

Complete TRS Specification for Abstract Collection Types

Anthony Hunt

May 29, 2025

Contents

1	Introduction	2
2	High Level Strategy	2
3	Supported Operations	3
4	Rules	5
4.1	Sets	5
4.2	Relations	7
4.3	Bags	8
5	Implementation Representation	8

1 Introduction

This document serves as a living specification of the underlying term rewriting system used in the compiler for a modelling-focused programming language.

2 High Level Strategy

General Strategy A basic strategy to optimize set and relational expressions is:

1. Normalize the expression as a set comprehensions
2. Simplify and reorganize conjuncts of the set comprehension body

Intuition The TRS for this language primarily involves lowering collection data type expressions into pointwise boolean quantifications. Breaking down each operation into set builder notation enables a few key actions:

- Quantifications over sets ($\{x \cdot G \mid P\}$) are naturally separated into generators (G) and (non-generating) predicates (P). For sets, at least one membership operator per top-level conjunction in G will serve as a concrete element generator in generated code. Then, top level disjunctions will select one membership operation to act as a generator, relegating all others to the predicate level. For example, if the rewrite system observes an intersection of the form $\{x \cdot x \in S \wedge x \in T\}$, the set construction operation must iterate over at least one of S and T . Then, the other will act as a condition to check every iteration (becoming $\{x \cdot x \in S \mid x \in T\}$).
- By definition of generators in quantification notation, operations in G must be statements of the form $x \in S$, where x is used in the “element” portion of the set construction. Statements like $x \notin T$ or checking a property $p(x)$ must act like conditions since they do not produce any iterable elements.
- Any boolean expression for conditions may be rewritten as a combination of \neg , \vee , and \wedge expressions. Therefore, by converting all set notation down into boolean notation and then generating code based on set constructor booleans, we can accommodate any form of predicate function.

Granular Strategy (Sets)

3 Supported Operations

Table 1: Summary table: a few operators on sets and relations.

Sets		Relations	
Syntax	Label/Description	Syntax	Label/Description
$set(T)$	Unordered, unique collection	$S \mapsto T$	Partial function
$S \leftrightarrow T$	Relation, $set(S \times T)$	$S \mapsto T$	Total injection
\emptyset	Empty set	$a \mapsto b$	Pair (relational element)
$\{a, b, \dots\}$	Set enumeration	$dom(S)$	Domain
$\{x \cdot x \in S \mid P\}$	Set comprehension	$ran(S)$	Range
$S \cup T$	Union	$R[S]$	Relational image
$S \cap T$	Intersection	$R \Leftarrow Q$	Relational overriding
$S \setminus T$	Difference	$R \circ Q$	Relational composition
$S \times T$	Cartesian Product	$S \triangleleft R$	Domain restriction
$S \subseteq T$	Subset	R^{-1}	Relational inverse

Table 2: Collection of operators on set data types.

Name	Definition
Empty Set	Creates a set with no elements.
Set Enumeration	Literal collection of elements to create a set.
Set Membership	The term $x \in S$ is True if x can be found somewhere in S .
Union	$S \cup T = \{x \cdot x \in S \vee x \in T\}$
Intersection	$S \cap T = \{x \cdot x \in S \wedge x \in T\}$
Difference	$S \setminus T = \{x \cdot x \in S \mid x \notin T\}$
Cartesian Product	$S \times T = \{x \mapsto y \cdot x \in S \wedge y \in T\}$
Powerset	$\mathbb{P}(S) = \{s \cdot s \subseteq S\}$
Magnitude	$\#S = \sum_{x \in S} 1$
Subset	$S \subseteq T \equiv \forall x \in S : x \in T$
Strict Subset	$S \subset T \equiv S \subseteq T \wedge S \neq T$
Superset	$S \supseteq T \equiv \forall x \in T : x \in S$
Strict Superset	$S \supset T \equiv S \supseteq T \wedge S \neq T$
Set Mapping	$f * S = \{f(x) \cdot x \in S\}$
Set Filter	$p \triangleleft S = \{x \cdot x \in S \mid p(x)\}$
Set Quantification (Folding)	$\oplus x \cdot x \in S \mid P$
Cardinality	$card(S) = \sum 1 \cdot x \in S$

Table 3: Collection of operators on bag/multiset data types.

Name	Definition
Empty Set	Creates a set with no elements.
Bag Enumeration	Literal collection of elements to create a set (for now, stored as a tuple of elements and number of occurrences).
Bag Membership	The term $x \in S$ is True if S contains one or more occurrences of x .
Union	$S \cup T = \{ (x, a + b) \cdot (x, a) \in S \wedge (x, b) \in T \mid a, b \geq 0 \}$
Intersection	$S \cap T = \{ (x, \min(a, b)) \cdot (x, a) \in S \wedge (x, b) \in T \mid a, b \geq 0 \}$
Difference	$S - T = \{ (x, a - b) \cdot (x, a) \in S \wedge (x, b) \in T \mid a, b \geq 0 \wedge a - b > 0 \}$
Bag Mapping	$f * S = \{ (f(x), r) \cdot (x, r) \in S \}$
Bag Filter	$p \triangleleft S = \{ (x, r) \cdot (x, r) \in S \mid p(x) \}$
Size	$size(S) = \sum r \cdot (x, r) \in S$
Zero Occurrences	$(x, 0) \in S \implies x \notin S$

Table 4: Collection of operators on sequence data types.

Name	Definition
Empty List	Creates a list with no elements.
List Enumeration	Literal collection of elements to create a list.
Construction	Alternative form of List Enumeration.
List Membership	The term $x \text{ in } S$ is True if x can be found somewhere in S .
Append	$[s_1, s_2, \dots, s_n] + t = [s_1, s_2, \dots, s_n, t]$
Concatenate	$[s_1, \dots, s_n] ++ [t_1, \dots, t_n] = [s_1, \dots, s_n, t_1, \dots, t_n]$
Length	$\#S = \sum 1 \cdot x \text{ in } S$
List Mapping	$f * S = [f(x) \cdot x \text{ in } S]$
List Filter	$p \triangleleft S = [f(x) \cdot x \text{ in } S \mid p(x)]$
Associative Reduction	$\oplus/[s_1, s_2, \dots, s_n] = s_1 \oplus s_2 \oplus \dots \oplus s_n$
Right Fold	$\text{foldr}(f, e, [s_1, s_2, \dots, s_n]) = f(s_1, f(s_2, f(\dots, f(s_n, e))))$
Left Fold	$\text{foldl}(f, e, [s_1, s_2, \dots, s_n]) = f(f(f(f(e, s_1), s_2), \dots), s_n)$

Table 5: Collection of operators on relation data types.

Name	Definition
Empty Relation	Creates a relation with no elements.
Relation Enumeration	Literal collection of elements to create a relation.
Identity	$id(S) = \{x \mapsto x \cdot x \in S\}$
Domain	$dom(R) = \{x \cdot x \mapsto y \in R\}$
Range	$ran(R) = \{y \cdot x \mapsto y \in R\}$
Relational Image	$R[S] = \{y \cdot x \mapsto y \in R \mid x \in S\}$
Overriding	$R \triangleleft Q = Q \cup (dom(Q) \triangleleft R)$
(Forward) Composition	$Q \circ R = \{x \mapsto z \cdot x \mapsto y \in R \wedge y \mapsto z \in Q\}$
Inverse	$R^{-1} = \{y \mapsto x \cdot x \mapsto y \in R\}$
Domain Restriction	$S \triangleleft R = \{x \mapsto y \cdot x \mapsto y \in R \mid x \in S\}$
Domain Subtraction	$S \triangleleft R = \{x \mapsto y \cdot x \mapsto y \in R \mid x \notin S\}$
Range Restriction	$R \triangleright S = \{x \mapsto y \cdot x \mapsto y \in R \mid y \in S\}$
Range Subtraction	$R \triangleright S = \{x \mapsto y \cdot x \mapsto y \in R \mid y \notin S\}$

4 Rules

Below is a list of rewrite rules for key abstract data types.

4.1 Sets

$S \rightarrow \{x \cdot x \in S\}$	(Set Construction)
$x \in e \rightarrow x = e$	(Singleton Membership)
$f(x) \cdot x \in \{g(y) \cdot y \in S \mid P\} \rightarrow f(g(y)) \cdot y \in S \mid P$	(Membership Collapse)
$x \cdot x \in S \wedge p(x) \rightarrow x \cdot x \in S \mid p(x)$	(Predicate Promotion)
<hr/>	
$x \in S \cup T \rightarrow x \in S \vee x \in T$	(Union)
$x \in S \cap T \rightarrow x \in S \wedge x \in T$	(Intersection)
$x \in S \setminus T \rightarrow x \in S \wedge x \notin T$	(Difference)
$card(S) \rightarrow \sum x \in S \cdot 1$	(Cardinality)
<hr/>	
$filter(p, S) \oplus filter(q, T) \rightarrow \{x \cdot (x \in S \wedge p(x)) \oplus_{bool} (x \in T \wedge q(y))\}$	(Predicate Operations)
$filter(p, S) \oplus filter(p, T) \rightarrow filter(p, S \oplus T)$	(Filter Over Operators)
$map(f, S) \oplus map(g, T) \rightarrow \{z \cdot (x \in S \wedge f(x) = z) \oplus_{bool} (x \in T \wedge g(y) = z)\}$	(Mapping Operations)

$$map(f, S) \oplus map(f, T) \rightarrow map(f, S \oplus T) \quad (\text{Mapping Over Operators})$$

$$map(f, S) \rightarrow \{ f(x) \cdot x \in S \} \quad (\text{Map})$$

$$filter(p, S) \rightarrow \{ x \cdot x \in S \mid p(x) \} \quad (\text{Filter})$$

$$S \cup T \text{ where } S \subseteq T \rightarrow T \quad (\text{Union Subset Simplification})$$

$$S \cap T \text{ where } S \subseteq T \rightarrow S \quad (\text{Intersection Subset Simplification})$$

$$S \setminus T \text{ where } S \subseteq T \rightarrow \emptyset \quad (\text{Difference Subset Simplification})$$

$$filter(p, S) \cup filter(q, T) \text{ where } S \subseteq T \rightarrow filter(p \vee q, T) \quad (\text{Predicate Union with Subset})$$

$$filter(p, S) \cap filter(q, T) \text{ where } S \subseteq T \rightarrow filter(p \wedge q, S) \quad (\text{Predicate Intersection with Subset})$$

$$\{ E \cdot x \in S \mid P \} \rightarrow \begin{array}{l} ret := \emptyset \\ \textbf{for } x \in S \textbf{ do} \\ \quad \textbf{if } P \textbf{ then } ret.append(E) \end{array} \quad (\text{Set generation})$$

$$\sum x \in S \mid P \cdot E \rightarrow \begin{array}{l} c := 0 \\ \textbf{for } x \in S \textbf{ do} \\ \quad \textbf{if } P \textbf{ then } c := c + E \end{array} \quad (\text{Summation})$$

$$\begin{array}{l} \textbf{if } P \wedge x \in S \wedge Q \textbf{ then } p \\ \text{where } x \text{ not free in } P \end{array} \rightarrow \begin{array}{l} \textbf{if } P \textbf{ then} \\ \quad \textbf{for } x \in S \textbf{ do} \\ \quad \quad \textbf{if } Q \textbf{ then } p \end{array} \quad (\text{Composed conditional})$$

$$\begin{array}{l} \textbf{if } P \wedge ((x \in S \wedge V) \vee W) \wedge Q \\ \textbf{then } p \\ \text{where } x \text{ not free in } P \end{array} \rightarrow \begin{array}{l} \textbf{if } P \textbf{ then} \\ \quad \textbf{for } x \in S \textbf{ do} \\ \quad \quad \textbf{if } V \wedge Q \textbf{ then } p \\ \quad \quad \textbf{if } W \wedge Q \textbf{ then } p \end{array} \quad (\text{Conjunct conditional})$$

Additional notes and extended context:

General When the intention is clear, operands are written in a concise variable form instead of constructor notation, although the implementation should assume that variables are fully expanded by the (Set Construction) phase. The functions *map* and *filter* correspond to applying a function to all elements of a set and adding a predicate to the set respectively, so $map(f, S) = \{ f(x) \cdot x \in S \}$ and $filter(p, S) = \{ x \cdot x \in S \mid p(x) \}$.

Membership Collapse Assumes context of the initial term results in a valid expression. For example, this rewrite would be reasonable within a quantification statement or set construction. The intention of this rule is to

simplify nested set constructions and should only be used after the (Set Construction) phase.

Predicate Promotion If a predicate within set construction cannot be considered as a loop iterator, it is promoted to a purely conditional role. The meaning of the before and after term does not necessarily change, but it clarifies the difference between candidate generators and regular predicates.

Operations involving \oplus The operator \oplus represents any union, intersection, or difference operation, with \oplus_{bool} representing the corresponding construction operation.

Summation The summation operation may be replaced with any reduction/-folding operation, provided the correct low-level operation. More specific forms may be needed to consider non-associative operations.

Composed Conditional x not free in P means x must already be bound. More clearly, no candidate generator variables that have not been defined should occur in P . We differentiate P from Q so that if-statements may leave behind predicates that are not needed in the nested loop. In the case of top-level AND-separated generators that bind to the same variable, the larger generator should be promoted to a conditional.

4.2 Relations

$R[S] \rightarrow \{x \mapsto y \in R \cdot x \in S \mid y\}$	(Image)
$x \mapsto y \in S \times T \rightarrow x \in S \wedge y \in T$	(Product)
$x \mapsto y \in R^{-1} \rightarrow y \mapsto x \in R$	(Inverse)
$x \mapsto y \in (Q \circ R) \rightarrow x \mapsto z \in Q \wedge z' \mapsto y \in R \wedge z = z'$	(Composition)
$R \triangleleft Q \rightarrow Q \cup (dom(Q) \triangleleft R)$	(Override)
<hr/>	
$dom(R) \rightarrow map(fst, R)$	(Domain)
$ran(R) \rightarrow map(snd, R)$	(Range)
<hr/>	
$S \triangleleft R \rightarrow filter(fst \in S, R)$	(Domain Restriction)
$S \triangleleft R \rightarrow filter(fst \notin S, R)$	(Domain Subtraction)
$S \triangleright R \rightarrow filter(snd \in S, R)$	(Range Restriction)
$S \triangleright R \rightarrow filter(snd \notin S, R)$	(Range Subtraction)

4.3 Bags

Since bags can be interpreted as a set of tuples (*element, repetitions*), all set operations apply, except for the overriding operations below.

$$S \cup T \rightarrow \{ (x, r) \cdot x \in \text{set}(S) \cup \text{set}(T) \mid r = \max(\#(x, S), \#(x, T)) \} \quad (\text{Union})$$

$$S \cap T \rightarrow \{ (x, \min(a, b)) \cdot (x, a) \in S \wedge (x, b) \in T \mid a, b \geq 0 \} \quad (\text{Intersection})$$

$$S + T \rightarrow \{ (x, r) \cdot x \in \text{set}(S) \cup \text{set}(T) \mid r = \#(x, S) + \#(x, T) \} \quad (\text{Sum})$$

$$S - T \rightarrow \{ (x, r) \cdot (x, a) \in S \mid r = a - \#(x, T) \wedge r > 0 \} \quad (\text{Difference})$$

$$\text{size}(S) \rightarrow \sum (x, r) \in S \cdot r \quad (\text{Size})$$

Additional notes and extended context:

The # Operator Defined as the number of occurrences of an element in a bag. If bags are represented by a relation, this corresponds to a direct lookup $\#(x, S) = S[x]$.

Intersection, Difference Since the intersection and difference operators are always decreasing (ex. $S \cap T \subseteq S \wedge S \cap T \subseteq T$ and $S - T \subseteq S$), we can short-circuit operations that would require looping over both sets instead of just S . *But how do we define this short-circuiting behaviour?* Intersections can make use of this property for both operands, but difference will always iterate over the first operand.

Difference $a - b > 0 \implies a > 0$.

Sum, Union The cast to set of $\text{set}(S)$ can be implemented by taking the domain of the bag-representing relations.

5 Implementation Representation

Different implementations of each data type will have varying strengths and weaknesses, not only in theoretical asymptotic time and space, but in concrete real-world tests. Cache usage and additional information through object metadata may prove influential on smaller tests. Since this document is only concerned with the theoretical compiler specification, we analyze the theoretical time and space complexity, then pair gathered examples with a test plan for hardware considerations.

A first approach to tackling these type representations would likely constitute a linked list. The space requirements for enumeration are straightforward, with extra allocations for link pointers. Insertions for unordered collections or append/concat operations are $O(1)$, but $O(n)$ for indexed insertion and union with one element. Lookups for all collections are $O(n)$, but this running time is undesirable for the often-used `in` operator for set-generated code. Since linked

lists naturally enforce element order, this structure may be suitable for fast-changing sequences. Although, a limited-size sequence may be better suited for a contiguous array for $O(1)$ indexing. *TODO: For sequences, we should also see if trees/heaps or bloom filters could provide efficient membership checking. Bloom filters are probabilistic but can determine \neq operations.*

On the other hand, hashmaps with $O(1)$ membership and element lookups are useful for all unordered collections. Relations may need bidirectional hashmaps that can efficiently handle many-to-many relations.

Compressed bitmaps may be used for sets, but require a lot of space for sparse elements.

Bags may be implemented either as a (linked) list, a set of tuples where the number of element occurrences is stored in the second tuple component, or a relation where the number of occurrences is the codomain.