -tuples  $\{t \mid t \in q(P(D))\}$ .

**Example 4.** Consider the program  $P_{graph}$  given in Example 3, and the database

$$D = \{edge(v_1, v_3), edge(v_2, v_3), edge(v_3, v_4), edge(v_4, v_5), edge(v_5, v_3), sp\_node\{v_1\}\}$$

for the schema  $EDB(P_{graph})$ ; the graph  $G$  encoded by  $D$  is depicted in Figure 1, where the special node is shaded. The BCQ  $q$  over  $SCH(P)$

$$ans() \wedge edge(X, Y), edge(Y, Z), edge(Z, X), sp\_reachable(X), sp\_reachable(Y), sp\_reachable(Z),$$

$$\forall X_1 \dots \forall X_m (a_0 \leftarrow a_1 \wedge \dots \wedge a_n),$$

$$D = \{edge(v_1, v_3), edge(v_2, v_3), edge(v_3, v_4), edge(v_4, v_5), edge(v_5, v_3), sp\_node(v_1)\}$$

$$ans() \leftarrow edge(X, Y), edge(Y, Z), edge(Z, X), sp\_reachable(X), sp\_reachable(Y), sp\_reachable(Z),$$



asks whether a directed triangle in G is reachable from the set of special nodes of G. It is easy to verify that

$$P_{graph}(D) = D \cup \bigcup_{i \in \{1,3,4,5\}} \{sp\_reachable(v_i)\}.$$

It is straightforward to see that there exists a homomorphism that maps body(q) to Pgraph {D}:  $P_{graph}(D) \models q$ . ■

It is important to say that for a given program P and a database D for EDB (P), there exists always a unique C-minimal model of SP containing D, i.e., P(D) exists. Moreover, it is always possible to construct it. In particular, P(D) is the unique database for SCH(P) which satisfies the following conditions:

- $P(D) = \{o \mid \text{for each model } M \text{ of } P \text{ containing } D, o \in M\}$ ,
- $dom(P(D)) \subseteq dom(P) \cup dom(D)$ ,
- for each  $r \in EDB(P)$ ,  $r(P(D)) = r(D)$ .

The above conditions actually provide an algorithm for computing the semantics of Datalog programs. However, this is clearly an inefficient procedure. A more reasonable algorithm is provided by the fixpoint semantics.

Fixpoint Semantics. The fixpoint semantics of Datalog programs relies on an operator called the immediate consequence operator. In fact, this operator produces new facts starting from known facts. It is possible to show that, given a Datalog program P and a database D, P(D) can be defined as the smallest solution of a fixpoint equation involving that operator. This approach can be seen

$$P_{graph}(D) = D \cup \bigcup_{i \in \{1,3,4,5\}} \{sp\_reachable(v_i)\}.$$

as an implementation of the model-theoretic semantics.

Let  $P$  be a Datalog program and  $D$  a database for  $SCH(P)$ . A fact  $a$  is an immediate consequence for  $D$  and  $P$  if either  $a \in G_r(D)$  for some predicate  $r \in G_{EDP}(P)$ , or there exists a rule  $a_0 \leftarrow a_1; \dots, a_n$  in  $P$  and a homomorphism  $h$  such that  $h(\{a_1; \dots, a_n\}) \subseteq D$  and  $a = h(a_0)$ . The immediate consequence operator of  $P$ , denoted  $TP$ , is the mapping from databases for  $SCH(P)$  to databases for  $SCH(P)$  defined as follows: for each database  $D$ ,  $Tp(D) = \{a \mid a \text{ is an immediate consequence for } D \text{ and } P\}$ . We write  $Tp^i(D)$  for the result obtained by applying  $i$  times  $TP$  starting from  $D$ . Formally,  $T_0(D) = D$  and  $T_{i+1}(D) = TP(Tp^i(D))$ ; let  $Tp^*(D) = \bigcup_{i \geq 0} T_{i+1}(D)$ . From the fact that  $D \subseteq T_{i+1}(D)$  and the monotonicity of  $TP$ , it is clear that  $Tp^i(D) \subseteq T_{i+1}(D)$ , for each  $i \geq 0$ . It is well-known that for a Datalog program  $P$  and a database  $D$  for  $EDB(P)$ ,  $Tp^*(D)$  is the minimum fixpoint of  $TP$  containing  $D$  which coincides with  $P(D)$ . Interestingly,  $Tp^*(D)$  - and thus  $P(D)$  - can be obtained by applying finitely many times  $TP$ . More precisely, given a Datalog program  $P$ , for each database  $D$  there exists an integer  $k_D$  (which depends on  $D$ ) such that  $Tp^*(D) = T_{k_D}(D) = P(D)$ . The smallest such integer is called the stage for  $D$  and  $P$ , denoted  $stage(D, P)$ .

Example 5. Recall the program  $P_{graph}$  given in Example 3, and the database  $D$  given in Example 4. Clearly,

$$T_{P_{graph}}(D) = D \cup \{sp.reachable(v_i)\}$$

$$T_{P_{graph}}^2(D) = T_{P_{graph}}(D) \cup \{up-$$

$\text{reachable}(v_3)\}$   
 $\text{TPSruPh}^{\wedge}D) = \text{TPSr}^*\text{ph}(D) \cup \{\text{sp\_reachable}(v_4)\}$   
 $\text{TPSruPh}^{\wedge}D) = \text{TP}^9\text{ruPh}^{\wedge}D) \cup \{\text{sp\_reachable}(v_5)\}$   
 $\text{ThpPh}(D) = \text{TPSaPhh}(D).$

Hence,  $\text{Tp}(D) = \text{Tp}(D)$ . Notice that  $\text{Pgraph}\{D\} = \text{Tp}(D)$ .

### 3.3 Complexity of Datalog

We conclude this section by discussing the complexity of query answering under Datalog programs. Formally, given a Datalog program  $P$ , a database  $D$  for  $\text{EDB}(P)$ , an  $n$ -ary CQ  $q$  over  $\text{SCH}(P)$ , and an  $n$ -tuple  $t \in G^n$ ,  $\text{CQAns}$  is defined as the problem of deciding whether  $t \in G^n$  ( $q(P(D))$ ). In case that  $q$  is a BCQ (and thus  $t$  is the empty tuple), the above problem is called  $\text{BCQAns}$ . Following Vardi's taxonomy [37], the data complexity of the above problems is the complexity calculated taking only the database as input, while the Datalog program and the query are fixed. The combined complexity is the complexity calculated considering as input, together with the database, also the program and the query. It is well-known that  $\text{CQAns}$  and  $\text{BCQAns}$  are  $\text{LOGSPACE}$ -equivalent; this result is implicit in [15], and stated explicitly in [8]. Henceforth, for brevity, we focus only on  $\text{BCQAns}$ , and all complexity results carry over to  $\text{CQAns}$ .

Given a Datalog program  $P$  and a database for  $\text{EDB}(P)$ , the total number of atoms that can appear in  $P(D)$  is bounded by  $|\text{SCH}(P)| \cdot (|\text{Idom}(D) \cup \text{dom}(P)|)^w$ ,

$$\begin{aligned}
 T_{P_{\text{graph}}}(D) &= D \cup \{\text{sp\_reachable}(v_1)\} \\
 T_{P_{\text{graph}}}^2(D) &= T_{P_{\text{graph}}}(D) \cup \{\text{sp\_reachable}(v_3)\} \\
 T_{P_{\text{graph}}}^3(D) &= T_{P_{\text{graph}}}^2(D) \cup \{\text{sp\_reachable}(v_4)\} \\
 T_{P_{\text{graph}}}^4(D) &= T_{P_{\text{graph}}}^3(D) \cup \{\text{sp\_reachable}(v_5)\} \\
 T_{P_{\text{graph}}}^5(D) &= T_{P_{\text{graph}}}^4(D).
 \end{aligned}$$



$$|\text{SCH}(P)| \cdot (|\text{dom}(D) \cup \text{dom}(P)|)^w,$$

where  $w$  is the maximum arity over all predicates of SCH (P). Clearly, if the program is fixed (resp., non-fixed), then P(D) can be constructed by applying polynomially (resp., exponentially) many times the operator TP, while each application is feasible in polynomial (resp., exponential) time. This implies that query answering under Datalog programs is in polynomial time in data complexity, and in exponential time in combined complexity.

Interestingly, the above upper bounds are also lower bounds. As shown in [18], the PTIME-hardness can be established by constructing a simple Datalog metainterpreter for the class of propositional Logic Programs whose clauses have at most three body-atoms, denoted LP(3). More precisely, given an LP(3) program P, a database DP can be constructed in LOGSPACE as follows: each rule  $a$  is encoded by an atom  $t(a)$ , while each rule  $a_0 \leftarrow a_1; \dots; a_i$  is encoded by an atom  $r_j(a_0, \dots, a_j)$ , where  $r^*$  is an  $(i + 1)$ -ary predicate. Then, P which is actually stored as described above in DP can be evaluated by a fixed Datalog program PMI constructed as follows: for each  $i \in [3]$ , add a rule  $t(X_0) \wedge t(X_1), \dots, t(X_i), r_i(X_0, \dots, X_i)$ .

Intuitively,  $t(b)$  means that the propositional atom  $b$  is true. It is not difficult to see that, given a propositional atom  $a$ ,  $P \models a$  iff  $\text{PMI}(D_p) \models q$ , where  $q$  is the BCQ  $t(a)$ . The EXPTIME-hardness can be established by simulating the behavior of an exponential time Turing machine by means of a Datalog program. In fact, this holds

$$t(X_0) \leftarrow t(X_1), \dots, t(X_i), r_i(X_0, \dots, X_i).$$

even if the input database is empty, and only the constants 0 and 1 are allowed in the program; for details we refer the reader to [18].

From the above discussion we immediately get the following complexity characterization of query answering under Datalog programs.

Theorem 1. BCQAns under Datalog programs is PTIME-complete in data complexity, and EXPTIME-complete in combined complexity.

An interesting (semantic) property of Datalog programs is boundedness. A Datalog program is bounded if there exists a constant  $k$  such that, for each database  $D$ ,  $\text{stage}(D, P) \leq k$ ; recall that in general  $\text{stage}(D, P)$  depends on  $D$ . If a program is bounded, then it is actually non-recursive, even if syntactically appears to be recursive. It is well-known that every bounded Datalog program can be transformed into an equivalent non-recursive program. For example, the (recursive) program

$$\text{buys}(X, Y) \leftarrow \text{trendy}(X), \text{buys}(Z, Y)$$
$$\text{buys}(X, Y) \leftarrow \text{likes}(X, Y)$$

is bounded, and the atom  $\text{buys}(Z, Y)$  in the body of the first rule can be replaced by  $\text{likes}(Z, Y)$ .

Given a bounded Datalog program  $P$  and a BCQ  $q : p \wedge \exists X$  over  $\text{SCH}(P)$ , which can be seen as a Datalog rule, the program  $P \cup \{p \wedge \exists X\}$  is still bounded since  $p \in \text{SCH}(P)$ . This observation implies that a non-recursive Datalog program  $Pq$  can be constructed such that  $P(D) = q$  iff  $p \in Pq(D)$ , for every database  $D$  for  $\text{EDP}(P)$ . The problem of deciding whether the propositional

*buys*(X, Y) ← *trendy*(X), *buys*(Z, Y)  
*buys*(X, Y) ← *likes*(X, Y)

atom  $p$  belongs to  $Pq$  ( $D$ ) can be easily reduced to the problem of evaluating a positive existential first-order sentence, obtained by “unfolding” (trăi ra) the intensional predicates of  $Pq$ , over  $D$ ; for details see [1]. For example, if  $Pq$  is the program

$p \wedge s(X), r(X, Y), t(Y, Z) \wedge r(X, Y) \wedge t(Y, X), u(X), v(Y, Z) \wedge r(X, Y) \wedge t(X, Y)$ , then the first-order sentence that must be evaluated is

$\exists X \exists Y \exists Z s(X) \wedge ((\exists Z' t(Y, X) \wedge u(X) \wedge v(Y, Z')) \vee t(X, Y)) \wedge t(Y, Z)$ .

Since evaluation of first-order sentences is feasible in  $AC_0$  in data complexity [40], the next result follows immediately.

**Theorem 2.** BCQAns under bounded Datalog programs is in  $AC_0$  in data complexity.

As said, boundedness is a semantic (and not a syntactic) property. In fact, the problem of deciding whether a Datalog program is bounded is undecidable (see, e.g., [24]). This implies that only sufficient syntactic conditions, which will not detect boundedness in some cases, can be defined. Decidability results for particular subclasses are established, e.g., in [17,38].

#### 4 Datalog Extensions

Datalog was designed and intended for reasoning over finite databases, assuming that only the values explicitly mentioned in the extensional database exist. For ontological reasoning, however, it would be desirable that an extended version of Datalog could be able to express the existence of certain values that are not necessarily from the EDB domain. This can be

$\exists X \exists Y \exists Z s(X) \wedge ((\exists Z' t(Y, X) \wedge u(X) \wedge v(Y, Z')) \vee t(X, Y)) \wedge t(Y, Z)$ .

achieved by allowing existentially quantified variables in rule heads. Other useful features for representing ontologies, which are not expressible by Datalog rules, are functionality and disjointness constraints. Such assertions can be supported by allowing the equality predicate and the truth constant false to appear in rule heads.

Example 6. As said, ontology languages are based on DLs. Such formalisms can express, for instance, that every person has exactly one father/mother who, moreover, is himself/herself a person, and also that fathers and mothers are disjoint sets, by the following axioms:  $\text{Person} \subseteq \exists \text{Father}$ ,  $\text{Person} \subseteq \exists \text{Mother}$ ,  $\exists \text{Father} \sim \text{C Person}$ ,  $\exists \text{Mother} \sim \text{C Person}$ ,  $\exists \text{Father} \sim \text{disj } \exists \text{Mother} \sim$ , (funct Father) and (funct Mother). In an appropriate extended version of Datalog, the same can be expressed as:

$\exists Y \text{ father}(X, Y) \wedge \text{person}(X)$   
 $\exists Y \text{ mother}(X, Y) \wedge \text{person}(X)$   
 $\text{person}(Y) \wedge \text{father}(X, Y) \text{ person}(Y)$   
 $\wedge \text{mother}(X, Y)$   
 $L \wedge \text{father}(X, Y), \text{mother}(Z, Y)$   
 $Y = Z \wedge \text{father}(X, Y), \text{father}(X, Z).$

Observe that all the predicates occurring in the above program appear both in the body and in the head of some rule. Therefore, in the desired extended version of Datalog, we no longer require the distinction between extensional and intensional predicates. ■

Recently a family of Datalog-based languages, called Datalog $_{\pm}$ , which is a new framework for tractable ontology querying, has been introduced [10]. The precise aim of

$\exists Y \text{ father}(X, Y) \leftarrow \text{person}(X)$   
 $\exists Y \text{ mother}(X, Y) \leftarrow \text{person}(X)$   
 $\text{person}(Y) \leftarrow \text{father}(X, Y)$   
 $\text{person}(Y) \leftarrow \text{mother}(X, Y)$   
 $\perp \leftarrow \text{father}(X, Y), \text{mother}(Z, Y)$   
 $Y = Z \leftarrow \text{father}(X, Y), \text{father}(X, Z).$

Datalog $\pm$  is to extend Datalog with existential quantification, the equality predicate, and the truth constant false, while preserving, not only decidability, but also tractability of query answering in data complexity. In view of the fact that query answering under Datalog extended with existential quantification is already undecidable [7], certain syntactic restrictions were needed (hence the symbol  $\pm$ ).

#### 4.1 Syntax of Datalog $_{\exists,=,\pm}$

A Datalog $_{\exists,=,\pm}$  rule  $p$  is an expression of the form

$h \text{ au } \dots, a_n,$

where  $a_1; \dots, a_n$  are atoms over a relational schema which contain constants of  $r$  and variables of  $rV$ , and  $h$  is either

— an expression  $\exists Y_1 \dots \exists Y_m a_0$ , where  $\{Y_1, \dots, Y_m\} \subseteq IV$ ,  $\{a_1, \dots, a_n\}$  and  $\text{var}(\{a_1; \dots, a_n\})$  are disjoint, and  $a_0$  is an atom which contains constants of  $r$  and variables of  $(\{Y_1, \dots, Y_m\} \cup \text{var}(\{a_1; \dots, a_n\}))$ , or

— an expression  $X = Y$ , where  $\{X, Y\} \subseteq \text{var}(\{a_1; \dots, a_n\})$ , or

— the symbol  $L$  which denotes the truth constant false.

If  $h$  is either an equality among variables or the symbol  $L$ , then  $n \geq 1$ ; otherwise,  $n \geq 0$ . The expression  $h$  is called the head, denoted  $\text{head}(p)$ , while the set of atoms  $\{a_1, \dots, a_n\}$  is called the body, denoted  $\text{body}(p)$ . A Datalog $_{\exists,=,\pm}$  program  $P$  is a finite set of Datalog $_{\exists,=,\pm}$  rules. The schema of  $P$ , written  $\text{SCH}(P)$ , is the set of predicates occurring in  $P$ . Notice that  $\text{SCH}(P)$  is not partitioned, as in plain Datalog programs, into extensional and

$$h \leftarrow \underline{a}_1, \dots, \underline{a}_n,$$



intensional predicates. Obviously, the program given in Example 6 is a Datalog<sub>3,=,±</sub> program.

#### 4.2 Semantics of Datalog<sub>3,=,x</sub>

**A Model-Theoretic Approach.** To define the semantics of Datalog<sub>3,=,±</sub> we follow an approach similar to that of the model-theoretic semantics of Datalog, i.e., we consider the given program as a first-order theory. More precisely, to a rule  $p : h \ a_1; \dots, \ a_n$ , where  $\text{var}(\text{body}(p)) = \{X_i, \dots, \ X_j\}$ , if  $h$  is of the form  $\exists Y_1 \dots \exists Y_m \ a_0$ , then we associate the first-order sentence  $\forall X_1 \dots \forall X_k \exists Y_1 \dots \exists Y_m (a_0 \leftarrow a_1 \wedge \dots \wedge a_n)$ .

An instance  $I$  satisfies  $p$  if for each homomorphism  $\theta$ ,  $\{p(a_i), \dots, \ p(a_n)\} \subseteq I$  implies that there exists an extension  $\pi$  of  $\theta$ , i.e.,  $\pi \supseteq \theta$ , such that  $\pi(a_0) \in I$ . Sentences of the above form are known as tuple-generating dependencies (TGDs) in the database literature - essentially they say that the presence of some tuples in the instance implies the presence of some other tuple in the same instance. Notice that in general the head of a TGD can be a conjunction of atoms, and not just a single atom. In this way, however, we do not extend the expressive power of the language, since a set of TGDs can be reduced in LOGSPACE into an equivalent set of TGDs with a single atom as a head [8]. Now, if  $h$  is of the form  $X = Y$ , where  $\{X, Y\} \subseteq \{X_1, \dots, \ X_k\}$ , then we associate to  $p$  the sentence  $\forall X_1 \dots \forall X_k (X = Y \leftarrow a_1 \wedge \dots \wedge a_n)$ .

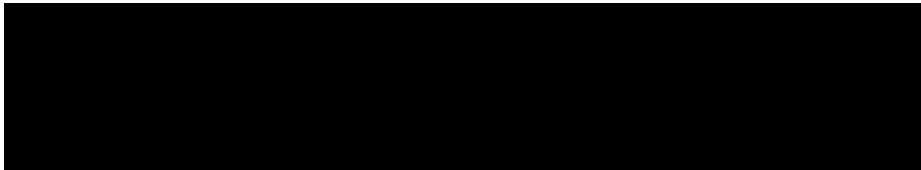
An instance  $I$  satisfies  $p$  if for each homomorphism  $\theta$ ,  $\{\theta(a_1), \dots, \ \theta(a_n)\} \subseteq I$  implies  $\theta(X) = \theta(Y)$ . Sentences of the above form are known as



$$\forall X_1 \dots \forall X_k \exists Y_1 \dots \exists Y_m (a_0 \leftarrow a_1 \wedge \dots \wedge a_n).$$



$$\forall X_1 \dots \forall X_k (X = Y \leftarrow a_1 \wedge \dots \wedge a_n).$$



equality- generating dependencies in the database literature - basically they assert that the presence of some tuples in the instance implies that certain tuple components are equal. Finally, if  $h = \pm$ , then we associate to  $p$  the sentence

$\forall X_1 \dots \forall X_k (\perp \leftarrow a_1 \wedge \dots \wedge a_n).$   
 An instance  $I$  satisfies  $p$  if for each homomorphism  $p$ ,  $\{p(a_1), \dots, p(a_n)\} \subseteq I$ . Such sentences are known as negative constraints, and actually they forbid the existence of some tuples in the instance.

The conjunction of the first-order sentences associated to the rules of a Datalog $_{\exists, =, \pm}$  program  $P$  is denoted  $SP$ . A model of  $P$  is an instance for  $SCH(P)$  that satisfies  $SP$ . The semantics of  $P$  on an input database  $D$ , denoted  $P(D)$ , is an instance  $M$  for  $SCH(P)$  such that  $M \subseteq D$ ,  $M$  is a model of  $P$ , and for each model  $M'$  of  $P$  containing  $D$ , there exists a homomorphism  $p$  such that  $p(M) \subseteq M'$ ; such an instance is called universal model of  $P$  w.r.t.  $D$ . Intuitively speaking, a universal model it contains no more and no less information than what the given program requires. It is not difficult to verify that, in general, there exist more than one universal models of  $P$  w.r.t.  $D$ .

However, by definition, they are the same up to homomorphic equivalence, i.e., for each pair of universal models  $M_1$  and  $M_2$ , there exist homomorphisms  $h_1$  and  $h_2$  such that  $h_1(M_1) \subseteq M_2$  and  $h_2(M_2) \subseteq M_1$ . Thus,  $P(D)$  is unique up to homomorphic equivalence. Given an  $n$ -ary CQ  $q : p(X) \wedge \exists Y \theta(X, Y)$  over  $SCH(P)$ , the answer to  $q$  w.r.t.  $P$  and

$$\forall X_1 \dots \forall X_k (\perp \leftarrow a_1 \wedge \dots \wedge a_n).$$

D is the set of n-tuples  $\{t \mid t \in q(P(D))\}$ . It is important to say that for query answering, homomorphically equivalent instances are indistinguishable, i.e., given two instances I and I' which are the same up to homomorphic equivalence,  $q(I)$  and  $q(I')$  coincide. From the above discussion we conclude that, for query answering purposes,  $P(D)$  is defined uniquely.

**A Fixpoint Approach for Datalog3.** Interestingly, analogous to the fixpoint semantics of Datalog programs, given a Datalog3 program P - in the rest of the paper, we denote Datalog3 the formalism obtained by allowing only existential quantification in rule heads - and an input database D,  $P(D)$  can be (constructively) defined as the least fixpoint of a monotonic operator (modulo homomorphic equivalence). This can be achieved by exploiting the well-known chase procedure, introduced for checking implication of dependencies [31], and later for checking query containment [28]. The chase is an algorithm that, roughly speaking, executes the rules of P starting from D in a forward chaining manner by inferring new atoms, and inventing new null values whenever an existential quantifier needs to be satisfied. The formal definition follows.

Consider a Datalog3 rule  $p : \exists Y_1 \dots \exists Y_m a_0 \leftarrow a_1, \dots, a_n$ , and an instance I. We say that p is applicable to I if there exists a homomorphism h such that  $h(\{a_1, \dots, a_n\}) \subseteq I$ , but there is no  $h' \subseteq h$  such that  $h'(a_0) \in I$ . Let I' be the instance  $I \cup \{h'(a_0)\}$ , where  $h' \subseteq h$  is such that

$h(Y_i)$  is a fresh labeled null of  $rN$  not occurring in  $I$ , for each  $i \in [m]$ , and  $h(Y_i) = h(Y_j)$ , for each pair  $(i, j)$  of distinct integers of  $[m]$ . We say that the result of applying  $p$  to  $I$  with  $h$  is  $I'$ ,

and write  $I \rightarrow I'$ ; in fact,  $I \rightarrow I'$  defines one single chase step. Consider now a Datalog<sup>3</sup> program  $P$ , and an input database  $D$ . A chase sequence of  $D$  w.r.t.  $P$  is a sequence of chase steps  $I_i \rightarrow I_{i+1}$ , where  $i \geq 0$ ,  $I_0 = D$ , and  $p_i \in P$ . The result of the chase of  $D$  w.r.t.  $P$ , denoted  $\text{chase}(D, P)$ , is defined as follows:

— A finite chase of  $D$  w.r.t.  $P$  is a finite chase sequence  $I_i \rightarrow I_{i+1}$ , where  $0 \leq i < m$ , and there is no  $p \in P$  which is applicable to  $I_m$ ; let  $\text{chase}(D, P) = I_m$ .

— An infinite chase sequence  $I_i \rightarrow I_{i+1}$ , where  $i \geq 0$ , is **fair (công bằng, hợp lý)** if whenever a rule  $p \in P$  is applicable to  $I_i$  with homomorphism  $h$ , then there exist  $h' \in D$  and  $k > i$  such that  $h'(\text{head}(p)) \in I_k$ . An infinite chase of  $D$  w.r.t.  $P$  is a fair infinite chase sequence  $I_i \rightarrow I_{i+1}$ , where  $i \geq 0$ ; let  $\text{chase}(D, P) = \bigcup_{i=0}^{\infty} I_i$ .

A useful technical notion is the depth of the chase. Roughly, the lower the depth of an atom, the earlier the atom has been obtained during the construction of  $\text{chase}(D, P)$ . For an atom  $a \in D$ , let  $\text{depth}(a) = 0$ , and for an atom  $a \in \text{chase}(D, P)$  obtained during the chase application  $I_i \rightarrow I_{i+1}$ , let  $\text{depth}(a) =$

$\max_{f \in \text{fl}(\text{body}(p))} \{\text{depth}(b)\} + 1$ . The chase of  $D$  w.r.t.  $P$  up to depth  $k \geq 0$ , denoted  $\text{chase}_k(D, P)$ , is defined as the instance  $\{a \mid a \in \text{chase}(D, P) \text{ and } \text{depth}(a) \leq k\}$ .

Example 7. Consider the Datalog<sup>3</sup>

program P:

P1 :  $\exists Y \text{ father}(X, Y) \wedge \text{person}(X)$  p2 :  $\text{person}(Y) \wedge \text{father}(X, Y)$ ,

which is a subset of the program given in Example 6, and the input database  $D = \{\text{person}(\text{john})\}$ . A fair infinite chase of  $D$  w.r.t.  $P$  follows:

$D$   
 $D \cup \{\text{father}(\text{john}, z_1)\}$   
 $D \cup \{\text{father}(\text{john}, z_1), \text{person}(z_1)\}$   
 $D \cup \{\text{father}(\text{john}, z_1), \text{person}(z_1), \text{father}(z_1, z_2)\}$

$D \cup \{\text{father}(\text{john}, z_1)\} \cup \{\text{father}(z_j, z_{j+1}), \text{person}(z_j)\}$   
 where  $z_1, z_2, \dots$  are labeled nulls of  $r_N$ . Thus,  $\text{chase}(D, P)$  is the infinite instance  $D \cup \{\text{father}(\text{john}, Z_1), \text{person}(Z_1)\} \cup \bigcup_{j=1}^{\infty} \{\text{father}(z_j, z_{j+1}), \text{person}(z_j)\}$  - ■

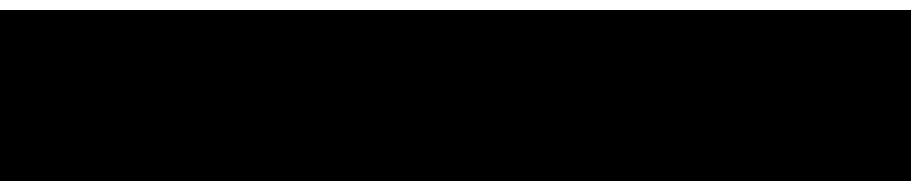
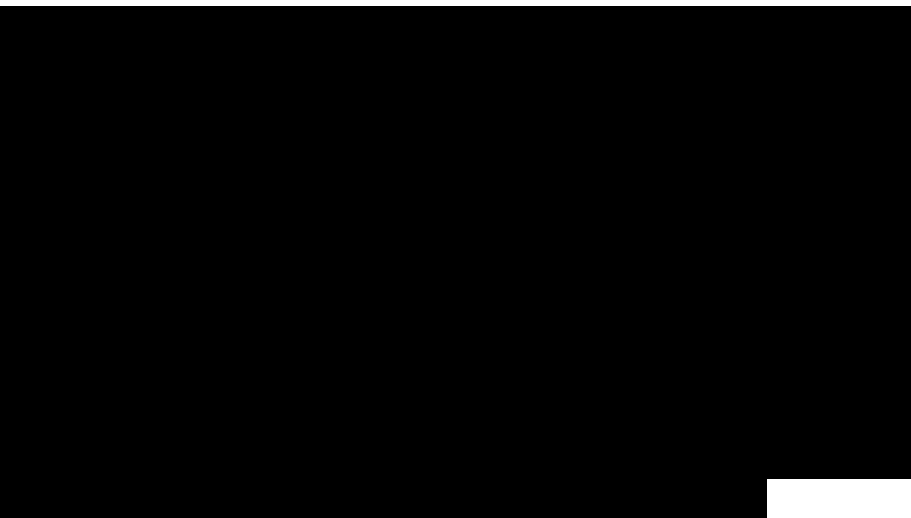
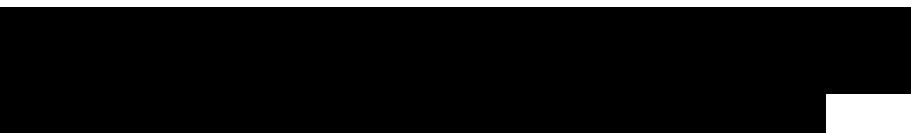
It turns out that the result of the chase of an input database  $D$  w.r.t. a Datalog<sup>3</sup> program  $P$  is a universal model of  $P$  w.r.t.  $D$  (see, e.g., [19,23]). This can be shown by induction on the number of applications of the chase step. In fact, it is possible to show that, for each  $k \geq 0$ , there exists a homomorphism  $h_k$  that maps the initial part of  $\text{chase}(D, P)$  obtained after  $k$  chase step applications to each model of  $P$  containing  $D$ . Moreover, it holds that  $h_0 \subseteq \dots \subseteq h_k \subseteq \dots$ , and therefore the homomorphism  $h_P$  maps  $\text{chase}(D, P)$  to each model of  $P$  containing  $D$ . Since, by definition,  $\text{chase}(D, P)$  is a model of  $P$  containing  $D$ , we conclude that  $\text{chase}(D, P)$  is a universal model of  $P$  w.r.t.  $D$ .

It is easy to verify that the result of the chase of  $D$  w.r.t.  $P$  is not unique since it depends on the order that the rules of  $P$  are executed. In other

$\rho_1 : \exists Y \text{ father}(X, Y) \leftarrow \text{person}(X)$      $\rho_2 : \text{person}(Y) \leftarrow \text{father}(X, Y)$ ,



$D$   
 $\rho_1, h_1 = \{X \rightarrow \text{john}\} \quad D \cup \{\text{father}(\text{john}, z_1)\}$   
 $\rho_2, h_2 = \{X \rightarrow \text{john}, Y \rightarrow z_1\} \quad D \cup \{\text{father}(\text{john}, z_1), \text{person}(z_1)\}$   
 $\rho_1, h_3 = \{X \rightarrow z_1\} \quad D \cup \{\text{father}(\text{john}, z_1), \text{person}(z_1), \text{father}(z_1, z_2)\}$   
 $\vdots$   
 $\rho_1, h_{2i-1} = \{X \rightarrow z_{i-1}\} \quad D \cup \{\text{father}(\text{john}, z_1)\} \cup \bigcup_{j=1}^{i-1} \{\text{person}(z_j), \text{father}(z_j, z_{j+1})\}$   
 $\vdots$



words, different chase sequences may yield different results. However, each result is a universal model of  $P$  w.r.t.  $D$ , and thus  $\text{chase}(D,P)$  is unique up to homomorphic equivalence. From the above discussion we conclude that  $P(D)$  and  $\text{chase}(D,P)$  coincide (up to homomorphic equivalence).

**The Challenge of Infinity.** Recall that for a Datalog program  $P$  and an input database  $D$ ,  $P(D)$  is finite and it is always possible to construct it. In fact, the fixpoint semantics of Datalog programs provide an efficient (w.r.t. the size of the data) algorithm, based on the immediate consequence operator, which constructs  $P(D)$ . Unfortunately, the situation changes dramatically if  $P$  is a Datalog<sub>3</sub> (and thus, a Datalog<sub>3,=,±</sub>) program. Due to the existentially quantified variables in rule heads,  $P(D)$  is in general infinite, and thus not explicitly computable. As illustrated in Example 7, the result of the chase procedure is in general infinite even for very simple Datalog<sub>3</sub> programs, and for extremely small input databases.

Already for Datalog<sub>3</sub> query answering is undecidable [7]. Worse than that, undecidability holds even if both the program and the query are fixed, and only the database is given as input [8]. It is thus necessary to identify expressive fragments of Datalog<sub>3,=,±</sub> for which query answering is decidable, and also tractable in data complexity. In what follows we present such fragments of Datalog<sub>3,=,±</sub>. First, in Subsection 4.3, we focus on Datalog<sub>3</sub>, and we

discuss how we can regain decidability, and also tractability in data complexity, of query answering. Then, in Subsection 4.5, we discuss how  $=$  can be safely added, and also that  $L$  can be easily treated.

### 4.3 Guarded Datalog<sub>3</sub>

Guardedness, proposed by Andreka et al. [2], is a well-known restriction of first-order logic that ensures decidability of satisfiability, i.e., the problem of deciding whether a first-order theory has at least one model. Inspired by the guarded fragment of first-order logic, guarded Datalog<sub>3</sub> has been proposed recently by Cali et al. [8,9]. A guarded Datalog<sub>3</sub> rule  $p$  is a Datalog<sub>3</sub> rule of the form  $h \wedge a_1, \dots, a_n$ , where at least one atom  $a$  in  $\text{body}(p)$  contains all the variables occurring in  $\text{body}(p)$ , i.e.,  $\text{var}(a) = \text{var}(\text{body}(p))$ . The rightmost such atom is called the guard of  $p$ . A guarded Datalog<sub>3</sub> program  $P$  is a finite set of guarded Datalog<sub>3</sub> rules. The sentences of the first-order theory  $SP$  are known as guarded TGDs [8,9]. Notice that  $SP$  falls in the guarded fragment of first-order logic.

The decidability of BCQAns follows from the fact that the result of the chase of a database w.r.t. to a guarded Datalog<sub>3</sub> program has finite treewidth, i.e., is a treelike structure; for more details we refer the interested reader to [8]. The data complexity of BCQAns under guarded Datalog<sub>3</sub> programs was investigated in [9]; in fact, it was shown that it is PTIME-complete. Let us briefly discuss how this result was obtained.

First, it was established that guarded

Datalog<sup>3</sup> enjoys the so-called bounded guard-depth property (BGDP). The chase graph of a database  $D$  w.r.t. a guarded Datalog<sup>3</sup> program  $P$  is a labeled directed graph  $(V, E, X)$ , where  $V$  is the node set,  $E$  is the edge set, and  $X$  is a node labeling function  $V \rightarrow \text{chase}(D, P)$ . For each atom  $a \in D$ , there exists a node  $v \in V$  such that  $X(v) = a$ . For each atom  $a \in \text{chase}(D, P)$  obtained during a chase step  $p_j \leftarrow p_{j-1}$ , where  $p_j$  is of the form  $\exists Y_1 \dots \exists Y_m a_0 \leftarrow a_1 \wedge \dots \wedge a_n$  and the guard of  $p_j$  is  $\bullet$ , there exist edges  $(v_1, u), \dots, (v_n, u)$ , where  $X(v_j) = a_j$  and  $X(u) = a_0$ ; the node  $v_j$  is marked as guard. The guarded chase forest of  $D$  w.r.t.  $P$  is the restriction of the chase graph of  $D$  w.r.t.  $P$  to all nodes marked as guards and their children.

Example 8. Consider the guarded Datalog<sup>3</sup> program  $P$ :

$p_1 : \exists Z r(Z, X) \wedge r(X, Y), s(Y)$   $p_2 : s(X) \wedge r(X, Y)$ ,

and the database  $D = \{r(a, b), s(b)\}$ .

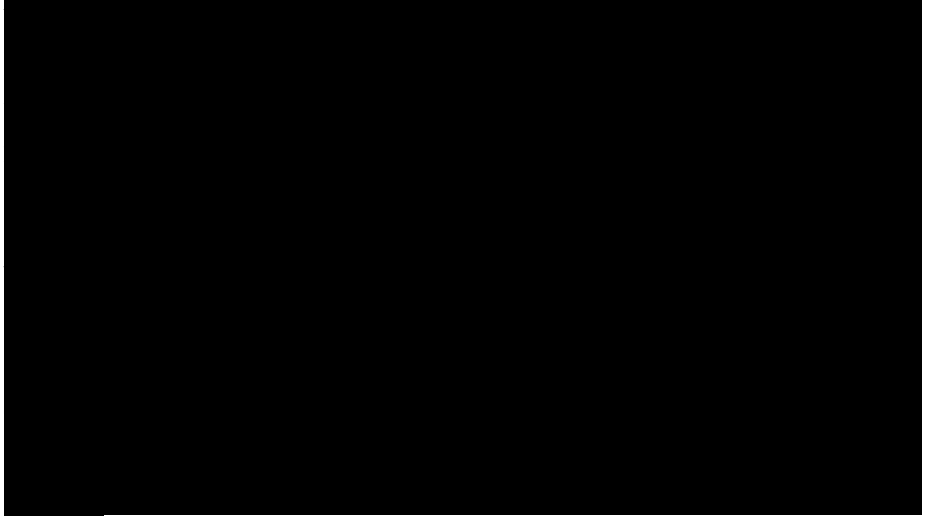
An initial part of the chase graph (resp., guarded chase forest) of  $D$  w.r.t.  $P$  is depicted in Figure 2(a) (resp., 2(b)). The

$r(a, b) \circ \quad s(b) \quad r(a, b) \circ \quad s(b)$

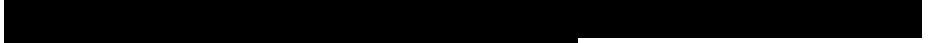
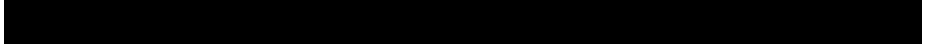
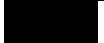
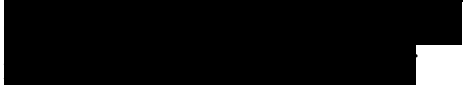
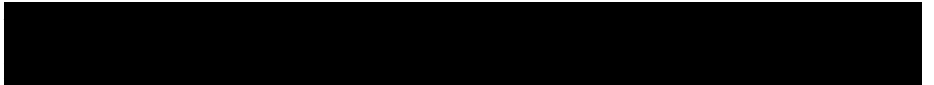
(a) (b)

Fig. 2. Chase graph and guarded chase forest for Example 8. Edges are labeled by the applied rules, and the number on the upper right side of each atom indicates its depth; formally not part of the graph (resp., forest). ■

The BGDP guarantees the following: whenever the given BCQ  $q$  is entailed by  $\text{chase}(D, P)$ , and thus



$\rho_1 : \exists Z r(Z, X) \leftarrow r(X, Y), s(Y) \quad \rho_2 : s(X) \leftarrow r(X, Y)$ ,

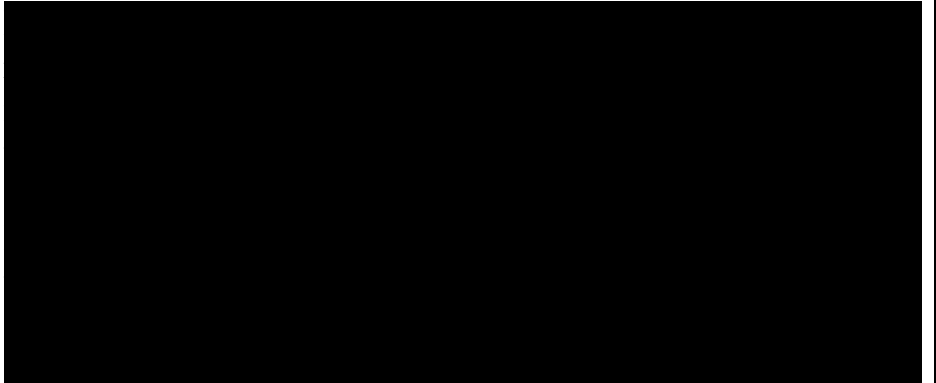




$P(D) = q$ , then the whole chase derivation of the query atoms is contained in a finite, initial part  $C$  of the guarded chase forest of  $D$  w.r.t  $P$ , whose depth depends only the query and the program, but not on the database. Thus, the BGD<sub>P</sub> allows us to construct the finite part  $C$  of the guarded chase forest (instead of considering the infinite instance  $\text{chase}(D,P)$ ), and then evaluate  $q$  over  $C$ . Finally, was shown that the construction of  $C$  can be done in polynomial time in the size of the data, and the PTIME upperbound of BCQAns under guarded Datalog<sub>3</sub> programs in data complexity follows. The PTIME-hardness of BCQAns under guarded Datalog<sub>3</sub> in data complexity follows immediately from the fact that the single rule used in the Datalog metainterpreter for LP(3) (see Subsection 3.3) is actually a guarded Datalog<sub>3</sub> rule (without existentially quantified variables). The combined complexity of BCQAns under guarded Datalog<sub>3</sub> has been investigated by Cali et al. in [8]; in fact, it is 2EXPTIME-complete.

Theorem 3. BCQAns under guarded Datalog<sub>3</sub> programs is PTIME-complete in data complexity, and 2EXPTIME-complete in combined complexity.

Guarded Datalog<sub>3</sub> strictly captures the DL ECHI, that is, the extension of EC with role hierarchies and inverse roles [4], without losing tractability of BCQAns in data complexity. In particular, in a normalized ECHI ontology we can have axioms of the form:  $A \sqsubseteq B$ ,  $A \sqsupseteq B$ ,  $A \sqsubseteq B \sqsubseteq C$ ,  $A \sqsubseteq B \sqsubseteq E$ ,  $A \sqsubseteq B \sqsubseteq E.B$ ,  $B \sqsubseteq A$ ,



BE.A Q B and E Q S, where A, B, C are atomic concepts, and E, S are basic roles. It is easy to verify that DL axioms of the above form can be translated into guarded Datalog<sub>3</sub> rules.

#### 4.4 Linear Datalog<sub>3</sub>

Interestingly, a (semantic) property analogous to boundedness of Datalog programs can be defined for Datalog<sub>3</sub> programs. A Datalog<sub>3</sub> program P has the bounded derivation-depth property (BDDP) if for every input database D, and for every BCQ q,  $P(D) = q$  implies  $\text{chase}_k(D, P) = q$ , where k depends only on P and q. As shown in [9], the problem of deciding whether  $\text{chase}_k(D, P) = q$  can be reduced to the problem of evaluating a first-order query Pq over D. Roughly, Pp can be constructed as follows. Suppose that  $P(D) = q$  which implies that  $\text{chase}_k(D, P) = q$ . Hence,  $\text{body}(q)$  is mapped to  $\text{chase}_k(D, P)$  via a homomorphism h. Since the derivation-depth and the number of body-atoms in rules of P are bounded, the cardinality of the set S of atoms at depth zero that are needed to entail the query is bounded. Thus, the number of all non-isomorphic sets S is also bounded. Consider the existentially quantified conjunction for every such set where q is answered positively. The first-order query Pq is the disjunction of all these formulas. Since evaluation of first-order queries is feasible in AC<sub>0</sub> in data complexity [40], the next result follows.

Theorem 4. BCQAns under Datalog<sub>3</sub> programs that have the BDDP is in

AC0 in data complexity.

Recall that BDDP is a semantic property. The question that comes up is whether we can identify a syntactic fragment of Datalog3 programs that enjoys the aforementioned property. Let us first say that every class of Datalog3 programs under which query answering is PTIME-hard in data complexity does not have the BDDP. Assume that such a class C has the BDDP. By Theorem 4, we get that BCQAns under C is in AC0 in data complexity, and thus AC0 coincides with PTIME. But this is a contradiction since AC0 C ptime (see, e.g., [34]). Consequently, guarded Datalog3 does not have the BDDP. If we further restrict guarded Datalog3 by allowing only one body-atom in the rules (which is trivially a guard), then we get the class linear Datalog3 which enjoys the BDDP. The PSPACE-completeness of BCQAns (in combined complexity) under linear Datalog3 is implicit in the seminal work of Johnson and Klug [28], where the problem of conjunctive query containment under inclusion dependencies has been investigated. The next result follows.

Theorem 5. BCQAns under linear Datalog3 programs is in AC0 in data complexity, and PSPACE-complete in combined complexity.

#### 4.5 Addition of = and $\pm$

We now discuss how guarded and linear Datalog3 can be enriched with rules that have the equality predicate and the truth constant false in the head, without increasing the complexity. For simplicity and technical clarity, we consider only

$\rho : X = Y \leftarrow \underline{a}_1, \underline{a}_2,$

equality rules of the form

$$P \cdot X == \wedge, (X \langle a_1, a_2 \rangle)$$

where  $a_1$  and  $a_2$  have the same predicate, there are no constants in  $p$ , each variable of  $\text{var}(a_j)$  occurs in  $\wedge$  at most once, and each variable of  $\text{var}(a_1)$  fl  $\text{var}(a_2)$  occurs in both atoms and  $a_2$  at the same position. Notice that rules of this form, despite their simplicity, they are expressive enough to capture functionality constraints. For example, the rule  $Y = Z \wedge \text{father}(X, Y), \text{father}(X, Z)$ , given in Example 6, which asserts that each person has at most one father, is exactly of this form. In the sequel, we call guarded/linear Datalog<sub>3</sub> the formalism obtained by combining guarded/linear Datalog<sub>3</sub> with equality rules of the above form.

It is known that query answering under linear Datalog<sub>3</sub> (and thus, also under guarded Datalog<sub>3</sub>) programs is undecidable; this is implicit in [16]. The reason for this is the interaction between existential and equality rules. Moreover, while the result of the chase procedure under Datalog<sub>3</sub> programs is well-defined as the least fixpoint of a monotonic operator (modulo homomorphic equivalence) - see Subsection 4.2 - this is not true if we consider also equality rules. Despite the fact that the chase procedure under Datalog<sub>3</sub> programs can be formally defined, it is not clear how the result of an infinite chase involving both existential and equality rules should be defined; the sequence of sets obtained in such an infinite chase is, in general, neither monotonic nor convergent.

Separability. It would be useful to identify syntactic restrictions which guarantee that, for query answering purposes, the semantics of Datalog<sub>3</sub>,= programs coincide with the semantics of Datalog<sub>3</sub>, providing that the input database satisfies the equality rules. In other words, we need sufficient syntactic conditions for the property of separability. Consider a Datalog<sub>3</sub>,= program P, and let P<sub>3</sub> (resp., P=) be the rules of P with an existential quantifier (resp., equality) in the head. P is separable if, for every input database D for SCH(P): if D = P=, then P(D) exists, and, for every BCQ q over SCH(P), P(D) = q iff P<sub>3</sub>(D) = q.

Separability allows us to answer queries under guarded/linear Datalog<sub>3</sub>,= programs by exploiting the algorithms for guarded/linear Datalog<sub>3</sub>. More precisely, given a program P, a database D for SCH(P), and a BCQ q for SCH(P), we can decide whether P(D) = q by applying the following algorithm: (1) if D = P=, then accept; (2) if P<sub>3</sub>(D) = q, then accept; otherwise, reject. Interestingly, as shown in [9], the problem of deciding whether the input database satisfies P= it is not harder than conjunctive query evaluation. In fact, this problem is feasible in ACo if P= is fixed, and in NP in general. From the above discussion, we conclude that the notion of separability allows us to transfer, not only the decidability results, but also the upper complexity bounds derived for guarded/linear Datalog<sub>3</sub>.

Theorem 6. BCQAns under separable guarded/linear Datalog<sub>3</sub>,

programs has the same data/combined complexity as BCQAns under guarded/linear Datalog3 programs.

Non-conflicting Condition.

Obviously, separability of Datalog3,= programs is a semantic (and not a syntactic) property. In what follows we present an efficiently checkable syntactic condition which is sufficient for separability of Datalog3,= programs. For brevity, given a rule  $p$  with an existential quantifier in the head, we define the set  $Up$  of universal positions of  $p$  as the set of positions in  $\text{head}(p)$  at which a variable of body ( $p$ ) appears. Moreover, given an equality rule  $p$ , we define the set  $Jp$  of joined positions of  $p$  as the set of positions at which a variable that occurs in both atoms of body ( $p$ ) appears; recall that we consider a special kind of equality rules with just two body-atoms. Formally, a Datalog3,= program  $P$ , where  $P_3$  (resp.,  $P_=$ ) are the rules of  $P$  with an existential quantifier (resp., equality) in the head, is non-conflicting if, for each  $(p_1, p_2) \in P_3 \times P_=$ , the following holds: (i)  $Up_1 \cap Jp_2 = \emptyset$ , and (ii) if  $Up_1 = Jp_2$ , then each existentially quantified variable in  $\text{head}(p_1)$  occurs just once.

Example 9. Consider the Datalog3,= program  $P$ :

$p_1 : \exists Y \text{ father}(X, Y) \wedge \text{person}(X)$   $p_2 : Y = Z \wedge \text{father}(X, Y), \text{father}(X, Z)$ , which asserts that each person has exactly one father. Clearly,  $BP_1 = \{\text{father}[1]\}$ , since the variable  $X$  is the first argument of  $\text{head}(p_1)$ , and  $JP_2 = \{\text{father}[1]\}$  due to the variable  $X$  in body ( $p_2$ ). Since  $BP_1 = JP_2$ ,

$\rho_1 : \exists Y \text{ father}(X, Y) \leftarrow \text{person}(X) \quad \rho_2 : Y = Z \leftarrow \text{father}(X, Y), \text{father}(X, Z)$ ,

and the existentially quantified variable  $Y$  occurs in  $\text{head}(p_i)$  only once, we get that  $P$  is non-conflicting. ■

Given a non-conflicting Datalog<sub>3</sub>,= program  $P$ , and an input database  $D$ , if  $D = P=$ , then during the construction of  $\text{chase}(D, P_3)$  it is not possible to violate an equality rule of  $P=$ . Let us explain briefly why this holds. Consider a pair  $(p_1, p_2) \in P_3 \times P=$  such that the head-predicate of  $p_1$  and the predicate of the atoms of body  $(p_2)$  is the same. In other words, the atom generated by applying  $p_1$  during the chase it is possible to violate  $p_2$ . Since  $P$  is non-conflicting

— either  $JP_2 \setminus UPI \neq \emptyset$ , and the application of  $p_1$  generates an atom  $a$  with a “fresh” null at some position of  $JP_2$ ; hence,  $a$  does not violate  $p_2$ ,

— or  $UPI = JP_2$ , and any newly generated atom  $a$  must have “fresh” distinct nulls (since each existentially quantified variable occurs in  $p_1$  just once) at all positions but those of  $JP_2$ . Therefore, if  $a$  coincides with some existing atom in the chase at the positions of  $JP_2$ , then  $a$  would not be added since  $p_1$  is not applicable; hence, again it is not possible to violate  $p_2$ .

From the above informal discussion we conclude that  $\text{chase}(D, P_3)$  is a universal model of  $P$  w.r.t.  $D$ . Consequently, if  $D = P=$ , then  $P(D)$  exists, and  $P(D) = q$  iff  $\text{chase}(D, P_3) = q$  iff  $P_3(D) = q$ ; hence,  $P$  is separable. The next result follows.

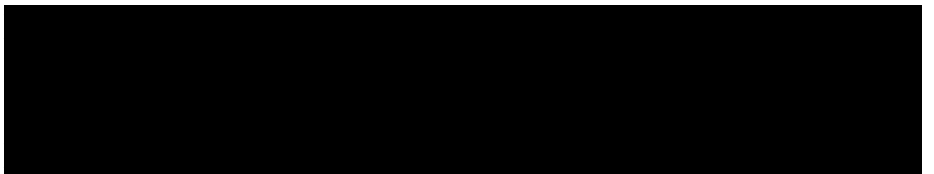
**Theorem 7.** BCQAns under non-conflicting guarded/linear Datalog<sub>3</sub>, programs has the same data/combined complexity as

BCQAns under guarded/linear Datalog<sub>3</sub> programs.

Adding L. We now show that the addition of L can be done effortlessly. Nonconflicting guarded/linear Datalog<sub>3,=,±</sub> is obtained by combining non-conflicting guarded/linear Datalog<sub>3,=</sub> with rules of the form  $\_L a_1, \dots, a_n$ . As shown in [9], the addition of this new feature does not increase the complexity of query answering. Let us explain why this holds. Consider a non-conflicting guarded/linear Datalog<sub>3,=,±</sub> program P, where  $P_{\exists}$  (resp.,  $P_{=}, P_{\pm}$ ) is the set of rules of P with existential quantifiers (resp., equality, the truth constant false) in the head, an input database D for SCH(P), and a BCQ q over SCH(P). For each rule  $p \in P_{\pm}$ , let  $q_p$  be the BCQ  $p \wedge \text{body}(p)$ , where  $p \in P_{\pm}$ . Also, let  $P_{\exists,=} = P_{\exists} \cup P_{=}$ . Clearly, if  $P_{\exists,=}(D) = q$ , for some  $p \in P_{\pm}$ , then there is no model of P w.r.t. D, and thus query answering is trivial since  $P(D)$  entails every BCQ over SCH(P). Consequently,  $P(D) = q$  iff  $P_{\exists,=}(D) = q$ , or  $P_{\exists,=}(D) = q_p$ , for some  $p \in P_{\pm}$ . Since P is non-conflicting, we immediately get that  $P(D) = q$  iff  $P_{\exists}(D) = q$ , or  $P_{\exists}(D) = q_p$ , for some  $p \in P_{\pm}$ . The next result follows.

Theorem 8. BCQAns under non-conflicting guarded/linear Datalog<sub>3</sub> programs has the same data/combined complexity as BCQAns under guarded/linear Datalog<sub>3</sub> programs.

DL-Lite is a well-known family of DLs under which query answering is highly tractable, i.e., in AC<sub>0</sub>, in data complexity [13,36]. Notice that DL-





Lite forms the OWL 2 QL profile of OWL 2 DL . It is easy to verify that a DL-Litex, where  $X \in \{F, R, A\}$ , ontology can be translated into a non-conflicting linear Datalog<sub>3,=,±</sub> program. Thus, non-conflicting linear Datalog<sub>3,=,±</sub> strictly captures the main languages of DL-Lite without losing the AC0 data complexity of BCQAns.

#### 5 Further Reading

Several classes of Datalog<sub>3</sub> programs have been investigated. Here are some examples. The class of weakly-guarded programs, an extension of guarded programs where the guard atom must cover only the body-variables that occur at affected positions, i.e., positions at which a null value can appear during the construction of the chase, has been investigated in [8]. The class of frontier-guarded programs, an extension of guarded programs where the guard atom must contain only the frontier, i.e., the set of variables that appear both in the body and the head, has been studied in [6]. As show in [9], guarded programs can be enriched with stratified negation, a simple non-monotonic form of negation often used in the context of plain Datalog. The class of weakly-acyclic programs has been proposed in the context of data exchange [23]. Notice that weak acyclicity guarantees the termination of the chase procedure. The more general class of super-weakly-acyclic programs has been studied in [32]. The class of sticky programs has been investigated in [11]. Stickiness is a decidability paradigm which

differs significantly from guardedness and weak acyclicity. The combination of stickiness with linearity and weak acyclicity is performed in [12]. The question whether weak acyclicity and guardedness can be combined in order to obtain more expressive decidable classes has been investigated in [29]. The work [22] presents the class of FDNC logic programs that allows for function symbols (F), disjunction (D), non-monotonic negation under the answer set semantics (N), and constraints (C), while standard reasoning tasks are decidable.