

Bottom-Up Evaluation of Datalog: Preliminary Report

Stefan Brass Heike Stephan

Institut für Informatik, Martin-Luther-Universität Halle-Wittenberg, Germany

brass@informatik.uni-halle.de

stephan@informatik.uni-halle.de

Bottom-up evaluation of Datalog has been studied for a long time, and is standard material in textbooks. However, if one actually wants to develop a deductive database system, it turns out that there are many implementation options. For instance, the sequence in which rule instances are applied is not given. In this paper, we study a method that immediately uses a derived tuple to derive more tuples (called the Push method). In this way, storage space for intermediate results can be reduced. The main contribution of our method is the way in which we minimize the copying of values at runtime, and do much work already at compile-time.

1 Introduction

The efficient evaluation of queries expressed as logic programs remains an everlasting problem. Of course, big achievements have been made, but at the same time problem size and complexity grows. Any further progress can increase the practical applicability of logic-based, declarative programming.

Our long-term goal is to develop a new deductive database system. This has many aspects, for instance, language design. However, in the current paper, we exclude all special language features, including negation, and focus on efficient query evaluation for basic Datalog.

The magic set method is the standard textbook method for making bottom-up evaluation goal-directed. Many optimizations have been proposed, including our own SLDMagic method [1] and a method based on Earley deduction [3]. We assume in the current paper that such a rewriting of the program has been done, so we can concentrate on pure bottom-up evaluation.

As we understand it, bottom-up evaluation is an implementation of the T_P -operator that computes the minimal model of the program. However, an implementation is free in the order in which it applies the rule instances, while the T_P -operator first derives all facts that are derivable with a given set of known facts, before the derived facts are used (in the next iteration). Furthermore, facts do not have to be stored until the end of query evaluation, but can be deleted as soon as all possible derivations using them have been done, except for the facts that form the answer to the query. Therefore, the sequence of rule instance application becomes important. If one computes predicate by predicate as the standard textbook method, one of course needs to store the entire extension of the predicates. However, if one uses derived tuples immediately, it might be possible to store only one tuple of the predicate during the evaluation. Of course, for duplicate elimination and termination, it might still be necessary to store extensions of a few selected predicates. It is also not given that tuples (facts) must be represented explicitly as records or objects in the program. It suffices if one knows where the values of single columns (predicate arguments) can be found. In this way, a lot of copying can be saved because tuples for the rule head are typically constructed from values bound in the rule body. Of course, one must ensure that the values are not changed before all usages are finished.

Our plan is to translate Datalog to C++, and to generate executable code from the resulting program. This permits to use existing compilers for low-level optimizations and gives an interface for defining built-in predicates. In [2], we already discussed implementation alternatives for bottom-up evaluation and did performance comparisons for a few example programs. Now we will improve the “push method” from that paper by changing the set of variables used to represent intermediate facts. This is essential for reducing the amount of copying. It also enables us to do more precomputation at “compile time”.

The idea of immediately using derived facts to derive more facts is not new. For instance, variants of semi-naive evaluation have been studied which work in this way [10, 12]. It also seems to be related to the propagation of updates to materialized views. However, the representation of tuples at runtime and the code structure is different from [10] (and this is essential for the reduction of copying values). The paper [12] translates from a temporal Datalog extension to Prolog, which makes any efficiency comparison dependent on implementation details of the used Prolog compiler. We also believe that the rule application graph introduced in our paper is a useful concept. Further literature about the implementation of deductive database systems is, for instance, [8, 4, 9, 11, 7, 13]. A current commercial deductive DB system is LogicBlox [5]. A benchmark collection is OpenRuleBench [6].

2 Basic Definitions

In this paper, we consider basic Datalog, i.e. pure Prolog without negation and without function symbols (i.e. terms can only be variables or constants). We also assume without loss of generality that all rules have at most two body literals. The output of our rewriting methods [1, 3] has this property. (But in any case, it is no restriction since one can introduce intermediate predicates.) Finally, we require range-restriction (allowedness), i.e. all variables in the head of the rule must also appear in a body literal. For technical purposes, we assume that each rule has a unique rule number.

As usual in deductive databases, we assume that EDB and IDB predicates are distinguished (“extensional” and “intensional database”). EDB predicates are defined by facts only, e.g. stored in a relational database or specially formatted files. Also program input is represented in this way. IDB predicates are defined by rules. There is a special IDB-predicate answer that only appears in the head of one or more rules. The task is to compute the extension of this predicate in the minimal model of the program, i.e. all derivable answer-facts.

We assume that the logic program for the IDB predicates as well as the query (i.e. the answer-rules) are given at “compile time”, whereas the database for the EDB predicates is only known at “runtime”. Since the same program can be executed several times with different database states, any optimization or precomputation we can do at compile time will pay off in most cases. It might even be advantageous in a single execution because the database is large.

Since we want to generate C++ code, we assume that a data type known for every argument of an EDB predicate. The method does not need type information for IDB predicates (this is implicitly computed). Data structures for storing relations for EDB predicates can be characterized with binding patterns: A binding pattern for a predicate p with n arguments is a string of length n over the alphabet $\{b, f\}$. The letter b (“bound”) means that a value for the corresponding argument is known when the predicate relation is accessed (input), f (“free”) means that a value needs to be looked up (output).

As mentioned above, our rewriting methods [1, 3] produce rules that have at most two body literals. Furthermore the case of two IDB-literals is rare — it is only used in special cases for translating complex recursions. Most rules have one body literal with IDB-predicate and one with EDB-predicate. Of course, there are also rules with only one body literal (EDB or IDB).

3 Accessing Database Relations

The approach we want to follow is to translate Datalog into C++, which can then be compiled to machine code. Of course, we need an interface to access relations for the EDB predicates. These relations can be stored in a standard relational database, but it is also possible to program this part oneself (at the moment, we do not consider concurrent updates and multi-user access).

We assume that it is possible to open a cursor (scan, iterator) over the relation, which permits to loop over all tuples. We assume that for every EDB predicate p there is a class p_cursor with the following methods:

- `void open()`: Open a scan over the relation, i.e. place the cursor before the first tuple.
- `bool fetch()`: Move the cursor to the next tuple. This function must also be called to access the first tuple. It returns `true` if there is a first/next tuple, or `false` if the cursor is at the end of the relation.
- `T col_i()`: Get the value of the i -th column (attribute) in the current tuple. Here T is the type of the i -th column.
- `close()`: Close the cursor.

For recursive rules, we will also need

- `push()`: Save the state of the cursor on a global stack.
- `pop()`: Restore the state of the cursor.

A relation may have special access structures (e.g. it might be stored in a B-tree, hash table or array). Then not only a full scan (corresponding to binding pattern `ff...f`) is possible, but also scans only over tuples with given values for certain arguments. We assume that in such cases there are additional cursor classes called p_cursor_β , with a binding pattern β . These classes have the same methods as the other cursor classes, only the `open`-method has parameters for the bound arguments. E.g. if p is a predicate of arity 3 that permits particularly fast access to tuples with a given value of the first argument, and if this argument has type `int`, the class `p_cursor_bff` would have the method `open(int x)`.

4 Duplicate Elimination and Termination

The main contribution of this paper is the way in which copying and materialization of tuples is avoided. Our method basically pushes newly derived facts to body literals where they can be used to derive further facts.

However, in the presence of recursion, we must be able to notice whether a derived tuple is new or not. Therefore, in each recursive cycle, at least one predicate must be materialized (“tabled”) to ensure termination. A simple solution is to create hash tables for the predicates in question.

This solution means that we materialize the extensions of some IDB predicates (hopefully, only a few) and copy all data values for the tuples of these predicates. In some cases, information about order or acyclicity might help to avoid this. Information about keys and data distribution could be used to make sensible optimization decisions. Furthermore, if tuples are produced in a sort order, the duplicate check can be done very efficiently and without storing the predicate extension. All this is subject of our future work.

It is also interesting that the data values in a derived tuple are stored at different times in program variables. For instance, we might know that when $p(X, Y)$ is generated, X only seldom changes, and Y

```

<Declaration Section>;
<Initialization Section>; // Initializes backtrack_stack
while(!backtrack_stack.is_empty()) {
    switch(backtrack_stack.pop()) {
        case L1:
            l1:
                <Code Piece 1>;
                // break or goto at end of Code Piece
        case L2:
            l2:
                <Code Piece 2>;
            ...
    }
}

```

Figure 1: Overall structure of the generated code

changes much more often. Then a nested relation might be best for tabling the predicate for the purpose of duplicate detection.

Of course, breaking each recursive cycle with a duplicate detection is only the minimum we have to do to ensure termination. Also non-recursive rules can generate duplicates, and in some cases it might be more efficient to detect these duplicates early in order to avoid duplicate computations (since the price for duplicate detection is quite high, in other cases it might be more efficient to simply do the duplicate work).

5 Code Generation: Overall Structure

The result of the translation looks basically as shown in Figure 1. So there are many small code pieces, each with a label that is suitable for a `goto`. Furthermore, when there are several things to do, e.g. a generated fact can be used in more than one rule, a backtrack point is set up for the second rule, and then a `goto` is done for the first. When an execution path reaches an end, the `switch` is left with `break`, and one of the delayed tasks is taken from the stack. Therefore, each code piece also has a unique number, which can be stored on the backtrack stack, and used in the `switch` to reach the code piece.

Optimizations are possible, e.g. one can order the code pieces such that some jumps can be eliminated, because the target is immediately following. Some backtrack points can be avoided by finding a suitable code sequence.

5.1 Declaration Section

Data not known at compile time always originates from the database. In order to minimize copying, we (usually) introduce a C++ variable only for Datalog variables which

- occur in an EDB body literal,
- but do not occur in an IDB body literal of that rule (because then the value comes from another rule, where a variable has been created, if the value is not known at compile time),

- and occur in the head of that rule (because otherwise the value does not really have to be processed in the program).

For instance, consider the following rule:

$$p(X, Y, a) \leftarrow q(Y) \wedge r(X, Y, Z, Z).$$

If q is an IDB predicate and r an EDB predicate, we create a C++ variable only for X . A variable or constant for Y exists already when the rule is activated.

In seldom cases of recursive rule applications (see Section 5.4 below) we create C++ variables for all variables of the rule.

If the above condition shows that we must create a C++ variable for variable X in rule ρ , we generate the following code line in the declaration section:

```
T v $\rho$ _X;
```

We use the prefix with the rule number so that there can be no name conflicts between variables of different rules. T is the C++ data type for the database column in which X occurs.

5.2 Symbolic Facts

A symbolic fact consists of an IDB predicate p and a tuple (t_1, \dots, t_n) of C++ variables (i.e. their identifiers) and constants, where n is the arity of p . So a symbolic fact represents what is known at compile time about a fact that will be derived at runtime. For some arguments, we might know the exact value (a constant), for other arguments, we know the C++ variable which will contain the value.

An initial set of symbolic facts is derived by rules without IDB body literals. Then our task is to pass each derived symbolic fact to matching IDB body literals and to derive a symbolic fact for the rule head. For each such rule application, a code piece is generated which does the remaining computation at runtime. The computation of symbolic facts is similar to the standard fixpoint iteration to compute the minimal model (but it is done at “compile time”, when the data for the EDB predicates are not yet known).

“Matching” between a symbolic fact and a body literal means that they are unifiable. In general a full unification must be done (at compile time). Consider e.g. the body literal $p(X, X, a)$ and the symbolic fact $p(b, v1_Y, v1_Y)$. The rule cannot be applied to the symbolic fact, so no code is generated for this case.

5.3 Rule Application Graph

As explained above, we assume that all rules have at most two body literals. A “Symbolic Rule Application” is represented by

- a rule from the logic program with one IDB body literal, together with a symbolic fact matching this body literal, or
- a rule without IDB body literals, or
- a rule with two IDB body literals, with one of the two selected, together with a symbolic fact matching this body literal. (In the rare case of two IDB body literals, we use temporary tables for facts matching each body literal. The symbolic fact in this rule application describes the situation that we just computed a new fact for one of the IDB body literals. For the other body literal we use the table with previously computed facts.)

The result of a symbolic rule application is a symbolic fact. Let $p(t_1, \dots, t_n)$ be the head of the rule, and ρ be its rule number. If the rule has an IDB body literal, let θ be a most general unifier with the input symbolic fact. We require that variable-to-variable bindings are done such that logic variables are replaced by C++ variables. Then the derived symbolic fact is $p(u_1, \dots, u_n)$, where u_i is

- t_i if t_i is a constant.
- $t_i\theta$ if t_i is a variable which appears in the IDB body literal (if there is one).
- $\nu\rho_X$ if t_i is a variable X which does not appear in the IDB body literal.

Now we can do a standard fixpoint computation to compute all symbolic facts which are derivable from the program. This process will come to an end, because the number of symbolic facts is bounded: There is only a finite number of C++ variables (at most the number of variables in the given logic program, where variables with the same name in different rules count as distinct). Furthermore, only a finite number of constants occurs in the given logic program (constants which appear only in the database are not known at “compile time” and not used for computing symbolic facts).

The structure of the computation can be shown in a “rule application graph”. It has two types of nodes, namely symbolic facts (“fact nodes”), and symbolic rule applications (“rule nodes”). There is an edge from every symbolic fact to every symbolic rule application which uses the symbolic fact. Furthermore, there is an edge from every symbolic rule application to the symbolic fact it generates.

Of course, it is possible to show only the rule in nodes for symbolic rule applications (since the symbolic fact is identified by the incoming edge, except in the case of two IDB body literals). However, then there can be several nodes marked with the same rule: It is possible that a single rule is compiled several times for different symbolic facts matching its IDB body literal.

Note also that not every application of a recursive rule to a symbolic fact is actually recursive: Only if the same symbolic fact can be generated by applying this rule (maybe indirectly via other rules), we have to be prepared for recursive invocations of the code piece for the symbolic rule application. This can be seen from cycles in the graph.

Finally, nodes in the graph from which there is no path to an answer-node can be eliminated: They do not contribute to the computation of the answer. If the program is the result of a program transformation like magic sets, this path will not be followed at runtime, but it is better not to generate code for it. An example of such a program is

```

answer(X) ← q(X,a).
q(X,Y)    ← p(Y,X).
p(a,X)    ← r(X).
p(b,X)    ← s(X).

```

The rule application graph is shown in Figure 2. The right path is useless. In the code generation below, we assume that such useless computation paths have been removed, i.e. from every node a fact node with predicate `answer` is reachable. This in particular means that every fact node with a predicate different from “`answer`” has an outgoing edge.

5.4 Variable Conflicts

In rare cases of recursive rule applications, it is possible that a rule is applied to a symbolic fact which contains already a variable generated for that rule. An example is

```

p(X,Y) ← r(X,Y).
p(Y,Z) ← p(X,Y) ∧ r(Y,Z).

```

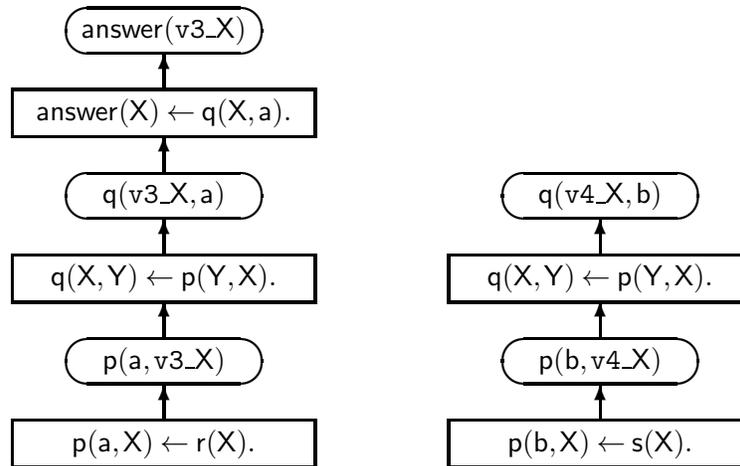


Figure 2: Rule Application Graph with Useless Part (to be eliminated).

The first rule generates the symbolic fact $p(v1_X, v1_Y)$. When we insert this into the second rule, we get $p(v1_Y, v2_Z)$. Now we have to insert this again into the second rule: $v2_Z$ contains the input value for Y , but must also be set with a new data value from r . In this case, some copying seems unavoidable. While there are optimizations possible, the simplest solution is to create a C++ variable for each logical variable of the rule, and to copy first the values from the input fact to the right variable (which might need temporary variables, e.g. for swapping the values of two variables). For recursive rule applications, the previous variable values are also stored on a stack (see Section 5.6 below).

5.5 Labels for Code Pieces

We need a `goto` label and/or a case selector value (a unique number) for each code piece implementing a symbolic rule application. We write

$$\text{l_start}(p(t_1, \dots, t_n), \rho, p(u_1, \dots, u_n))$$

for the `goto`-label of the code piece for application of rule ρ with body literal $p(u_1, \dots, u_n)$ to the symbolic fact $p(t_1, \dots, t_n)$. Of course, instead of listing the body literal $p(u_1, \dots, u_n)$ explicitly, one could also use its position number in the rule ρ . In any case, the implementation will replace this by `l_start_n` with some unique number n . The symbolic constant for the case-value is written as `L_START(...)` (and also made a legal C++ identifier by using the same unique number). Sometimes there are continuations or other code pieces, therefore the label is marked as “start”.

5.6 Protection of Variable Values

Of course, when a code piece corresponding to a symbolic rule application is executed, the C++ variables in the symbolic fact $p(t_1, \dots, t_n)$ must still have the same value as when this task was generated. It is possible that the ID/label of the code piece was pushed on the backtrack stack and it is executed only later.

However, for every C++ variable, a new value is assigned only in code pieces for the single rule for which the variable was introduced (to hold a data value for an EDB literal in that rule).

Furthermore, it is important that the backtrack points are kept on a stack. So we will return to that rule only after all backtrack points which use the value (and are thus generated later) have been processed—unless the rule is recursive. In this case, the variable value must be saved (on another stack suitable for the data type), and we put the ID of a code piece on the backtrack stack which restores the variable value. This is done whenever we enter a recursive rule, and only for variables set in this rule (the derived fact might contain also variables passed from elsewhere and not changed in the rule).

If the backtrack stack shrinks below this point, all usages of the new variable value are done, and the old value is restored, so that older backtrack points find the value which was current when the backtrack point was created.

6 Code Pieces

In this section, we define a number of code pieces which are translations of different types of rules. Each code piece corresponds to a symbolic rule application. For simplicity, we do not consider variable conflicts (Section 5.4) here.

6.1 IDB-Facts

Suppose the program contains an IDB-fact $p(c_1, \dots, c_n)$. For each body literal $p(t_1, \dots, t_n)$ of a rule ρ that unifies with the fact $p(c_1, \dots, c_n)$, the case selector value

$$\text{L_START}(p(c_1, \dots, c_n), \rho, p(t_1, \dots, t_n))$$

is pushed on the backtrack stack during initialization.

6.2 One EDB-Body Literal

Consider the rule $p(t_1, \dots, t_n) \leftarrow r(u_1, \dots, u_m)$ where r is an EDB predicate. Let ρ be the rule number. Let $p(\bar{t}_1, \dots, \bar{t}_n)$ be the symbolic fact generated by the rule ($\bar{t}_i := t_i$ if t_i is a constant, and $\bar{t}_i := \forall \rho _X$ if t_i is the variable X).

Among all possible cursors $\text{cursor_r_}\beta$ for r choose one such that for all bound argument positions i (i.e. $\beta_i = b$), u_i is a constant. This is always possible because every relation supports a full table scan, i.e. an access path with all argument positions “free”. But obviously, if there are constants among the u_i , and there are available indexes, it is best to choose one with the smallest estimated result size. In the declaration section, generate

$$\text{cursor_r_}\beta \text{ } c\rho;$$

Define symbolic constants $\text{L_INIT_}\rho$ and $\text{L_CONT_}\rho$ as unique numbers for cases in the switch. Generate the following code in the initialization section:

$$\text{backtrack_stack.push}(\text{L_INIT_}\rho);$$

All following code is generated in the switch:

1. Generate

$$\text{case L_INIT_}\rho:$$

2. Let i_1, \dots, i_k be the bound argument positions in β . Generate:

$$c\rho.\text{open}(u_{i_1}, \dots, u_{i_k});$$

(Note that although another case follows, execution simply continues.)

3. Generate:

```
case L_CONT_ρ:
```

The following loop (item 4) is left with `goto` when the first fact is generated. But before the jump, this case label is pushed on the backtrack stack, so that the loop is continued later.

4. Generate

```
while(cρ.fetch()) {
```

5. Let u_{i_1}, \dots, u_{i_k} be the constants among the u_1, \dots, u_m which correspond to free argument positions in β . If $k \geq 1$, generate

```
if(cρ.col_i1() != u_{i_1} || ... || cρ.col_i_k() != u_{i_k})
  continue;
```

I.e. if the current tuple of the EDB-predicate does not have the required values for the constant arguments, we immediately start the next iteration of the `while`-loop (i.e. fetch the next tuple).

6. For every variable Y , which appears more than once among the u_1, \dots, u_m : Let u_{i_1}, \dots, u_{i_k} be all equal to Y (note that $k \geq 2$). Generate a test that the same value appears in these columns:

```
if(cρ.col_i1() != cρ.col_i2 || ... || cρ.col_i_{k-1}() != cρ.col_i_k())
  continue;
```

7. For every variable X_i in the head let u_j be any occurrence of this variable among the u_1, \dots, u_m . Because of the range restriction (allowedness) condition on the rules, X_i must occur in the body. Generate for each X_i :

```
vρ_X_i = cρ.col_j();
```

8. In case the predicate p was selected for a duplicate check, the following must be done here: The result tuple $p(\bar{t}_1, \dots, \bar{t}_n)$ with the current values of the C++ variables is entered into a hash table or other data structure. If the tuple was already present, one simply does “`continue;`” to skip it.

9. Generate:

```
backtrack_stack.push(L_CONT_ρ);
```

This ensures that the `while`-loop above will be continued later. Since then the values of the variables introduced in the rule will change, this label must be on the stack below every task using the generated tuple.

10. Let ρ_1, \dots, ρ_k be all rules with an IDB body literal B_i , $i := 1, \dots, k$, which matches the generated symbolic fact $p(\bar{t}_1, \dots, \bar{t}_n)$. For $i := 2, \dots, k$, generate

```
backtrack_stack.push(L_START(p(\bar{t}_1, \dots, \bar{t}_n), ρ_i, B_i));
```

Finally, generate

```
goto l_start(p(\bar{t}_1, \dots, \bar{t}_n), ρ_1, B_1);
```

11. Generate

```
} // End of while-loop
break;
```

The `break;` is important if the `while`-loop ends because no (further) matching fact is found in the relation r . Otherwise, the loop is left with `goto` when the first/next matching fact is found.

6.3 Two EDB-Body Literals

In the output of SLDMagic, this case does not occur. However, it is easy to extend the above program code. One uses two cursors, one for each body literal, and two nested `while`-loops. For simplicity, we implement all joins as “nested loop join” (or “index join” if the data structure for the relation supports the corresponding binding pattern). Later, sort orders might be used, so that also a “merge join” can be generated.

6.4 One IDB-Body Literal

Consider the rule

$$p(t_1, \dots, t_n) \leftarrow q(u_1, \dots, u_m),$$

where q is an IDB-predicate. Let ρ be the number of this rule. Due to partial evaluation done at compile time, several specializations of the same rule might be generated. There is one code piece per symbolic fact $q(\bar{u}_1, \dots, \bar{u}_m)$ which matches the body literal. Let θ be a most general unifier, where variable-to-variable bindings are done such that logic variables are replaced by C++ variables i.e. u_i is replaced by \bar{u}_i , if both are variables. The generated symbolic fact is $p(t_1\theta, \dots, t_n\theta)$. Note that because of the range restriction requirement, every variable among the t_i also appears as an u_j , and then it is unified with a constant or a C++ variable. Thus, no new C++ variables are introduced in this case.

1. Generate

```
case L_START(q( $\bar{u}_1, \dots, \bar{u}_m$ ),  $\rho$ , q( $u_1, \dots, u_m$ )):
  l_start(q( $\bar{u}_1, \dots, \bar{u}_m$ ),  $\rho$ , q( $u_1, \dots, u_m$ )):
```

If this rule application is activated via backtracking, the `case` label is used. If it is activated as the first usage of a generated fact, a jump to the `goto` label is done (as a slight optimization of pushing something on the backtrack stack and immediately popping it again).

2. Now the part of the unification which can only be done at runtime must be generated. Let V_1, \dots, V_k be all C++ variables which θ replaces by constants or a different variable (i.e. $V_i\theta \neq V_i$). If $k > 0$, generate:

```
if ( $V_1 \neq V_1\theta$  || ... ||  $V_k \neq V_k\theta$ )
  break;
```

So we simply stop executing this code piece if the current fact for q does not unify with the body literal. Then another task will be taken from the backtrack stack in the main loop.

3. Next, if the predicate p was selected for a duplicate check, the code to enter the result tuple $p(t_1\theta, \dots, t_n\theta)$ with the current values of the C++ variables into a hash table is generated here. If the tuple was already present (so we just computed a duplicate), one simply does “`break;`” to end the code piece (as under 2 above).
4. Let ρ_1, \dots, ρ_k be all rules with an IDB body literal B_i , $i := 1, \dots, k$, which matches the generated symbolic fact $p(t_1\theta, \dots, t_n\theta)$. For $i := 2, \dots, k$, generate

```
backtrack_stack.push(L_START(p( $t_1\theta, \dots, t_n\theta$ ),  $\rho_i$ ,  $B_i$ ));
```

Finally, generate

```
goto l_start(p( $t_1\theta, \dots, t_n\theta$ ),  $\rho_1$ ,  $B_1$ );
```

6.5 One IDB- and one EDB-Body Literal

Consider the rule

$$p(t_1, \dots, t_n) \leftarrow q(u_1, \dots, u_m) \wedge r(v_1, \dots, v_l),$$

where q is an IDB-predicate and r is an EDB-predicate. Let ρ be the number of this rule. Again, there is one code piece per symbolic fact $q(\bar{u}_1, \dots, \bar{u}_m)$ which matches the IDB body literal. Let θ be a most general unifier and $p(\bar{t}_1, \dots, \bar{t}_n)$ be the generated symbolic fact as defined in Section 5.3. As in Section 6.2, select a binding pattern β for accessing the EDB-relation r . A value for v_i is known (i.e. β_i can be “bound”) if v_i is a constant or a variable which also appears in $q(u_1, \dots, u_m)$ (when execution reaches this code piece, a concrete fact is given for the IDB body literal). In the declaration section, generate

```
cursor_r_β cρ;
```

All following code is generated in the switch:

1. Generate

```
case L_START(q(ū1, ..., ūm), ρ, q(u1, ..., um)):
  l_start(q(ū1, ..., ūm), ρ, q(u1, ..., um)):
```

2. Now the part of the unification of the given fact with the IDB body literal, which can only be done at runtime, must be generated. Let V_1, \dots, V_k be all C++ variables with $V_i\theta \neq V_i$. If $k > 0$, generate:

```
if (V1 != V1θ || ... || Vk != Vkθ)
  break;
```

This ends the execution of this code piece if the rule is not applicable.

3. Now we must use a cursor to access the tuples for the EDB body literal $r(v_1, \dots, v_l)$. In case this rule is recursive, it might be possible that the state of the cursor and the values of the variables $v\rho_X$ set in this rule are still needed by backtrack points on the stack (unless we know that there are no such backtrack points, e.g. because the recursive rule application is the last use of the fact). Therefore, we generate

```
cρ.push();
```

And for each variable $v\rho_X$ we generate

```
value_stack.push(vρ_X);
```

(There are probably several value stacks for different data types.)

Finally we generate a backtrack point:

```
backtrack_stack.push(L_RESTORE_ρ);
```

The code for this case in the switch simply restores the variable values and the cursor state by popping them (in the inverse order). In this way, all earlier backtrack points (below the one just generated) find the old cursor state and variable values.

4. Now we open the cursor over the EDB-relation r with the selected binding pattern β . Let i_1, \dots, i_k be the bound argument positions in β . Generate:

```
cρ.open(ṽi1, ..., ṽik);
```

where \bar{v}_{i_j} is

- v_{i_j} if this is a constant,
- $v_{i_j}\theta$ if v_{i_j} is a variable which appears in the IDB body literal $q(\dots)$.

5. Generate:

```
case L_CONT_ρ:
```

When we are finished with using the computed fact, backtracking returns here to continue the following loop.

6. Generate

```
while(cρ.fetch()) {
```

7. Now we do that part of the section for $r(v_1, \dots, v_l)$ that was not supported by the data structure for r with binding pattern β (e.g., hash table). Let i_1, \dots, i_k be the free argument positions in β such that v_{i_j} is a constant or a variable which appears in the IDB body literal $q(\dots)$. Let \bar{v}_{i_j} be defined as in 4 above. If $k \geq 1$, generate

```
if(cρ.col_i1() !=  $\bar{v}_{i_1}$  || ... || cρ.col_i_k() !=  $\bar{v}_{i_k}$ )
  continue;
```

I.e. if the current tuple in the EDB relation does not have the required values, we continue with the next iteration of the while-loop under 6.

8. For every variable Y , which appears more than once among the v_1, \dots, v_l , but not in the IDB literal $q(\dots)$: Let u_{i_1}, \dots, u_{i_k} be all equal to Y (note that $k \geq 2$). Generate:

```
if(cρ.col_i1() != cρ.col_i2 || ... || cρ.col_i_{k-1}() != cρ.col_i_k())
  continue;
```

9. For every variable X_i in the head, which does not appear in the IDB body literal $q(\dots)$, let v_j be any occurrence of this variable among the v_1, \dots, v_l . Because of the range restriction (allowedness) condition on the rules, it must occur there. Generate for each X_i :

```
vρ_Xi = cρ.col_j();
```

10. In case predicate p was selected for a duplicate check, we again enter the result tuple $p(\bar{t}_1, \dots, \bar{t}_n)$ with the current values of the C++ variables into a hash table. If the tuple was already present, one simply does “continue;” to compute the next tuple.

11. Generate:

```
backtrack_stack.push(L_CONT_ρ);
```

12. Let ρ_1, \dots, ρ_k be all rules with an IDB body literal B_i , $i := 1, \dots, k$ that matches the generated symbolic fact $p(\bar{t}_1, \dots, \bar{t}_n)$. For $i := 2, \dots, k$, generate

```
backtrack_stack.push(L_START(p( $\bar{t}_1, \dots, \bar{t}_n$ ),  $\rho_i$ ,  $B_i$ ));
```

Then generate

```
goto l_start(p( $\bar{t}_1, \dots, \bar{t}_n$ ),  $\rho_1$ ,  $B_1$ );
```

13. Finally, we must close the open while-loop (6. above) and finish the code piece in case the loop does not find any (further) matching tuple. Generate:

```
}
break;
```

In addition, there is a code piece for case $L_RESTORE_ρ$ as explained under item 3 above. It is executed when this rule application is finished. It pops everything pushed there and then does `break;` to continue with the next task from the backtrack stack.

6.6 Two IDB-Body Literals

This is a complicated case and needs intermediate storage of tuples generated for the body literals. Fortunately, this occurs rarely in the output of the SLDMagic method (only when translating recursions that are not tail recursions).

A general solution, which does not need information about the order of generated tuples, is to manage one set of tuples for each body literal. In the code piece for the case that a new tuple has been derived for the left body literal, this tuple is joined with all tuples in the current set for the right body literal. In the same way, when a new tuple is generated for the right body literal, it is joined with all existing tuples for the left body literal.

This means that we now need cursors also for the intermediate storage of generated IDB facts, and these cursors must keep information about the last fact when they were created (since new facts can be appended to the list while the cursor is active—these facts must not be returned by the cursor).

Recursion can be handled in the same way as before: When a new fact is generated for a body literal, we save the state of the cursor and all variables of that code piece, continue with derivations using the new fact, and later return to the old fact. However, since we anyway have intermediate storage now, it is also possible to create a queue of facts for each body literal, which must still be used in derivations.

If it is possible to generate all facts for the left body literal before the first fact for the right body literal, one obviously needs intermediate storage only for the left body literal. In this case it can later be treated like an EDB literal.

If one should be able generate facts for both literals in the sort order of the common variables (the join attributes), one would need intermediate storage only for a single tuple for each body literal (we would basically do a merge join).

6.7 Generated Answer-Facts

For rules about the predicate answer, one can print the generated tuple or insert it into a result relation whenever the above code would jump to a body literal which uses the generated fact (there are no body literals with predicate answer). One can also offer the cursor interface of Section 3.

7 Conclusion

We have presented a detailed description of the push method, an efficient bottom-up evaluation algorithm for pure Datalog programs. It is implemented as a translation from Datalog to C++.

In the push method, derived facts are immediately used to derive new facts without generally materializing immediate results. A specific feature introduced in this paper is the representation of derived tuples which significantly reduces the amount of copying. This is the result of the partial evaluation to do as much as possible work already at “compilation time” (at the price of producing several specializations of the same rule, i.e. the generated code might grow while runtime is saved). The rule application graph defined here is useful for planning the evaluation.

First performance tests show some improvement over the previous version of the push method from [2]. We plan to develop a more complete implementation and to investigate further optimizations. The current state of the project is reported at

<http://www.informatik.uni-halle.de/~brass/push/>.

References

- [1] Stefan Brass (2000): *SLDMagic — The Real Magic (with Applications to Web Queries)*. In W. Lloyd et al., editors: *First International Conference on Computational Logic (CL'2000/DOOD'2000)*, LNCS 1861, Springer, pp. 1063–1077, doi:10.1007/3-540-44957-4_71. Available at <http://users.informatik.uni-halle.de/~brass/sldmagic/>.
- [2] Stefan Brass (2010): *Implementation Alternatives for Bottom-Up Evaluation*. In Manuel Hermenegildo & Torsten Schaub, editors: *Technical Communications of the 26th International Conference on Logic Programming (ICLP'10), Leibniz International Proceedings in Informatics (LIPIcs) 7*, Schloss Dagstuhl, pp. 44–53, doi:10.4230/LIPIcs.ICLP.2010.44. Available at <http://users.informatik.uni-halle.de/~brass/botup/>.
- [3] Stefan Brass & Heike Stephan (2013): *A Variant of Earley Deduction with Partial Evaluation*. In Wolfgang Faber & Domenico Lembo, editors: *Web Reasoning and Rule Systems - 7th International Conference, RR 2013, LNCS 7994*, Springer-Verlag, pp. 35–49, doi:10.1007/978-3-642-39666-3_4. Available at <http://dbs.informatik.uni-halle.de/Earley/>.
- [4] Marcia A. Derr, Shinichi Morishita & Geoffrey Phipps (1994): *The Glue-Nail Deductive Database System: Design, Implementation and Evaluation*. *The VLDB Journal* 3, pp. 123–160, doi:10.1007/BF01228879. Available at <http://www.vldb.org/journal/VLDBJ3/P123.pdf>.
- [5] Todd J. Green, Molham Aref & Grigoris Karvounarakis (2012): *LogicBlox, Platform and Language: A Tutorial*. In Pablo Barceló & Reinhard Pichler, editors: *Datalog in Academia and Industry, 2nd Int. Workshop, Datalog 2.0, LNCS 7494*, Springer-Verlag, pp. 1–8, doi:10.1007/978-3-642-32925-8_1. Available at <https://developer.logicblox.com/2012/08/>.
- [6] Senlin Liang, Paul Fodor, Hui Wan & Michael Kifer (2009): *OpenRuleBench: An Analysis of the Performance of Rule Engines*. In: *Proceedings of the 18th International Conference on World Wide Web (WWW'09)*, ACM, pp. 601–610, doi:10.1145/1526709.1526790. Available at <http://rulebench.projects.semwebcentral.org/>.
- [7] Mengchi Liu (2000): *Design and Implementation of the ROL System*. *Journal of Intelligent Information Systems* 14, pp. 1–21, doi:10.1023/A:1008774006482. Available at <http://www.scs.carleton.ca/~mengchi/papers/rol-JIIS00.ps>.
- [8] Kotagiri Ramamohanarao (1993): *An Implementation Overview of the Aditi Deductive Database System*. In Stefano Ceri, Katsumi Tanaka & Shalom Tsur, editors: *Deductive and Object-Oriented Databases, Third Int. Conf., (DOOD'93), LNCS 760*, Springer, pp. 184–203, doi:10.1007/3-540-57530-8_12.
- [9] Konstantinos Sagonas, Terrance Swift & David S. Warren (1994): *XSB as an Efficient Deductive Database Engine*. In Richard T. Snodgrass & Marianne Winslett, editors: *Proc. of the 1994 ACM SIGMOD Int. Conf. on Management of Data (SIGMOD'94)*, pp. 442–453, doi:10.1145/191843.191927. Available at <http://user.it.uu.se/~kostis/Papers/xsbddb.html>.
- [10] Heribert Schütz (1993): *Tupelweise Bottom-up-Auswertung von Logikprogrammen (Tuple-wise bottom-up evaluation of logic programs)*. Ph.D. thesis, TU München.
- [11] Praveen Seshadri, Shaun Flisakowski & Seymour Hersh (1996): *CORAL: The Inside Story. Shocking Hacks Revealed*. Technical Report, Department of Computer Sciences, The University of Wisconsin-Madison. Available at <http://ftp.cs.wisc.edu/coral/doc/Inside.ps>.
- [12] Donald A. Smith & Mark Utting (1999): *Pseudo-Naive Evaluation*. In: *Australasian Database Conference*, pp. 211–223. Available at <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.177.5047>.
<http://www.cs.waikato.ac.nz/research/jstar/1999-ADC-pseudo-naive-eval.pdf>.
- [13] Guizhen Yang & Michael Kifer (2000): *FLORA: Implementing an Efficient DOOD System Using a Tabling Logic Engine*. In W. Lloyd et al., editors: *First International Conference on Computational Logic (CL'2000/DOOD'2000)*, LNCS 1861, Springer, pp. 1078–1093, doi:10.1007/3-540-44957-4_72. Available at <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.110.4304>.