

# **Modular and Polymorphic Set-Based Analysis: Theory and Practice**

*Cormac Flanagan*  
*Matthias Felleisen*

**Rice COMP TR96-266**

**November 1996**

Department of Computer Science  
Rice University  
P.O. Box 1892  
Houston, TX 77251-1892

Copyright ©1996 by  
Cormac Flanagan and Matthias Felleisen

# Modular and Polymorphic Set-Based Analysis: Theory and Practice

(revision 1.4)

Cormac Flanagan      Matthias Felleisen

Department of Computer Science,  
Rice University,  
Houston, TX 77251-1892

Contact email: `cormac@cs.rice.edu`

## Abstract

Set-based analysis (SBA) produces good predictions about the behavior of functional and object-oriented programs. The analysis proceeds by inferring *constraints* that characterize the data flow relationships of the analyzed program. Experiences with Rice's program development environment, which includes a static debugger based on SBA, indicate that SBA can deal with programs of up to a couple of thousand lines of code. However, SBA does not cope with larger programs because it generates large systems of constraints for these programs. These constraint systems are at least linear, and possibly quadratic, in the size of the analyzed program.

This paper presents theoretical and practical results concerning methods for reducing the size of constraint systems. The theoretical results include a complete proof-theoretic characterization of the observable behavior of a constraint system, which we use to establish a close connection between the observable equivalence of constraint systems and the equivalence of regular-tree grammars. We then exploit this connection to adapt a variety of algorithms for simplifying grammars to the practical problem of simplifying constraint systems. Based on the resulting algorithms, we develop a *componential set-based analysis*, a modular and polymorphic variant of SBA. Experimental results validate the practicality of the simplification algorithms and the analysis. The simplified constraint systems are typically at least an order of magnitude smaller than the original systems, which significantly reduces both the time and space required by the analysis.

# Contents

<b>1</b>	<b>The Effectiveness of Set-Based Analysis</b>	<b>1</b>
<b>2</b>	<b>The Source Language</b>	<b>2</b>
<b>3</b>	<b>Set-Based Analysis</b>	<b>3</b>
3.1	The Source Language . . . . .	4
3.1.1	Syntax . . . . .	4
3.1.2	Semantics . . . . .	4
3.2	The Constraint Language . . . . .	5
3.3	The Meaning of Set Constraints . . . . .	6
3.3.1	The Semantic Domain . . . . .	6
3.3.2	The Semantics of Constraints . . . . .	7
3.4	Deriving Constraints . . . . .	8
3.5	Soundness of the Derived Constraints . . . . .	10
3.6	Solving Set Constraints . . . . .	11
3.6.1	Computing the Least Solution . . . . .	13
<b>4</b>	<b>Constraint Simplification</b>	<b>13</b>
4.1	Conditions for Constraint Simplification . . . . .	14
4.2	The Proof Theory of Observable Equivalence . . . . .	15
4.3	Deciding Observable Equivalence . . . . .	18
4.4	Regular Grammars . . . . .	18
4.5	Regular Tree Grammars . . . . .	20
4.6	Staging . . . . .	21
4.7	The Entailment Algorithm . . . . .	21
4.8	Practical Constraint System Simplification . . . . .	23
4.8.1	Empty Constraint Simplification . . . . .	24
4.8.2	Unreachable Constraint Simplification . . . . .	24
4.8.3	Removing $\epsilon$ -Constraints . . . . .	25
4.8.4	Hopcroft's Algorithm . . . . .	25
4.9	Simplification Benchmarks . . . . .	26
<b>5</b>	<b>Componential Set-Based Analysis</b>	<b>27</b>
5.1	Componential Set-Based Analysis . . . . .	27
5.2	Experimental Results . . . . .	28
5.3	User Interface for Multi-File Programs . . . . .	30
5.4	Efficient Polymorphic Analysis . . . . .	30
<b>6</b>	<b>Competitive Work</b>	<b>31</b>
<b>7</b>	<b>Future Work</b>	<b>32</b>

<b>A</b>	<b>Proofs</b>	<b>32</b>
A.1	Subject Reduction Proof . . . . .	32
A.2	Proofs for Computing Set-Based Analysis . . . . .	39
A.3	Proofs for Conditions for Constraint Simplification . . . . .	42
A.4	Proofs for Proof Theory of Observable Equivalence . . . . .	43
A.5	Proofs for Deciding Observable Equivalence . . . . .	49
A.6	Correctness of the Entailment Algorithm . . . . .	52
A.7	Correctness of the Hopcroft Algorithm . . . . .	55
<b>B</b>	<b>Notations</b>	<b>56</b>

# 1 The Effectiveness of Set-Based Analysis

Rice’s Scheme program development environment provides a static debugger, MrSpidey, which statically analyzes a program and, using the results of this analysis, checks the soundness of all primitive operations [8]. If an operation may fault due to a violation of its invariant, MrSpidey highlights the program operation so that the programmer can investigate the potential fault site before running the program. Using the graphical explanation facilities of MrSpidey, the programmer can determine whether this fault will really happen or whether the corresponding correctness proof is beyond the analysis’s capabilities.

MrSpidey’s program analysis is a constraint-based system similar to Heintze’s set-based analysis of ML programs [10]. The analysis consists of two co-mingled phases: a *specification* phase, during which MrSpidey derives constraints describing the data flow relationships, and a *solution* phase, during which MrSpidey solves the constraints. The solution conservatively approximates the set of values that may result from each program expression.

In practice, MrSpidey has proven highly effective for pedagogic programming, which includes programs of several hundreds to 2,000 lines of code. It also works reasonably well on some programs of up to several thousand lines in length. However, it becomes less useful for debugging large programs for two reasons:

- Set-based analysis has an  $O(n^3)$  worst-case time bound. Although the constant on the cubic element in the polynomial is small, it is noticeable for programs of several thousand lines.
- Large programming projects tend to re-use functions in a polymorphic fashion. To avoid merging information between unrelated calls to such functions, the analysis must duplicate the constraints for each corresponding call site. This duplication is expensive because the size of the constraint set is at best linear, and possibly quadratic, in the size of the function.

A closer look at these two obstacles quickly reveals that the major limitation of set-based analysis is the size of the constraint system describing the data flow relationships of a program. If we could reduce the size of constraint systems without affecting the solution space that they denote, we could simplify constraint sets at intermediate stages during the specification phase and thus reduce the analysis time. In particular, after simplifying the constraint set of a polymorphic function, MrSpidey would duplicate a much smaller constraint set at each polymorphic reference than the plain analysis; similarly, producing a simplified constraint set for a module would substantially reduce the cost of solving the combined set of constraints for a modularized program.

The simplification of constraint sets raises both interesting theoretical and practical questions. On the theoretical side, we need to ensure that simplification preserves the observable behavior of a constraint system. In this paper, we provide a complete characterization of observable behavior and, in the course of this development, establish a close connection between this observable equivalence of constraint systems and the equivalence of regular tree grammars (RTGs).<sup>1</sup> Exploiting this connection, we develop a complete al-

---

<sup>1</sup>A number of researchers, including Reynolds [18], Jones and Muchnick [14], Heintze [10], Aiken [2], and Cousot and Cousot [3] previously exploited the relationship between RTGs and the *least solution* of a

gorithm for deciding the equivalence of constraint systems. Unfortunately, the problem is PSPACE-hard.

Fortunately, a minimized constraint set is only optimal but not necessary for practical purposes. The practical question concerns finding approximate algorithms for simplifying constraint sets that would make MrSpidey more useful. To answer this question, we exploit the correspondence between the minimization problem for RTGs and constraint sets to adapt a variety of algorithms for simplifying RTG to the problem of simplifying constraint systems. Based on these simplification algorithms, we develop a *componential*,<sup>2</sup> or component-wise, variant of set-based analysis. Experimental results verify the effectiveness of these algorithms and the corresponding flavors of the analysis. The simplified constraint systems are typically at least an order of magnitude smaller than the original constraint set. Moreover, these reductions in size result in significant gains in the speed of the analysis.

We expect that some of our theoretical and practical results as well as the techniques will carry over to other constraint-based systems, such as the conditional type system of Aiken *et al.* [2], Eifrig *et al.*'s object-oriented type system [5], or Pottier's or Smith *et al.*'s subtyping simplification algorithms [17, 20].

The presentation of our results proceeds as follows. Section 2 describes an idealized source language. Section 3 describes our set-based analysis. Section 4 characterizes the relationship between constraint systems and RTGs, and section 5 exploits this connection to derive a number of practical constraint simplification algorithms. Sections 6 and 7 discuss how these algorithms perform in a realistic program analysis system. Section 8 discusses related work, and section 9 describes directions for future research. Appendix A contains proof of various theorems and lemmas. Appendix B contains an index of notations used in the paper.

## 2 The Source Language

---

Syntax:

$M \in \Lambda$	$= x \mid V \mid (M \ M) \mid (\mathbf{let} \ (x \ M) \ M) \mid M^l$	(Expressions)
$V \in \text{Value}$	$= b \mid (\lambda^t x. M)$	(Values)
$x \in \text{Vars}$	$= \{x, y, z, \dots\}$	(Variables)
$b \in \text{BasicConst}$		(Basic Constants)
$t \in \text{Tag}$		(Function Tags)
$l \in \text{Label}$		(Expression Labels)

Figure 1: The source language  $\Lambda$

---

constraint system. We present an additional result, namely a connection between RTGs and the observable behavior (*i.e.*, the *entire solution space*) of constraint systems.

<sup>2</sup>**componential** *a.* of or pertaining to components; *spec. (Ling.)* designating the analysis of distinctive sound units or grammatical elements into phonetic or semantic components (*New Shorter Oxford English Dictionary*, Clarendon Press, 1993)

For simplicity, we present our results for a  $\lambda$ -calculus-like language with constants and labeled expressions. It is straightforward to extend the analysis to a realistic language including assignments, recursive data structures, objects and modules along the lines described in a related report [?].

Expressions in the language are either variables, values, function applications, **let**-expressions, or labeled expressions. We use labels to identify those program expressions whose values we wish to predict. Values include basic constants and functions. Functions have identifying tags so that MrSpidey can reconstruct a call-graph from the results of the analysis. We use **let**-expressions to introduce polymorphic bindings, and hence restrict these bindings to syntactic values [22]. We work with the usual conventions and terminology of the  $\lambda_v$ -calculus when discussing syntactic issues. In particular, the substitution operation  $M[x \leftarrow V]$  replaces all free occurrences of  $x$  within  $M$  by  $V$ , and  $\Lambda^0$  denotes the set of closed terms, also called *programs*.

We specify the meaning of programs based upon the *notions of reduction*  $\beta_v$ ,  $\beta_{let}$  and *unlabel*:

$$\begin{array}{lll} ((\lambda^t x.M) V) & \longrightarrow & M[x \mapsto V] & (\beta_v) \\ (\mathbf{let} (x V) M) & \longrightarrow & M[x \mapsto V] & (\beta_{let}) \\ V^l & \longrightarrow & V & (unlabel) \end{array}$$

The  $\beta_v$  and  $\beta_{let}$  rules are the conventional rules for the  $\lambda$ -calculus. The *unlabel* rule simply removes the label from an expression once its value is needed.

An *evaluation context*  $\mathcal{E}$  is an expression containing a hole  $[]$  in place of the next subterm to be evaluated:

$$\mathcal{E} = [] \mid (\mathcal{E} M) \mid (V \mathcal{E}) \mid (\mathbf{let} (x \mathcal{E}) M) \mid \mathcal{E}^l$$

For example, in the term  $(N M)$ , the next expression to be evaluated lies within  $N$ , and thus the definition of evaluation contexts includes the clause  $(\mathcal{E} M)$ . An evaluation context always contains a single hole  $[]$ , and we use the notation  $\mathcal{E}[M]$  to denote term produced by filling the hole in  $\mathcal{E}$  with the term  $M$ .

The standard reduction relation  $\mapsto$  is the compatible closure of  $\longrightarrow$  with respect to evaluation contexts:

$$\mathcal{E}[M] \mapsto \mathcal{E}[N] \quad \Longleftrightarrow \quad M \longrightarrow N$$

The relation  $\mapsto^*$  is the reflexive, transitive closure of  $\mapsto$ . The semantics of the language is defined via the partial function *eval* on programs:

$$\begin{array}{lll} eval : \Lambda^0 & \longrightarrow_p & Value \\ eval(M) & = & V \quad \text{if } M \mapsto^* V \end{array}$$

### 3 Set-Based Analysis

Conceptually, set-based analysis consists of two phases: a *specification* phase and a *solution* phase.<sup>3</sup> During the specification phase, the analysis tool derives *constraints* on the sets of values that program expressions may assume. These constraints describe the data flow

---

<sup>3</sup>Cousot and Cousot showed that set-based analysis can alternatively be formulated as an abstract interpretation computed by chaotic iteration [3].



relationships of the analyzed program. During the solution phase, the analysis produces finite descriptions of the potentially infinite sets of values that satisfy these constraints. The result provides an approximate set of values for each labeled expression in the program.

### 3.1 The Source Language

We develop the analysis for an idealized,  $\lambda$ -calculus-like language  $\Lambda$  with constants and labeled expressions. This section introduces the syntax and semantics of  $\Lambda$ .

#### 3.1.1 Syntax

---

Syntax:

$M \in \Lambda$	$= x \mid V \mid (M\ M) \mid (\mathbf{let}\ (x\ V)\ M) \mid M^l$	(Expressions)
$V \in \text{Value}$	$= b \mid (\lambda^t x.M)$	(Values)
$x \in \text{Var}$	$= \{x, y, z, \dots\}$	(Variables)
$b \in \text{Basic Const}$		(Basic Constants)
$t \in \text{FnTag}$		(Function Tags)
$l \in \text{Label}$		(Expression Labels)

Figure 2: The source language  $\Lambda$

---

Expressions in the language are either variables, values, applications, **let**-expressions, or labeled expressions: see figure 2. Values include basic constants and functions. Each function has an identifying tag so that MrSpidey can reconstruct the textual source of function values from the results of the analysis. We use **let**-expressions to introduce polymorphic bindings, and hence restrict these bindings to syntactic values [22]. We use labels to identify those program expressions whose values we wish to predict.

We work with the usual conventions and terminology of the  $\lambda_v$ -calculus when discussing syntactic issues. In particular, the substitution operation  $M[x \leftarrow V]$  replaces all free occurrences of  $x$  within  $M$  by  $V$ , and  $\Lambda^0$  denotes the set of closed terms, also called *programs*.

#### 3.1.2 Semantics

We specify the meaning of programs based upon three *notions of reduction*:

$$\begin{array}{lll}
 ((\lambda^t x.M)\ V) & \longrightarrow & M[x \mapsto V] & (\beta_v) \\
 (\mathbf{let}\ (x\ V)\ M) & \longrightarrow & M[x \mapsto V] & (\beta_{let}) \\
 V^l & \longrightarrow & V & (unlabel)
 \end{array}$$

The  $\beta_v$  and  $\beta_{let}$  rules are the conventional rules for the  $\lambda$ -calculus. The *unlabel* rule simply removes the label from an expression once its value is needed.

An *evaluation context*  $\mathcal{E}$  is an expression containing a hole  $[ ]$  in place of the next sub-term to be evaluated:

$$\mathcal{E} = [ ] \mid (\mathcal{E}\ M) \mid (V\ \mathcal{E}) \mid \mathcal{E}^l$$

For example, in the term  $(N \ M)$ , the next expression to be evaluated lies within  $N$ , and thus the definition of evaluation contexts includes the clause  $(\mathcal{E} \ M)$ . An evaluation context always contains a single hole  $[\ ]$ , and we use the notation  $\mathcal{E}[ \ M \ ]$  to denote term produced by filling the hole in  $\mathcal{E}$  with the term  $M$ .

The *standard reduction relation*  $\mapsto$  is the compatible closure of  $\longrightarrow$  with respect to evaluation contexts:

$$\mathcal{E}[ \ M \ ] \mapsto \mathcal{E}[ \ N \ ] \quad \text{iff} \quad M \longrightarrow N$$

The relation  $\mapsto^*$  is the reflexive, transitive closure of  $\mapsto$ . The semantics of the language is defined via the partial function *eval* on programs:

$$\begin{aligned} eval : \Lambda^0 &\longrightarrow_p Value \\ eval(M) &= V \quad \text{if } M \mapsto^* V \end{aligned}$$

### 3.2 The Constraint Language

To simplify the later derivation of the constraint simplification algorithms (see chapter ??), we express the constraint language in terms of selectors, instead of the more usual constructors. Specifically, a *set expression*  $\tau$  is either a set variable; a constant; or one of the “selector” expressions  $\mathbf{dom}(\tau)$  or  $\mathbf{rng}(\tau)$ :

$$\begin{aligned} \tau &\in SetExp = \alpha \mid c \mid \mathbf{dom}(\tau) \mid \mathbf{rng}(\tau) \\ \alpha, \beta &\in SetVar \supset Label \\ c &\in Const = BasicConst \cup FnTag \end{aligned}$$

By using selector expressions, we can specify each “quantum” of the program’s data-flow behavior independently; using constructors would combine several of these quanta into one constraint. For example, we specify a function’s behavior via the two constraints  $\{\mathbf{dom}(\alpha) \leq \alpha_1, \alpha_2 \leq \mathbf{rng}(\alpha)\}$  instead of the combined constraint  $\{(\alpha_1 \rightarrow \alpha_2) \leq \alpha\}$ .

The meta-variables  $\alpha, \beta, \gamma$  range over set variables, and we include program labels in the collection of set variables. Constants include both basic constants and function tags. A *constraint*  $C \in Constraint$  is an inequality  $\tau_1 \leq \tau_2$  relating two set expressions.

$$C \in Constraint = \tau_1 \leq \tau_2$$

We sometimes enclose constraints inside square brackets for clarity:  $[\tau_1 \leq \tau_2]$ . A *constraint system*  $S \in ConstraintSystem$  is a collection of constraints.

$$S \in ConstraintSystem = \mathcal{P}_{\text{fin}}(Constraint)$$

We use  $SetVar(S)$  to denote the collection of set variables in a constraint system  $S$ . In some cases, the relevant constraints in a constraint system are those that only mention certain set variables. The *restriction* of a constraint system to a collection of set variables  $E$  is:

$$S \mid_E = \{C \in S \mid C \text{ only mentions set variables in } E\}$$

### 3.3 The Meaning of Set Constraints

Intuitively, each set expression  $\tau$  denotes a set of run-time values, and each constraint  $\tau_1 \leq \tau_2$  denotes a corresponding set containment relationship. We formalize the meaning of set constraints by mapping syntactic set expressions onto a semantic domain. The next subsection describes the precise structure of the semantic domain, and the second subsection describes the mapping from set expressions to that domain.

#### 3.3.1 The Semantic Domain

A set expression denotes a collection of values. For our sample language, the collection consists of basic constants and functions and is therefore best represented as a triple  $X = \langle C, D, R \rangle$ . The first component  $C \in \mathcal{P}(\text{Const})$ <sup>4</sup> is a set of basic constants and function tags. The second and third components of  $X$  denote the possible arguments and results of functions in  $X$ , respectively. Since these two components also denote value sets, the appropriate model for set expressions is the solution of the equation:

$$\mathcal{D} = \mathcal{P}(\text{Const}) \times \mathcal{D} \times \mathcal{D}$$

The solution  $\mathcal{D}$  is equivalent to the set of all infinite binary trees<sup>5</sup> with each node labeled with an element of  $\mathcal{P}(\text{Const})$ . This set can be formally defined as the set of total functions  $f : \{\text{dom}, \text{rng}\}^* \rightarrow \mathcal{P}(\text{Const})$ , and the rest of the development can be adapted *mutandis mutatis* [16]. For clarity, we present our results using the more intuitive equational definition instead.

We use the functions  $\text{const} : \mathcal{D} \rightarrow \mathcal{P}(\text{Const})$  and  $\text{dom}, \text{rng} : \mathcal{D} \rightarrow \mathcal{D}$  to extract the respective components of an element of  $\mathcal{D}$ .

Each element of  $\mathcal{D}$  represents a set of run-time values (relative to a given program) according to the set of basic constants and function tags in its first component. The set of values represented by an element  $X \in \mathcal{D}$  is defined through the relation  $V$  in  $X$ :

$$\begin{aligned} b \text{ in } \langle C, D, R \rangle &\text{ iff } b \in C \\ (\lambda^t x. M) \text{ in } \langle C, D, R \rangle &\text{ iff } t \in C \end{aligned}$$

We order the elements of  $\mathcal{D}$  according to a relation that is anti-monotonic in the domain position:

$$\langle C_1, D_1, R_1 \rangle \sqsubseteq \langle C_2, D_2, R_2 \rangle \text{ iff } C_1 \subseteq C_2, D_2 \subseteq D_1, R_1 \subseteq R_2$$

This ordering is anti-monotonic in the domain position because information about argument values at an application needs to flow *backward* along data-flow paths to the formal parameter of the corresponding function definitions. To illustrate this idea, consider a program that binds a function  $f$  to a program variable  $g$ . This behavior is described in the semantic domain as the inequality  $X_f \sqsubseteq X_g$ , where  $X_f$  and  $X_g$  describe the values sets for  $f$  and  $g$  respectively. Since the possible argument set for  $f$  must contain all values to which  $g$  is

---

<sup>4</sup> $\mathcal{P}$  denotes the power-set constructor.

<sup>5</sup>We can analyze languages with additional data structures by extending  $\mathcal{D}$  to infinite  $n$ -ary trees, where  $n$  is the number of selectors (*e.g.*, `dom`, `rng`) corresponding to the extended language.

applied, the inequality  $\text{dom}(X_g) \sqsubseteq \text{dom}(X_f)$  must also hold. Thus the domain  $\mathcal{D}$  should satisfy the inference rule:

$$\frac{X_f \sqsubseteq X_g}{\text{dom}(X_g) \sqsubseteq \text{dom}(X_f)}$$

which is why the ordering  $\sqsubseteq$  needs to be anti-monotonic in the domain element.

Under the defined ordering, the set  $\mathcal{D}$  forms a complete lattice; the top and bottom elements are the solutions to the equations:

$$\begin{aligned} \top &= \langle \text{Const}, \perp, \top \rangle \\ \perp &= \langle \emptyset, \top, \perp \rangle \end{aligned}$$

respectively. The least upper bound and greatest lower bound operations are recursively defined as:

$$\begin{aligned} \langle C_1, D_1, R_1 \rangle \sqcup \langle C_2, D_2, R_2 \rangle &= \langle C_1 \cup C_2, D_1 \sqcap D_2, R_1 \sqcup R_2 \rangle \\ \langle C_1, D_1, R_1 \rangle \sqcap \langle C_2, D_2, R_2 \rangle &= \langle C_1 \cap C_2, D_1 \sqcup D_2, R_1 \sqcap R_2 \rangle \end{aligned}$$

### 3.3.2 The Semantics of Constraints

The semantics of a constraint system is defined with respect to a *set environment*  $\rho$ , which specifies an element of  $\mathcal{D}$  for each set variable in a constraint system:

$$\rho \in \text{SetEnv} = \text{SetVar} \longrightarrow \mathcal{D}$$

The collection of set environments forms a complete lattice  $\text{SetEnv}$  under the pointwise extension of the ordering relation  $\sqsubseteq$  on  $\mathcal{D}$ .

For each set environment  $\rho$ , we define the following unique extension  $\rho^*$  that specifies a meaning for set expressions:

$$\begin{aligned} \rho^* : \text{SetExp} &\longrightarrow \mathcal{D} \\ \rho^*(\alpha) &= \rho(\alpha) \\ \rho^*(c) &= \langle \{c\}, \top, \perp \rangle \\ \rho^*(\text{dom}(\tau)) &= \text{dom}(\rho^*(\tau)) \\ \rho^*(\text{rng}(\tau)) &= \text{rng}(\rho^*(\tau)) \end{aligned}$$

Where there is no confusion, we remove the asterisk and simply use  $\rho$  to denote  $\rho^*$ .

A set environment  $\rho$  *satisfies* a constraint  $C = [\tau_1 \leq \tau_2]$  (written  $\rho \models C$ ) if  $\rho(\tau_1) \sqsubseteq \rho(\tau_2)$ . Similarly,  $\rho$  *satisfies*  $S$ , or  $\rho$  is a *solution* of  $S$  (written  $\rho \models S$ ) if  $\rho \models C$  for each  $C \in S$ . The relation  $\models$  is obviously reflexive and transitive. The *solution space* of a constraint system  $S$  is:

$$\text{Soln}(S) = \{\rho \mid \rho \models S\}$$

A constraints set  $S_1$  *entails*  $S_2$  (written  $S_1 \models S_2$ ) iff  $\text{Soln}(S_1) \subseteq \text{Soln}(S_2)$ , and  $S_1$  is *observably equivalent to*  $S_2$  (written  $S_1 \cong S_2$ ) iff  $S_1 \models S_2$  and  $S_2 \models S_1$ .

The *restriction* of a solution space to a collection of set variables  $E$  is:

$$\text{Soln}(S) \upharpoonright_E = \{\rho \mid \exists \rho' \in \text{Soln}(S) \text{ such that } \rho(\alpha) = \rho'(\alpha) \ \forall \alpha \in E\}$$

---

$\Gamma \cup \{x : \alpha\} \vdash x : \alpha, \emptyset$	$(var)$
$\Gamma \vdash b : \alpha, \{b \leq \alpha\}$	$(const)$
$\frac{\Gamma \vdash M : \alpha, \mathcal{S}}{\Gamma \vdash M^l : \alpha, \mathcal{S} \cup \{\alpha \leq l\}}$	$(label)$
$\frac{\Gamma \cup \{x : \alpha_1\} \vdash M : \alpha_2, \mathcal{S}}{\Gamma \vdash (\lambda^t x.M) : \alpha, \mathcal{S} \cup \{t \leq \alpha, \text{dom}(\alpha) \leq \alpha_1, \alpha_2 \leq \text{rng}(\alpha)\}}$	$(abs)$
$\frac{\Gamma \vdash M_i : \beta_i, \mathcal{S}_i}{\Gamma \vdash (M_1 \ M_2) : \alpha, \mathcal{S}_1 \cup \mathcal{S}_2 \cup \{\beta_2 \leq \text{dom}(\beta_1), \text{rng}(\beta_1) \leq \alpha\}}$	$(app)$
$\frac{\begin{array}{l} \Gamma \vdash V : \alpha, \mathcal{S}_V \\ A = \text{Vars}(\mathcal{S}_V) \setminus (FV[\text{rng}(\Gamma)] \cup \text{Label}) \\ \Gamma \cup \{x : \forall A. (\alpha, \mathcal{S}_V)\} \vdash M : \beta, \mathcal{S} \end{array}}{\Gamma \vdash (\text{let } (x \ V) \ M) : \beta, \mathcal{S}}$	$(let)$
$\frac{\psi \text{ is a substitution of fresh vars for } A}{\Gamma \cup \{x : \forall A. (\alpha, \mathcal{S}_V)\} \vdash x : \psi(\alpha), \psi(\mathcal{S}_V)}$	$(inst)$

Figure 3: Constraint derivation rules.

---

There are actually *more* set environments in the restricted solution space, since these additional environments can specify arbitrary domain elements for all set variables that are not in  $E$ .

We extend the notion of restriction to the entailment and observable equivalence of constraint systems.

**Definition 3.1.** (*Restricted Entailment, Restricted Observable Equivalence*)

- If  $\text{Soln}(\mathbf{S}_1) \mid_E \subseteq \text{Soln}(\mathbf{S}_2) \mid_E$ , then  $\mathbf{S}_1$  *entails*  $\mathbf{S}_2$  *with respect to*  $E$  (written  $\mathbf{S}_1 \models_E \mathbf{S}_2$ ).
- If  $\mathbf{S}_1 \models_E \mathbf{S}_2$  and  $\mathbf{S}_2 \models_E \mathbf{S}_1$  then that  $\mathbf{S}_1$  and  $\mathbf{S}_2$  are *observably equivalent with respect to*  $E$  (written  $\mathbf{S}_1 \cong_E \mathbf{S}_2$ ).

■

### 3.4 Deriving Constraints

The specification phase of set-based analysis derives constraints on the sets of values that program expressions may assume. Following Aiken *et al.* [2] and Palsberg and O’Keefe [16], we formulate this derivation as a proof system.

The derivation proceeds in a syntax-directed manner according to the constraint derivation rules presented in figure 3. Each rule infers a judgment of the form  $\Gamma \vdash M : \alpha, \mathcal{S}$ , where:

1. the *derivation context*  $\Gamma$  maps the free variables of the expression  $M$  to either set variables or *constraint schemas* (see below);
2.  $\alpha$  names the value set of  $M$ ; and
3. the constraint system  $\mathcal{S}$  is a *simple constraint system* (see below) describes the data-flow relationships of  $M$ , using  $\alpha$ .

The constraint derivation rules only generates a certain subset of the constraint language, called *simple constraints*. Simple constraints have the form:

$$\mathcal{C} \in \text{SimpleCon} = \begin{array}{c|c|c} c \leq \beta & & \alpha \leq \beta \\ \hline \alpha \leq \text{dom}(\beta) & | & \text{dom}(\alpha) \leq \beta \\ \hline \alpha \leq \text{rng}(\beta) & | & \text{rng}(\alpha) \leq \beta \end{array}$$

$$\mathcal{S} \in \text{SimpleConSystem} = \mathcal{P}_{\text{fn}}(\text{SimpleCon})$$

A collection of such simple constraints forms a *simple constraint system*. We use the calligraphic letters  $\mathcal{C}$  and  $\mathcal{S}$  as meta-variables ranging over simple constraints and simple constraint systems, respectively.

The constraint derivation rule (*var*) derives appropriate constraints for a variable reference  $x$ . This rule generates the constraint  $\beta \leq \alpha$ , where  $\beta$  describes the value set of  $x$ , and  $\alpha$  denoting the value set for this reference to  $x$ . The constraint derivation rule (*const*) generates the constraint  $b \leq \alpha$ , which ensures that the value set for a constant expression contains that constant. The rule (*label*) records the possible values of a labeled expression  $M^l$  in the label  $l$ .

The rule (*abs*) for functions records the function's tag, and also propagates values from the function's domain into its formal parameter and from the function's body into its range. The rule (*app*) for applications propagates values from the argument expression into the domain of the applied function and from the range of that function into the result of the application expression. The correctness of the rules (*abs*) and (*app*) relies on the anti-monotonicity of the underlying ordering  $\sqsubseteq$  in the domain position.

The rule (*let*) produces a *constraint schema*  $\sigma = \forall \bar{\alpha}. (\beta, \mathcal{S})$  for polymorphic, **let**-bound values [2, 19]. The set variable  $\beta$  names the result of the value; the simple constraint system  $\mathcal{S}$  describes the data-flow relationships of the value, using  $\beta$ ; and the set  $\bar{\alpha} = \{\alpha_1, \dots, \alpha_n\}$  contains those internal set variables of the constraint system that must be duplicated at each reference to the **let**-bound variable via the rule (*inst*).

The derivation context  $\Gamma$  maps program variables to either set variables or constraint schemas:

$$\begin{aligned} \Gamma &\in \text{DerivCtxt} = \text{Var} \longrightarrow_p \text{SetVar} \cup \text{ConSchema} \\ \sigma &\in \text{ConSchema} = \forall \bar{\alpha}. (\beta, \mathcal{S}) \end{aligned}$$

We use  $FV[\text{range}(\Gamma)]$  to denote the free set variables in the range of  $\Gamma$ . The free set variables of a constraint schema  $\forall \bar{\alpha}. (\beta, \mathcal{S})$  are those in  $\mathcal{S}$  but not in  $\bar{\alpha}$ , and the free variables of a set variable is simply the set variable itself.

Many of the constraint derivation rules contain *meta* set variables. For example, the rule (*const*):

$$\Gamma \vdash b : \alpha, \{b \leq \alpha\} \quad (\text{const})$$

mentions the meta set variable  $\alpha$ . Any time this rule is applied, we need to choose an appropriate set variable for this meta variable. Choosing a fresh set variable not used elsewhere in the derivation yields a more accurate analysis. A *most general constraint derivation* is one that always uses fresh set variable for these meta variables, and a *most general constraint system* for an expression is one produced by a most general constraint derivation. However, the use of fresh variables is not strictly necessary for the correctness of the analysis. As an extreme example, we could perform the entire analysis using a single set variable, although this would yield extremely coarse results, and would be of no practical use. But the ability to consider constraint derivations that re-use certain set variables significantly simplifies the subject reduction proofs of the following section.

### 3.5 Soundness of the Derived Constraints

Let  $P$  be a program such that  $\emptyset \vdash P : \alpha, \mathcal{S}$ . Typically,  $\mathcal{S}$  has many solutions. Each solution  $\rho$  of  $\mathcal{S}$  correctly approximates the value sets of labeled expressions in  $P$ . That is, if  $\rho$  is a solution of  $\mathcal{S}$  and  $V$  is a possible value of some expression  $M^l$  in  $P$ , then  $V$  in  $\rho(l)$ . We prove this property using a subject reduction proof [7], following Wright and Felleisen [21] and Palsberg [15].

**Main Lemma 3.2 (Soundness of the Derived Constraints)** *If  $\emptyset \vdash P : \alpha, \mathcal{S}$  and  $\rho \models \mathcal{S}$  and  $P \mapsto^* \mathcal{E}[V^l]$  then  $V$  in  $\rho(l)$ .*

**Proof:** The Subject Reduction for  $\mapsto$  Lemma (3.3) shows that standard reduction steps preserve entailment. Hence, since  $P \mapsto^* \mathcal{E}[V^l]$ , there exists some  $\mathcal{S}'$  such that  $\emptyset \vdash \mathcal{E}[V^l] : \alpha, \mathcal{S}'$  and  $\mathcal{S} \models \mathcal{S}'$ . The derivation of this judgment must contain a sub-derivation concluding:

$$\frac{\Gamma \vdash V : \beta, \mathcal{S}_V}{\Gamma \vdash V^l : \beta, \mathcal{S}_V \cup \{\beta \leq l\}} \quad (label)$$

Except for the rule (*let*), each application of a constraint derivation rule can only extend the constraint system produced by its sub-derivation. Since definition of evaluation contexts does not contain a clause for **let**-expressions, there cannot be any **let**-expressions on the spine from  $V^l$  to  $\mathcal{E}[V^l]$ . Hence  $\mathcal{S}_V \cup \{\beta \leq l\} \subseteq \mathcal{S}'$ .

Since  $\rho \models \mathcal{S}$ ,  $\mathcal{S} \models \mathcal{S}'$ , and  $\mathcal{S}' \supseteq \mathcal{S}_V \cup \{\beta \leq l\}$ , we have that  $\rho \models \mathcal{S}_V \cup \{\beta \leq l\}$ . Hence  $V$  in  $\rho(\beta)$  by the Value Typing Lemma 3.7. But  $\rho(\beta) \sqsubseteq \rho(l)$ , hence  $V$  in  $\rho(l)$ , as required. ■

The proof of the above result relies on the following lemma showing that standard reduction steps preserves the entailment of the derived constraint systems.

**Lemma 3.3 (Subject Reduction for  $\mapsto$ )** *If  $\Gamma \vdash M_1 : \alpha, \mathcal{S}_1$  and  $M_1 \mapsto M_2$ , then  $\Gamma \vdash M_2 : \alpha, \mathcal{S}_2$  where  $\mathcal{S}_1 \models \mathcal{S}_2$ .*

**Proof:** Follows from the Subject Reduction Lemma 3.4 and the Replacement Lemma 3.5. ■

**Lemma 3.4 (Subject Reduction for  $\rightarrow$ )** *If  $\Gamma \vdash M_1 : \alpha, \mathcal{S}_1$  and  $M_1 \rightarrow M_2$ , then  $\Gamma \vdash M_2 : \alpha, \mathcal{S}_2$  such that  $\mathcal{S}_1 \models \mathcal{S}_2$ .*

**Proof:** See Appendix A.1. ■

**Lemma 3.5 (Replacement)** *If:*

1.  $D$  is a deduction concluding  $\Gamma \vdash \mathcal{E}[M_1] : \alpha, \mathcal{S}_1$ ,
2.  $D_1$  is a sub-deduction of  $D$  concluding  $\Gamma' \vdash M_1 : \beta, \mathcal{S}'_1$ ,
3.  $D_1$  occurs in  $D$  in the position corresponding to the hole  $([ ])$  in  $\mathcal{E}$ , and
4.  $\Gamma' \vdash M_2 : \beta, \mathcal{S}'_2$  where  $\mathcal{S}'_1 \models \mathcal{S}'_2$ ,

then  $\Gamma \vdash \mathcal{E}[M_2] : \alpha, \mathcal{S}_2$  where  $\mathcal{S}_1 \models \mathcal{S}_2$ .

**Proof:** Follows the proof idea of Hindley and Seldin [11:page 181]. ■

The Flow Lemma describes conditions under which we can replace the result set variable returned by the constraint derivation rules.

**Lemma 3.6 (Flow)** *If  $\Gamma \vdash M : \alpha, \mathcal{S}$  then for all  $\gamma \in \text{SetVar}$ ,  $\Gamma \vdash M : \gamma, \mathcal{S}'$  with  $\mathcal{S} \cup \{\alpha \leq \gamma\} \models \mathcal{S}'$ .*

**Proof:** See Appendix A.1. ■

The Value Typing Lemma simply states that any solution to the constraint system for a syntactic value always corresponds to a value set invariant that includes that value.

**Lemma 3.7 (Value Typing)** *If  $\Gamma \vdash V : \alpha, \mathcal{S}$  and  $\rho \models \mathcal{S}$ , then  $V$  in  $\rho(\alpha)$ .*

**Proof:** By considering the two cases  $V = b$  and  $V = (\lambda^t x.M)$ . ■

### 3.6 Solving Set Constraints

Every simple constraint system admits the trivial solution  $\rho^{\top_s}$  defined by  $\rho^{\top_s}(\alpha) = \top_s$  where  $\top_s$  is recursively defined as  $\top_s = \langle \text{Const}, \top_s, \top_s \rangle$ . The domain element  $\top_s$  represents the set of all run-time values, including functions that can take any value as an argument, and return any value as a result.

**Lemma 3.8** *If  $\mathcal{S}$  is a simple constraint system then  $\rho^{\top_s} \models \mathcal{S}$ .*

**Proof:** By a case analysis showing that  $\rho^{\top_s} \models \mathcal{C}$  for any simple constraint  $\mathcal{C}$ . ■

Since  $\top_s$  represents all run-time values, this solution is highly approximate and thus utterly useless. Fortunately, simple constraint systems yield many additional solutions that more accurately characterize the value sets of program expressions.

To illustrate this idea, consider the program  $P = (\lambda^t x.x)$ . According to the constraint derivation rules of figure 3, this program yields the constraint system:

$$\{t \leq \alpha_P, \text{dom}(\alpha_P) \leq \alpha_x, \alpha_x \leq \alpha_M, \alpha_M \leq \text{rng}(\alpha_P)\}$$



In addition to the trivial solution  $\rho^{\top_s}$ , this constraint system admits a number of other solutions, including:

$$\begin{aligned}\rho_1 &= \{\alpha_P \mapsto \langle \{t\}, \perp, \perp \rangle, \alpha_x \mapsto \perp, \alpha_M \mapsto \perp\} \\ \rho_2 &= \{\alpha_P \mapsto \langle \{t\}, \top, \top \rangle, \alpha_x \mapsto \top, \alpha_M \mapsto \top\} \\ \rho_3 &= \{\alpha_P \mapsto \langle \{t, c_1\}, X, X \rangle, \alpha_x \mapsto X, \alpha_M \mapsto X\}\end{aligned}$$

where  $X = \langle \{c_2\}, \perp, \perp \rangle$ , and  $c_1$  and  $c_2$  are arbitrary constants. Because we assume  $P$  to be the entire program, the function tagged  $t$  is never applied, and hence the set of run-time values for  $x$  is simply the empty set. The solution  $\rho_1$  describes this (empty) set of run-time values of  $x$  more accurately than either  $\rho_2$  or  $\rho_3$ . Yet these three solutions are incomparable under the ordering  $\sqsubseteq$ , since the ordering models the flow of values through a program, but does not rank set environments according to their accuracy.

Therefore we introduce an alternative ordering  $\sqsubseteq_s$  on  $\mathcal{D}$  that ranks environments according to their accuracy. This ordering is monotonic in the domain position:

$$\langle C_1, D_1, R_1 \rangle \sqsubseteq_s \langle C_2, D_2, R_2 \rangle \text{ iff } C_1 \subseteq C_2, D_1 \sqsubseteq_s D_2, R_1 \sqsubseteq_s R_2$$

The maximal and minimal elements of  $\mathcal{D}$  under  $\sqsubseteq_s$  are the solutions to the equations:

$$\begin{aligned}\top_s &= \langle \text{Const}, \top_s, \top_s \rangle \\ \perp_s &= \langle \emptyset, \perp_s, \perp_s \rangle\end{aligned}$$

respectively. The least upper bound and greatest lower bound operations are recursively defined as:

$$\begin{aligned}\langle C_1, D_1, R_1 \rangle \sqcup_s \langle C_2, D_2, R_2 \rangle &= \langle C_1 \cup C_2, D_1 \sqcup_s D_2, R_1 \sqcup_s R_2 \rangle \\ \langle C_1, D_1, R_1 \rangle \sqcap_s \langle C_2, D_2, R_2 \rangle &= \langle C_1 \cap C_2, D_1 \sqcap_s D_2, R_1 \sqcap_s R_2 \rangle\end{aligned}$$

Under the ordering  $\sqsubseteq_s$ , a simple constraint system has both a maximal solution ( $\rho^{\top_s}$  above) and a minimal solution. The minimal solution exists because the greatest lower bound  $\sqcap_s$  with respect to  $\sqsubseteq_s$  of two solutions is also a solution [10].

**Lemma 3.9 (Least Solution of Simple Constraint Systems)** *Every simple constraint system has a solution that is least with respect to  $\sqsubseteq_s$ .*

**Proof:** See Appendix A.2. ■

Using Lemma 3.9, it makes sense to define  $LeastSoln(\mathcal{S})$  as the least solution of the simple constraint system  $\mathcal{S}$  under the ordering  $\sqsubseteq_s$ . Since this solution yields the most accurate invariants consistent with the constraints  $\mathcal{S}$ , we define set-based analysis as the function that extracts the possible values for each labeled expression from this least solution.

**Definition 3.10.** ( $sba : \Lambda^0 \longrightarrow (Label \longrightarrow \mathcal{P}(Value))$ ) If  $\emptyset \vdash P : \alpha, \mathcal{S}$  is a most general derivation, then

$$sba(P)(l) = \{V \mid V \text{ in } LeastSoln(\mathcal{S})(l)\}$$

■

By Lemma 3.2,  $sba(P)$  correctly characterizes the possible value sets for each labeled expression.

**Theorem 3.11** *If  $P \longmapsto^* \mathcal{E}[V^l]$  then  $V \in sba(P)(l)$ .*

**Proof:** Follows from Lemma 3.2. ■

### 3.6.1 Computing the Least Solution

To compute  $sba(P)$ , we derive the most general constraint system for  $P$  and close that constraint system under the rules  $\Theta$  described in figure 4. Intuitively, these rules infer all the data-flow paths in the program, which are described by constraints of the form  $\beta \leq \gamma$  (for  $\beta, \gamma \in SetVar$ ), and propagate values along those data-flow paths. Specifically, the rules  $(s_1)$ ,  $(s_2)$ , and  $(s_3)$  propagate information about constants, function domains and function ranges forward along the data-flow paths of the program. The rule  $(s_4)$  constructs the data-flow paths from actual to formal parameters for each function call, and the rule  $(s_5)$  similarly constructs data-flow paths from function bodies to corresponding call sites. We write  $\mathcal{S} \vdash_{\Theta} \mathcal{C}$  if  $\mathcal{S}$  proves  $\mathcal{C}$  via the rules  $\Theta$ , and use  $\Theta(\mathcal{S})$  to denote the closure of  $\mathcal{S}$  under  $\Theta$ , *i.e.*, the set  $\{\mathcal{C} \mid \mathcal{S} \vdash_{\Theta} \mathcal{C}\}$ . An algorithm for computing  $\Theta(\mathcal{S})$  is included in the next section.

---


$$\begin{array}{c}
 \frac{c \leq \beta \quad \beta \leq \gamma}{c \leq \gamma} \quad (s_1) \\
 \frac{\alpha \leq \mathbf{rng}(\beta) \quad \beta \leq \gamma}{\alpha \leq \mathbf{rng}(\gamma)} \quad (s_2) \\
 \frac{\mathbf{dom}(\beta) \leq \alpha \quad \beta \leq \gamma}{\mathbf{dom}(\gamma) \leq \alpha} \quad (s_3) \\
 \frac{\alpha \leq \mathbf{rng}(\beta) \quad \mathbf{rng}(\beta) \leq \gamma}{\alpha \leq \gamma} \quad (s_4) \\
 \frac{\alpha \leq \mathbf{dom}(\beta) \quad \mathbf{dom}(\beta) \leq \gamma}{\alpha \leq \gamma} \quad (s_5)
 \end{array}$$

Figure 4: The rules  $\Theta = \{s_1, \dots, s_5\}$

---

This closure process propagates all information concerning the possible constants for labeled expressions into constraints of the form  $c \leq l$ . Hence, we can infer  $sba(P)$  from  $\Theta(\mathcal{S})$  according to the following theorem.

**Theorem 3.12** *If  $P \in \Lambda^0$  and  $\emptyset \vdash P : \alpha$ ,  $\mathcal{S}$  is a most general constraint derivation then:*

$$\begin{aligned}
 sba(P)(l) = & \quad \{b \mid \mathcal{S} \vdash_{\Theta} b \leq l\} \\
 & \cup \quad \{(\lambda^t x.M) \mid \mathcal{S} \vdash_{\Theta} t \leq l\}
 \end{aligned}$$

**Proof:** See Appendix A.2. ■

## 4 Constraint Simplification

The traditional set-based analysis we have just described has proven highly effective for programs of up to a couple of thousand lines of code. Unfortunately, it is useless for

larger programs due to its nature as a whole-program analysis and due to the size of the constraint systems it produces, which are quadratic in the size of (large) programs. Storing and manipulating these constraint systems can be extremely expensive.

To overcome this problem, we have developed algorithms for *simplifying* constraint systems. A simplified version of a constraint system contains significantly fewer constraints, yet still preserves the essential characteristics of the original system. Applying these simplification algorithms at strategic points during the constraint derivation, *e.g.*, to the constraint system for a module definition or a polymorphic function definition, significantly reduces both the time and space required by the overall analysis.

The following section shows that constraint simplification does not affect the analysis results, provided the simplified system is *observably equivalent* to the original system. Section 4.2 presents a proof-theoretic formulation of observable equivalence, and section 4.3 exploits this formulation to develop a complete algorithm for deciding the observable equivalence of constraint systems. The insights provided by this theoretical development lead to the practical constraint simplification algorithms of section 4.8.

#### 4.1 Conditions for Constraint Simplification

Let us consider a program  $P$  containing a program component  $M$ , where  $M$  may be a module definition or a polymorphic function definition. Suppose the constraint derivations for  $M$  concludes:

$$\Gamma \vdash M : \alpha, \mathcal{S}_1$$

where  $\mathcal{S}_1$  is the constraint system for  $M$ . Our goal is to replace  $\mathcal{S}_1$  by a simpler constraint system, say  $\mathcal{S}_2$ , without changing the results of the analysis.

Let the context surrounding  $M$  be  $C$ , *i.e.*,  $P = C[M]$ . Since the constraint derivation process is compositional, the constraint derivation for the entire program concludes:

$$\emptyset \vdash P : \beta, \mathcal{S}_C \cup \mathcal{S}_1$$

where  $\mathcal{S}_C$  is the constraint system for  $C$ . The union of the sets  $\mathcal{S}_C$  and  $\mathcal{S}_1$  describes the space of solutions for the entire program, which is the same as the intersection of the two respective solution spaces:

$$\text{Soln}(\mathcal{S}_C \cup \mathcal{S}_1) = \text{Soln}(\mathcal{S}_C) \cap \text{Soln}(\mathcal{S}_1)$$

Hence  $\text{Soln}(\mathcal{S}_1)$  describes at least all the properties of  $\mathcal{S}_1$  relevant to the analysis, but it may also describe solutions for set variables that are not relevant to the analysis. In particular:

- $\text{sbv}(P)$  only references the solutions for label variables; and
- an inspection of the constraint derivation rules shows that the only interactions between  $\mathcal{S}_C$  and  $\mathcal{S}_1$  are due to the set variables in  $\{\alpha\} \cup \text{FV}[\text{range}(\Gamma)]$ .

In short, the only properties of  $\mathcal{S}_1$  relevant to the analysis is the solution space for its *external set variables*:

$$E = \text{Label} \cup \{\alpha\} \cup \text{FV}[\text{range}(\Gamma)]$$

For our original problem, this means that we want a constraint system  $\mathcal{S}_2$  whose solution space restricted to  $E$  is equivalent to that of  $\mathcal{S}_1$  restricted to  $E$ :

$$\text{Soln}(\mathcal{S}_1) \upharpoonright_E = \text{Soln}(\mathcal{S}_2) \upharpoonright_E$$

or, with the notation from section 3.2,  $\mathcal{S}_1$  and  $\mathcal{S}_2$  are observably equivalent on  $E$ :

$$\mathcal{S}_1 \cong_E \mathcal{S}_2 .$$

We can translate this compaction idea into an additional rule for the constraint derivation system:

$$\frac{\Gamma \vdash_{\cong} M : \alpha, \mathcal{S}_1 \quad \mathcal{S}_1 \cong_E \mathcal{S}_2 \text{ where } E = \text{Label} \cup \text{FV}[\text{range}(\Gamma)] \cup \{\alpha\}}{\Gamma \vdash_{\cong} M : \alpha, \mathcal{S}_2} \quad (\cong)$$

This rule is *admissible* in that any derivation (denoted using  $\vdash_{\cong}$ ) in the extended constraint derivation system produces information that is equivalent to the information produced by the original analysis.

**Lemma 4.1 (Admissibility of  $(\cong)$ )** *If  $\emptyset \vdash_{\cong} P : \alpha, \mathcal{S}$  is a most general constraint derivation then:*

$$\text{sba}(P)(l) = \text{const}(\text{LeastSoln}(\mathcal{S})(l))$$

**Proof:** See Appendix A.3. ■

## 4.2 The Proof Theory of Observable Equivalence

Since the new derivation rule  $(\cong)$  involves the *semantic* notion of observably equivalent constraint systems, it cannot be used in an algorithmic manner. To make this rule useful, we must first reformulate the observable equivalence relation (or some conservative approximation thereof) as a *syntactic* proof system.

The key properties of the observational equivalence relation are reflections of the reflexivity and transitivity of the ordering relation ( $\sqsubseteq$ ) and the monotonicity and anti-monotonicity of the functions *rng* and *dom*, respectively. We can reify these properties into a syntactic proof system via the following inference rules  $\Delta$ :

$$\begin{array}{lll} \alpha \leq \alpha & (\text{reflex}) & \frac{\tau_1 \leq \tau \quad \tau \leq \tau_2}{\tau_1 \leq \tau_2} \quad (\text{trans}_{\tau}) \quad \frac{\kappa_1 \leq \kappa_2}{\text{rng}(\kappa_1) \leq \text{rng}(\kappa_2) \quad \text{dom}(\kappa_2) \leq \text{dom}(\kappa_1)} \quad (\text{compat}) \end{array}$$

The meta-variables  $\kappa, \kappa_1, \kappa_2$  range over non-constant set expressions:

$$\kappa, \kappa_1, \kappa_2 = \alpha \mid \text{dom}(\kappa) \mid \text{rng}(\kappa)$$

This restriction avoids inferring useless tautologies. For example, without this restriction, the constraint  $c \leq \alpha$  would yield the constraint  $\text{rng}(c) \leq \text{rng}(\alpha)$  via  $(\text{compat})$ , which is a tautology since  $\text{rng}(c) = \perp$ .

The rules (*reflex*) and (*trans* <sub>$\tau$</sub> ) capture the reflexivity and transitivity of the ordering relation  $\sqsubseteq$ ; (*compat*) reflects the monotonicity and anti-monotonicity of the functions *rng* and *dom*, respectively. Since many of the inferred constraints lie outside of the original language of simple constraints, we define an extended *compound constraint* language that includes all of the inferred constraints:

$$\begin{aligned} \mathbf{C} &\in \text{CmpdConstraint} = c \leq \kappa \mid \kappa \leq \kappa \\ \mathbf{S} &\in \text{CmpdConSystem} = \mathcal{P}_{\text{fin}}(\text{CmpdConstraint}) \end{aligned}$$

We use the boldface roman letters  $\mathbf{C}$  and  $\mathbf{S}$  as meta-variables ranging over compound constraints and compound constraint systems, respectively.

The proof system  $\Delta$  completely captures the relevant properties of the ordering  $\sqsubseteq$  and the functions *rng* and *dom*. That is,  $\Delta$  is both sound and complete.

**Lemma 4.2 (Soundness and Completeness of  $\Delta$ )** *For a compound constraint system  $\mathbf{S}$  and a compound constraint  $\mathbf{C}$ :*

$$\mathbf{S} \vdash_{\Delta} \mathbf{C} \iff \mathbf{S} \models \mathbf{C}$$

**Proof:** See Appendix A.4. ■

This lemma implies that  $\Delta(\mathcal{S})$  contains exactly those compound constraints that hold in all environments in  $\text{Soln}(\mathcal{S})$ . Hence, if we consider a collection of external set variables  $E$ , then  $\Delta(\mathcal{S}) \mid_E$  contains all compound constraints that hold in all environments in  $\text{Soln}(\mathcal{S}) \mid_E$ . Therefore the following lemma holds.

**Lemma 4.3** *For a compound constraint system  $\mathbf{S}$ ,  $\mathbf{S} \cong_E \Delta(\mathbf{S}) \mid_E$ .*

**Proof:** See Appendix A.4. ■

We could use this result to define a proof-theoretic equivalent of restricted entailment as follows:

$$\mathcal{S}_1 \vdash_{\Delta}^E \mathcal{S}_2 \text{ if and only if } \Delta(\mathcal{S}_1) \mid_E \supseteq \Delta(\mathcal{S}_2) \mid_E$$

and then show that  $\mathcal{S}_1 \vdash_{\Delta}^E \mathcal{S}_2$  if and only if  $\mathcal{S}_1 \models_E \mathcal{S}_2$ . However, this definition based on the proof system  $\Delta$  does not lend itself to an efficient implementation. Specifically, checking if two potential antecedents of (*trans* <sub>$\tau$</sub> ) contain the same set expression  $\tau$  involves comparing two potentially large set expressions. Hence we develop an alternative proof system that is more suitable for an implementation, yet infers the same constraints as  $\Delta$ .

The alternative system consists of the inference rules  $\Psi$  described in figure 5, together with the rules  $\Theta$  from figure 4. The rules (*compose*<sub>1...4</sub>) of  $\Psi$  replace a reference to a set variable by an upper or lower (non-constant) bound for that variable, as appropriate. The rules (*reflex*) and (*compat*) of  $\Psi$  are as described above. The rule (*trans* <sub>$\alpha$</sub> ) of  $\Psi$  provides a weaker characterization of transitivity than the previous rule (*trans* <sub>$\tau$</sub> ), but, provided we start from with a simple constraint system, the additional rules,  $\Theta$  and (*compose*<sub>1...4</sub>), compensate for this weakness. That is, suitable combinations of these additional rules allow us to infer any constraint that could be inferred by the rule (*trans* <sub>$\tau$</sub> ).

---

$\frac{\alpha \leq \mathbf{rng}(\beta) \quad \beta \leq \kappa}{\alpha \leq \mathbf{rng}(\kappa)}$	$(compose_1)$	$\frac{\alpha \leq \mathbf{dom}(\beta) \quad \beta \geq \kappa}{\alpha \leq \mathbf{dom}(\kappa)}$	$(compose_2)$
$\frac{\alpha \geq \mathbf{rng}(\beta) \quad \beta \geq \kappa}{\alpha \geq \mathbf{rng}(\kappa)}$	$(compose_3)$	$\frac{\alpha \geq \mathbf{dom}(\beta) \quad \beta \leq \kappa}{\alpha \geq \mathbf{dom}(\kappa)}$	$(compose_4)$
$\alpha \leq \alpha$	$(reflex)$	$\frac{\tau_1 \leq \alpha \quad \alpha \leq \tau_2}{\tau_1 \leq \tau_2}$	$(trans_\alpha)$
$\frac{\kappa_1 \leq \kappa_2}{\mathbf{rng}(\kappa_1) \leq \mathbf{rng}(\kappa_2) \quad \mathbf{dom}(\kappa_2) \leq \mathbf{dom}(\kappa_1)}$	$(compat)$		

---

Figure 5: The Inference Rule System  $\Psi$ 

**Lemma 4.4 (Equivalence of Proof Systems)** *For a simple constraint system  $\mathcal{S}$ :*

$$\Delta(\mathcal{S}) = \Psi\Theta(\mathcal{S})$$

**Proof:** See Appendix A.4. ■

We could use this result to define a proof-theoretic equivalent of restricted entailment as follows:

$$\mathcal{S}_1 \vdash_{\Psi\Theta}^E \mathcal{S}_2 \text{ if and only if } \Psi\Theta(\mathcal{S}_1) \mid_E \supseteq \Psi\Theta(\mathcal{S}_2) \mid_E$$

and then show that  $\mathcal{S}_1 \vdash_{\Psi\Theta}^E \mathcal{S}_2$  if and only if  $\mathcal{S}_1 \models_E \mathcal{S}_2$ . However, this definition is needlessly inefficient. Because  $(compat)$  does not eliminate any variables, any  $(compat)$ -consequent in  $\Psi\Theta(\mathcal{S}_2) \mid_E$  is subsumed by its antecedent. If we define:

$$\Pi = \Psi \setminus \{compat\}$$

then this argument implies that  $\Psi\Theta(\mathcal{S}_2) \mid_E \cong \Pi\Theta(\mathcal{S}_2) \mid_E$ . Hence we get the following lemma.

**Lemma 4.5**  $\Psi\Theta(\mathcal{S}) \mid_E \cong \Pi\Theta(\mathcal{S}) \mid_E$ .

**Proof:** See Appendix A.4. ■

Together, Lemmas 4.3, 4.4 and 4.5 provide the basis to introduce proof-theoretic equivalents of restricted entailment and observable equivalence.

**Definition 4.6.**  $(\vdash_{\Psi\Theta}^E, =_{\Psi\Theta}^E)$

- $\mathcal{S}_1 \vdash_{\Psi\Theta}^E \mathcal{S}_2$  if and only if  $\Psi\Theta(\mathcal{S}_1) \mid_E \supseteq \Pi\Theta(\mathcal{S}_2) \mid_E$ ,
- $\mathcal{S}_1 =_{\Psi\Theta}^E \mathcal{S}_2$  if and only if  $\mathcal{S}_1 \vdash_{\Psi\Theta}^E \mathcal{S}_2$  and  $\mathcal{S}_2 \vdash_{\Psi\Theta}^E \mathcal{S}_1$ .

■

The two relations  $\vdash_{\Psi\Theta}^E$  and  $=_{\Psi\Theta}^E$  completely characterize restricted entailment and observable equivalence of constraint systems.

**Theorem 4.7 (Soundness and Completeness of  $\vdash_{\Psi\Theta}^E$  and  $=_{\Psi\Theta}^E$ )**

1.  $\mathcal{S}_1 \vdash_{\Psi\Theta}^E \mathcal{S}_2$  if and only if  $\mathcal{S}_1 \models_E \mathcal{S}_2$ .
2.  $\mathcal{S}_1 =_{\Psi\Theta}^E \mathcal{S}_2$  if and only if  $\mathcal{S}_1 \cong_E \mathcal{S}_2$ .

**Proof:** See Appendix A.4. ■

### 4.3 Deciding Observable Equivalence

While the relation  $=_{\Psi\Theta}^E$  completely characterizes the model-theoretic observable equivalence relation  $\cong_E$ , an implementation of the extended constraint derivation system needs a decision algorithm for  $=_{\Psi\Theta}^E$ .

Given two simple constraint systems  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , this algorithm needs to verify that  $\Psi\Theta(\mathcal{S}_1) \mid_E = \Psi\Theta(\mathcal{S}_2) \mid_E$ . If  $\mathcal{S}_1$  and  $\mathcal{S}_2$  are first closed under  $\Theta$ , then the algorithm only needs to verify that  $\Psi(\mathcal{S}_1) \mid_E = \Psi(\mathcal{S}_2) \mid_E$ . The naive approach to enumerate and to compare the two constraint systems  $\Psi(\mathcal{S}_1) \mid_E$  and  $\Psi(\mathcal{S}_2) \mid_E$  does not work, since they are typically infinite. For example, if  $\mathcal{S} = \{\alpha \leq \mathbf{rng}(\alpha)\}$ , then  $\Psi(\mathcal{S})$  is the infinite set  $\{\alpha \leq \mathbf{rng}(\alpha), \alpha \leq \mathbf{rng}(\mathbf{rng}(\alpha)), \dots\}$ .

Fortunately, the infinite constraint systems inferred by  $\Psi$  exhibit a regular structure, which we exploit to decide observable equivalence as follows:

1. We generate regular grammars describing the upper and lower bounds for each set variable.
2. We extend these regular grammars to regular tree grammars (RTGs) describing all constraints in  $\Pi(\mathcal{S}_1) \mid_E$  and  $\Pi(\mathcal{S}_2) \mid_E$ . This representation allows us to use a standard RTG containment algorithm to decide if  $\Pi(\mathcal{S}_1) \mid_E \supseteq \Pi(\mathcal{S}_2) \mid_E$ .
3. Based on the RTG containment algorithm, we develop an extended algorithm that decides the more difficult entailment question  $\Psi(\mathcal{S}_1) \mid_E \supseteq \Pi(\mathcal{S}_2) \mid_E$  by allowing for the additional (*compat*) inferences on  $\mathcal{S}_1$ .

By checking entailment in both directions, we can decide if two constraint systems are observably equivalent. These steps are described in more detail below.

### 4.4 Regular Grammars

Our first step is to describe, for each set variable  $\alpha$  in a simple constraint system  $\mathcal{S}$ , the following two languages of the lower and upper non-constant bounds of  $\alpha$ :

$$\begin{aligned} &\{\kappa \mid [\kappa \leq \alpha] \in \Pi(\mathcal{S}) \text{ and } \mathit{SetVar}(\kappa) \subseteq E\} \\ &\{\kappa \mid [\alpha \leq \kappa] \in \Pi(\mathcal{S}) \text{ and } \mathit{SetVar}(\kappa) \subseteq E\} \end{aligned}$$

These languages are generated by a regular grammar. For each set variable  $\alpha$ , the grammar contains the non-terminals  $\alpha_L$  and  $\alpha_U$ , which generate the above lower and upper bounds of  $\alpha$ , respectively.

To illustrate this idea consider the program component  $P = (\lambda^g y.((\lambda^f x.1) y))$ , where  $f$  and  $g$  are function tags, and take  $E = \{\alpha^P\}$ . The constraint system  $\mathcal{S}_P$  for  $P$  (closed under  $\Theta$ ) is described in figure 6, together with the productions in the corresponding regular grammar. This grammar describes the upper and lower non-constant bounds for each set variable in  $\Pi(\mathcal{S}_P) \upharpoonright_E$ . For example, the productions:

$$\begin{aligned} \alpha_L^x &\mapsto \alpha_L^r \\ \alpha_L^r &\mapsto \mathbf{dom}(\alpha_U^P) \\ \alpha_U^P &\mapsto \alpha^P \end{aligned}$$

imply that  $\alpha_L^x \mapsto^* \mathbf{dom}(\alpha_P)$ . This lower bound for  $x$  means that the value set for  $x$  must include all values to which the function  $P$  is applied.

Constraints $\mathcal{S}_P$	Grammar $G_r(\mathcal{S}_P, E)$	Additional productions in $G_t(\mathcal{S}_P, E)$
$f \leq \alpha^f$		$R \mapsto [f \leq \alpha_U^f]$
$\mathbf{dom}(\alpha^f) \leq \alpha^x$	$\alpha_L^x \mapsto \mathbf{dom}(\alpha_U^f)$	
$1 \leq \alpha^1$		$R \mapsto [1 \leq \alpha_U^1]$
$\alpha^1 \leq \mathbf{rng}(\alpha^f)$	$\alpha_U^1 \mapsto \mathbf{rng}(\alpha_U^f)$	
$\mathbf{rng}(\alpha^f) \leq \alpha^a$	$\alpha_L^a \mapsto \mathbf{rng}(\alpha_L^f)$	
$\alpha^y \leq \alpha^r$	$\alpha_U^y \mapsto \alpha_U^r \quad \alpha_L^r \mapsto \alpha_L^y$	
$\alpha^r \leq \mathbf{dom}(\alpha^f)$	$\alpha_U^r \mapsto \mathbf{dom}(\alpha_U^f)$	
$g \leq \alpha^P$		$R \mapsto [g \leq \alpha_U^P]$
$\mathbf{dom}(\alpha^P) \leq \alpha^y$	$\alpha_L^y \mapsto \mathbf{dom}(\alpha_U^P)$	
$\alpha^a \leq \mathbf{rng}(\alpha^P)$	$\alpha_U^a \mapsto \mathbf{rng}(\alpha_U^P)$	
$\alpha^r \leq \alpha^x$	$\alpha_U^r \mapsto \alpha_U^x \quad \alpha_L^x \mapsto \alpha_L^r$	
$\alpha^1 \leq \alpha^a$	$\alpha_U^1 \mapsto \alpha_U^a \quad \alpha_L^a \mapsto \alpha_L^1$	
$1 \leq \alpha^a$		$R \mapsto [1 \leq \alpha_U^a]$
	$\alpha_L^P \mapsto \alpha^P \quad \alpha_U^P \mapsto \alpha^P$	$R \mapsto [\alpha_L \leq \alpha_U] \forall \alpha \in \text{SetVar}(\mathcal{S}_P)$

Figure 6: The constraint system, regular grammar, and regular tree grammar for  $P = (\lambda^g y.((\lambda^f x.1) y))$

The productions of the grammar are determined by  $\mathcal{S}_P$  and  $\Pi$ . For example, the constraint  $[\alpha^1 \leq \mathbf{rng}(\alpha^f)] \in \mathcal{S}_P$  implies that for each upper bound  $\kappa$  of  $\alpha^f$ , the rule (*compose*<sub>1</sub>) infers the upper bound  $\mathbf{rng}(\kappa)$  of  $\alpha^1$ . Since, by induction,  $\alpha^f$ 's upper bounds are generated by  $\alpha_U^f$ , the production  $\alpha_U^f \mapsto \mathbf{rng}(\alpha_U^f)$  generates the corresponding upper bounds of  $\alpha^1$ .

More generally, the collection of productions:

$$\{\alpha_U \mapsto \mathbf{rng}(\beta_U) \mid [\alpha \leq \mathbf{rng}(\beta)] \in \mathcal{S}\}$$

describes all bounds inferred via (*compose*<sub>1</sub>) on a simple constraint system  $\mathcal{S}$ . Bounds inferred via the remaining (*compose*) rules can be described in a similar manner. Bounds



inferred via the rule (*reflex*) imply the productions  $\alpha_U \mapsto \alpha$  and  $\alpha_L \mapsto \alpha$  for  $\alpha \in E$ . Finally, consider the rule (*trans <sub>$\alpha$</sub>* ), and suppose this rule infers an upper bound  $\tau$  on  $\alpha$ . This bound must be inferred from an upper bound  $\tau$  on  $\beta$ , using the additional antecedent  $[\alpha \leq \beta]$ . Hence the productions  $\{\alpha_U \mapsto \beta_U \mid [\alpha \leq \beta] \in \mathcal{S}\}$  generate all upper bounds inferred via (*trans <sub>$\alpha$</sub>* ). In a similar fashion, the productions  $\{\beta_L \mapsto \alpha_L \mid [\alpha \leq \beta] \in \mathcal{S}\}$  generate all lower bounds inferred via (*trans <sub>$\alpha$</sub>* ).

**Definition 4.8.** (*Regular Grammar  $G_r(\mathcal{S}, E)$* ) Let  $\mathcal{S}$  be a simple constraint system and  $E$  a collection of set variables. The regular grammar  $G_r(\mathcal{S}, E)$  consists of the non-terminals  $\{\alpha_L, \alpha_U \mid \alpha \in \text{SetVar}(\mathcal{S})\}$  and the following productions:

$$\begin{array}{ll} \alpha_U \mapsto \alpha, \alpha_L \mapsto \alpha & \forall \alpha \in E \\ \alpha_U \mapsto \beta_U, \beta_L \mapsto \alpha_L & \forall [\alpha \leq \beta] \in \mathcal{S} \\ \alpha_U \mapsto \text{dom}(\beta_L) & \forall [\alpha \leq \text{dom}(\beta)] \in \mathcal{S} \\ \alpha_U \mapsto \text{rng}(\beta_U) & \forall [\alpha \leq \text{rng}(\beta)] \in \mathcal{S} \\ \beta_L \mapsto \text{dom}(\alpha_U) & \forall [\text{dom}(\alpha) \leq \beta] \in \mathcal{S} \\ \beta_L \mapsto \text{rng}(\alpha_L) & \forall [\text{rng}(\alpha) \leq \beta] \in \mathcal{S} \end{array}$$

■

The grammar  $G_r(\mathcal{S}, E)$  generates two languages for each set variable that describe the upper and lower non-constant bounds for that set variable. Specifically, if  $\mapsto_G^*$  denotes a derivation in the grammar  $G$ , and  $\mathcal{L}_G(x)$  denotes the language  $\{\tau \mid x \mapsto_G^* \tau\}$  generated by a non-terminal  $x$ , then the following lemma holds.

**Lemma 4.9** *Let  $G = G_r(\mathcal{S}, E)$ . Then:*

$$\begin{aligned} \mathcal{L}_G(\alpha_L) &= \{\kappa \mid [\kappa \leq \alpha] \in \Pi(\mathcal{S}) \text{ and } \text{SetVar}(\kappa) \subseteq E\} \\ \mathcal{L}_G(\alpha_U) &= \{\kappa \mid [\alpha \leq \kappa] \in \Pi(\mathcal{S}) \text{ and } \text{SetVar}(\kappa) \subseteq E\} \end{aligned}$$

**Proof:** See Appendix A.5. ■

## 4.5 Regular Tree Grammars

The grammar  $G_r(\mathcal{S}, E)$  does not describe all constraints in  $\Pi(\mathcal{S}) \mid_E$ . In particular:

- $G_r(\mathcal{S}, E)$  does not describe constraints of the form  $[c \leq \tau]$ . Thus, for example, the regular grammar for the example program component  $P$  does not describe the constraint  $[1 \leq \text{rng}(\alpha^P)]$  in  $\Pi(\mathcal{S}_P) \mid_E$ .
- $G_r(\mathcal{S}, E)$  does not describe constraints inferred by the (*trans <sub>$\alpha$</sub>* ) rule that are not bounds of the form  $[\kappa \leq \alpha]$  or  $[\alpha \leq \kappa]$ . To illustrate this idea, consider the program  $Q = (\lambda^t x. x)$  whose constraint system is:

$$\mathcal{S}_Q = \{t \leq \alpha^Q, \text{dom}(\alpha^Q) \leq \alpha^x, \alpha^x \leq \text{rng}(\alpha^Q)\}.$$

The regular grammar  $G_r(\mathcal{S}_Q, E)$  for  $Q$  describes the constraints  $\{\text{dom}(\alpha^Q) \leq \alpha^x, \alpha^x \leq \text{rng}(\alpha^Q)\}$  in  $\Pi(\mathcal{S}_Q) \mid_E$ , but it does not describe the *trans <sub>$\alpha$</sub>*  consequent  $[\text{dom}(\alpha^Q) \leq \text{rng}(\alpha^Q)]$  of those constraints, which is also in  $\Pi(\mathcal{S}_Q) \mid_E$ .

For an arbitrary constraint system  $\mathcal{S}$ , we represent the constraint system  $\Pi(\mathcal{S}) \mid_E$  by extending the grammar  $G_r(\mathcal{S}, E)$  to a *regular tree grammar*  $G_t(\mathcal{S}, E)$ . The extended grammar combines upper and lower bounds for set variables in the same fashion as the  $(trans_\alpha)$  rule, and also generates constraints of the form  $[c \leq \tau]$  where appropriate.

**Definition 4.10.** (*Regular Tree Grammar  $G_t(\mathcal{S}, E)$* ) The regular tree grammar  $G_t(\mathcal{S}, E)$  extends the grammar  $G_r(\mathcal{S}, E)$  with the root non-terminal  $R$  and the additional productions:

$$\begin{aligned} R &\mapsto [\alpha_L \leq \alpha_U] & \forall \alpha \in SetVar(\mathcal{S}) \\ R &\mapsto [c \leq \alpha_U] & \forall [c \leq \alpha] \in \mathcal{S} \end{aligned}$$

where  $[\cdot \leq \cdot]$  is viewed as a binary constructor. ■

The extended regular tree grammar  $G_t(\mathcal{S}, E)$  describes all constraints in  $\Pi(\mathcal{S}) \mid_E$ .

**Lemma 4.11** *Let  $G = G_t(\mathcal{S}, E)$ . Then  $\Pi(\mathcal{S}) \mid_E = \mathcal{L}_G(R)$ .*

**Proof:** See Appendix A.5. ■

The grammar  $G_t(\mathcal{S}_P, E)$  for the example program component  $P$  is described in figure 6. This grammar yields all constraints in  $\Pi(\mathcal{S}_P) \mid_E$ . For example, the productions:

$$R \mapsto [1 \leq \alpha_U^a] \quad \alpha_U^a \mapsto \mathbf{rng}(\alpha_U^P) \quad \alpha_U^P \mapsto \alpha^P$$

imply that  $R \mapsto^* [1 \leq \mathbf{rng}(\alpha^P)]$ , or that the constant 1 is returned as a possible result of the function  $P$ .

## 4.6 Staging

Before we can exploit the grammar representation of  $\Pi(\mathcal{S}) \mid_E$ , we must prove that the closure under  $\Theta \cup \Pi \cup \{compat\}$  can be performed in a staged manner. The following lemma justifies this staging of the closure algorithm. In particular, it states that  $\Pi$  does not create any additional opportunities for rules in  $\Theta$ , and  $(compat)$  does not create any additional opportunities for  $\Pi$  or  $\Theta$ .

**Lemma 4.12 (Staging)** *For any simple constraint system  $\mathcal{S}$ :*

$$\Psi\Theta(\mathcal{S}) = \Psi(\Theta(\mathcal{S})) = compat(\Pi(\Theta(\mathcal{S})))$$

**Proof:** See Appendix A.5. ■

## 4.7 The Entailment Algorithm

We can check entailment based on lemmas 4.11 and 4.12 as follows. Given  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , we close them under  $\Theta$  and then have:

$$\begin{aligned} & \mathcal{S}_2 \vdash_{\Psi\Theta}^E \mathcal{S}_1 \\ \iff & \Psi\Theta(\mathcal{S}_2) \mid_E \supseteq \Pi\Theta(\mathcal{S}_1) \mid_E & \text{by defn } \vdash_{\Psi\Theta}^E \\ \iff & \Psi(\Theta(\mathcal{S}_2)) \mid_E \supseteq \Pi(\Theta(\mathcal{S}_1)) \mid_E & \text{by lemma 4.12} \\ \iff & \Psi(\mathcal{S}_2) \mid_E \supseteq \Pi(\mathcal{S}_1) \mid_E & \text{as } \mathcal{S}_i = \Theta(\mathcal{S}_i) \\ \iff & compat(\Pi(\mathcal{S}_2) \mid_E) \supseteq \Pi(\mathcal{S}_1) \mid_E & \text{by lemma 4.12} \\ \iff & compat(\mathcal{L}_{G_2}(R)) \supseteq \mathcal{L}_{G_1}(R) & \text{by lemma 4.11} \\ & \text{where } G_i = G_t(\mathcal{S}_i, E) \end{aligned}$$

---

### The Entailment Algorithm

In the following,  $\mathcal{P}_{\text{fin}}$  denotes the finite power-set constructor.

Let:

$$\begin{aligned} G_1 &= G_r(\mathcal{S}_1, E) & L_i &= \{\alpha_L \mid \alpha \in \text{Vars}(\mathcal{S}_i)\} \\ G_2 &= G_t(\mathcal{S}_2, E) & U_i &= \{\alpha_U \mid \alpha \in \text{Vars}(\mathcal{S}_i)\} \end{aligned}$$

Assume  $G_1$  and  $G_2$  are pre-processed to remove  $\epsilon$ -transitions. For  $C \in \mathcal{P}_{\text{fin}}(L_2 \times U_2)$ , define:

$$\mathcal{L}(C) = \{[\tau_L \leq \tau_U] \mid \langle \alpha_L, \beta_U \rangle \in C, \alpha_L \mapsto_{G_2} \tau_L, \beta_U \mapsto_{G_2} \tau_U\}$$

The relation  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\cdot, \cdot, \cdot, \cdot]$  is defined as the largest relation on  $L_1 \times U_1 \times \mathcal{P}_{\text{fin}}(L_2 \times U_2) \times \mathcal{P}_{\text{fin}}(L_2 \times U_2)$  such that if:

$$\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D] \quad \alpha_L \mapsto_{G_1} X \quad \beta_U \mapsto_{G_1} Y$$

then one of the following cases hold:

1.  $\mathcal{L}([X \leq Y]) \subseteq \mathcal{L}(C \cup D)$ .
2.  $X = \text{rng}(\alpha'_L)$ ,  $Y = \text{rng}(\beta'_U)$  and  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha'_L, \beta'_U, C, D']$ , where

$$D' = \{(\gamma'_L, \delta'_U) \mid \langle \gamma_L, \delta_U \rangle \in C \cup D, \gamma_L \mapsto_{G_2} \text{rng}(\gamma'_L), \delta_U \mapsto_{G_2} \text{rng}(\delta'_U)\}$$

3.  $X = \text{dom}(\alpha'_U)$ ,  $Y = \text{dom}(\beta'_L)$  and  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\beta'_L, \alpha'_U, C, D']$ , where

$$D' = \{(\delta'_L, \gamma'_U) \mid \langle \gamma_L, \delta_U \rangle \in C \cup D, \gamma_L \mapsto_{G_2} \text{dom}(\gamma'_U), \delta_U \mapsto_{G_2} \text{dom}(\delta'_L)\}$$

The *computable entailment relation*  $\mathcal{S}_2 \vdash_{\text{alg}}^E \mathcal{S}_1$  holds if and only if  $\forall \alpha \in \text{Vars}(\mathcal{S}_1)$ :

$$\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \alpha_U, \{\langle \gamma_L, \gamma_U \rangle \mid \gamma \in \text{Vars}(\mathcal{S}_2)\}, \emptyset]$$

Figure 7: The computable entailment relation  $\vdash_{\text{alg}}^E$

---

The containment question:

$$\mathcal{L}_{G_2}(R) \supseteq \mathcal{L}_{G_1}(R)$$

can be decided via a standard RTG containment algorithm [?]. To decide the more difficult question:

$$\text{compat}(\mathcal{L}_{G_2}(R)) \supseteq \mathcal{L}_{G_1}(R)$$

we extend the RTG containment algorithm to allow for constraints inferred via (*compat*) on the language  $\mathcal{L}_{G_2}(R)$ .

The extended algorithm is presented in figure 7. It first computes the largest relation  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}$  such that  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D]$  holds if and only if:

$$\mathcal{L}([\alpha_L \leq \beta_U]) \subseteq \text{compat}(\mathcal{L}(C)) \cup \mathcal{L}(D)$$

where  $\alpha_L, \beta_U$  describe collections of types;  $C, D$  describe collections of constraints; and  $\mathcal{L}([\alpha_L \leq \beta_U])$  denotes the language  $\{[\tau_L \leq \tau_U] \mid \alpha_L \mapsto^* \tau_L, \beta_U \mapsto^* \tau_U\}$ . The first case

in the definition of  $\mathcal{R}$  uses an RTG containment algorithm to detect if  $\mathcal{L}([\alpha_L \leq \beta_U]) \subseteq \mathcal{L}(C) \cup \mathcal{L}(D)$ . The two remaining cases handle constraints of the form  $[\text{rng}(\alpha'_L) \leq \text{rng}(\beta'_U)]$  or  $[\text{dom}(\alpha'_U) \leq \text{dom}(\beta'_L)]$ , and allow for inferences via (*compat*). The relation  $\mathcal{R}$  can be computed by starting with a maximal relation (true at every point), and then iteratively setting entries to false, until the largest relation satisfying the definition is reached.

Based on this relation, the algorithm then defines a *computable entailment relation*  $\vdash_{\text{alg}}^E$  on constraint systems. This relation is equivalent to  $\vdash_{\Psi}^E$ .

**Theorem 4.13 (Correctness of the Entailment Algorithm)**  $\mathcal{S}_2 \vdash_{\Psi}^E \mathcal{S}_1$  if and only if  $\mathcal{S}_2 \vdash_{\text{alg}}^E \mathcal{S}_1$ .

**Proof:** See Appendix A.6. ■

The entailment algorithm takes exponential time, since the size of  $\mathcal{R}$  is exponential in the number of set variables in  $\mathcal{S}_2$ . Although faster algorithms for the entailment may exist, these algorithms must all be in PSPACE, because the containment problem on NFA's, which is PSPACE-complete [1], can be polynomially reduced to the entailment problem on constraint systems.

By using the entailment algorithm in both directions, we can now decide if two constraint systems are observable equivalent. Thus, given a constraint system, we can find a minimal, observably equivalent system by systematically generating *all* constraint systems in order of increasing size, until we find one observably equivalent to the original system. Of course, the process of computing the minimal equivalent system with this algorithm is far too expensive for use in practical program analysis systems.

## 4.8 Practical Constraint System Simplification

Fortunately, to take advantage of the rule ( $\cong$ ) in program analysis algorithms, we do not need a completely minimized constraint system. Any *simplifications* in a constraint system produces corresponding reductions in the time and space required for the overall analysis.

For this purpose, we exploit the connection between constraint systems and RTGs. By Lemma 4.11, any transformation on constraint systems that preserves the language:

$$\mathcal{L}_{G_t(\Theta(\mathcal{S}), E)}(R)$$

also preserves the observable behavior of  $\mathcal{S}$  with respect to  $E$ . Based on this observation, we have transformed a variety of existing algorithms for simplifying regular tree grammars to algorithms for simplifying constraint systems. In the following subsections, we present the four most promising algorithms found so far. We use  $G$  to denote  $G_t(\mathcal{S}, E)$ , and we let  $X$  range over non-terminals and  $p$  over *paths*, which are sequences of the constructors **dom** and **rng**. Each algorithm assumes that the constraint system  $\mathcal{S}$  is closed under  $\Theta$ . Computing this closure corresponds to propagating data-flow information locally within a program component. This step is relatively cheap, since program components are typically small (less than a thousand lines of code).

Constraints	Production Rules	Non-empty	Reachable
$f \leq \alpha^f$	$R \mapsto [f \leq \alpha_U^f]$		
$\text{dom}(\alpha^f) \leq \alpha^x$	$\alpha_L^x \mapsto \text{dom}(\alpha_U^f)$		
$1 \leq \alpha^1$	$R \mapsto [1 \leq \alpha_U^1]$	$1 \leq \alpha^1$	$1 \leq \alpha^1$
$\alpha^1 \leq \text{rng}(\alpha^f)$	$\alpha_U^1 \mapsto \text{rng}(\alpha_U^f)$		
$\text{rng}(\alpha^f) \leq \alpha^a$	$\alpha_L^a \mapsto \text{rng}(\alpha_L^f)$		
$\alpha^y \leq \alpha^r$	$\alpha_U^y \mapsto \alpha_U^r$ $\alpha_L^r \mapsto \alpha_L^y$	$\alpha^y \leq \alpha^r$	
$\alpha^r \leq \text{dom}(\alpha^f)$	$\alpha_U^r \mapsto \text{dom}(\alpha_U^f)$		
$g \leq \alpha^P$	$R \mapsto [g \leq \alpha_U^P]$	$g \leq \alpha^P$	$g \leq \alpha^P$
$\text{dom}(\alpha^P) \leq \alpha^y$	$\alpha_L^y \mapsto \text{dom}(\alpha_U^P)$	$\text{dom}(\alpha^P) \leq \alpha^y$	
$\alpha^a \leq \text{rng}(\alpha^P)$	$\alpha_U^a \mapsto \text{rng}(\alpha_U^P)$	$\alpha^a \leq \text{rng}(\alpha^P)$	$\alpha^a \leq \text{rng}(\alpha^P)$
$\alpha^r \leq \alpha^x$	$\alpha_U^r \mapsto \alpha_U^x$ $\alpha_L^x \mapsto \alpha_L^r$	$\alpha^r \leq \alpha^x$	
$\alpha^1 \leq \alpha^a$	$\alpha_U^1 \mapsto \alpha_U^a$ $\alpha_L^a \mapsto \alpha_L^1$	$\alpha^1 \leq \alpha^a$	$\alpha^1 \leq \alpha^a$
$1 \leq \alpha^a$	$R \mapsto [1 \leq \alpha_U^a]$	$1 \leq \alpha^a$	$1 \leq \alpha^a$
	$\alpha_L^P \mapsto \alpha^P$ $\alpha_U^P \mapsto \alpha^P$		

Figure 8: The constraint system, grammar and simplified systems for  $P = (\lambda^g y.((\lambda^f x.1) y))$ 

#### 4.8.1 Empty Constraint Simplification

A non-terminal  $X$  is *empty* if  $\mathcal{L}_G(X) = \emptyset$ . Similarly, a production is *empty* if it refers to empty non-terminals, and a constraint is *empty* if it only induces empty productions. Since empty productions have no effect on the language generated by  $G$ , an empty constraint in  $\mathcal{S}$  can be deleted without changing  $\mathcal{S}$ 's observable behavior.

Let us illustrate this idea with the program component  $P = (\lambda^g y.((\lambda^f x.1) y))$  we considered earlier. Although this example is unrealistic, it illustrates the behavior of our simplification algorithms. The solved constraint system  $\mathcal{S}_P$  for  $P$  is shown in figure 8, together with the corresponding grammar  $G_t(\mathcal{S}_P, E)$  where  $E = \{\alpha^P\}$ . An inspection of this grammar shows that the set of non-empty non-terminals is:

$$\{\alpha_L^P, \alpha_U^P, \alpha_U^y, \alpha_U^a, \alpha_L^r, \alpha_U^1, \alpha_L^x, R\}$$

Five of the constraints in  $\mathcal{S}_P$  are empty, and are removed by this simplification algorithm, yielding a simplified system of eight non-empty constraints.

#### 4.8.2 Unreachable Constraint Simplification

A non-terminal  $X$  is *unreachable* if there is no production  $R \mapsto [Y \leq Z]$  or  $R \mapsto [Z \leq Y]$  such that  $\mathcal{L}_G(Y) \neq \emptyset$  and  $Z \rightarrow_G^* p(X)$ . Similarly, a production is *unreachable* if it refers to unreachable non-terminals, and a constraint is *unreachable* if it only induces unreachable productions. Unreachable productions have no effect on the language  $\mathcal{L}_G(R)$ , and hence unreachable constraints in  $\mathcal{S}$  can be deleted without changing the observable behavior of  $\mathcal{S}$ .

In the above example, the reachable non-terminals are  $\alpha_U^1$ ,  $\alpha_U^a$  and  $\alpha_U^g$ . Three of the constraints are unreachable, and are removed by this algorithm, yielding a simplified system with five reachable constraints.

- 
1. Use a variant of Hopcroft's algorithm [12] to compute an equivalence relation  $\sim$  on the set variables of  $\mathcal{S}$  that satisfies the following conditions:
    - (a) Each set variable in  $E$  is in an equivalence class by itself.
    - (b) If  $[\alpha \leq \beta] \in \mathcal{S}$  then  $\forall \alpha \sim \alpha' \exists \beta \sim \beta'$  such that  $[\alpha' \leq \beta'] \in \mathcal{S}$ .
    - (c) If  $[\alpha \leq \mathbf{rng}(\beta)] \in \mathcal{S}$  then  $\forall \alpha \sim \alpha' \exists \beta \sim \beta'$  such that  $[\alpha' \leq \mathbf{rng}(\beta')] \in \mathcal{S}$ .
    - (d) If  $[\mathbf{rng}(\alpha) \leq \beta] \in \mathcal{S}$  then  $\forall \alpha \sim \alpha' \exists \beta \sim \beta'$  such that  $[\mathbf{rng}(\alpha') \leq \beta'] \in \mathcal{S}$ .
    - (e) If  $[\alpha \leq \mathbf{dom}(\beta)] \in \mathcal{S}$  then  $\forall \alpha \sim \alpha' \forall \beta \sim \beta'$  such that  $[\alpha' \leq \mathbf{dom}(\beta')] \in \mathcal{S}$ .
  2. Merge set variables according to their equivalence class.

Figure 9: The *Hopcroft* algorithm

---

#### 4.8.3 Removing $\epsilon$ -Constraints

A constraint of the form  $[\alpha \leq \beta] \in \mathcal{S}$  is an  $\epsilon$ -constraint. Suppose  $\alpha \notin E$  and the only upper bound on  $\alpha$  in  $\mathcal{S}$  is the  $\epsilon$ -constraint  $[\alpha \leq \beta]$ , i.e., there are no other constraints of the form  $\alpha \leq \tau$ ,  $\mathbf{rng}(\alpha) \leq \gamma$ , or  $\gamma \leq \mathbf{dom}(\alpha)$  in  $\mathcal{S}$ . Then, for any solution  $\rho$  of  $\mathcal{S}$ , the set environment  $\rho'$  defined by:

$$\rho'(\delta) = \begin{cases} \rho(\delta) & \text{if } \delta \not\equiv \alpha \\ \rho(\beta) & \text{if } \delta \equiv \alpha \end{cases}$$

is also a solution of  $\mathcal{S}$ . Therefore we can replace all occurrences of  $\alpha$  in  $\mathcal{S}$  by  $\beta$  while still preserving the observable behavior  $\text{Soln}(\mathcal{S})|_E$ . This substitution transforms the constraint  $[\alpha \leq \beta]$  to the tautology  $[\beta \leq \beta]$ , which can be deleted. Dually, if  $[\alpha \leq \beta] \in \mathcal{S}$  with  $\beta \notin E$  and  $\beta$  having no other lower bounds, then we can replace  $\beta$  by  $\alpha$ , again eliminating the constraint  $[\alpha \leq \beta]$ .

To illustrate this idea, consider the remaining constraints for  $P$ . In this system, the only upper bound for the set variable  $\alpha^1$  is the  $\epsilon$ -constraint  $[\alpha^1 \leq \alpha^a]$ . Hence this algorithm replaces all occurrences of  $\alpha^1$  by  $\alpha^a$ , which further simplifies this constraint system into:

$$\{1 \leq \alpha^a, \alpha^a \leq \mathbf{rng}(\alpha^P), g \leq \alpha^P\}$$

This system is the smallest simple constraint system observably equivalent to the original system  $\Theta(\mathcal{S})$ .

#### 4.8.4 Hopcroft's Algorithm

The previous algorithm *merges* set variables under certain circumstances, and only when they are related by an  $\epsilon$ -constraint. We would like to identify more general circumstances under which set variables can be merged. To this end, we define a *valid unifier* for  $\mathcal{S}$  to be an equivalence relation  $\sim$  on the set variables of  $\mathcal{S}$  such that we can merge the set variables in each equivalence class of  $\sim$  without changing the observable behavior of  $\mathcal{S}$ . Using a model-theoretic argument, we can show that an equivalence relation  $\sim$  is a valid unifier for  $\mathcal{S}$  if for all solutions  $\rho \in \text{Soln}(\mathcal{S})$  there exists another solution  $\rho' \in \text{Soln}(\mathcal{S})$  such that  $\rho'$  agrees with  $\rho$  on  $E$  and  $\rho'(\alpha) = \rho'(\beta)$  for all  $\alpha \sim \beta$ .

A natural strategy for generating  $\rho'$  from  $\rho$  is to map each set variable to the least upper bound of the set variables in its equivalence class:

$$\rho'(\alpha) = \bigsqcup_{\alpha' \sim \alpha} \rho(\alpha')$$

Figure 9 describes sufficient conditions to ensure that  $\rho'$  is a solution of  $\mathcal{S}$ , and hence that  $\sim$  is a valid unifier for  $\mathcal{S}$ . To produce an equivalence relation satisfying these conditions, we use a variant of Hopcroft’s  $O(n \lg n)$  time algorithm [12] for computing an equivalence relation on states in a DFA and then merge set variables according to their equivalence class.<sup>6</sup>

The following theorem states that this simplification algorithm preserves the observable behavior of constraint systems.

**Theorem 4.14 (Correctness of the Hopcroft Algorithm)** *Let  $\mathcal{S}$  be a simple constraint system with external variables  $E$ ; let  $\sim$  be an equivalence relation on the set variables in a constraint system  $\mathcal{S}$  satisfying conditions (a) to (e) from figure 9; let the substitution  $f$  map each set variable to a representation element of its equivalence class; and let  $\mathcal{S}' = f(\mathcal{S})$ , i.e.,  $\mathcal{S}'$  denotes the constraint system  $\mathcal{S}$  with set variables merged according to their equivalence class. Then  $\mathcal{S} \cong_E \mathcal{S}'$ .*

**Proof:** See Appendix A.7. ■

## 4.9 Simplification Benchmarks

Definition	lines	size	empty		unreachable		$\epsilon$ -removal		Hopcroft	
			factor	time	factor	time	factor	time	factor	time
map	5	221	3	<10	6	20	11	30	13	30
reverse	6	287	4	<10	8	20	20	10	20	30
substring	8	579	12	10	64	10	64	10	96	20
qsort	41	1387	15	<10	15	30	58	50	66	40
unify	89	2921	10	10	11	80	55	120	65	150
hopcroft	201	8429	25	10	42	100	118	100	124	200
check	237	21854	4	50	4	1150	26	370	168	510
escher-fish	493	30509	187	10	678	40	678	40	678	80
scanner	1209	59215	3	180	17	840	45	2450	57	2120

Figure 10: Behavior of the constraint simplification algorithms.

To test the effectiveness of the simplification algorithms, we extended MrSpidey with the four algorithms that we have just described: *empty*, *unreachable*,  $\epsilon$ -removal, and *Hopcroft*. Each algorithm also implements the preceding simplification strategies. The first three algorithms are linear in the number of non-empty constraints in the system, and *Hopcroft* is log-linear.

<sup>6</sup>A similar development based on the definition  $\rho'(\alpha) = \sqcap \{\rho(\alpha') \mid \alpha \sim \alpha'\}$  results in an alternative algorithm, which is less effective in practice.

We tested the algorithms on the constraint systems for nine program components on a 167MHz Sparc Ultra 1 with 326M of memory, using the MzScheme compiler [9]. The results are described in figure 10. The second column gives the number of lines in each program component, and the third column gives the number of constraints in the original (unsimplified) constraint system after closing it under the rules  $\Theta$ . The remaining columns describe the behavior of each simplification algorithm, presenting the factor by which the number of constraints was reduced, and the time (in milliseconds) required for this simplification. Since MzScheme is a byte code compiler, porting the simplification algorithms to a native code compiler could be expected to produce a speed-up of roughly a factor of 5.

The results demonstrate the effectiveness and efficiency of our simplification algorithms. The resulting constraint systems are typically at least an order of magnitude smaller than the original system. The cost of these algorithms is reasonable, particularly considering that they were run on a byte code compiler. As expected, the more sophisticated algorithms are more effective, but are also more expensive.

## 5 Componential Set-Based Analysis

Equipped with the simplification algorithms, we can now return to our original problem of extending set-based analysis to handle significantly larger programs. These programs are typically constructed as a collection of program components (*e.g.* modules, packages or files). Exploiting this component-based structure is the key to analyzing such programs efficiently.

The following section describes componential set-based analysis. Section 5.2 presents experimental results on the effectiveness of the analysis, and section 5.3 describes how MrSpidey presents the analysis results for multi-component programs to the programmer.

The constraint simplification algorithms also enables an efficient polymorphic, or context-sensitive, analysis that only duplicates a simplified constraint system for each reference to a polymorphic function. A description of this polymorphic analysis is presented in section ??, together with experimental results on the behavior of the analysis.

### 5.1 Componential Set-Based Analysis

Componential set-based analysis processes programs in three steps.

1. For each component in the program, the analysis derives and simplifies the constraint system for that component and saves the simplified system in a *constraint file*, for use in later runs of the analysis. The simplification is performed with respect to the external variables of the component, *excluding* expression labels, in order to minimize the size of the simplified system. Thus, the simplified system only needs to describe how the component interacts with the rest of the program, and the simplification algorithm can discard constraints that are only necessary to infer local value set invariants. These discarded constraints are reconstructed later as needed.

This step can be skipped for each program component that has not changed since the last run of the analysis, and the component's constraint file can be used instead.



2. The analysis combines the simplified constraint systems of the *entire* program and closes the combined collection of constraints under  $\Theta$ , thus propagating data-flow information between the constraint systems for the various program components.
3. Finally, to reconstruct the full analysis results for the program component that the programmer is focusing on, the analysis tool combines the constraint system from the second step with the unsimplified constraint system for that component. It closes the resulting system under  $\Theta$ , which yields appropriate value set invariants for each labeled expression in the component.

The new analysis can easily process programs that consist of many components. For its first step, it eliminates all those constraints that have only local relevance, thus producing a small combined constraint system for the entire program. As a result, the analysis tool can solve the combined system more quickly and using less space than traditional set-based analysis [10]. Finally, it recreates as much precision as traditional set-based analysis as needed on a per-component basis.

The new analysis performs extremely well in an interactive setting because it exploits the saved constraint files where possible and thus avoids re-processing many program components unnecessarily.

## 5.2 Experimental Results

We implemented four variants of componential set-based analysis. Each analysis uses a particular simplification algorithm from chapter ?? to simplify the constraint systems for the program components. We tested these analyses with five benchmark programs, ranging from 1,200 to 17,000 lines. For comparison purposes, we also analyzed each benchmark with the *standard* set-based analysis that performs no simplification. The analyses handled library functions in a context-sensitive, polymorphic manner according to the constraint derivation rules (*let*) and (*inst*) to avoid merging information between unrelated calls to these functions. The remaining functions were analyzed in a context-insensitive, monomorphic manner. The results are documented in figure 11.

The fourth column in the figure shows the maximum size of the constraint system generated by each analysis, and also shows this size as a percentage of the constraint system generated by the *standard* analysis. The analyses based on the simplification algorithms produce significantly smaller constraint systems, and can also analyze more programs, such as *sba* and *poly*, for which the *standard* analysis exhausted heap space.

The fifth column shows the time required to analyze each program from scratch, without using any existing constraint files.<sup>7</sup> The analyses that exploit constraint simplification yield significant speed-ups over the *standard* analysis because they manipulate much smaller constraint systems. The results indicate that, for these benchmarks, the  $\epsilon$ -removal algorithm yields the best trade-off between efficiency and effectiveness of the simplification algorithms. The additional simplification performed by the more expensive *Hopcroft* algorithm is outweighed by the overhead of running the algorithm. The tradeoff may change as we analyze larger programs.

---

<sup>7</sup>These times exclude scanning and parsing time.

Program	lines	Analysis	Num. of constraints	Analysis time	Re-analysis time	Constraint file size (bytes)
<b>scanner</b>	1253	<i>standard</i>	61K	14.1s	7.7s	572K
		<i>empty</i>	24K (39%)	12.0s	3.1s	189K
		<i>unreachable</i>	15K (25%)	9.7s	2.0s	39K
		<i><math>\epsilon</math>-removal</i>	14K (23%)	9.5s	1.7s	28K
		<i>Hopcroft</i>	14K (23%)	10.4s	1.7s	25K
<b>zodiac</b>	3419	<i>standard</i>	704K	133.4s	110.6s	1634K
		<i>empty</i>	62K (9%)	34.1s	8.1s	328K
		<i>unreachable</i>	21K (3%)	28.8s	4.5s	169K
		<i><math>\epsilon</math>-removal</i>	13K (2%)	28.8s	3.8s	147K
		<i>Hopcroft</i>	11K (2%)	31.4s	3.8s	136K
<b>nucleic</b>	3432	<i>standard</i>	333K	83.9s	51.2s	2882K
		<i>empty</i>	90K (27%)	52.8s	17.8s	592K
		<i>unreachable</i>	68K (20%)	48.4s	14.6s	386K
		<i><math>\epsilon</math>-removal</i>	56K (17%)	48.3s	13.1s	330K
		<i>Hopcroft</i>	56K (17%)	60.9s	13.2s	328K
<b>sba</b>	11560	<i>standard</i>	*, >5M	*	*	*
		<i>empty</i>	1908K (<38%)	181.5s	65.5s	1351K
		<i>unreachable</i>	105K (<2%)	149.5s	43.3s	920K
		<i><math>\epsilon</math>-removal</i>	76K (<2%)	147.1s	42.2s	770K
		<i>Hopcroft</i>	65K (<1%)	156.8s	41.1s	716K
<b>mod-poly</b>	17661	<i>standard</i>	*, >5M	*	*	*
		<i>empty</i>	*, >5M	*	*	*
		<i>unreachable</i>	201K (<4%)	259.6s	26.9s	1517K
		<i><math>\epsilon</math>-removal</i>	68K (<1%)	239.6s	13.3s	1038K
		<i>Hopcroft</i>	38K (<1%)	254.1s	10.9s	907K

\* indicates the analysis exhausted heap space

Figure 11: Behavior of the modular analyses.

To test the responsiveness of the componential analyses in an interactive setting based on an analyze-debug-edit cycle, we re-analyzed each benchmark after changing a randomly chosen component in that benchmark. The re-analysis times are shown in the sixth column of figure 11. These times show an order-of-magnitude improvement in analysis times over the original, *standard* analysis, since the saved constraint files are used to avoid reanalyzing all of the unchanged program components. For example, the analysis of **zodiac**, which used to take over two minutes, now completes in under four seconds. Since practical debugging sessions using MrSpidey typically involve repeatedly analyzing the project each time the source code of one module is modified, *e.g.*, when a bug is identified and eliminated, using separate analysis substantially improves the usability of MrSpidey.

The disk-space required to store the constraint files is shown in column seven. Even though these files use a straight-forward, text-based representation, their size is typically within a factor of two or three of the corresponding source file.

---

Program	lines	<i>copy</i> analysis	Relative time of <i>smart</i> polymorphic analyses				Mono. analysis
			<i>empty</i>	<i>unreachable</i>	<i>ε-removal</i>	<i>Hopcroft</i>	
<b>lattice</b>	215	4.2s	39%	36%	35%	38%	42%
<b>browse</b>	233	2.5s	76%	76%	76%	81%	75%
<b>splay</b>	265	7.9s	75%	73%	70%	72%	83%
<b>check</b>	281	50.1s	21%	23%	14%	14%	23%
<b>graphs</b>	621	2.8s	85%	85%	82%	87%	82%
<b>boyer</b>	624	4.3s	46%	46%	49%	50%	40%
<b>matrix</b>	744	7.5s	64%	57%	51%	52%	45%
<b>maze</b>	857	6.2s	64%	59%	58%	61%	54%
<b>nbody</b>	880	39.6s	57%	25%	25%	26%	28%
<b>nucleic</b>	3335	*	* 243s	* 42s	* 42s	* 44s	* 36s

\* indicates the *copy* analysis exhausted heap space,  
and the table contains absolute times for the other analyses

Figure 12: Times for the smart polymorphic analyses, relative to the *copy* analysis.

---

### 5.3 User Interface for Multi-File Programs

We extended MrSpidey’s user interface to cope with programs consisting of multiple source files, or components. MrSpidey first analyses the program, using the componential set-based analysis described above, and then displays an annotated version of the program’s main file with the usual static debugging mark-ups. The programmer can also view annotated versions of any other source file by using the **File|Open . . .** dialog box (shown in Figure ??) to select the source file of interest.

In multi-file programs, the source (or destination) of an arrow may sometimes refer to a program point in a separate file. In this case MrSpidey draws an arrow originating (or terminating) in the left margin of the program, as shown in figure ?. Clicking on the arrow provides the option to zoom to and highlight the term at the other end of the arrow, as shown in figures ? and ?. These facilities are useful for following the flow of values through multi-file programs.

### 5.4 Efficient Polymorphic Analysis

The constraint simplification algorithms also enables an efficient polymorphic, or context-sensitive, analysis. To avoid merging information between unrelated calls to functions that are used in a polymorphic fashion, a polymorphic analysis duplicates the function’s constraints at each call site. We extended MrSpidey with five polymorphic analyses. The first analysis is *copy*, which duplicates the constraint system for each polymorphic reference via a straightforward implementation of the rules (*let*) and (*inst*).<sup>8</sup> The remaining four analyses are *smart* analyses that simplify the constraint system for each polymorphic definition.

We tested the analyses using a standard set of benchmarks [13]. The results of the test runs are documented in figure 12. The second column shows the number of lines in each

---

<sup>8</sup>We also implemented a polymorphic analysis that re-analyzes a definition at each reference, but found its performance to be comparable to, and sometimes worse than, the *copy* analysis.

benchmark; the third column presents the time for the *copy* analysis; and columns four to seven show the times for each smart polymorphic analysis, as a percentage of the *copy* analysis time. For comparison purposes, the last column shows the relative time of the original, but less accurate, monomorphic analysis.

The results again demonstrate the effectiveness of our constraint simplification algorithms. The smart analyses that exploit constraint simplification are always significantly faster and can analyze more programs than the *copy* analysis. For example, while *copy* exhausts heap space on the **nucleic** benchmark, all smart analyses successfully analyzed this benchmark.

Again, it appears that the  $\epsilon$ -removal analysis yields the best trade-off between efficiency and effectiveness of the simplification algorithms. This analysis provides the additional accuracy of polymorphism without much additional cost over the coarse, monomorphic analysis. With the exception of the benchmarks **browse**, **splay** and **graphs**, which do not re-use many functions in a polymorphic fashion, this analysis is a factor of 2 to 4 times faster than the *copy* analysis, and it is also capable of analyzing larger programs.

## 6 Competitive Work

A number of researchers have investigated the problem of constraint simplification in order to derive faster and more scalable analyses and type systems.

Deutsch and Heintze [?] examine constraint simplification for set-based analysis. They discover two simplification algorithms, which are analogous to our empty and unreachable constraint simplification algorithms, but do not present results on the cost or effectiveness of these simplification algorithms.

Fähndrich and Aiken [6] examine constraint simplification for an analysis based on a more complex constraint language. They develop a number of heuristic algorithms for constraint simplification, which they test on programs of up to 6000 lines. Their fastest approach yields a factor of 3 saving in both time and space, but is slow in absolute times compared to other program analyses.

Pottier [17] studies an ML-style language with subtyping. Performing type inference on this language produces subtype constraints that are similar to our constraints. Pottier defines an entailment relation on constraints, and presents an incomplete algorithm for deciding entailment. In addition, he proposes some *ad hoc* algorithms for simplifying constraints. He does not report any results on the cost or effectiveness of these algorithms.

Trifonov and Smith [20] describe a subtyping relation between constrained types, which are similar to our constraint systems, and they present an incomplete decision algorithm for subtyping. They do not discuss constraint simplification. Eifrig, Smith and Trifonov [5] discuss constraint simplification in the context of type inference for objects. They present three algorithms for simplifying constraint systems, two of which are similar to the *empty* and  $\epsilon$ -removal algorithms, and the third is a special case of the *Hopcroft* algorithm. They do not present results on the cost or effectiveness of these algorithms.

Duesterwald *et al* [4] describe algorithms for simplifying data flow equations. These algorithms are similar to the  $\epsilon$ -removal and *Hopcroft* algorithms. Their approach only preserves the greatest solution of the equation system and assumes that the control flow graph is already known. Hence it cannot be used to analyze programs in a componential manner

or to analyze programs with advanced control-flow mechanisms such as first-class functions and virtual methods. The paper does not present results on the cost or effectiveness of these algorithms.

## 7 Future Work

All our constraint simplification algorithms preserve the observable behavior of constraint systems, and thus do not effect the accuracy of the analysis. If we were willing to tolerate a less accurate analysis, we could choose a compressed constraint system that does not preserve the observable behavior of the original, but only *entails* that behavior. This approach allows the use of much smaller constraint systems, and hence yields a faster analysis.

A promising approach for deriving such approximate constraint systems is to rely on a programmer-provided *signature* describing the behavior of a polymorphic function or module, and to derive the new constraint system from that signature. After checking the entailment condition to verify that signature-based constraints correctly approximates the behavior of the module, we could use those constraints in the remainder of the analysis. Since the signature-based constraints are expected to be smaller than the derived ones, this approach could significantly reduce analysis times for large projects. We are investigating this approach for developing a typed module language on top of Scheme.

## A Proofs

### A.1 Subject Reduction Proof

**Lemma 3.4** (*Subject Reduction for  $\longrightarrow$* ). *If  $\Gamma \vdash M_1 : \alpha, \mathcal{S}_1$  and  $M_1 \longrightarrow M_2$ , then  $\Gamma \vdash M_2 : \alpha, \mathcal{S}_2$  such that  $\mathcal{S}_1 \models \mathcal{S}_2$ .*

**Proof:** The proof proceeds by case analysis according to the relation  $M_1 \longrightarrow M_2$ .

- Suppose  $M_1 \longrightarrow M_2$  via  $(\beta_v)$ . Then:

$$\begin{aligned} M_1 &= ((\lambda^t x. N) V) \\ M_2 &= N[x \mapsto V] \end{aligned}$$

The typing derivation on  $M_1$  is of the form:

$$\frac{\frac{\Gamma \cup \{x : \beta_x\} \vdash N : \beta_N, \mathcal{S}_N}{\Gamma \vdash (\lambda^t x. N) : \beta_t, \mathcal{S}_t} \quad (abs) \quad \Gamma \vdash V : \beta_V, \mathcal{S}_V}{\Gamma \vdash M_1 : \alpha, \mathcal{S}_1} \quad (app)$$

where:

$$\begin{aligned} \mathcal{S}_t &= \mathcal{S}_N \cup \{t \leq \beta_t, \text{dom}(\beta_t) \leq \beta_x, \beta_N \leq \text{rng}(\beta_t)\} \\ \mathcal{S}_1 &= \mathcal{S}_t \cup \mathcal{S}_V \cup \{\beta_V \leq \text{dom}(\beta_t), \text{rng}(\beta_t) \leq \alpha\} \end{aligned}$$

By the Substitution Lemma A.1:

$$\Gamma \vdash M_2 : \beta_N, \mathcal{S}'_2 \quad \text{where} \quad \mathcal{S}_N \cup \mathcal{S}_V \cup \{\beta_V \leq \beta_x\} \models \mathcal{S}'_2$$

By the Flow Lemma 3.6:

$$\Gamma \vdash M_2 : \alpha, \mathcal{S}_2 \quad \text{where} \quad \mathcal{S}'_2 \cup \{\beta_N \leq \alpha\} \models \mathcal{S}_2$$

Since  $\mathcal{S}_1 \supseteq \mathcal{S}_N \cup \mathcal{S}_V$ , we have that  $\mathcal{S}_1 \models \mathcal{S}_N \cup \mathcal{S}_V$ . Also:

$$\begin{aligned} \mathcal{S}_1 &\supseteq \{\beta_V \leq \text{dom}(\beta_t), \text{dom}(\beta_t) \leq \beta_x\} \models \{\beta_V \leq \beta_x\} \\ \mathcal{S}_1 &\supseteq \{\beta_N \leq \text{rng}(\beta_t), \text{rng}(\beta_t) \leq \alpha\} \models \{\beta_N \leq \alpha\} \end{aligned}$$

Hence  $\mathcal{S}_1 \models \mathcal{S}_2$ , as required.

- Suppose  $M_1 \longrightarrow M_2$  via  $(\beta_{let})$ . Then:

$$\begin{aligned} M_1 &= (\mathbf{let} \ (x \ V) \ N) \\ M_2 &= N[x \mapsto V] \end{aligned}$$

The typing derivation on  $M_1$  is of the form:

$$\frac{\begin{array}{c} \Gamma \vdash V : \alpha_V, \mathcal{S}_V \\ \overline{\alpha} = \text{SetVar}(\mathcal{S}_V) \setminus (FV[\text{range}(\Gamma)] \cup \text{Label}) \\ \Gamma \cup \{x : \forall \overline{\alpha}. (\alpha_V, \mathcal{S}_V)\} \vdash N : \beta, \mathcal{S} \end{array}}{\Gamma \vdash M_1 : \beta, \mathcal{S}} \quad (let)$$

By the Subject Reduction for **let** Lemma A.2:

$$\Gamma \vdash M_2 : \beta, \mathcal{S}$$

as required.

- Suppose  $M_1 \longrightarrow M_2$  via  $(unlabel)$ . Then:

$$\begin{aligned} M_1 &= V^l \\ M_2 &= V \end{aligned}$$

The typing derivation on  $M_1$  is of the form:

$$\frac{\Gamma \vdash V : \alpha, \mathcal{S}_V}{\Gamma \vdash V^l : \beta, \mathcal{S}_V \cup \{\alpha \leq l, \alpha \leq \beta\}} \quad (label)$$

Hence  $\Gamma \vdash V : \alpha, \mathcal{S}_V$ , and by the Flow Lemma 3.6,  $\Gamma \vdash V : \beta, \mathcal{S}_2$  where  $\mathcal{S}_1 \models \mathcal{S}_2$ .

■

**Lemma A.1 (Substitution)** *If*

$$\begin{aligned} \Gamma \cup \{x : \alpha_x\} &\vdash N : \alpha_N, \mathcal{S}_N \\ \Gamma &\vdash V : \alpha_V, \mathcal{S}_V \end{aligned}$$

*then*

$$\Gamma \vdash N[x \mapsto V] : \alpha_N, \mathcal{S} \quad \text{where} \quad \mathcal{S}_N \cup \mathcal{S}_V \cup \{\alpha_V \leq \alpha_x\} \models \mathcal{S}$$

**Proof:** The proof proceeds by induction on the number of **let**-expressions in  $N$ , and on the size of  $N$ :

$$\Gamma \cup \{x : \alpha_x\} \vdash N : \alpha_N, \mathcal{S}_N$$

If  $x \notin FV[N]$ , then  $N[x \mapsto V] = N$ ,  $\mathcal{S} = \mathcal{S}_N$  and the lemma trivially holds.

Otherwise we proceed by case analysis on the constraint derivation rule used in the last step in the derivation.

- (*var*): Since  $x \in FV[N]$ ,  $N = x$ . Hence  $\mathcal{S}_N = \{\alpha_x \leq \alpha_N\}$ ,  $N[x \mapsto V] = V$  and this case holds via the Flow Lemma 3.6.
- (*const*): This case cannot occur since  $x \in FV[N]$ .
- (*label*): This case is straightforward.
- (*abs*): In this case  $N = (\lambda^t y. M)$ , and the constraint derivation is:

$$\frac{\Gamma \cup \{x : \alpha_x, y : \alpha_y\} \vdash M : \alpha_M, \mathcal{S}_M}{\Gamma \cup \{x : \alpha_x\} \vdash (\lambda^t y. M) : \alpha_N, \mathcal{S}_N} \quad (abs)$$

where:

$$\mathcal{S}_N = \mathcal{S}_M \cup \{t \leq \alpha_N, \text{dom}(\alpha_N) \leq \alpha_y, \alpha_M \leq \text{rng}(\alpha_N)\}$$

Since  $x \in FV[N]$ ,  $x \neq y$ . Hence  $N[x \mapsto V] = (\lambda^t y. M[x \mapsto V])$ . By induction:

$$\Gamma \cup \{x : \alpha_x, y : \alpha_y\} \vdash M[x \mapsto V] : \alpha_M, \mathcal{S}'_M$$

where:

$$\mathcal{S}_M \cup \mathcal{S}_V \cup \{\alpha_V \leq \alpha_x\} \models \mathcal{S}'_M$$

Hence, via (*abs*):

$$\Gamma \vdash N[x \mapsto V] : \alpha_N, \mathcal{S}$$

where:

$$\mathcal{S} = \mathcal{S}'_M \cup \{t \leq \alpha_N, \text{dom}(\alpha_N) \leq \alpha_y, \alpha_M \leq \text{rng}(\alpha_N)\}$$

Since  $\mathcal{S}_N \cup \mathcal{S}_V \cup \{\alpha_V \leq \alpha_x\} \models \mathcal{S}$ , the lemma holds for this case.

- (*app*): In this case  $M = (M_1 \ M_2)$ , and the constraint derivation is:

$$\frac{\Gamma \cup \{x : \alpha_x\} \vdash M_i : \beta_i, \mathcal{S}_i}{\Gamma \cup \{x : \alpha_x\} \vdash (M_1 \ M_2) : \alpha_N, \mathcal{S}_N} \quad (app)$$

where

$$\mathcal{S}_N = \mathcal{S}_1 \cup \mathcal{S}_2 \cup \{\beta_2 \leq \text{dom}(\beta_1), \text{rng}(\beta_1) \leq \alpha_N\}$$

By induction,

$$\Gamma \vdash M_i[x \mapsto V] : \beta_i, \mathcal{S}'_i$$

where

$$\mathcal{S}_i \cup \mathcal{S}_V \cup \{\alpha_V \leq \alpha_x\} \models \mathcal{S}'_i$$

Hence

$$\Gamma \vdash M[x \mapsto V] : \alpha_N, \mathcal{S}$$

where

$$\mathcal{S} = \mathcal{S}'_1 \cup \mathcal{S}'_2 \cup \{\beta_2 \leq \text{dom}(\beta_1), \text{rng}(\beta_1) \leq \alpha_N\}$$

Obviously,  $\mathcal{S}_N \cup \mathcal{S}_V \cup \{\alpha_V \leq \alpha_x\} \models \mathcal{S}$ , and the lemma holds for this case.

- (*let*): In this case  $N = (\text{let } (y \ W) \ M)$ , where  $W \in \text{Value}$ . Hence:

$$\Gamma \cup \{x : \alpha_x\} \vdash (\text{let } (y \ W) \ M) : \alpha_N, \mathcal{S}_N$$

and therefore by the following Subject Reduction for **let** Lemma A.2:

$$\Gamma \cup \{x : \alpha_x\} \vdash M[y \mapsto W] : \alpha_N, \mathcal{S}_N$$

By induction:

$$\Gamma \vdash M[y \mapsto W][x \mapsto V] : \alpha_N, \mathcal{S}$$

where  $\mathcal{S}_N \cup \mathcal{S}_V \cup \{\alpha_V \leq \alpha_x\} \models \mathcal{S}$ . Since:

$$M[y \mapsto W][x \mapsto N] = M[x \mapsto N][y \mapsto W[x \mapsto N]]$$

we have that:

$$\Gamma \vdash M[x \mapsto V][y \mapsto W[x \mapsto V]] : \alpha_N, \mathcal{S}$$

and therefore, by the Subject Reduction for **let** Lemma A.2:

$$\Gamma \vdash (\text{let } (y \ W[x \mapsto V]) \ M[x \mapsto V]) : \alpha_N, \mathcal{S}$$

or, equivalently:

$$\Gamma \vdash (\text{let } (y \ W) \ M)[x \mapsto V] : \alpha_N, \mathcal{S}$$

and thus the lemma holds in this case.

- (*inst*): This case cannot occur since  $x \in FV[N]$  and  $x$  is bound to a set variable in the derivation context.

■

The constraint derivation rules uses constraint schemas to accurately analyze polymorphic **let**-expressions. The constraint system for a **let**-expression is actually *equivalent* to the constraint system for the corresponding  $\beta_{let}$ -expanded expression<sup>9</sup>. The following lemma demonstrates this equivalence of constraint systems.

**Lemma A.2 (Subject Reduction for let)**

$$\Gamma \vdash (\text{let } (x \ V) \ N) : \alpha_N, \mathcal{S}_N$$

*if and only if*

$$\Gamma \vdash N[x \mapsto V] : \alpha_N, \mathcal{S}_N$$

---

<sup>9</sup>This equivalence contrasts with the situation for the other reduction rules, where the constraint system for the redex only *entails* the constraint system for the contractum, as shown in the Subject Reduction for  $\rightarrow$  Lemma 3.4.



**Proof:** The derivation  $\Gamma \vdash (\mathbf{let} (x \ V) \ N) : \alpha_N, \mathcal{S}_N$  holds if and only if:

$$\begin{aligned} & \Gamma \vdash V : \alpha_V, \mathcal{S}_V \\ & \overline{\alpha} = \text{SetVar}(\mathcal{S}_V) \setminus (FV[\text{range}(\Gamma)] \cup \text{Label}) \\ & \sigma = \forall \overline{\alpha}. (\alpha_V, \mathcal{S}_V) \\ & \Gamma \cup \{x : \sigma\} \vdash N : \alpha_N, \mathcal{S}_N \end{aligned}$$

The proof of both directions proceeds by induction on the number of **let**-expressions in  $N$ , and on the size of  $N$ :

If  $x \notin FV[N]$ , then  $N[x \mapsto V] = N$  and the lemma trivially holds.

Otherwise we proceed by case analysis on the constraint derivation rule used in the last step in the derivation:

$$\Gamma \cup \{x : \sigma\} \vdash N : \alpha_N, \mathcal{S}_N$$

- (*var*): This case cannot occur since  $x \in FV[N]$  implies  $N = x$ , but  $x$  is bound to a schema in the derivation context and so the rule (*inst*) applies.
- (*const*): This case cannot occur since  $x \in FV[N]$ .
- (*abs*): In this case  $N = (\lambda^t y. M)$ , and the typing derivation is:

$$\frac{\Gamma \cup \{x : \sigma, y : \alpha_y\} \vdash M : \alpha_M, \mathcal{S}_M}{\Gamma \cup \{x : \sigma\} \vdash (\lambda^t y. M) : \alpha_N, \mathcal{S}_N} \quad (\text{abs})$$

where:

$$\mathcal{S}_N = \mathcal{S}_M \cup \{t \leq \alpha_N, \text{dom}(\alpha_N) \leq \alpha_y, \alpha_M \leq \text{rng}(\alpha_N)\}$$

Since  $x \in FV[N]$ ,  $x \neq y$ . Hence  $N[x \mapsto V] = (\lambda^t y. M[x \mapsto V])$ . By induction,

$$\Gamma \cup \{x : \sigma, y : \alpha_y\} \vdash M[x \mapsto V] : \alpha_M, \mathcal{S}_M$$

Hence, via (*abs*):

$$\Gamma \vdash N[x \mapsto V] : \alpha_N, \mathcal{S}_N$$

The reasoning for the converse direction is similar.

- (*app*): In this case  $M = (M_1 \ M_2)$ , and the typing derivation is:

$$\frac{\Gamma \cup \{x : \sigma\} \vdash M_i : \beta_i, \mathcal{S}_i}{\Gamma \cup \{x : \sigma\} \vdash (M_1 \ M_2) : \alpha_N, \mathcal{S}_N} \quad (\text{app})$$

where:

$$\mathcal{S}_N = \mathcal{S}_1 \cup \mathcal{S}_2 \cup \{\beta_2 \leq \text{dom}(\beta_1), \text{rng}(\beta_1) \leq \alpha_N\}$$

By induction,

$$\Gamma \vdash M_i[x \mapsto V] : \beta_i, \mathcal{S}_i$$

Hence

$$\Gamma \vdash M[x \mapsto V] : \alpha_N, \mathcal{S}_N$$

as required. The reasoning for the converse direction is similar.

- (*label*): This case is straightforward.
- (*let*): In this case  $N = (\mathbf{let} (y W) M)$ , where  $W \in Value$ . Hence:

$$\Gamma \cup \{x : \sigma\} \vdash (\mathbf{let} (y W) M) : \alpha_N, \mathcal{S}_N$$

Since  $M$  has fewer **let**-expressions than  $N$ , by induction:

$$\Gamma \cup \{x : \sigma\} \vdash M[y \mapsto W] : \alpha_N, \mathcal{S}_N$$

Since  $M[y \mapsto W]$  has fewer **let**-expressions than  $N$ , by induction:

$$\Gamma \vdash M[y \mapsto W][x \mapsto V] : \alpha_N, \mathcal{S}_N$$

Since:

$$M[y \mapsto W][x \mapsto N] = M[x \mapsto N][y \mapsto W[x \mapsto N]]$$

we have that:

$$\Gamma \vdash M[x \mapsto V][y \mapsto W[x \mapsto V]] : \alpha_N, \mathcal{S}_N$$

Since  $M[x \mapsto V][y \mapsto W[x \mapsto V]]$  is smaller than  $N$ , by induction:

$$\Gamma \vdash (\mathbf{let} (y W[x \mapsto V]) M[x \mapsto V]) : \alpha_N, \mathcal{S}_N$$

or, equivalently:

$$\Gamma \vdash (\mathbf{let} (y W) M)[x \mapsto V] : \alpha_N, \mathcal{S}_N$$

and thus the lemma holds in this case.

- (*inst*): Since  $x \in FV[N]$ ,  $N = x$ , and the derivation on  $N$  must be:

$$\Gamma \cup \{x : \sigma\} \vdash x : \alpha_N, \mathcal{S}_N$$

where  $\psi$  is a substitution of fresh variables for  $\bar{\alpha}$  and  $\mathcal{S}_N = \psi(\mathcal{S}_V) \cup \{\psi(\alpha_V) \leq \alpha_N\}$ .

If  $\mathcal{D}$  is the derivation concluding

$$\Gamma \vdash V : \alpha_V, \mathcal{S}_V$$

then  $\psi(\mathcal{D})$  is an analogous derivation concluding

$$\Gamma \vdash V : \psi(\alpha_V), \psi(\mathcal{S}_V)$$

Now  $N[x \mapsto V] = V$ , and by the Flow Lemma 3.6:

$$\Gamma \vdash V : \alpha_N, \mathcal{S}_N$$

as required.

■

**Lemma 3.6 (Flow).** *If  $\Gamma \vdash M : \alpha, \mathcal{S}$  then for all  $\gamma \in SetVar$ ,  $\Gamma \vdash M : \gamma, \mathcal{S}'$  with  $\mathcal{S} \cup \{\alpha \leq \gamma\} \models \mathcal{S}'$ .*

**Proof:** By induction on the derivation  $\Gamma \vdash M : \alpha, \mathcal{S}$  and by case analysis on the last step in this derivation.

- (*var*): In this case the derivation for  $M = x$  is:

$$\Gamma' \cup \{x : \beta\} \vdash x : \alpha, \mathcal{S}$$

where  $\mathcal{S} = \{\beta \leq \alpha\}$ . For any  $\gamma \in \text{SetVar}$ , let  $\mathcal{S}' = \{\beta \leq \gamma\}$ , and then:

$$\Gamma' \cup \{x : \beta\} \vdash x : \gamma, \mathcal{S}'$$

with  $\mathcal{S} \cup \{\alpha \leq \gamma\} \models \mathcal{S}'$ , as required.

- (*const*): This case follows by reasoning similar to the (*var*) case.
- (*label*): The derivation for  $M = N^l$  must conclude:

$$\frac{\Gamma \vdash N : \beta, \mathcal{S}_N}{\Gamma \vdash N^l : \alpha, \mathcal{S}} \quad (\text{label})$$

where  $\mathcal{S} = \mathcal{S}_N \cup \{\beta \leq l, \beta \leq \alpha\}$ . Let  $\mathcal{S}' = \{\beta \leq l, \beta \leq \gamma\}$ , and then:

$$\frac{\Gamma \vdash N : \beta, \mathcal{S}_N}{\Gamma \vdash N^l : \gamma, \mathcal{S}'} \quad (\text{label})$$

with  $\mathcal{S} \cup \{\alpha \leq \gamma\} \models \mathcal{S}'$ .

- (*abs*): The derivation for  $M = (\lambda^t x. N)$  must conclude:

$$\frac{\Gamma \cup \{x : \alpha_1\} \vdash N : \alpha_2, \mathcal{S}_N}{\Gamma \vdash (\lambda^t x. N) : \alpha, \mathcal{S}} \quad (\text{abs})$$

where  $\mathcal{S} = \{t \leq \alpha, \text{dom}(\alpha) \leq \alpha_1, \alpha_2 \leq \text{rng}(\alpha)\}$ .

Let  $\mathcal{S} = \{t \leq \gamma, \text{dom}(\gamma) \leq \alpha_1, \alpha_2 \leq \text{rng}(\gamma)\}$ , and then:

$$\frac{\Gamma \cup \{x : \alpha_1\} \vdash N : \alpha_2, \mathcal{S}_N}{\Gamma \vdash (\lambda^x N.) : \gamma, \mathcal{S}} \quad (\text{abs})$$

with  $\mathcal{S} \cup \{\alpha \leq \gamma\} \models \mathcal{S}'$ .

- (*app*): This case follows by reasoning similar to the (*app*) case.
- (*let*): This case follows by induction.
- (*abs*): The derivation for  $M = x$  must be:

$$\Gamma \cup \{x : \forall \bar{\alpha}. (\alpha_V, \mathcal{S}_V)\} \vdash x : \alpha, \mathcal{S}$$

where  $\mathcal{S} = \psi(\mathcal{S}_V) \cup \{\psi(\alpha_V) \leq \alpha\}$ , and  $\psi$  is a substitution of set variables for  $\bar{\alpha}$ .

Let  $\mathcal{S}' = \psi(\mathcal{S}_V) \cup \{\psi(\alpha_V) \leq \gamma\}$ . Then

$$\Gamma \cup \{x : \forall \bar{\alpha}. (\alpha_V, \mathcal{S}_V)\} \vdash x : \gamma, \mathcal{S}'$$

with  $\mathcal{S} \cup \{\alpha \leq \gamma\} \models \mathcal{S}'$ .

■

## A.2 Proofs for Computing Set-Based Analysis

**Lemma 3.9** (*Least Solution of Simple Constraint Systems*). *Every simple constraint system has a solution that is least with respect to  $\sqsubseteq_s$ .*

**Proof:** Let  $\mathcal{S}$  be a simple constraint system, and define  $\rho = \sqcap_s \text{Soln}(\mathcal{S})$ , using the pointwise extension of  $\sqcap_s$  to set environments. We prove that  $\rho \in \text{Soln}(\mathcal{S})$  by showing that  $\rho$  satisfies any constraint  $\mathcal{C} \in \mathcal{S}$ . The proof proceeds by case analysis on  $\mathcal{C}$ .

- The case where  $\mathcal{C} = [\alpha \leq \beta]$  follows from Lemma A.3.
- Suppose  $\mathcal{C} = [c \leq \beta]$ . Then  $c \in \text{const}(\rho'(\beta))$  for all  $\rho' \in \text{Soln}(\mathcal{S})$ , therefore  $c \in \text{const}(\rho(\beta))$  and  $\rho \models \mathcal{C}$ .
- Suppose  $\mathcal{C} = [\alpha \leq \text{rng}(\beta)]$ . Then:

$$\begin{aligned}
 \rho(\alpha) &= \sqcap_{\rho' \in \text{Soln}(\mathcal{S})} \rho'(\alpha) \\
 &\sqsubseteq \sqcap_{\rho' \in \text{Soln}(\mathcal{S})} \rho'(\text{rng}(\beta)) \\
 &\quad \text{by Lemma A.3, since } \rho'(\alpha) \sqsubseteq \rho'(\text{rng}(\beta)) \\
 &= \text{rng} \left( \sqcap_{\rho' \in \text{Soln}(\mathcal{S})} \rho'(\beta) \right) \quad \text{by definition of } \sqcap_s \\
 &= \text{rng}(\rho(\beta)) \\
 &= \rho(\text{rng}(\beta))
 \end{aligned}$$

Hence  $\rho \models \alpha \leq \text{rng}(\beta)$ .

The remaining cases are similar. Hence the set of environments satisfying  $\mathcal{S}$  has a least element  $\sqcap_s \{\rho \mid \rho \models \mathcal{S}\}$ . ■

The following lemma describes some properties about how the two orderings  $\sqsubseteq$  and  $\sqsubseteq_s$  defined on  $\mathcal{D}$  interact.

**Lemma A.3** *Let  $I$  be an index set, and let  $x_i, y_i \in \mathcal{D}$  for all  $i \in I$ .*

- If  $x_i \sqsubseteq y_i$  for all  $i \in I$ , then:

$$\sqcap_{i \in I} x_i \sqsubseteq \sqcap_{i \in I} y_i$$

$$\sqcup_{i \in I} x_i \sqsubseteq \sqcup_{i \in I} y_i$$

- If  $x_i \sqsubseteq_s y_i$  for all  $i \in I$ , then:

$$\sqcap_{i \in I} x_i \sqsubseteq_s \sqcap_{i \in I} y_i$$

$$\sqcup_{i \in I} x_i \sqsubseteq_s \sqcup_{i \in I} y_i$$

**Proof:** The proof is based on the interpretation of  $\mathcal{D}$  as the set of total functions

$$f : \{\mathbf{dom}, \mathbf{rng}\}^* \longrightarrow \mathcal{P}(\mathit{Const})$$

and proceeds by showing the appropriate relation holds between the sets of constant elements at any path in  $\{\mathbf{dom}, \mathbf{rng}\}^*$ .

To prove the first relation, assume  $x_i \sqsubseteq y_i$  for all  $i \in I$ , and let  $p$  be a path in  $\{\mathbf{dom}, \mathbf{rng}\}^*$ . If  $p$  is monotonic, then  $p(x_i) \sqsubseteq p(y_i)$ . Hence:

$$\begin{aligned} p \left( \bigcap_{i \in I} x_i \right) &= \bigcap_{i \in I} p(x_i) \\ &\sqsubseteq \bigcap_{i \in I} p(y_i) \\ &= p \left( \bigcap_{i \in I} y_i \right) \end{aligned}$$

Conversely, if  $p$  is anti-monotonic, then  $p(x_i) \supseteq p(y_i)$ . Hence:

$$\begin{aligned} p \left( \bigcap_{i \in I} x_i \right) &= \bigcap_{i \in I} p(x_i) \\ &\supseteq \bigcap_{i \in I} p(y_i) \\ &= p \left( \bigcap_{i \in I} y_i \right) \end{aligned}$$

Hence

$$\bigcap_{i \in I} x_i \sqsubseteq \bigcap_{i \in I} y_i$$

as required. ■

**Theorem 3.12** *If  $P \in \Lambda^0$  and  $\emptyset \vdash P : \alpha, \mathcal{S}$  is a most general constraint derivation then:*

$$\begin{aligned} sba(P)(l) = & \{b \mid \mathcal{S} \vdash_{\Theta} b \leq l\} \\ & \cup \{(\lambda^t x.M) \mid \mathcal{S} \vdash_{\Theta} t \leq l\} \end{aligned}$$

**Proof:**

$$\begin{aligned} \mathcal{S} \vdash_{\Theta} c \leq \alpha &\iff \mathcal{S} \models c \leq \alpha && \text{by lemma A.4} \\ &\iff \forall \rho \models \mathcal{S}. \rho \models c \leq \alpha \\ &\iff \forall \rho \in \mathit{Soln}(\mathcal{S}). c \in \mathit{const}(\rho(\alpha)) \\ &\iff c \in \bigcap (\{\mathit{const}(\rho(\alpha)) \mid \rho \in \mathit{Soln}(\mathcal{S})\}) \\ &\iff c \in \mathit{const}(\bigcap (\{\rho(\alpha) \mid \rho \in \mathit{Soln}(\mathcal{S})\})) \\ &\iff c \in \mathit{const}(\bigcap (\{\rho \mid \rho \in \mathit{Soln}(\mathcal{S})\})(\alpha)) \\ &\iff c \in \mathit{const}(\mathit{LeastSoln}(\mathcal{S})(\alpha)) \end{aligned}$$

The correctness of this theorem then follows from definition 3.10. ■

**Lemma A.4 (Soundness and Completeness of  $\Theta$ )** For any simple constraint system  $\mathcal{S}$ :

$$\mathcal{S} \models c \leq \alpha \iff \mathcal{S} \vdash_{\Theta} c \leq \alpha$$

**Proof:** The soundness of  $\Theta$  is straightforward. To prove the completeness of  $\Theta$ , assume  $\mathcal{S} \models c \leq \alpha$ . Let  $\rho$  be any fixpoint of the functional  $F$  defined as:

$$\begin{aligned} F : \text{SetEnv} &\longrightarrow \text{SetEnv} \\ F(\rho)(\alpha) &= \langle \{c \mid \mathcal{S} \vdash_{\Theta} c \leq \alpha\}, \\ &\quad \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \alpha \leq^* \delta, \gamma \leq \text{dom}(\delta)\}, \\ &\quad \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\alpha)\} \rangle \end{aligned}$$

where the notation  $\mathcal{S} \vdash_{\Theta} \beta \leq^* \delta$  means there exists some  $\delta_1, \dots, \delta_n$  with  $\beta \equiv \delta_1$  and  $\delta_n \equiv \delta$  such that:

$$\mathcal{S} \vdash_{\Theta} \{\delta_i \leq \delta_{i+1}, \mid 1 \leq i < n\}$$

The asymmetry between the definition of the domain and range components  $F(\rho)(\alpha)$  arises from the rules  $\Theta$ . These rules propagate set variables denoting the result of functions in  $\alpha$  forward along data-flow paths into constraints of the form  $\gamma \leq \text{rng}(\alpha)$ . However, the same propagation does not occur for set variables denoting argument values to functions in  $\alpha$ , and hence this propagation is performed in the definition of  $F(\rho)(\alpha)$  by finding all  $\gamma$  such that  $\gamma \leq \text{dom}(\delta)$  and  $\alpha \leq^* \delta$ .

If  $\rho \models \mathcal{S}$ , then  $\rho \models c \leq \alpha$  and hence  $\mathcal{S} \vdash_{\Theta} c \leq \alpha$  by the definition of  $\rho$ , as required. Thus it just remains to prove that  $\rho \models \mathcal{S}$ . We proceed by case analysis on constraints  $\mathcal{C} \in \mathcal{S}$ .

- Suppose  $\mathcal{C} = [\alpha \leq \beta]$ . We need to show that the correct ordering holds between the corresponding components of  $\rho(\alpha)$  and  $\rho(\beta)$ . For the first component, by  $(s_1)$ , which is the first rule in  $\Theta$ :

$$\begin{aligned} \{c \mid \mathcal{S} \vdash_{\Theta} c \leq \alpha\} &\subseteq \{c \mid \mathcal{S} \vdash_{\Theta} c \leq \beta\} \\ \therefore \text{const}(\rho(\alpha)) &\sqsubseteq \text{const}(\rho(\beta)) \end{aligned}$$

For the second (domain) component, by  $(s_3)$ :

$$\begin{aligned} [\mathcal{S} \vdash_{\Theta} \beta \leq^* \delta] &\Rightarrow [\mathcal{S} \vdash_{\Theta} \alpha \leq^* \delta] \\ \therefore \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \beta \leq^* \delta, \gamma \leq \text{dom}(\delta)\} &\subseteq \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \alpha \leq^* \delta, \gamma \leq \text{dom}(\delta)\} \\ \therefore \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \beta \leq^* \delta, \gamma \leq \text{dom}(\delta)\} &\sqsubseteq \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \alpha \leq^* \delta, \gamma \leq \text{dom}(\delta)\} \\ \therefore \text{dom}(\rho(\beta)) &\sqsubseteq \text{dom}(\rho(\alpha)) \end{aligned}$$

For the third (range) component, by  $(s_2)$ :

$$\begin{aligned} [\mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\alpha)] &\Rightarrow [\mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\beta)] \\ \therefore \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\alpha)\} &\subseteq \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\beta)\} \\ \therefore \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\alpha)\} &\sqsubseteq \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\beta)\} \\ \therefore \text{rng}(\rho(\alpha)) &\sqsubseteq \text{rng}(\rho(\beta)) \end{aligned}$$

Hence  $\rho(\alpha) \sqsubseteq \rho(\beta)$ .

- Suppose  $\mathcal{C} = [c \leq \beta]$ :

$$\begin{aligned} \rho(\beta) &\sqsupseteq \langle \{c \mid \mathcal{S} \vdash_{\Theta} c \leq \beta\}, \top, \perp \rangle \\ &\sqsupseteq \langle \{c\}, \top, \perp \rangle \\ &= \rho(c) \end{aligned}$$

- Suppose  $\mathcal{C} = [\alpha \leq \text{rng}(\beta)]$ .

$$\begin{aligned} \{\rho(\alpha)\} &\subseteq \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\beta)\} \\ \therefore \rho(\alpha) &\sqsubseteq \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\beta)\} \\ &= \rho(\text{rng}(\beta)) \end{aligned}$$

- Suppose  $\mathcal{C} = [\text{rng}(\alpha) \leq \beta]$ . Then

$$\begin{aligned} \rho(\beta) &\sqsupseteq \rho(\gamma) && \forall \mathcal{S} \vdash_{\Theta} \gamma \leq \beta \\ \therefore \rho(\beta) &\sqsupseteq \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \beta\} \\ &\sqsupseteq \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{rng}(\alpha)\} && \text{by } (s_4) \\ &= \rho(\text{rng}(\alpha)) \end{aligned}$$

- Suppose  $\mathcal{C} = [\alpha \leq \text{dom}(\beta)]$ . Then

$$\begin{aligned} \rho(\alpha) &\sqsubseteq \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \text{dom}(\beta)\} \\ &\sqsubseteq \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \beta \leq^* \delta, \gamma \leq \text{dom}(\delta)\} \\ &= \rho(\text{dom}(\beta)) \end{aligned}$$

- Suppose  $\mathcal{S} = [\text{dom}(\alpha) \leq \beta]$ .

$$\begin{aligned} \rho(\text{dom}(\alpha)) &= \bigsqcup \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \alpha \leq^* \delta, \gamma \leq \text{dom}(\delta)\} && \text{by } (s_3) \\ &\sqsubseteq \bigsqcup_s \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \text{dom}(\delta) \leq \beta, \gamma \leq \text{dom}(\delta)\} \\ &\sqsubseteq \bigsqcup_s \{\rho(\gamma) \mid \mathcal{S} \vdash_{\Theta} \gamma \leq \beta\} \\ &\sqsubseteq \rho(\beta) \end{aligned}$$

Hence  $\rho \models \mathcal{S}$ , and the lemma holds. ■

### A.3 Proofs for Conditions for Constraint Simplification

The following lemma demonstrates that the rule  $(\cong)$  is *admissible* in that any derivation in the extended constraint derivation system produces information equivalent to that produced by the original analysis.

**Lemma 4.1** (*Admissibility of  $(\cong)$* ). *If  $\emptyset \vdash_{\cong} P : \alpha, \mathcal{S}$  is a most general constraint derivation then:*

$$sba(P)(l) = \text{const}(\text{LeastSoln}(\mathcal{S})(l))$$

**Proof:** This lemma follows from the induction hypothesis:

If  $\Gamma \vdash_{\cong} M : \alpha, \mathcal{S}_1$ , and  $E = FV[\text{range}(\Gamma)] \cup \{\alpha\} \cup \text{Label}$ , then there exists  $\mathcal{S}_2$  such that  $\Gamma \vdash M : \alpha, \mathcal{S}_2$  and  $\mathcal{S}_1 \cong_E \mathcal{S}_2$ .

We prove this hypothesis by induction on the derivation  $\Gamma \vdash_{\cong} M : \alpha, \mathcal{S}_1$ , and by case analysis on the last step in the derivation.

- If the last step in the derivation  $\Gamma \vdash_{\cong} M : \alpha, \mathcal{S}_1$  uses a derivation rule other than  $(\cong)$ , then the lemma holds based on the induction hypothesis.
- Suppose  $\Gamma \vdash_{\cong} M : \alpha, \mathcal{S}_1$  via  $(\cong)$  because  $\Gamma \vdash_{\cong} M : \alpha, \mathcal{S}_3$  and  $\mathcal{S}_3 \cong_E \mathcal{S}_1$ . By induction,  $\Gamma \vdash M : \alpha, \mathcal{S}_4$  where  $\mathcal{S}_3 \cong_E \mathcal{S}_4$ . Since  $\cong_E$  is an equivalence relation,  $\mathcal{S}_1 \cong_E \mathcal{S}_2$ , and hence the lemma holds.

■

#### A.4 Proofs for Proof Theory of Observable Equivalence

The following proofs require a number of auxiliary definitions.

**Definition A.5.** (*Paths*)

- A *path*  $p, q \in \text{Path}$  is a sequence of the constructors **dom** and **rng**. We use  $\epsilon$  to denote the empty sequence, and  $p.q$  to denote the concatenation of the paths  $p$  and  $q$ .
- The *arity* of a path  $p$ , denoted  $\pi p$ , is the number of **dom**'s in  $p$ , taken modulo 2. If  $\pi p$  is 0, we say  $p$  is *monotonic*, otherwise  $p$  is *anti-monotonic*.
- For a path  $p$ , the notation  $p(\tau)$  denotes the set expression  $\tau$  enclosed in the **dom**'s and **rng**'s of  $p$ , *i.e.*, if  $p \equiv \text{rng.dom}$ , then  $p(\alpha) \equiv \text{rng}(\text{dom}(\alpha))$ .
- The relations  $\leq_0$  and  $\leq_1$  denote  $\leq$  and  $\geq$ , respectively.
- The relations  $\sqsubseteq_0$  and  $\sqsubseteq_1$  denote  $\sqsubseteq$  and  $\sqsupseteq$ , respectively.
- The relations  $\sqsupseteq_0$  and  $\sqsupseteq_1$  denote  $\sqsupseteq$  and  $\sqsubseteq$ , respectively.
- The relations  $\subseteq_0$  and  $\subseteq_1$  denote  $\subseteq$  and  $\supseteq$ , respectively.
- The operations  $\sqcup_0$  and  $\sqcup_1$  denote  $\sqcup$  and  $\sqcap$ , respectively.
- For a path  $p$  and a domain element  $X \in \mathcal{D}$ , the notation  $p(X)$  extracts the component of  $X$  at the position  $p$ . This notation is formalized as follows:

$$\begin{aligned} \epsilon(X) &= X \\ (\text{rng}.p)(X) &= \text{rng}(p(X)) \\ (\text{dom}.p)(X) &= \text{dom}(p(X)) \end{aligned}$$



- For a path  $p$  and a domain element  $X \in \mathcal{D}$ , the notation  $X@p$  is defined as follows:

$$\begin{aligned}
\cdot @ \cdot : \mathcal{D} \times path &\longrightarrow \mathcal{D} \\
X @ \epsilon &= X \\
X @ (\text{dom}.p) &= \langle \emptyset, X, \perp_s \rangle @ p \\
X @ (\text{rng}.p) &= \langle \emptyset, \perp_s, X \rangle @ p
\end{aligned}$$

■

**Lemma 4.2** (*Soundness and Completeness of  $\Delta$* ). For a compound constraint system  $\mathbf{S}$  and a compound constraint  $\mathbf{C}$ :

$$\mathbf{S} \vdash_{\Delta} \mathbf{C} \iff \mathbf{S} \models \mathbf{C}$$

**Proof:** The soundness of  $\Delta$  is straightforward. To demonstrate the completeness of  $\Delta$ , we assume  $\mathbf{S} \models \mathbf{C}$  and prove that  $\mathbf{S} \vdash_{\Delta} \mathbf{C}$  by case analysis on  $\mathbf{C}$ .

- Suppose  $\mathbf{C} = [c \leq \kappa]$ . Define  $\rho$  by:

$$\forall p \in Path. \forall \alpha \in SetVar. const(p(\rho(\alpha))) = \{c \mid \mathbf{S} \vdash_{\Delta} c \leq p(\alpha)\}$$

We prove  $\rho \models \mathbf{S}$  by a case analysis showing that  $\rho$  satisfies every constraint  $\mathbf{C}' \in \mathbf{S}$ .

- Suppose  $\mathbf{C}' = [c \leq q(\beta)]$ . Then, by the definition of  $\rho$ ,  $c \in const(\rho(q(\beta)))$ , and hence  $\rho \models c \leq q(\beta)$ .
- Suppose  $\mathbf{C}' = [p(\alpha) \leq q(\beta)]$ . We need to show that  $\rho(p(\alpha)) \sqsubseteq \rho(q(\beta))$ . We prove this inequality by showing that for any path  $r$ :

$$const(r(\rho(p(\alpha)))) \subseteq_{\pi_r} const(r(\rho(q(\beta))))$$

If  $r$  is monotonic, then:

$$\begin{aligned}
&const(r(\rho(p(\alpha)))) \\
&= const(r(p(\rho(\alpha)))) \\
&= \{c \mid \mathbf{S} \vdash_{\Delta} c \leq r(p(\alpha))\} \\
&\subseteq \{c \mid \mathbf{S} \vdash_{\Delta} c \leq r(q(\beta))\} \\
&\quad \text{via } (trans_{\tau}), \text{ since } [p(\alpha) \leq q(\beta)] \in \mathbf{S} \\
&\quad \text{and hence } \mathbf{S} \vdash_{\Delta} r(p(\alpha)) \leq r(q(\beta)) \text{ via } (compat) \\
&= const(r(q(\rho(\beta)))) \\
&= const(r(\rho(q(\beta))))
\end{aligned}$$

The case where  $r$  is anti-monotonic follows by a similar argument.

Hence  $\rho \models \mathbf{S}$ . But since  $\mathbf{S} \models c \leq \kappa$ ,  $\rho \models c \leq \kappa$ . Since  $\kappa = p(\alpha)$  for some  $p$  and  $\alpha$ , then we have that:

$$\begin{aligned}
c &\in const(\rho(p(\alpha))) \\
&= const(p(\rho(\alpha))) \\
&= \{c \mid \mathbf{S} \vdash_{\Delta} c \leq p(\alpha)\}
\end{aligned}$$

Hence,  $\mathbf{S} \vdash_{\Delta} c \leq \kappa$ , as required.

- Suppose  $\mathbf{C} = [\kappa_1 \leq \kappa_2]$ . Let  $c$  be a constant not used in  $\mathbf{S}$  or  $\mathbf{C}$ ; let  $\mathbf{S}' = \mathbf{S} \cup \{c \leq \kappa_1\}$ ; and let  $\rho = \text{LeastSoln}(\mathbf{S}')$ . Since  $\rho \models \mathbf{C}$ , we have that:

$$\rho \models \{c \leq \kappa_1, \kappa_1 \leq \kappa_2\}$$

Hence  $\rho \models c \leq \kappa_2$  and by the first part of this proof,  $\mathbf{S}' \vdash_{\Delta} c \leq \kappa_2$ .

We now show that for any  $\kappa'$ ,  $\mathbf{S}' \vdash_{\Delta} c \leq \kappa'$  if and only if  $\mathbf{S} \vdash_{\Delta} \kappa_1 \leq \kappa'$ . We prove this hypothesis by induction on the derivation of  $\mathbf{S}' \vdash_{\Delta} c \leq \kappa'$ .

- Suppose  $\mathbf{S}' \vdash_{\Delta} c \leq \kappa'$  because  $[c \leq \kappa'] \in \mathbf{S}'$ . Then  $\kappa' = \kappa_1$ , and by the rule (*reflex*),  $\mathbf{S} \vdash_{\Delta} \kappa_1 \leq \kappa_1$ , as required.
- If  $[c \leq \kappa'] \notin \mathbf{S}'$ , then  $\mathbf{S}' \vdash_{\Delta} c \leq \kappa'$  must be derived via the rule (*trans <sub>$\tau$</sub>* ) based on the antecedents  $\mathbf{S}' \vdash_{\Delta} \{c \leq \kappa'', \kappa'' \leq \kappa'\}$ . By induction,  $\mathbf{S} \vdash_{\Delta} \kappa_1 \leq \kappa''$ . Hence  $\mathbf{S} \vdash_{\Delta} \kappa_1 \leq \kappa'$  via (*trans <sub>$\tau$</sub>* ), as required.

Since  $\mathbf{S}' \vdash_{\Delta} c \leq \kappa_2$ , the above induction hypothesis implies that  $\mathbf{S} \vdash_{\Delta} \kappa_1 \leq \kappa_2$ , as required.

■

**Lemma 4.3** *For a compound constraint system  $\mathbf{S}$ ,  $\mathbf{S} \cong_E \Delta(\mathbf{S}) \mid_E$ .*

**Proof:** We need to show that  $\mathbf{S} \cong_E \Delta(\mathbf{S}) \mid_E$ , i.e.:

$$\text{Soln}(\mathbf{S}) \mid_E = \text{Soln}(\Delta(\mathbf{S}) \mid_E) \mid_E$$

Since the rules  $\Delta$  are sound:

$$\begin{aligned} \text{Soln}(\mathbf{S}) \mid_E &= \text{Soln}(\Delta(\mathbf{S})) \mid_E \\ &\subseteq \text{Soln}(\Delta(\mathbf{S}) \mid_E) \mid_E \end{aligned}$$

because the solution space increases as the constraints  $\Delta(\mathbf{S})$  are restricted to  $E$ .

To show the containment in the other direction, assume  $\rho \models \Delta(\mathbf{S}) \mid_E$ . Without loss of generality, assume  $\rho(\alpha) = \perp_s$  for all  $\alpha \notin E$ . We extend  $\rho$  to a super-environment  $\rho'$  that satisfies  $\mathbf{S}$  as follows:

$$\forall p \in \text{Path}. \forall \alpha \in \text{SetVar}. \text{const}(p(\rho'(\alpha))) = \bigcup \{ \text{const}(\rho(\tau)) \mid \mathbf{S} \vdash_{\Delta} \tau \leq p(\alpha) \}$$

We show that  $\rho' \models \mathbf{S}$  by case analysis on the constraints  $\mathbf{C} \in \mathbf{S}$ .

- Suppose  $\mathbf{C} = [c \leq q(\beta)]$ . Then

$$\begin{aligned} \text{const}(q(\beta)) &= \bigcup \{ \text{const}(\rho(\tau)) \mid \mathbf{S} \vdash_{\Delta} \tau \leq q(\beta) \} \\ &\supseteq \{c\} \end{aligned}$$

as required.

- Suppose  $\mathbf{C} = [p(\alpha) \leq q(\beta)]$ . Then for any path  $r$ ,  $\mathbf{S} \vdash_{\Delta} r(p(\alpha)) \leq_{\pi r} r(q(\beta))$ . Hence:

$$\begin{aligned} & \bigcup \{ \text{const}(\rho(\tau)) \mid \mathbf{S} \vdash_{\Delta} \tau \leq r(p(\alpha)) \} \\ \subseteq_{\pi r} & \bigcup \{ \text{const}(\rho(\tau)) \mid \mathbf{S} \vdash_{\Delta} \tau \leq r(q(\beta)) \} \end{aligned}$$

Therefore:

$$\text{const}(\rho'(r(p(\alpha)))) \subseteq_{\pi r} \text{const}(\rho'(r(q(\beta))))$$

Hence:

$$\rho'(p(\alpha)) \sqsubseteq \rho'(q(\beta))$$

And hence  $\rho' \models \mathbf{C}$ , as required.

Thus  $\rho' \models \mathbf{S}$ . It remains to show that  $\rho$  and  $\rho'$  agree on  $E$ . Let  $\alpha \in E$  and  $r \in \text{Path}$ . Then:

$$\begin{aligned} \text{const}(\rho'(r(\alpha))) &= \bigcup \{ \text{const}(\rho(\tau)) \mid \mathbf{S} \vdash_{\Delta} \tau \leq r(\alpha) \} \\ &\quad \text{by definition of } \rho' \\ &= \bigcup \{ \text{const}(\rho(\tau)) \mid \mathbf{S} \vdash_{\Delta} \tau \leq r(\alpha), \text{SetVar}(\tau) \subseteq E \} \\ &\quad \text{since } \rho(\beta) = \perp_s \text{ for } \beta \notin E \\ &\quad \text{and hence } \rho(\tau) = \perp_s \text{ for } \text{SetVar}(\tau) \not\subseteq E \\ &= \bigcup \{ \text{const}(\rho(\tau)) \mid \tau \leq r(\alpha) \in \Delta(\mathbf{S}) \mid_E \} \\ &= \text{const}(\rho(r(\alpha))) \end{aligned}$$

since  $[r(\alpha) \leq r(\alpha)] \in \Delta(\mathbf{S}) \mid_E$  by *(reflex)* and *(compat)*, and for  $[\tau \leq r(\alpha)] \in \Delta(\mathbf{S}) \mid_E$ ,  $\text{const}(\rho(\tau)) \subseteq \text{const}(\rho(r(\alpha)))$ . Thus  $\rho$  and  $\rho'$  agree on  $E$ , and the lemma holds. ■

**Lemma A.6** *For any  $p \in \text{Path}$  and  $X \in \mathcal{D}$ ,  $p(X @ p) = X$ .*

**Proof:** By induction on the length of  $p$ , and by case analysis on the top constructor in  $p$ . ■

**Lemma 4.4** *(Equivalence of Proof Systems). For a simple constraint system  $\mathcal{S}$ :*

$$\Delta(\mathcal{S}) = \Psi\Theta(\mathcal{S})$$

**Proof:** We show that  $\Psi\Theta(\mathcal{S}) \subseteq \Delta(\mathcal{S})$  by induction on the derivation of  $\mathbf{C} \in \Psi\Theta(\mathcal{S})$ . For the base case, if  $\mathbf{C} \in \Psi\Theta(\mathcal{S})$  because  $\mathbf{C} \in \mathcal{S}$ , then  $\mathbf{C} \in \Delta(\mathcal{S})$ . Otherwise we proceed by case analysis on the last rule used in the derivation of  $\mathbf{C} \in \Psi\Theta(\mathcal{S})$ .

- *(compose<sub>1</sub>)*: In this case  $\mathbf{C} = [\alpha \leq \text{rng}(\kappa)]$  is derived from the antecedents  $\{\alpha \leq \text{rng}(\beta), \beta \leq \kappa\} \in \Psi\Theta(\mathcal{S})$ . By induction, these antecedents are also in  $\Delta(\mathcal{S})$ , and hence the following derivation shows that  $\mathbf{C} \in \Delta(\mathcal{S})$ :

$$\frac{\alpha \leq \text{rng}(\beta) \quad \frac{\beta \leq \kappa}{\text{rng}(\beta) \leq \text{rng}(\kappa)} \quad (\text{compat})}{\alpha \leq \text{rng}(\kappa)} \quad (\text{trans}_{\tau})$$

- $(compose_2)$ ,  $(compose_3)$ ,  $(compose_3)$ : These cases follow by similar reasoning.
- $(reflex)$ ,  $(trans_\alpha)$ ,  $(compat)$ : These rules are either equivalent to or subsumed by corresponding rules in  $\Delta$ .
- $(s_1)$ ,  $(s_4)$ ,  $(s_5)$ : For these cases  $\mathbf{C} \in \Delta(\mathcal{S})$  via  $(trans_\tau)$ .
- $(s_2)$ ,  $(s_3)$ : These rules are special cases of the rules  $(compose_1)$  and  $(compose_4)$ , respectively.

There are no other possibilities for the derivation  $\mathbf{C} \in \Psi\Theta(\mathcal{S}) \setminus \Delta(\mathcal{S})$ , and hence  $\Psi\Theta(\mathcal{S}) \subseteq \Delta(\mathcal{S})$ .

We prove the converse inclusion  $\Delta(\mathcal{S}) \subseteq \Psi\Theta(\mathcal{S})$  by induction on the derivation of  $\mathbf{C} \in \Delta(\mathcal{S})$ . Again, for the base case, if  $\mathbf{C} \in \Delta(\mathcal{S})$  because  $\mathbf{C} \in \mathcal{S}$ , then  $\mathbf{C} \in \Psi\Theta(\mathcal{S})$ . Otherwise we proceed by case analysis on the last rule used in the derivation of  $\mathbf{C} \in \Delta(\mathcal{S})$ .

- $(reflex)$ ,  $(compat)$ : These rules are also in  $\Psi$  and, by induction, the antecedents are in  $\Psi\Theta(\mathcal{S})$ , hence  $\mathbf{C} \in \Psi\Theta(\mathcal{S})$ .
- $(trans_\tau)$ : The last step in the derivation must be:

$$\frac{\tau_1 \leq \tau \quad \tau \leq \tau_2}{\tau_1 \leq \tau_2} \quad (trans_\tau)$$

We proceed by case analysis on  $\tau$  to show that  $[\tau_1 \leq \tau_2] \in \Psi\Theta(\mathcal{S})$ .

- The case  $\tau = c$  is impossible, since  $[\tau_1 \leq c]$  is not a compound constraint.
- If  $\tau \in SetVar$ , then  $[\tau_1 \leq \tau_2] \in \Psi\Theta(\mathcal{S})$  via  $(trans_\alpha)$ .
- Suppose  $\tau = \mathbf{rng}(\tau')$ . If  $\tau' \in SetVar$  then  $[\tau_1 \leq \tau_2] \in \Psi\Theta(\mathcal{S})$  via  $(s_4)$ .

Otherwise  $\tau_1 \leq \mathbf{rng}(\tau')$  and  $\mathbf{rng}(\tau') \leq \tau_2$  are not simple constraints, and we proceed by considering the derivation of these constraints in  $\Delta(\mathcal{S})$ . The last step in the derivation of  $\tau_1 \leq \mathbf{rng}(\tau')$  is either via:

$$\frac{\alpha_1 \leq \mathbf{rng}(\beta_1) \quad \beta_1 \leq \tau'}{\alpha_1 \leq \mathbf{rng}(\tau')} \quad (compose_1)$$

where  $\tau_1 = \alpha_1$ , or:

$$\frac{\tau'_1 \leq \tau'}{\mathbf{rng}(\tau'_1) \leq \mathbf{rng}(\tau')} \quad (compat)$$

where  $\tau_1 = \mathbf{rng}(\tau'_1)$ . Similarly, the last step in the derivation of  $\mathbf{rng}(\tau') \leq \tau_2$  is either via:

$$\frac{\tau' \leq \beta_2 \quad \mathbf{rng}(\beta_2) \leq \alpha_2}{\mathbf{rng}(\tau') \leq \alpha_2} \quad (compose_3)$$

where  $\tau_2 = \alpha_2$ , or:

$$\frac{\tau' \leq \tau'_2}{\mathbf{rng}(\tau') \leq \mathbf{rng}(\tau'_2)} \quad (compat)$$

where  $\tau_2 = \mathbf{rng}(\tau'_2)$ . We consider the four possible combinations for the derivations of  $\tau_1 \leq \mathbf{rng}(\tau')$  and  $\mathbf{rng}(\tau') \leq \tau_2$ :

- \* Suppose  $\tau_1 \leq \mathbf{rng}(\tau')$  is inferred via  $(compose_1)$  and  $\mathbf{rng}(\tau') \leq \tau_2$  is inferred via  $(compose_3)$ . Then  $\{\beta_1 \leq \tau', \tau' \leq \beta_2\} \subseteq \Delta(\mathcal{S})$ , and therefore  $[\beta_1 \leq \beta_2] \in \Delta(\mathcal{S})$  via  $(trans_\alpha)$ . By induction,  $[\beta_1 \leq \beta_2] \in \Psi\Theta(\mathcal{S})$ , and the following derivation then shows that  $[\tau_1 \leq \tau_2] \in \Psi\Theta(\mathcal{S})$ .

$$\frac{\frac{\alpha_1 \leq \mathbf{rng}(\beta_1) \quad \beta_1 \leq \beta_2}{\alpha_1 \leq \mathbf{rng}(\beta_2)} \quad (s_2) \quad \mathbf{rng}(\beta_2) \leq \alpha_2}{\alpha_1 \leq \alpha_2} \quad (s_4)$$

- \* Suppose  $\tau_1 \leq \mathbf{rng}(\tau')$  is inferred via  $(compose_1)$  and  $\mathbf{rng}(\tau') \leq \tau_2$  is inferred via  $(compat)$ . Then  $\{\beta_1 \leq \tau', \tau' \leq \tau'_2\} \subseteq \Delta(\mathcal{S})$ , and therefore  $[\beta_1 \leq \tau'_2] \in \Delta(\mathcal{S})$  via  $(trans_\alpha)$ . By induction,  $[\beta_1 \leq \tau'_2] \in \Psi\Theta(\mathcal{S})$ , and the following derivation shows that  $[\tau_1 \leq \tau_2] \in \Psi\Theta(\mathcal{S})$ .

$$\frac{\alpha_1 \leq \mathbf{rng}(\beta_1) \quad \beta_1 \leq \tau'_2}{\alpha_1 \leq \mathbf{rng}(\tau'_2)} \quad (compose_1)$$

- \* Suppose  $\tau_1 \leq \mathbf{rng}(\tau')$  is inferred via  $(compat)$  and  $\mathbf{rng}(\tau') \leq \tau_2$  is inferred via  $(compose_3)$ . This case holds by similar reasoning to the previous case.
- \* Suppose  $\tau_1 \leq \mathbf{rng}(\tau')$  is inferred via  $(compat)$  and  $\mathbf{rng}(\tau') \leq \tau_2$  is inferred via  $(compat)$ . Then  $\{\tau'_1 \leq \tau', \tau' \leq \tau'_2\} \subseteq \Delta(\mathcal{S})$ , and therefore  $[\tau'_1 \leq \tau'_2] \in \Delta(\mathcal{S})$  via  $(trans_\alpha)$ . By induction,  $[\tau'_1 \leq \tau'_2] \in \Psi\Theta(\mathcal{S})$ , and therefore a  $(compat)$ -inference shows that  $[\tau_1 \leq \tau_2] \in \Psi\Theta(\mathcal{S})$ .

There are no other possibilities for the derivations of  $\tau_1 \leq \mathbf{rng}(\tau')$  and  $\mathbf{rng}(\tau') \leq \tau_2$ .

- Suppose  $\tau = \mathbf{dom}(\tau')$ . This case holds by similar reasoning to the previous case where  $\tau = \mathbf{rng}(\tau')$ .

There are no other possibilities for  $\tau$ .

There are no other possibilities for the derivation of  $\mathbf{C} \in \Delta(\mathcal{S})$ , and hence  $\Delta(\mathcal{S}) \subseteq \Psi\Theta(\mathcal{S})$ . ■

**Lemma 4.5**  $\Psi\Theta(\mathcal{S}) \mid_E \cong \Pi\Theta(\mathcal{S}) \mid_E$ .

**Proof:** Since the rule  $(compat)$  does not create any  $\Pi$  or  $\Theta$  opportunities,  $\Psi\Theta(\mathcal{S}) = compat(\Pi\Theta(\mathcal{S}))$ , and hence we just need to show that:

$$compat(\Pi\Theta(\mathcal{S})) \mid_E \cong \Pi\Theta(\mathcal{S}) \mid_E$$

Now:

$$\begin{aligned} compat(\Pi\Theta(\mathcal{S})) &\supseteq \Pi\Theta(\mathcal{S}) \\ \therefore compat(\Pi\Theta(\mathcal{S})) \mid_E &\supseteq \Pi\Theta(\mathcal{S}) \mid_E \\ \therefore compat(\Pi\Theta(\mathcal{S})) \mid_E &\models \Pi\Theta(\mathcal{S}) \mid_E \end{aligned}$$

To prove the converse direction, let  $\rho \models \Pi\Theta(\mathcal{S}) \mid_E$ . If  $\rho \not\models compat(\Pi\Theta(\mathcal{S})) \mid_E$ , then let  $\mathbf{C}$  be the constraint in  $compat(\Pi\Theta(\mathcal{S})) \mid_E$  with the smallest derivation such that  $\rho \not\models \mathbf{C}$ . Then the last step in the derivation of  $\mathbf{C}$  must be via  $(compat)$ . Let  $\mathbf{C}'$  be the antecedent of this rule in

$compat(\Pi\Theta(\mathcal{S}))$ . Then  $SetVar(\mathbf{C}') = SetVar(\mathbf{C}) \subseteq E$ , and hence  $\mathbf{C}' \in compat(\Pi\Theta(\mathcal{S})) \mid_E$  with a smaller derivation. Therefore  $\rho \models \mathbf{C}'$ , and hence since  $(compat)$  is sound,  $\rho \models \mathbf{C}$ . Thus  $\rho \not\models compat(\Pi\Theta(\mathcal{S})) \mid_E$ , as required. ■

**Theorem 4.7 [Soundness and Completeness of  $\vdash_{\Psi\Theta}^E$  and  $=_{\Psi\Theta}^E$ ]**

1.  $\mathcal{S}_1 \vdash_{\Psi\Theta}^E \mathcal{S}_2$  if and only if  $\mathcal{S}_1 \models_E \mathcal{S}_2$ .
2.  $\mathcal{S}_1 =_{\Psi\Theta}^E \mathcal{S}_2$  if and only if  $\mathcal{S}_1 \cong_E \mathcal{S}_2$ .

1. Suppose  $\mathcal{S}_1 \models_E \mathcal{S}_2$ . Then

$$\begin{aligned}
& \mathcal{S}_1 \models_E \Delta(\mathcal{S}_2) && \text{by the soundness of } \Delta \\
\therefore & \mathcal{S}_1 \models_E \Delta(\mathcal{S}_2) \mid_E \\
\therefore & \forall \mathcal{C} \in \Delta(\mathcal{S}_2) \mid_E. \mathcal{S}_1 \models \mathcal{C} \\
\therefore & \forall \mathcal{C} \in \Delta(\mathcal{S}_2) \mid_E. \mathcal{S}_1 \vdash_{\Delta} \mathcal{C} && \text{by Lemma 4.2} \\
\therefore & \forall \mathcal{C} \in \Pi\Theta(\mathcal{S}_2) \mid_E. \mathcal{S}_1 \vdash_{\Delta} \mathcal{C} && \text{by Lemma 4.4} \\
\therefore & \forall \mathcal{C} \in \Pi\Theta(\mathcal{S}_2) \mid_E. \mathcal{S}_1 \vdash_{\Psi\Theta} \mathcal{C} && \text{by Lemma 4.4} \\
\therefore & \forall \mathcal{C} \in \Pi\Theta(\mathcal{S}_2) \mid_E. \mathcal{C} \in \Psi\Theta(\mathcal{S}_1) \mid_E \\
\therefore & \Psi\Theta(\mathcal{S}_1) \mid_E \supseteq \Pi\Theta(\mathcal{S}_2) \mid_E \\
\therefore & \mathcal{S}_1 \vdash_{\Psi\Theta}^E \mathcal{S}_2
\end{aligned}$$

Conversely, suppose  $\mathcal{S}_1 \vdash_{\Psi\Theta}^E \mathcal{S}_2$ . Then

$$\begin{aligned}
& \Psi\Theta(\mathcal{S}_1) \mid_E \supseteq \Pi\Theta(\mathcal{S}_2) \mid_E \\
\therefore & Soln(\Psi\Theta(\mathcal{S}_1) \mid_E) \subseteq Soln(\Pi\Theta(\mathcal{S}_2) \mid_E) \\
\therefore & Soln(\Psi\Theta(\mathcal{S}_1) \mid_E) \mid_E \subseteq Soln(\Pi\Theta(\mathcal{S}_2) \mid_E) \mid_E \\
\therefore & Soln(\mathcal{S}_1) \mid_E \subseteq Soln(\mathcal{S}_2) \mid_E \\
& \text{by Lemmas 4.3, 4.4 and 4.5, since } Soln(\mathcal{S}_i) \mid_E = Soln(\Pi\Theta(\mathcal{S}_i) \mid_E) \mid_E \\
\therefore & \mathcal{S}_1 \models_E \mathcal{S}_2
\end{aligned}$$

2. Follows from part 1.

## A.5 Proofs for Deciding Observable Equivalence

We repeat definition 4.8 here, to avoid having to refer back to the original definition earlier in the text.

**Definition 4.8** (Regular Grammar  $G_r(\mathcal{S}, E)$ ) Let  $\mathcal{S}$  be a simple constraint system and  $E$  a collection of set variables. The regular grammar  $G_r(\mathcal{S}, E)$  consists of the non-terminals  $\{\alpha_L, \alpha_U \mid \alpha \in SetVar(\mathcal{S})\}$  and the following productions:

$$\begin{array}{ll}
\alpha_U \mapsto \alpha, \alpha_L \mapsto \alpha & \forall \alpha \in E \\
\alpha_U \mapsto \beta_U, \beta_L \mapsto \alpha_L & \forall [\alpha \leq \beta] \in \mathcal{S} \\
\alpha_U \mapsto \text{dom}(\beta_L) & \forall [\alpha \leq \text{dom}(\beta)] \in \mathcal{S} \\
\alpha_U \mapsto \text{rng}(\beta_U) & \forall [\alpha \leq \text{rng}(\beta)] \in \mathcal{S} \\
\beta_L \mapsto \text{dom}(\alpha_U) & \forall [\text{dom}(\alpha) \leq \beta] \in \mathcal{S} \\
\beta_L \mapsto \text{rng}(\alpha_L) & \forall [\text{rng}(\alpha) \leq \beta] \in \mathcal{S}
\end{array}$$

■

**Lemma 4.9** *Let  $G = G_r(\mathcal{S}, E)$ . Then:*

$$\begin{aligned}\mathcal{L}_G(\alpha_L) &= \{\kappa \mid [\kappa \leq \alpha] \in \Pi(\mathcal{S}) \text{ and } \text{SetVar}(\kappa) \subseteq E\} \\ \mathcal{L}_G(\alpha_U) &= \{\kappa \mid [\alpha \leq \kappa] \in \Pi(\mathcal{S}) \text{ and } \text{SetVar}(\kappa) \subseteq E\}\end{aligned}$$

**Proof:** We prove the left-to-right inclusion by induction on the derivation of the constraint  $\mathbf{C} \in \Psi(\mathcal{S})$ , and by case analysis on the last step in that derivation.

- Suppose  $\mathbf{C} \in \Psi(\mathcal{S})$  because  $\mathbf{C} \in \mathcal{S}$ . We proceed by case analysis on  $\mathbf{C}$ .
  - Suppose  $\mathbf{C} = [\alpha \leq \beta]$ . Then  $\alpha \in E$ , so  $\beta_L \mapsto \alpha_L$  and  $\alpha_L \mapsto \alpha$  are productions in the grammar. Hence  $\alpha \in \mathcal{L}_G(\beta_L)$ . Similarly,  $\beta \in \mathcal{L}_G(\alpha_U)$ .

The remaining cases for  $\mathbf{C}$  follow by similar reasoning.

- Suppose  $\mathbf{C} = [\alpha \leq \mathbf{rng}(\kappa)]$  is inferred via  $(compose_1)$  from the antecedents  $[\alpha \leq \mathbf{rng}(\beta)]$  and  $[\beta \leq \kappa]$ . Then  $\alpha_U \mapsto \mathbf{rng}(\beta_U)$ , and by induction  $\beta_U \mapsto^* \kappa$ . Hence  $\alpha_U \mapsto \mathbf{rng}(\kappa)$ , as required.

The remaining cases follow by similar reasoning.

We prove the right-to-left inclusion by induction on the derivation  $\alpha_L \mapsto^* \kappa$  or  $\alpha_U \mapsto^* \kappa$ , and by case analysis on the last step in the derivation. The reasoning for each case is straightforward. ■

We repeat definition 4.10 here, to avoid having to refer back to the original definition earlier in the text.

**Definition 4.10** (Regular Tree Grammar  $G_t(\mathcal{S}, E)$ ) The regular tree grammar  $G_t(\mathcal{S}, E)$  extends the grammar  $G_r(\mathcal{S}, E)$  with the root non-terminal  $R$  and the additional productions:

$$\begin{aligned}R &\mapsto [\alpha_L \leq \alpha_U] & \forall \alpha \in \text{SetVar}(\mathcal{S}) \\ R &\mapsto [c \leq \alpha_U] & \forall [c \leq \alpha] \in \mathcal{S}\end{aligned}$$

where  $[\cdot \leq \cdot]$  is viewed as a binary constructor. ■

**Lemma 4.11** *Let  $G = G_t(\mathcal{S}, E)$ . Then  $\Pi(\mathcal{S}) \upharpoonright_E = \mathcal{L}_G(R)$ .*

**Proof:** We prove the left-to-right inclusion by case analysis on  $\mathbf{C} \in \Pi(\mathcal{S}) \upharpoonright_E$ .

- Suppose  $\mathbf{C} = [\alpha \leq \kappa]$ . Then by Lemma 4.9,  $\alpha_U \mapsto_G^* \kappa$ . Since  $\text{SetVar}(\mathbf{C}) \subseteq E$ ,  $\alpha \in E$ , and hence  $\alpha_L \mapsto_G \alpha$ . Thus  $R \mapsto_G [\alpha_L \leq \alpha_U] \mapsto_G^* [\alpha \leq \kappa]$ , and hence  $[\alpha \leq \kappa] \in \mathcal{L}_G(R)$ .
- The case where  $\mathbf{C} = [\kappa \leq \alpha]$  follows by similar reasoning.

- Suppose  $\mathbf{C} = [c \leq \kappa]$ . If  $\mathbf{C} \in \mathcal{S}$ , then  $\kappa = \alpha$ ,  $\alpha \in E$ , and

$$R \mapsto_G [c \leq \alpha_U] \mapsto_G [c \leq \alpha]$$

as required.

If  $\mathbf{C} \notin \mathcal{S}$ , then an examination of the inference rules in  $\Pi$  shows that  $\mathbf{C}$  can only be inferred via  $(trans_\alpha)$ , based on the antecedents  $[c \leq \alpha]$  and  $[\alpha \leq \kappa]$ . By Lemma 4.9,  $\alpha_L \mapsto_G^* c$  and  $\alpha_U \mapsto_G^* \kappa$ . Hence  $R \mapsto [c \leq \kappa]$ , and hence  $[c \leq \kappa] \in \mathcal{L}_G(R)$ , as required.

- Otherwise  $\mathbf{C} = [\kappa_1 \leq \kappa_2]$ , where  $\kappa_1, \kappa_2 \notin SetVar$ . An examination of the inference rules in  $\Pi$  shows that  $\mathbf{C}$  can only be inferred via  $(trans_\alpha)$ , based on the antecedents  $[\kappa_1 \leq \alpha]$  and  $[\alpha \leq \kappa_2]$ . By Lemma 4.9,  $\alpha_L \mapsto_G^* \kappa_1$  and  $\alpha_U \mapsto_G^* \kappa_2$ . Hence  $R \mapsto [\kappa_1 \leq \kappa_2]$ , and hence  $[\kappa_1 \leq \kappa_2] \in \mathcal{L}_G(R)$ , as required.

We prove the right-to-left inclusion by case analysis on  $\mathbf{C} \in \mathcal{L}_G(R)$ .

- Suppose  $\mathbf{C} = [\kappa_1 \leq \kappa_2]$ . Then for some  $\alpha$ ,  $\alpha_L \mapsto_G^* \kappa_1$  and  $\alpha_U \mapsto_G^* \kappa_2$ . By Lemma 4.9,  $\{\kappa_1 \leq \alpha, \alpha \leq \kappa_2\} \subseteq \Psi(\mathcal{S})$  and  $SetVar(\kappa_i) \subseteq E$ . By Lemma 4.12,  $\{\kappa_1 \leq \alpha, \alpha \leq \kappa_2\} \subseteq \Pi(\mathcal{S})$ . Hence  $[\kappa_1 \leq \kappa_2] \in \Pi(\mathcal{S}) \mid_E$ , as required.
- Otherwise  $\mathbf{C} = [c \leq \kappa]$ . Then for some  $\alpha$ ,  $[c \leq \alpha] \in \mathcal{S}$  and  $\alpha_U \mapsto_G^* \kappa$ . By Lemma 4.9,  $\{\alpha \leq \kappa\} \subseteq \Psi(\mathcal{S})$  and  $SetVar(\kappa) \subseteq E$ . By Lemma 4.12,  $\{\alpha \leq \kappa\} \subseteq \Pi(\mathcal{S})$ . Hence  $[c \leq \kappa] \in \Pi(\mathcal{S}) \mid_E$ , as required.

■

**Lemma 4.12 (Staging)** *For any simple constraint system  $\mathcal{S}$ :*

$$\Psi\Theta(\mathcal{S}) = \Psi(\Theta(\mathcal{S})) = compat(\Pi(\Theta(\mathcal{S})))$$

**Proof:** The equality  $\Psi\Theta(\mathcal{S}) = compat(\Pi(\Theta(\mathcal{S})))$  holds since  $(compat)$  does not create any  $\Pi$  or  $\Theta$  opportunities.

The inclusion  $\Psi\Theta(\mathcal{S}) \supseteq \Psi(\Theta(\mathcal{S}))$  obviously holds. To prove the inclusion  $\Psi\Theta(\mathcal{S}) \subseteq \Psi(\Theta(\mathcal{S}))$  holds, we suppose  $\mathcal{S} \vdash_{\Psi\Theta} \mathbf{C}$ , and prove  $\Theta(\mathcal{S}) \vdash_{\Psi} \mathbf{C}$  by induction on the derivation  $\mathcal{S} \vdash_{\Psi\Theta} \mathbf{C}$ , and by case analysis on the last step in this derivation.

- Suppose  $\mathcal{S} \vdash_{\Psi\Theta} \mathbf{C}$  via some rule in  $\Psi$ . By induction, the antecedents of this rule are in  $\Psi(\Theta(\mathcal{S}))$ , and hence  $\mathbf{C}$  is also in  $\Psi(\Theta(\mathcal{S}))$ .
- Suppose  $\mathcal{S} \vdash_{\Psi\Theta} \mathbf{C}$  via one of the rules  $(s_1)$ ,  $(s_2)$  or  $(s_3)$ . These rules are subsumed by  $(trans_\alpha)$ ,  $(compose_1)$  and  $(compose_4)$ , and hence this case is subsumed by the previous case.
- Suppose  $\mathcal{S} \vdash_{\Psi\Theta} \mathbf{C}$  via  $(s_4)$ , based on the antecedents  $\{\alpha \leq \mathbf{rng}(\beta), \mathbf{rng}(\beta) \leq \gamma\}$ . By induction, these antecedents are in  $\Psi(\Theta(\mathcal{S}))$ . An examination of  $\Psi$  shows that  $\Psi$  can only infer  $[\alpha \leq \mathbf{rng}(\beta)]$  if there exists  $\alpha', \beta'$  such that  $\Theta(\mathcal{S})$  contains the constraints:

$$\alpha \leq^* \alpha' \quad \alpha' \leq \mathbf{rng}(\beta') \quad \beta' \leq^* \beta$$



### The Entailment Algorithm

In the following,  $\mathcal{P}_{\text{fin}}$  denotes the finite power-set constructor.

Let:

$$\begin{aligned} G_1 &= G_r(\mathcal{S}_1, E) & L_i &= \{\alpha_L \mid \alpha \in \text{Vars}(\mathcal{S}_i)\} \\ G_2 &= G_t(\mathcal{S}_2, E) & U_i &= \{\alpha_U \mid \alpha \in \text{Vars}(\mathcal{S}_i)\} \end{aligned}$$

Assume  $G_1$  and  $G_2$  are pre-processed to remove  $\epsilon$ -transitions. For  $C \in \mathcal{P}_{\text{fin}}(L_2 \times U_2)$ , define:

$$\mathcal{L}(C) = \{[\tau_L \leq \tau_U] \mid \langle \alpha_L, \beta_U \rangle \in C, \alpha_L \mapsto_{G_2} \tau_L, \beta_U \mapsto_{G_2} \tau_U\}$$

The relation  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\cdot, \cdot, \cdot, \cdot]$  is defined as the largest relation on  $L_1 \times U_1 \times \mathcal{P}_{\text{fin}}(L_2 \times U_2) \times \mathcal{P}_{\text{fin}}(L_2 \times U_2)$  such that if:

$$\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D] \quad \alpha_L \mapsto_{G_1} X \quad \beta_U \mapsto_{G_1} Y$$

then one of the following cases hold:

1.  $\mathcal{L}([X \leq Y]) \subseteq \mathcal{L}(C \cup D)$ .
2.  $X = \text{rng}(\alpha'_L)$ ,  $Y = \text{rng}(\beta'_U)$  and  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha'_L, \beta'_U, C, D']$ , where

$$D' = \{\langle \gamma'_L, \delta'_U \rangle \mid \langle \gamma_L, \delta_U \rangle \in C \cup D, \gamma_L \mapsto_{G_2} \text{rng}(\gamma'_L), \delta_U \mapsto_{G_2} \text{rng}(\delta'_U)\}$$

3.  $X = \text{dom}(\alpha'_L)$ ,  $Y = \text{dom}(\beta'_U)$  and  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\beta'_L, \alpha'_U, C, D']$ , where

$$D' = \{\langle \delta'_L, \gamma'_U \rangle \mid \langle \gamma_L, \delta_U \rangle \in C \cup D, \gamma_L \mapsto_{G_2} \text{dom}(\gamma'_U), \delta_U \mapsto_{G_2} \text{dom}(\delta'_L)\}$$

The *computable entailment relation*  $\mathcal{S}_2 \vdash_{\text{alg}}^E \mathcal{S}_1$  holds if and only if  $\forall \alpha \in \text{Vars}(\mathcal{S}_1)$ :

$$\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \alpha_U, \{\langle \gamma_L, \gamma_U \rangle \mid \gamma \in \text{Vars}(\mathcal{S}_2)\}, \emptyset]$$

Figure 13: The computable entailment relation  $\vdash_{\text{alg}}^E$

Similarly,  $\Psi$  can only infer  $[\text{rng}(\beta) \leq \gamma]$  if there exists  $\beta'', \gamma'$  such that  $\Theta(\mathcal{S})$  contains the constraints:

$$\beta \leq^* \beta'' \quad \text{rng}(\beta'') \leq \gamma' \quad \gamma' \leq^* \gamma$$

Hence:

$$\begin{aligned} \mathcal{S} \vdash_{\Theta} \alpha' &\leq \text{rng}(\beta'') && \text{via multiple applications of } (s_2) \\ \mathcal{S} \vdash_{\Theta} \alpha' &\leq \gamma' && \text{via } (s_4) \\ \Theta(\mathcal{S}) \vdash_{\Psi} \alpha &\leq \gamma && \text{via multiple applications of } (trans_{\alpha}) \end{aligned}$$

- The case for  $(s_5)$  holds by similar reasoning.

■

### A.6 Correctness of the Entailment Algorithm

**Theorem 4.13** (*Correctness of the Entailment Algorithm*).  $\mathcal{S}_2 \vdash_{\Psi}^E \mathcal{S}_1$  if and only if  $\mathcal{S}_2 \vdash_{\text{alg}}^E \mathcal{S}_1$ .

**Proof:** The definitions of the computable entailment relation and the relation  $\mathcal{R}$  are shown in figure 13. We prove this theorem based on the following invariant concerning the relation  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\cdot, \cdot, \cdot, \cdot]$ :

$$\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D] \iff \mathcal{L}([\alpha_L \leq \beta_U]) \subseteq \text{compat}(\mathcal{L}(C)) \cup \mathcal{L}(D)$$

Assume this relation holds, and that  $\mathcal{S}_2 \vdash_{\Psi}^E \mathcal{S}_1$ . Then  $\Pi(\mathcal{S}_1) \upharpoonright_E \subseteq \text{compat}(\Pi(\mathcal{S}_2)) \upharpoonright_E$ . By lemma 4.11,  $\Pi(\mathcal{S}_i) \upharpoonright_E = \mathcal{L}_{G_i}(R)$ , and hence:

$$\mathcal{L}_{G_1}(R) \subseteq \text{compat}(\mathcal{L}_{G_2}(R))$$

Thus, for all  $R \mapsto_{G_1} [\alpha_L \leq \alpha_U]$ :

$$\begin{aligned} \mathcal{L}_{G_1}([\alpha_L \leq \alpha_U]) &\subseteq \text{compat}(\mathcal{L}_{G_2}(R)) \\ \therefore \mathcal{L}_{G_1}([\alpha_L \leq \alpha_U]) &\subseteq \text{compat}(\mathcal{L}_{G_2}(\{\langle \gamma_L, \gamma_U \rangle \mid \gamma \in \text{SetVar}(\mathcal{S}_2)\})) \end{aligned}$$

Hence:

$$\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \alpha_U, \{\langle \gamma_L, \gamma_U \rangle \mid \gamma \in \text{SetVar}(\mathcal{S}_2)\}, \emptyset]$$

Also, from  $\mathcal{L}_{G_1}(R) \subseteq \text{compat}(\mathcal{L}_{G_2}(R))$ , we have that for all  $R \mapsto_{G_1} [c \leq \alpha_U]$ :

$$\begin{aligned} \mathcal{L}_{G_1}([c \leq \alpha_U]) &\subseteq \text{compat}(\mathcal{L}_{G_2}(R)) \\ \therefore \mathcal{L}_{G_1}([c \leq \alpha_U]) &\subseteq \mathcal{L}_{G_2}(R) \\ \therefore \mathcal{L}_{G_1}(\alpha_U) &\subseteq \mathcal{L}_{G_2}(\{\gamma_U \mid R \mapsto_{G_2} [c \leq \gamma_U]\}) \end{aligned}$$

Hence  $\mathcal{S}_2 \vdash_{\text{alg}}^E \mathcal{S}_1$  holds. The proof of the converse implication that  $\mathcal{S}_2 \vdash_{\text{alg}}^E \mathcal{S}_1$  implies  $\mathcal{S}_2 \vdash_{\Psi}^E \mathcal{S}_1$  proceeds by a similar argument.

It remains to show that the invariant concerning  $\mathcal{R}$  holds. To prove the left-to-right direction, suppose  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D]$  and:

$$\begin{aligned} \alpha_L &\mapsto_{G_1} X \mapsto_{G_1}^* \tau_L \\ \beta_U &\mapsto_{G_1} Y \mapsto_{G_1}^* \tau_U \end{aligned}$$

We prove by induction on  $\tau_L$  that

$$\mathcal{L}([\tau_L \leq \tau_U]) \subseteq \text{compat}(\mathcal{L}(C)) \cup \mathcal{L}(D)$$

One of three cases in the definition of  $\mathcal{R}$  must hold.

1.  $\mathcal{L}([\alpha_L \leq \beta_U]) \subseteq \mathcal{L}(C \cup D)$ . This case is trivial.
2. In this case:

$$\begin{aligned} X &= \text{rng}(\alpha'_L) & \alpha'_L &\mapsto_{G_1}^* \tau'_L & \tau_L &= \text{rng}(\tau'_L) \\ Y &= \text{rng}(\beta'_U) & \beta'_U &\mapsto_{G_1}^* \tau'_U & \tau_U &= \text{rng}(\tau'_U) \end{aligned}$$

and  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha'_L, \beta'_U, C, D']$ , where

$$D' = \{\langle \gamma'_L, \delta'_U \rangle \mid \langle \gamma_L, \delta_U \rangle \in C \cup D, \gamma_L \mapsto_{G_2} \text{rng}(\gamma'_L), \delta_U \mapsto_{G_2} \text{rng}(\delta'_U)\}$$

By induction,  $[\tau'_L \leq \tau'_U] \in \text{compat}(\mathcal{L}(C)) \cup \mathcal{L}(D')$ .

- If  $[\tau'_L \leq \tau'_U] \in \mathcal{L}(D')$  then there exists  $\langle \gamma'_L, \delta'_U \rangle \in D$  such that  $\gamma'_L \mapsto_{G_2}^* \tau'_L$  and  $\delta'_U \mapsto_{G_2}^* \tau'_U$ . By the definition of  $D'$ , there exists  $\langle \gamma_L, \delta_U \rangle \in C \cup D$  such that  $\gamma_L \mapsto_{G_2}^* \tau_L$  and  $\delta_U \mapsto_{G_2}^* \tau_U$ . Therefore  $[\tau_L \leq \tau_U] \in \mathcal{L}(C \cup D)$ , as required.
  - If  $[\tau'_L \leq \tau'_U] \in \text{compat}(\mathcal{L}(C))$  then  $[\tau_L \leq \tau_U] \in \text{compat}(\mathcal{L}(C))$ , as required.
3. The proof for the third case of the definition of  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\cdot, \cdot, \cdot, \cdot]$  is similar to that for the second case.

To prove the right-to-left direction, suppose:

$$\mathcal{L}([\alpha_L \leq \beta_U]) \subseteq \text{compat}(\mathcal{L}(C)) \cup \mathcal{L}(D)$$

and that the relation  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D]$  does not hold. Hence there exists  $X, Y$  such that  $\alpha_L \mapsto_{G_1} X$  and  $\beta_U \mapsto_{G_1} Y$  and none of the three conditions in figure 13 hold. Furthermore, since  $\mathcal{R}$  is the largest relation satisfying the conditions in figure 13, there exists a finite proof that none of the three conditions hold.

Of all possible such *counter-examples*  $\langle \alpha_L, \beta_U, X, Y, C, D \rangle$ , we pick the one with the smallest proof that the relation  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D]$  does not hold, and proceed by case analysis on the last step in this proof.

- Suppose  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D]$  does not hold because of condition one. Then  $\mathcal{L}([X \leq Y]) \not\subseteq \mathcal{L}(C \cup D)$ , which contradicts the assumptions above.
- Suppose  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D]$  does not hold because of condition 2. Then  $X = \text{rng}(\alpha'_L)$  and  $Y = \text{rng}(\beta'_U)$ . Consider any pair of set expressions  $\tau_L$  and  $\tau_U$  such that  $\alpha'_L \mapsto_{G_1}^* \tau_L$  and  $\beta'_U \mapsto_{G_1}^* \tau_U$ . We consider the two possibilities for  $[\text{rng}(\tau_L) \leq \text{rng}(\tau_U)] \in \text{compat}(\mathcal{L}(C)) \cup \mathcal{L}(D)$  separately.

- If  $[\text{rng}(\tau_L) \leq \text{rng}(\tau_U)] \in \mathcal{L}(C) \cup \mathcal{L}(D)$ , then there exists  $\langle \gamma_L, \delta_U \rangle \in C \cup D$  such that:

$$\begin{aligned} \gamma_L \mapsto_{G_2} \text{rng}(\gamma'_L) &\mapsto_{G_2}^* \text{rng}(\tau_L) \\ \delta_U \mapsto_{G_2} \text{rng}(\delta'_U) &\mapsto_{G_2}^* \text{rng}(\tau_U) \end{aligned}$$

Hence  $[\tau_L \leq \tau_U] \in \mathcal{L}(D')$ , where:

$$D' = \{ \langle \gamma'_L, \delta'_U \rangle \mid \langle \gamma_L, \delta_U \rangle \in C \cup D, \gamma_L \mapsto_{G_2} \text{rng}(\gamma'_L), \delta_U \mapsto_{G_2} \text{rng}(\delta'_U) \}$$

- Otherwise  $[\text{rng}(\tau_L) \leq \text{rng}(\tau_U)] \in \text{compat}(\mathcal{L}(C)) \setminus \mathcal{L}(C)$ , and hence  $[\tau_L \leq \tau_U] \in \text{compat}(\mathcal{L}(C))$ .

Hence

$$\mathcal{L}([\alpha'_L \leq \beta'_U]) \subseteq \text{compat}(\mathcal{L}(C)) \cup \mathcal{L}(D')$$

The proof that  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D]$  does not hold cannot rely on a smaller proof that  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha'_L, \beta'_U, C, D]$  does not hold, since that would yield a counter-example with a smaller proof.

- The case where  $\mathcal{R}_{\mathcal{S}_1, \mathcal{S}_2}[\alpha_L, \beta_U, C, D]$  does not hold because of condition 3 is also impossible via similar reasoning.

Thus the invariant on  $\mathcal{R}$  is true, and thus the lemma holds. ■

## A.7 Correctness of the Hopcroft Algorithm

**Theorem 4.14** (*Correctness of the Hopcroft Algorithm*). *Let  $\mathcal{S}$  be a simple constraint system with external variables  $E$ ; let  $\sim$  be an equivalence relation on the set variables in a constraint system  $\mathcal{S}$  satisfying conditions (a) to (e) from figure 9; let the substitution  $f$  map each set variable to a representation element of its equivalence class; and let  $\mathcal{S}' = f(\mathcal{S})$ , i.e.,  $\mathcal{S}'$  denotes the constraint system  $\mathcal{S}$  with set variables merged according to their equivalence class. Then  $\mathcal{S} \cong_E \mathcal{S}'$ .*

**Proof:** Let  $\rho$  be a solution of  $\mathcal{S}$ . Define  $\rho'$  by:

$$\rho'(\alpha) = \bigsqcup_{\alpha' \sim \alpha} \rho(\alpha')$$

Obviously  $\rho, \rho'$  agree on  $E$  by condition (a) on  $\sim$ . We claim that  $\rho' \models \mathcal{C}$  for all  $\mathcal{C} \in \mathcal{S}$  by case analysis on  $\mathcal{C}$ .

- Suppose  $\mathcal{C} = [\alpha \leq \text{rng}(\beta)]$ . Then for all  $\alpha'$  such that  $\alpha \sim \alpha'$  there exists  $\beta'$  such that  $\beta \sim \beta'$  and:

$$\rho(\alpha') \sqsubseteq \rho(\text{rng}(\beta'))$$

Hence for all  $\alpha \sim \alpha'$ :

$$\rho(\alpha') \sqsubseteq \bigsqcup_{\beta' \sim \beta} \rho(\text{rng}(\beta'))$$

and therefore:

$$\bigsqcup_{\alpha' \sim \alpha} \rho(\alpha') \sqsubseteq \bigsqcup_{\beta' \sim \beta} \rho(\text{rng}(\beta'))$$

Hence:

$$\begin{aligned} \rho'(\alpha) &= \bigsqcup_{\alpha' \sim \alpha} \rho(\alpha') \\ &\sqsubseteq \bigsqcup_{\beta' \sim \beta} \rho(\text{rng}(\beta')) \\ &= \text{rng} \left( \bigsqcup_{\beta' \sim \beta} \rho(\beta') \right) \\ &= \text{rng}(\rho'(\beta)) \\ &= \rho'(\text{rng}(\beta)) \end{aligned}$$

and thus  $\rho' \models \mathcal{C}$ , as required.

- The cases where  $\mathcal{C} = [\alpha \leq \beta]$  and  $\mathcal{C} = [\text{rng}(\alpha) \leq \beta]$  follow by similar reasoning.
- Suppose  $\mathcal{C} = [\alpha \leq \text{dom}(\beta)]$ . Then  $\forall \alpha \sim \alpha' \forall \beta \sim \beta'$  such that:

$$\rho(\alpha') \sqsubseteq \rho(\text{dom}(\beta'))$$

Hence  $\forall \alpha \sim \alpha'$ :

$$\rho(\alpha') \sqsubseteq \sqcap_{\beta' \sim \beta} \rho(\text{dom}(\beta'))$$

and therefore:

$$\bigsqcup_{\alpha' \sim \alpha} \rho(\alpha') \sqsubseteq \sqcap_{\beta' \sim \beta} \rho(\text{dom}(\beta'))$$

Hence:

$$\begin{aligned} \rho'(\alpha) &= \bigsqcup_{\alpha' \sim \alpha} \rho(\alpha') \\ &\sqsubseteq \sqcap_{\beta' \sim \beta} \rho(\text{dom}(\beta')) \\ &= \text{dom} \left( \bigsqcup_{\beta' \sim \beta} \rho(\beta') \right) \\ &= \text{dom}(\rho'(\beta)) \\ &= \rho'(\text{dom}(\beta)) \end{aligned}$$

- Suppose  $\mathcal{C} = [\text{dom}(\beta) \leq \alpha]$ . Then:

$$\begin{aligned} \rho'(\text{dom}(\beta)) &= \text{dom}(\rho'(\beta)) \\ &= \text{dom} \left( \bigsqcup_{\beta' \sim \beta} \rho(\beta') \right) \\ &= \sqcap_{\beta' \sim \beta} \rho(\text{dom}(\beta')) \\ &\sqsubseteq \rho(\text{dom}(\beta)) \\ &\sqsubseteq \rho(\alpha) \\ &\sqsubseteq \bigsqcup_{\alpha' \sim \alpha} \rho(\alpha') \\ &= \rho'(\alpha) \end{aligned}$$

- Suppose  $\mathcal{C} = [c \leq \alpha]$ . Then

$$\begin{aligned} \rho'(\alpha) &= \bigsqcup_{\alpha' \sim \alpha} \rho(\alpha') \\ &\sqsupseteq \rho(\alpha) \\ &= c \end{aligned}$$

■

## B Notations

Symbol	Meaning	Section	Page
$M \in \Lambda$	Terms	2	4
$V \in \text{Value}$	Values	2	4
$x \in \text{Vars}$	Variables	2	4
$b \in \text{BasicConst}$	Basic constants	2	4

$t \in Tag$	Function tags	2	4
$l \in Label$	Labels	2	4
$\beta_v, \beta_{let}, unlabeled$	Reduction rules	3.1.2	4
$\mathcal{E}$	Evaluation contexts	3.1.2	4
$\longrightarrow, \longrightarrow^*$	Standard reduction relation	3.1.2	5
$eval$	Evaluator	3.1.2	5
$\tau \in SetExp$	Set expressions	3.2	5
$\alpha, \beta, \dots \in SetVar$	Set variables	3.2	5
$c \in Const$	Constants	3.2	5
$dom, rng$	Type expression constructors	3.2	5
$\mathcal{C}$	Constraints	3.2	5
$\mathcal{S}$	Constraint systems	3.2	5
$\mathcal{S} \mid_E$	Restriction of a constraint system	3.2	5
$\mathcal{D}, \sqsubseteq, \top, \perp, \sqcup, \sqcap$	Domain, ordering, elements and operations	3.3.1	6
$\mathcal{P}$	Power set constructor	3.2	5
$in$	Values described by constants	3.3.1	6
$const, dom, rng$	Extract components of element of $\mathcal{D}$	??	??
$\rho \in SetEnv$	Set environment	3.3.2	7
$Soln(\mathcal{S})$	Solution space	3.3.2	7
$\models, \models_E$	Satisfies, or entails	3.3.2	7
$Soln(\mathcal{S}) \mid_E$	Restriction	3.3.2	7
$\cong_E$	Observable equivalence	3.1	8
$\vdash$	Constraint derivation rules	3.4	8
$\Gamma$	Set variable context	3.4	8
$\sigma$	Constraint schema	3.4	9
$\sqsubseteq_s, \top_s, \perp_s, \sqcup_s, \sqcap_s$	Alternative ordering on domain	3.4	9
$sba$	Analysis function	3.10	12
$LeastSoln$	Least Solution	3.6	12
$\Theta = \{s_1, \dots, s_n\}$	Inference rules	3.6.1	13
$\vdash$	Deduction	3.6.1	13
$E$	External variables	4.1	14
$FV[rng(\Gamma)]$	Free variables in $\Gamma$	3.4	9
$\Psi, \Pi$	Inference rules on constraint systems	4.2	17
$G$	Grammar	4.8	20
$G_r$	Function producing regular grammar	4.8	20
$\alpha_L, \alpha_U$	Grammar non-terminals	4.8	20
$\mathcal{L}_G(X)$	Language for $X$ in $G$	4.8	20
$G_t$	Function producing RTG	4.10	21

$R$	Root non-terminal	4.10	21
$p, q$	Paths	A.4	43
$\pi$	Arity function	A.4	43
$\leq_i$	$\leq$ or $\geq$	A.4	43
<b>sel</b>	Either <b>dom</b> or <b>rng</b>	A.4	43
$\overline{C}$	Converts compound to simple constraints	A.4	43
$\cdot @ \cdot$	Injection function	A.4	43
$\leq^*$	Transitive closure of $\leq$	A.2	41

## References

- [1] AHO, A., J. HOPCROFT AND J. ULLMAN. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, Reading, Mass., 1974.
- [2] AIKEN, A., WIMMERS, E. L., AND LAKSHMAN, T. K. Soft typing with conditional types. In *Proceedings of the ACM Sigplan Conference on Principles of Programming Languages* (1994), pp. 163–173.
- [3] COUSOT, P., AND COUSOT, R. Formal language, grammar, and set-constraint-based program analysis by abstract interpretation. In *Proceedings of the 1995 Conference on Functional Programming and Computer Architecture* (1995), pp. 170–181.
- [4] DUESTERWALD, E., GUPTA, R., AND SOFFA, M. L. Reducing the cost of data flow analysis by congruence partitioning. In *International Conference on Compiler Construction* (April 1994).
- [5] EIFRIG, J., SMITH, S., AND TRIFONOV, V. Sound polymorphic type inference for objects. In *Conference on Object-Oriented Programming Systems, Languages, and Applications* (1995).
- [6] FÄHNDRICH, M., AND AIKEN, A. Making set-constraint based program analyses scale. Technical Report UCB/CSD-96-917, University of California at Berkeley, 1996.
- [7] FLANAGAN, C., AND FELLEISEN, M. Set-based analysis for full Scheme and its use in soft-typing. Technical Report TR95-254, Rice University, 1995.
- [8] FLANAGAN, C., FLATT, M., KRISHNAMURTHI, S., WEIRICH, S., AND FELLEISEN, M. Finding bugs in the web of program invariants. In *Proceedings of the ACM Conference on Programming Language Design and Implementation* (1996), pp. 23–32.
- [9] FLATT, M. *MzScheme Reference Manual*. Rice University.
- [10] HEINTZE, N. Set-based analysis of ML programs. In *Proceedings of the ACM Conference on Lisp and Functional Programming* (1994), pp. 306–317.
- [11] HINDLEY, R. J., AND SELDIN, J. P. *Introduction to Combinators and  $\lambda$ -Calculus*. Cambridge University Press, 1986.

- [12] HOPCROFT, J. E. An  $n \log n$  algorithm for minimizing the states of a finite automaton. *The Theory of Machines and Computations* (1971), 189–196.
- [13] JAGANNATHAN, S., AND WRIGHT, A. K. Effective flow analysis for avoiding run-time checks. In *Proc. 2nd International Static Analysis Symposium, LNCS 983* (September 1995), Springer-Verlag, pp. 207–224.
- [14] JONES, N., AND MUCHNICK, S. A flexible approach to interprocedural data flow analysis and programs with recursive data structures. In *Conference Record of the Ninth Annual ACM Symposium on Principles of Programming Languages* (January 1982), pp. 66–74.
- [15] PALSBERG, J. Closure analysis in constraint form. *Transactions on Programming Languages and Systems* 17, 1 (1995), 47–62.
- [16] PALSBERG, J., AND O'KEEFE, P. A type system equivalent to flow analysis. In *Proceedings of the ACM SIGPLAN '95 Conference on Principles of Programming Languages* (1995), pp. 367–378.
- [17] POTTIER, F. Simplifying subtyping constraints. In *Proceedings of the 1996 ACM SIGPLAN International Conference on Functional Programming* (1996), pp. 122–133.
- [18] REYNOLDS, J. Automatic computation of data set definitions. *Information Processing '68* (1969), 456–461.
- [19] TOFTE, M. Type inference for polymorphic references. *Information and Computation* 89, 1 (November 1990), 1–34.
- [20] TRIFONOV, V., AND SMITH, S. Subtyping constrained types. In *Third International Static Analysis Symposium (LNCS 1145)* (1996), pp. 349–365.
- [21] WRIGHT, A., AND FELLEISEN, M. A syntactic approach to type soundness. *Information and Computation* 115, 1 (1994), 38–94.
- [22] WRIGHT, A. K. Simple imperative polymorphism. *Lisp and Symbolic Computation* 8, 4 (Dec. 1995), 343–356.