

Type Inference with Non-structural Subtyping*

Jens Palsberg[†] Mitchell Wand^{‡§} Patrick O’Keefe[¶]

November 27, 2004

Abstract

We present an $O(n^3)$ time type inference algorithm for a type system with a largest type \top , a smallest type \perp , and the usual ordering between function types. The algorithm infers type annotations of least shape, and it works equally well for recursive types. For the problem of typability, our algorithm is simpler than the one of Kozen, Palsberg, and Schwartzbach for type inference *without* \perp . This may be surprising, especially because the system with \perp is strictly more powerful.

**Formal Aspects of Computing*, 9:49–67, 1997.

[†]Massachusetts Institute of Technology, NE43–340, 545 Technology Square, Cambridge, MA 02139, USA. E-mail: palsberg@theory.lcs.mit.edu.

[‡]Work supported by the National Science Foundation under grants CCR-9304144 and CCR-9404646.

[§]College of Computer Science, Northeastern University, 360 Huntington Avenue, 161CN, Boston, MA 02115, USA. E-mail: wand@ccs.neu.edu.

[¶]151 Coolidge Avenue #211, Watertown, MA 02172, USA. E-mail: pmo@world.std.com.

1 Introduction

This paper concerns types for the λ -calculus that can be generated from the grammar:

$$t ::= \perp \mid \top \mid t_1 \rightarrow t_2 .$$

Intuitively, \perp is a least type containing only the divergent computation; \top is a maximal or universal type containing all values; and $t_1 \rightarrow t_2$ is the usual function space.

Types are partially ordered by \leq which is the smallest binary relation on types such that

1. $\perp \leq t \leq \top$ for all types t ; and
2. $s \rightarrow t \leq s' \rightarrow t'$ if and only if $s' \leq s$ and $t \leq t'$.

Following [10], we denote this type system by **PTB** (**PTB** indicates *Partial Types with Bottom*). Thatte's system of partial types [8] did not include \perp ; the fragment of **PTB** without \perp will be denoted **PT** [10].

Subtyping in **PTB** and **PT** is *non-structural*, that is, two types can be related by \leq even in cases where they have different structures. For example $\top \rightarrow \top \leq \top$ and $\top \rightarrow \top \leq (\top \rightarrow \top) \rightarrow \top$. In contrast, type systems with atomic subtyping [4] have the property that if s is a subtype of t , then s and t have the same structure; only the leaves of s and t can be different. We are interested in non-structural subtyping because this is the form of subtyping that is often used in connection with records and objects. The type system **PTB** is a convenient setting for studying certain aspects of non-structural subtyping, especially the ordering of function types.

The type system **PTB** can be extended with recursive types, yielding Amadio and Cardelli's system [1], here denoted **PTB _{μ}** . The further extension of **PTB _{μ}** with the type **Int** of integers will be denoted **PTB _{μ} \cup {Int}**.

We have earlier shown that **PTB** is strictly more powerful than **PT** [10], and recursive types add further power.

It is known that type inference for **PT** is computable in $O(n^3)$ time [2]. It is well known that all λ -terms are typable in **PTB _{μ}** by the type $\mu\alpha.\alpha \rightarrow \alpha$. It is also known that type inference for **PTB _{μ} \cup {Int}** is computable in $O(n^3)$ time, by reduction to a flow analysis problem [6].

In this paper we show that type inference for **PTB** and **PTB _{μ}** is also computable in $O(n^3)$ time, by an algorithm similar to that in [2]. The algorithm

infers type annotations of least shape. A type is a tree (which can be infinite if we consider recursive types), so its shape is the set of paths from the root. In general, there can be several types of least shape. For example, the λ -term $\lambda x.\lambda y.(\lambda f.f(fx))(\lambda v.vy)$ can be typed such that

$$\begin{aligned} x & : \perp \\ y & : \perp \\ f & : \perp \rightarrow \perp \\ v & : \perp . \end{aligned}$$

These type annotations are of least shape. If we replace the type of y by \top , then we obtain another type annotation of least shape. The above λ -term and the two typings themselves are not particularly interesting, but it is the smallest example we have found for which there are more than one least-shape typing.

By contrast, in the type system PT , every typable λ -term has exactly one type annotation of least shape. The algorithm by Kozen, Palsberg, and Schwartzbach [2] infers this least-shape type annotation.

The type inference algorithm by Palsberg and O’Keefe [6] for $\text{PTB}_\mu \cup \{\text{Int}\}$ does in general not infer type annotations of least shape. It remains open if our result can be extended to $\text{PTB}_\mu \cup \{\text{Int}\}$.

For the problem of typability, our algorithm is simpler than the one of Kozen, Palsberg, and Schwartzbach for type inference *without* \perp . This may be surprising, especially because the system with \perp is strictly more powerful. In our view, the simpler algorithm motivates having \top and \perp together in a type system.

Our algorithm proceeds by first calculating a compact representation of the shape of a type for each subterm, and then assigning \perp or \top to each leaf. Our method for choosing between \perp and \top is simple and it seems possible that it can be generalized to work for richer type systems.

In the following section we define the type inference problem and summarize some known results that will be helpful in the later sections. In Section 3 we define an automaton that will be used when defining possible solutions to the type inference problem, and in Section 4 we define a useful notion of upper bound. In Section 5 we prove our main theorem, which characterizes a least-shape solution, and in Section 6 we present our type inference algorithm. Finally, in Section 7 we give an example of how our type inference

algorithm works.

2 The type system

2.1 Type rules

If E is a λ -term, t is a type, and A is a type environment, *i.e.* a partial function assigning types to variables, then the judgement

$$A \vdash E : t$$

means that E has the type t in the environment A . Formally, this holds when the judgement is derivable using the following four rules:

$$A \vdash x : t \quad (\text{provided } A(x) = t)$$

$$\frac{A[x \leftarrow s] \vdash E : t}{A \vdash \lambda x. E : s \rightarrow t}$$

$$\frac{A \vdash E : s \rightarrow t \quad A \vdash F : s}{A \vdash EF : t}$$

$$\frac{A \vdash E : s \quad s \leq t}{A \vdash E : t}$$

The first three rules are the usual rules for simple types and the last rule is the rule of *subsumption*.

The type system has the subject reduction property, that is, if $A \vdash E : t$ is derivable and E β -reduces to E' , then $A \vdash E' : t$ is derivable. This is proved by straightforward induction on the structure of the derivation of $A \vdash E : t$ [5].

This system types λ -terms which are not typable in the simply-typed λ -calculus. For example, consider $\lambda f.(fK(fI))$, where K and I are the usual combinators. This is not typable in the ordinary calculus, since K and I have different types, but it is typable under partial typing: assign f the type $\top \rightarrow (\top \rightarrow \top)$. Both the K and I can be coerced to type \top , and the result (fI) , of type $\top \rightarrow \top$, can be coerced to \top to form the second argument of the first f . Therefore the entire λ -term has type $(\top \rightarrow (\top \rightarrow \top)) \rightarrow \top$.

Similarly, some self-application is possible: $(\lambda x.xx)$ has type $(\top \rightarrow t) \rightarrow t$ for all t , since the final x can be coerced to \top .

However, not all λ -terms are typable in this system. Any λ -term typable in PTB is strongly normalizing [10]. To indicate the flavor of those strongly normalizing λ -terms which are not typable in this system, we have earlier [10] showed that $(\lambda x.xxx)(\lambda y.y)$ is not typable in PTB. We also showed that $(\lambda f.f(fx))(\lambda v.vy)$ is typable in PTB but *not* in PT. This demonstrates that \perp adds power to a type system.

2.2 Constraints

Given a λ -term E , the type inference problem can be rephrased in terms of solving a system of type constraints. Assume that E has been α -converted so that all bound variables are distinct. Let X_E be the set of λ -variables x occurring in E , and let Y_E be a set of variables disjoint from X_E consisting of one variable $\llbracket F \rrbracket$ for each occurrence of a subterm F of E . (The notation $\llbracket F \rrbracket$ is ambiguous because there may be more than one occurrence of F in E . However, it will always be clear from context which occurrence is meant.) We generate the following system of inequalities over $X_E \cup Y_E$:

- for every occurrence in E of a subterm of the form $\lambda x.F$, the inequality

$$x \rightarrow \llbracket F \rrbracket \leq \llbracket \lambda x.F \rrbracket ;$$

- for every occurrence in E of a subterm of the form GH , the inequality

$$\llbracket G \rrbracket \leq \llbracket H \rrbracket \rightarrow \llbracket GH \rrbracket ;$$

- for every occurrence in E of a λ -variable x , the inequality

$$x \leq \llbracket x \rrbracket .$$

Denote by $T(E)$ the system of constraints generated from E in this fashion. For every λ -term E , let $\mathbf{Tmap}(E)$ be the set of total functions from $X_E \cup Y_E$ to the set of types. The function $\psi \in \mathbf{Tmap}(E)$ is a *solution* of $T(E)$ if and only if it is a solution of each constraint in $T(E)$. Specifically, for $V, V', V'' \in X_E \cup Y_E$:

The constraint:	has solution ψ if:
$V \rightarrow V' \leq V''$	$\psi(V) \rightarrow \psi(V') \leq \psi(V'')$
$V \leq V' \rightarrow V''$	$\psi(V) \leq \psi(V') \rightarrow \psi(V'')$
$V \leq V'$	$\psi(V) \leq \psi(V')$

The solutions of $T(E)$ correspond to the possible type annotations of E in a sense made precise by Theorem 1.

Let A be a type environment assigning a type to each λ -variable occurring freely in E . If ψ is a function assigning a type to each variable in $X_E \cup Y_E$, we say that ψ *extends* A if A and ψ agree on the domain of A .

Theorem 1 *The judgement $A \vdash E : t$ is derivable if and only if there exists a solution ψ of $T(E)$ extending A such that $\psi(\llbracket E \rrbracket) = t$. In particular, if E is closed, then E is typable with type t if and only if there exists a solution ψ of $T(E)$ such that $\psi(\llbracket E \rrbracket) = t$.*

Proof. Similar to the proof of Theorem 2.1 in the journal version of [2], in outline as follows. Given a solution of the constraint system, it is straightforward to construct a derivation of $A \vdash E : t$. Conversely, observe that if $A \vdash E : t$ is derivable, then there exists a derivation of $A \vdash E : t$ such that each use of one of the ordinary rules is followed by exactly one use of the subsumption rule. The approach in for example [9, 7] then gives a set of inequalities of the desired form. \square

We now begin the running example of this paper. It is small and does not fully illustrate all constructions, but it may be helpful for understanding the basic ideas. A larger example is presented in Section 7. Consider the λ -term $E = (\lambda x.x)y$. We have

$$\begin{aligned} X_E &= \{x, y\} \\ Y_E &= \{\llbracket x \rrbracket, \llbracket y \rrbracket, \llbracket \lambda x.x \rrbracket, \llbracket (\lambda x.x)y \rrbracket\} . \end{aligned}$$

The constraint system $T(E)$ looks as follows.

$$\begin{array}{ll} \text{From } \lambda x.x & x \rightarrow \llbracket x \rrbracket \leq \llbracket \lambda x.x \rrbracket \\ \text{From } (\lambda x.x)y & \llbracket \lambda x.x \rrbracket \leq \llbracket y \rrbracket \rightarrow \llbracket (\lambda x.x)y \rrbracket \\ \text{From } x & x \leq \llbracket x \rrbracket \\ \text{From } y & y \leq \llbracket y \rrbracket \end{array}$$

To the left of the constraints, we have indicated from where they arise. The example will be continued later.

2.3 A representation of types

We will use a standard representation of types that will be convenient in Sections 3–6. Intuitively, a type is a binary tree where its interior nodes are labeled with \rightarrow and its leaves are labeled with \top or \perp . A path from the root of a type can be thought of as a string of 0's and 1's. Starting at the root, as long as the nodes along the path are labeled with an arrow, a 0 in the path selects the left subtree, a 1 selects the right subtree. In the following we formalize this intuition, along the lines of for example [3, 6].

Let $\Sigma = \{\rightarrow, \perp, \top\}$ be the ranked alphabet where \rightarrow is binary and \perp, \top are nullary. A *type* is a finite tree over Σ . A *path* from the root of such a tree is a string over $\{0, 1\}$, where 0 indicates “left subtree” and 1 indicates “right subtree”. This set of types is the same as the one defined by the grammar in Section 1. A set S of strings is *prefix closed* if $\alpha\beta \in S$ implies $\alpha \in S$. We represent a type by a *term*, which is a function mapping each path from the root of the type to the symbol at the end of the path, as follows.

Definition 2 A *term* is a partial function

$$t : \{0, 1\}^* \rightarrow \Sigma$$

with domain $\mathcal{D}(t)$ such that $\mathcal{D}(t)$ is non-empty and prefix-closed, and such that if $t(\alpha) = \rightarrow$, then $\{i \mid \alpha i \in \mathcal{D}(t)\} = \{0, 1\}$; and if $t(\alpha) \in \{\perp, \top\}$, then $\{i \mid \alpha i \in \mathcal{D}(t)\} = \emptyset$. The set of all such terms is denoted T_Σ . The set of terms with finite domain is denoted F_Σ . They represent finite trees, by König's Lemma. In our application, F_Σ is the set of types.

For $s, t \in T_\Sigma$, define $\rightarrow^{T_\Sigma} : T_\Sigma \times T_\Sigma \rightarrow T_\Sigma$ and $\perp^{T_\Sigma}, \top^{T_\Sigma} \in T_\Sigma$ with domains

$$\begin{aligned} \mathcal{D}(s \rightarrow^{T_\Sigma} t) &= \{\epsilon\} \cup \{0\alpha \mid \alpha \in \mathcal{D}(s)\} \cup \{1\alpha \mid \alpha \in \mathcal{D}(t)\} \\ \mathcal{D}(\perp^{T_\Sigma}) &= \{\epsilon\} \\ \mathcal{D}(\top^{T_\Sigma}) &= \{\epsilon\} \end{aligned}$$

by

$$\begin{aligned} (s \rightarrow^{T_\Sigma} t)(0\alpha) &= s(\alpha) \\ (s \rightarrow^{T_\Sigma} t)(1\alpha) &= t(\alpha) \\ (s \rightarrow^{T_\Sigma} t)(\epsilon) &= \rightarrow \\ \perp^{T_\Sigma}(\epsilon) &= \perp \\ \top^{T_\Sigma}(\epsilon) &= \top. \end{aligned}$$

For $t \in T_\Sigma$ and $\alpha \in \{0, 1\}^*$, define the partial function $t \downarrow \alpha : \{0, 1\}^* \rightarrow \Sigma$ by

$$t \downarrow \alpha(\beta) = t(\alpha\beta) .$$

If $t \downarrow \alpha$ has non-empty domain, then it is a term, and is called the subterm of t at position α .

For $x \subseteq \{0, 1\}^*$ and $\alpha \in \{0, 1\}^*$, define $x \downarrow \alpha = \{\beta \mid \alpha\beta \in x\}$. Notice that the notation $_ \downarrow _$ is overloaded. Which definition is being used will be clear from context. \square

The following properties are immediate from the definitions:

- (i) $(s \rightarrow^{T_\Sigma} t) \downarrow 0 = s$
- (ii) $(s \rightarrow^{T_\Sigma} t) \downarrow 1 = t$
- (iii) $(t \downarrow \alpha) \downarrow \beta = t \downarrow \alpha\beta$

Following [3], we henceforth omit the superscript T_Σ on the operators $\rightarrow^{T_\Sigma}, \perp^{T_\Sigma}, \top^{T_\Sigma}$.

We will say that a term s has smaller shape than a term t if $\mathcal{D}(s) \subseteq \mathcal{D}(t)$. Our type inference algorithm will infer types that are of least shape.

Terms are ordered by the subtype relation \leq , as follows.

Definition 3 The *parity* of $\alpha \in \{0, 1\}^*$ is the number mod 2 of 0's in α . The parity of α is denoted $\pi\alpha$. A string α is said to be *even* if $\pi\alpha = 0$ and *odd* if $\pi\alpha = 1$. Let \leq_0 be the partial order on Σ given by

$$\perp \leq_0 \rightarrow \quad \text{and} \quad \rightarrow \leq_0 \top$$

and let \leq_1 be its reverse

$$\top \leq_1 \rightarrow \quad \text{and} \quad \rightarrow \leq_1 \perp .$$

For $s, t \in T_\Sigma$, define $s \leq t$ if $s(\alpha) \leq_{\pi\alpha} t(\alpha)$ for all $\alpha \in \mathcal{D}(s) \cap \mathcal{D}(t)$. \square

Kozen, Palsberg, and Schwartzbach [3] showed that the relation \leq is equivalent to the order defined by Amadio and Cardelli [1]. The relation \leq is a partial order with $\perp \leq t \leq \top$ for all $t \in T_\Sigma$, and $s \rightarrow t \leq s' \rightarrow t'$ if and only if $s' \leq s$ and $t \leq t'$ [1, 3]. Moreover, on F_Σ , \leq coincides with the type ordering defined in Section 1 [3].

2.4 Graphs

In this section we reduce the problem of solving constraint systems to a problem of solving a certain kind of constraint graph. The graphs yield a convenient setting for applying algorithms like transitive closure.

Definition 4 A *constraint graph* is a directed graph $G = (S, L, R, \leq)$ consisting of a set of nodes S and three sets of directed edges L, R, \leq , where L indicates “left” and R indicates “right”. We write $s \xrightarrow{0} t$ to indicate that the pair (s, t) is in the edge set L , we write $s \xrightarrow{1} t$ to indicate that the pair (s, t) is in the edge set R , and we write $s \xrightarrow{\leq} t$ to indicate that the pair (s, t) is in the edge set \leq . A constraint graph must satisfy the properties:

- any node has at most one outgoing L edge and at most one outgoing R edge;
- a node has an outgoing L edge if and only if it has an outgoing R edge.

A *solution* for G is any map $h : S \rightarrow T_\Sigma$ such that

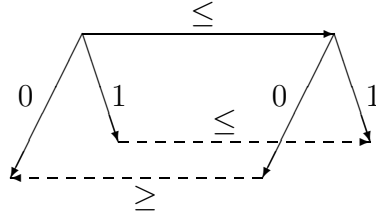
- (i) if $u \xrightarrow{0} v$ and $u \xrightarrow{1} w$, then $h(u) = h(v) \rightarrow h(w)$;
- (ii) if $u \xrightarrow{\leq} v$, then $h(u) \leq h(v)$.

Notice that we consider solutions in $S \rightarrow T_\Sigma$, not just in $S \rightarrow F_\Sigma$. The solution h is *finite* if $h(s)$ is a finite set for all s . \square

A system of type constraints as described above gives rise to a constraint graph by associating a unique node with every subexpression occurring in the system of constraints, defining L and R edges from an occurrence of an expression to its left and right subexpressions, and defining \leq edges for the inequalities. Clearly, a solution of a system of type constraints can be transformed into a finite solution of the corresponding constraint graph, and vice versa.

Definition 5 A constraint graph is *closed* if the edge relation \leq is reflexive, transitive, and closed under the following rule which says that the dashed

edges exist whenever the solid ones do:



The *closure* of a constraint graph G is the smallest closed graph containing G as a subgraph. \square

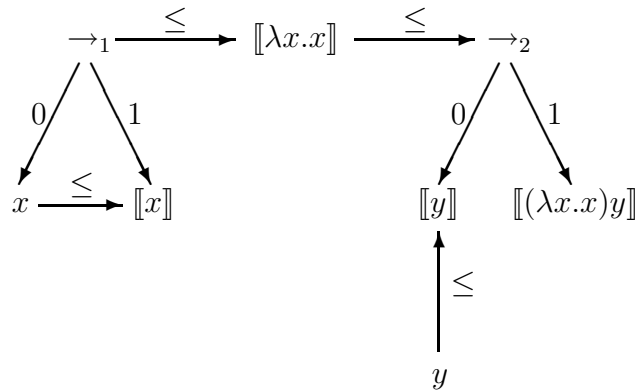
Lemma 6 *A constraint graph and its closure have the same set of solutions.*

Proof. Any solution of the closure of G is also a solution of G , since G has fewer constraints. Conversely, the closure of G can be constructed from G by iterating the closure rules, and it follows inductively that any solution of G satisfies the additional constraints added by this process. \square

Consider again the running example $E = (\lambda x.x)y$. We use the following abbreviations

The symbol:	abbreviates:
\rightarrow_1	$x \rightarrow \llbracket x \rrbracket$
\rightarrow_2	$\llbracket y \rrbracket \rightarrow \llbracket (\lambda x.x)y \rrbracket$

The constraint graph generated from $T(E)$ looks as follows.



To close the constraint graph, it is sufficient to fill in the edges

$$\begin{aligned} \llbracket x \rrbracket &\stackrel{\leq}{\rightarrow} \llbracket (\lambda x.x)y \rrbracket \\ \llbracket y \rrbracket &\stackrel{\leq}{\rightarrow} x \end{aligned}$$

and make \leq reflexive and transitive.

2.5 On solving graphs

In Sections 3–6 we give an algorithm for solving constraint graphs. Here follows an informal account of the basic idea.

Every constraint graph is solvable. The basic problem is: can we solve a given graph using only finite types, that is, using only elements of F_Σ ? Our approach to doing this is to compute a least-shape solution. The idea is that if a least-shape solution uses at least one infinite type then so does every other solution.

A least-shape solution assigns arrow types to as few nodes as possible, and it assigns \top or \perp to the remaining nodes. For example, consider again the constraint graph from the running example. Suppose h is a solution of the graph. Rule (i) in Definition 4 forces both $h(\rightarrow_1)$ and $h(\rightarrow_2)$ to be an arrow type. Moreover, rule (ii) in Definition 4 forces $h(\rightarrow_1) \leq h(\llbracket \lambda x.x \rrbracket) \leq h(\rightarrow_2)$. Now, since $h(\llbracket \lambda x.x \rrbracket)$ is a type, it must be either \perp , \top or an arrow type. It cannot be \perp , because $h(\rightarrow_1) \leq \perp$ is false. Likewise $h(\llbracket \lambda x.x \rrbracket)$ cannot be \top , because $\top \leq h(\rightarrow_2)$ is false. Therefore $h(\llbracket \lambda x.x \rrbracket)$ must be an arrow type.

A naive algorithm might work by creating two fresh nodes to hold the left and right subexpressions of $h(\llbracket \lambda x.x \rrbracket)$. Such a strategy may in the worst case require exponential space for storing the nodes of the graph. To see this, consider the constraint system

$$\begin{aligned} V_1 &= V_0 \rightarrow V_0 \\ V_2 &= V_1 \rightarrow V_1 \\ &\dots \\ V_n &= V_{n-1} \rightarrow V_{n-1} \end{aligned}$$

where each constraint of the form $X = Y$ denotes the two constraints $X \leq Y$ and $Y \leq X$. If we use the naive algorithm mentioned before, then at least

2^n nodes will be generated. Moreover, care needs to be taken to ensure that the naive algorithm terminates in cases where an infinite type is required.

To sidestep these problems we shift the focus from the question

“which nodes have to be assigned an arrow type?”

to the question

“what set of paths have to be in the type for a given node?”

The first question is easy to answer, but it has to be repeated and answered many times before a solution is constructed. The second question seems more difficult, but it turns out that we can encode the necessary paths using time and space which is polynomial in the size of the graph. The basic idea is to use a nondeterministic automaton, as shown in the following section.

3 An automaton

We now present the key construction (Definition 7) that we will use to efficiently solve constraint graphs. Given a closed constraint graph, the idea is to construct an automaton that

- keeps track of upper and lower bounds in the graph, and
- accepts paths that must exist in every solution of the graph.

In the following sections we will use the automaton to construct a solution of least shape. Each state in the automaton is a pair of nodes from the graph. Moreover, we ensure that if a state (u', v') is reachable from a state of the form (u, u) , then $u' \preceq v'$ is in the graph (Lemma 8). The language accepted by the automaton is a set of paths. If s is a node in the constraint graph and we take (s, s) as the start state of the automaton, then the accepted paths *must* be in the domain of every type for s (Lemma 9). The present section includes three useful results (Lemmas 10, 11, 12) about the automaton; they will be applied in later sections.

Definition 7 Let a constraint graph $G = (S, L, R, \preceq)$ be given. The automaton \mathcal{M} is defined as follows. The input alphabet of \mathcal{M} is $\{0, 1\}$. The

set of states of \mathcal{M} is $S \times S$. A state is written (s, t) . The transitions are defined as follows.

$$\begin{aligned} (u, v) &\xrightarrow{\epsilon} (u, v') && \text{if } v \xrightarrow{\leq} v' \text{ in } G \\ (u, v) &\xrightarrow{\epsilon} (u', v) && \text{if } u' \xrightarrow{\leq} u \text{ in } G \\ (u, v) &\xrightarrow{1} (u', v') && \text{if } u \xrightarrow{1} u' \text{ and } v \xrightarrow{1} v' \text{ in } G \\ (u, v) &\xrightarrow{0} (v', u') && \text{if } u \xrightarrow{0} u' \text{ and } v \xrightarrow{0} v' \text{ in } G \end{aligned}$$

If p and q are states of \mathcal{M} and $\alpha \in \{0, 1\}^*$, we write $p \xrightarrow{\alpha} q$ if the automaton can move from state p to state q under input α , including possible ϵ -transitions.

The automaton \mathcal{M}_s is the automaton \mathcal{M} with start state (s, s) . All states are accept states; thus the language accepted by \mathcal{M}_s is the set of strings α for which there exists a state (u, v) such that $(s, s) \xrightarrow{\alpha} (u, v)$. We denote this language by $\mathcal{L}(s)$. \square

Informally, we can think of the automaton \mathcal{M}_s as follows. We start with two pebbles, one green and one red, on the node s of the constraint graph G . We can move the green pebble forward along a \leq edge at any time, and we can move the red pebble backward along a \leq edge at any time. We can move both pebbles simultaneously along R edges leading out of the nodes they occupy. We can also move them simultaneously along outgoing L edges, but in the latter case we switch their colors. The sequence of 0's and 1's that were seen gives a string in $\mathcal{L}(s)$, and all strings in $\mathcal{L}(s)$ are obtained in this way.

For comparison, the automaton that was defined in [2] for doing type inference without \perp had eight rather than four rules for generating transitions.

It turns out that once we have \mathcal{M} , we can construct a solution with domain $\mathcal{L}(s)$ for each s .

The fundamental relation between a closed constraint graph G and the automaton \mathcal{M} is expressed by Lemma 8.

Lemma 8 *If G is closed and $(u, u) \xrightarrow{\alpha} (u', v')$, then $u' \xrightarrow{\leq} v'$.*

Proof. By induction on the number of transitions in $(u, u) \xrightarrow{\epsilon} (u', v')$. \square

Lemma 9 *If G is closed, $h : S \rightarrow T_\Sigma$ is any solution of G , and $(s, s) \xrightarrow{\alpha} (u, v)$, then $\alpha \in \mathcal{D}(h(s))$. Moreover, $h(u) \leq h(s) \downarrow \alpha \leq h(v)$.*

Proof. We proceed by induction on the number of transitions. If this is zero, then $(u, v) = (s, s)$ and $\alpha = \epsilon$, and the result is immediate. Otherwise, assume that $(s, s) \xrightarrow{\alpha} (u, v)$ and the lemma holds for this sequence of transitions. We argue by cases, depending on the form of the next transition out of (u, v) .

If $(u, v) \xrightarrow{\epsilon} (u', v')$, then $u' \xrightarrow{\leq} u$ and $v \xrightarrow{\leq} v'$, so $\alpha\epsilon = \alpha \in \mathcal{D}(h(s))$ and

$$h(u') \leq h(u) \leq h(s) \downarrow \alpha \leq h(v) \leq h(v').$$

If $(u, v) \xrightarrow{1} (u', v')$, then $u \xrightarrow{1} u'$ and $v \xrightarrow{1} v'$, so $h(u') = h(u) \downarrow 1$ and $h(v') = h(v) \downarrow 1$. Then $1 \in \mathcal{D}(h(u))$ and $1 \in \mathcal{D}(h(v))$, so $1 \in \mathcal{D}(h(s) \downarrow \alpha)$ and $\alpha 1 \in \mathcal{D}(h(s))$, and

$$h(u') = h(u) \downarrow 1 \leq h(s) \downarrow \alpha 1 \leq h(v) \downarrow 1 = h(v').$$

If $(u, v) \xrightarrow{0} (v', u')$, then $u \xrightarrow{0} u'$ and $v \xrightarrow{0} v'$, so $h(u') = h(u) \downarrow 0$ and $h(v') = h(v) \downarrow 0$. Then $0 \in \mathcal{D}(h(u))$ and $0 \in \mathcal{D}(h(v))$, so $0 \in \mathcal{D}(h(s) \downarrow \alpha)$ and $\alpha 0 \in \mathcal{D}(h(s))$, and

$$h(u') = h(u) \downarrow 0 \geq h(s) \downarrow \alpha 0 \geq h(v) \downarrow 0 = h(v').$$

□

When certain paths in \mathcal{M} exist, others must exist too, as expressed in Lemmas 10, 11, and 12.

Define

$$\mathcal{N} = \{u \in S \mid u \text{ has outgoing } L \text{ and } R \text{ edges in } G\}$$

Intuitively, \mathcal{N} is the set of “ \rightarrow nodes.”

Lemma 10 *If $(u, v) \xrightarrow{\epsilon} (u', v')$, then for any $w \in S$, $(w, v) \xrightarrow{\epsilon} (w, v')$ and $(u, w) \xrightarrow{\epsilon} (u', w)$.*

Proof. By induction on the number of transitions in $(u, v) \xrightarrow{\epsilon} (u', v')$. □

Lemma 11 *If $(u_1, v_1) \xrightarrow{\alpha} (u'_1, v'_1)$ and $(u_2, v_2) \xrightarrow{\alpha} (u'_2, v'_2)$, then $(u_1, v_2) \xrightarrow{\alpha} (u'_1, v'_2)$.*

Proof. We proceed by induction on the length of α . In the base case, consider α with length 0. We then have $\alpha = \epsilon$, so by Lemma 10,

$$(u_1, v_2) \xrightarrow{\epsilon} (u'_1, v_2) \xrightarrow{\epsilon} (u'_1, v'_2) .$$

In the induction step, consider first

$$\begin{aligned} (u_1, v_1) &\xrightarrow{\alpha 1} (u'_1, v'_1) \\ (u_2, v_2) &\xrightarrow{\alpha 1} (u'_2, v'_2) . \end{aligned}$$

In both cases, let us make explicit the 1-transition, that is, find $x_1, y_1, x_2, y_2 \in \mathcal{N}$ and $x'_1, y'_1, x'_2, y'_2 \in S$ such that

$$\begin{aligned} (u_1, v_1) &\xrightarrow{\alpha} (x_1, y_1) \xrightarrow{1} (x'_1, y'_1) \xrightarrow{\epsilon} (u'_1, v'_1) \\ (u_2, v_2) &\xrightarrow{\alpha} (x_2, y_2) \xrightarrow{1} (x'_2, y'_2) \xrightarrow{\epsilon} (u'_2, v'_2) \end{aligned}$$

and $x_1 \xrightarrow{1} x'_1, y_1 \xrightarrow{1} y'_1, x_2 \xrightarrow{1} x'_2, y_2 \xrightarrow{1} y'_2$. From the induction hypothesis and Lemma 10 we get

$$(u_1, v_2) \xrightarrow{\alpha} (x_1, y_2) \xrightarrow{1} (x'_1, y'_2) \xrightarrow{\epsilon} (u'_1, y'_2) \xrightarrow{\epsilon} (u'_1, v'_2) .$$

Consider then

$$\begin{aligned} (u_1, v_1) &\xrightarrow{\alpha 0} (u'_1, v'_1) \\ (u_2, v_2) &\xrightarrow{\alpha 0} (u'_2, v'_2) . \end{aligned}$$

Again, let us make explicit the 0-transition, that is, find $x_1, y_1, x_2, y_2 \in \mathcal{N}$ and $x'_1, y'_1, x'_2, y'_2 \in S$ such that

$$\begin{aligned} (u_1, v_1) &\xrightarrow{\alpha} (x_1, y_1) \xrightarrow{0} (y'_1, x'_1) \xrightarrow{\epsilon} (u'_1, v'_1) \\ (u_2, v_2) &\xrightarrow{\alpha} (x_2, y_2) \xrightarrow{0} (y'_2, x'_2) \xrightarrow{\epsilon} (u'_2, v'_2) \end{aligned}$$

and $x_1 \xrightarrow{0} x'_1, y_1 \xrightarrow{0} y'_1, x_2 \xrightarrow{0} x'_2, y_2 \xrightarrow{0} y'_2$. From the induction hypothesis and Lemma 10 we get

$$(u_1, v_2) \xrightarrow{\alpha} (x_2, y_1) \xrightarrow{0} (y'_1, x'_2) \xrightarrow{\epsilon} (u'_1, x'_2) \xrightarrow{\epsilon} (u'_1, v'_2) .$$

□

Lemma 12 Suppose $i \in \{0, 1\}$, $u \xrightarrow{i} v$, and $\alpha \in \mathcal{L}(v)$. If G is closed, then $\mathcal{L}(v) \downarrow \alpha = \mathcal{L}(u) \downarrow i\alpha$.

Proof. Clearly, $\mathcal{L}(v) \downarrow \alpha \subseteq \mathcal{L}(u) \downarrow i\alpha$. To prove the converse, suppose

$$(u, u) \xrightarrow{i} (u_1, v_1) \xrightarrow{\alpha} (x, y) \xrightarrow{\beta} (x', y') .$$

From $(u, u) \xrightarrow{i} (v, v)$ and Lemma 11 we get

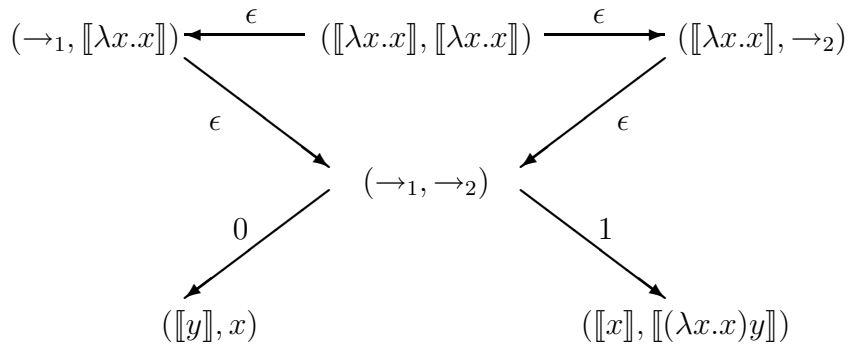
$$\begin{aligned} (u, u) &\xrightarrow{i} (u_1, v) \\ (u, u) &\xrightarrow{i} (v, v_1) . \end{aligned}$$

From Lemma 8 we get $u_1 \xrightarrow{\leq} v \xrightarrow{\leq} v_1$, hence

$$(v, v) \xrightarrow{\epsilon} (u_1, v_1) \xrightarrow{\alpha} (x, y) \xrightarrow{\beta} (x', y') ,$$

so $\mathcal{L}(v) \downarrow \alpha \supseteq \mathcal{L}(u) \downarrow i\alpha$. □

Consider again the running example $E = (\lambda x.x)y$. The constraint graph generated from $T(E)$ has 8 nodes, so the automaton \mathcal{M} has $8^2 = 64$ states. The following picture shows some of the states reachable from the state $(\llbracket \lambda x.x \rrbracket, \llbracket \lambda x.x \rrbracket)$.



It turns out that $\mathcal{L}(x) = \mathcal{L}(y) = \{\epsilon\}$ and $\mathcal{L}(\llbracket \lambda x.x \rrbracket) = \{\epsilon, 0, 1\}$.

4 Upper bounds

We want to construct a solution $\psi : S \rightarrow T_\Sigma$ such that $\mathcal{D}(\psi(s)) = \mathcal{L}(s)$ for all $s \in S$. Suppose $\mathcal{L}(s) \downarrow \alpha = \{\epsilon\}$ for some $s \in S$. Then we want $(\psi(s))(\alpha) \in \{\perp, \top\}$. The question is: should we pick \perp or \top ? Let us examine this question by considering

$$(s, s) \xrightarrow{\alpha} (u, w).$$

Recall that \mathcal{N} is the set of “ \rightarrow nodes.” Clearly, at most one of u and w is in \mathcal{N} ; otherwise (u, w) would have outgoing 0 and 1 transitions, contradicting $\mathcal{L}(s) \downarrow \alpha = \{\epsilon\}$. If either of u and w is in \mathcal{N} , then that puts a constraint on the choice for the value of $(\psi(s))(\alpha)$. Intuitively, if w is in \mathcal{N} , then we must define $(\psi(s))(\alpha) = \perp$, and if u is in \mathcal{N} , then we must define $(\psi(s))(\alpha) = \top$. If neither of u and w are in \mathcal{N} , then no constraint is put on $(\psi(s))(\alpha)$. To capture the set of imposed constraints, we define $\text{Up}(s, \alpha)$, which intuitively is the set of upper bounds of s at level α (Definition 13).

We let $2^{\mathcal{N}}$ denote the powerset of \mathcal{N} .

Definition 13 Define $\text{Up} : S \times \{0, 1\}^* \rightarrow 2^{\mathcal{N}}$ as follows:

$$\text{Up}(s, \alpha) = \{w \in \mathcal{N} \mid (s, s) \xrightarrow{\alpha} (u, w) \text{ for some } u \in S\}$$

□

Next, we prove two lemmas that give conditions for when $\text{Up}(s, \alpha)$ is empty.

Lemma 14 Suppose $u \xrightarrow{\leq} v$ and $\alpha \in \mathcal{L}(u) \cap \mathcal{L}(v)$.

1. If $\text{Up}(u, \alpha) = \emptyset$, then $\text{Up}(v, \alpha) = \emptyset$.
2. If $\text{Up}(v, \alpha) = \emptyset$ and $\pi\alpha = 1$, then $\text{Up}(u, \alpha) = \emptyset$.
3. If $\mathcal{L}(u) \downarrow \alpha \neq \{\epsilon\}$ and $\mathcal{L}(v) \downarrow \alpha = \{\epsilon\}$, then $\text{Up}(v, \alpha) = \emptyset$.
4. If $\mathcal{L}(v) \downarrow \alpha \neq \{\epsilon\}$, $\mathcal{L}(u) \downarrow \alpha = \{\epsilon\}$, and $\pi\alpha = 1$, then $\text{Up}(u, \alpha) = \emptyset$.

Proof. To prove (1), suppose $\text{Up}(v, \alpha) \neq \emptyset$. Choose $w \in \mathcal{N}$ and $u' \in S$ such that $(v, v) \xrightarrow{\alpha} (u', w)$. Choose also $u_1, v_1 \in S$ such that $(u, u) \xrightarrow{\alpha} (u_1, v_1)$. From $u \xrightarrow{\leq} v$ and Lemma 11 we then get

$$(u, u) \xrightarrow{\epsilon} (u, v) \xrightarrow{\alpha} (u_1, w) ,$$

contradicting $\text{Up}(u, \alpha) = \emptyset$.

To prove (2), suppose $\text{Up}(u, \alpha) \neq \emptyset$. Since $\pi\alpha = 1$, we can write $\alpha = \beta 0 \gamma$ where $\beta \in \{0, 1\}^*$ and $\gamma \in \{1\}^*$. Choose $x_1, y_1, w \in \mathcal{N}$ and $x_2, y_2, u' \in S$ such that

$$(u, u) \xrightarrow{\beta} (x_1, y_1) \xrightarrow{0} (y_2, x_2) \xrightarrow{\gamma} (u', w)$$

and $x_1 \xrightarrow{0} x_2$ and $y_1 \xrightarrow{0} y_2$. Choose $p_1, q_1 \in \mathcal{N}$ and $p_2, q_2, u'', v'' \in S$ such that

$$(v, v) \xrightarrow{\beta} (p_1, q_1) \xrightarrow{0} (q_2, p_2) \xrightarrow{\gamma} (u'', v'')$$

and $p_1 \xrightarrow{0} p_2$ and $q_1 \xrightarrow{0} q_2$. From $u \xrightarrow{\leq} v$ and Lemma 11 we then get

$$(v, v) \xrightarrow{\epsilon} (u, v) \xrightarrow{\beta} (x_1, q_1) \xrightarrow{0} (q_2, x_2) \xrightarrow{\gamma} (u'', w) ,$$

contradicting $\text{Up}(v, \alpha) = \emptyset$.

To prove (3), suppose $\text{Up}(v, \alpha) \neq \emptyset$. Choose $w \in \mathcal{N}$ and $u' \in S$ such that $(v, v) \xrightarrow{\alpha} (u', w)$. Choose also $x, y \in \mathcal{N}$ such that $(u, u) \xrightarrow{\alpha} (x, y)$. From $u \xrightarrow{\leq} v$ and Lemma 11 we get

$$(v, v) \xrightarrow{\epsilon} (u, v) \xrightarrow{\alpha} (x, w) ,$$

contradicting $\mathcal{L}(v) \downarrow \alpha = \{\epsilon\}$.

To prove (4), suppose $\text{Up}(u, \alpha) \neq \emptyset$. Since $\pi\alpha = 1$, we can write $\alpha = \beta 0 \gamma$ where $\beta \in \{0, 1\}^*$ and $\gamma \in \{1\}^*$. Choose $x_1, y_1, w \in \mathcal{N}$ and $x_2, y_2, u' \in S$ such that

$$(u, u) \xrightarrow{\beta} (x_1, y_1) \xrightarrow{0} (y_2, x_2) \xrightarrow{\gamma} (u', w)$$

and $x_1 \xrightarrow{0} x_2$ and $y_1 \xrightarrow{0} y_2$. Since $\mathcal{L}(v) \downarrow \alpha \neq \{\epsilon\}$, choose $p_1, q_1, u'', v'' \in \mathcal{N}$ and $p_2, q_2 \in S$ such that

$$(v, v) \xrightarrow{\beta} (p_1, q_1) \xrightarrow{0} (q_2, p_2) \xrightarrow{\gamma} (u'', v'')$$

and $p_1 \xrightarrow{0} p_2$ and $q_1 \xrightarrow{0} q_2$. From $u \xrightarrow{\leq} v$ and Lemma 11 we then get

$$(u, u) \xrightarrow{\epsilon} (u, v) \xrightarrow{\beta} (x_1, q_1) \xrightarrow{0} (q_2, x_2) \xrightarrow{\gamma} (u'', w) ,$$

contradicting $\mathcal{L}(u) \downarrow \alpha = \emptyset$. □

Lemma 15 *Suppose $i \in \{0, 1\}$, $u \xrightarrow{i} v$, and $\alpha \in \mathcal{L}(v)$.*

1. *If $\text{Up}(u, i\alpha) = \emptyset$, then $\text{Up}(v, \alpha) = \emptyset$.*
2. *Suppose G is closed. If $\text{Up}(v, \alpha) = \emptyset$, then $\text{Up}(u, i\alpha) = \emptyset$.*

Proof. To prove (1), suppose $\text{Up}(v, \alpha) \neq \emptyset$. Choose $w \in \mathcal{N}$ and $u' \in S$ such that

$$(u, u) \xrightarrow{i} (v, v) \xrightarrow{\alpha} (u', w) ,$$

contradicting $\text{Up}(u, i\alpha) = \emptyset$.

To prove (2), suppose $\text{Up}(u, i\alpha) \neq \emptyset$. Choose $w \in \mathcal{N}$ and $u', u_1, v_1 \in S$ such that

$$(u, u) \xrightarrow{i} (u_1, v_1) \xrightarrow{\alpha} (u', w) .$$

From $(u, u) \xrightarrow{i} (v, v)$ and Lemma 11 we get

$$\begin{aligned} (u, u) &\xrightarrow{i} (u_1, v) \\ (u, u) &\xrightarrow{i} (v, v_1) . \end{aligned}$$

From Lemma 8 we get $u_1 \xrightarrow{\leq} v \xrightarrow{\leq} v_1$, hence

$$(v, v) \xrightarrow{\epsilon} (u_1, v_1) \xrightarrow{\alpha} (u', w) ,$$

contradicting $\text{Up}(v, \alpha) = \emptyset$. □

Consider again the running example $E = (\lambda x.x)y$. We have

$$\begin{aligned} \text{Up}(x, \epsilon) &= \text{Up}(y, \epsilon) = \text{Up}(\llbracket \lambda x.x \rrbracket, 0) = \text{Up}(\llbracket \lambda x.x \rrbracket, 1) = \emptyset \\ \text{Up}(\llbracket \lambda x.x \rrbracket, \epsilon) &= \{\rightarrow_2\} . \end{aligned}$$

5 Main result

We can now define the least-shape solution of a constraint graph. The definition assigns \perp or \top to each leaf of $\mathcal{L}(s)$, with the help of $\text{Up}(s, \alpha)$. If neither \perp nor \top is forced for a certain leaf, then we arbitrarily choose \top .

Definition 16 Let a closed constraint graph $G = (S, L, R, \leq)$ be given. Define $\psi : S \rightarrow T_\Sigma$ as follows:

$$\psi(s) = \lambda\alpha. \begin{cases} \top & \text{if } \alpha \in \mathcal{L}(s) \wedge \mathcal{L}(s) \downarrow \alpha = \{\epsilon\} \wedge \text{Up}(s, \alpha) = \emptyset \\ \perp & \text{if } \alpha \in \mathcal{L}(s) \wedge \mathcal{L}(s) \downarrow \alpha = \{\epsilon\} \wedge \text{Up}(s, \alpha) \neq \emptyset \\ \rightarrow & \text{if } \alpha \in \mathcal{L}(s) \wedge \mathcal{L}(s) \downarrow \alpha \neq \{\epsilon\} \\ \text{undefined} & \text{if } \alpha \notin \mathcal{L}(s) \end{cases}$$

Notice that each $\psi(s)$ has domain $\mathcal{L}(s)$ for all $s \in S$. □

Theorem 17 *If G is a closed constraint graph, then the function ψ is a solution of G . Moreover, if $h : S \rightarrow T_\Sigma$ is any other solution of G , then $\mathcal{D}(\psi(s)) \subseteq \mathcal{D}(h(s))$ for any s .*

Proof. For any $s \in S$, notice that $\psi(s)$ is nonempty since $(s, s) \xrightarrow{\epsilon} (s, s)$, that it is prefix-closed since all states of \mathcal{M}_s are accept states, that if $(\psi(s))(\alpha) = \rightarrow$ then $\{i \mid \alpha i \in \mathcal{D}(\psi(s))\} = \{0, 1\}$, and that if $(\psi(s))(\alpha) \in \{\perp, \top\}$ then $\{i \mid \alpha i \in \mathcal{D}(\psi(s))\} = \emptyset$, so $\psi(s) \in T_\Sigma$.

We first prove that

$$\text{If } u \xrightarrow{0} v_0 \text{ and } u \xrightarrow{1} v_1, \text{ then } \psi(u) = \psi(v_0) \rightarrow \psi(v_1).$$

We begin by reformulating $\psi(u) = \psi(v_0) \rightarrow \psi(v_1)$. Notice that it is equivalent to the conjunction of

$$\begin{aligned} \psi(u) &\leq \psi(v_0) \rightarrow \psi(v_1) \text{ and} \\ \psi(u) &\geq \psi(v_0) \rightarrow \psi(v_1). \end{aligned}$$

From the definition of \leq (Definition 3) we get that this conjunction is equivalent to the conjunction of

$$\begin{aligned} (\psi(u))(\alpha) &\leq_{\pi\alpha} (\psi(v_0) \rightarrow \psi(v_1))(\alpha) \text{ for all } \alpha \in D, \text{ and} \\ (\psi(u))(\alpha) &\geq_{\pi\alpha} (\psi(v_0) \rightarrow \psi(v_1))(\alpha) \text{ for all } \alpha \in D, \end{aligned}$$

where $D = \mathcal{D}(\psi(u)) \cap \mathcal{D}(\psi(v_0) \rightarrow \psi(v_1))$. Since $\leq_{\pi\alpha}$ is a partial order, we can fold the second conjunction into the single condition

$$(\psi(u))(\alpha) = (\psi(v_0) \rightarrow \psi(v_1))(\alpha) \text{ for all } \alpha \in D.$$

To prove that $u \xrightarrow{0} v_0$ and $u \xrightarrow{1} v_1$ implies this condition, notice first that $(\psi(u))(\epsilon) = (\psi(v_0) \rightarrow \psi(v_1))(\epsilon) \Rightarrow$. Then, let $i \in \{0, 1\}$ and $i\alpha \in \mathcal{D}(\psi(u)) \cap \mathcal{D}(\psi(v_0) \rightarrow \psi(v_1))$. It is sufficient to prove

$$(\psi(u))(i\alpha) = (\psi(v_i))(\alpha) .$$

There are three cases. Suppose $(\psi(v_i))(\alpha) \Rightarrow$. Then $\mathcal{L}(v_i) \downarrow \alpha \neq \{\epsilon\}$. From Lemma 12 we get $\mathcal{L}(u) \downarrow i\alpha \neq \{\epsilon\}$, so $(\psi(u))(i\alpha) \Rightarrow$. Suppose then $(\psi(v_i))(\alpha) = \perp$. Then $\mathcal{L}(v_i) \downarrow \alpha = \{\epsilon\}$ and $\mathbf{Up}(v_i, \alpha) \neq \emptyset$. From Lemma 15.1 we get $\mathbf{Up}(u, i\alpha) \neq \emptyset$. From Lemma 12 we get $\mathcal{L}(u) \downarrow i\alpha = \{\epsilon\}$, so $(\psi(u))(i\alpha) = \perp$. Suppose finally $(\psi(v_i))(\alpha) = \top$. Then $\mathcal{L}(v_i) \downarrow \alpha = \{\epsilon\}$ and $\mathbf{Up}(v_i, \alpha) = \emptyset$. From Lemma 15.2 we get $\mathbf{Up}(u, i\alpha) = \emptyset$. From Lemma 12 we get $\mathcal{L}(u) \downarrow i\alpha = \{\epsilon\}$, so $(\psi(u))(i\alpha) = \top$.

We then prove that

$$\text{If } u \xrightarrow{\leq} v, \text{ then } \psi(u) \leq \psi(v).$$

Suppose $\alpha \in \mathcal{D}(\psi(u)) \cap \mathcal{D}(\psi(v)) = \mathcal{L}(u) \cap \mathcal{L}(v)$. We then need to prove $(\psi(u))(\alpha) \leq_{\pi\alpha} (\psi(v))(\alpha)$. There are two cases.

Suppose first $\pi\alpha = 0$. If $(\psi(u))(\alpha) = \perp$, then the result is immediate. If $(\psi(u))(\alpha) = \top$, then $\mathcal{L}(u) \downarrow \alpha = \{\epsilon\}$ and $\mathbf{Up}(u, \alpha) = \emptyset$. By Lemma 14.1, $\mathbf{Up}(v, \alpha) = \emptyset$. Hence, $\mathcal{L}(v) \downarrow \alpha = \{\epsilon\}$, so $(\psi(v))(\alpha) = \top$, from which the result follows. If $(\psi(u))(\alpha) \Rightarrow$, then there are two cases. If $\mathcal{L}(v) \downarrow \alpha \neq \{\epsilon\}$, then $(\psi(v))(\alpha) \Rightarrow$, from which the result follows. If $\mathcal{L}(v) \downarrow \alpha = \{\epsilon\}$, then it follows from Lemma 14.3 that $\mathbf{Up}(v, \alpha) = \emptyset$. Hence, $(\psi(v))(\alpha) = \top$, from which the result follows.

Suppose then $\pi\alpha = 1$. If $(\psi(v))(\alpha) = \perp$, then the result is immediate. If $(\psi(v))(\alpha) = \top$, then $\mathcal{L}(v) \downarrow \alpha = \{\epsilon\}$ and $\mathbf{Up}(v, \alpha) = \emptyset$. By Lemma 14.2, $\mathbf{Up}(u, \alpha) = \emptyset$. Hence, $\mathcal{L}(u) \downarrow \alpha = \{\epsilon\}$, so $(\psi(u))(\alpha) = \top$, from which the result follows. If $(\psi(v))(\alpha) \Rightarrow$, then there are two cases. If $\mathcal{L}(u) \downarrow \alpha \neq \{\epsilon\}$, then $(\psi(u))(\alpha) \Rightarrow$, from which the result follows. If $\mathcal{L}(u) \downarrow \alpha = \{\epsilon\}$, then it follows from Lemma 14.4 that $\mathbf{Up}(u, \alpha) = \emptyset$. Hence, $(\psi(u))(\alpha) = \top$, from which the result follows.

To show that ψ is a solution of least shape, we need to show that for any other solution $h : S \rightarrow T_\Sigma$ of G , $\mathcal{D}(\psi(s)) \subseteq \mathcal{D}(h(s))$ for any s . This follows directly from Lemma 9. \square

Consider for the last time the running example $E = (\lambda x.x)y$. We get

$$\begin{aligned} \psi(x) = \psi(\llbracket x \rrbracket) = \psi(y) = \psi(\llbracket y \rrbracket) = \psi(\llbracket (\lambda x.x)y \rrbracket) &= \top \\ \psi(\llbracket \lambda x.x \rrbracket) &= \top \rightarrow \top . \end{aligned}$$

6 An Algorithm

We have shown that the type inference problem for PTB can be reduced to solving constraint graphs. Using the characterization of Theorem 17, we get a straightforward type inference algorithm.

Theorem 18 *One can decide in time $O(n^3)$ whether a constraint graph of size n has a finite solution.*

Proof. By Theorem 17, there exists a finite solution if and only if the constructed solution ψ has the property that $\psi(s)$ is finite for all $s \in \mathcal{D}(\psi)$. To determine this, we need only check whether any $\mathcal{L}(s)$ contains an infinite path. We first form the constraint graph, then close it; this gives a graph with n vertices and $O(n^2)$ edges. This can be done in time $O(n^3)$. We then form the automaton \mathcal{M} , which has n^2 states but only $O(n^3)$ transitions, at most $O(n)$ from each state. We then check for a cycle with at least one non- ϵ transition reachable from some (s, s) . This can be done in linear time in the size of the graph using depth-first search. The entire algorithm requires time $O(n^3)$. \square

A λ -term of size n yields a constraint graph with $O(n)$ nodes and $O(n)$ edges. Allowing types to be represented succinctly by the automata \mathcal{M}_s , we get our main result.

Corollary 19 *The type inference problem for PTB is solvable in $O(n^3)$ time.*

Recursive types are just regular trees [1]. The least-shape solution we have constructed, although possibly infinite, is a regular tree for every node in the constraint graph. Thus, we get the following result.

Corollary 20 For PTB_μ , we can infer a least-shape type annotation in $O(n^3)$ time.

It remains open how to extend this result to $\text{PTB}_\mu \cup \{\text{Int}\}$.

7 Example

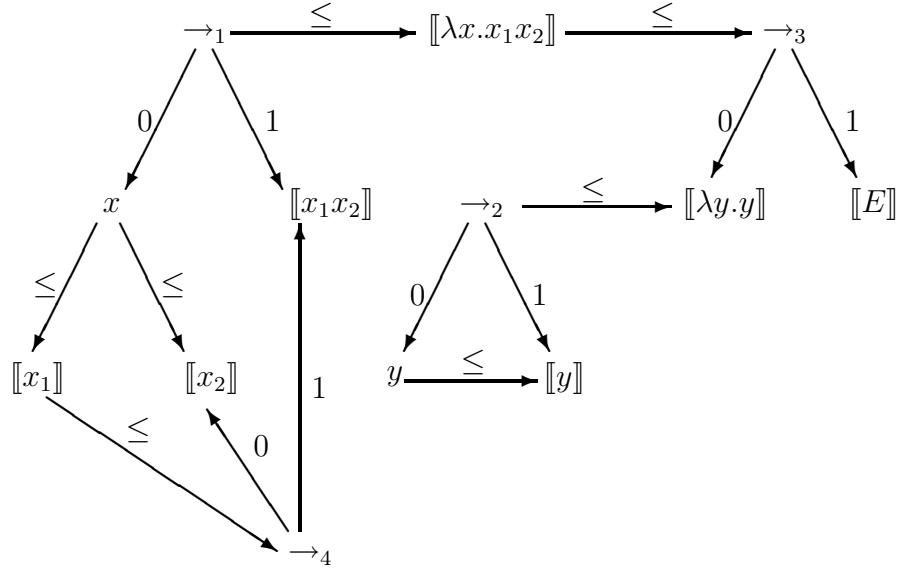
We now give an example of how our type inference algorithm works. Consider the λ -term $E = (\lambda x.xx)(\lambda y.y)$ which was also treated in [6]. We give each of the two occurrences of x a label so that the λ -term reads $(\lambda x.x_1x_2)(\lambda y.y)$. The constraint system $T(E)$ looks as follows:

$$\begin{array}{ll}
\text{From } \lambda x.x_1x_2 & x \rightarrow \llbracket x_1x_2 \rrbracket \leq \llbracket \lambda x.x_1x_2 \rrbracket \\
\text{From } \lambda y.y & y \rightarrow \llbracket y \rrbracket \leq \llbracket \lambda y.y \rrbracket \\
\text{From } E & \llbracket \lambda x.x_1x_2 \rrbracket \leq \llbracket \lambda y.y \rrbracket \rightarrow \llbracket E \rrbracket \\
\text{From } x_1x_2 & \llbracket x_1 \rrbracket \leq \llbracket x_2 \rrbracket \rightarrow \llbracket x_1x_2 \rrbracket \\
\text{From } x_1 & x \leq \llbracket x_1 \rrbracket \\
\text{From } x_2 & x \leq \llbracket x_2 \rrbracket \\
\text{From } y & y \leq \llbracket y \rrbracket
\end{array}$$

To the left of the constraints, we have indicated from where they arise. We will use the following abbreviations:

The symbol:	abbreviates:
\rightarrow_1	$x \rightarrow \llbracket x_1x_2 \rrbracket$
\rightarrow_2	$y \rightarrow \llbracket y \rrbracket$
\rightarrow_3	$\llbracket \lambda y.y \rrbracket \rightarrow \llbracket E \rrbracket$
\rightarrow_4	$\llbracket x_2 \rrbracket \rightarrow \llbracket x_1x_2 \rrbracket$

The constraint graph derived from $T(E)$ looks as follows:

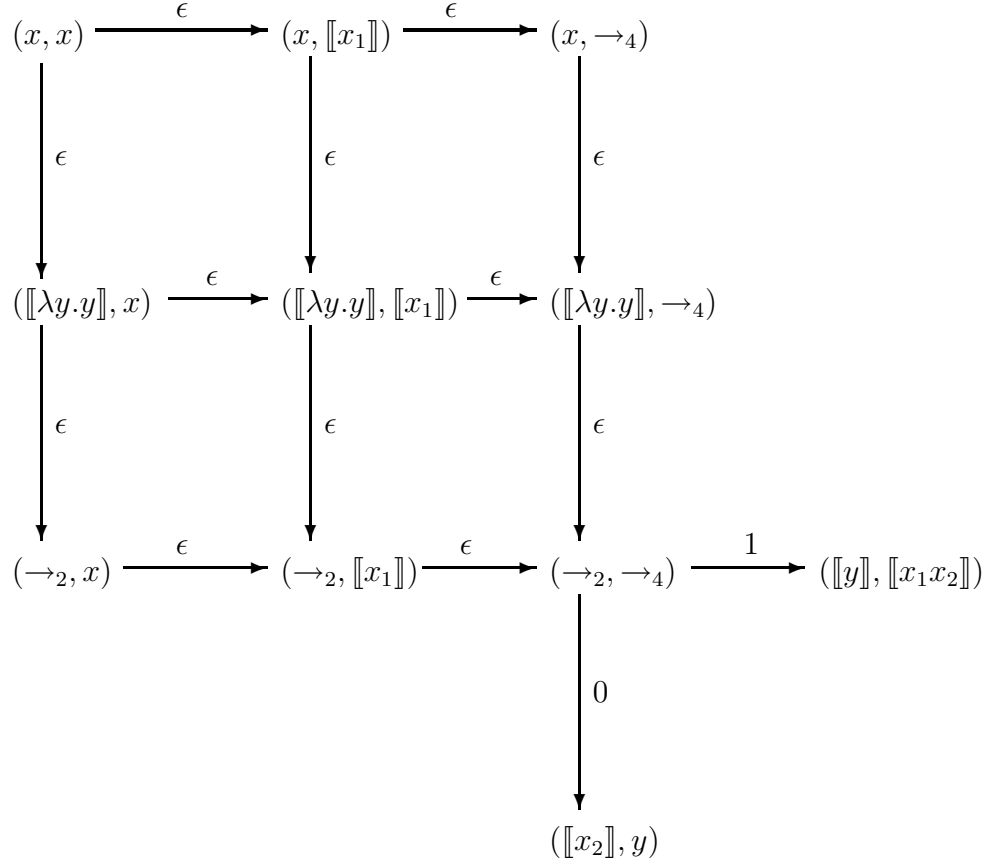


To close the constraint graph, it is sufficient to first fill in the edges

$$\begin{aligned}
 [[x_1x_2]] &\stackrel{\leq}{\rightarrow} [[E]] \\
 [[\lambda y.y]] &\stackrel{\leq}{\rightarrow} x \\
 [[y]] &\stackrel{\leq}{\rightarrow} [[x_1x_2]] \\
 [[x_2]] &\stackrel{\leq}{\rightarrow} y ,
 \end{aligned}$$

and make \leq reflexive and transitive.

The constraint graph has 13 nodes, so the automaton \mathcal{M} has $13^2 = 169$ states. The following picture shows some of the states reachable from the state (x, x) .



It turns out that

$$\begin{aligned}
\mathcal{L}(x) &= \{\epsilon, 0, 1\} \\
\mathcal{L}(y) &= \{\epsilon\} \\
\text{Up}(x, \epsilon) &= \{\rightarrow_4\} \\
\text{Up}(x, 0) &= \emptyset \\
\text{Up}(x, 1) &= \emptyset \\
\text{Up}(y, \epsilon) &= \emptyset .
\end{aligned}$$

(Recall from Definition 13 that $\text{Up}(s, \alpha) \subseteq \mathcal{N}$ for all s, α , and recall that in this example we have $\mathcal{N} = \{\rightarrow_1, \rightarrow_2, \rightarrow_3, \rightarrow_4\}$.) We then get

$$\begin{aligned}
\psi(x) &= \top \rightarrow \top \\
\psi(y) &= \top ,
\end{aligned}$$

and it turns out that E is typable.

For comparison, we can apply the algorithm of [6] for type inference for PTB extended with recursive types to E . The result is that both x and y get annotated by the infinite type $\mu\alpha.\alpha \rightarrow \alpha$ [6].

References

- [1] Roberto M. Amadio and Luca Cardelli. Subtyping recursive types. *ACM Transactions on Programming Languages and Systems*, 15(4):575–631, 1993. Also in Proc. POPL’91.
- [2] Dexter Kozen, Jens Palsberg, and Michael I. Schwartzbach. Efficient inference of partial types. *Journal of Computer and System Sciences*, 49(2):306–324, 1994. Preliminary version in Proc. FOCS’92, 33rd IEEE Symposium on Foundations of Computer Science, pages 363–371, Pittsburgh, Pennsylvania, October 1992.
- [3] Dexter Kozen, Jens Palsberg, and Michael I. Schwartzbach. Efficient recursive subtyping. *Mathematical Structures in Computer Science*, 5(1):113–125, 1995. Preliminary version in Proc. POPL’93, Twentieth Annual SIGPLAN–SIGACT Symposium on Principles of Programming Languages, pages 419–428, Charleston, South Carolina, January 1993.
- [4] John Mitchell. Coercion and type inference. In *Eleventh Symposium on Principles of Programming Languages*, pages 175–185, 1984.
- [5] John C. Mitchell. Type inference with simple subtypes. *Journal of Functional Programming*, 1:245–285, 1991.
- [6] Jens Palsberg and Patrick M. O’Keefe. A type system equivalent to flow analysis. *ACM Transactions on Programming Languages and Systems*, 17(4):576–599, July 1995. Preliminary version in Proc. POPL’95, 22nd Annual SIGPLAN–SIGACT Symposium on Principles of Programming Languages, pages 367–378, San Francisco, California, January 1995.
- [7] Jens Palsberg and Michael I. Schwartzbach. Safety analysis versus type inference for partial types. *Information Processing Letters*, 43:175–180, 1992.

- [8] Satish Thatte. Type inference with partial types. In *Proc. International Colloquium on Automata, Languages, and Programming 1988*, pages 615–629. Springer-Verlag (LNCS 317), 1988.
- [9] Mitchell Wand. Type inference for record concatenation and multiple inheritance. *Information and Computation*, 93(1):1–15, 1991.
- [10] Mitchell Wand, Patrick M. O’Keefe, and Jens Palsberg. Strong normalization with non-structural subtyping. *Mathematical Structures in Computer Science*, 5(3):419–430, 1995.