

# MIF14-BDD: PROJECT

## Implementing a Top-down Query Evaluation Engine for Datalog

May 14, 2021

### Abstract

This project will let you delve into the details of the implementing a query evaluation engine for Datalog. The project consists of two main parts.

The first part (about 12 points) consists of **implementing** an engine that, given as input a set of facts, a set of rules and a query, computes the answers, through a *top-down* evaluation heuristic.

The second part (about 8 points) consists of **analyzing** the performance of your engine.

To this end, you should:

- (about 4 points) *benchmark*, i.e., write a set of 10 example Datalog programs and queries to be used as sample inputs, together with their respective translations into DES syntax;
- (about 2 points) *test*, i.e., run your engine on the benchmark you designed and check that it outputs the same result as DES;
- (about 2 points) *evaluate*, i.e., compare the runtimes you obtain against those of DES and/or of other open-source Datalog engines and plot the results.

This project is to be performed by one or two students. The programming language is **Java**. A small report (at most 2 pages, without appendix) must be added to your project.

*The project is due on June 10, 2021 at noon (12pm).*

*You have to submit it on Tomuss by this deadline (no email will be accepted).*

## 1 The Datalog Language

A Datalog clause is of the form:

$$\forall \bar{x}. \forall \bar{y}. (\phi(\bar{x}, \bar{y}) \rightarrow H(\bar{y})).$$

The formula  $\phi(\bar{x}, \bar{y})$  is built from the conjunction of positive relational atoms having  $\bar{x}$  and  $\bar{y}$  as free variables. In the remainder, quantifiers and term arguments will be omitted.

A clause  $C$  is thus of the form:  $A_1, \dots, A_n \rightarrow H$ , where  $head(C) = H$ ,  $body(C) = \{A_1, \dots, A_n\}$ .

The following safety condition is imposed on all clauses of a program  $P$ :

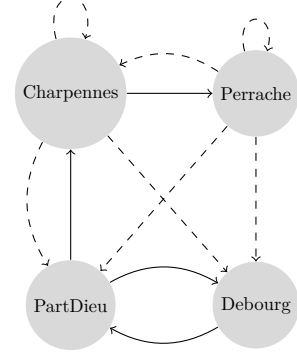
*all head variables should appear among those in the body.*

**Example 1.1** (Reachability Queries). Consider the following Datalog program. It consists of a set of **facts**, denoting that there is a *link* — marked by an arrow — between two stations, and of a set of **rules** indicating the *metro* names and the *reachability* between two stations (via multiple links) — marked by a dotted arrow. The program **query** computes all *metro* stations that are *reachable* from *Charpennes*.

```

link ( Charpennes , Perrache )
link ( PartDieu , Charpennes )
link ( Debourg , PartDieu )
link ( PartDieu , Debourg )
link ( X,Y )→metro ( X ) .
link ( X,Y )→metro ( Y ) .
link ( X,Y )→reachable ( X,Y ) .
link ( X,Z ) , reachable ( Z,Y )→reachable ( X,Y ) .
reachable ( Charpennes , Y ) → query ( Y ) .

```



The *bottom-up* evaluation of the program above, proceeds as follows to compute its fixpoint model:

$$\mathbf{T}_P^0(\emptyset) = \{link(Charpennes, Perrache), link(PartDieu, Charpennes), link(Debourg, PartDieu), link(PartDieu, Debourg)\}$$

$$\mathbf{T}_P^1(T_P^0) = T_P^0 \cup \{metro(Charpennes), metro(Perrache), metro(PartDieu), metro(Debourg), \textcolor{red}{reachable(Charpennes, Perrache)}, reachable(PartDieu, Charpennes), reachable(Debourg, PartDieu), reachable(PartDieu, Debourg)\}$$

$$\mathbf{T}_P^2(T_P^1) = T_P^1 \cup \{reachable(PartDieu, Perrache), reachable(Debourg, Charpennes), reachable(Debourg, Debourg), reachable(PartDieu, PartDieu), \textcolor{red}{query(Perrache)}\}$$

$$\mathbf{T}_P^3(T_P^2) = T_P^2 \cup \{reachable(Debourg, Perrache)\}$$

Note that this heuristic is *sub-optimal*, as the engine infers additional reachability facts that are *not relevant* for computing the needed query. In particular, as this example illustrates, it is sometimes more efficient to “push” the constants from the query definition (i.e., *Charpennes*) down into the rules that are relevant for the query, i.e., those deriving into the predicates being queried. This creates more (sub)queries from the atoms belonging to the bodies of these rules. In turn, subqueries are answered in a similar, “top-down” fashion. We will implement such *top-down* heuristic, namely the **Query-Subquery (QSQ)** algorithm.

## 2 Example Top-Down Query-Subquery Evaluation

The idea of **QSQ** evaluation is to *unify* the distinguished *query* atom, i.e., *reachable(Charpennes, Y)*, with the *relevant* rule *heads* and then propagate the acquired information inside the respective rule *bodies*. Each *body* atom becomes a *subquery* that is answered similarly, by unifying them with *relevant* rules, possibly resulting in more subqueries. When all subqueries pertaining to a rule *body* are answered, **QSQ** produces an answer set for the (sub)query pertaining to the rule *head*. The production and answering of queries/subqueries repeat until no more answer sets and no more subqueries are obtained.

We explain this procedure on the previous example. For the *reachable(Charpennes, Y)* query, the *relevant* rules from our previous program are  $R_1$  and  $R_2$  given below.

$$R_1 \quad link(X, Y) \rightarrow \textcolor{blue}{reachable}(X, Y) .$$

$$R_2 \quad link(X, Z), \textcolor{red}{reachable}(Z, Y) \rightarrow \textcolor{blue}{reachable}(X, Y) .$$

**Top-Down Information Passing** ↓ Passing the variable binding information to a rule *head* is called *top-down information passing*. In particular, we seek to determine what are the *constants* with which we can replace *head variables*. Let us understand how this works by focusing on  $R_1$ .

$$R_1 \quad link(X, Y) \rightarrow \textcolor{blue}{reachable}(X, Y) .$$

Unifying the *query* atom *reachable(Charpennes, Y)* with the head of the first rule below, *reachable(X, Y)*, first introduces the variable binding  $\{X \mapsto \textcolor{red}{Charpennes}\}$ . In this case, we say that the variable  $X$  is *bound*, as we have found a constant with which we can instantiate it.

**Sideways Information Passing**  $\leftarrow$  Passing the variable binding information from atom to atom in the same rule body is called *sideways information passing*. “Pushing” the variable binding information  $\{X \mapsto \text{Charpennes}\}$  into the body of  $R_1$ , we obtain:

$R_1$   $\text{link}(\text{Charpennes}, Y) \rightarrow \text{reachable}(\text{Charpennes}, Y)$ .

We can now take as a **subquery** the atom  $\text{link}(\text{Charpennes}, Y)$ . As  $\text{link}(\text{Charpennes}, \text{Perrache})$  is in our set of facts, we obtain the binding information  $\{Y \mapsto \text{Perrache}\}$ . Hence, we now know that  $\text{reachable}(\text{Charpennes}, \text{Perrache})$ .

Let us illustrate the  $\downarrow$  and  $\leftarrow$  information passing steps on rule  $R_2$  as well.

$R_2$   $\text{link}(X, Z), \text{reachable}(Z, Y) \rightarrow \text{reachable}(X, Y)$ .

From the  $\downarrow$  step, we obtain, as previously,  $\{X \mapsto \text{Charpennes}\}$ . Applying a  $\leftarrow$  step, we get:

$R_2$   $\text{link}(\text{Charpennes}, Z), \text{reachable}(Z, Y) \rightarrow \text{reachable}(\text{Charpennes}, Y)$ .

We now take  $\text{link}(\text{Charpennes}, Z)$  to be our **subquery**. As before, since  $\text{link}(\text{Charpennes}, \text{Perrache})$  is in our set of facts, we obtain the binding information  $\{Z \mapsto \text{Perrache}\}$ . We can now *backtrack* and continue to propagate the information sideways in another  $\leftarrow$  step. We thus have:

$R_2$   $\text{link}(\text{Charpennes}, \text{Perrache}), \text{reachable}(\text{Perrache}, Y) \rightarrow \text{reachable}(\text{Charpennes}, Y)$ .

We now take  $\text{reachable}(\text{Perrache}, Y)$  to be our **subquery**. We repeat the same procedure. The *relevant* rules are still  $R_1$  and  $R_2$ . Consider  $R_1$ :

$R_1$   $\text{link}(X, Y) \rightarrow \text{reachable}(X, Y)$ .

From the  $\downarrow$  step, we get  $\{X \mapsto \text{Perrache}\}$ . Applying the  $\leftarrow$  step:

$R_1$   $\text{link}(\text{Perrache}, Y) \rightarrow \text{reachable}(\text{Perrache}, Y)$ .

When taking  $\text{link}(\text{Perrache}, Y)$  to be our **subquery**, we notice there is no fact it can be unified against, so no new facts can be inferred using  $R_1$ . We *backtrack* to the  $\text{reachable}(\text{Perrache}, Y)$  **subquery** and are left to consider  $R_2$ :

$R_2$   $\text{link}(X, Z), \text{reachable}(Z, Y) \rightarrow \text{reachable}(X, Y)$ .

From the  $\downarrow$  step, we get  $\{X \mapsto \text{Perrache}\}$ . Applying the  $\leftarrow$  step:

$R_2$   $\text{link}(\text{Perrache}, Z), \text{reachable}(Z, Y) \rightarrow \text{reachable}(\text{Perrache}, Y)$ .

When taking  $\text{link}(\text{Perrache}, Z)$  to be our **subquery**, we notice there is no fact it can be unified against, so no new facts can be inferred using  $R_2$  either. Hence, no new reachability information can be obtained.

### 3 Subquery Components

To be able to process a **subquery** we need two ingredients:

- *argument information* – what are the *bound* variables that should be replaced by constants ?
- *binding information* – what is the *set of substitutions* that we can pass  $\downarrow$  or  $\rightarrow$  ?

**Argument Information** To explicitly mark which arguments are ready to receive binding information, for every **IDB predicate**  $P$  of our program, we compute its *adorned* version  $P^\gamma$ . The parameter  $\gamma$  is a list, whose length corresponds to the *arity* of  $P$  and whose elements are boolean flags, denoted  $b$  and  $f$ , marking whether the predicate argument in the respective position is *bound* or *free*.

*Adornments* propagate to rules, as every rule body atom  $B_i$  is rewritten to refer to an *adorned* predicate. To compute an adornment, we determine if the argument at each position is bound, by checking if:

- it is a *constant*
- it is a variable bound in the rule **head**, according to the adornment of the head atom
- it is a variable bound by some other atom  $B_j$ , where  $j < i$

Note that different atom orderings in the rule body can yield different adornments.

For example, the *adornment* of the rules we considered previously is:

$R_1$   $\text{link}(X, Y) \rightarrow \text{reachable}^{bf}(X, Y)$ .  
 $R_2$   $\text{link}(X, Z), \text{reachable}^{bf}(Z, Y) \rightarrow \text{reachable}^{bf}(X, Y)$ .  
 $R_3$   $\text{reachable}^{bf}(\text{Charpennes}, Y) \rightarrow \text{query}^f(Y)$ .

- In  $R_3$ , since the first body atom has a constant argument the corresponding adornment is  $\text{reachable}^{bf}$ , marking that the first position is *bound*, while the second position contains the *free* variable  $Y$ .
- In  $R_2$ , since  $X$  is bound in the head, it is also bound in the first argument of the first body atom; also, since any match for the first body atom bounds  $Z$ , we have that  $Z$  is also bound when evaluating the second atom, due to the left-to-right evaluation.
- In  $R_1$ , since  $X$  is bound in the head, it is also bound in the first argument of the first body atom.

**Binding Information/Auxiliary Relations** A *auxiliary relation*  $R$  consists of the set of all possible substitutions for the *bound arguments* of an adorned predicate  $P^\gamma$  and represents the *binding information* to be passed  $\downarrow$  or  $\leftarrow$ . For example,  $R = \{\{X \mapsto \text{Charpennes}, Y \mapsto \text{Perrache}\}\}$  contains a possible binding for  $X$  and  $Y$ .

Depending on the direction in which its bindings will applied,  $R$  can either be:

- an *input relation* –  $R \triangleq \text{input\_}P^\gamma(V)$  – representing information passed  $\downarrow$  into the head of adorned rules deriving  $P^\gamma$ , where  $V$  is a set of bound arguments in  $P^\gamma$  and the *arity* of  $\text{input\_}P^\gamma(V)$  is equal to the number of  $b$  arguments in  $\gamma$
- a *supplementary relation* –  $R \triangleq \text{sup}_j^i(V)$  – representing information passed  $\leftarrow$  into the  $j$ -th atom of the  $i$ -th rule, where  $V$  is a set of arguments that are:
  - bound by atoms before  $j$
  - later referenced
- a *output relation* –  $R \triangleq \text{output\_}P^\gamma(V)$  – representing the “completed” input, with added values for all unbound variables added, where the *arity* of  $\text{output\_}P^\gamma(V)$  is equal to the length of  $\gamma$

Note that when evaluating the body atoms from left to right, in a given rule  $i$ , with head predicate  $P$ , the *first* set of bindings  $\text{sup}_0^i$  come from  $\text{input\_}P^\gamma$  – namely, they are the set of substitutions for *the bound variables in the head* – and the *last* set of bindings  $\text{sup}_n^i$  – namely, the set of substitutions for *all variables in the head* – go to  $\text{output\_}P^\gamma$ .

For example, we mark the auxiliary relations in the previously considered rules as follows:

$R_1$   $[\text{sup}_0^1(X)] \text{link}(X, Y) [\text{sup}_1^1(X, Y)] \rightarrow \text{reachable}^{bf}(X, Y) [\text{output\_reachable}^{bf}(X, Y)]$ .  
 $R_2$   $[\text{sup}_0^2(X)] \text{link}(X, Z), [\text{sup}_1^2(X, Z)] \text{reachable}^{bf}(Z, Y) [\text{sup}_2^2(X, Y)] \rightarrow \text{reachable}^{bf}(X, Y) [\text{output\_reachable}^{bf}(X, Y)]$ .  
 $R_3$   $[\text{input\_reachable}^{bf}(X)] \text{reachable}^{bf}(\text{Charpennes}, Y) \rightarrow \text{query}^f(Y)$ .

## 4 QSQR Evaluation Steps

**Preprocessing and Initialization** Before the evaluation stage, we first add the needed *argument and binding information* to our program.

- We create the *adorned program* by recursively creating adorned rules for all adorned predicates.
- We initialize all *auxiliary relations* to empty sets.

The QSQR evaluation procedure consists of the following four steps:

**Step 1** The **query** atom – with predicate  $P$  – is *unified* with the *relevant* adorned rules – with head predicates  $P^\gamma$ , according to which of its arguments are *bound*. The corresponding substitutions are added to the **input relation**, namely **input\_** $P^\gamma$ .

**Step 2** For a *relevant* adorned rule  $i$ , we compute  $\text{supp}_0^i$ , by projecting out the necessary variables from the input relation.

**Step 3** We produce and evaluate **subqueries** by computing subsequent *auxiliary relations*. For a *relevant* adorned rule  $i$  – with body  $B \triangleq A_1, \dots, A_l$ , we do the following:

- if the predicate  $S$  of  $A_k$  is an EDB, given  $\text{supp}_{k-1}^i$ , we compute  $\text{supp}_k^i$ , by adding to each set of substitutions in  $\text{supp}_{k-1}^i$ , corresponding bindings that instantiate the atom's unbound variables to the appropriate constants of the facts with predicate  $S$ .
- if the predicate  $S^{\gamma_k}$  of  $A_k$  is an IDB, we add to **input\_** $S^{\gamma_k}$ , new bindings from  $\text{supp}_k^i$  combined with constants in  $A_k$
- if **input\_** $S^{\gamma_k}$  changed, *recursively* evaluate the **subqueries** corresponding to all rules with head predicate  $S^{\gamma_k}$
- compute  $\text{supp}_k^i$  from  $\text{supp}_{k-1}^i$  and **output\_** $S^{\gamma_k}$ .

**Step 4** Add the results of  $\text{supp}_l^i$  to **output\_** $P^\gamma$ .

We illustrate the evaluation, by revisiting our running example, in which we have our query relation, with the **input relation** explicitly marked:

$R_3$  [**input\_reachable**<sup>bf</sup>( $X$ )] **reachable**<sup>bf</sup>(*Charpennes*,  $Y$ )  $\rightarrow$  **query**<sup>f</sup>( $Y$ ).

**Step 1** We compute **input\_reachable**<sup>bf</sup>( $X$ ), by *unifying* the query body atom **reachable**<sup>bf</sup>(*Charpennes*,  $Y$ ), against the adorned head **reachable**<sup>bf</sup>( $X$ ,  $Y$ ) of the relevant rule  $R_1$ .

$R_1$  *link*( $X$ ,  $Y$ )  $\rightarrow$  **reachable**<sup>bf</sup>( $X$ ,  $Y$ ).

We obtain: **input\_reachable**<sup>bf</sup>( $X$ ) =  $\{\{X \mapsto \text{Charpennes}\}\}$ .

**Step 2** We then compute  $\text{sup}_0^1(X)$ , by projecting from **input\_reachable**<sup>bf</sup>( $X$ ) the bindings corresponding to the variables marked as bound in the adornment of the head predicate (**reachable**<sup>bf</sup>), i.e., for  $X$ . Hence,  $\text{sup}_0^1(X) = \{\{X \mapsto \text{Charpennes}\}\}$ .

$R_1$  [ $\text{sup}_0^1(X)$ ] *link*( $X$ ,  $Y$ )  $\rightarrow$  **reachable**<sup>bf</sup>( $X$ ,  $Y$ ).

**Step 3** Next, we compute  $\text{sup}_1^1(X, Y)$ . To this end, we notice that *link* is a EDB predicate and that we can directly compute the binding  $\{Y \mapsto \text{Perrache}\}$  from  $\text{sup}_0^1(X, Y)$ .

Hence,  $\text{sup}_1^1(X, Y) = \{\{X \mapsto \text{Charpennes}, Y \mapsto \text{Perrache}\}\}$ .

It follows that **output\_reachable**( $X$ ,  $Y$ ) =  $\{\{X \mapsto \text{Charpennes}, Y \mapsto \text{Perrache}\}\}$

and  $\text{sup}_0^2(X)$ , by projecting from **input\_reachable**<sup>bf</sup>( $X$ ) the bindings corresponding to the variables marked as bound in the adornment of each of the head predicates in  $R_1$  and  $R_2$ . In our setting, this bound variable corresponds to  $X$  in both cases and  $\text{sup}_0^1(X) = \text{sup}_0^2(X) = \{\{X \mapsto \text{Charpennes}\}\}$ .

$R_1$  [ $\text{sup}_0^1(X)$ ] *link*( $X$ ,  $Y$ )  $\rightarrow$  **reachable**<sup>bf</sup>( $X$ ,  $Y$ ).

$R_2$  [ $\text{sup}_0^2(X)$ ] *link*( $X$ ,  $Z$ ), **reachable**<sup>bf</sup>( $Z$ ,  $Y$ )  $\rightarrow$  **reachable**<sup>bf</sup>( $X$ ,  $Y$ ).

In  $R_1$ , we use the *binding information* provided by  $\text{sup}_0^1(X)$  to

We obtain that **input\_reachable**<sup>bf</sup>( $X$ ) =  $\{\{X \mapsto \text{Charpennes}\}\}$ ,  $\text{sup}_0^1 = \text{sup}_0^2 = \{X \mapsto \text{Charpennes}\}$ .

## 5 Input File Format

In the first part, we assume an input file containing a stratified Datalog programs in textual form. Their syntax is specified in BNF in Annex A.1. A parser (using JavaCC <sup>1</sup>) that checks the syntactic correctness of the program, according to the provided BNF grammar is available on **Google drive** for your convenience (you do not need to implement it).

Note that for handling constants in the atoms of a mapping, we have to declare a special **unique**, unary atom for that constant, as illustrated in the below example, i.e., to express:

---

```
EDB
link (Charpennes , Perrache)
link (PartDieu , Charpennes)
link (Debourg , PartDieu)
link (PartDieu , Debourg)

IDB
metro ($x)
reachable ($x , $y)
query ($y)

MAPPING
reachable (Charpennes , $y) -> query ($y).
```

---

...we will write:

---

```
EDB
link (Charpennes , Perrache)
link (PartDieu , Charpennes)
link (Debourg , PartDieu)
link (PartDieu , Debourg)
cst (Charpennes)

IDB
metro ($x)
reachable ($x , $y)
query ($y)

MAPPING
cst ($x) , reachable ($x , $y) -> query ($y).
```

---

---

<sup>1</sup><https://javacc.java.net/>

**Example 5.1.** The full Example 1.1 is written according to the BNF grammar as follows:

EDB

```
link(Charpennes,Perrache)
link(PartDieu,Charpennes)
link(Debourg,PartDieu)
link(PartDieu,Debourg)
cst(Charpennes)
```

IDB

```
metro($x)
reachable($x,$y)
query($y)
```

MAPPING

```
link($x,$y) -> metro($x).
link($x,$y) -> metro($y).
link($x,$y) -> reachable($x,$y).
link($x,$z), reachable($z,$y) -> reachable($x,$y).
cst($x), reachable($x,$y) -> query($y).
```

# A Formal Grammar

## A.1 BNF for stratified Datalog clauses

The BNF for input files is given below. Note that, for sake of simplicity, usual skippable characters (spaces, tabulations, carriage returns) have been omitted. Comments (also omitted and to be skipped) start with two dashes (--) and finish at the end of the line.

```
START      ::= "EDB" SCHEMA "IDB" SCHEMA "MAPPING" TGDS

SCHEMA     ::= RELATION SCHEMA
              | RELATION

RELATION   ::= NAME "(" ATTS ")"

ATTS       ::= NAME "," ATTS
              | NAME

TGDS       ::= TGD
              | TGDS TGD

TGD        ::= QUERY "->" ATOM "."

LITERAL    ::= "NEG" ATOM
              | ATOM

QUERY      ::= LITERAL "," QUERY
              | LITERAL

ATOM       ::= NAME "(" ARGS ")"

ARGS       ::= VALUE "," ARGS
              | VALUE

VALUE      ::= VARIABLE
              | CONSTANT

VARIABLE   ::= "$" NAME

NAME       ::= LETTER
              | LETTER NAME2

NAME2      ::= LETTER_OR_DIGIT
              | LETTER_OR_DIGIT NAME2

CONSTANT   ::= DIGITS

DIGITS     ::= DIGIT DIGITS
              | DIGIT

LETTER_OR_DIGIT ::= LETTER
                | DIGIT

LETTER     ::= "a" | "b" | "c" | "d" | "e" | "f" | "g" | "h" | "i" | "j"
              | "k" | "l" | "m" | "n" | "o" | "p" | "q" | "r" | "s" | "t"
              | "u" | "v" | "w" | "x" | "y" | "z"
              | "A" | "B" | "C" | "D" | "E" | "F" | "G" | "H" | "I" | "J"
              | "K" | "L" | "M" | "N" | "O" | "P" | "Q" | "R" | "S" | "T"
              | "U" | "V" | "W" | "X" | "Y" | "Z"
              | "_"

DIGIT      ::= "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8" | "9"
```