

Optimizing Inequality Joins in Datalog with Approximated Constraint Propagation

Dario Campagna¹, Beata Sarna-Starosta², and Tom Schrijvers³

¹ Dept. of Mathematics and Computer Science, University of Perugia, Italy
dario.campagna@dmi.unipg.it

² LogicBlox Inc., Atlanta, Georgia, USA
bss@logicblox.com

³ Dept. of Applied Mathematics and Computer Science, UGent, Belgium
tom.schrijvers@ugent.be

Abstract. Datalog systems evaluate joins over arithmetic (in)equalities as a naive generate-and-test of Cartesian products. We exploit aggregates in a source-to-source transformation to reduce the size of Cartesian products and to improve performance. Our approach approximates the well-known propagation technique from Constraint Programming. Experimental evaluation shows good run time speed-ups on a range of non-recursive as well as recursive programs. Furthermore, our technique improves upon the previously reported in the literature constraint magic set transformation approach.

1 Introduction

Datalog [13,1] is a syntactic subset of Prolog introduced in the 1980s for database processing. By supporting a limited, safe form of recursion, Datalog considerably extends the expressive power of traditional database query languages like SQL. At the same time, unlike Prolog, Datalog allows SQL's set-at-a-time evaluation. Also similarly to SQL, the programs in Datalog are guaranteed to terminate. Hence, extra-logical constructs such as Prolog's cut ('!') operator are not needed.

After its original introduction as a smarter version of SQL, Datalog lost the interest of researchers for a time, until recently re-gaining attention in applications falling outside of the realm of traditional database reasoning, which include: program analysis [8], networks [12], security protocols [10], knowledge representation [9], robotics [2] and gaming [19]. Our industrial partner, LogicBlox Inc. [11], uses a variant of Datalog, called Datalog^{LB}, as the basis for implementing decision automation and business planning systems.

Many of the above application domains rely on processing numerical data with arithmetic operations, in Datalog available as built-in relations (predicates). We focus in particular on built-in arithmetic (in)equality predicates ($>$, $<$, etc.) which we also refer to as *(arithmetic) constraints*. Existing Datalog compilers do not exploit the constraining properties of arithmetic predicates, but rather implement them as ordinary tests. As a result, evaluation of programs with arithmetic constraints follows the naive *generate-and-test* approach, where ordinary

predicates act as generators, and the entire search space they produce is enumerated before the constraints can be applied to prune the candidate solutions. In database terminology, the full $\mathcal{O}(n^2)$ Cartesian product of two tables is computed. This is in stark contrast with $\mathcal{O}(n \log n)$ equality-based joins, for which current Datalog systems are optimized.

The research area of Constraint Programming (CP) offers approaches that prune the search space more eagerly, e.g., *constrain-and-generate*, as well as the constraint implementation technique, called *constraint propagators*, which allows to prune the domains of the variables involved in the constraints to narrow down the sets of candidate values even before the values are enumerated.

We adapt the CP constraint propagator technique to filter the individual Datalog generators in $\mathcal{O}(n)$ time before they are used in, potentially much smaller, Cartesian products. For this purpose we extended the Datalog^{LB} system with an automatic program transformation framework. Experimental evaluation shows that our technique enables good run-time improvements for a variety of test programs.

2 Datalog^{LB}

LogicBlox is a commercial Datalog-based platform for building enterprise-scale corporate planning and pricing applications. LogicBlox is currently used in several commercial decision automation applications, including retail supply-chain management [14] and software program analysis [3,4,16]. A typical LogicBlox application involves computational analyses that require aggregation across very large data sets combined with simulation and modeling techniques. The platform accommodates these features by means of its custom query language Datalog^{LB}, a type-safe variant of Datalog, based on incremental evaluation, with trigger-like functionality and support for dynamic updates, declarative specification of functional dependencies, non-deterministic choice, stratified negation, meta-programming, and a wide range of extra-logical computations, including aggregation utilized by our optimization approach. In the following paragraphs we outline the main features of Datalog^{LB} and the LogicBlox run-time system. A more exhaustive description of Datalog^{LB} can be found in [21]. Readers familiar with Datalog may want to use this section as a reference when reading the remainder of the paper.

The Datalog^{LB} Language. Figure 1 shows a Datalog^{LB} encoding of the cryptarithmic puzzle LP+FP=PL, the goal of which is to find an assignment of digits to letters that satisfies the equation LP+FP = PL.

The basic programming construct in Datalog^{LB} is the implication ‘ \leftarrow ’, denoting derivation rules of the form:

$$\text{Head} \leftarrow \text{Body}.$$

where **Head** and **Body** are conjunctions of *atoms*. An atom can be either a predicate with variable or constant arguments, a comparison expression, an arithmetic

```

1 digit(_) ->.
2 digit(d), val(d:v) -> uint[8](v), v<=9.
3
4 solution(l,p,f) -> digit(l), digit(p), digit(f).
5 solution(l,p,f) <-
6     digit(l), val(l:vl),
7     digit(p), val(p:vp),
8     digit(f), val(f:vf),
9     vl != 0, vp != 0, vf != 0,
10    vl != vp, vl != vf, vp != vf
11    10*vl+vp + 10*vf+vp = 10*vp+vl.

```

Fig. 1: The Datalog^{LB} encoding of the LP+FP=PL cryptarithmic puzzle.

expressions, or a negated atom. The above rule means that the atoms constituting **Head** can be derived from the atoms constituting **Body**. The example program in Figure 1 contains only one rule (lines 5-11), which derives the facts of the predicate `solution` based on the facts of the predicates `digit` and `val`, and the constraints represented as comparisons and arithmetic expressions on their arguments.

Datalog^{LB} extends Datalog with the notion of an *integrity constraint* of the form:

Lhs -> **Rhs**.

Informally, the above constraint means that if **Lhs** is true, then **Rhs** must also be true, where **Lhs** and **Rhs** are conjunctions of atoms. The difference between a constraint and a rule is that a rule derives data for the atoms in its head, whereas a constraint checks that for the existing data matching its left-hand side, the right-hand side holds. The integrity constraints constitute the basis of Datalog^{LB}'s static type system, which guarantees at compile-time that certain kinds of constraints always hold for all possible instantiations of a given schema. Our approach uses integrity constraints to declare *filter types* which allow to reduce the domains of predicates subjected to arithmetic constraints.

Datalog^{LB} types are represented as unary predicates. Custom types may be defined using *entity predicates*. For instance, in Figure 1, the constraint in line 1 declares the entity predicate `digit`. The constraint in line 4 is a *type declaration* for the predicate `solution`, which states that for every tuple `solution(l,p,f)`, the arguments `l`, `p`, and `f` must be `digit` entities. An entity predicate P may be associated with a *reference mode predicate*, which uniquely identifies each element in P with a value of a primitive type, thus allowing to access the specific entity elements from user applications. For instance, line 2 of Figure 1 declares a reference mode predicate `val`, which associates each entity element `d` in `digit` with `v`, an 8-bit unsigned integer value no greater than 9, thus binding the `digit` type to represent single-digit integers. The syntactic form `val(d:v)` denotes the

one-to-one functional relation between d and v , and is reserved for declaring reference mode predicates. The decision to express digits as entities is dictated by one of the mechanisms contributing to Datalog^{LB}'s termination guarantee, which restricts the use of primitive types as arguments to built-in predicates such as arithmetic operations.

The extra-logical operations supported by LogicBlox, including aggregation computations, are represented by special-syntax rules of the form:

$$\mathbf{result}[x_1, \dots, x_n]=v \leftarrow Op \langle\langle v=Method \rangle\rangle Body.$$

The head of the rule uses Datalog^{LB}'s shorthand notation for declaring functional dependencies: $\mathbf{result}[x_1, \dots, x_n]=v$ declares the predicate \mathbf{result} to be a function from x_1, \dots, x_n to v . The notation also allows declaring singleton (constant) values: $p[]=v$ declares the predicate p to be a singleton that contains only the value v . The value can be retrieved through $p[]$. The right-hand side of the above rule, in addition to the conjunction of atoms in $Body$, includes a directive which specifies the type of the operation to be performed (e.g., aggregation), and the particular method (e.g., finding the minimum value) to be used. For instance, in Section 3.1 we show the following rule which finds the lower bound for the \mathbf{val} predicate:

$$\mathbf{lb_digit}[]=n \leftarrow \mathbf{agg} \langle\langle n=\min(v) \rangle\rangle \mathbf{digit}(d), \mathbf{val}(d:v).$$

Above, \mathbf{agg} states that the rule computes an aggregation, and \mathbf{min} names the specific operation to be applied to the values referenced by v .

3 The Filter Predicates Transformation

This section describes the details of our transformation, beginning with non-recursive programs, and then considering the impact of recursion.

3.1 Non-Recursive Programs

Recall the LP+FP=PL program from the previous section. Our goal is to reduce the number of different candidate values that are used for producing answers. Thus, we exploit the equality constraint

$$10 * v_l + v_p + 10 * v_f + v_p = 10 * v_p + v_l$$

from the program rule to filter candidate values in the generator predicate \mathbf{digit} . Specifically, for each generator predicate atom appearing in the constraint, we consider the value generated by this atom in the context of the upper and lower bounds of the values produced by other generator atoms.

For instance, for the generator atom $\mathbf{digit}(1)$, the original constraint, which is equivalent to the pair of inequalities:

$$\begin{cases} 10 * v_l + v_p + 10 * v_f + v_p \leq 10 * v_p + v_l \\ 10 * v_l + v_p + 10 * v_f + v_p \geq 10 * v_p + v_l \end{cases}$$

yields the pair of inequalities:

$$\begin{cases} 10 * v_l + l_d + 10 * l_d + l_d \leq 10 * u_d + v_l \\ 10 * v_l + u_d + 10 * u_d + u_d \geq 10 * l_d + v_l \end{cases}$$

where u_d and l_d represent the upper and lower bound of the generator predicate `digit`, respectively. We use these inequalities in the Datalog definition of the filter predicate for `digit(l)`, which is linear in the size of the `digit` set.

```
digit_filtered_l(l) <-
  digit(l),
  val(l:v_l),
  lb_digit[]=l_d,
  ub_digit[]=u_d,
  10*v_l+l_d+10*l_d+l_d <= 10*u_d+v_l,
  10*l_d+v_l <= 10*v_l+u_d+10*u_d+u_d.
```

Similar filter predicates are generated for the remaining generator atoms.

The bounds for the generator predicates are computed in separate aggregate predicates, again adding only linear overhead, and reused in all filter predicates:

```
lb_digit[]=n <- agg<<n=min(v)>> digit(d), val(d:v).
ub_digit[]=n <- agg<<n=max(v)>> digit(d), val(d:v).
```

In the last step of the transformation we replace the generator predicate atoms in the body of the `solution/3` rule by atoms representing corresponding filter predicates:

```
solution(l,p,f) <-
  digit_filtered_l(l),
  digit_filtered_p(p),
  digit_filtered_f(f),
  ... % rest of the original LP+FP=PL program
```

As the transformation adds only linear overhead, the overall worst-case time complexity is not increased. Moreover, the filtered generator sets are potentially much smaller than the original sets, thus resulting in a Cartesian product much smaller than the original one. In this small example the filtered generator sets for `l`, `p` and `f` are all reduced from $[0, 9]$ to respectively $[1, 8]$, $[2, 9]$ and $[1, 8]$.

Our approach is inspired by the well-known *bounds consistency* technique [5], in CP implemented by finite-domain constraint propagators. We simplify constraint propagation in two ways: (1) by computing filtered domains on the original domains rather than as a fixed point of the filtering process, and (2) by filtering only at the beginning of the evaluation rather than repeatedly after every enumeration step (in CP terminology known as *labeling*). As a consequence of these simplifications, (1) we cannot encode unbounded fixpoint computations,

```

p(t,w) -> string(t), int[64](w).
s(t,w) -> string(t), int[64](w).

e(t,w) -> string(t), int[64](w).
e(t,w) <- p(t,w).
e(t,w) <- s(t,w),
        e(tp,wp),
        w - wp <= 100,
        w + wp >= 19500.

```

Fig. 2: The `Engine` program.

```

e(t,w) <- p(t,w).
e(t,w) <- s_filtered(t,w),
        e_filtered(tp,wp),
        w-wp <= 100,
        w+wp >= 19500.

s_filtered(t,w) <-
s(t,w),
w-ub_e[] <= 100,
19500 <= w+ub_e[].

e_filtered(tp,wp) <-
e(tp,wp),
lb_s[]-wp <= 100,
19500 <= ub_s[]+wp.

lb_s[]=n <- agg<<n=min(v)>> s(_,v).
ub_s[]=n <- agg<<n=max(v)>> s(_,v).
ub_e[]=n <- agg<<n=max(v)>> e(_,v). % ERROR

```

Fig. 3: Ill-formed `Engine` program after naive transformation.

and (2) computing and storing many successively filtered tables for the same variable adds considerable time and space overhead. Nevertheless, our approach yields a light-weight technique that is easily provided on top of the existing `Data-log` implementations, offering a satisfactory compromise between the anticipated speed-up and the overhead.

3.2 Recursive Programs

Recursion considerably complicates our transformation. Consider the `Engine` program listed in Figure 2. The program selects suitable engines for an engine yard. In the predicates $p(t,w)$, $s(t,w)$ and $e(t,w)$, t corresponds to the engine type and w to the produced wattage. The predicate p represents the primary engines, and the predicate s represents the potentially spare engines. A suitable engine for the engine compound $e(t,w)$ is either a primary engine, or a spare engine that can assist another engine in the compound. The difference in wattage between the assisting engine and the assistee should not exceed 100, and the total wattage of the compound should be no less than 19,500.

The naive application of our technique yields the ill-formed program shown in Figure 3. The program involves recursion through aggregation: in order to

compute the set of $\epsilon/2$ we need to know the upper bound of $\epsilon/2$. Such recursion is not supported by Datalog^{LB} (nor by any other LP system we are aware of).

Since it is not possible to effectively compute the *exact* upper bound on $\epsilon/2$, we approximate it as the upper bound of the approximated upper bounds of the two rules defining $\epsilon/2$. For the first, non-recursive rule, such an approximated (and exact) upper bound is $\text{ub_p}[]$. A crudely approximated upper bound for the second, recursive rule, is $\text{ub_s}[]$. Hence:

$$\text{ub_e}[] = n \leftarrow n = \max(\text{ub_p}[], \text{ub_s}[]).$$

where

$$\text{ub_p}[] = n \leftarrow \text{agg}\langle\langle n = \max(v) \rangle\rangle p(_, v).$$

We may attempt to tighten the upper bound of the second rule, based on the observation that it is bounded from above by $\text{ub_e}[] + 100$:

$$\text{ub_e}[] = n \leftarrow n = \max(\text{ub_p}[], \min(\text{ub_s}[], \text{ub_e}[] + 100)).$$

Alas, this step reintroduces recursion through aggregation. We eliminate it in the same way as before, by substituting the cruder approximation derived earlier:

$$\text{ub_e}[] = n \leftarrow n = \max(\text{ub_p}[], \min(\text{ub_s}[], \max(\text{ub_p}[], \text{ub_s}[]) + 100)).$$

We further simplify the above expression by noticing that

$$\forall x, y, c \in \mathbb{N}. \min(x, \max(y, x) + c) = x$$

This step brings us back to the first approximation, thus proving that the refinement attempt was unsuccessful. Nevertheless, as we show in Section 5, our approximation is still quite effective at pruning the predicate domains and improving the performance of the programs.

4 Implementation

Most of the Datalog^{LB} syntax is compatible with the term syntax of standard Prolog. The discrepancies in the particular notations, such as the functional dependency syntax, can be easily accommodated by simple processing steps. Hence, we chose Prolog (specifically, SWI-Prolog [20]) to implement the transformations of Datalog^{LB} programs. Our analyzer consists of three Prolog modules, for the total of about 1,500 lines of Prolog code, including comments.

4.1 LogicBlox/SWI-Prolog Interface

Figure 4 shows the LogicBlox compilation scheme and outlines the communication between the LogicBlox engine and the SWI-Prolog analyzer.

The LogicBlox compiler rewrites a source Datalog^{LB} program into a core representation, which is then encoded using Google's protocol buffers (GPBs)

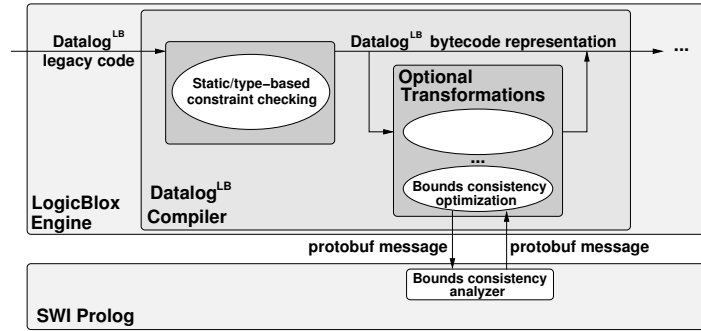


Fig. 4: LogicBlox Compiler and its communication with SWI-Prolog

interface for further use by a number of tools, including an interpreter. GPBs [7] is a platform-independent, extensible mechanism for serializing structured data to binary form. GPBs allow the programmers to determine how to structure their data by defining simple data structures (messages) in a dedicated specification language, and then to compile those data structures into the language and platform of their choice. As shown in Figure 4, the core program representation generated by the LogicBlox compiler is either passed directly to the subsequent phase of run-time processing, or subjected to one or more optional transformations aimed at optimizing the compiled code or collecting information to be used in further evaluation steps. This infrastructure makes the GPBs the medium of choice for interfacing LogicBlox with external analysis modules. The interface for our application comprises three Prolog modules for the total of about 1,700 lines of code (including comments), and five new modules in the LogicBlox code, for the total of 1,800 lines.

The interface on the SWI-Prolog side is based on SWI-Prolog’s native GPBs library [15]. We extended the library with the capability to represent recursively structured data (which is essential to encode Datalog^{LB} programs), and optimized it to linear run-time complexity. Our version of the library is available in a dedicated branch of the SWI-Prolog code repository.

The communication between LogicBlox and SWI-Prolog proceeds as follows. The output of the Datalog^{LB} compiler is received by a new LogicBlox module which extracts from it the information relevant to our analysis, encodes it as a collection of GPBs messages, and opens a socket connection with SWI-Prolog. Once the connection is established, the messages are supplied to our analyzer. The analyzer decodes the messages into a program representation, applies the transformation, encodes the resulting program, and sends it back to the LogicBlox side, where another dedicated module retrieves the transformation results and updates the core representation of the program accordingly.

We illustrate the use of GPBs on Datalog^{LB} rule bodies. A rule body is a formula defined as an atom, a disjunction, a conjunction, or a negation. LogicBlox serializes and deserializes formulas with GPBs messages of the following form:


```

message Formula {
    optional Atom atom = 1;
    optional Negation negation = 2;
    optional Conjunction conjunction = 3;
    optional Disjunction disjunction = 4;
}
message Conjunction { repeated Formula formula = 1; }
...

```

Note the mutually recursive nature of the `Formula` and `Conjunction` definitions. On the SWI-Prolog side, the messages are defined in `message/2` clauses:

```

protobufs:message(formula, [ optional(1,message(atom))
                           , optional(2,message(negation))
                           , optional(3,message(conjunction))
                           , optional(4,message(disjunction)) ]).

protobufs:message(conjunction, [ repeated(1,message(formula)) ]).

```

The predicate `message/2`, which we added to the SWI-Prolog GPBs library, enables naming message templates. It is essential for recursive and repeated embedded messages. The `protobuf_message/2` predicate serializes and deserializes messages to and from binary form, like the representation of the single-atom formula `digit(d)`.

```

?- protobuf_message(message(formula,
    [ optional(1, message(atom,
        [ string(1,"digit")
          , repeated(2,[ /* variable d */ ],term) ]))
    ]),Bytes).

```

4.2 The Transformation

Given a representation of a `DatalogLB` program, our transformation processes in turn each of its rules. For every rule with one or more arithmetic constraints, it identifies the generator predicate atoms, exploits the constraints to produce corresponding filter predicates, and replaces the generator atoms accordingly. It also extends the program with the definitions for the auxiliary predicates performing bounds computations.

Implementation of the code that generates the bounds-computing predicates turned out to be one of the more involved aspects of our project. The numerical data appearing in the arithmetic constraints pertinent to our transformation is often represented as the values of the `DatalogLB` reference mode predicates where the keys are the entities produced by the predicates serving as generators. To access these data, it is necessary to reconstruct the chain of functional dependencies connecting each value with the appropriate entity generator. For instance, to compute predicate bounds for the atom set:

$p(x), \text{val}_1(x:vx), q(y), \text{val}_2(y:vy), vx > vy$

we need to reconstruct the chain connecting vx with p and vy with q . Additional complications arise when the reference mode predicates (and the corresponding generators) have non-unary keys, in which case the reconstructed dependencies are trees with the functional dependencies as nodes and the generators as leaves.

As mentioned in Section 2, Datalog^{LB}'s static type system relies on the type information in the form of the integrity constraints. To ensure completeness of the type information in the transformed programs, we need to provide type declarations for the predicates generated by the analyzer (i.e., filter and aggregate predicates). It turns out that we can conveniently derive these directly from the original predicates, with no additional bookkeeping during the transformation.

5 Evaluation

We now present the results of applying our transformation to a variety of programs. All experiments were performed on a machine with a 2.83 GHz Intel® Core™ 2 Quad CPU and 4 GB of RAM, running Ubuntu 10.10 (Linux kernel 2.6.35-24-server). For each experiment we show the run times, in seconds, for the original programs (*Original* column), and the relative performance change after the filter predicates transformations (*FP* column).

For LogicBlox (v 3.7), in the *Opt* column, we additionally measure the impact of the system's optimizer [17] aimed at improving the performance of equality-based joins by reordering the goals and applying a variant of magic-set rewrite.

In order to have a point of reference, we also report the results of tabled top-down evaluation of our test suites using XSB Prolog 3.3.1. The changes required to accommodate Datalog^{LB} programs in XSB are minimal and mainly syntactic in nature: we omit type declarations, replace '<->' arrows with ':-', capitalize variable names, change functional dependencies to ordinary arguments, and provide Prolog implementations for aggregates. To guarantee termination, we declare all predicates as tabled.

5.1 Non-Recursive Programs

Cryptarithmic Puzzles. Table 1 shows the evaluation run times for a set of cryptarithmic puzzles building on the idea of the LP+FP=PL program from Section 2. In almost all cases the transformation yields drastic performance improvements (3× to 10×) over both original and optimized LogicBlox evaluation. There are two exceptions. In the first case, the overhead of the auxiliary predicates introduced by the transformation dominates the extremely short run time of the original program. In the second case, the transformed program prunes very few values from the initial domains, and consequently shows performance similar to that of the original program.

The XSB evaluation yields similar results both in terms of the original program performance, and the benefits from the transformation.

Puzzle	Datalog ^{LB}			XSB	
	Original	Opt	FP	Original	FP
Puzzle 1	0.01 sec.	100.00 %	140.90 %	0.01 sec.	100.00 %
LP+FP=PL	0.01 sec.	100.00 %	100.00 %	0.01 sec.	100.00 %
Puzzle 2	0.80 sec.	72.50 %	14.02 %	0.65 sec.	15.38 %
Puzzle 3	3.10 sec.	25.16 %	11.42 %	2.60 sec.	11.92 %
Puzzle 4	2.67 sec.	104.49 %	12.79 %	2.73 sec.	12.09 %
Puzzle 5	6.39 sec.	114.71 %	15.02 %	7.70 sec.	12.60 %
Puzzle 6	3.90 sec.	82.56 %	27.05 %	8.75 sec.	25.26 %
Puzzle 7	17.54 sec.	50.85 %	105.01 %	17.20 sec.	107.62 %
Puzzle 8	20.63 sec.	92.05 %	11.71 %	19.99 sec.	52.53 %

Table 1: Benchmark results for cryptarithmic puzzles.

Tons range	Datalog ^{LB}			XSB	
	Original	Opt	FP	Original	FP
[1,500]	0.60 sec.	101.67 %	103.64 %	0.29 sec.	96.55 %
[1,1000]	2.81 sec.	46.26 %	100.75 %	1.10 sec.	98.18 %
[1,2500]	12.37 sec.	42.52 %	40.60 %	5.01 sec.	92.41 %
[1,5000]	13.71 sec.	43.69 %	41.27 %	5.90 sec.	88.30 %

Table 2: Benchmark results for the Production problem.

The Production Problem. The `Production` program⁴ models the mathematical programming problem of optimizing the profit from manufacturing several types of products, subject to a set of constraints such as production costs and maximum number of items to be manufactured for each product type, or the availability of the factory line. From a technical point of view this program is interesting because it contains multi-key functional dependencies that drive the filter predicates. Another non-standard feature is the use of the aggregates for computing the optimized profit.

Table 2 reports the results of evaluating the original and transformed program with four data sets differing in the range of the generator predicate indicating the number of tons of products being manufactured. Clearly, for LogicBlox evaluation, the transformation has no significant effect on the program for the small tons ranges, but enables a lot of pruning, and thus considerable performance improvement, when the tons ranges are large. On XSB the effects of the transformation are more uniform across the different data sets, with slightly better performance improvements for the larger tons ranges.

5.2 Recursive programs

The Engine Program. To evaluate the effects of our transformation on the recursive `Engine` program from Figure 2, we used four different data sets. Each

⁴ We refer to <http://users.ugent.be/~tschrijv/Datalog> for the source code.

Data set	Datalog ^{LB}			XSB	
	Original	Opt	FP	Original	FP
<i>Set</i> ₁	26.87 sec.	106.43 %	21.41 %	43.40 sec.	6.11 %
<i>Set</i> ₂	9.82 sec.	106.92 %	4.65 %	8.29 sec.	0.84 %
<i>Set</i> ₃	172.47 sec.	100.93 %	84.18 %	119.82 sec.	64.91 %
<i>Set</i> ₄	53.61 sec.	100.39 %	104.75 %	20.30 sec.	97.93 %

Table 3: Benchmark results for the **Engine** program.

data set defines the sets of couples produced by $p/2$ (denoted P in the following), and $s/2$ (denoted S). Let

$$\mathcal{T} = \{\text{Steam engine, Internal combustion engine, Gas Turbine}\}$$

The four data sets define the sets P and S as follows.

$$\begin{aligned}
- \textit{Set}_1: & \begin{cases} P = \mathcal{T} \times [1100, 11500] \\ S = \mathcal{T} \times [1, 10000] \end{cases} & - \textit{Set}_3: & \begin{cases} P = \mathcal{T} \times [500, 16000] \\ S = \mathcal{T} \times [1000, 14000] \end{cases} \\
- \textit{Set}_2: & \begin{cases} P = \mathcal{T} \times [500, 5000] \\ S = \mathcal{T} \times [1, 6000] \end{cases} & - \textit{Set}_4: & \begin{cases} P = \mathcal{T} \times [10000, 16000] \\ S = \mathcal{T} \times [8, 12000] \end{cases}
\end{aligned}$$

The results of the evaluation are shown in Table 3. There is a visible correlation between the particular data set and the effects of the transformation. With little pruning comes modest speed-up or even a slow-down, whereas considerable pruning yields large performance improvements. Again our transformation achieves drastic improvements where the LogicBlox optimizer does not.

Multi-Legged Flights Program. The **Flights** program (Figure 5) models multi-legged flights and their travel distance. More abstractly, it captures the transitive closure of a directed weighted graph. The Datalog^{LB} encoding consists of the basic variant of the program, based on that studied by Stuckey and Sudarshan [18], together with a sample query to compute all possible destinations no further than 10,000 miles from Sydney.

Predicate $e(x, y, d)$ (line 1) denotes a flight leg, i.e., a direct connection between cities x and y with the distance d . The data of this predicate are given as facts. The predicate f (lines 3-7) defines a multi-legged flight as the transitive closure of the predicate e . Since the second rule for f contains recursion, to be expressible in Datalog^{LB}, it needs to be bounded. Hence, we have added the constraint ‘ $d \leq 10000$ ’ (line 7), which is not present in the encoding of [18]. Lines 9-10 define the **query** predicate.

It turns out that our transformation has no significant effect on the performance of the **Flights** program; it does not provide additional pruning. Fortunately, to our aid comes the *constraint magic set transformation* [18]. Not only is the constraint magic set rewritten (CMR) variant of the program (Figure 6) faster than the original, but also it is amenable to our transformation.

```

1 e(x,y,d) -> string(x), string(y), int[64](d).
2
3 f(x,y,d) -> string(x), string(y), int[64](d).
4 f(x,y,d) <- e(x,y,d), d >= 0.
5 f(x,y,d) <- e(x,z,d1), d1 >= 0,
6           f(z,y,d2), d2 >= 0,
7           d = d1 + d2, d <= 10000.
8
9 query(x,y,d) -> string(x), string(y), int[64](d).
10 query("Sydney",y,d) <- f("Sydney",y,d), d >= 0, d <= 10000.

```

Fig. 5: The Datalog^{LB} encoding of the `Flights` program.

Table 4 shows the results of evaluating the CMR variant of the `Flights` program without (CMR) and with (CMR+FP) filter predicate transformation for a collection of 19 different data graphs, with different structures.

For the LogicBlox evaluation, Table 4 reports performance decrease for three transformed programs with corresponding original run times below 0.1s, and visible improvement for all other benchmarks. The speed-up varies roughly between $2\times$ for the original programs with the shorter run times and $8\times$ for those with longer run times. Interestingly, the performance in XSB is very different. First, we observe that the run times for programs without the transformation are considerably shorter than in LogicBlox. Furthermore, applying the transformation has no effect on the three programs with the shortest original run times, whereas it significantly slows down the evaluation of all other programs. We attribute this negative effect to the ordering of constraints—imposed by our transformation when introducing filter predicates—which forces overhead computations in the order-sensitive XSB.

6 Conclusion and Future Work

We presented a technique exploiting Datalog with aggregates to improve the performance of Datalog^{LB} programs with arithmetic (in)equalities. Our approach employs a source-to-source program transformation that approximates the propagation technique from Constraint Programming. The experimental evaluation of the approach shows good run time speed-ups on a range of non-recursive as well as recursive programs. Furthermore, our technique improves upon the constraint magic set transformation approach proposed by Stuckey and Sudarshan.

In the future we plan to investigate ways to integrate finite domain solvers with the Datalog’s semi-naive bottom-up evaluation mechanism to enable further benefits from constraint propagation. We would also like to compare our transformation-based approach to the tabled constraint programming approach proposed by Cui and Warren [6], applied to a finite domain constraint solver.

```

answer_f(x,y,d) -> string(x), string(y), int[64](d).
answer_f(x,y,d) <- x = "Sydney", f_a(x,y,d), d >= 0, d <= 10000.

f_a(x,y,d) -> string(x), string(y), int[64](d).
f_a(x,y,d) <- query_f_a(x,ld,ud), ld <= ud,
              e(x,y,d), d >= 0, d >= ld, d <= ud.
f_a(x,y,d) <- query_f_a(x,ld,ud), ld <= ud,
              e(x,z,d1), d1 >= 0,
              f_a(z,y,d2), d2 >= 0,
              d = d1 + d2, d >= ld, d <= ud.

query_f_a(x,ld,ud) -> string(x), int[64](ld), int[64](ud).
query_f_a("Sydney",0,10000).
query_f_a(y,ld2,ud2) <- query_f_a(x,ld,ud), ld <= ud,
                      e(x,y,d), d >= 0,
                      ud2 = ud - d, ld2 = max(ld-d,0).

e(x,y,d) -> string(x), string(y), int[64](d).

```

Fig. 6: Constraint magic rewritten variant of the `Flights` program.

References

1. S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
2. M. P. Ashley-Rollman, M. De Rosa, S. S. Srinivasa, P. Pillai, S. C. Goldstein, and J. D. Campbell. Declarative Programming for Modular Robots. In *Workshop on Self-Reconfigurable Robots/Systems and Applications at IROS*, 2007.
3. M. Bravenboer and Y. Smaragdakis. Exception Analysis and Points-To Analysis: Better Together. In *ISSTA*, pages 1–12, 2009.
4. M. Bravenboer and Y. Smaragdakis. Strictly Declarative Specification of Sophisticated Points-To Analyses. In *OOPSLA*, pages 243–262, 2009.
5. C. Choi, W. Harvey, J. Lee, and P. Stuckey. Finite domain bounds consistency revisited. In *AI 2006: Advances in Artificial Intelligence*, volume 4304 of *Lecture Notes in Computer Science*, pages 49–58. Springer Berlin / Heidelberg, 2006.
6. B. Cui and D. S. Warren. A System for Tabled Constraint Logic Programming. In *CL*, pages 478–492, 2000.
7. Google’s Protocol Buffers. <http://code.google.com/apis/protocolbuffers/>.
8. M. S. Lam, J. Whaley, V. B. Livshits, M. C. Martin, D. Avots, M. Carbin, and C. Unkel. Context-sensitive program analysis as database queries. In *PODS*, pages 1–12, 2005.
9. N. Leone, G. Pfeifer, W. Faber, T. Eiter, G. Gottlob, S. Perri, and F. Scarcello. The DLV system for knowledge representation and reasoning. *ACM Trans. Comput. Logic*, 7(3):499–562, 2006.
10. N. Li and J. C. Mitchell. Datalog with constraints: A foundation for trust management languages. In *PADL*, pages 58–73, 2003.
11. LogicBlox. <http://logicblox.com/>.

Graph	Datalog ^{LB}		XSB	
	CMR	CMR+FP	CMR	CMR+FP
Graph 1	0.01 sec.	191.1 %	0.01 sec.	100.0 %
Graph 2	0.03 sec.	162.0 %	0.01 sec.	100.0 %
Graph 3	0.02 sec.	117.2 %	0.01 sec.	100.0 %
Graph 4	0.19 sec.	54.5 %	0.02 sec.	250.0 %
Graph 5	4.47 sec.	21.2 %	0.51 sec.	468.6 %
Graph 6	0.24 sec.	63.5 %	0.04 sec.	925.0 %
Graph 7	0.76 sec.	41.7 %	0.12 sec.	2266.7 %
Graph 8	2.91 sec.	22.1 %	0.31 sec.	442.8 %
Graph 9	65.79 sec.	13.3 %	5.28 sec.	988.8 %
Graph 10	5.76 sec.	42.0 %	1.26 sec.	504.8 %
Graph 11	1.94 sec.	21.3 %	0.19 sec.	163.1 %
Graph 12	2.40 sec.	38.0 %	0.39 sec.	2761.5 %
Graph 13	2.83 sec.	22.1 %	0.29 sec.	320.7 %
Graph 14	4.99 sec.	25.5 %	0.73 sec.	291.8 %
Graph 15	66.93 sec.	13.0 %	5.14 sec.	1010.9 %
Graph 16	1.92 sec.	22.9 %	0.17 sec.	170.6 %
Graph 17	2.85 sec.	21.5 %	0.27 sec.	340.7 %
Graph 18	1.92 sec.	21.4 %	0.16 sec.	181.2 %
Graph 19	67.60 sec.	13.3 %	5.06 sec.	1030.0 %

Table 4: Benchmark results for the `Flights` program.

12. B. Thau Loo, T. Condie, M. N. Garofalakis, D. E. Gay, J. M. Hellerstein, P. Maniatis, R. Ramakrishnan, T. Roscoe, and I. Stoica. Declarative networking: language, execution and optimization. In *SIGMOD*, pages 97–108, 2006.
13. D. Maier and D. S. Warren. *Computing with Logic: Logic Programming with Prolog*. Benjamin/Cummings, 1988.
14. Predictix. <http://www.predictix.com/>.
15. J. Rosenwald. SWI-Prolog Google’s Protocol Buffers library. <http://www.swi-prolog.org/pldoc/package/protobufs.html>.
16. Semmle. <http://semmle.com/>.
17. D. Sereni, P. Avgustinov, and O. de Moor. Adding magic to an optimising Datalog compiler. In *SIGMOD*, pages 553–565, 2008.
18. P. J. Stuckey and S. Sudarshan. Compiling query constraints (extended abstract). In *PODS*, pages 56–67, 1994.
19. W. White, A. Demers, C. Koch, J. Gehrke, and R. Rajagopalan. Scaling games to epic proportions. In *SIGMOD*, pages 31–42, 2007.
20. J. Wielemaker. SWI-Prolog 5.10 Reference Manual. <http://www.swi-prolog.org>, April 2010.
21. D. Zook, E. Pasalic, and B. Sarna-Starosta. Typed Datalog. In *PADL*, pages 168–182, 2009.