

Binding-Time Analysis for Standard ML

LARS BIRKEDAL*

*School of Computer Science, Carnegie Mellon University
5000 Forbes Avenue, Pittsburgh, PA 15213, USA*

birkedal@cs.cmu.edu

MORTEN WELINDER

*DIKU, Department of Computer Science, University of Copenhagen
DK-2100 Copenhagen Ø, Denmark*

terra@diku.dk

Abstract. We present an efficient base algorithm for binding-time analysis based on constraint solving and the union-find algorithm. In practice it has been used to handle all of Standard ML except modules and we show the principles of how constraints can be used for binding-time analysis of Standard ML; in particular we show how to binding-time analyse nested pattern matching. To the best of our knowledge no previous binding-time analysis has treated nested pattern matching.

Keywords: binding-time analysis, partial evaluation, Standard ML.

1. Introduction

There are two common ways of doing binding-time analysis: by fixed-point iteration and by constraint solving. Analyses based on the former method tend to be stronger as proven by Palsberg and Schwartzbach [11], but have so far been rather slow — cubic complexity for higher-order languages — and thus impractical for large programs. Analyses based on the latter method are as fast as we can expect — essentially linear-time complexity — and have been treated by Henglein [7], for the untyped lambda calculus, by Andersen [1], for a subset of C, by Bondorf and Jørgensen [4], for a subset of Scheme, and by the present authors [3], for a subset of Standard ML.

Except for the latter, constraint-based binding-time analyses have either been for dynamically typed languages or have not used the available type information. We show that the use of type information is advantageous because it leads to a simple algorithm. This algorithm has been proven correct and very efficient, and it has been implemented.

We have used our algorithm to handle all of Core Standard ML as defined in [10] which clearly demonstrates that usefulness has not been sacrificed for simplicity.

1.1. Constraint-Based Binding-Time Analysis

Figure 1 outlines the stages of constraint-based binding-time analysis. First, the

* The present work was conducted while the author was at DIKU.

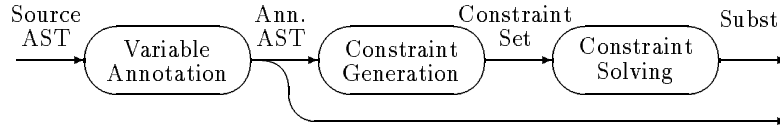


Figure 1: The structure of a BTA based on constraints.

abstract syntax tree of the source program is annotated with one or more variables for every node in the tree. These variables are placeholders for the binding times that will be inferred and for certain control information to be described in Section 3.

In the next phase a set of constraints on the variables is generated from the annotated syntax tree. The set of constraints is called a constraint system and the constraints describe the necessary relations between the binding-time types that the variables will take.

We present our constraint systems in Section 2 and we shall see that our constraint systems are more general than those used by Henglein [7]; in particular, we allow more expressive logical side conditions in the constraint systems.

In Section 3 we explain how the constraints can be used in a binding-time analysis for Standard ML by describing which constraints are generated for typical constructs in Standard ML.

After the constraint generation, the set of constraints is solved by assigning a binding-time type to every variable. Constraint solving is done by rewriting constraint systems until a so-called normal form is reached. From the normal form it is easy to obtain the solution mapping the binding-time variables to binding-time types. In Section 4 we define the normalization process and how to extract a solution from a normal form and we prove the correctness of the constraint-solving algorithm. Moreover, we describe and solve a problem with Henglein’s constraint-solving algorithm [7].

A naive implementation of the constraint-solving algorithm would be extremely slow. In Section 5 we therefore show how it can be implemented efficiently using union-find techniques in two levels.

The result of the binding-time analysis is the solution mapping the binding-time variables to binding-time types together with the annotated syntax tree. Often, these two parts are combined resulting in a so-called two-level syntax tree, but for simplicity we will not do that here.

2. Constraint Systems

Definition. The binding-time types are generated by the following grammar:

$$\tau ::= \mathbf{D} \mid \mathbf{S} \mid [\tau_1, \dots, \tau_n]$$

where $n \geq 0$.

These will be used to describe respectively dynamic entities, static first-order values, and structured values. (The square brackets do *not* imply any connection with lists.) We assume the existence of an infinite set of variables called binding-time variables, and introduce the convention that α 's range over binding-time variables, that β 's range over binding-time variables or the binding-time type \mathbf{D} , and that τ 's range over binding-time types.

Definition. A constraint system is a finite set of formal constraints, each of which must take one of the forms

$$\beta_1 = \beta_2 \quad (\beta_1, \dots, \beta_n) \triangleright \beta \quad [\beta_1, \dots, \beta_n] \leq \beta \quad \beta_1 \rightsquigarrow \beta_2$$

where $n \geq 0$. These constraints are respectively called “equality”, “dependency”, “structure”, and “lift” constraints.

Note that the notation here is just syntax and β 's. There is no implied connection between say $[\beta_1, \dots, \beta_n]$ and structured binding-time types; the interpretation is given below. Constraints of the form $(\beta_1) \triangleright \beta$ are abbreviated as $\beta_1 \triangleright \beta$. The binding-time type \mathbf{S} is not allowed in any constraint system.

To the above formal constraints we attach semantics¹. Each constraint is either satisfied or not, so constraints are predicates in the binding-time variables. The semantics is relative to a substitution, S , that maps binding-time variables to binding-time types, and \mathbf{D} to \mathbf{D} .

Definition. Let τ_i stand for $S(\beta_i)$ and so on. The constraints have the following semantics with respect to S :

- The equality constraint $\beta_1 = \beta_2$ is satisfied if and only if τ_1 and τ_2 consist of the same string of symbols.
- The dependency constraint $(\beta_1, \dots, \beta_n) \triangleright \beta$ is satisfied if $\tau = \mathbf{D}$ or if $\tau_i \neq \mathbf{D}$ for some i . The informal meaning is that “if τ_1, τ_2, \dots , and τ_n all are \mathbf{D} , then τ must also be \mathbf{D} .”
- The structure constraint $[\beta_1, \dots, \beta_n] \leq \beta$ is satisfied either by equality or if τ and all τ_i 's are \mathbf{D} . In words, structured values that are dynamic have dynamic components.
- The lift constraint $\beta \rightsquigarrow \beta'$ is satisfied either by equality, or if $\tau = \mathbf{S}$ and $\tau' = \mathbf{D}$. This is used to express that static first-order values are allowed in a dynamic context.

Example: All the following five constraints are satisfied under the substitution $\{\beta_1 \mapsto \mathbf{S}, \beta_2 \mapsto \mathbf{D}, \beta_3 \mapsto [\mathbf{S}, \mathbf{D}], \mathbf{D} \mapsto \mathbf{D}\}$:

$$\beta_1 \rightsquigarrow \mathbf{D} \quad \beta_1 \rightsquigarrow \beta_2 \quad [\beta_1, \beta_2] \leq \beta_3 \quad (\beta_2, \beta_3) \triangleright \beta_1 \quad (\mathbf{D}, \beta_2) \triangleright \beta_2$$

whereas none of the following are

$$\beta_2 \rightsquigarrow \beta_1 \quad (\mathbf{D}, \beta_2) \triangleright \beta_1 \quad [\beta_2, \beta_1] \leq \beta_3 \quad [\mathbf{S}, \mathbf{D}, \mathbf{D}] \leq \beta_2 \quad \beta_1 = \beta_2$$

Note that the order of components is important in structure constraints. \square

2.1. Well-Typedness

To facilitate efficient and simple solving of constraint systems we impose a well-typedness condition on constraint systems. The intention is that constraint systems generated from well-typed programs should be well-typed.

Definition. Let C be a constraint system and let $R(C)$ be the subset of C consisting of all constraints not having the form $(\dots) \triangleright \beta$. Let $R'(C)$ be $R(C)$ with all occurrences of \mathbf{D} replaced by different fresh binding-time variables.

We define C to be *well-typed* if $R'(C)$ can be solved equationally (i.e., if each structure and lift constraint can be replaced by a corresponding equality constraint, and a most general unifier, $U(C)$, exists for the resulting equational system). The solution may be circular to handle recursive data types.

The exact purpose of this (admittedly rather obscure) definition is to make sure that whenever we have two structure constraints with the same right-hand side then they have the same number of components on the left-hand side.

Example: Consider the constraint set $C = \{[\beta_1, \beta_2] \leq \beta_3, [\beta_4, \beta_5, \beta_6] \leq \beta_3\}$. Here $R'(C)$ cannot be solved by equality and thus C is not well-typed. \square

2.2. Solutions

A constraint system is a kind of inequality system to be solved by replacing binding-time variables with binding-time types.

Definition. Let C be a well-typed constraint system and let S be a substitution mapping C 's binding-time variables into binding-time types. Let S' be the extension of S to all binding-time variables by mapping variables not occurring in C to themselves and let S'' be the component-wise extension of S' to constraints. We say that S is a *solution* to C if for every $c \in C$ it holds that $S''(c)$ is satisfied.

It is easy to see that any constraint system has a solution, $S_{\mathbf{D}}$, which maps all variables present to the constant \mathbf{D} . We need more interesting solutions so we introduce the following partial ordering on solutions.

Definition. We define a partial order, \preceq , on binding-time types by $\tau \preceq \mathbf{D}$, $\tau \preceq \tau$, and $[\tau_1, \dots, \tau_n] \preceq [\tau'_1, \dots, \tau'_n]$ if $\forall i : \tau_i \preceq \tau'_i$. We extend this to variables by defining $\alpha \preceq \mathbf{D}$ and $\alpha \preceq \alpha$. We finally extend it to solutions by defining $S \preceq S'$ if $\forall \alpha : S(\alpha) \preceq S'(\alpha)$, and $S \prec S'$ if $S \preceq S' \wedge S \neq S'$. A solution S of a constraint system C is said to be *minimal* if there is no solution S' of C that satisfies $S' \prec S$.

Minimality for a solution means that it maps as few variables as possible to \mathbf{D} .

2.3. Logical Conditions

If we model truth by \mathbf{D} and falsehood by any other binding-time type we see that it is possible to express logical conditions in the constraint systems. These logical conditions can use both conjunctions (modelled by having more than one component on the left-hand side of dependency constraints) and disjunctions (modelled by multiple dependency constraints with the same right-hand side). Negation cannot be modelled (a corollary of Theorem 6) so the logic is in no sense complete. In the following section we show how such logical conditions can be used.

3. Binding-Time Analysis for Core SML

We have used the constraint systems in a binding-time analysis for programs written in Standard ML (without modules) [3]. In this section we describe, by way of example, the main ideas of how the constraints are to be used in such a binding-time analysis. In particular we argue that the small and fixed set of formal constraints we have introduced is expressive enough to make a non-trivial binding-time analysis of a full-scale, realistic, statically typed, higher-order language like Standard ML.

For Standard ML three kinds of structured values are considered: functions, tuples, and constructed values. In all three cases structured binding-time types are used and well-typedness of the source program implies that *binding-time variables associated with different kinds of structured values are never mixed*. This includes the situation of tuples having different lengths or constructed values from different Standard ML types. To emphasize: from a well-typed Standard ML program a well-typed constraint system in the sense defined above is generated. The reader is referred to [3] for details.

In the following section we describe how the constraints can be used to express binding-time requirements of typical constructs in Standard ML.

3.1. Function Application

Consider an application expression $e \equiv e' e''$ of one expression e' to another expression e'' . In our setting, functional values such as e' either have binding-time type $[\tau_1, \tau_2]$ (which is our notation for Henglein's $\tau_1 \rightarrow \tau_2$) or type \mathbf{D} . The structured binding-time type is used if the function can be applied during specialization whereas \mathbf{D} is used if the function is unknown during specialization and thus cannot be applied. If the binding time of the functional expression e' is \mathbf{D} then the argument expression must be residualized (as the application cannot be performed during specialization) and, moreover, the binding time of the whole expression must also be dynamic. On the other hand, if the application can be performed during specialization then the binding-time types of e'' and e must agree with the binding-time types of τ_1 and τ_2 respectively. These requirements can be specified by the constraint $[\beta'', \beta] \leq \beta'$ where the binding-time variable β holds the binding-time

type of the application expression e , β' the binding-time type of the functional expression e' , and β'' the binding-time type of the argument expression e'' . For every application expression, such a constraint is generated.

3.2. Tuples and Records

Standard ML records can be treated as tuples by sorting the record fields by labels, so we will restrict ourselves to tuples. Consider the tuple expression (e_1, \dots, e_n) and assume that e_i has binding-time type τ_i . We will deal with partially static tuples to the extent that we will allow the tuple to have binding-time value $[\tau_1, \dots, \tau_n]$ no matter what the τ_i 's are, or to have binding-time type \mathbf{D} if and only if $\tau_i = \mathbf{D}$ for all i . In words, we will only build code for a tuple when all components are code. These requirements are specified by the constraint $[\beta_1, \dots, \beta_n] \leq \beta$ where β_i is a binding-time variable holding the binding-time type for e_i and β holds the binding-time type for the entire expression.

3.3. User-Defined Algebraic Data Types

Now consider the case of Standard ML user-defined algebraic data types. We will allow constructors of a data type to be either known during specialization or not; the arguments of constructors may have any binding-time type. Thus we allow partially static structures where, for instance, the spine of a list may be known and the elements may be dynamic. As an example we consider the following Standard ML data type of binary trees.

```
datatype T = L of int | N of int * T * T;
```

For this data type we use six binding-time variables: $\beta_{\mathbf{T}}$ is used for the whole data type; it is \mathbf{D} if the constructors are not all known during specialization; $\beta_{\mathbf{L}}$ and $\beta_{\mathbf{N}}$ are used for the \mathbf{L} and \mathbf{N} constructors and express the constructor arguments' binding-time types, and $\beta_1, \beta_2, \beta_3$ are used for the binding-times of the components of the argument of the \mathbf{N} constructor. The following constraints are used to express the requirements among the binding-time variables:

$$[\beta_{\mathbf{L}}, \beta_{\mathbf{N}}] \leq \beta_{\mathbf{T}} \quad [\beta_1, \beta_2, \beta_3] \leq \beta_{\mathbf{N}} \quad \beta_2 = \beta_{\mathbf{T}} \quad \beta_3 = \beta_{\mathbf{T}}$$

Notice how cycles in the constraint system are used to express recursiveness of data types. The binding-time types of the data type depend, of course, on how the constructors are used in the analyzed program. Constructors occur in expressions and patterns in Standard ML. For every occurrence of the constructor \mathbf{N} as an expression e we generate the following two constraints

$$[\beta_{\mathbf{N}}, \beta_{\mathbf{T}}] \leq \beta_e \quad \beta_{\mathbf{T}} \triangleright \beta_e$$

Intuitively, the first constraint expresses that the constructor \mathbf{N} is in essence a function, i.e., the expression \mathbf{N} takes something of binding-time type $\beta_{\mathbf{N}}$ as argument

and returns something of binding-time type $\beta_{\mathbf{T}}$. The second constraint expresses that if the constructors of the data type are not known during specialization, then the constructor application cannot be performed during specialization and thus the functional expression e must be residualized. For every occurrence of the constructor \mathbf{N} in a pattern p of the form $\mathbf{N}x$ we generate the following constraints

$$\beta_{\mathbf{T}} = \beta_p \quad \beta_{\mathbf{N}} = \beta_x$$

where β_p is used to hold the binding time of the pattern (corresponds to the binding time of the expression the pattern will be matched against) and β_x holds the binding time of the variable x .

We have now seen how we can use the *fixed*, small set of formal constraints to express binding-time properties of *user-defined* data types in well-typed programs. In fact, the constraints can also be used to express other more advanced aspects of binding-time analysis as we shall see now.

3.4. Pattern Matching

3.4.1. Choice Between Branches

Pattern matching in Standard ML is combined with lambda abstraction. Thus, given a lambda abstraction with several branches it is not trivial statically to analyse whether one during specialization can determine which of the branches to choose. This problem is part of the analysis of where to insert memoization points in [3] — if one cannot choose the right branch then a memoization point is to be inserted to avoid infinite unfolding. We now consider this problem.

For each lambda abstraction we will use a binding-time variable, β , as a flag describing whether we can choose a single branch during specialization. If β becomes \mathbf{D} it means that we cannot choose the right branch. Let us assume that the branches are always exhaustive; this can always be ensured by preprocessing. If in one of these branches there is a constant and the corresponding binding-time type is \mathbf{D} , then we cannot determine which branch to choose and the flag β must be set to \mathbf{D} . Consider the following example where we have three branches:

```

fn (1, x)      => ...
  | (2, (1,3)) => ...
  | _          => ...

```

Assume that the binding time for the patterns is $[\mathbf{S}, [\mathbf{S}, \mathbf{D}]]$. Then we cannot determine which branch to choose as the binding time corresponding to the constant $\mathbf{3}$ is \mathbf{D} .

The idea is to decompose the binding time for the patterns according to the patterns and generate $(\cdot) \triangleright \beta$ constraints for all constants. As can be seen from the example, it is not obvious how to decompose the binding time for the patterns — the second pattern would intuitively yield a more fine-grained decomposition

than the first. The point is, however, that it is not necessary to find *one* most fine-grained decomposition of the binding time of the patterns. We can instead consider each pattern in turn, decomposing the binding time each time. It is not only for constants that $(\cdot) \triangleright \beta$ constraints should be generated. If the pattern is a constructed pattern then such a constraint should also be generated with the binding-time variable for the constructed pattern substituted for the dot; because if we do not know the constructor, we cannot determine the branch (constructors correspond here to ordinary constants).

Assume that β_p holds the binding time of the pattern in our example under consideration. Then the following constraints will be generated:

$$\begin{aligned} [\beta_1, \beta_2] &\leq \beta_p & \beta_1 &\triangleright \beta \\ [\beta_3, \beta_4] &\leq \beta_p & \beta_3 &\triangleright \beta \\ [\beta_5, \beta_6] &\leq \beta_4 & \beta_5 &\triangleright \beta & \beta_6 &\triangleright \beta \end{aligned}$$

Notice how \leq constraints are used to decompose the binding time of the pattern — this is safe due to well-typedness. The binding time for the whole lambda abstraction will be $[\beta_p, \beta_e]$ for some β_e . If one compares with the constraints generated for a function application expression (see above), one sees that if the lambda abstraction is applied to something with binding-time type $[\mathbf{S}, [\mathbf{S}, \mathbf{D}]]$, then the constraints generated do express that we cannot choose between the branches. In this case, β will be \mathbf{D} after solving the constraints.

3.4.2. Nested Pattern Matching

We have now seen how we can generate constraints expressing whether we can choose between several branches. We now consider the problem of statically analyzing which parts of a nested pattern matching can be performed during specialization and which parts must be residualized. To the best of our knowledge, no other binding-time analysis have dealt with nested pattern matching — other binding-time analyses have only dealt with non-nested pattern matching, e.g. [9]. Of course, nested pattern matching can be compiled into non-nested pattern matching, but we conjecture that valuable information may be lost in that process. As an example, consider the following lambda abstraction.

```
fn (1, 1, u) => ...
  | (1, y, z) => ...
  | _       => ...
```

In case we *can* choose between the branches during specialization, i.e., if the first two components of the patterns are both static, then the whole pattern matching can be performed during specialization, no matter what the binding time of the third component is. If on the other hand we *cannot* choose the right branch during specialization then things get complicated. Thus assume that the pattern has binding-time type $[\mathbf{D}, \mathbf{S}, \mathbf{D}]$. Now

1. we may perform the matching on the second component during specialization (since this component is static), and
2. no matching, beta-reduction, should be performed on the third component (since a memoization point is to be inserted as mentioned above).

In other words, the pattern matching of a component should be residualized only if we cannot choose between the branches *and* the pattern component is dynamic. The generalized dependency constraint $(\beta_1, \beta_2) \triangleright \beta$ is used to express this conjunction. Note that it is not expressible in Henglein’s system [7]. For the above example, the following constraints are generated, assuming that β is used as a flag saying whether we can choose between the branches (see above), and $\beta'_1, \beta'_2, \beta'_3$ are used for each component respectively to say whether the matching should be residualized, and β_p is used to hold the binding time for the whole pattern.

$$[\beta_1, \beta_2, \beta_3] \leq \beta_p \quad [\beta, \beta_1] \triangleright \beta'_1 \quad [\beta, \beta_2] \triangleright \beta'_2 \quad [\beta, \beta_3] \triangleright \beta'_3$$

The issue of how the binding-time information is used to actually stage the pattern matching in partial evaluation is beyond the scope of this paper; the interested reader is referred to [3].

4. Normalization

This section describes the process of normalizing a set of constraints. The purpose is to reduce the set to a form that can be solved essentially by “treating the structure and lift constraints as equalities.” For this to succeed, equality-conflicting constraints must be rewritten or eliminated. The normalization process also eliminates all non-trivial equality constraints by substitution.

All reductions that we will make on constraint systems preserve solutions. This will be proved in Theorem 5.

Figure 2 outlines how normalization works. The initial constraint system C_0 undergoes a finite number of rewritings spawning off a substitution for each step. The final constraint system C_n is in normal form and can easily be solved in a best-possible way. (“Best-possible” will be defined later.) Finally, the solution to C_n is combined with all the substitutions from the rewritings to form a best-possible solution to C_0 .

4.1. Variable Equivalence

The constraint-solving algorithm presented by Henglein in [7] handles constraint systems similar to those presented here except that the only structures considered there are functions. Our algorithm solves some problems with the algorithm in [7]. Using our notation, a constraint system like

$$\{[\beta_1, \beta_2] \leq \beta_3, \quad [\mathbf{D}, \mathbf{D}] \leq \beta_4, \quad \beta_5 \rightsquigarrow \beta_3, \quad \beta_5 \rightsquigarrow \beta_4, \quad \beta_1 \triangleright \beta_3\}$$

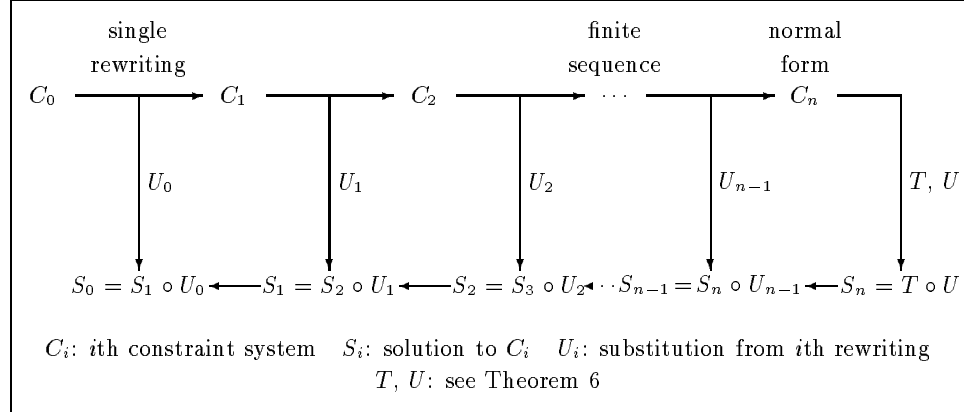


Figure 2: Solving a constraint system by normalization.

is a normal form constraint system in the sense defined in [7]. The first four constraints in the system can be solved by equality but contrary to the claim in Henglein’s Proof 4 some variables (β_1, β_2) are mapped to \mathbf{D} and the fifth constraint is not satisfied by the “solution.” The same problem can be found in [8], Section 8.7, [1], Chapter 5.2, and [2], Section 5.4. A binding-time analysis based on an algorithm with this problem can yield incorrect results leading to a crash of the partial evaluator.

The problem is that the rewritings performed on the constraint systems in [7] only propagate the information that β_3 occurs on the right-hand side of a \leq_f -constraint through lift constraints *from left to right*, so the analysis does not discover that β_3 and β_4 are related. Unfortunately a purely local rewriting rule propagating information the other way through lift constraints would not be solution preserving.

The problem can be overcome by introducing the notion of variable equivalence, which can be conceived of as weakened equality. The concept of equivalence is global in the sense that it takes the entire constraint system into account.

Definition. For a constraint system C we define \approx to be the smallest equivalence relation on the set of binding-time variables and \mathbf{D} satisfying

$$\beta_1 \approx \beta_2 \quad \text{if} \quad \beta_1 = \beta_2 \in C \text{ or } \beta_1 \rightsquigarrow \beta_2 \in C$$

Informally, this means that two binding-time variables are equivalent if they are connected by a chain of lift and equality constraints. In the example above β_3 and β_4 are equivalent. The symmetry of the equivalence relation ensures that information about structure can be propagated not only left-to-right but also right-to-left through lift constraints. Considering also equality constraints makes Theorem 1 (below) hold, which in turn is important for the efficient implementation.

The introduction of variable equivalence leads to an implementation that is rather different from Henglein’s [7], and also from Jones et al.’s presentation [8]. Bondorf

#	Constraint(s)	Condition	Replacement	Substitution
SU-1	$\alpha = \mathbf{D}$ or $\mathbf{D} = \alpha$	-	-	$\{\alpha \mapsto \mathbf{D}\}$
SU-2	$\alpha_1 = \alpha_2$	-	-	$\{\alpha_2 \mapsto \alpha_1\}$
SI-1	$() \triangleright \beta$	-	$\beta = \mathbf{D}$	I
SI-2	$(\beta_1, \dots, \mathbf{D}, \dots, \beta_n) \triangleright \beta$	-	$(\beta_1, \dots, \beta_n) \triangleright \beta$	I
SI-3	$\mathbf{D} \rightsquigarrow \beta$	-	$\beta = \mathbf{D}$	I
C-1	$\begin{cases} [\beta_1, \dots, \beta_n] \leq \alpha \\ [\beta'_1, \dots, \beta'_n] \leq \alpha' \end{cases}$	$\begin{cases} \alpha \approx \alpha' \\ \alpha \not\approx \mathbf{D} \end{cases}$	$\begin{cases} [\beta_1, \dots, \beta_n] \leq \alpha \\ \beta_i = \beta'_i, \alpha = \alpha' \end{cases}$	I
C-2	$[\beta_1, \dots, \beta_n] \leq \beta$	$\beta \approx \mathbf{D}$	$\beta_i = \mathbf{D}, \beta = \mathbf{D}$	I

Figure 3: Rewriting Rules

and Jørgensen’s binding-time analysis [4] uses a separate analysis called a “flow analysis” to capture a form of variable equivalence.

4.2. Normal Forms

Definition. A normal form constraint system is a constraint system where only the following types of constraints occur

$$\mathbf{D} = \mathbf{D}, \quad \alpha \rightsquigarrow \beta, \quad (\alpha_1, \dots, \alpha_m) \triangleright \beta, \quad [\beta_1, \dots, \beta_n] \leq \alpha$$

where $m \geq 1$, $n \geq 0$, and for which the following conditions hold

1. If $\alpha \approx \alpha'$ then there are no two different \leq constraints with α on the right-hand side of the first and α' on the right-hand side of the second.
2. If $\alpha \approx \mathbf{D}$ then α does not occur on the right-hand side of any \leq constraint.

We shall later see that a normal form constraint system can be solved by equality yielding a minimal solution, and that every well-typed constraint system can be “normalized” in a solution-preserving way. Allowing $\mathbf{D} = \mathbf{D}$ in a normal form may look strange, but it saves a rewriting rule below.

4.3. Rewriting Rules

We will turn constraint systems into normal form by repeatedly applying some rewriting rule until no rule applies. The rewriting rules are shown in Figure 3.

Every rewriting rule takes one or two constraints and replaces them by a number of new constraints. The column “Condition” contains a precondition that must be satisfied before a rule can be used, and the column “Substitution” contains what we will call the resulting substitution. I is the identity substitution. In rule C-1 the two source constraints must be different.

There are three classes of rewriting rules: *Substitution rules* (SU- n), *Simplification rules* (SI- n), and *Combination rules* (C- n).

Definition. Normalization is the process of exhaustive use of the above rewriting rules in any order.

4.4. Properties

We are now ready to prove a number of properties of constraint systems. Some of these are almost trivial and others are not, but they are all needed to understand why the algorithm we will eventually present works.

THEOREM 1 (VARIABLE EQUIVALENCE IS PRESERVED) *Let C be a well-typed constraint system, let C' be the resulting system after use of one of the rewriting rules, and let S be the resulting substitution. If $\beta_1 \approx \beta_2$ with respect to C , then $S(\beta_1) \approx S(\beta_2)$ with respect to C' .*

Proof: Only the substitution rules and rule SI-3 remove any equality or lift constraint from C . In the substitution case the equivalence is moved into the substitution and in rule SI-3 one equivalence-causing constraint is replaced by another. In all other cases the theorem is trivially satisfied since the new relation is larger or equal. ■

We do not simply preserve equivalence by making too many variables equivalent in the rewriting process — this fact is a corollary of Theorem 5. Preservation of variable equivalence is important because it means that we need not recalculate the equivalence relation from scratch after each rewriting. Instead we can make do with taking unions of equivalence classes.

THEOREM 2 (WELL-TYPEDNESS IS PRESERVED) *For every well-typed constraint system, C , any of the rewriting rules yields a well-typed constraint system, C' .*

Proof: We prove this rule by rule:

For rule SU-1 we have a simpler unification problem for $R'(C')$ since different occurrences of β have been changed to different fresh variables. For SU-2 we have $\alpha_1 = \alpha_2 \in R'(C)$. This guarantees the existence of a $U(C')$.

Since all occurrences of D 's are replaced by different variables at unification time it is easy to see that all simplification rules preserve well-typedness.

For rule C-1 we have that $U(\alpha) = U(\alpha')$ since $\alpha \approx \alpha'$ and $\alpha \not\approx \mathbf{D}$. We then have $U(\beta_i) = U(\beta'_i)$ so U is also a unifier for C' . For rule C-2 the unification in C' is simpler since the i th \mathbf{D} matches anything that β_i matched. ■

Again it follows from Theorem 5 that we do not obtain preservation of well-typedness by removing too much information. That well-typedness is preserved means that the theorems below which depend on well-typedness can be used not only on the original constraint system but also on rewritten constraint systems.

THEOREM 3 (TERMINATION) *For any well-typed constraint system normalization terminates.*

Proof: This is true because all rules lexicographically decrease (a, b) , where a is the total number of symbols (being “(”, “]”, variables, “ \mathbf{D} ”, ...) on left-hand sides of non-equality constraints, and b is the number of constraints. ■

THEOREM 4 (SOUNDNESS) *Let C be a well-typed constraint system, and let C' be the constraint system obtained by normalization. C' is a normal form constraint system.*

Proof: First we prove that only the allowed types of constraints occur in C' . Concerning lift constraints: the claim is that \mathbf{D} does not occur on the left-hand side; this holds by rule SI-3. Concerning inequality constraints: the claim is that \mathbf{D} does not occur on the right-hand side; this holds by rule C-2. Concerning dependency constraints: the claim is that \mathbf{D} does not occur on the left-hand side and that the left-hand side is not empty; this holds by rules SI-1 and SI-2. Concerning equality constraint: there are none with variables by rules SU-1 and SU-2.

Second we prove that the claims 1-2 in the definition of normal form holds in C' . (1) holds by well-typedness (only compatible entities meet) and by rule C-1. (2) holds by rule C-2. ■

Theorems 3 and 4 tell us that normalization is indeed an algorithm: it does what it is supposed to in a finite number of steps.

LEMMA 1 *Let C be a well-typed constraint system, let C' be the result of applying one of the rewriting rules on C , and let U be the resulting substitution. Then*

$$\begin{aligned} S \text{ is a solution for } C &\Rightarrow \exists S' : S = S' \circ U \quad \text{and} \quad S' \text{ is a solution for } C' \\ S' \text{ is a solution for } C' &\Rightarrow S \stackrel{\text{def}}{=} S' \circ U \text{ is a solution for } C. \end{aligned}$$

Proof: Observe that the sets of variables used are unchanged except for the substitution rules and the requirement for substitutions to map variables not occurring in a constraint system to themselves is therefore preserved in those cases. In the same cases we further have that U is the identity substitution so that the above claims reduce to showing that solutions coincide. We prove the two claims by examining the rewriting rules in turn.

SU-1: The first claim is satisfied with $S' = V^{-1} \circ S \circ V$, where V maps α to a variable not used elsewhere and passes everything else. S' maps α onto itself and no other variables disappear or arrive. The second claim is satisfied since S must substitute \mathbf{D} for α .

SU-2: The first claim is satisfied as for the preceding rule, but with α_2 instead of α . The second claim is satisfied since S must map α_1 and α_2 to the same thing.

SI- n : These preserve solutions both ways as they are direct application of the semantics of constraints.

C-1: The first claim holds since if $S(\alpha) = \mathbf{D}$ then also $S(\alpha') = \mathbf{D}$ and S clearly solves C' . Otherwise $S(\alpha) = [\tau_1, \dots, \tau_n] = S(\alpha')$ since structured values cannot be lifted and because $\alpha \approx \alpha'$. The second claim holds obviously.

C-2: The first claim holds since $\beta \approx \mathbf{D}$ and because structured values cannot be lifted so we have $S(\beta) = \mathbf{D}$ and thus $S(\beta_i) = \mathbf{D}$. The second claim holds obviously. ■

THEOREM 5 (SOLUTIONS ARE PRESERVED) *Let C be a well-typed constraint system, let C' be the result of applying one of the rewriting rules on C . Then there is a one-to-one correspondence between solutions of C and solutions of C' .*

Proof: For non-substitution rules this was proved in the previous lemma. For the substitution rules we have the correspondence $S' = V^{-1} \circ S \circ V$ and $S = V \circ S' \circ V^{-1}$ using notation from the lemma.² ■

THEOREM 6 *Every normal form constraint system C has a minimal solution.*

Proof: We will prove this constructively by exhibiting a minimal solution. Let C' be the subset of C consisting of constraints having the form $[\beta_1, \dots, \beta_n] \leq \alpha$ or $\alpha \rightsquigarrow \alpha'$ (i.e., leave out dependency constraints and lift constraints with \mathbf{D} on the right-hand side). Solve all constraints in C' by equality. This is possible due to well-typedness and the definition of normal form. Call the resulting substitution U and let T be the substitution that maps all variables in $U(C)$ to the constant value \mathbf{S} . We claim that $V = T \circ U$ solves the entire system C and that V is minimal.

Notice first that neither U nor T maps any variable to \mathbf{D} . Every dependency constraint is therefore satisfied by V . For every $\alpha \rightsquigarrow \mathbf{D}$ in C we have by the second condition of normal form that $V(\alpha) = \mathbf{S}$ so the constraint is satisfied. All other remaining constraints are equality-solved by U and therefore satisfied by V . As variables not in C are mapped to themselves by both U and T we conclude that V is indeed a solution to C .

Now consider some other solution $V' \neq V$. If V' solves all the above mentioned constraints by equality it must map at least one of the remaining variables to something different from \mathbf{S} , and the two solutions are then either incommensurable or we have $V \preceq V'$. If on the other hand V' does not solve all the above mentioned

constraints by equality then it maps some α to \mathbf{D} where V does not, and we have $V' \not\leq V$. We finally conclude that V is minimal. ■

THEOREM 7 (MINIMALITY IS PRESERVED) *Let C be a well-typed constraint system, let C' be the result of applying one of the rewriting rules on C with the resulting substitution U . If S' is a minimal solution to C' then $S \equiv S' \circ U$ is a minimal solution to C .*

Proof: Assume that T is a solution to C and that $T \prec S$. It follows from Lemma 1 that S really is a solution and that we can choose a solution T' to C' such that $T = T' \circ U$.

For an α where $S(\alpha) \neq T(\alpha)$ we now have $T(\alpha) \neq \mathbf{D}$ and $S(\alpha) = \mathbf{D}$ from which it follows that $T'(\alpha) \neq \mathbf{D}$, $U(\alpha) \neq \mathbf{D}$, and $S'(U(\alpha)) = \mathbf{D}$. Thus $T'(\alpha) \preceq S'(U(\alpha))$.

For an α where $S(\alpha) = T(\alpha)$ we have by choice of T' that $S'(U(\alpha)) = T'(U(\alpha))$. If $U(\alpha) = \alpha$ then this means that $T'(\alpha) = S'(\alpha)$. Otherwise, by the definition of a solution, we directly have $T'(\alpha) = S'(\alpha)$.

Since $S' \neq T'$ we have thus proved $T' \prec S'$ which is a contradiction to the minimality of S' . ■

All in all, the previous theorems prove that every well-typed constraint system C has a minimal solution: since C is well-typed it can be normalized yielding a normal form constraint system C' . A minimal solution exists for C' and this gives rise to a minimal solution for C .

5. Efficient Implementation

A naïve implementation of the normalization process described above would be very slow since it would require searching for some constraints every time one of the combination rules should be used. This section presents an efficient implementation that runs in almost linear time.

The key to efficiency is to consider the constraints one by one and to record enough information about already-seen constraints to make it possible to decide which rules can be applied when a new constraint is considered. A concrete Standard ML implementation can be found in [3].

Phase 1 — initialization

- Apply simplification rules SI-1 and SI-2 exhaustively. For every $\beta_1 \rightsquigarrow \beta_2$ constraint add a $\beta_1 \triangleright \beta_2$ constraint. This handles rule SI-3 in a simple way.
- Using union-find techniques an evolving equivalence relation can be constructed on the set of binding-time variables and \mathbf{D} . This relation is used to trace variable substitutions (*not* variable equivalence!) The class in which β occurs is written $\bar{\beta}$ and we associate a set of dependencies, $\text{dps}(\bar{\beta})$, with every class.

- The working list is a list of constraints that we have not dealt with yet. We initialize the working list to all non-dependency constraints. Note that a copy of this list must be kept for phase 3.
- For each constraint $(\beta_1, \dots, \beta_k) \triangleright \beta$ create a pair (k, β) . Insert a reference to this pair in all $\text{dps}(\overline{\beta}_i)$. Since k may be decremented later it is important that all the β_i 's reference the same pair. The meaning of such a pair is that k is decremented every time some β_i turns \mathbf{D} and if k eventually becomes zero then the constraint $\beta = \mathbf{D}$ will be inserted into the working list.
- Construct a term graph for all members of the working list using nodes with labels \mathbf{S} (arity 0), \mathbf{D} (arity 0), $=$ (arity 2), \rightsquigarrow (arity 2), $[]$ (some finite arity), and \leq (arity 2). The binding-time variables are the other leaves of the graph.
- On the set of classes, $\overline{\beta}$, create another evolving equivalence relation for handling variable equivalence. Call the resulting classes $\overline{\overline{\beta}}$. With each class $\overline{\overline{\beta}}$ we associate a field, $\text{memory}(\overline{\overline{\beta}})$, which is to hold at most one unprocessed inequality constraint with β on the right-hand side. Notice by inspection of the second combination rule that it does not matter which $\beta \in \overline{\overline{\beta}}$ or $\overline{\overline{\beta}} \in \overline{\overline{\beta}}$ is used. The memory field of $\overline{\overline{\mathbf{D}}}$ will always be kept empty.

Phase 2 — main loop

While the working list is not empty remove a constraint, c , from it. If there is a lift or an equality constraint in the list, c must be such a constraint. (In other words we handle structure constraints only when no other constraints are available.)

If c is an equality or a lift constraint we now recognize the two sides to be variable equivalent and therefore unify the two $\overline{\overline{\beta}}$ classes if they are not already equivalent. This may lead to application of one of the combination rules since two conflicting memory fields may collide or one of c 's sides may be variable equivalent to \mathbf{D} .

If c is an equality constraint with different sides we also unify the two $\overline{\overline{\beta}}$ classes thus recording the substitution. If both sides are variables then the dps -sets are concatenated. If on the other hand one of the sides is \mathbf{D} we update the dps -sets as described above.

If c is a structure constraint with right-hand side β we have by the choice of c that the variable equivalence relation is up to date. If $\text{memory}(\overline{\overline{\beta}})$ is empty simply record c there. Otherwise apply one of the combination rules.

Phase 3 — postprocessing

- For each $[...] \leq \beta$ or $\cdot \rightsquigarrow \beta$ in the saved copy of the working list solve by equality if $\overline{\overline{\beta}} \neq \overline{\overline{\mathbf{D}}}$. This is done efficiently by changing β to a link to the left-hand side when the sides are different.
- Unify any remaining free variable with \mathbf{S} .

Complexity

The algorithm described above is very efficient and our implementation confirms this. All in all the algorithm can be made to run in time $O(n \cdot \alpha(n, n))$, where n is the total number of components in the constraint system, and α is an inverse of Ackermann’s function (originating from the union-find structures) [7], [3].

6. Related Work

As already noted the algorithm described here is closely based on the one by Henglein in [7]. In particular, the process outlined in Figure 2 comes from that paper. The ideas in that paper have proven very useful as it can be seen from the fact that they play an important rôle in [1], [2], [4], [3]. These works also show how to use constraint systems for doing binding-time analysis.

Launchbury [9] used projections and abstract interpretation to perform binding-time analysis on a small strongly typed first-order language. His methods were later used by de Niel in [6] for a similar language. The methods are well-suited for both polymorphic programs and for polyvariant specialization, but it seems to be difficult to use the methods for higher-order programs. Recently, however, Kei Davis have proposed to use projections for higher-order monomorphic binding-time analysis [5].

The constraint systems used in [4] are significantly different from those in this article because, for historical reasons, a completely different set of binding-time types are used. Also, as noted in Section 4.1, a separate “flow analysis” is used instead of variable equivalence.

7. Conclusion and Future Work

We have presented and proved the correctness of a general framework for doing binding-time analysis of statically typed languages. We have also shown that the method presented can be used for languages with a variety of different data types, such a Standard ML’s functions, records, user-defined algebraic data types, and pattern matching.

We believe that the fixed set of constraint types presented here is strong enough to express most needs for realistic binding-time analyses of statically typed languages. This belief is supported by the fact they are strong enough for (monomorphic) Standard ML with partially static data structures and nested pattern matching.

Future research should be conducted in two directions: the application of constraint-based binding-time analysis to polymorphically typed programs, and polyvariant binding-time analysis by constraint solving.

Acknowledgments

The authors wish to thank the DIKU Topps group for valuable discussions. The idea of solving the normalization problem from Section 4.1 with a (slightly different) equivalence relation as well as the initial conjecture that type information might prove useful in binding-time analyses originate from informal discussions with Fritz Henglein. Further thanks go to Peter Sestoft, Olivier Danvy, Robert Glück, Karoline Malmkjær, and the referees for comments on earlier versions of this article.

Notes

1. Following mathematical tradition we do not distinguish notationally between formal and semantic constraints.
2. For the formally inclined reader we note that the definition of a solution basically is the same in [7] as in this paper; unused variables are required to be mapped to themselves. Strictly speaking this property together with Henglein's rewriting rule 3a invalidates the claim in [7], Theorem 2 and its proof. To see this, one just has to consider the constraint system $\{\beta \triangleright \mathbb{D}\}$ which will be normalized to $\{\}$. The former system has a countably infinite number of solutions (because β can be mapped to any binding-time type), the latter has exactly one solution. The same problem can be found in [8], Section 8.7, [1], Chapter 5, and [2], Figure 7.

References

1. Lars Ole Andersen. C program specialization. Technical Report 92/14, DIKU, Department of Computer Science, University of Copenhagen, May 1992.
2. Lars Ole Andersen. Binding time analysis and the taming of C pointers. In *Proceedings of the ACM SIGPLAN Symposium on Partial Evaluation and Semantics-Based Program Manipulation*, pages 47–58. ACM Press, June 1993.
3. Lars Birkedal and Morten Welinder. Partial evaluation of Standard ML. Technical Report 93/22, DIKU, October 1993.
4. Anders Bondorf and Jesper Jørgensen. Efficient analyses for realistic off-line partial evaluation. *Journal of Functional Programming, special issue on partial evaluation*, 3, July 1993.
5. Kei Davis. Pers from projections for binding-time analysis. In *Lisp and Symbolic Computation, special issue on Partial Evaluation and Semantics-Based Program Manipulation (this volume)*. Klüwer Academic Publishers, 1995.
6. Anne de Niel. *Self-applicable Partial Evaluation of Polymorphically Typed Functional Languages*. PhD thesis, Katholieke Universiteit Leuven, January 1993.
7. Fritz Henglein. Efficient type inference for higher-order binding-time analysis. In John Hughes, editor, *Functional Programming Languages and Computer Architecture, Cambridge, Massachusetts, August 1991. (Lecture Notes in Computer Science, vol. 523)*, pages 448–472. Springer-Verlag, 1991.
8. Neil D. Jones, Carsten K. Gomard, and Peter Sestoft. *Partial Evaluation and Automatic Program Generation*. Prentice-Hall, 1993.
9. John Launchbury. *Projection Factorisations in Partial Evaluation*. Cambridge University Press, 1991.
10. Robin Milner, Mads Tofte, and Robert Harper. *The Definition of Standard ML*. MIT Press, 1990.
11. Jens Palsberg and Michael I. Schwartzbach. Binding-time analysis: Abstract interpretation versus type inference. In *Proc. ICCL'94, Fifth IEEE International Conference on Computer Languages, Toulouse, France, May 1994*.