

# Evaluation of Logic Programs with Built-Ins and Aggregation: A Calculus for Bag Relations

Matthew Francis-Landau<sup>[0000–0002–4139–1111]</sup>, Tim Vieira<sup>[0000–0002–2043–1073]</sup>,  
and Jason Eisner<sup>[0000–0002–8861–0772]</sup>

Johns Hopkins University, Baltimore, MD 21218 USA  
`{mfl,timv,jason}@cs.jhu.edu`

**Abstract.** We present a scheme for translating logic programs with built-ins and aggregation into algebraic expressions that denote bag relations over ground terms of the Herbrand universe. To evaluate queries against these relations, we develop an operational semantics based on term rewriting of the algebraic expressions. This approach can exploit arithmetic identities and recovers a range of useful strategies, including lazy strategies that defer work until it becomes both possible and necessary. Code is available at <https://github.com/matthewfl/dyna-R>.

**Keywords:** Logic Programming, Relational Algebra, Rewriting Systems

## 1 Introduction

We are interested in developing execution strategies for deductive databases whose defining rules make use of aggregation, recursion, and arithmetic. Languages for specifying such deductive databases are expressive enough that it can be challenging to answer queries against a given database. Term rewriting systems are an attractive approach because they can start with the query-database pair itself as an intensional description of the answer, and then attempt to rearrange it into a more useful extensional description such as a relational table. We will give a notation for algebraically constructing potentially infinite bag relations from constants and built-in relations by unions, joins, projection, aggregation, and recursion. We show how programs in existing declarative languages such as Datalog, pure Prolog, and Dyna can be converted into this notation.

We will then sketch a term rewriting system for simplifying these algebraic expressions, which can in principle be used to answer queries against a bag relation. Term rewriting naturally handles materialization, delayed constraints, constraint propagation, short-circuit evaluation, lazy iteration, and enumerative strategies such as nested-iterator join.

There remains a practical challenge (which is beyond the scope of this paper): determining which rewrites to apply and when, much as in automated theorem proving. Our current implementation is essentially a priority rewrite

system with heuristic priorities and memoization.<sup>1</sup> While an improved execution engine is work in progress, our design includes fair (breadth-first) nondeterministic search, polyvariant specialization through generating new rewrite rules and programmable reduction strategies, static analysis by abstract interpretation, and guess-check-update strategies to solve recursive systems of constraints.

## 1.1 Approach

Dyna [4] is a generalization of Datalog [1,7] and pure Prolog [3,2]. Our methods apply to all three of these logic programming languages (and more trivially to languages like SQL that can be written in standard relational algebra).

We are given a Herbrand universe  $\mathcal{G}$  of ground terms. A Dyna program serves to define a partial map from  $\mathcal{G}$  to  $\mathcal{G}$ , which may be regarded as a set of key-value pairs. Datalog and Prolog are similar, but they can map a key only to **true**, so the program serves only to define the set of keys.

Given a program, a user may query the value of a specific key (ground term). More generally, a user may use a non-ground term to query the values of all keys that match it, so that the answer is itself a set of key-value pairs.

A set of key-value pairs—either a program or the answer to a query—is a relation on two  $\mathcal{G}$ -valued variables. Our method in this paper will be to describe the desired relation algebraically, building up the description from simpler relations using a relational algebra. These simpler relations can be over any number of variables; they may or may not have functional dependencies as the final key-value relation does; and they may be bag relations, i.e., a given tuple may appear in the relation more than once, or even infinitely many times. Given this description of the desired relation, we will use rewrite rules to simplify it into a form that is more useful to the user. Although we use a **term rewriting system**, we refer to the descriptions being rewritten as **R-exprs** (relational expressions) to avoid confusion with the terms of the object language, Dyna.

## 1.2 Dyna Examples

To illustrate the task, we first briefly give a couple of examples (adapted from [9]) of useful Dyna programs and queries against them. In §4, we will sketch how to translate Dyna programs into our algebra.

First, we start with a canonical Datalog program written in standard Dyna notation to compute the shortest path in a graph. Additional rules not shown here define the values of **start** and the various **edge** terms.

```

1 | path(start) min= 0.
2 | path(Sink)  min= path(Source) + edge(Source, Sink).
```

<sup>1</sup> This is already more flexible than Prolog’s SLD resolution, which confines its attention at any time to a specific subgoal term and will crash with an “instantiation fault” if that subgoal cannot be rewritten.

This program defines a map with keys such as `edge("albany", "buffalo")`, whose value is the distance from Albany to Buffalo, and keys such as `path("chicago")`, whose value is the total length of the shortest path from `start` to Chicago.

The second rule implies that the value of `path("chicago")` is the minimum value achieved by `path(Source) + edge(Source, "chicago")` for any instantiation of the variable `Source`. If there are no instantiations of `Source` such that `path(Source)` and `edge(Source, "chicago")` both have values, then `path("chicago")` is not a key in the database at all, meaning that Chicago is not reachable at all from `start`.<sup>2</sup>

How is this database used? A user who is located at `start` might query `path("chicago")` to find out how far away it is, or they might query `path(Y)` to find all reachable cities `Y` along with how far away they are. Other queries would return other objects in the database, such as edges.

In the above example, `Sink` and `Source` may range over only a finite set of cities. However, we can easily encounter cases that define infinite relations, as in the following program that runs a simple convolutional neural network:

```

3 |  $\sigma(X) = 1/(1+\exp(-X))$ . % define sigmoid function
4 | out(J) =  $\sigma$ (in(J)). % apply sigmoid function
5 | in(J) += out(I) * edge(I,J). % vector-matrix product
6 | in(input(X,Y)) += pixel_brightness(X,Y) % external input
7 | loss += (out(J) - target(J))**2. % L2 loss of the predictions
8 | edge(input(X,Y),hidden(X+DX,Y+DY)) = weight_conv(DX,DY). % layer 1
9 | edge(hidden(X,Y),output(Z)) = weight_output(X,Y,Z). % layer 2
10 | weight_output(X,Y,Z) := random(*,-1,1). % init with random
11 | weight_conv(DX,DY) := random(*,-1,1) for DX:-4..4, DY:-4..4.

```

Without giving a detailed exposition of this program, we point out that it again has edge keys, which specify the weighted edges of a neural network. This time, however, the rules that specify such edges define infinitely many of them, in a convolutional structure on an infinite set of neurons.

As a result, a query `edge(I,J)` must return a description of an infinite set of edges. Even so, only finitely many of these edges contribute to the value of `loss`, provided that the input image specifies `pixel_brightness(X,Y)` at only finitely many `(X,Y)` coordinates, and the loss function specifies `target` values for only finitely many neurons. As a result, a clever system will be able to answer a query for `loss` in finite time, essentially by finding the paths between the pixels and the targets (which requires addition/subtraction of `DX` and `DY`) and by determining which of the infinitely many keys  `$\sigma(X)$`  need to compute their values in order to find the out activations of the neurons on these paths.

A version of this program indeed runs in our present implementation, although there is not space here to work through such a detailed example. Broadly speaking, the implementation unrolls the definition of `loss` until it is possible to apply rewrites that can make progress on simplifying the definition, for example, by performing arithmetic on known quantities.

<sup>2</sup> Unless "chicago" happens to be the value of the `start` key, in which case the first rule comes into play as well. In this case 0 is included in the minimum, so Chicago is always reachable with total distance of at most 0.

## 2 Syntax and Semantics of **R**-exprs

Let  $\mathcal{G}$  be the Herbrand universe of **ground terms** built from a given set  $\mathcal{F}$  of ranked functors. We treat constants (including numeric constants) as 0-ary functors. Let  $\mathcal{M} = \mathbb{N} \cup \{\infty\}$  be the set of **multiplicities**. A simple definition of a **bag relation** [8] would be a map  $\mathcal{G}^n \rightarrow \mathcal{M}$  for some  $n$ . Such a map would assign a multiplicity to each possible *ordered  $n$ -tuple* of ground terms. However, we will use names rather than positions to distinguish the roles of the  $n$  objects being related: in our scheme, the  $n$  tuple elements will be named by variables.

Let  $\mathcal{V}$  be an infinite set of distinguished **variables**. A **named tuple**  $E$  is a function mapping some of these variables to ground terms. For any  $\mathcal{U} \subseteq \mathcal{V}$ , we can write  $\mathcal{G}^{\mathcal{U}}$  for the environments  $E : \mathcal{U} \mapsto \mathcal{G}$  with domain  $\mathcal{U}$ . These are just the possible  **$\mathcal{U}$ -tuples**: that is, tuples over  $\mathcal{G}$  whose elements are named by  $\mathcal{U}$ . This set  $\mathcal{G}^{\mathcal{U}}$  of named  $\mathcal{U}$ -tuples replaces the set  $\mathcal{G}^n$  of ordered  $n$ -tuples from above.

Below, we will inductively define the set  $\mathcal{R}$  of **R**-exprs. The reader may turn ahead to later sections to see some examples. Each **R**-expr  $R$  has a finite set of **free variables**  $\text{vars}(R) \subseteq \mathcal{V}$ , namely the variables that appear in  $R$  in positions where they are not bound by an operator such as **proj** or **sum**. The idea is for  $R$  to specify a bag relation over domain  $\mathcal{G}$ , with columns named by  $\text{vars}(R)$ .

The **denotation function**  $\llbracket \cdot \rrbracket_E$  interprets **R**-exprs in the **environment**  $E$ . It defines a multiplicity  $\llbracket R \rrbracket_E$  for any **R**-expr  $R$  whose  $\text{vars}(R) \subseteq \text{domain}(E)$ .

If  $\mathcal{U} \supseteq \text{vars}(R)$ , we can dually regard  $R$  as inducing the map  $\mathcal{G}^{\mathcal{U}} \rightarrow \mathcal{M}$  defined by  $E \mapsto \llbracket R \rrbracket_E$ . In other words, given  $\mathcal{U}$ ,  $R$  specifies a bag relation whose column names are  $\mathcal{U}$ . This is true for *any*  $\mathcal{U} \supseteq \text{vars}(R)$ , but the relation constrains only the columns  $\text{vars}(R)$ . The other columns can take any values in  $\mathcal{G}$ . A tuple's multiplicity never depends on its values in those other columns, since our definition of  $\llbracket \cdot \rrbracket_E$  will ensure that  $\llbracket R \rrbracket_E$  depends only on the restriction of  $E$  to  $\text{vars}(R)$ .

We say that  $T$  is a **term** if  $T \in \mathcal{V}$  or  $T = f(T_1, \dots, T_n)$  where  $f \in \mathcal{F}$  has rank  $n$  and  $T_1, \dots, T_n$  are also terms. Terms typically appear in the object language (e.g., Dyna) as well as in our meta-language (**R**-exprs). Let  $\mathcal{T} \supseteq \mathcal{G}$  be the set of terms. Let  $\text{vars}(T)$  be the set of vars appearing in  $T$ , and extend  $E$  in the natural way over terms  $T$  for which  $\text{vars}(T) \subseteq \text{domain}(E)$ :  $E(f(T_1, \dots, T_n)) = f(E(T_1), \dots, E(T_n))$ .

We now define  $\llbracket R \rrbracket_E$  for each type of **R**-expr  $R$ , thus also presenting the different types of **R**-exprs in  $\mathcal{R}$ . First, we have **equality constraints** between non-ground terms, which are true in an environment that grounds those terms to be equal. True is represented by multiplicity 1, and false by multiplicity 0.

1.  $\llbracket T=U \rrbracket_E = \text{if } E(T) = E(U) \text{ then } 1 \text{ else } 0, \quad \text{where } T, U \in \mathcal{T}$

We also have **built-in constraints**, such as

2.  $\llbracket \text{plus}(I, J, K) \rrbracket_E = \text{if } E(I) + E(J) = E(K) \text{ then } 1 \text{ else } 0, \quad \text{where } I, J, K \in \mathcal{T}$

The above **R**-exprs are said to be **constraints** because they always have multiplicity 1 or 0 in any environment. Taking the **union** of **R**-exprs via  $+$  may yield larger multiplicities:

3.  $\llbracket R+S \rrbracket_E = \llbracket R \rrbracket_E + \llbracket S \rrbracket_E, \quad \text{where } R, S \in \mathcal{R}$

The **R**-expr  $\emptyset$  denotes the empty bag relation, and more generally,  $\mathbf{M} \in \mathcal{M}$  denotes the bag relation that contains  $\mathbf{M}$  copies of every  $\mathcal{U}$ -tuple:

$$4. \llbracket \mathbf{M} \rrbracket_E = \mathbf{M}, \quad \text{where } \mathbf{M} \in \mathcal{M}$$

To **intersect** two bag relations, we must use multiplication  $*$  to combine their multiplicities [8]:

$$5. \llbracket \mathbf{R} * \mathbf{S} \rrbracket_E = \llbracket \mathbf{R} \rrbracket_E \cdot \llbracket \mathbf{S} \rrbracket_E, \quad \text{where } \mathbf{R}, \mathbf{S} \in \mathcal{R}$$

Here we regard both  $\mathbf{R}$  and  $\mathbf{S}$  as bag relations over columns  $\mathcal{U} \supseteq \text{vars}(\mathbf{R}) \cup \text{vars}(\mathbf{S}) = \text{vars}(\mathbf{R} * \mathbf{S}) = \text{vars}(\mathbf{R} + \mathbf{S})$ . The names **intersection**, **join** (or **equijoin**), and **Cartesian product** are conventionally used for the cases of  $\mathbf{R} * \mathbf{S}$  where (respectively)  $\text{vars}(\mathbf{R}) = \text{vars}(\mathbf{S})$ ,  $|\text{vars}(\mathbf{R}) \cap \text{vars}(\mathbf{S})| = 1$ , and  $|\text{vars}(\mathbf{R}) \cap \text{vars}(\mathbf{S})| = 0$ . As a special case of Cartesian product, notice that  $\mathbf{R} * 3$  denotes the same bag relation as  $\mathbf{R} + \mathbf{R} + \mathbf{R}$ .

We next define **projection**, which removes a named column (variable) from a bag relation, summing the multiplicities of rows (tuples) that have thus become equal. When we translate a logic program into an **R**-expr (§4), we will generally apply projection operators to each rule to eliminate that rule's local variables.

$$6. \llbracket \text{proj}(\mathbf{X}, \mathbf{R}) \rrbracket_E = \sum_{x \in \mathcal{G}} \llbracket \mathbf{R} \rrbracket_{E[\mathbf{X}=x]} \\ \text{where } \mathbf{X} \in \mathcal{V}, \mathbf{R} \in \mathcal{R}, \text{ and where } E[\mathbf{X}=x] \text{ means a version of } E \text{ that has been} \\ \text{modified to put } E(\mathbf{X}) = x \text{ (adding } \mathbf{X} \text{ to its domain if it is not already there)}$$

Projection collapses each group of rows that are identical except for their value of  $\mathbf{X}$ . **Summation** does the same, but instead of adding up the multiplicities of the rows in each possibly empty group, it adds up their  $\mathbf{X}$  values to get a  $\mathbf{Y}$  value for the new row, which has multiplicity 1. Thus, it removes column  $\mathbf{X}$  but introduces a new column  $\mathbf{Y}$ . The summation operator is defined as follows (note that  $\text{sum} \notin \mathcal{F}$ ):

$$7. \llbracket \mathbf{A} = \text{sum}(\mathbf{X}, \mathbf{R}) \rrbracket_E = \text{if } E(\mathbf{A}) = \sum_{x \in \mathcal{G}} x * \llbracket \mathbf{R} \rrbracket_{E[\mathbf{X}=x]} \text{ then } 1 \text{ else } 0 \\ \text{where } \mathbf{A}, \mathbf{X} \in \mathcal{V}, \mathbf{R} \in \mathcal{R}, \text{ and the } * \text{ in the summand means that we sum up} \\ \llbracket \mathbf{R} \rrbracket_{E[\mathbf{X}=x]} \text{ copies of the value } x \text{ (perhaps } \infty \text{ copies). If there are no summands,} \\ \text{the result of the } \sum \cdots \text{ is defined to be the identity element } \text{id}_{\text{sum}}.$$

Notice that an **R**-expr of this form is a constraint.  $\text{sum}$  is just one type of **aggregation operator**, based on the binary  $+$  operation and its identity element  $\text{id}_{\text{sum}} = 0$ . In exactly the same way, one may define other aggregation operators such as  $\text{min}$ , based on the binary  $\text{min}$  operation and its identity element  $\text{id}_{\text{min}} = \infty$ . Variants of these will be used in §4 to implement the aggregations in  $\text{+=}$  and  $\text{min=}$  rules like those in §1.2.

In projections  $\text{proj}(\mathbf{X}, \mathbf{R})$  and aggregations such as  $\text{sum}(\mathbf{X}, \mathbf{R})$  and  $\text{min}(\mathbf{X}, \mathbf{R})$ , we say that occurrences of  $\mathbf{X}$  within  $\mathbf{R}$  are **bound** by the projection or aggregation operator,<sup>3</sup> so that they are not in the vars of the resulting **R**-expr. However, the most basic aggregation operator does not need to bind a variable:

$$8. \llbracket \mathbf{M} = \text{count}(\mathbf{R}) \rrbracket_E = \text{if } E(\mathbf{M}) = \llbracket \mathbf{R} \rrbracket_E \text{ then } 1 \text{ else } 0$$

<sup>3</sup> But notice that  $\text{sum}(\mathbf{X}, \mathbf{R})$  does not correspond to  $\sum_{\mathbf{X}} \cdots$  but rather to  $\sum_{\text{row} \in \mathbf{R}} \text{row}[\mathbf{X}]$ .

In effect,  $\mathbf{M}=\text{count}(\mathbf{R})$  is a version of  $\mathbf{R}$  that changes every tuple's multiplicity to 1 but records its original multiplicity in a new column  $\mathbf{M}$ . It is equivalent to  $\mathbf{M}=\text{sum}(\mathbf{N}, (\mathbf{N}=1)*\mathbf{R})$  (where  $\mathbf{N} \notin \text{vars}(\mathbf{R})$ ), but we define it separately here so that it can later serve as an intermediate form in the operational semantics.

Finally, it is convenient to augment the built-in constraint types with **user-defined** relation types. Choose a new functor of rank  $n$  that is  $\notin \mathcal{F}$ , such as  $\mathbf{f}$ , and choose some  $\mathbf{R}$ -expr  $\mathbf{R}_f$  with  $\text{vars}(\mathbf{R}_f) \subseteq \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  (which are  $n$  distinct variables) to serve as the **definition** (macro expansion) of  $\mathbf{f}$ . Now define

9.  $\llbracket \mathbf{f}(\mathbf{T}_1, \dots, \mathbf{T}_n) \rrbracket_E = \llbracket \mathbf{R}_f\{\mathbf{x}_1 \mapsto \mathbf{T}_1, \dots, \mathbf{x}_n \mapsto \mathbf{T}_n\} \rrbracket_E$  where  $\mathbf{T}_1, \dots, \mathbf{T}_n \in \mathcal{T} \cup \mathcal{R}$   
 The  $\{\mapsto\}$  notation denotes substitution for variables, where bound variables of  $\mathbf{R}_f$  are renamed to avoid capturing free variables of the  $\mathbf{T}_i$ .

With user-defined relation types, it is possible for a user to write  $\mathbf{R}$ -exprs that are circularly defined in terms of themselves or one another (similarly to a `let rec` construction in functional languages). Indeed, a Dyna program normally does this. In this case, the definition of  $\llbracket \cdot \rrbracket_E$  is no longer a well-founded inductive definition. Nonetheless, we can still interpret the  $\llbracket \cdot \rrbracket_E = \dots$  equations in the numbered points above as *constraints* on  $\llbracket \cdot \rrbracket_E$ , and attempt to solve for a denotation function  $\llbracket \cdot \rrbracket_E$  that satisfies these constraints [4]. Some circularly defined  $\mathbf{R}$ -exprs may be constructed so as to have unique solutions, but this is not the case in general.

### 3 Rewrite Rules

Where the previous section gave a denotational semantics for  $\mathbf{R}$ -exprs, we now sketch an operational semantics. The basic idea is that we can use rewrite rules to simplify an  $\mathbf{R}$ -expr until it is either a finite materialized relation—a list of tuples—or has some other convenient form. All of our rewrite rules are semantics-preserving, but some may be more helpful than others. For some  $\mathbf{R}$ -exprs that involve infinite bag relations, there may be no way to eliminate all built-in constraints or aggregation operations. The reduced form then includes delayed constraints (just as in constraint logic programming) or delayed aggregators. Even so, conjoining this reduced form with a query can permit further simplification; therefore, some queries may still yield simple answers.

#### 3.1 Finite Materialized Relations

We may express the finite bag relation shown at left by a simple **sum-of-products**  $\mathbf{R}$ -expr, shown at the right. In this example, the ground values being related are integers.

$$\begin{array}{ccc}
 g = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 \\ 1 & 1 \\ 2 & 6 \\ 2 & 7 \\ 2 & 7 \\ 5 & 7 \end{bmatrix} & \xrightarrow{\text{to } \mathbf{R}\text{-expr}} & \mathbf{R}_g = \begin{array}{l} ( (\mathbf{x}_1 = 1) * (\mathbf{x}_2 = 1) \\ + (\mathbf{x}_1 = 2) * (\mathbf{x}_2 = 6) \\ + (\mathbf{x}_1 = 2) * (\mathbf{x}_2 = 7) \\ + (\mathbf{x}_1 = 2) * (\mathbf{x}_2 = 7) \\ + (\mathbf{x}_1 = 5) * (\mathbf{x}_2 = 7) ) \end{array} \quad (1)
 \end{array}$$

We see that each individual row of the table (tuple) translates into a product (\*) of several (**Variable** = value) expressions, where the variable is the column's name and the value is the cell in the table. To encode multiple rows, the **R**-expr simply adds the **R**-exprs for the individual rows. When evaluated in a given environment, the **R**-expr is a sum of products of multiplicities. But abstracting over possible environments, it represents a union of Cartesian products of 1-tuples, yielding a bag relation.

We may use this **R**-expr as the basis of a new user-defined **R**-expr type  $g$  (case 9 of §2) by taking its definition  $R_g$  to be this **R**-expr. Our **R**-exprs can now include constraints such as  $g(J,K)$  or  $g(J,7)$ . When adding a new case in the denotational semantics in this way, we always match it in the operational semantics by introducing a corresponding rewrite rule  $g(X_1, X_2) \rightarrow R_g$ .

A sum-of-products **R**-expr simply enumerates the tuples of a bag relation, precisely as a boolean expression in disjunctive normal form (that is, an “or-of-ands” expression) enumerates the satisfying assignments. Just as in the boolean case, a disjunct does not have to constrain every variable: it may have “don't care” elements. For example,  $(J=1)*(K=1) + (J=2)$  describes an infinite relation because the second disjunct  $J=2$  is true for any value of  $K$ .

### 3.2 Equality Constraints and Multiplicity Arithmetic

We may wish to query whether the relation  $g$  in the above section relates 2 to 7. Indeed it does—and twice. We may discover this by considering  $g(2,7)$ , which rewrites immediately via substitution to an **R**-expr that has no variables at all  $((2=1)*(7=1) + \dots)$ , and finding that this further reduces to the multiplicity 2.

How does this reduction work? First, we need to know how to evaluate the equality constraints: we need to rewrite  $2=1 \rightarrow 0$  but  $2=2 \rightarrow 1$ . The necessary rewrite rules are special cases of the following structural unification rules:

$$\begin{aligned}
 (f(U_1, \dots, U_n) = g(V_1, \dots, V_m)) &\rightarrow 0 \text{ if } f, g \in \mathcal{F} \text{ and } (f, n) \neq (g, m) &> \text{functor clash} \\
 (f(U_1, \dots, V_n) = f(V_1, \dots, V_n)) &\rightarrow (U_1 = V_1) * \dots * (U_n = V_n) \text{ if } f \in \mathcal{F} \text{ and } n \geq 0 \\
 (T = X) &\rightarrow (X = T) \text{ if } X \in \mathcal{V} \text{ and } T \in \mathcal{T} &> \text{put var on left to match rules below} \\
 (X = X) &\rightarrow 1 \text{ if } X \in \mathcal{V} &> \text{because true for every value of } X \\
 (X = T) &\rightarrow 0 \text{ if } X \in \mathcal{V} \text{ and } T \in \mathcal{T} \text{ and } X \in \text{vars}(T) &> \text{not true for any } X \text{ (occurs check)}
 \end{aligned}$$

We now have an arithmetic expression, which we can simplify to the multiplicity 2 via rewrites that implement basic arithmetic on multiplicities  $\mathcal{M}$ :

$$M + N \rightarrow L \text{ if } M, N \in \mathcal{M} \text{ and } M+N=L; \quad M * N \rightarrow L \text{ if } M, N \in \mathcal{M} \text{ and } M*N=L$$

Above, we relied on the definition of the new relation type  $g$ , which allowed us to request a specialization of  $R_g$ . Do we need to make such a definition in order to query a given bag relation? No: we may do so by conjoining the **R**-expr with additional constraints. For example, to get the multiplicity of the pair (2, 7) in  $R_g$ , we may write  $R_g*(X_1=2)*(X_2=7)$ . This filters the original relation to just the pairs that match (2, 7), and simplifies to  $2*(X_1=2)*(X_2=7)$ . To accomplish this simplification, we need to use the following crucial rewrite:

$$(X=T)*R \rightarrow (X=T)*R\{X \mapsto T\} \text{ if } X \in \mathcal{V} \text{ and } T \in \mathcal{T} \quad > \text{equality propagation}$$

As a more interesting example, the reader may consider simplifying the query  $R_g*\text{lessthan}(X_2, 7)$ , which uses a built-in inequality constraint (see §3.4).

### 3.3 Joining Relations

Analogous to eq. (1), we define a second tabular relation  $f$  with a rewrite rule  $f(X_1, X_2) \rightarrow R_f$ .

$$f = \begin{bmatrix} X_1 & X_2 \\ 1 & 2 \\ 3 & 4 \end{bmatrix} \xrightarrow{\text{to } \mathbf{R}\text{-expr}} R_f = \begin{matrix} (X_1 = 1) * (X_2 = 2) \\ + (X_1 = 3) * (X_2 = 4) \end{matrix} \quad (2)$$

We can now consider simplifying the  $\mathbf{R}$ -expr  $f(I, J) * g(J, K)$ , which the reader may recognize as an equijoin on  $f$ 's second column and  $g$ 's first column.<sup>4</sup> Notice that the  $\mathbf{R}$ -expr has *renamed* these columns to its own free variables ( $I, J, K$ ). Reusing the variable  $J$  in each factor is what gives rise to the join on the relevant column. (Compare  $f(I, J) * g(J', K)$ , which does not share the variable and so gives the Cartesian product of the  $f$  and  $g$  relations.)

We can “materialize” the equijoin by reducing it to a sum-of-products form as before, if we wish:  $(I=1)*(J=2)*(K=6) + 2*(I=1)*(J=2)*(K=7)$ .

To carry out such simplifications, we use the fact that multiplicities form a commutative semiring under  $+$  and  $*$ . Since any  $\mathbf{R}$ -expr evaluates to a multiplicity, these rewrites can be used to rearrange unions and joins of  $\mathbf{R}$ -exprs  $Q, R, S \in \mathcal{R}$ :

$1 * R \leftrightarrow R$	$\triangleright$ multiplicative identity
$0 * R \leftrightarrow 0$	$\triangleright$ multiplicative annihilation
$0 + R \leftrightarrow R$	$\triangleright$ additive identity
$\infty * R \rightarrow \infty$ if $R \in \mathcal{M}$ and $R > 0$	$\triangleright$ absorbing element
$\infty + R \rightarrow \infty$	$\triangleright$ absorbing element
$R + S \leftrightarrow S + R;$	
$R * S \leftrightarrow S * R$	$\triangleright$ commutativity
$Q + (R + S) \leftrightarrow (Q + R) + S;$	$\triangleright$ associativity
$Q * (R + S) \leftrightarrow Q * R + Q * S$	$\triangleright$ distributivity
$R * M \rightarrow R + (R * N)$ if $M, N \in \mathcal{M}$ and $(1+N \rightarrow M)$	$\triangleright$ implicitly does $M \rightarrow 1+N$
$R \leftrightarrow R * R$ if $R$ is a constraint	$\triangleright$ as defined in §2

We can apply some of these rules to simplify our example as follows:

$$\begin{aligned}
& f(I, J) * g(J, K) \\
& \rightarrow ((I=1)*(J=2) + (I=3)*(J=4)) * g(J, K) && \triangleright \text{eq. (2)} \\
& \rightarrow (I=1)*(J=2)*g(J, K) + (I=3)*(J=4)*g(J, K) && \triangleright \text{distributivity} \\
& \rightarrow (I=1)*(J=2)*g(2, K) + (I=3)*(J=4)*g(4, K) && \triangleright \text{equality propagation} \\
& \rightarrow (I=1)*(J=2)*((K=6)+(K=7)*2) + (I=3)*(J=4)*0 && \triangleright \text{via eq. (1)} \\
& \rightarrow (I=1)*(J=2)*((K=6)+(K=7)*2) && \triangleright \text{annihilation} \\
& \rightarrow (I=1)*(J=2)*(K=6) + (I=1)*(J=2)*(K=7)*2 && \triangleright \text{distributivity}
\end{aligned}$$

Notice that the factored intermediate form  $(I=1)*(J=2)*((K=6)*1 + (K=7)*2)$  is more compact than the final sum of products, and may be preferable in some settings. In fact, it is an example of a **trie** representation of a bag relation. Like the root node of a trie, the expression partitions the bag of  $(I, J, K)$  tuples into disjuncts according to the value of  $I$ . Each possible value of  $I$  (in this case only  $I=1$ ) is multiplied by a trie representation of the bag of  $(J, K)$  tuples that can co-occur

<sup>4</sup> This notation may be familiar from Datalog, except that we are writing the conjunction operation as  $*$  rather than with a comma, to emphasize the fact that we are multiplying multiplicities rather than merely conjoining booleans.



with this  $I$ . That representation is a sum over possible values of  $J$  (in this case only  $J=2$ ), which recurses again to a sum over possible values of  $K$  ( $K=6$  and  $K=7$ ). Finally, the multiplicities 1 and 2 are found at the leaves. A trie-shaped  $R$ -expr generally has a smaller branching factor than a sum-of-products  $R$ -expr. As a result, it is comparatively fast to query it for all tuples that strongly restrict  $I$  or  $(I, J)$  or  $(I, J, K)$ , by narrowing down to the matching summand(s) at each node. For example, multiplying our example trie by the query  $I=5$  gives an  $R$ -expr that can be immediately simplified to  $\emptyset$ , as the single disjunct (for  $I=1$ ) does not match.

That example query also provides an occasion for a larger point. This trie simplification has the form  $(I=5)*(I=1)*R$ , an expression that in general may be simplified to  $\emptyset$  on the basis of the first two factors, without spending any work on simplifying the possibly large expression  $R$ . This is an example of **short-circuiting** evaluation—the same logic that allows a SAT solver or Prolog solver to backtrack immediately upon detecting a contradiction.

### 3.4 Rewrite Rules for Built-In Constraints

Built-in constraints are an important ingredient in constructing infinite relations. While they are not the only method,<sup>5</sup> they have the advantage that libraries of built-in constraints such as  $\text{plus}(I, J, K)$  (case 2 of §2) usually come with rewrite rules for reasoning about these constraints [6]. Some of the rewrite rules invoke opaque procedural code.

Recall that the arguments to a  $\text{plus}$  constraint are terms, typically either variables or numeric constants. Not all  $\text{plus}$  constraints can be rewritten, but a library should provide at least the following cases:

$$\begin{aligned} \text{plus}(I, J, K) &\rightarrow \underbrace{I(I + J = K)}_{\in \{0,1\}} \text{ if } I, J, K \in \mathbb{R} \\ \text{plus}(I, J, X) &\rightarrow (X = I + J) \text{ if } I, J \in \mathbb{R} \text{ and } X \in \mathcal{V} \\ \text{plus}(I, X, K) &\rightarrow (X = K - I) \text{ if } I, K \in \mathbb{R} \text{ and } X \in \mathcal{V} \\ \text{plus}(X, J, K) &\rightarrow (X = \underbrace{K - J}_{\in \mathbb{R}}) \text{ if } J, K \in \mathbb{R} \text{ and } X \in \mathcal{V} \end{aligned}$$

The  $R$ -expr  $R = \text{proj}(J, \text{plus}(I, 3, J) * \text{plus}(J, 4, K))$  represents the infinite set of  $(I, K)$  pairs such that  $K = (I + 3) + 4$  arithmetically. (The intermediate temporary variable  $J$  is projected out.) The rewrite rules already presented (plus a rewrite rule from §3.5 below to eliminate  $\text{proj}$ ) suffice to obtain a satisfactory answer to the query  $I=2$  or  $K=9$ , by reducing either  $(I=2)*R$  or  $R*(K=9)$  to  $(I=2)*(K=9)$ .

On the other hand, if we wish to reduce  $R$  itself, the above rules do not apply. In the jargon, the two  $\text{plus}$  constraints within  $R$  remain as **delayed constraints**, which cannot do any work until more of their variable arguments are replaced by constants (e.g., due to equality propagation from a query, as above).

We can do better in this case with a library of additional rewrite rules that implement standard axioms of arithmetic [6], in particular the associative law. With these,  $R$  reduces to  $\text{plus}(I, 7, K)$ , which is a simpler description of this infinite relation. Such rewrite rules are known as **constraint propagators**. Other

<sup>5</sup> Others are structural equality constraints and recursive user-defined constraints.

useful examples concerning `plus` include  $\text{plus}(\emptyset, J, K) \rightarrow K=J$  and  $\text{plus}(I, J, J) \rightarrow (I=\emptyset)$ , since unlike the rules at the start of this section, they can make progress even on a single `plus` constraint whose arguments include more than one variable. Similarly, some useful constraint propagators for the `lessthan` relation include  $\text{lessthan}(J, J) \rightarrow \emptyset$ ; the transitivity rule  $\text{lessthan}(I, J) * \text{lessthan}(J, K) \rightarrow \text{lessthan}(I, J) * \text{lessthan}(J, K) * \text{lessthan}(I, K)$ ; and  $\text{lessthan}(\emptyset, I) * \text{plus}(I, J, K) \rightarrow \text{lessthan}(\emptyset, I) * \text{plus}(I, J, K) * \text{lessthan}(J, K)$ . The integer domain can be **split** by rules such as  $\text{int}(I) \rightarrow \text{int}(I) * (\text{lessthan}(\emptyset, I) + \text{lessthan}(I, 1))$  in order to allow case analysis of, for example,  $\text{int}(I) * \text{myconstraint}(I)$ . All of these rules apply even if their arguments are variables, so they can apply early in a reduction before other rewrites have determined the values of those variables. Indeed, they can sometimes short-circuit the work of determining those values.

Like all rewrites, built-in rewrites  $R \rightarrow S$  must not change the denotation of  $R$ : they ensure  $\llbracket R \rrbracket_E = \llbracket S \rrbracket_E$  for all  $E$ . For example,  $\text{lessthan}(X, Y) * \text{lessthan}(Y, X) \rightarrow^* \emptyset$  is semantics-preserving because both forms denote the empty bag relation.

### 3.5 Projection

Projection is implemented using the following rewrite rules. The first two rules make it possible to push the  $\text{proj}(X, \dots)$  operator down through the sums and products of  $R$ , so that it applies to smaller subexpressions that mention  $X$ :

$$\begin{aligned} \text{proj}(X, R+S) &\leftrightarrow \text{proj}(X, R) + \text{proj}(X, S) && \triangleright \text{distributivity over } + \\ \text{proj}(X, R*S) &\leftrightarrow R * \text{proj}(X, S) \quad \text{if } X \notin \text{vars}(R) && \triangleright \text{see also the } R * \infty \text{ rule below} \end{aligned}$$

Using the following rewrite rules, we can then eliminate the projection operator from smaller expressions whose projection is easy to compute. (In other cases, it must remain as a delayed operator.) How are these rules justified? Observe that  $\text{proj}(X, R)$  in an environment  $E$  denotes the number of  $X$  values that are consistent with  $E$ 's binding of  $R$ 's other free variables. Thus, we may safely rewrite it as another expression that always achieves the same denotation.

$$\begin{aligned} \text{proj}(X, (X=T)) &\rightarrow 1 \quad \text{if } T \in \mathcal{T} \text{ and } X \notin \text{vars}(T) && \triangleright \text{occurs check} \\ \text{proj}(A, (A=\text{sum}(X, R))) &\rightarrow 1 \quad \text{if } A \notin \text{vars}(R) && \triangleright \text{cardinality of an aggregated variable} \\ \text{proj}(X, R) &\rightarrow R * \infty \quad \text{if } X \notin \text{vars}(R) && \triangleright \text{cardinality of an unconstrained variable} \\ \text{proj}(X, \text{bool}(X)) &\rightarrow 2 && \triangleright \text{cardinality of a variable given a certain unary constraint} \\ \text{proj}(X, \text{int}(X)) &\rightarrow \infty && \triangleright \text{cardinality of a variable given a certain unary constraint} \\ \text{proj}(X, \text{proj}(Y, \text{nand}(X, Y))) &\rightarrow 3 && \triangleright \text{card. of a pair given a certain binary constraint} \end{aligned}$$

As a simple example, let us project column  $K$  out of the table  $g(J, K)$  from eq. (1).

$$\begin{aligned} \text{proj}(K, ((J=1) * (K=1) &+ (J=2) * (K=6) &+ (J=2) * (K=7) &+ (J=2) * (K=7) &+ (J=5) * (K=7))) \\ \rightarrow \text{proj}(K, g(J, K)) &\rightarrow ( (J=1) * \text{proj}(K, (K=1)) &+ (J=2) * \text{proj}(K, (K=6)) &+ (J=2) * \text{proj}(K, (K=7)) &+ (J=2) * \text{proj}(K, (K=7)) &+ (J=5) * \text{proj}(K, (K=7))) \\ &\rightarrow (J=1) + (J=2)*3 + (J=5) \end{aligned}$$

When multiple projection operators are used, we may push them down independently of each other, since they commute:

$$\text{proj}(X, \text{proj}(Y, R)) \rightarrow \text{proj}(Y, \text{proj}(X, R))$$

### 3.6 Aggregation

The simple count aggregator from §2 is implemented with the following rewrite rules, which resemble those for proj:

```

M=count(R+S) → proj(L, (L=count(R)) * proj(N, (N=count(S)) * plus(L,N,M)))
M=count(R*S) → proj(L, (L=count(R)) * proj(N, (N=count(S)) * times(L,N,M))) if
    vars(R) ∩ vars(S) = ∅
M=count(N) → (M=N) if N ∈ M

```

In the first two rules,  $L$  and  $N$  are new bound variables introduced by the right-hand side of the rule. (When the rewrite rule is applied, they will—as is standard—potentially be renamed to avoid capturing free variables in the arguments to the left-hand side.) They serve as temporary registers. The third rule covers the base case where the expression has been reduced to a constant multiplicity: e.g.,

```

M=count(5=5) → M=count(1) → M=1
M=count(plus(I,J)*(I=5)) → M=count((I=0)*(I=5)) →* M=count(0) → M=0

```

The following rewrite rules implement `sum`. (The rules for other aggregation operators are isomorphic.) The usual strategy is to rewrite  $A=\text{sum}(X,R)$  as a chain of `plus` constraints that maintain a running total. The following rules handle cases where  $R$  is expressed as a union of 0, 1, or 2 bag relations, respectively. (A larger union can be handled as a union of 2 relations, e.g.,  $(Q+R)+S$ .)

```

A=sum(X, 0) → (A=idsum)
A=sum(X, (X=T)) → (A=T) if T ∈ T and X ∉ vars(T) ▷occurs check
A=sum(X, R+S) → proj(B, (B=sum(X,R)) * proj(C, (C=sum(X,S)) * plus(B,C,A)))

```

The second rule above handles only one of the base cases of 1 bag relation. We must add rules to cover other base cases, such as these:<sup>6</sup>

```

A=sum(X, (X=sum(Y,R))) → (A=sum(Y,R)) if X ∉ vars(R)
A=sum(X, (X=min(Y,R))) → (A=min(Y,R)) if X ∉ vars(R)

```

Most important of all is this case, which is analogous to the second rule of §3.5 and is needed to aggregate over sum-of-products constructions:

```

A=sum(X, R*S) ↔ sum_copies(R,B,A)*(B=sum(X,S)) if X ∉ vars(R)

```

Here, `sum_copies(M,B,A)` for  $M \in \mathcal{M}$  constrains  $A$  to be the aggregation of  $M$  copies of the aggregated value  $B$ . The challenge is that in the general case we actually have `sum_copies(R,B,A)`, so the multiplicity  $M$  may vary with the free variables of  $R$ . The desired denotational semantics are

$$\begin{aligned}
\llbracket \text{sum\_copies}(R,B,A) \rrbracket_E &= \text{if } \llbracket R \rrbracket_E \cdot \llbracket B \rrbracket_E = \llbracket A \rrbracket_E \text{ then } 1 \text{ else } 0 \\
\llbracket \text{min\_copies}(R,B,A) \rrbracket_E &= \text{if } (\llbracket R \rrbracket_E = 0 \text{ and } \llbracket A \rrbracket_E = \text{id}_{\min}) \\
&\quad \text{or } (\llbracket R \rrbracket_E > 0 \text{ and } \llbracket A \rrbracket_E = \llbracket B \rrbracket_E) \text{ then } 1 \text{ else } 0
\end{aligned}$$

where  $R \in \mathcal{R}$  and  $B, A \in \mathcal{T}$

<sup>6</sup> As in §3.5, we could also include special rewrites for certain aggregations that have a known closed-form result, such as certain series sums.

where we also show the interesting case of  $\text{min\_copies}(\mathbf{M}, \mathbf{B}, \mathbf{A})$ , which is needed to help define the min aggregator. We can implement these by the rewrite rules

$$\begin{aligned} \text{sum\_copies}(\mathbf{R}, \mathbf{B}, \mathbf{A}) &\rightarrow \text{proj}(\mathbf{M}, (\mathbf{M} = \text{count}(\mathbf{R})) * \text{times}(\mathbf{M}, \mathbf{B}, \mathbf{A})) && \triangleright \text{assumes } \text{id}_{\text{sum}} = 0 \\ \text{min\_copies}(\mathbf{R}, \mathbf{B}, \mathbf{A}) &\rightarrow \text{proj}(\mathbf{M}, (\mathbf{M} = \text{count}(\mathbf{R})) * ((\mathbf{M} = 0) * (\mathbf{A} = \text{id}_{\text{min}}) + \text{lessthan}(\emptyset, \mathbf{M}) * (\mathbf{A} = \mathbf{B}))) \end{aligned}$$

Identities concerning aggregation yield additional rewrite rules. For example, since multiplication distributes over  $\sum$ , summations can be merged and factored via  $(\mathbf{B} = \text{sum}(\mathbf{I}, \mathbf{R})) * (\mathbf{C} = \text{sum}(\mathbf{J}, \mathbf{S})) * \text{times}(\mathbf{B}, \mathbf{C}, \mathbf{A}) \leftrightarrow \mathbf{A} = \text{sum}(\mathbf{K}, \mathbf{R} * \mathbf{S} * \text{times}(\mathbf{I}, \mathbf{J}, \mathbf{K}))$  provided that  $\mathbf{I} \in \text{vars}(\mathbf{R})$ ,  $\mathbf{J} \in \text{vars}(\mathbf{S})$ ,  $\mathbf{K} \notin \text{vars}(\mathbf{R} * \mathbf{S})$ , and  $\text{vars}(\mathbf{R}) \cap \text{vars}(\mathbf{S}) = \emptyset$ . Other distributive properties yield more rules of this sort. Moreover, projection and aggregation operators commute if they are over different free variables.

To conclude this section, we now attempt aggregation over infinite streams. We wish to evaluate  $\mathbf{A} = \text{exists}(\mathbf{B}, \text{proj}(\mathbf{I}, \text{peano}(\mathbf{I}) * \text{myconstraint}(\mathbf{I}) * (\mathbf{B} = \text{true})))$  to determine whether there exists any Peano numeral that satisfies a given constraint. Here  $\text{exists}$  is the aggregation operator based on the binary or operation.

$\text{peano}(\mathbf{I})$  represents the infinite bag of Peano numerals, once we define a user constraint via the rewrite rule  $\text{peano}(\mathbf{I}) \rightarrow (\mathbf{I} = \text{zero}) + \text{proj}(\mathbf{J}, (\mathbf{I} = \mathbf{s}(\mathbf{J})) * \text{peano}(\mathbf{J}))$ . Rewriting  $\text{peano}(\mathbf{I})$  provides an opportunity to apply the rule again (to  $\text{peano}(\mathbf{J})$ ). After  $k \geq 0$  rewrites we obtain a representation of the original bag that explicitly lists the first  $k$  Peano numerals as well as a continuation that represents all Peano numerals  $\geq k$ :

$$\rightarrow (\mathbf{I} = \text{zero}) + \dots + (\mathbf{I} = \underbrace{\mathbf{s}(\dots \mathbf{s}(\text{zero}) \dots)}_{k-1 \text{ times}}) + \overbrace{\text{proj}(\mathbf{J}, (\mathbf{I} = \underbrace{\mathbf{s}(\dots \mathbf{s}(\mathbf{J}) \dots)}_{k-1 \text{ times}})) * \text{peano}(\mathbf{J})}^{\text{continuation}}$$

Rewriting the  $\text{exists}$  query over this  $(k + 1)$ -way union results in a chain of  $k$  or constraints. If one of the first  $k$  Peano numerals in fact satisfies  $\text{myconstraint}$ , then we can “short-circuit” the infinite regress and determine our answer without further expanding the continuation, thanks to the useful rewrite  $\text{or}(\text{true}, \mathbf{C}, \mathbf{A}) \rightarrow (\mathbf{A} = \text{true})$ , which can apply even while  $\mathbf{C}$  remains unknown.

In general, one can efficiently aggregate a function of the Peano numerals by alternating between expanding  $\text{peano}$  to draw the next numeral from the iterator, and rewriting the aggregating  $\mathbf{R}$ -expr to aggregate the value of the function at that numeral into a running “total.” If the running total ever reaches an absorbing element  $a$  of the aggregator’s binary operation—such as  $\text{true}$  for the or operation—then one will be able to simplify the expression to  $\mathbf{A} = a$  and terminate. We leave the details as an exercise.

## 4 Translation of Dynabases to $\mathbf{R}$ -exprs

The translation of a Dyna program to a single recursive  $\mathbf{R}$ -expr can be performed mechanically. We will illustrate the basic construction on the small contrived example below. We will focus on the first three rules, which define  $\mathbf{f}$  in terms of  $\mathbf{g}$ . The final rule, which defines  $\mathbf{g}$ , will allow us to take note of a few subtleties.

```

12 | f(X) += X * X.
13 | f(4) += 3.
14 | f(X) += g(X, Y).
15 | g(4 * C, Y) += C - 1 for Y > 99.

```

Recall that a Dyna program represents a set of key-value pairs. `is(Key,Val)` is the conventional name for the key-value relation. The above program translates into the following user-defined constraint, which recursively invokes itself when the value of one key is defined in terms of the values of other keys.

```
is(Key,Val) → (Val=sum(Result,                                ▷sum represents the += aggregator
  proj(X, (Key=f(X))*times(X,X,Result) )                      ▷f(X) += X*X.
+      (Key=f(4))*(Result=3)                                  ▷f(4) += 3.
+proj(X, (Key=f(X))*proj(Y,is(g(X,Y),Result))) )              ▷f(X) += g(X,Y).
+proj(C, proj(Y, proj(Temp,(Key=g(Temp,Y))*times(4,C,Temp))    ▷g(4*C,Y) +=
  *minus(C,1,Result)*lessthan(99,Y)                            ▷... C-1 for Y > 99.
) * notnull(Val)        ▷notnull discards any pair whose Val aggregated nothing
```

Each of the 4 Dyna rules translates into an **R**-expr (as indicated by the comments above) that describes a bag of `(Key,Result)` pairs of ground terms. In each pair, `Result` represents a possible ground value of the rule's body (the Dyna expression to the right of the aggregator `+=`), and `Key` represents the corresponding grounding of the rule head (the term to the left of the aggregator), which will receive `Result` as an aggregand. Note that the same `(Key,Result)` pair may appear multiple times. Within each rule's **R**-expr, we project out the variables such as `X` and `Y` that appear locally within the rule, so that the **R**-expr's free variables are only `Key` and `Result`.<sup>7</sup>

Dyna mandates in-place evaluation of Dyna expressions that have values [4]. For each such expression, we create a new local variable to bind to the result. Above, the expressions such as `X*X`, `g(X,Y)`, `4*C`, and `C-1` were evaluated using `times`, `is`, `times`, and `minus` constraints, respectively, and their results were assigned to new variables `Result`, `Result`, `Temp`, and `Result`. Importantly, `g(X,Y)` refers to a key of the Dyna program itself, so the Dyna program translated into an `is` constraint whose definition recursively invokes `is(g(X,Y),Result)`.

The next step is to take the bag union (+) of these 4 bag relations. This yields the bag of all `(Key,Result)` pairs in the program. Finally, we wrap this combined bag in `Val=sum(Result,...)` to aggregate the `Results` for each `Key` into the `Val` for that key. This yields a set relation with exactly one value for each key. For `Key=f(4)`, for example, the first and second rules will each contribute one `Result`, while the third rule will contribute as many `Results` as the map has keys of the form `g(4,Y)`.<sup>8</sup>

We use `sum` as our aggregation operator because all rules in this program specify the `+=` aggregator. One might expect `sum` to be based on the binary operator that is implemented by the `plus` builtin, as described before, with identity element `idsum = 0`. There is a complication, however: in Dyna, a `Key` that has

<sup>7</sup> It is always legal to project out the local variables at the top level of the rule, e.g., `proj(X,proj(Y,...))` for rule 3. However, we have already seen rewrite rules that can narrow the scope of `proj(Y,...)` to a sub-**R**-expr that contains all the copies of `Y`. Here we display the version after these rewrites have been done.

<sup>8</sup> The reader should be able to see that the third Dyna rule will contribute infinitely many copies of `0`, one for each `Y > 99`. This is an example of multiplicity  $\infty$ . Fortunately, `sum_copies` (invoked when rewriting the `sum` aggregation operator) knows that summing any positive number of `0` terms—even infinitely many—will give `0`.

no **Results** should not in fact be paired with **Val**=0. Rather, this **Key** should not appear as a key of the final relation at all! To achieve this, we augment  $\mathcal{F}$  with a new constant **null** (similar to Prolog’s **no**), which represents “no results.” We define  $\text{id}_{\text{sum}} = \text{null}$ , and we base **sum** on a modified version of **plus** for which **null** is indeed the identity (so the constraints  $\text{plus}(\text{null}, \mathbf{x}, \mathbf{x})$  and  $\text{sum\_copies}(\emptyset, \mathbf{x}, \text{null})$  are both true for all  $\mathbf{x}$ ). All aggregation operators in Dyna make use of **null** in this way. Our  $\text{Val}=\text{sum}(\text{Result}, \dots)$  relation now obtains **null** (rather than  $\emptyset$ ) as the unique **Val** for each **Key** that has no aggregands. As a final step in expressing a Dyna program, we always remove from the bag all  $(\text{Key}, \text{Val})$  pairs for which  $\text{Val}=\text{null}$ , by conjoining with a  $\text{notnull}(\text{Val})$  constraint. This is how we finally obtain the **R**-expr above for  $\text{is}(\text{Key}, \text{Val})$ .

To query the Dyna program for the value of key  $f(4)$ , we can reduce the expression  $\text{is}(f(4), \text{Val})$ , using previously discussed rewrite rules. We can get as far as this before it must carry out its own query  $g(4, Y)$ :

```
proj(C, (C=sum(Result, proj(Y, is(g(4, Y), Result))))
      * plus(19, C, Val) ) * notnull(Val)
```

where the local variable **C** captures the total contribution from rule 3, and 19 is the total contribution of the other rules. To reduce further, we now recursively expand  $\text{is}(g(4, Y), \text{Result})$  and ultimately reduce it to  $\text{Result}=\emptyset$  (meaning that  $g(4, Y)$  turns out to have value  $\emptyset$  for all ground terms  $Y$ ).  $\text{proj}(Y, \text{Result}=\emptyset)$  reduces to  $(\text{Result}=\emptyset)*\infty$ —a bag relation with an infinite multiplicity—but then  $\text{C}=\text{sum}(\text{Result}, (\text{Result}=\emptyset)*\infty)$  reduces to  $\text{C}=\emptyset$  (via  $\text{sum\_copies}$ , as footnote 8 noted). The full expression now easily reduces to  $\text{Val}=19$ , the correct answer.

What if the Dyna program has different rules with different aggregators? Then our translation takes the form

```
is(Key, Val) → Val=only(Val1, (Val1=sum(Result, RSum))
                             + (Val1=min(Result, RMin))
                             + (Val1=only(Result, REq))
                             + ... ) * notnull(Val)
```

where **RSum** is the bag union of the translated  $\text{+=}$  rules as in the previous example, **RMin** is the bag union of the translated  $\text{min=}$  rules, **REq** is the bag union of the translated  $\text{=}$  rules, and so on. The new aggregation operator **only** is based on a binary operator that has identity  $\text{id}_{\text{only}} = \text{null}$  and combines any pair of non-null values into error. For each **Key**, therefore, **Val** is bound to the aggregated result **Val1** of the *unique* aggregator whose rules contribute results to that key. If multiple aggregators contribute to the same key, the value is **error**.<sup>9</sup>

A Dyna program may consist of multiple *dynabases* [4]. Each dynabase defines its own key-value mapping, using aggregation over only the rules in that dynabase, which may refer to the values of keys in other dynabases. In this case, instead of defining a single constraint **is** as an **R**-expr, we define a different named constraint for each dynabase as a different **R**-expr, where each of the **R**-exprs may call any of these named constraints.

<sup>9</sup> A Dyna program is also supposed to give an **error** value to the key **a** if the program contains both the rules  $a = 3$  and  $a = 4$ , or the rule  $a = f(X)$  when both  $f(\emptyset)$  and  $f(1)$  have values. So we also used **only** above as the aggregation operator for  $\text{=}$  rules.

## 5 Conclusions and Ongoing Work

We have shown how to algebraically represent the bag-relational semantics of any program written in a declarative logic-based language like Dyna. A query against a program can be evaluated by joining the query to the program and simplifying the resulting algebraic expression with appropriate term rewriting rules. It is congenial that this approach allows evaluation to flexibly make progress, propagate information through the expression, exploit arithmetic identities, and remove irrelevant subexpressions rather than wasting possibly infinite work on them.

In ongoing work, we are considering methods for practical interpretation and compilation of rewrite systems. Our goal is to construct an evaluator that will both perform static analysis and “learn to run fast” [9] by constructing or discovering reusable strategies for reducing expressions of frequently encountered sorts. In the case of cyclic rewrite systems, the system should be able to guess portions of a solution and then verify that the guesses are consistent with the rewrite rules. A forward chaining mechanism can be used to invalidate or correct inconsistent guesses, as is common in Datalog, and this mechanism can also be used for change propagation when the program or the query is externally updated [5].

**Acknowledgements.** This material is based upon work supported by the National Science Foundation under Grant No. 1629564. We thank Scott Smith for helpful comments, and the WRLA program chairs Santiago Escobar and Narciso Martí Oliet for allowing generous time to revise the paper.

## References

1. Ceri, S., Gottlob, G., Tanca, L.: What you always wanted to know about Datalog (and never dared to ask). In: IEEE Transactions on Knowledge and Data Engineering (1989)
2. Clocksin, W.F., Mellish, C.S.: Programming in Prolog. Springer-Verlag (1984)
3. Colmerauer, A., Roussel, P.: The Birth of Prolog, chap. 7. Association for Computing Machinery (1996)
4. Eisner, J., Filardo, N.W.: Dyna: Extending Datalog for modern AI. In: de Moor, O., Gottlob, G., Furche, T., Sellers, A. (eds.) Datalog Reloaded. Lecture Notes in Computer Science, Springer (2011)
5. Filardo, N.W., Eisner, J.: A flexible solver for finite arithmetic circuits. In: Technical Communications of the International Conference on Logic Programming. Leibniz International Proceedings in Informatics (LIPIcs) (2012)
6. Frühwirth, T.: Theory and practice of constraint handling rules. The Journal of Logic Programming **37**(1) (1998)
7. Gallaire, H., Minker, J., Nicolas, J.M.: Logic and databases: A deductive approach. ACM Comput. Surv. **16**(2) (1984)
8. Green, T.J.: Bag semantics. In: LIU, L., ÖZSU, M.T. (eds.) Encyclopedia of Database Systems. Springer (2009)
9. Vieira, T., Francis-Landau, M., Filardo, N.W., Khorasani, F., Eisner, J.: Dyna: Toward a self-optimizing declarative language for machine learning applications. In: Proceedings of the ACM SIGPLAN Workshop on Machine Learning and Programming Languages (MAPL). ACM (2017)