

Complete Specification and Design for Simile

Anthony Hunt

October 25, 2025

Contents

1	Overview	3
1.1	High Level Ideas	3
2	Motivation	3
3	Language	3
3.1	Lexicon	4
3.1.1	Whitespace	4
3.1.2	Identifiers	4
3.1.3	Primitives	4
3.1.4	Expressions	4
3.1.5	Reserved Keywords	6
3.2	Syntax	7
3.3	Static Analysis - Type System	10
3.3.1	Base Environment Interactions	10
3.3.2	The Set Type	10
3.3.3	Primitive Types	10
3.3.4	Collection Types	11
3.3.5	Relation Subtypes	12
3.3.6	Simple Commands	12
3.3.7	Quantification Expressions	13
3.3.8	Boolean Operations	14
3.3.9	Equality and Membership	14
3.3.10	Set Operations	15
3.3.11	Bag Operations	16
3.3.12	Relation Operations	16
3.3.13	Sequence Operations	17
3.3.14	Numerical Operations	17
3.3.15	Record Types	17
3.3.16	Compound Commands	18
3.3.17	Relational Subtypes After Operations	18
3.3.18	Built-in Functions	19
3.4	Definedness	20

3.5	Language Examples	20
4	Features and Behaviours	20
5	User-Facing Specification	20
5.1	(Abstract) Data Types	20
5.1.1	Primitives	20
5.1.2	Sets	21
5.1.3	Relations	22
5.1.4	Bags	23
5.1.5	Sequences	24
5.2	Data Type Operations	24
5.2.1	Primitives	24
5.2.2	Sets	24
5.2.3	Relations	24
5.2.4	Bags	24
5.2.5	Sequences	24
6	Implementation Specification	24
6.1	Abstract Data Types	24
6.2	Atomic Operations	24
6.3	Optimization Frontend	24
6.4	Optimization Backends	24
7	AST Lowering and Optimization Rules	24
7.1	Syntactic Sugar for Bags	24
7.2	Syntactic Sugar for Sequences	25
7.3	Builtin Functions	25
7.4	Comprehension Construction	25
7.5	Disjunctive Normal Form	26
7.6	Or-wrapping	27
7.7	Generator Selection	27
7.8	GSP to Loop	28
7.9	Relational Subtyping Loop Simplification	29
7.10	Loop Code Generation	30
7.11	Replace and Simplify	31

1 Overview

This document serves as a living specification of the underlying term rewriting system used in the compiler for a modelling-focused programming language. It includes an overview of the language Simile, a user-facing reference, and implementation plans.

1.1 High Level Ideas

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.

2 Motivation

3 Language

This section denotes the design and behaviour of Simile.

Much of the language design is inherited from Event-B and the mathematical language definition in [1].

3.1 Lexicon

Since Simile makes extensive use of Event-B-like set notation and otherwise non-ASCII characters, this section contains a list of such characters, translations to ASCII symbols where suitable, and other useful lexing information.

3.1.1 Whitespace

Similar to Python, indentation and newline tokens are significant in Simile. In particular, nested code blocks, as in if-statements, for-loops, and procedure definitions, require the use of indentations to signify the beginning and end. Indentation/newlines may also be used within collection-type comprehensions or enumerations, similar to the style seen in JSON. Comments are single-lined and start with a pound symbol (#).

3.1.2 Identifiers

Identifiers in Simile follow the standard regular expression, using only ASCII letters, underscores, and numbers: [a – zA – Z_][a – zA – Z_0 – 9].

3.1.3 Primitives

Integers and floating point numbers follow a similar style to modern programming languages. The dash character (–) may be used in either subtraction or as a unary negation. Superscript numbers may be used for constant exponentiation. Booleans (True/False) follow Python style, with capitalized first letters. Strings are delimited with double-quotes ("") and may be multilined.

3.1.4 Expressions

ASCII tokens will be provided as an alternative to unicode format where suitable.

The following are symbols used for expressions:

Token	ASCII Token	Name
.	.	Center Dot (Quantifier Separator)
,	,	Comma
:	:	Colon
;	;	Semicolon
		Vertical Bar
λ	lambda	Lambda
(Left Parenthesis
)		Right Parenthesis
<	<<	Left Angle Bracket
	[Left Angle Bracket (Alternative)
>	>>	Right Angle Bracket
]	Right Angle Bracket (Alternative)

Token	ASCII Token	Name
{		Left Brace
}		Right Brace
[[[Left Double Bracket
]]]	Right Double Bracket
=		Equals
\neq	!=	Not Equals
\neg	not	Not
$!$!	Not (Alternative)
\wedge	and	And
\vee	or	Or
\Rightarrow	=>	Implies
\equiv	==	Equivalent
$\not\equiv$!==	Not Equivalent
\forall	forall	For All
\exists	exists	There Exists
$+$	+	Plus
$-$	-	Minus
$*$	*	Multiplication
\times	><	Cartesian Product
	/	Division
	div	Integer Division
	mod	Modulo
	^	Exponent
$<$		Less Than
\leq	<=	Less Than or Equal
$>$		Greater Than
\geq	>=	Greater Than or Equal
\in	in	In (Membership)
\notin	not in	Not In (Membership)
\cup	\vee	Union
\cap	\wedge	Intersection
\setminus	\setminus	Backslash (Set Minus)
\subset	<<:	Subset
\subseteq	<:	Subset or Equal
\supset	:>>	Superset
\supseteq	:>	Superset or Equal
$\not\subset$		Not Subset
$\not\subseteq$		Not Subset or Equal
$\not\supset$		Not Superset
$\not\supseteq$		Not Superset or Equal
\cup	union	General Union
\cap	intersection	General Intersection
\mapsto	->	Maplet
\oplus	<+>	Relation Overriding
\circ	circ	Composition

Token	ASCII Token	Name
++	++	Concatenation
-1	~	Inverse
▷	<	Domain Restriction
▷	<<	Domain Subtraction
▷	>	Range Restriction
▷	>>	Range Subtraction
↔	<->	Relation
↔	<<->	Total Relation ¹
↔	<->>	Surjective Relation ²
↔	<<->>	Total Surjective Relation ³
→	+->	Partial Function
→	->	Total Function
⤠	>+>	Partial Injection
⤠	>->	Total Injection
⤠	+->>	Partial Surjection
⤠	->>	Total Surjection
⤠	>->>	Bijection
..		Upto
℘	powerset	Powerset
℘₁	powerset1	Non-Empty Powerset

The following are symbols used for programming constructs:

(Unicode) Token	(ASCII) Token	Name
:=	:=	Assignment
:∈	::	Non-deterministic Membership Assignment
→	->	Right Arrow

3.1.5 Reserved Keywords

The following list of symbols and identifiers are reserved for programming constructs:

true	false	if	else	for
while	record	enum	procedure	is
is not	return	exists	break	continue
from	import	skip	with	type
min	max	choice		

¹Brave Sans Mono unicode: U+E100

²Brave Sans Mono unicode: U+E101

³Brave Sans Mono unicode: U+E102

3.2 Syntax

We follow the standard EBNF notation for describing the language structure.

Start ::= Statements

Statements ::= {SimpleStmt | CompoundStmt}

SimpleStmt ::= Expr [[‘:’ Expr] (‘:=’ | ‘∈’) Expr] NEWLINE [WithStmt]
| ControlFlowStmt NEWLINE
| ImportStmt NEWLINE

WithStmt ::= INDENT ‘with’ Expr NEWLINE {Expr NEWLINE} DEDENT

ControlFlowStmt ::= ‘return’ [Expr]
| ‘break’
| ‘continue’
| ‘skip’

ImportStmt ::= ‘import’ STRING
| ‘from’ STRING ‘import’ ImportList

ImportList ::= ‘*’
| FlatTupleIdentifier
| ‘(’ FlatTupleIdentifier ‘)’

FlatTupleIdentifier ::= IDENTIFIER {‘,’ IDENTIFIER }

Expr ::= Quantification | Predicate

Quantification ::= ‘λ’ FlatTupleIdentifier ‘.’ Predicate ‘|’ Expr
| ‘U’ QuantificationBody
| ‘∩’ QuantificationBody
| ‘min’ QuantificationBody
| ‘max’ QuantificationBody

QuantificationBody ::= IdentList ‘.’ Predicate ‘|’ Expr
| Expr ‘|’ Predicate

IdentList ::= IdentListItem {‘,’ IdentListItem}

IdentListItem ::= IDENTIFIER
| IdentListItem ‘→’ IdentListItem
| ‘(’ IdentList ‘)’

Predicate ::= UnquantifiedPredicate
| ‘V’ IdentList ‘.’ Predicate
| ‘E’ IdentList ‘.’ Predicate

```

UnquantifiedPredicate ::= Implication
| UnquantifiedPredicate ('≡' | '≢') Implication

Implication ::= Disjunction
| Disjunction {'⇒' Disjunction}

Disjunction ::= Conjunction {'∨' Conjunction}

Conjunction ::= Negation {'∧' Negation}

Negation ::= AtomBool
| '¬' Negation

AtomBool ::= PairExpr [('=' | '≠' | '<' | '>' | '≤' | '≥' | '∈' | '∉' | '⊂' | '⊆' |
    '⊃' | '⊇' | '⊄' | '⊅' | '⊅' | '⊅') PairExpr]

PairExpr ::= RelSetExpr
| PairExpr '↪' RelSetExpr

RelSetExpr ::= SetExpr [(↔ | ↔↔ | ↔↔↔ | ↔↔↔ | → | →↔ | →↔↔ | →↔↔↔ |
    ↔↔ | ↔↔↔) RelSetExpr]

SetExpr ::= IntervalExpr
| IntervalExpr {'∪' IntervalExpr}
| IntervalExpr {'×' IntervalExpr}
| IntervalExpr {'⊕' IntervalExpr}
| IntervalExpr {'◦' IntervalExpr}
| IntervalExpr {'∩' IntervalExpr}
| IntervalExpr {'++' IntervalExpr}
| IntervalExpr ('⟨' | '⟨') IntervalExpr {'∩' IntervalExpr} [RelSubExpr]

RelSubExpr ::= ('\\' | '▷' | '▷') IntervalExpr

IntervalExpr ::= ArithmeticExpr [..] ArithmeticExpr

ArithmeticExpr ::= Term
| ArithmeticExpr ('+' | '-') Term

Term ::= Factor
| Term ('*' | '/') Factor

Factor ::= [(+' | -')] Power

Power ::= Primary ['^' Factor | ['-']('INTEGER' | 'FLOAT')]

Primary ::= Atom
| Atom '.' IDENTIFIER

```

```

| Atom `(` Expr {` ,` Expr} `)`
| Atom ` [` Expr `]`

Atom ::= INTEGER | FLOAT | STRING | BOOLEAN | IDENTIFIER
| Set | Sequence | Bag | Tuple
| `P` `(` Expr `)` | `P1` `(` Expr `)` | `(` Expr `)`

Set ::= ` {` CollectionBody ` }`

Bag ::= ` [` CollectionBody ` ]`

Sequence ::= ` <` CollectionBody ` >`

Tuple ::= `(` `)` | `(` [NEWLINE INDENT] Expr `,` [NEWLINE [INDENT]]
[EnumerationBody] `)`

CollectionBody ::= QuantificationBody | EnumerationBody

EnumerationBody ::= [NEWLINE INDENT] Expr {` ,` [NEWLINE [INDENT]]
Expr} [NEWLINE DEDENT]

CompoundStmt ::= IfStmt | ForStmt | WhileStmt
| RecordStmt | ProcedureStmt

IfStmt ::= `if` Predicate `:` Block [ElseStmt]

ElseStmt ::= `else` `:` Block
| `else if` Predicate `:` Block [ElseStmt]

ForStmt ::= `for` IdentList `in` Expr `:` Block

WhileStmt ::= `while` Expr `:` Block

RecordStmt ::= `record` IDENTIFIER `:` NEWLINE INDENT (
`skip` | TypedName {(` ,` [NEWLINE] | NEWLINE) TypedName }
) NEWLINE DEDENT

ProcedureStmt ::= `procedure` IDENTIFIER `(` TypedName {` ,` TypedName
}`)` `→` Expr `:` Block

TypedName ::= IDENTIFIER [` :` Expr]

Block ::= SimpleStmt
| NEWLINE INDENT Statements DEDENT

```

3.3 Static Analysis - Type System

Let Γ be the type environment set consisting of elements either $identifier : Type$ or $Type$. The typing rules have been organized in distinct categories for readability. The notation $Type[x, y, z]$ denotes refinements on $Type$ either through instantiating parametric (generic) types or restricting values by way of `with` clauses.

Type refinements may be applicable to types with certain traits (ex., min/-max values for ordered types, definedness of expressions with certain values, set sizes).

3.3.1 Base Environment Interactions

Includes: Mechanisms for adding/extracting variables from the type environment.

$$\frac{\overline{\emptyset \vdash \diamond}}{\Gamma \vdash A \quad I \notin \text{dom}(\Gamma)} \quad (\text{Env } \emptyset)$$

$$\frac{}{\Gamma, I : A \vdash \diamond} \quad (\text{Env } I)$$

$$\frac{\Gamma_1, I : A, \Gamma_2 \vdash \diamond}{\Gamma_1, I : A, \Gamma_2 \vdash I : A} \quad (\text{Fetch Identifier})$$

3.3.2 The Set Type

Includes: Universal type; everything is a set. The empty set is a part of every type (aside from primitives).

$$\frac{\Gamma \vdash T}{\Gamma \vdash \text{set}[T]} \quad (\text{Powerset})$$

$$\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{EmptySet}} \quad (\text{Emptyset})$$

$$\frac{\Gamma \vdash \text{EmptySet}, \text{set}[T]}{\text{EmptySet} \leq \text{set}[T]} \quad (\text{Emptyset Top})$$

$$\frac{\Gamma \vdash e_1, \dots, e_n : A}{\Gamma \vdash \{e_1, \dots, e_n\} : \text{set}[A]} \quad (\text{Set Enumeration})$$

3.3.3 Primitive Types

Includes: Basic types built in to the language, like numbers (int and float), strings, and booleans. Although these could all be represented as sets, it is useful to distinguish these common types from other collections or user-defined types.

$$\begin{array}{c}
\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{bool}} \quad (\text{Primitives - bool}) \\
\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{int}} \quad (\text{Primitives - int}) \\
\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{float}} \quad (\text{Primitives - float}) \\
\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{str}} \quad (\text{Primitives - str}) \\
\frac{\Gamma \vdash \text{int}}{\Gamma \vdash \text{nat} \quad \text{nat} = \{ n \geq 0 \cdot n \in \text{int} \}} \quad (\text{Nat from int}) \\
\frac{\Gamma \vdash \text{int}, \text{float}}{\text{int} \leq \text{float}} \quad (\text{Int from float})
\end{array}$$

3.3.4 Collection Types

Includes: Sequences, relations, enums (frozen sets), and bags (multisets); Syntactic sugar for abstract data types defined in terms of sets.

$$\begin{array}{c}
\frac{\Gamma \vdash T_1, \dots, T_n}{\Gamma \vdash \text{tuple}[T_1, \dots, T_n] \quad \text{tuple}[T_1, \dots, T_n] = T_1 \times \dots \times T_n} \quad (\text{Tuples}) \\
\frac{\Gamma \vdash \diamond}{\Gamma \vdash () : \text{EmptyTuple}} \quad (\text{Empty Tuple}) \\
\frac{\Gamma \vdash \text{EmptyTuple}, T}{\Gamma \vdash \text{EmptyTuple} \subseteq \text{tuple}[T]} \quad (\text{Empty Tuple Top}) \\
\frac{\Gamma \vdash e_1 : A_1, \dots, e_n : A_n}{\Gamma \vdash (e_1, \dots, e_n) : \text{tuple}[A_1, \dots, A_n]} \quad (\text{Tuple Enumeration}) \\
\frac{\Gamma \vdash \text{set}[T_1], \text{set}[T_2]}{\Gamma \vdash T_1 \leftrightarrow T_2 \quad T_1 \leftrightarrow T_2 = \text{set}[\text{tuple}[T_1, T_2]]} \quad (\text{Relation Type}) \\
\frac{\Gamma \vdash T, \text{nat}}{\Gamma \vdash \text{bag}[T] = T \rightarrowtail \text{nat}} \quad (\text{Bag Type}) \\
\frac{\Gamma \vdash e_1, \dots, e_n : A}{\Gamma \vdash [e_1, \dots, e_n] : \text{bag}[A]} \quad (\text{Bag Enumeration}) \\
\frac{\Gamma \vdash T, \text{nat}}{\Gamma \vdash \text{bag}[T] = \text{nat} \rightarrowtail T} \quad (\text{Sequence Type}) \\
\frac{\Gamma \vdash e_1, \dots, e_n : A}{\Gamma \vdash \langle e_1, \dots, e_n \rangle : \text{sequence}[A]} \quad (\text{Sequence Enumeration}) \\
\frac{\Gamma \vdash s_1, \dots, s_n : \text{str} \quad A = \{s_1, \dots, s_n\} \quad \text{immutable}(A)}{\Gamma \vdash \text{enum}[A]} \quad (\text{Enum from static set})
\end{array}$$

3.3.5 Relation Subtypes

Includes: Relation subtypes examine four key properties of relations: totality (t), surjectivity (t_r), one-to-manyness ($1-$), and many-to-oneness ($1-r$). Syntactic sugar to produce refined relations with these properties follows the definitions presented in Event-B. See Section 5.1.3 for more details.

$$\begin{array}{c}
 \frac{\Gamma \vdash set[T_1], set[T_2]}{\Gamma \vdash T_1 \leftrightarrow T_2 \quad T_1 \leftrightarrow T_2 = (T_1 \leftrightarrow T_2)_{\text{Relation Subtype}}[\text{with } t] - \text{Total Relation}} \\
 \frac{\Gamma \vdash set[T_1], set[T_2]}{\Gamma \vdash T_1 \leftrightarrow T_2 \quad T_1 \leftrightarrow T_2 = (T_1 \leftrightarrow T_2)_{\text{Relation Subtype}}[\text{with } t_r] - \text{Surjective Relation}} \\
 \frac{\Gamma \vdash set[T_1], set[T_2]}{\Gamma \vdash T_1 \leftrightarrow T_2 \quad T_1 \leftrightarrow T_2 = (T_1 \leftrightarrow T_2)_{\text{Relation Subtype}}[\text{with } t, \text{with } t_r] - \text{Total Surjective Relation}} \\
 \frac{\Gamma \vdash set[T_1], set[T_2]}{\Gamma \vdash T_1 \rightarrow T_2 \quad T_1 \rightarrow T_2 = (T_1 \rightarrow T_2)_{\text{Relation Subtype}}[\text{with } t] - \text{Partial Function}} \\
 \frac{\Gamma \vdash set[T_1], set[T_2]}{\Gamma \vdash T_1 \rightarrow T_2 \quad T_1 \rightarrow T_2 = (T_1 \rightarrow T_2)_{\text{Relation Subtype}}[\text{with } t, \text{with } 1] - \text{Total Function}} \\
 \frac{\Gamma \vdash set[T_1], set[T_2]}{\Gamma \vdash T_1 \rightarrow T_2 \quad T_1 \rightarrow T_2 = (T_1 \rightarrow T_2)_{\text{Relation Subtype}}[\text{with } t, \text{with } 1, \text{with } 1] - \text{Injection}} \\
 \frac{\Gamma \vdash set[T_1], set[T_2]}{\Gamma \vdash T_1 \rightarrow T_2 \quad T_1 \rightarrow T_2 = (T_1 \leftrightarrow T_2)_{\text{Relation Subtype}}[\text{with } t, \text{with } 1, \text{with } 1] - \text{Total Surjection}} \\
 \frac{\Gamma \vdash set[T_1], set[T_2]}{\Gamma \vdash T_1 \rightarrow T_2 \quad T_1 \rightarrow T_2 = (T_1 \leftrightarrow T_2)_{\text{Relation Subtype}}[\text{with } t, \text{with } t_r, \text{with } 1] - \text{Total Surjection}}
 \end{array}$$

3.3.6 Simple Commands

Includes: Assignment, type assignment, type refinements, simple programming language commands/statements.

$$\begin{array}{c}
 \frac{\Gamma \vdash E : A \quad I \notin \text{dom}(\Gamma) \vee \Gamma \vdash I : A}{\Gamma \vdash I : A := E} \quad (\text{Variable Assignment}) \\
 \frac{\Gamma \vdash T \quad I \notin \text{dom}(\Gamma)}{\Gamma \vdash I := T} \quad (\text{Type Alias Assignment})
 \end{array}$$

$$\begin{array}{c}
\frac{\Gamma \vdash I := T}{\Gamma, I \vdash \diamond \quad I = T} \quad (\text{Type Alias}) \\
\frac{\Gamma \vdash E : A, E_1 : A_1, \dots, E_n : A_n \quad (I \notin \text{dom}(\Gamma) \vee \Gamma \vdash I : A)}{\Gamma \vdash I : A[\text{with } E_1 : A_1, \dots, \text{with } E_n : A_n \text{ Refined } \bar{E} \text{ Variable Assignment}]} \\
\frac{\Gamma \vdash T, E_1 : A_1, \dots, E_n : A_n \quad I \notin \text{dom}(\Gamma)}{\Gamma \vdash I := T[\text{with } E_1 : A_1, \dots, \text{with } E_n : A_n]} \quad (\text{Refined Type Alias}) \\
\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{break} : C} \quad (\text{Command - break}) \\
\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{continue} : C} \quad (\text{Command - continue}) \\
\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{skip} : C} \quad (\text{Command - skip}) \\
\frac{\Gamma \vdash E : A}{\Gamma \vdash (\text{return } E) : C} \quad (\text{Command - return})
\end{array}$$

3.3.7 Quantification Expressions

Includes: Set, sequence, relation, and bag comprehensions, in addition to (set based) lambda expressions. Comprehension variables may be bound through either an *Identifier* or an arbitrarily nested tuple of *Identifiers*

$$\begin{array}{c}
\frac{\Gamma \vdash E : A, P : \text{bool} \quad I_1 : T_1, \dots, I_n : T_n \text{ are vars in } E}{\Gamma \vdash (\lambda I_1, \dots, I_n \cdot P \mid E) : (T_1, \dots, T_n) \rightarrow A} \quad (\text{Lambda Expression}) \\
\frac{\begin{array}{c} \Gamma \vdash E : A, P : \text{bool} \\ I_1, \dots, I_n : \text{TupleIdentifier} \\ Ids = [\text{identifiers}(I_i) \cdot 0 \leq i \leq n] \\ \forall e \mapsto \text{count} \in Ids \cdot \text{count} = 1 \wedge e \notin \text{dom}(\Gamma) \\ M = \{I_i \mapsto T \cdot 0 \leq i \leq n \\ \wedge \exists g \cdot \Gamma \vdash g : \text{set}[T] \\ \wedge \text{structuralMatch}(e, T) \quad (\text{Quantification Body } ^4) \\ \wedge \text{hasGenerator}(P, e \in g) \\ \wedge \text{structuralMatchInAllORs}(e, T, P)\} \end{array}}{\Gamma \vdash (I_1, \dots, I_n \cdot P \mid E) : \text{QuantificationBody}[M, A]} \\
\frac{\Gamma \vdash (I \cdot P \mid E) : \text{QuantificationBody}[M, A]}{\Gamma \vdash (\bigcup I \cdot P \mid E) : A} \quad (\text{General Union}) \\
\frac{\Gamma \vdash (I \cdot P \mid E) : \text{QuantificationBody}[M, A]}{\Gamma \vdash (\bigcap I \cdot P \mid E) : A} \quad (\text{General Intersection}) \\
\frac{\Gamma \vdash (I \cdot P) : \text{QuantificationBody}[M, \text{bool}]}{\Gamma \vdash (\forall I \cdot P) : \text{bool}} \quad (\text{Forall}) \\
\frac{\Gamma \vdash (I \cdot P) : \text{QuantificationBody}[M, \text{bool}]}{\Gamma \vdash (\exists I \cdot P) : \text{bool}} \quad (\text{Exists})
\end{array}$$

$$\begin{array}{c}
 \frac{\Gamma \vdash (I \cdot P \mid E) : QuantificationBody[M, A]}{\Gamma \vdash \{ I \cdot P \mid E \} : set[A]} \quad (\text{Set Comprehension}) \\
 \frac{\Gamma \vdash (I \cdot P \mid E) : QuantificationBody[M, A]}{\Gamma \vdash \llbracket I \cdot P \mid E \rrbracket : bag[A]} \quad (\text{Bag Comprehension}) \\
 \frac{\Gamma \vdash (I \cdot P \mid E) : QuantificationBody[M, A]}{\Gamma \vdash [I \cdot P \mid E] : sequence[A]} \quad (\text{Sequence Comprehension})
 \end{array}$$

3.3.8 Boolean Operations

Includes: Common predicate logic operations (equivalence, implication, and, or, etc.).

$$\begin{array}{c}
 \frac{\Gamma \vdash b_1, b_2 : bool}{\Gamma \vdash b_1 \equiv b_2 : bool} \quad (\text{Boolean Operations - Equivalence}) \\
 \frac{\Gamma \vdash b_1, b_2 : bool}{\Gamma \vdash b_1 \not\equiv b_2 : bool} \quad (\text{Boolean Operations - Not Equivalence}) \\
 \frac{\Gamma \vdash b_1, b_2 : bool}{\Gamma \vdash b_1 \implies b_2 : bool} \quad (\text{Boolean Operations - Implication}) \\
 \frac{\Gamma \vdash b_1, b_2 : bool}{\Gamma \vdash b_1 \wedge b_2 : bool} \quad (\text{Boolean Operations - And}) \\
 \frac{\Gamma \vdash b_1, b_2 : bool}{\Gamma \vdash b_1 \vee b_2 : bool} \quad (\text{Boolean Operations - Or}) \\
 \frac{\Gamma \vdash b_1 : bool}{\Gamma \vdash \neg b_1 : bool} \quad (\text{Boolean Operations - Negation})
 \end{array}$$

3.3.9 Equality and Membership

Includes: Equality and ordering comparisons, set membership.

$$\begin{array}{c}
 \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2}{\Gamma \vdash a_1 = a_2 : bool} \quad (\text{Equals}) \\
 \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2}{\Gamma \vdash a_1 \neq a_2 : bool} \quad (\text{Not Equals}) \\
 \frac{\Gamma \vdash a_1, a_2 : T \quad T \in \{int, float, nat\}}{\Gamma \vdash a_1 < a_2 : bool} \quad (\text{Less Than}) \\
 \frac{\Gamma \vdash a_1, a_2 : T \quad T \in \{int, float, nat\}}{\Gamma \vdash a_1 \leq a_2 : bool} \quad (\text{Less Than or Equal}) \\
 \frac{\Gamma \vdash a_1, a_2 : T \quad T \in \{int, float, nat\}}{\Gamma \vdash a_1 > a_2 : bool} \quad (\text{Greater Than}) \\
 \frac{\Gamma \vdash a_1, a_2 : T \quad T \in \{int, float, nat\}}{\Gamma \vdash a_1 \geq a_2 : bool} \quad (\text{Greater Than or Equal})
 \end{array}$$

$$\begin{array}{c}
 \frac{\Gamma \vdash a_1 : T_1, a_2 : set[T_2]}{\Gamma \vdash a_1 \in a_2 : bool} \quad (\text{Set Membership}) \\
 \frac{\Gamma \vdash a_1 : T_1, a_2 : set[T_2]}{\Gamma \vdash a_1 \notin a_2 : bool} \quad (\text{Not Set Membership})
 \end{array}$$

3.3.10 Set Operations

Includes: Subset comparisons, union, intersection, cartesian product, etc.
Also includes Maplet and Range.

$$\begin{array}{c}
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \subset a_2 : bool} \quad (\text{Set Operations - Subset}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \subseteq a_2 : bool} \quad (\text{Set Operations - Subset Equal}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \supset a_2 : bool} \quad (\text{Set Operations - Superset}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \supseteq a_2 : bool} \quad (\text{Set Operations - Superset Equal}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \not\subset a_2 : bool} \quad (\text{Set Operations - Not Subset}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \not\subseteq a_2 : bool} \quad (\text{Set Operations - Not Subset Equal}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \not\supset a_2 : bool} \quad (\text{Set Operations - Not Superset}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \not\supseteq a_2 : bool} \quad (\text{Set Operations - Not Superset Equal}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \not\ni a_2 : set[T]} \quad (\text{Set Operations - Union}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \cap a_2 : set[T]} \quad (\text{Set Operations - Intersection}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \setminus a_2 : set[T]} \quad (\text{Set Operations - Difference}) \\
 \frac{\Gamma \vdash a_1 : set[T_1], a_2 : set[T_2]}{\Gamma \vdash a_1 \times a_2 : T_1 \leftrightarrow T_2} \quad (\text{Cartesian Product}) \\
 \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2}{\Gamma \vdash a_1 \mapsto a_2 : tuple[T_1, T_2]} \quad (\text{Maplet}) \\
 \frac{\Gamma \vdash a_1, a_2 : int}{\Gamma \vdash a_1 .. a_2 : set[int, \text{with } exclude(n) = a_1 < n, \text{ with } exclude(n) = a_2 \geq n]} \quad (\text{Numerical Range})
 \end{array}$$

$$\frac{\Gamma \vdash a : T}{\Gamma \vdash \mathbb{P}(a) : set[T]} \quad (\text{Set Operations - Powerset})$$

3.3.11 Bag Operations

Includes: Bag union, bag intersection, etc.

$$\begin{array}{c} \frac{\Gamma \vdash a_1, a_2 : bag[T]}{\Gamma \vdash a_1 \cup a_2 : bag[T]} \quad (\text{Bag Operations - (Max) Union}) \\ \frac{\Gamma \vdash a_1, a_2 : bag[T]}{\Gamma \vdash a_1 \cap a_2 : bag[T]} \quad (\text{Bag Operations - (Min) Intersection}) \\ \frac{\Gamma \vdash a_1, a_2 : bag[T]}{\Gamma \vdash a_1 + a_2 : bag[T]} \quad (\text{Bag Operations - Addition}) \\ \frac{\Gamma \vdash a_1, a_2 : bag[T]}{\Gamma \vdash a_1 - a_2 : bag[T]} \quad (\text{Bag Operations - Subtraction}) \\ \frac{\Gamma \vdash a : T_1 \leftrightarrow T_2, b : bag[T_1]}{\Gamma \vdash a[b] : bag[T_2] \quad a[b] = a^{-1} \circ b} \quad (\text{Bag Operations - Image}) \end{array}$$

3.3.12 Relation Operations

Includes: Relation overriding, composition, domain/range manipulation. See Section 5.1.3 for changes in relational subtypes after applying operations.

$$\begin{array}{c} \frac{\Gamma \vdash a : T_1 \leftrightarrow T_2, b : T_1}{\Gamma \vdash a(b) : T_2} \quad (\text{Relation Operations - Function Call}) \\ \frac{\Gamma \vdash a : T_1 \leftrightarrow T_2, b : set[T_1]}{\Gamma \vdash a[b] : set[T_2]} \quad (\text{Relation Operations - Image}) \\ \frac{\Gamma \vdash a_1, a_2 : T_1 \leftrightarrow T_2}{\Gamma \vdash a_1 \oplus a_2 : T_1 \leftrightarrow T_2} \quad (\text{Relation Operations - Overriding}) \\ \frac{\Gamma \vdash a_1 : T_1 \leftrightarrow T_2, a_2 : T_2 \leftrightarrow T_3}{\Gamma \vdash a_1 \circ a_2 : T_1 \leftrightarrow T_3} \quad (\text{Relation Operations - Composition}) \\ \frac{\Gamma \vdash a_1 : T_1, a_2 : T_1 \leftrightarrow T_2}{\Gamma \vdash a_1 \lhd a_2 : T_1 \leftrightarrow T_2} \quad (\text{Relation Operations - Domain Restriction}) \\ \frac{\Gamma \vdash a_1 : T_1, a_2 : T_1 \leftrightarrow T_2}{\Gamma \vdash a_1 \lhd a_2 : T_1 \leftrightarrow T_2} \quad (\text{Relation Operations - Domain Subtraction}) \\ \frac{\Gamma \vdash a_1 : T_1 \leftrightarrow T_2, a_2 : T_2}{\Gamma \vdash a_1 \rhd a_2 : T_1 \leftrightarrow T_2} \quad (\text{Relation Operations - Range Restriction}) \\ \frac{\Gamma \vdash a_1 : T_1 \leftrightarrow T_2, a_2 : T_2}{\Gamma \vdash a_1 \rhd a_2 : T_1 \leftrightarrow T_2} \quad (\text{Relation Operations - Range Subtraction}) \end{array}$$

3.3.13 Sequence Operations

Includes: Concatenation

$$\frac{\Gamma \vdash a_1, a_2 : \text{sequence}[T]}{\Gamma \vdash a_1 ++ a_2 : \text{sequence}[T]} \quad (\text{Sequence Operations - Concatenation})$$

3.3.14 Numerical Operations

Includes: Addition, subtraction, multiplication, etc. on natural numbers, integers, and floats.

$$\begin{array}{c} \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2 \quad T_i \in \{\text{int}, \text{nat}\}}{\Gamma \vdash a_1 \text{ div } a_2 : \max(T_1, T_2)} \quad (\text{Integer Operations - Division}) \\ \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2 \quad T_i \in \{\text{int}, \text{nat}\}}{\Gamma \vdash a_1 \text{ mod } a_2 : \max(T_1, T_2)} \quad (\text{Integer Operations - Modulo}) \\ \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2 \quad T_i \in \{\text{int}, \text{nat}, \text{float}\}}{\Gamma \vdash a_1 + a_2 : \max(T_1, T_2)} \quad (\text{Numerical Operations - Addition}) \\ \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2 \quad T_i \in \{\text{int}, \text{nat}, \text{float}\}}{\Gamma \vdash a_1 - a_2 : \max(T_1, T_2, \text{int})} \quad (\text{Numerical Operations - Subtraction}) \\ \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2 \quad T_i \in \{\text{int}, \text{nat}, \text{float}\}}{\Gamma \vdash a_1 / a_2 : \max(T_1, T_2, \text{float})} \quad (\text{Numerical Operations - Floating Division}) \\ \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2 \quad T_i \in \{\text{int}, \text{nat}, \text{float}\}}{\Gamma \vdash a_1 * a_2 : \max(T_1, T_2)} \quad (\text{Numerical Operations - Multiplication}) \\ \frac{\Gamma \vdash a_1 : T \quad T \in \{\text{int}, \text{nat}, \text{float}\}}{\Gamma \vdash -a_1 : \max(\text{int}, T)} \quad (\text{Numerical Operations - Negation}) \\ \frac{\Gamma \vdash a_1 : T_1, a_2 : T_2 \quad T_i \in \{\text{int}, \text{nat}, \text{float}\}}{\Gamma \vdash a_1 ^ a_2 : \max(T_1, T_2)} \quad (\text{Numerical Operations - Exponentiation}) \end{array}$$

3.3.15 Record Types

Includes: Types for named tuples.

$$\begin{array}{c} \frac{\Gamma \vdash a : \text{Record}[I : A, \dots]}{\Gamma \vdash a.I : A} \quad (\text{Records - Access}) \\ \frac{\Gamma \vdash A_1, \dots, A_n \quad \forall I_i, I_j \cdot I_i \neq I_j}{\Gamma \vdash \text{Record}[I_1 : A_1, \dots, I_n : A_n]} \quad (\text{Records - Type Definition}) \\ \frac{\Gamma \vdash a_1 : A_1, \dots, a_n : A_n}{\Gamma \vdash \text{record}(a_1 = A_1, \dots, a_n = A_n) : \text{Record}[a_1 \text{ }\overset{A_1}{\text{Records}}\text{ } a_n \text{ }\overset{A_n}{\text{Initialization}}]} \end{array}$$

3.3.16 Compound Commands

Includes: If, While, For, and Import statements, along with Procedure definitions and calls.

$$\begin{array}{c}
 \frac{\Gamma \vdash C_1, C_2}{\Gamma \vdash C_1 \setminus n C_2} \quad (\text{Command - Composition}) \\
 \frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{import } M : C} \quad (\text{Command - Valid Import Module}) \\
 \frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{from } M \text{ import } I_1, \dots, I_n : C} \quad (\text{Command - Valid Import Names}) \\
 \frac{\Gamma \vdash \text{import } M \text{ (with identifiers and types } I_1 : A_1, \dots, I_n : A_n\text{)}}{\Gamma, M = \text{record}(I_1 = A_1, \dots, I_n = A_n) \vdash \text{Command - Import Module}} \\
 \frac{\Gamma \vdash \text{from } M \text{ import } I_1, \dots, I_n \text{ (with types } A_1, \dots, A_n\text{)}}{\Gamma, I_1 : A_1, \dots, I_n : A_n \vdash \text{Command - Import Names}} \\
 \frac{\Gamma \vdash C, E, A_1, \dots, A_n}{\Gamma \vdash \text{procedure } I = C :: I : \text{procedure}[A_1, \dots, A_n, E]} \quad (\text{Command - Procedure Definition}) \\
 \frac{\Gamma \vdash \text{procedure}[A_1, \dots, A_n, E], a_1 : A_1, \dots, a_n : A_n}{\Gamma \vdash I(a_1, \dots, a_n) : E} \quad (\text{Command - Procedure Call}) \\
 \frac{\Gamma \vdash C, b : \text{bool}}{\Gamma \vdash \text{if } b : C} \quad (\text{Command - If}) \\
 \frac{\Gamma \vdash C_1, C_2, b : \text{bool}}{\Gamma \vdash \text{if } b : C_1 \text{ else : } C_2} \quad (\text{Command - If Else}) \\
 \frac{\Gamma \vdash C, b : \text{bool}}{\Gamma \vdash \text{while } b : C} \quad (\text{Command - While}) \\
 \frac{\Gamma \vdash C, p : \text{bool} \quad p = i \in G}{\Gamma \vdash \text{for } p : C} \quad (\text{Command - For}) \\
 \frac{\Gamma \vdash D :: I : A \quad \Gamma, I : A \vdash C}{\Gamma \vdash C \text{ substitute vars in } D} \quad (\text{Command - Block})
 \end{array}$$

3.3.17 Relational Subtypes After Operations

Includes: Changes in relational subtyping (totality, surjectivity, etc.) after applying operations. For simplicity, we implement subtype changes in a few binary operators as a boolean mask representing: [with total (t), with total on range (t_r), with one-to-manyness ($1-$), with many-to-oneness ($1-r$)]. Subtyping is done conservatively: we keep subtype properties only when they are guaranteed. See Section 5.1.3 for more on relational subtyping.

$$\begin{array}{c}
\frac{\Gamma \vdash s : set[T_1], r : (T_1 \leftrightarrow T_2)[True, x, y, z]}{\Gamma \vdash (s \lhd r) : (T_1 \leftrightarrow T_2)[False, x, y, z] \text{ Domain Restriction}} \\
\frac{\Gamma \vdash s : set[T_1], r : (T_1 \leftrightarrow T_2)[True, x, y, z]}{\Gamma \vdash (s \lhd r) : (T_1 \leftrightarrow T_2)[False, x, y, z] \text{ Domain Subtraction}} \\
\frac{\Gamma \vdash s : set[T_2], r : (T_1 \leftrightarrow T_2)[x, True, y, z]}{\Gamma \vdash (r \triangleright s) : (T_1 \leftrightarrow T_2)[x, False, y, z] \text{ Range Restriction}} \\
\frac{\Gamma \vdash s : set[T_2], r : (T_1 \leftrightarrow T_2)[x, True, y, z]}{\Gamma \vdash (r \triangleright s) : (T_1 \leftrightarrow T_2)[x, False, y, z] \text{ Range Subtraction}} \\
\frac{\Gamma \vdash r : (T_1 \leftrightarrow T_2)[w, x, y, z]}{\Gamma \vdash r^{-1} : (T_1 \leftrightarrow T_2)[x, w, z, y]} \quad (\text{Relational Subtype - Inverse}) \\
\frac{\Gamma \vdash r : (T_1 \leftrightarrow T_2)[a, b, c, d], t : (T_1 \leftrightarrow T_2)[w, x, y, z]}{\Gamma \vdash (r \oplus t) : (T_1 \leftrightarrow T_2)[w, x, y, z] \text{ Relational Subtype - Overriding}} \\
\frac{\Gamma \vdash r : (T_1 \leftrightarrow T_2)[a, b, c, d], t : (T_2 \leftrightarrow T_3)[w, x, y, z]}{\Gamma \vdash (r \circ t) : (T_1 \leftrightarrow T_3)[a, b, c, d] \text{ Relational Subtype - Composition}}
\end{array}$$

3.3.18 Built-in Functions

Includes: Common set-related functions.

$$\begin{array}{c}
\frac{\Gamma \vdash T[\text{with } order \in traits(T)]}{\Gamma \vdash min : set[T] \rightarrow T} \quad (\text{Built-in - Minimum}) \\
\frac{\Gamma \vdash T, int}{\Gamma \vdash min : set[T] \times (T \rightarrow int) \rightarrow T} \quad (\text{Built-in - Mapped Minimum}) \\
\frac{\Gamma \vdash T[\text{with } order \in traits(T)]}{\Gamma \vdash max : set[T] \rightarrow T} \quad (\text{Built-in - Maximum}) \\
\frac{\Gamma \vdash T, int}{\Gamma \vdash max : set[T] \times (T \rightarrow int) \rightarrow T} \quad (\text{Built-in - Mapped Maximum}) \\
\frac{\Gamma \vdash T}{\Gamma \vdash choice : set[T] \rightarrow T} \quad (\text{Built-in - Choice}) \\
\frac{\Gamma \vdash T}{\Gamma \vdash dom : T_1 \leftrightarrow T_2 \rightarrow set[T_1]} \quad (\text{Built-in - Domain}) \\
\frac{\Gamma \vdash T}{\Gamma \vdash ran : T_1 \leftrightarrow T_2 \rightarrow set[T_2]} \quad (\text{Built-in - Range}) \\
\frac{\Gamma \vdash T}{\Gamma \vdash card : set[T] \rightarrow nat} \quad (\text{Built-in - Cardinality})
\end{array}$$

$$\frac{\Gamma \vdash T}{\Gamma \vdash \text{size} : \text{bag}[T] \rightarrow \text{nat}} \quad (\text{Built-in - Bag Size})$$

$$\frac{\Gamma \vdash T \quad T \in \{\text{int}, \text{float}, \text{nat}\}}{\Gamma \vdash \text{sum} : \text{set}[T] \rightarrow T} \quad (\text{Built-in - Sum})$$

3.4 Definedness

3.5 Language Examples

4 Features and Behaviours

5 User-Facing Specification

5.1 (Abstract) Data Types

For the most part, Simile’s semantics resemble those of set theory and specification languages (specifically, Event-B). Aside from a few primitive types that abstract over bare set theory and programming language specific structures, like integer arithmetic and procedure definitions, all objects inherit properties from sets. For example, a sequence can be modelled as a set of pairs from natural numbers to the sequence element type.

In the forthcoming sections, we explicitly state the design of these set-theory-inspired types, useful properties, and underlying implementations.

5.1.1 Primitives

With consideration to execution efficiency and the extensive amount of existing optimization in modern languages, Simile makes use of the following basic types:

Booleans As the cornerstone of logical reasoning, booleans deserve an explicit type with reserved identifiers: $\mathbb{B} = \{\text{True}, \text{False}\}$. Further, Simile provides a healthy supply of common operations beyond \wedge and \vee : $\implies, \equiv, \Leftarrow, \forall, \exists$, etc.

Integers While integers are a very mature type in the world of computing, trade-offs between arbitrary-precision and fixed-point arithmetic may require testing and careful review.

If Simile chooses to represent numbers through conventional fixed point arithmetic, struggles with overflow and extremely large numbers are compounded by the innate non-determinism within Simile. For example, calculating the sum of a set of very large numbers can be done in arbitrary order, according to the semantics of mathematics (and Simile). But in fixed-point arithmetic, the resulting value can very well differ based on the order of operations. In other words, commutative math operations are no longer commutative under fixed precision. At some point in the computation pipeline, Simile will need to decide the order in which to evaluate the set, which may cause hard-to-spot bugs and unexpected runtime

failures.

Conversely, choosing arbitrary precision integers may bring Simile operators closer to mathematics at the cost of speed. The core features of Simile, while targeting formal methods users, attempt to efficiently execute set and relation operations. Adding checks for overflow on every mathematical operation may prove detrimental to performance.

The implementation decision essentially is one of context vs. safety. To benchmark Simile with these types, we gather a list of examples and judge whether the performance hit is worth the reliability benefits. Specific questions to examine: How often does overflow occur in these examples? Do numbers approach the upper limits of fixed precision often enough to warrant accommodating infrastructure? What is the performance of fixed vs. arbitrary precision calculations.

Floating Point Numbers In a similar vein to integers, floats are an imperfect translation of real numbers into efficiently calculable representations. Since floats (or even reals) are not used as often in specifications as integers, we leave the implementation as the conventional 64-bit representation. However, we note that non-determinism on these operations will prove especially tricky, since the result of commutative complements may differ without explicit over/underflow.

Strings Strings are syntactic sugar over Sequences of characters, represented internally as a sequence of (fixed precision) integers. Eventually, Simile aims to support built-in string operations as seen in many common programming languages.

5.1.2 Sets

Foundational Simile type: an unordered, unique collection of objects.

Sets may be used in two ways:

1. Enumerations by way of $\{x, y, z\}$ (read: The set of elements x, y , and z)
2. Comprehensions of the form $\{x \cdot x \in S \mid f(x)\}$ (read: The set of function f applied to all elements in S). A small effort is made to accommodate a shorthand form, wherein $\{x \cdot x \in S \mid x\}$ may be written as $\{x \cdot x \in S\}$, but this syntax easily confuses bound/unbound variables and is not recommended for complex comprehension expressions.

Sets, as mathematical objects with only two major restrictions, lend themselves well to a variety of computational implementations. Arrays, ordered arrays, hash maps, bit sets, roaring bit sets, path maps, and tries can all be used to satisfy the uniqueness property and hide an internally ordered representation from users.

No matter the underlying implementation, the optimizer may make use of syntactically-analyzed set sizes to direct the choice of lowered iterator. Further, recording the cardinality of a set at runtime may prove beneficial to the

performance of common cardinality-based specification expressions at a very small memory cost. Other properties based on element type, like max/min for numeric types, limits for numeric ranges, maximum sizes, may be useful.

5.1.3 Relations

Relations are sets that make use of paired/maplet elements (i.e., they define a relationship between two sets). Like sets, implementations may vary greatly, with the potential to increase efficiency of relation-specific operations.

Like sets, relations may benefit from compile-time computed properties: size of the domain, size of the range, totality, many-to-oneness, $O(1)$ access to the domain and range, and access to associated domain/range properties.

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 3: 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)

Subtype 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:

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

5.1.4 Bags

Bags (or Multisets) are a refinement of relations wherein, for a bag with element type T , the bag is a relational function from $T \rightarrow \mathbb{Z}^+$, denoting the count of items in a collection. Once the count of an item reaches 0, the item is removed from the bag.

Bags must be many-to-one and carry specialized behaviours for conventional relation operations. For example, relational image ($R[B]$) with a bag argument will now return a bag, where the number of occurrences of the resulting item replaces what would otherwise be a regular set of those items.

Additional fields to record the size of the bag (sum of the range) may prove useful.

5.1.5 Sequences

Sequences are yet another refinement of relations where a sequence with element type T represents a relational function from $\mathbb{N} \rightarrow T$ that is total on some defined bounds.

The dense nature of sequence domains, along with knowledge of sequence sizes, may open up further data type refinements.

5.2 Data Type Operations

5.2.1 Primitives

5.2.2 Sets

5.2.3 Relations

5.2.4 Bags

5.2.5 Sequences

6 Implementation Specification

6.1 Abstract Data Types

6.2 Atomic Operations

6.3 Optimization Frontend

6.4 Optimization Backends

7 AST Lowering and Optimization 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.

7.1 Syntactic Sugar for Bags

$$\begin{aligned} R[B] &\rightsquigarrow R^{-1} \circ B && \text{(Bag Image)} \\ S \cup T &\rightsquigarrow \llbracket x \cdot x \in \text{dom}(S) \cup \text{dom}(T) \wedge n = \max(S[x] \cup T[x]) \mid x \mapsto n \rrbracket && \text{(Bag Predicates - Union)} \\ S \cap T &\rightsquigarrow \llbracket x \cdot x \in \text{dom}(S) \cap \text{dom}(T) \wedge n = \min(S[x] \cup T[x]) \wedge n > 0 \mid x \mapsto n \rrbracket && \text{(Bag Predicates - Intersection)} \\ S + T &\rightsquigarrow \llbracket x \cdot x \in \text{dom}(S) \cup \text{dom}(T) \wedge n = \text{sum}(S[x] \cup T[x]) \mid x \mapsto n \rrbracket && \text{(Bag Predicates - Sum)} \\ S - T &\rightsquigarrow \llbracket x \cdot x \in \text{dom}(S) \wedge \text{pop}(S[x]) - \text{popDefault}(T[x], 0) \wedge n > 0 \mid x \mapsto n \rrbracket && \text{(Bag Predicates - Difference)} \end{aligned}$$

7.2 Syntactic Sugar for Sequences

$$\begin{aligned}
L + M &\rightsquigarrow L \cup \{x \mapsto y \cdot x \mapsto y \in M \mid x + \max(\text{dom}(L)) \mapsto y\} && (\text{Concatenation}) \\
\text{flatten}(X) &\rightsquigarrow \text{foldL}(++, [], X) && (\text{Flatten}) \\
\text{foldL}(f, e, X) &\rightsquigarrow \text{foldL}(f, f(e, \text{first}(X)), \text{tail}(X)) && (\text{Fold Left Associative}) \\
\text{foldL}(f, e, []) &\rightsquigarrow e && (\text{Fold Left Associative - Empty}) \\
\text{first}(X) &\rightsquigarrow X(0) && (\text{First}) \\
\text{tail}(X) &\rightsquigarrow \{i \mapsto x \cdot i \mapsto x \in X \wedge i \neq 0 \mid i - 1 \mapsto x\} && (\text{Tail}) \\
\text{append}(X, y) &\rightsquigarrow X \cup \{\max(\text{dom}(X)) \mapsto y\} && (\text{Append}) \\
\text{prepend}(X, y) &\rightsquigarrow \{i \mapsto x \cdot i \mapsto x \in X \mid i + 1 \mapsto x\} \cup 0 \mapsto y && (\text{Prepend})
\end{aligned}$$

7.3 Builtin Functions

$$\begin{aligned}
f(x) := E &\rightsquigarrow f := f \oplus \{x \mapsto E\} && (\text{Functional Override}) \\
\text{card}(S) &\rightsquigarrow \sum x \cdot x \in S \mid 1 && (\text{Cardinality}) \\
\text{dom}(R) &\rightsquigarrow \{x \mapsto y \cdot x \mapsto y \in R \mid x\} && (\text{Domain}) \\
\text{ran}(R) &\rightsquigarrow \{x \mapsto y \cdot x \mapsto y \in R \mid y\} && (\text{Range}) \\
R \oplus Q &\rightsquigarrow Q \cup (\text{dom}(Q) \triangleleft R) && (\text{Override}) \\
S \triangleleft R &\rightsquigarrow \{x \mapsto y \cdot x \mapsto y \in R \wedge x \in S \mid x \mapsto y\} && (\text{Domain Restriction}) \\
S \triangleleft R &\rightsquigarrow \{x \mapsto y \cdot x \mapsto y \in R \wedge x \notin S \mid x \mapsto y\} && (\text{Domain Subtraction}) \\
R \triangleright S &\rightsquigarrow \{x \mapsto y \cdot x \mapsto y \in R \wedge y \in S \mid x \mapsto y\} && (\text{Range Restriction}) \\
R \triangleright S &\rightsquigarrow \{x \mapsto y \cdot x \mapsto y \in R \wedge y \notin S \mid x \mapsto y\} && (\text{Range Subtraction}) \\
\text{size}(S) &\rightsquigarrow \sum x \mapsto n \in S \cdot n && (\text{Bag Size})
\end{aligned}$$

7.4 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.

$R(x) \rightsquigarrow \text{pop}(R[x])$	(Functional Image ^a)
$R[S] \rightsquigarrow \{x \mapsto y \cdot x \mapsto y \in R \wedge x \in S \mid y\}$	(Image)
$x \mapsto y \in S \times T \rightsquigarrow x \in S \wedge y \in T$	(Product)
$x \mapsto y \in R^{-1} \rightsquigarrow y \mapsto x \in R$	(Inverse)
$x \mapsto y \in (Q \circ R) \rightsquigarrow x \mapsto z \in Q \wedge z' \mapsto y \in R \wedge z = z'$	(Composition)

^aMay fail if R is not total.

$S \cup T \rightsquigarrow \{x \cdot x \in S \vee x \in T\}$	(Predicate Operations - Union)
$S \cap T \rightsquigarrow \{x \cdot x \in S \wedge x \in T\}$	(Predicate Operations - Intersection)
$S \setminus T \rightsquigarrow \{x \cdot x \in S \wedge x \notin T\}$	(Predicate Operations - Difference)
$x \in \oplus(E \mid P) \rightsquigarrow P \wedge x = E$	(Membership Collapse ^a)

^aRule 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).

7.5 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.

$x_1 \wedge \dots \wedge (x_i \wedge x_{i+1}) \wedge \dots \rightsquigarrow x_1 \wedge \dots \wedge x_i \wedge x_{i+1} \wedge \dots$	(Flatten Nested \wedge)
$x_1 \vee \dots \vee (x_i \vee x_{i+1}) \vee \dots \rightsquigarrow x_1 \vee \dots \vee x_i \vee x_{i+1} \vee \dots$	(Flatten Nested \vee)
$\neg\neg x \rightsquigarrow x$	(Double Negation)
$\neg(x \vee y) \rightsquigarrow \neg x \wedge \neg y$	(Distribute De Morgan - Or)
$\neg(x \wedge y) \rightsquigarrow \neg x \vee \neg y$	(Distribute De Morgan - And)
$x \wedge (y \vee z) \rightsquigarrow (x \wedge y) \vee (x \wedge z)$	(Distribute \wedge over \vee)

7.6 Or-wrapping

Intuition Restructuring quantifier predicates with top-level ORs for later rules.

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.

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

^aTo 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.

7.7 Generator Selection

Intuition All nested ANDs are placed in a tree structure to better suit loop lowering.

Post-condition

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

Notes All matches of this phase occur only inside quantifier predicates. GSP is a tuple of form (generator chain, predicates) that can be flattened to a conjunctive clause. A CombGSP is a triple of form (shared generator, disjunctive children predicates/generators, predicates) that can likewise be flattened into a conjunctive clause. We keep the structure of these tuples distinct to ease the rewriting process.

Brainstorming selection heuristics Conditions to consider:

- Prefer composition 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?

$$\bigwedge P_i \rightsquigarrow GSP\left(\bigwedge P_{g_i}, \bigwedge_{P_i \notin P_g} P_i\right)$$

(GSP Wrapping ^a)

$$GSP(x \wedge xs, P) \vee GSP(x \wedge ys, Q) \rightsquigarrow \text{CombGSP}(x, GSP(xs, P) \vee GSP(ys, Q))$$

(Nested Generator Selection ^b)

^aGSP is a tuple-like structure, where all entries can be flattened into an AND. The LH term must occur inside a quantifier's predicate - one generator per top-level or-clause. Chained generators may be necessary to allow for nested loop operations (ie. non-top-level clauses may require generators). P_g is a list of chained generators, specific set membership clauses (of form $x \in S$) distinguished from the rest of $\bigwedge P_i$. Currently, the selection of P_g is arbitrary (and thus the rewrite system is not confluent), but heuristics may be added later to choose better generators.

^b*free* is a function that returns true when variables in the predicate are defined (either bound by the current generator or bound prior to the current statement)

7.8 GSP to Loop

Intuition Start lowering expressions into imperative-like loops.

Post-condition

- All quantifiers are transformed into *loop* structures.
- *loop* predicates are of the form $x \in S \wedge \bigwedge P_i$.

- All *loop* predicates have an assigned generator of the form $x \in S$.

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

${}^a\oplus$ 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.

$$\begin{array}{c} \text{loop } P_0 \\ \text{body} \\ \text{loop } \bigvee P_i : \\ \text{body} \rightsquigarrow \begin{array}{c} \text{loop } P_1 \wedge \neg P_0 \\ \text{body} \\ \text{loop } P_2 \wedge \neg \bigvee_{i < 2} P_i \\ \text{body} \end{array} \end{array} \quad (\text{Top-level Or-Loop})$$

$$\begin{array}{c} \text{loop } GSP(g \wedge gs, P) \\ \text{body} \rightsquigarrow \begin{array}{c} \text{loop} \\ \rightsquigarrow \begin{array}{c} \text{SingleGSP}(g, \text{free}(P)) \\ \text{loop } GSP(gs, \text{bound}(P)) \\ \text{body} \end{array} \end{array} \end{array} \quad (\text{Chained GSP Loop } {}^a)$$

${}^a\text{free}$ considers defined variables at this point in time.

$$\begin{array}{c} \text{loop } GSP([], P) \\ \text{body} \rightsquigarrow \begin{array}{c} \text{if } P \text{ then} \\ \text{body} \end{array} \end{array} \quad (\text{Empty GSP Loop})$$

$$\begin{array}{c} \text{loop } SingleGSP(g, P) \\ \text{loop } gs_0 \\ \text{body} \\ \text{loop } CombGSP(g, gs, P) \\ \text{body} \rightsquigarrow \begin{array}{c} \dots \\ \text{loop } gs_1 \\ \text{body} \\ \dots \end{array} \end{array} \quad (\text{Combined GSP Loop})$$

7.9 Relational Subtyping Loop Simplification

Intuition Simplify unnecessary iteration structures into direct accesses.

Post-condition

- All quantifiers are transformed into *loop* structures.
- *loop* predicates are of the form $x \in S \wedge \bigwedge P_i$.

- All *loop* predicates have an assigned generator of the form $x \in S$.

No change from previous, but the number of loop structures should not increase.

$$\begin{array}{ll} x \in \text{dom}(R) \rightsquigarrow \text{True} & (\text{Total membership elimination } {}^a) \\ x \in \text{ran}(R) \rightsquigarrow \text{True} & (\text{Surjective membership elimination } {}^b) \end{array}$$

^a R is total and x satisfies the type of R 's domain. Currently unused since it may eliminate some generators - move higher in the list?

^b R is surjective and x satisfies the type of R 's range. Currently unused since it may eliminate some generators - move higher in the list? Or

$$\begin{array}{ll} \text{loop } \text{SingleGSP}(x \mapsto y \in R, x = a \wedge P) \rightsquigarrow & \begin{array}{l} \text{if } a \in \text{dom}(R) \text{ then} \\ \quad \text{loop } \text{SingleGSP}(y \in R[a], P[x := a]) \\ \quad \text{body} \\ \quad (\text{Concrete Domain Image}) \end{array} \\ \text{body} & \\ \\ \text{loop } \text{SingleGSP}(x \mapsto y \in R, y = a \wedge P) \rightsquigarrow & \begin{array}{l} \text{if } a \in \text{ran}(R) \text{ then} \\ \quad \text{loop } \text{SingleGSP}(x \in R^{-1}[a], P[y := a]) \\ \quad \text{body} \\ \quad (\text{Concrete Range Image } {}^a) \end{array} \\ \text{body} & \end{array}$$

^aRequires efficient lookups for *ran*, inverse, and image

$$\begin{array}{ll} \text{loop } \text{SingleGSP}(y \in R[x], P) \rightsquigarrow & \begin{array}{l} \text{if } P[y := R(a)] \text{ then} \\ \quad \text{body}[y := R(a)] \end{array} \\ \text{body} & \\ \text{where } \text{free}(x) \text{ and } R \text{ is many-to-one} & \\ & (\text{Single Element Loop}) \end{array}$$

$$\begin{array}{ll} \text{loop } \text{SingleGSP}(x \in \{y\}, P) \rightsquigarrow & \begin{array}{l} \text{if } P[x := y] \text{ then} \\ \quad \text{body}[x := y] \end{array} \\ \text{body} & \\ & (\text{Singleton Membership Elimination}) \end{array}$$

7.10 Loop Code Generation

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.

$$\begin{array}{ccc}
 \text{loop } SingleGSP(P_g, \bigwedge P_i) & \rightsquigarrow & \begin{array}{c} \text{if } \bigwedge_{free(P_i)} P_i \text{ then} \\ \text{for } P_g \text{ do} \\ \quad \text{if } \bigwedge_{bound(P_i)} P_i \text{ then} \\ \quad \quad \text{body} \\ \quad \text{(Conjunct Conditional } ^a\text{)} \end{array}
 \end{array}$$

^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.

7.11 Replace and Simplify

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.

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

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

$$\begin{array}{ll}
 P \wedge True \rightsquigarrow P & (\text{Simplify And-true}) \\
 P \wedge False \rightsquigarrow False & (\text{Simplify And-false}) \\
 P \vee True \rightsquigarrow True & (\text{Simplify Or-true}) \\
 P \vee False \rightsquigarrow P & (\text{Simplify Or-false}) \\
 Statement(Statement(x)) \rightsquigarrow Statement(x) & (\text{Flatten Nested Statements})
 \end{array}$$

References

- [1] Christophe Métayer and Laurent Voisin. “The Event-B Mathematical Language”. 2007.