

Complete TRS Specification for Abstract Collection Types

Anthony Hunt

August 6, 2025

Contents

1	Introduction	2
2	High Level Strategy	2
3	Supported Operations	3
4	Relational Subtypes	5
4.1	Properties	6
5	Rules	7
5.1	Builtin Functions	7
5.2	Sets	7
5.2.1	Set Comprehension Construction	7
5.2.2	Disjunctive Normal Form	8
5.2.3	Predicate Simplification 1	8
5.2.4	Generator Selection	9
5.2.5	Predicate Simplification 2	10
5.2.6	Predicate Simplification 3	11
5.2.7	Code Generation 1	11
5.2.8	Code Generation 2	12
5.2.9	Equality Elimination	12
5.3	Relations	13
5.4	Bags	13
6	Implementation Representation	14

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)

Set Comprehension Construction Break down all qualifying sets into comprehension forms, collapsing and simplifying where needed.

DNF Predicates Revise comprehension predicates to top-level disjunctive normal form. Each or-clause should have at least one feasible generator. Each clause should record a list of candidate generators

Predicate Simplification Remove superfluous dummy variables, group or-clauses that use the exact same generator (ex. $\{x \cdot x \in S \wedge x \neq 0 \vee x \in S \wedge x = 0\} \rightarrow \{x \cdot x \in S \wedge (x \neq 0 \vee x = 0)\}$). Clauses should be group-able based on DNF, and generators should be selected and recorded.

Set Code Generation Converts quantifiers into for-loops and if-statements.

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	$S \leftrightarrow T$	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	$S \cup T = \{ \langle x, a+b \rangle \cdot (x, a) \in S \wedge (x, b) \in T \mid a, b \geq 0 \}$
Intersection	$S \cap T = \{ \langle x, \min(a, b) \rangle \cdot (x, a) \in S \wedge (x, b) \in T \mid a, b \geq 0 \}$
Difference	$S - T = \{ \langle x, a-b \rangle \cdot (x, a) \in S \wedge (x, b) \in T \mid a, b \geq 0 \wedge a-b > 0 \}$
Bag Mapping	$f * S = \{ \langle f(x), r \rangle \cdot (x, r) \in S \}$
Bag Filter	$p \triangleleft S = \{ \langle x, r \rangle \cdot (x, r) \in S \mid p(x) \}$
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 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 Relational Subtypes

The notion of a relation is a broad classification of a set of pairs, where the domain, range, and relationship are not bound by any conditions. However, when generating code based on purely mathematical relations, it can be advantageous to use notions of totality, injectivity, surjectivity, etc. In particular, these named relational subtypes can be described with four properties: total on the domain, total on the range, one-to-many, and many-to-one.

The below table converts conventional notation into these four subtypes:

Relational Subtype	Properties			
	T	TR	1-	1-R
Relation				✓
Partial Function			✓	
Partial Injection			✓	✓
Surjective Relation		✓		
		✓		✓
Partial Surjection		✓	✓	
		✓	✓	✓
Total Relation	✓			
	✓			✓
Total Function	✓		✓	
Total Injection	✓		✓	✓
Total Surjective Relation	✓	✓		
	✓	✓		✓
Total Surjection	✓	✓	✓	
Bijection	✓	✓	✓	✓

Table 6: Conventional mathematical notation of relational subtypes decomposed into four properties: domain totality (T), range totality (TR), domain one-to-manyness (1-), and range one-to-manyness (1-R)

4.1 Properties

Properties that can affect code generation may involve:

- New properties after applying operations on sets/relations
- Size of result
- Validity of operations
- Deterministic nature of operations
- Super/subset properties
- Relational identities

Below is a list of the impact of these properties.

Domain Totality:

- $R[S] \neq \emptyset$
- $R(S)$ is always valid
- R^{-1} is TR
- $|R| \geq \text{dom}(R)$

Range Totality:

- R^{-1} is T
- $|R| \geq \text{ran}(R)$

Domain One-to-manyness:

- $R(\{e\})$ is deterministic
- $|R[\{e\}]| = 1$
- $|R[S]| \leq |S|$
- $|R| \leq \text{dom}(R)$

Range One-to-manyness:

- $|R| \leq \text{ran}(R)$

5 Rules

Below is a list of rewrite rules for key abstract data types, syntactic sugar, and some builtin functions. Phases are intended to be executed in order; the post-condition of one phase serves as the pre-condition for the next.

5.1 Syntactic Sugar - Bag Image

$$R[B] \rightsquigarrow R^{-1} \circ B \quad (\text{Bag Image})$$

5.2 Builtin Functions

$$\begin{aligned} \text{card}(S) &\rightsquigarrow \sum x \cdot x \in S \mid 1 && (\text{Cardinality}) \\ \text{dom}(R) &\rightsquigarrow \{ x \mapsto y \cdot x \mapsto y \in R \mid x \} && (\text{Domain}) \\ \text{ran}(R) &\rightsquigarrow \{ x \mapsto y \cdot x \mapsto y \in R \mid y \} && (\text{Range}) \end{aligned}$$

5.3 Sets

Let S, T be sets, P, E expressions, and x, e any type.

5.3.1 Set Comprehension Construction

Intuition All set-like variables and literals are decomposed into set comprehensions.

Post-condition

- All terms with a set-like type (relations, bags, sets, sequences, etc.) must be in comprehension form.

$$S \cup T \rightsquigarrow \{x \cdot x \in S \vee x \in T\} \quad (\text{Predicate Operations - Union})$$

$$S \cap T \rightsquigarrow \{x \cdot x \in S \wedge x \in T\} \quad (\text{Predicate Operations - Intersection})$$

$$S \setminus T \rightsquigarrow \{x \cdot x \in S \wedge x \notin T\} \quad (\text{Predicate Operations - Difference})$$

$$x \in \{e\} \rightsquigarrow x = e \quad (\text{Singleton Membership } ^a)$$

$$x \in \oplus(E \mid P) \rightsquigarrow P \wedge x = E \quad (\text{Membership Collapse } ^b)$$

^aCurrently unused. We need to be careful to handle the case where x is a free variable.

^bRule only matches inside the predicate of a quantifier. Explicitly enumerating all matches for all quantification types and predicate cases (ANDs, ORs, etc.) would require too much boilerplate. x must be bound by the encasing quantifier.

The \oplus operator represents any quantifier that returns a set-like type (ex. generalized union/intersection, set comprehension, relation comprehension).

5.3.2 Disjunctive Normal Form

Intuition All quantifier predicates are expanded to DNF (i.e. \wedge -operations nested within top-level \vee -operations).

Notes All matches of this phase occur only inside quantifier predicates.

Post-condition

- All terms with a set-like type (relations, bags, sets, sequences, etc.) must be in comprehension form.
- Quantifier predicates are in disjunctive normal form - top level or-clauses with inner and-clauses.
- If no \vee operators exist within a quantifier predicate, the predicate must only contain \wedge operators.

$$x_1 \wedge \dots \wedge (x_i \wedge x_{i+1}) \wedge \dots \rightsquigarrow x_1 \wedge \dots \wedge x_i \wedge x_{i+1} \wedge \dots \quad (\text{Flatten Nested } \wedge)$$

$$x_1 \vee \dots \vee (x_i \vee x_{i+1}) \vee \dots \rightsquigarrow x_1 \vee \dots \vee x_i \vee x_{i+1} \vee \dots \quad (\text{Flatten Nested } \vee)$$

$$\neg \neg x \rightsquigarrow x \quad (\text{Double Negation})$$

$$\neg(x \vee y) \rightsquigarrow \neg x \wedge \neg y \quad (\text{Distribute De Morgan - Or})$$

$$\neg(x \wedge y) \rightsquigarrow \neg x \vee \neg y \quad (\text{Distribute De Morgan - And})$$

$$x \wedge (y \vee z) \rightsquigarrow (x \wedge y) \vee (x \wedge z) \quad (\text{Distribute } \wedge \text{ over } \vee)$$

5.3.3 Predicate Simplification 1

Intuition Restructuring quantifications to better suit code generation (ex. Some comprehensions require nested loops to generate all elements).

Post-condition

- All terms with a set-like type (relations, bags, sets, sequences, etc.) must be in comprehension form.
- Quantifier predicates are in disjunctive normal form - top level or-clauses with inner and-clauses.
- One bound variable per quantifier.

$$\{x, y \cdot P \wedge Q \mid E\} \rightsquigarrow \{x \cdot P \mid \{y \cdot Q \mid E\}\} \quad (\text{Nesting } ^a)$$

$$\{x \cdot \bigwedge P_i \mid E\} \rightsquigarrow \{x \cdot \bigvee \bigwedge P_i \mid E\} \quad (\text{Or-wrapping } ^b)$$

^a y cannot occur in P . Nesting selection is currently arbitrary; in the future we plan to find an optimal selection heuristic.

^bTo simplify the matching process later on, we wrap every top-level AND statement (which is guaranteed to be a ListOp by the dataclass field type definition) with an OR.

5.3.4 Generator Selection

Intuition Selecting element generators for use as iterables in generated for-loops.

Notes All matches of this phase occur only inside quantifier predicates.

Post-condition

- All terms with a set-like type (relations, bags, sets, sequences, etc.) must be in comprehension form.
- Quantifier predicates are in disjunctive normal form - top level or-clauses with inner and-clauses.
- One bound variable per quantifier.
- Each top-level or-clause within quantification predicates must have one selected ‘generator’ predicate of the form $x \in S$ that loops over its bound dummy variable x .

Brainstorming selection heuristics Conditions to consider:

- Prefer relation chains in case of nested loops (ex. $x \mapsto y \in R \wedge y \mapsto z \in Q$, but what about renaming? - $x \mapsto y \in R \wedge y' \mapsto z \in Q \wedge y = y'$). Perhaps this could just be a preference for relations over sets if we see a nested quantifier.
- Choose most common generator among a list of or-clauses (allows for greater simplification later on)

- Choose smallest generator (by set size). This information may not always be accessible (ex. if a set is created by reading values from standard input). Functions may need to choose this on a case-by-case basis (ie. one function call could have two args, with a small set on the left arg and a large set on the right arg. But what if the sizes are reversed later in the code? We've already statically lowered it).
- When should constant folding happen? Equality substitution necessary for this?
- This may need to work with nesting considerations.
- What if we try moving these optimizations to *loop* structures far later in the pipeline?

$$\bigwedge P_i \rightsquigarrow P_g \wedge \bigwedge_{P_i \neq P_g} P_i \quad (\text{Generator Selection } ^a)$$

^aThe LH term must occur inside a quantifier's predicate - one generator per or-clause. P_g is the generator, a single clause distinguished from the rest of $\bigwedge P_i$. P_g must be a set membership operation, of the form $x \in S$. Currently, the selection of P_g is arbitrary (and thus the rewrite system is not confluent), but heuristics may be added later to choose optimal generators.

5.3.5 Predicate Simplification 2

Intuition Group up similar or-separated clauses sharing the same generator (factoring that generator).

Post-condition

- All terms with a set-like type (relations, bags, sets, sequences, etc.) must be in comprehension form.
- Quantifier predicates are in disjunctive partially-normal form - top level or-clauses with inner and-clauses. Generator predicates can "factor" one or more or-clauses. The expected predicate shape is either $\vee(P_g \wedge (\vee P_i))$ or $\vee(P_g \wedge P_i)$.
- One bound variable per quantifier.
- Each top-level or-clause within quantification predicates must have one selected 'generator' predicate of the form $x \in S$ that loops over its bound dummy variable x .

$$(P_g \wedge \bigwedge P_i) \vee (P_g \wedge \bigwedge Q_i) \rightsquigarrow P_g \wedge (\bigwedge P_i \vee \bigwedge Q_i)$$

(Reduce duplicate generators ^a)

^aThe LH term must occur inside a quantifier's predicate. Combines clauses with the same generator. Requires Generator Selection to be run first. We leave this optimization out for now, given difficulties with counterexample $\{s \cdot (x \in \{1, 2\} \wedge (s = x + 1 \vee s = x + 2)) \vee (x \in \{4\} \wedge s = x + 1)\}$.

5.3.6 Predicate Simplification 3

Intuition Adds extra nested or-structures to ensure consistency within each top-level or-clause.

Post-condition

- All terms with a set-like type (relations, bags, sets, sequences, etc.) must be in comprehension form.
- Quantifier predicates are in disjunctive partially-normal form - top level or-clauses with inner and-clauses. Generator predicates can "factor" one or more or-clauses. The expected predicate shape is $\vee(P_g \wedge (\vee P_i))$.
- One bound variable per quantifier.
- Each top-level or-clause within quantification predicates must have one selected 'generator' predicate of the form $x \in S$ that loops over its bound dummy variable x .

$$P_g \wedge \bigwedge P_i \rightsquigarrow P_g \wedge (\bigvee \bigwedge P_i) \quad (\text{Partially-normalized Disjunctive Form})$$

5.3.7 Code Generation 1

Intuition Start lowering expressions into imperative-like loops.

Post-condition

- All quantifiers are transformed into *loop* structures.
- *loop* predicates are of the form $x \in S \wedge P \wedge \bigvee Q$
- All *loop* predicates have an assigned generator of the form $x \in S$.

$$\oplus E \mid P \rightsquigarrow \begin{array}{l} a := \text{identity}(\oplus) \\ \mathbf{loop} \ P : \\ \quad a := \text{accumulate}(a, E) \end{array} \quad (\text{Quantifier Generation } ^a)$$

^a \oplus works for any quantifier (but not \forall and \exists). The identity and accumulate functions are determined by the realized \oplus . For example, if $\oplus = \sum$, the identity is 0 and accumulate is addition.

$$\begin{array}{ccc}
a := accumulator(a, \oplus(E \mid P)) & \rightsquigarrow & \begin{array}{l} \mathbf{loop} \ P : \\ \quad a := accumulator(a, E) \\ \quad \text{(Nested Quantifier)} \end{array} \\
\\
\mathbf{loop} \ \bigvee P_i & \rightsquigarrow & \begin{array}{l} \mathbf{loop} \ P_0 : \\ \quad \text{body} \\ \mathbf{loop} \ P_1 \wedge \neg P_0 : \\ \quad \text{body} \\ \mathbf{loop} \ P_2 \wedge \neg(\bigvee_{i<2} P_i) : \\ \quad \text{body} \\ \dots \\ \text{(Disjunct conditional)} \end{array}
\end{array}$$

5.3.8 Code Generation 2

Intuition Eliminate all intermediate loop structures.

Post-condition

- AST is in imperative code style (for loops, if statements, etc).
- All *loop* and quantification constructs have been eliminated.
- Some variables may not be defined.

$$\begin{array}{ccc}
\mathbf{loop} \ P_g \wedge P \wedge \bigwedge Q_i & \rightsquigarrow & \begin{array}{l} \mathbf{if} \ \bigwedge_{free(P_i)} P_i \ \mathbf{then} \\ \quad \mathbf{for} \ P_g \ \mathbf{do} \\ \quad \quad \mathbf{if} \ \bigwedge_{bound(P_i)} P_i \ \mathbf{then} \\ \quad \quad \quad \mathbf{if} \ Q_0 \ \mathbf{then} \\ \quad \quad \quad \quad \text{body} \\ \quad \quad \quad \mathbf{if} \ Q_1 \wedge \neg Q_0 \ \mathbf{then} \\ \quad \quad \quad \quad \text{body} \\ \quad \quad \quad \mathbf{if} \ Q_2 \wedge \neg(\bigvee_{i<2} Q_i) \ \mathbf{then} \\ \quad \quad \quad \quad \text{body} \\ \quad \quad \dots \text{(Conjunct conditional } ^a) \end{array}
\end{array}$$

^aFunction *free* returns clauses in *P* that contain only free + defined variables. *bound* returns the clauses that contain the bound variable *x* or undefined variables. *P_g* is the selected generator.

5.3.9 Equality Elimination

Intuition Eliminate all undefined variables (with implicit \exists quantifiers).

Post-condition

- AST is in imperative code style (for loops, if statements, etc).
- All variables are defined.

$$\begin{array}{ccc} \text{if } Identifier(x) = E \wedge P \text{ then} & \rightsquigarrow & \text{if } P[x := E] \text{ then} \\ \text{body} & & \text{body}[x := E] \\ & & \text{(Equality Elimination } ^a) \end{array}$$

^a x is an undefined, unbound variable in the current scope. E is an expression that does not contain x

5.4 Relations

$$\begin{array}{ll} R[S] \rightarrow \{ x \mapsto y \in R \cdot x \in S \mid y \} & \text{(Image)} \\ x \mapsto y \in S \times T \rightarrow x \in S \wedge y \in T & \text{(Product)} \\ x \mapsto y \in R^{-1} \rightarrow y \mapsto x \in R & \text{(Inverse)} \\ x \mapsto y \in (Q \circ R) \rightarrow x \mapsto z \in Q \wedge z' \mapsto y \in R \wedge z = z' & \text{(Composition)} \\ \exists z \cdot x \mapsto z \in Q \wedge z \mapsto y \in R \rightarrow x \mapsto z \in Q \wedge z' \mapsto y \in R \wedge z = z' & \text{(Exists Elimination)} \\ R \triangleleft Q \rightarrow Q \cup (dom(Q) \triangleleft R) & \text{(Override)} \end{array}$$

$$\begin{array}{ll} S \triangleleft R \rightarrow filter(fst \in S, R) & \text{(Domain Restriction)} \\ S \triangleleft R \rightarrow filter(fst \notin S, R) & \text{(Domain Subtraction)} \\ R \triangleright S \rightarrow filter(snd \in S, R) & \text{(Range Restriction)} \\ R \triangleright S \rightarrow filter(snd \notin S, R) & \text{(Range Subtraction)} \end{array}$$

5.5 Bags

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

$$\begin{array}{ll} S \cup T \rightarrow \{ (x, r) \cdot x \in set(S) \cup set(T) \mid r = \max(\#(x, S), \#(x, T)) \} & \text{(Union)} \\ S \cap T \rightarrow \{ (x, \min(a, b)) \cdot (x, a) \in S \wedge (x, b) \in T \mid a, b \geq 0 \} & \text{(Intersection)} \\ S + T \rightarrow \{ (x, r) \cdot x \in set(S) \cup set(T) \mid r = \#(x, S) + \#(x, T) \} & \text{(Sum)} \\ S - T \rightarrow \{ (x, r) \cdot (x, a) \in S \mid r = a - \#(x, T) \wedge r > 0 \} & \text{(Difference)} \\ size(S) \rightarrow \sum (x, r) \in S \cdot r & \text{(Size)} \end{array}$$

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 $set(S)$ can be implemented by taking the domain of the bag-representing relations.

6 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.