

# Efficient ML Type Inference Using Ranked Type Variables

George Kuan    David MacQueen

University of Chicago  
{gkuan,dbm}@cs.uchicago.edu

## Abstract

The Damas-Milner type inference algorithm (commonly known as algorithm  $\mathcal{W}$ ) is at the core of all ML type checkers. Although the algorithm is known to have poor worst-case behavior [8], in practice well-engineered type checkers will run in approximately linear time. To achieve this efficiency, implementations need to improve on algorithm  $\mathcal{W}$ 's method of scanning the complete type environment to determine whether a type variable can be generalized at a let binding. Following a suggestion of Damas, most ML type checkers use an alternative method based on ranking unification variables to track their position in the type environment.

Here we formalize two such ranking systems, one based on lambda depth (used in the SML/NJ compiler), and the other based on let depth (used in OCaml, for instance). Each of these systems is formalized both with and without the value restriction, and they are proved correct relative to the classic algorithm  $\mathcal{W}$ . Our formalizations of the various algorithms use simple abstract machines that are similar to those derived from small-step evaluation semantics.

**Categories and Subject Descriptors** D.3.4 [Programming Languages]: Compilers

**General Terms** Languages

**Keywords** Type Inference, Type Checking, Algorithm  $\mathcal{W}$ , Standard ML, Functional Programming, Compilers

## 1. Introduction

ML compilers have adapted Milner's Algorithm  $\mathcal{W}$  [?] to type-check realistically rich languages with reasonable efficiency. But to achieve efficiency and the ability to scale to very large programs, implementations of ML type inference must overcome various engineering challenges. The challenge we focus on here is the basic problem of determining which type variables can be generalized to create polymorphic types at let bindings. In algorithm  $\mathcal{W}$ , the criterion for generalization is that the type variable does not appear free in the "context", represented by a type environment. Since type environments can grow quite large, both in length and the sizes of types they contain, this test can be expensive (quadratic in the size of the program). More efficient algorithms have been developed that use ranking of type variables and a simple rank bound to test for generalizability. The rank associated with a type variable indi-

	Pure	Ref
Classic $\mathcal{W}$	$\mathcal{W}$ (Fig. 4)	$\mathcal{W}_V$ (Fig. 17)
$\lambda$ -Ranking	$\mathcal{W}_\lambda$ (Fig. 7)	$\mathcal{W}_{\lambda V}$ (Fig. 19)
let-Ranking	$\mathcal{W}_L$ (Fig. 14)	$\mathcal{W}_{LV}$ (Fig. 20)

Figure 1. Summary of Type Inference Abstract Machines

cates its "position" in the type environment, where the type environment is viewed as ordered from outer to inner bindings.

There are two main variants of this variable ranking approach. The one originally suggested by Damas [1] is based on depth of lambda bindings, and is used in SML/NJ. An equivalent alternative is based on depth of let bindings [13]. Although these implementation techniques have been used in compilers for a couple of decades, their formalization and proof of correctness has been largely neglected. Our goal is to begin to address that gap by proving that these ranked variable variants are equivalent to the original algorithm  $\mathcal{W}$ . We start by considering type checkers for a pure form of microML (which we will call Pure ML), and then go on to introduce side-effects (ref) and the value restriction [18], calling this language variant Impure ML.

Our approach is to formalize the various type inference algorithms that we will be studying using simple abstract machines such as those derived from small-step operational semantics for evaluation [2]. We then prove equivalences between these abstract machines, whose operation differs mainly at the point where polymorphic generalization occurs for let bindings. Since we have three variants of the type inference algorithm  $\mathcal{W}$  (classic, lambda-ranked, and let-ranked) for two variants of the language (Pure and Impure ML), we have a total of six abstract machines, as shown in Fig. 1.

## 2. An Abstract Machine for Classic Algorithm $\mathcal{W}$

The abstract machines that we use here are simplified versions of the abstract machines used in an earlier paper [4], where we were viewing type checking and inference as a term-rewriting system that incrementally transformed expressions into their types. The abstract machines described below express type inference as a process of nonstandard evaluation of expressions, yielding their types.

The language we will start with is the usual minimal ML, presented in Fig. 3. This language can easily be extended with constant expressions (*e.g.*, numbers, booleans), and with additional basic type constructors (*e.g.*, Int, Bool, List), and we will feel free to do so in examples. Later on we will introduce references, which will force us to complicate type inference by introducing the value restriction.

Types ( $\tau \in \text{TYPE}$ ) are constructed over two forms of type variable. Polymorphic bound type variables ( $\alpha, \beta, \gamma$ ) will occur only in the bodies of polytypes (PTYPE) where they appear in the quantifier prefix – in other words, the types that occur during

$\mathcal{W}(\Gamma, x)$	=	$(id, \mathcal{I}(\Gamma(x)))$	
$\mathcal{W}(\Gamma, \lambda x.e)$	= let	$(\theta, \tau) = \mathcal{W}(\Gamma[x : \xi], e)$	$\xi$ fresh
	in	$(\theta, \theta\xi \rightarrow \tau)$	
$\mathcal{W}(\Gamma, e e')$	= let	$(\theta, \tau) = \mathcal{W}(\Gamma, e)$	
		$(\theta', \tau') = \mathcal{W}(\theta\Gamma, e')$	
		$\theta'' = \mathcal{U}(\theta'\tau, \tau' \rightarrow \xi)$	$\xi$ fresh
	in	$(\theta''\theta'\theta, \theta''\xi)$	
$\mathcal{W}(\Gamma, \text{let } x = e \text{ in } e')$	= let	$(\theta, \tau) = \mathcal{W}(\Gamma, e)$	
		$(\theta', \tau') = \mathcal{W}(\theta\Gamma[x : \sigma], e')$	
	in	$(\theta'\theta, \tau') \quad \sigma = \forall \bar{\alpha}. \{\bar{\alpha}/\mathcal{G}_\Gamma(\tau)\}\tau$	

**Figure 2.** Algorithm  $\mathcal{W}$  from Lee and Yi [5]

type inference will never contain free occurrences of such type variables. The other type variables ( $\xi, \zeta, \phi, \psi$ ) are called *unification variables*, or *univariables* for short. They should be considered part of the internal machinery of type inference and should generally be eliminated from the final type produced.<sup>1</sup>

It should also be noted that although neither types nor polytypes (PTYPE) will contain free occurrences of polymorphic type variables, types and polytypes can both contain occurrences of univariables (which are necessarily free occurrences). We will use the term *monotype* for types  $\tau \in \text{TYPE}$  to distinguish them from polytypes  $\sigma \in \text{PTYPE}$ , even though monotypes may contain free occurrences of univariables.

For reference, Fig. 2 shows a standard version of algorithm  $\mathcal{W}$  for Pure ML, following the modern presentation given in Lee and Yi [5]. This version differs from Milner’s original formulation [9] in that there are explicit generalization and generic instantiation steps that respectively introduce and eliminate polytypes.

Note that  $\mathcal{W}$  returns only monotypes  $\tau$  as results, while both types and polytypes occur in type environments. We see that polytypes are created and then immediately entered into the type environment while processing let bindings. Sometimes, when no univariables are generalizable, these let-bound polytypes will be degenerate in the sense that the quantifier prefix is empty.

New univariables are introduced at three points: (1) when entering a lambda abstraction, where we introduce a new univariable to represent the unknown type of the  $\lambda$ -bound variable, (2) when typing a variable, where we create a *generic instance* of the polytype bound to the variable in the type environment, and (3) in the application case, where we introduce a fresh univariable to represent the return type before unifying. Conversely, univariables can be eliminated (1) by being instantiated (*i.e.* replaced by substitution), or (2) by being generalized at a let binding.

Using standard methods [14], we can transform the classic algorithm of Fig. 2 into the abstract machine in Fig. 4.

The abstract machine states consist of three components, or *registers*. The control, the first register, can be either the next (sub)expression we are type checking or a monotype produced by completing the typing of a subexpression. The type environment, the second register, maps program variables to either monotypes (for  $\lambda$  bindings) or polytypes (for let bindings). The type environment context is an ordered sequence of  $\lambda$  and let-bindings that map expression variables to their types, with the outermost bindings first, so that the  $\lambda$  and let-nesting is reflected in its binding order. The distinction between let and lambda bindings is indicated

<sup>1</sup>A top level expression like  $\lambda x.x$  will be assigned a type  $\xi \rightarrow \xi$  containing a free univariable. We can eliminate this univariable by supplying an implicit top level let declaration context for the expression, typing it as let  $it = \lambda x.x$  where the scope is the “rest of the program”. Then the univariable can be eliminated by generalization. If the value restriction applies this will not always work (*e.g.*, ref nil), and we may have to eliminate the univariables by dummy instantiation with arbitrary ground types, as is done in the SML/NJ type checker.

$e ::=$	$x \mid \lambda x.e \mid e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2$	EXP
$\tau ::=$	$\alpha \mid \xi \mid \tau \rightarrow \tau$	TYPE
$\sigma ::=$	$\forall \bar{\alpha}. \tau$	PTYPE

**Figure 3.** Pure ML mini-language

by whether the type in the binding is a polytype (possibly degenerate) introduced by a let, or a monotype, introduced by a lambda.

The stack, the final register, represents a typing continuation that keeps track of what still needs to be done to complete the type checking.  $\lambda$  and let stack frames tell the machine when to pop off associated  $\lambda$ - and let-bindings from the type environment (*i.e.*, when the type checker is leaving the scope of the  $\lambda$  or let body).  $(\square e)$  and  $(\tau \square)$  frames save relevant context information when type checking the operator and operands of applications. Finally, let  $x$  in  $e$  frames tell the machine to generalize the control type, which must be the type of the definien, and to add the resulting polytype to the type environment when type checking the body  $e$ .

As in the original algorithm  $\mathcal{W}$ , polytypes are only introduced by [let-body], which puts them in the environment. The [var] rule looks up a variable in the type environment and immediately applies  $\mathcal{I}$  to generically instantiate the result, eliminating the polytype and (possibly) introducing fresh univariables. Consequently, polytypes only exist in the typing context and are never found in the control.

The [ $\lambda$ -in] rule pushes a new  $\lambda$ -binding onto the type environment, with a fresh univariable as the type of the bound variable. Simultaneously, a  $\lambda$  frame is pushed onto the stack as a reminder to pop off the associated  $\lambda$  binding when we have finished type checking the body. The [ $\lambda$ -out] rule pops off the  $\lambda$ -binding and constructs the function type to complete the typing of an abstraction. The [let-out] rule plays an analogous role for let expressions, popping the let-binding from the typing context and returning the body type. The remaining rules essentially serve as search rules.

The [app-out] rule sets up the function-argument agreement constraint, using a fresh univariable, and solves it by unification ( $\mathcal{U}$  defined in Fig. 5). In the process, the unification may eliminate some univariables. Unification fails (in a stuck state) if the function and argument cannot be made to agree, in which case a type error is signaled.

It is straightforward to show that the Algorithm  $\mathcal{W}$  abstract machine of Fig. 4 is equivalent to the classical algorithm of Fig. 2, in the sense that for any closed expression  $e$ ,

$$\mathcal{W}(e, \emptyset) = (\tau, \theta) \Leftrightarrow (e, \emptyset, \bullet) \mapsto (\tau, \emptyset, \bullet) \quad (1)$$

and we therefore take the abstract machine as our reference implementation from here on.

### 3. Depth Ranking Techniques

Let us make some key observations about the effect of substitutions on the occurrence of univariables in type environments, which is the critical factor in deciding whether a univariable can be generalized.

First note that for any  $\xi \in \text{ftv}(\Gamma)$  we can define the *rank* of  $\xi$  relative to the environment  $\Gamma$  as the index (counting from 1) of the first (*i.e.*, outermost) binding in  $\Gamma$  containing an occurrence of  $\xi$ . We will denote this rank by  $\mathcal{R}(\xi, \Gamma)$ . If  $\xi \notin \text{ftv}(\Gamma)$ , we define  $\mathcal{R}(\xi, \Gamma) = \infty$ .

So at a let binding, a univariable is generalizable if and only if its rank relative to the current type environment is infinite:

$$\xi \in \mathcal{G}_\Gamma(\tau) \Leftrightarrow \xi \in \text{ftv}(\tau) \text{ and } \mathcal{R}(\xi, \Gamma) = \infty$$

Another useful observation is that at a let binding, if we generalize relative to the current environment  $\Gamma$  to get a polytype  $\sigma$  and add the polytype binding  $[x : \sigma]$  to  $\Gamma$ , then any univariable  $\xi \in \text{ftv}(\sigma)$

$\Gamma ::= \emptyset \mid \Gamma[x : \tau] \mid \Gamma[x : \sigma]$	Types
$f ::= \lambda \mid (\square e) \mid (\tau \square) \mid \text{let } x \text{ in } e \mid \text{let}$	Frames
$k ::= \bullet \mid f :: k$	Stack
$c ::= e \mid \tau$	Control
$s ::= (c, \Gamma, k)$	States
Initial States $s_0 = (e, \emptyset, \bullet)$	
$(x, \Gamma, k) \mapsto (\mathcal{I}(\Gamma(x)), \Gamma, k)$	[var]
$(\lambda x.e, \Gamma, k) \mapsto (e, \Gamma[x : \xi], \lambda :: k)$	[λ-in]
	$\xi \text{ fresh}$
$(\tau, \Gamma[x : \tau'], \lambda :: k) \mapsto (\tau' \rightarrow \tau, \Gamma, k)$	[λ-out]
$(e_1 e_2, \Gamma, k) \mapsto (e_1, \Gamma, (\square e_2) :: k)$	[app-l]
$(\tau, \Gamma, (\square e_2) :: k) \mapsto (e_2, \Gamma, (\tau \square) :: k)$	[app-r]
$(\tau, \Gamma, (\tau' \square) :: k) \mapsto (\theta\xi, \theta\Gamma, \theta k)$	[app-out]
	$\text{where } \theta = \mathcal{U}(\tau', \tau \rightarrow \xi), \xi \text{ fresh}$
$(\text{let } x = e_1 \text{ in } e_2, \Gamma, k) \mapsto (e_1, \Gamma, \text{let } x \text{ in } e_2 :: k)$	[let-def]
$(\tau, \Gamma, \text{let } x \text{ in } e_2 :: k) \mapsto (e_2, \Gamma[x : \sigma], \text{let} :: k)$	[let-body]
	$\text{where } \sigma = \forall \bar{\alpha}. \{\bar{\alpha} / \mathcal{G}_\Gamma(\tau)\} \tau, \bar{\alpha} \text{ fresh}$
$(\tau, \Gamma[x : \sigma], \text{let} :: k) \mapsto (\tau, \Gamma, k)$	[let-out]
$\mathcal{G}_\Gamma(\tau) = \text{ftv}(\tau) - \text{ftv}(\Gamma)$	Generalization
$\mathcal{I}(\forall \bar{\alpha}.\tau) = \{\bar{\xi} / \bar{\alpha}\} \tau$ where $\bar{\xi}$ is fresh	Generic Instantiation
$\mathcal{I}(\tau) = \tau$	

Figure 4. Algorithm  $\mathcal{W}$  Abstract Machine

$\mathcal{U}(\tau, \tau) = \varepsilon$	$\mathcal{U}(\tau, \xi) = \mathcal{U}(\xi, \tau)$
	$\text{where } \tau \text{ is not a univariable}$
$\mathcal{U}(\xi, \tau) = \{\tau / \xi\}$ if $\xi \notin \text{ftv}(\tau)$	
$\mathcal{U}(\tau_1 \rightarrow \tau_2, \tau_3 \rightarrow \tau_4) = \mathcal{U}(\theta(\tau_2), \theta(\tau_4)) \circ \theta$	
	$\text{where } \theta = \mathcal{U}(\tau_1, \tau_3)$

Figure 5. Unification

already appeared in  $\Gamma$ , *i.e.*,  $\xi \in \text{ftv}(\Gamma)$ . This is because if  $\xi$  did not appear in  $\Gamma$ , it would have been generalized and replaced by a polymorphic bound type variable. Thus, at the point where it is added, a let binding introduces no new univariables into the environment. It also means that the first occurrence of a univariable in the environment is always at a  $\lambda$  binding. These properties are obviously preserved if a type substitution is applied to a type environment. So we have the invariant (for Pure ML):

1. *The first occurrence of a univariable in an environment is always at a  $\lambda$  binding.*

Now let us consider the effect of a substitution on the rank of a univariable in an environment. Consider a univariable  $\xi$ , an environment  $\Gamma$ , and a type substitution  $\theta$ . What is the relation between  $\mathcal{R}(\xi, \Gamma)$  and  $\mathcal{R}(\xi, \theta\Gamma)$ ? If  $\xi \in \text{dom}(\theta)$ , then  $\xi$  is eliminated by  $\theta$  (remember that all our substitutions are idempotent because of the occurrence check), so  $\mathcal{R}(\xi, \theta\Gamma) = \infty$ . But in this case, since  $\theta$  is applied uniformly to all components of the abstract machine,  $\xi$  is eliminated everywhere, and its rank is no longer relevant since we will never ask whether  $\xi$  can be generalized.

Next assume that  $\xi \notin \text{dom}(\theta)$ . Then applying  $\theta$  to  $\Gamma$  will not remove any occurrences of  $\xi$ , so the rank of  $\xi$  in  $\Gamma$  will not increase, but it could decrease if  $\xi$  occurs in the range of  $\theta$  and new occurrences of  $\xi$  are added in  $\theta(\Gamma)$  (including the case where  $\xi \notin \text{ftv}(\Gamma)$  but  $\xi \in \text{ftv}(\theta(\Gamma))$ ). Hence

$$\mathcal{R}(\xi, \theta(\Gamma)) \leq \mathcal{R}(\xi, \Gamma) \text{ if } \xi \notin \text{dom}(\theta)$$

Thus, except for the case where a univariable is eliminated by a substitution, a substitution can only decrease the rank of a univariable or leave it the same. This means that a substitution may “promote”

a univariable to an earlier binding and thus make it “less generalizable”.

As we have defined it, computing the rank of a univariable relative to an environment is no less expensive than computing whether it occurs in the environment. But the main idea of the ranked variable algorithms we will present below is that we can avoid scanning the environment by making the rank an attribute of the univariable and managing that attribute incrementally so that it is always equal to the intrinsic rank defined above, relative to the current environment in the machine state.

Since in the  $\mathcal{W}$  machine for Pure ML, only the  $\lambda$  bindings *count*, our first version of univariable ranking will assign ranks based on the depth of  $\lambda$  nesting (*i.e.* counting only  $\lambda$  bindings). We will then look at a variant for Pure ML that counts only let definiens marks (which demarcate the scopes of nested let definiens). When we introduce references and the consequent requirement for the value restriction, things become more complicated, because the invariant that univariables first occur in  $\lambda$  bindings will no longer be the case, and non-value let bindings will act like  $\lambda$  bindings.

## 4. $\lambda$ -Depth Ranking

In Figs. 6 and 7, we modify the Pure ML inference machine by adding ranks to the univariables. The rank associated with a univariable is either a natural number or  $\infty$ , and if a univariable occurs in the current type environment its rank is meant to correspond to the index of the earliest  $\lambda$ -binding in the environment in which it occurs. A rank greater than the number of  $\lambda$  bindings in the current environment will indicate that the univariable is generalizable, *i.e.* its intrinsic rank is  $\infty$ .

For convenience, we add a new register to the machine states containing the current  $\lambda$  nesting depth  $d$ . This register starts at 0 and counts how many  $\lambda$ -abstractions we have passed through to reach to the current control. Its value will always be equal to both the number of  $\lambda$ -bindings in the current environment and the number of  $\lambda$  frames in the context stack.

The  $\mathcal{W}$  machine rules are easily adapted to incorporate  $\lambda$ -depth ranking, and we will call this new machine  $\mathcal{W}_\lambda$ . The new [λ-in] rule increments the depth, introduces a fresh univariable whose initial rank is the new  $\lambda$  nesting depth, and pushes a new  $\lambda$  binding onto the environment. The [λ-out] rule now also decrements the depth.

We can categorize bindings in type environments according to whether they are introduced by the [let-body] rule, in which case the variable is bound to a polytype  $\sigma$ , or by the [λ-in] rule, in which case the variable is bound to a monotype  $\tau$  (and this distinction is preserved under substitution). This distinction on bindings is expressed in the definition of environments in Fig. 7.

As noted in the previous section, if a univariable occurs in an environment, the first binding it occurs in will always be a  $\lambda$ -binding, and let-bindings play a secondary role. So we can modify the notion of the occurrence index of a univariable  $\xi$  in an environment  $\Gamma$ , which we called the intrinsic rank  $\mathcal{R}(\xi, \Gamma)$ , by counting only  $\lambda$ -bindings from left to right when we calculate this index. The univariable ranks in the  $\mathcal{W}_\lambda$  machine are assigned and managed so that they always correspond to the intrinsic rank determined by the location of the univariable in the current environment (if any).

Note that because applying unification substitutions in [app-out] can “promote” univariables to earlier  $\lambda$ -bindings, the first  $\lambda$ -binding containing a univariable may not correspond to the program point where it was introduced. A substitution that promotes a univariable to an earlier binding must compensate for the change in the intrinsic rank of the univariable by adjusting the explicit rank of the univariable to reflect its new earliest occurrence in the environment.

The  $\mathcal{L}$  limit substitution defined in Fig. 8 performs this necessary adjustment to univariable ranks. When a unification substitu-

$e$	::=	$x \mid \lambda x.e \mid e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2$	EXP
$\tau$	::=	$\xi^d \mid \tau \rightarrow \tau \mid \alpha$	TYPE
$d$	::=	$0 \mid 1 \mid 2 \mid \dots \mid \infty$	Ranks
$\sigma$	::=	$\forall \bar{\alpha}.\tau$	PTYPE

**Figure 6.**  $\lambda$ -Ranking Language

$\Gamma$	::=	$\emptyset \mid \Gamma[x : \tau] \mid \Gamma[x : \sigma]$	Tyenv
$f$	::=	$\lambda \mid (\square e) \mid (\tau \square) \mid \text{let } x \text{ in } e \mid \text{let}$	Frames
$k$	::=	$\bullet \mid f :: k$	Stack
$c$	::=	$e \mid \tau$	Control
$s$	::=	$(c, \Gamma, d, k)$	States
Initial State $s_0 = (e, \emptyset, 0, \bullet)$			
$(x, \Gamma, d, k)$	$\mapsto_\lambda$	$(\mathcal{I}(\Gamma(x)), \Gamma, d, k)$	[var]
$(\lambda x.e, \Gamma, d, k)$	$\mapsto_\lambda$	$(e, \Gamma[x : \xi^{d+1}], d + 1, \lambda :: k)$	[ $\lambda$ -in]
$(\tau, \Gamma[x : \tau'], d, \lambda :: k)$	$\mapsto_\lambda$	$(\tau' \rightarrow \tau, \Gamma, d - 1, k)$	[ $\lambda$ -out]
$(e_1 e_2, \Gamma, d, k)$	$\mapsto_\lambda$	$(e_1, \Gamma, d, (\square e_2) :: k)$	[app-l]
$(\tau, \Gamma, d, (\square e_2) :: k)$	$\mapsto_\lambda$	$(e_2, \Gamma, d, (\tau \square) :: k)$	[app-r]
$(\tau, \Gamma, d, (\tau' \square) :: k)$	$\mapsto_\lambda$	$(\theta \xi^\infty, \theta \Gamma, d, \theta k)$	[app-out]
$\theta = \mathcal{U}(\tau', \tau \rightarrow \xi^\infty)$			
$(\text{let } x = e_1 \text{ in } e_2, \Gamma, d, k)$	$\mapsto_\lambda$	$(e_1, \Gamma, d, \text{let } x \text{ in } e_2 :: k)$	[let-def]
$(\tau, \Gamma, d, \text{let } x \text{ in } e_2 :: k)$	$\mapsto_\lambda$	$(e_2, \Gamma[x : \sigma], d, \text{let} :: k)$	[let-body]
$\sigma = \forall \bar{\alpha}.\{\bar{\alpha}/\mathcal{G}_d(\tau)\}\tau$ where $\bar{\alpha}$ is fresh			
$(\tau, \Gamma[x : \sigma], d, \text{let} :: k)$	$\mapsto_\lambda$	$(\tau, \Gamma, d, k)$	[let-out]
$\mathcal{G}_d(\tau)$	=	$\{\xi^m \in \text{ftv}(\tau) \mid m > d\}$	Generalizable
$\mathcal{I}(\forall \bar{\alpha}.\tau)$	=	$\{\xi^\infty/\bar{\alpha}\}\tau$	Generic Instantiation
$\mathcal{I}(\tau)$	=	$\tau$	

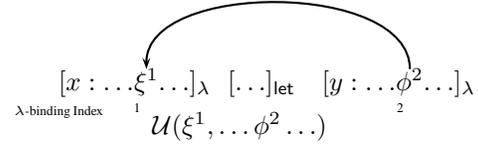
**Figure 7.** Algorithm  $\mathcal{W}$  Machine extended with  $\lambda$ -depth Ranking ( $\mathcal{W}_\lambda$ )

$\mathcal{U}(\tau, \tau) = \varepsilon$	$\mathcal{U}(\tau, \xi^m) = \mathcal{U}(\xi^m, \tau)$
where $\tau$ is not a univariable	
$\mathcal{U}(\xi^m, \tau) = \mathcal{L}_m(\tau) \circ \{\tau/\xi^m\}$ if $\xi^m \notin \text{ftv}(\tau)$	
$\mathcal{U}(\tau_1 \rightarrow \tau_2, \tau_3 \rightarrow \tau_4) = \mathcal{U}(\theta(\tau_2), \theta(\tau_4)) \circ \theta$	
where $\theta = \mathcal{U}(\tau_1, \tau_3)$	
$\mathcal{L}_d(\tau) = \{\xi^d/\xi^m \mid \xi^m \in \text{ftv}(\tau) \wedge m > d\}$	

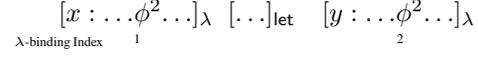
**Figure 8.** Unification with  $\lambda$ -rank Maintenance

tion maps a univariable  $\xi$  to a type  $\tau$ , univariables occurring in  $\tau$  may be introduced to  $\lambda$ -bindings to the left of their original first occurrences, if  $\xi$  occurs earlier in the environment. Thus, if the first occurrence of  $\xi$  was to the left of the first occurrence of some univariable  $\phi \in \text{ftv}(\tau)$ , then after unification, the first occurrence of  $\phi$  is now where the first occurrence of  $\xi$  was before as shown in Fig. 9. If the first occurrence of  $\xi$  was to the right of the first occurrence of  $\phi$ , then the first occurrence of  $\phi$  does not move after unification as shown in Fig. 10.

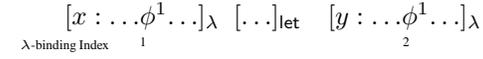
As a simple example of this issue, consider the expression  $\lambda x.\text{let } y = \lambda z.xz \text{ in } \dots$ . The  $\lambda$  bindings for  $x$  and  $z$  produce an environment  $\Gamma = [x : \xi_x^1][z : \xi_z^2]$  for the typing of the application  $xz$ . The typing of this application involves a unification producing the substitution  $\theta = \xi_x^1 \mapsto (\xi_z^2 \rightarrow \xi_r^\infty)$ . But if  $\theta$  is applied to the typing environment it would produce  $\Gamma' = [x : \xi_z^2 \rightarrow \xi_r^\infty][z : \xi_z^2]$  in which the ranks of the univariables in the  $x$  binding no longer correspond to its index in the environment, which is 1. So an additional rank limiting substitution is required to adjust the ranks of  $\xi_z$  and  $\xi_r$  downward:  $(\xi_z^2 \mapsto \xi_z^1, \xi_r^\infty \mapsto \xi_r^1)$ . When this is applied we get a properly ranked environment  $\Gamma'' = [x : \xi_z^1 \rightarrow \xi_r^1][z : \xi_z^1]$ . Now the adjusted rank of  $\xi_z^1$  will prevent the type of  $y$ , namely  $\xi_z^1 \rightarrow \xi_z^1$  from being generalized at the let binding.



(a) If  $\xi$  (a contextual variable) instantiates to  $y$ 's type, then  $\phi$  (a local variable) propagates into  $x$ 's  $\lambda$ -binding.

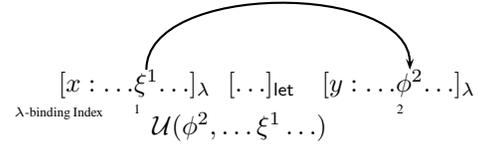


(b) The rank of  $\phi$  must now be limited to that of  $\xi$  because the first occurrence of  $\phi$  has shifted left, thus making it contextual.

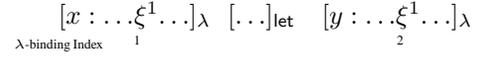


(c) After Applying the Limit Substitution

**Figure 9.** Limit Substitution ( $\mathcal{L}$ ) Example: Subscripts on bindings indicate the kind of binding.



(a) If  $\phi$  instantiates to  $x$ 's type, then  $\xi$  propagates into  $y$ 's  $\lambda$ -binding.



(b)  $\xi$ 's rank is fine because its first occurrence has not shifted. It remains contextual.

**Figure 10.** Example of when the limit substitution is unnecessary

### Definition 1 (Valid Machine States ( $\mathcal{VMS}$ ))

A machine state  $(c, \Gamma, d, k)$  is valid if it is reachable from some valid initial machine state  $(e, \emptyset, 0, \bullet)$  (i.e.,  $(e, \emptyset, 0, \bullet) \mapsto^* (c, \Gamma, d, k)$  where  $\mapsto$  is the machine transition relation).

From this point on, a metavariable with a tilde ( $\tilde{\phantom{x}}$ ) on top denotes the corresponding instance of that construct where the ranks of any univariables that occur are erased or ignored.

### Definition 2

$|\Gamma| = \text{the length of } \Gamma, \text{ and similarly for } |\tilde{\Gamma}|.$

### Definition 3

$\mathcal{D}(\Gamma) = \# \lambda \text{ bindings in } \Gamma, \text{ and similarly for } \mathcal{D}(\tilde{\Gamma}).$

### Definition 4

$\Gamma_\lambda$  denotes the list of  $\lambda$ -bindings in  $\Gamma$ , and similarly for  $\tilde{\Gamma}_\lambda$ .

**Definition 5 (Intrinsic Univariable Rank)**

$$\begin{aligned}\tilde{\mathcal{R}}(\xi, \emptyset) &= \infty \\ \tilde{\mathcal{R}}(\xi, \tilde{\Gamma}[x : \tilde{\tau}]) &= \min(\tilde{\mathcal{R}}(\xi, \tilde{\Gamma}), \mathcal{D}(\tilde{\Gamma}) + 1) \text{ if } \xi \in \text{ftv}(\tilde{\tau}) \\ &= \infty \text{ otherwise} \\ \tilde{\mathcal{R}}(\xi, \tilde{\Gamma}[x : \tilde{\sigma}]) &= \tilde{\mathcal{R}}(\xi, \tilde{\Gamma})\end{aligned}$$

This computes the index of the first  $\lambda$ -binding containing an occurrence of a univariable, ignoring let-bindings. If a univariable is not in the environment, it is assigned the intrinsic rank  $\infty$ , which ensures that such a univariable should be generalized because it is always greater than any let expression's finite  $\lambda$ -depth.

**Definition 6**

$\mathcal{R}(\xi^m, \Gamma) = \tilde{\mathcal{R}}(\xi, \mathcal{E}_r(\Gamma))$   
where  $\mathcal{E}_r$  erases the ranks of univariables in  $\Gamma$ .

By definition, type environments can be extended with more bindings. We formally define the notion of extension:

**Definition 7**

$\tilde{\Gamma}'$  is an **extension** of  $\tilde{\Gamma}$  if and only if there exists  $\tilde{\Gamma}''$  such that  $\tilde{\Gamma}'$  is the concatenation of  $\tilde{\Gamma}$  and  $\tilde{\Gamma}''$  (in that order).

Notice that if a type variable occurs in a type environment, then extension of that type environment does not change that type variable's leftmost occurrence:

**Lemma 1 ( $\tilde{\mathcal{R}}$  Invariant Under  $\tilde{\Gamma}$  Extension)**

If  $\xi \in \text{ftv}(\tilde{\Gamma})$  and  $\tilde{\Gamma}'$  is an extension of  $\tilde{\Gamma}$ , then  $\tilde{\mathcal{R}}(\xi, \tilde{\Gamma}') = \tilde{\mathcal{R}}(\xi, \tilde{\Gamma})$ .

**Proof:** By induction on the length of  $\tilde{\Gamma}$  ■

Because  $\tilde{\mathcal{R}}$  computes the first occurrence index for a univariable in an environment, if a univariable does occur in the environment  $\tilde{\mathcal{R}}$  will return the same index for any extension of that environment.

By definition,  $\tilde{\mathcal{R}}$  counts the number of  $\lambda$ -bindings up to a univariable's first occurrence. Naturally, this number cannot exceed the total number of  $\lambda$ -bindings on  $\tilde{\Gamma}$ . We observe this fact in Lemma 2.

**Lemma 2**

If  $\xi \in \text{ftv}(\tilde{\Gamma}_\lambda)$ , then  $\tilde{\mathcal{R}}(\xi, \tilde{\Gamma}) \leq \mathcal{D}(\tilde{\Gamma})$ .

**Proof:** By induction on the length of  $\tilde{\Gamma}$  ■

We determined by construction that the rank of a univariable that does not occur in the environment is  $\infty$ . This is formally stated in Lemma 3.

**Lemma 3**

If  $\xi \notin \text{ftv}(\tilde{\Gamma}_\lambda)$  then  $\tilde{\mathcal{R}}(\xi, \tilde{\Gamma}) = \infty$ .

**Proof:** By induction on the structure of  $\tilde{\Gamma}$  ■

To formalize the notion of first occurrence of a univariable, we appeal to Lemma 4.

**Lemma 4**

If  $\xi^m$  occurs in the  $n$ -th  $\lambda$ -binding in  $\Gamma$  (counting from the left) and  $\xi^m$  occurs in no  $\lambda$ -binding to the left of that  $n$ -th  $\lambda$ -binding, then  $\mathcal{R}(\xi^m, \Gamma) = n$ .

**Proof:** By induction on the length of  $\Gamma$  and Lemma 3 ■

We also make a simple observation about substitutions and environments that simplifies the proofs in Lemma 5.

**Lemma 5 (Substitution Range Preservation)**

If  $\tau \in \text{rng } \Gamma$ , then  $\sigma\tau \in \text{rng } \sigma\Gamma$ .

**Proof:** This follows from definition of type variable substitution on type environments. ■

**Definition 8 (Properties for  $\mathcal{W}_\lambda$ )**

The following are invariant properties on any valid  $\mathcal{W}_\lambda$  machine state  $s = (c, \Gamma, d, k)$ :

**Prop<sub>1</sub>**  $\forall \xi^m \in \text{ftv}(\Gamma), \mathcal{R}(\xi^m, \Gamma) = m$

**Prop<sub>2</sub>**  $\forall \xi^m \in \text{ftv}(c)$  such that  $m \leq \mathcal{D}(\Gamma)$ , then  $\xi^m \in \text{ftv}(\Gamma)$

**Prop<sub>3</sub>**  $\mathcal{D}(\Gamma) = \#$  of  $\lambda$ s on  $k$

**Prop<sub>4</sub>**  $\forall \xi^m \in \text{ftv}(k)$  such that  $m \leq r$  where  $r$  is the number of  $\lambda$  frames below the frame with the occurrence of  $\xi^m$  in the stack  $k$ ,  $\xi^m \in \text{ftv}(\Gamma)$

Prop<sub>1</sub> asserts that a univariable's rank corresponds to the index of the  $\lambda$ -binding containing the first occurrence of that univariable in the environment. This is a precursor to Corollary 1 below that provides a notion of completeness for the  $\lambda$ -ranking scheme. All univariables in the environment must be given a low rank.

Prop<sub>2</sub> says that all the "contextual" univariables occurring in the control are bound in the environment. Contextual refers to having a rank no greater than  $\mathcal{D}(\Gamma)$ . This specifically excludes univariables that would run off the end of the environment at the current scope because their ranks are too great, including univariables with rank  $\infty$ . Prop<sub>2</sub> asserts that all low ranked univariables must be in the environment. Therefore it provides a notion of the soundness of the ranking scheme.

Prop<sub>3</sub> tells us that the number of  $\lambda$ -bindings in the environment corresponds to the  $\lambda$ -nesting depth. This simple observation is necessary for establishing Prop<sub>1</sub> when only  $\lambda$ -depth information is available.

Prop<sub>4</sub> is a rather subtle property. It ensures that the a modified version of Prop<sub>2</sub> holds for those univariables occurring on the stack in  $(\tau \square)$  frames. Because the type  $\tau$  of the operator is in a sense "parallel" to the current operand expression that the machine is working on, it may contain some seemingly low ranked univariables (rank < current  $\lambda$ -depth) that are not in the current environment for the operand. These univariables are not truly low ranked because at the point when  $(\tau \square)$  was pushed onto the stack,  $\tau$  was in the control, so by Prop<sub>2</sub> any univariables not in the environment at that point had rank higher than the  $\mathcal{D}(\Gamma')$ , where  $\Gamma'$  was the then-current environment. But  $\mathcal{D}(\Gamma')$  was also equal to the number of  $\lambda$  frames on the stack before the  $(\tau \square)$  frame was pushed. Once the machine pops all the frames sitting on top of the  $(\tau \square)$  frame, the environment will once again have its original  $\lambda$ -depth.

For example, in  $(\lambda x.x)(\lambda y.\lambda z.z)$ , before the machine reaches the body of  $\lambda z.z$ ,  $\lambda x.x$  is given the type  $\xi^1 \rightarrow \xi^1$ . When typing the body of  $\lambda z.z$ , the  $\lambda$ -depth is 2. It seems that  $\xi^1$  has a low rank. However, at the point we were just finished with  $\lambda x.x$  and when we pop out of the operand, the  $\lambda$ -depth is 0. Therefore  $\xi^1$  is actually a high ranked univariable that we do not expect to find in the environment when we get back to it.

We proceed by first showing that these properties are invariant under the unification that results from [app-out].

**Lemma 6 (Invariants under unification)**

For all  $\Gamma$ , if the properties from def. 8 hold for  $s = (\tau, \Gamma, d, (\tau' \square) :: k) \in \mathcal{VMS}$  and  $\theta = \mathcal{U}(\tau', \tau \rightarrow \psi^\infty)$  (where  $\psi^\infty$  is fresh), then the properties must also hold for  $s' = (\theta\psi^\infty, \theta\Gamma, d, \theta k)$ .

**Proof (sketch):** This proof is by induction on sizes of  $\tau'$  and the other type argument of  $\mathcal{U}$ . The instantiation case is the interesting part of this proof. In that case, by definition of unification,  $\theta = \mathcal{U}(\phi^n, \tau) = \mathcal{L}_n(\tau) \circ \{\tau/\phi^n\}$ .

The interesting cases for  $\text{Prop}_1(s')$  are illustrated in Figs. 10 and 9. In the first case, assume  $\xi^m \in \text{ftv}(\theta\Gamma)$  was not limited, hence it must have a rank that is lower or equal to  $n$ . Because  $\phi^n$  was substituted away by  $\theta$ , it must occur in  $\Gamma$ . By, Lemma 2,  $n \leq \mathcal{D}(\Gamma)$ . Thus,  $m \leq \mathcal{D}(\Gamma)$ . By  $\text{Prop}_2(s)$ ,  $\xi^m$  must occur in  $\Gamma$ . By  $\text{Prop}_1(s)$ ,  $\mathcal{R}(\xi^m, \Gamma) = m$ . In the case where  $\xi^m$  was limited, the first occurrence of  $\xi^m$  shifted left to the former first occurrence of  $\phi^n$ . Thus, it should also receive rank  $n$ . Lem. 4 helps establish this.

For  $\text{Prop}_2(s')$ , the unlimited case is simple. The case where  $\xi^m$  was limited is somewhat tricky. If the univariable to be instantiated  $\phi^n$  has a rank  $n > \mathcal{D}(\theta\Gamma)$ , the case is vacuous, so assume that  $m \leq \mathcal{D}(\theta\Gamma)$ . Note that by the definition of  $\mathcal{L}$ ,  $m = n$ . We first have to show that  $\phi^n$  must be in  $\Gamma$ .  $\phi^n$  can come from  $\tau$ ,  $\tau'$ , or  $\psi^\infty$ . The  $\psi^\infty$  case is vacuous because  $\infty > \mathcal{D}(\Gamma)$ . The  $\tau$  case is handled by  $\text{Prop}_2(s)$ . The  $\tau'$  case takes a little more work because  $\text{Prop}_2(s)$  did not say anything about  $\tau'$ . This is the reason why we need  $\text{Prop}_4$ .  $\text{Prop}_4(s)$  gives us what we need. Assuming that  $\phi^n \in \text{ftv}(\Gamma)$ , it is straightforward to show that  $\xi^m \in \text{ftv}(\Gamma)$  with the help of Lemma 5.

$\text{Prop}_3(s')$  is easily to establish because substitution does not change the length of  $\Gamma$  or  $k$ .  $\text{Prop}_4$  is also relatively straightforward. ■

The properties trivially hold for the initial states. Now we proceed to show that the properties are indeed invariant under all machine transition rules. This would establish that the properties are invariants over all valid machine states.

**Lemma 7 (Properties invariant over all machine transitions)**

If the props. in def. 8 hold for  $s \in \mathcal{VMS}$  and  $s \mapsto_\lambda s'$ , then they must hold for  $s'$ .

**Proof (sketch):** We need  $\text{Prop}_3$  to push through  $\text{Prop}_1$ 's  $[\lambda\text{-in}]$  case, thus showing that the ranks do measure  $\lambda$  nesting depth. This fact combined with Lemma 3 establish that the ranks correspond to the  $\lambda$ -binding index of the first occurrence (actually the only occurrence) of the fresh univariable representing this  $\lambda$ 's parameter type in the environment.

For the  $[\lambda\text{-out}]$ ,  $[\text{let-body}]$ , and  $[\text{let-out}]$  cases, we use Lemma 1 to show  $\text{Prop}_1$ . For  $[\text{app-r}]$ ,  $\text{Prop}_4(s')$  simply inherits the property from  $\text{Prop}_2(s)$  because  $\tau$  merely shifted from the control to the top of the stack.  $\text{Prop}_1$  and  $\text{Prop}_2$  for the  $[\text{app-out}]$  case follow from Lemma 6.

For  $\text{Prop}_1$  in the  $[\text{let-body}]$  case, we use  $\text{Prop}_2$ 's prop. that we assume holds before transitioning by  $[\text{let-body}]$  to show that any univariable  $\xi$  in  $\text{ftv}(\Gamma[x : \sigma])$  must also occur in  $\text{ftv}(\Gamma)$ . That is to say, first occurrences of univariables in environments are never in let-bindings. We can then conclude by assumption and Lemma 1 that  $\mathcal{R}(\xi^m, \Gamma[x : \sigma]) = m$ .

For  $\text{Prop}_2$ , we note that  $\xi^m \in \text{ftv}(\tau' \rightarrow \tau)$  in  $[\lambda\text{-out}]$  cannot have a first occurrence in the last binding in  $\Gamma[x : \tau']$  because by assumption,  $m \leq \mathcal{D}(\Gamma)$ . If this were not the case, then we would be seeing unbound yet non-generic univariables after  $[\lambda\text{-out}]$  transitions.

The rest of the cases are fairly straightforward case analyses assuming that the properties holds before each transition rule. ■

From  $\text{Prop}_1$ , we can easily obtain a completeness corollary.

**Corollary 1 (Ranking is Complete)**

For all  $\Gamma$  from  $s = (c, \Gamma, d, k) \in \mathcal{VMS}$   
 $\xi^m \in \text{ftv}(\Gamma) \Rightarrow m \leq \mathcal{D}(\Gamma)$

**Proof (sketch):** By  $\text{Prop}_1$  from Lemma 7,  $\mathcal{R}(\xi^m, \Gamma) = m$ . Since  $\mathcal{D}(\Gamma) = |\Gamma_\lambda|$ ,  $m \leq \mathcal{D}(\Gamma)$  by definition of  $\mathcal{R}$  and Lemma 2. ■

We can put together Corollary 1 and  $\text{Prop}_2$  to establish that the ranking corresponds to searching through the environment for all valid machine states.

**Lemma 8**

For any  $(\tau, \Gamma, d, k) \in \mathcal{VMS}$  and  $\forall \xi^m \in \text{ftv}(\Gamma)$ ,  
 $m \leq d$  iff  $\xi^m \in \text{ftv}(\Gamma)$ .

**Proof:**

$\Rightarrow$  : This direction follows from  $\text{Prop}_2$  of Lemma 7.

$\Leftarrow$  : This direction follows from Corollary 1.  $m \leq \mathcal{D}(\Gamma) = d$  ■

The generalizability criteria are equivalent. Consequently, we always generalize the same set of univariables using the classical and  $\lambda$ -ranking algorithms.

**Lemma 9 (Equivalence of Generalizability Criteria)**

If  $(\tau, \Gamma, d, \text{let } x = \square \text{ in } e :: k) \mapsto (e, \Gamma[x : \sigma], d, \text{let } :: k)$  where  $\sigma = \forall \bar{\alpha}. \{\bar{\alpha}/\mathcal{G}_d(\Gamma)\}\tau$ , then  $\mathcal{G}_d(\Gamma) = \mathcal{G}_\Gamma(\Gamma)$ .

**Proof:** This proof falls out from the definitions of  $\mathcal{G}_d$  and  $\mathcal{G}_\Gamma$ , and Lemma 8. ■

Finally, we put everything together. The classical algorithm  $\mathcal{W}$  and the  $\lambda$ -ranking algorithms are equivalent.

**Theorem 1 (Equivalence to Classical Algorithm  $\mathcal{W}$ )**

$(e, \emptyset, 0, \bullet) \mapsto^* (\tau, \emptyset, 0, \bullet)$  iff  $(e, \emptyset, \bullet) \mapsto^* (\tau, \emptyset, \bullet)$ .

**Proof:** We only inspect the depth register in  $[\text{let-body}]$ . In all other cases, the  $\lambda$ -ranking and classical machines are identical. The depth register is only along for the ride. Because of Lemma 9, the  $[\text{let-body}]$  behaviors are also identical. ■

## 5. Let-Depth Ranking

The let-depth ranking scheme is similar to the  $\lambda$ -depth ranking except for a few significant adjustments. We can reuse the language and the same notation for ranks. The main observation that the let-depth ranking relies on is that the locality of univariables at different  $\lambda$ -depths is the same as long as we are in the same let definien's scope. In contrast,  $\lambda$ -ranks distinguish among  $\lambda$ 's at different  $\lambda$ -nestings even though they may all be in the same let definien scope. We have to be careful, though, because a let expression only extended the type environment when the type checker reaches the body. At that point, generalization of the definien type would have already occurred. So we cannot use let-bindings to derive a ranking for determining generalization. We have to determine a rank before generalization occurs. So, we turn to another implicit structure of the environment, the partitioning of bindings into contiguous

```

 $\lambda x. \text{let } a = \lambda y_1. \lambda y_2. \text{let } b = \text{let } c = \bullet$ 
  in ...
  in ...
  in ...

```

Environment at  $\bullet$ :  $[x : \xi^0][y_1 : \zeta^1][y_2 : \phi^1]$

**Figure 11.** When we reach  $\bullet$ , we pass into three let definiens but added no let-bindings to the environment. The type of  $x$  is less generalizable than that of the  $y$ 's. The type of the  $y$ 's should be generalizable at the same point.

```

let  $a = \lambda x. \text{let } b = \dots$ 
  in let  $c = \dots$ 
  in  $(\lambda y. y) \bullet$ 
in ...

```

Environment at  $\bullet$ :  $[x : \xi^1][b : \sigma_b][c : \sigma_c][y : \zeta^1]$

**Figure 12.** When we reach  $\bullet$ , we are in no let definiens but added two let-bindings to the environment. The types of  $x$  and  $y$  are both generalizable at the let  $a$ .

blocks according to the definien scope. We can extend the machine introduced in  $\mathcal{W}_\lambda$  to make this let definien partitioning manifest.

We cannot naively replace  $\lambda$ 's with let's in the development from the previous section to get a working let-ranking scheme. In particular,  $\lambda$ 's were simple. We only traverse a  $\lambda$ -abstraction in one pass. First, we extend the environment with the bound  $\lambda$ -variable. Then we traverse the body of the  $\lambda$  using this extended environment. Once we are done, we pop out of the  $\lambda$  never to revisit it. In contrast, for let's, we have to first traverse the definien without extending the environment, pop out of the definien, extend the environment with the generalized definien type, and only then traverse the body of the let with this extended environment. We have to increment the depth when we pass into a definien, as for instance in the expression  $\lambda x. \text{let } y = x \text{ in } \dots$ , where we need to prevent the type of  $x$  from being generalized where  $y$  is bound. On the other hand, the body should have the depth that held before we entered into the definien. This permits the type of  $\lambda y. y$  to generalize in  $\text{let } x = \text{let } \dots \text{ in } \lambda y. y \text{ in } \dots$ . The complication is that there is no simple connection between the number of let-bindings in the environment and the number of let definiens we have passed into (but have not finished typing). Figs. 11 and 12 give examples of this problem. In the environments for the two examples, we label the unvariables with the let definien depth at the point of introduction. Although in Fig. 11  $\xi$  has a lower let definien depth than  $\zeta$  and  $\phi$ , there is nothing in the environment to demarcate the two groups of bindings. In Fig. 12,  $\xi$  and  $\zeta$  have the same depth and yet they are separated by two let-bindings.

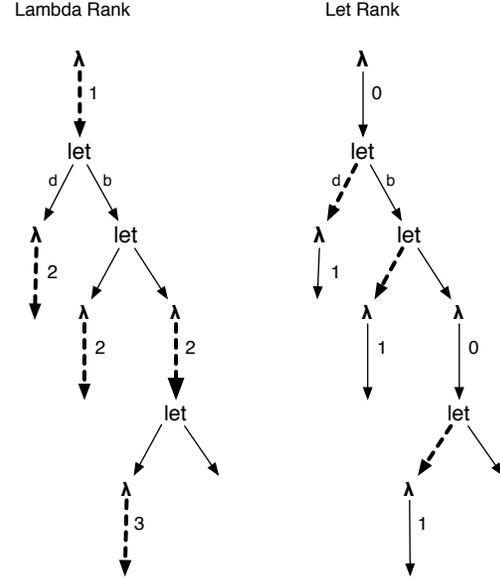
To overcome this problem, we rank by the number of let definien letd markers in an environment instead of  $\lambda$ -bindings. We want to identify contiguous blocks of  $\lambda$ -bindings in the environment that are within the same let definien scope. Thus, every time we pass into a let definien, we mark the (rightmost) end of the environment with a letd marker. This marker tells us that any  $\lambda$  to left of the marker will be less generalizable than any  $\lambda$  that will be to the right. We extend the definition of  $\Gamma$  with these letd marks in Fig. 14. The letd marked versions of the environments in Figs. 11 and 12 would be  $[x : \xi^0][\text{letd}][y_1 : \zeta^1][y_2 : \phi^1][\text{letd}][\text{letd}]$  and  $[\text{letd}][x : \xi^1][b : \sigma_b][c : \sigma_c][y : \zeta^1]$  respectively.

Fig. 13 shows both the  $\lambda$  and let ranking of the following program:

```

 $\lambda s. \text{let } x = \lambda t. \dots$ 
  in let  $y = \lambda u. \dots$ 

```



**Figure 13.** Comparison of  $\lambda$  and let ranks in expression structure

```

in  $\lambda v. \text{let } z = \lambda w. \dots$ 
  in ...

```

The left branch of each let node represents its definien (d) and the right branch its body (b). The dashed lines indicate where we increment the rank. In this example, there are no let definiens nested in another let definien, hence generalizability at a let definien is determined solely by whether a univariable is local to that let definien or outside of any let definien. Unvariables at that let definien could not have been introduced in another let definien because they would have already been generalized at that earlier (to the left) let definien or we would have not reached that let definien yet. Notice that the  $\lambda$  ranking method increments the rank much more frequently. In fact, the distinction between ranks on the right spine of the program does not make a difference with respect to generalizability. For example, the types of  $s$  and  $v$  are in the same scope of generalization (*i.e.*, they are both outside of all let definiens). The let ranking method only increments the rank at the let definiens where the rank distinction may make a difference with respect to generalizability.

#### Definition 9 (Let Depth)

$\mathcal{D}_{\text{let}}(\tilde{\Gamma}) = \# \text{ of letd markers in } \tilde{\Gamma}$

#### Definition 10 (Univariable Let-Rank)

$$\begin{aligned}
 \tilde{\mathcal{R}}_{\text{let}}(\xi, \emptyset) &= \infty \\
 \tilde{\mathcal{R}}_{\text{let}}(\xi, \tilde{\Gamma}[x : \tilde{\tau}]) &= \min(\tilde{\mathcal{R}}_{\text{let}}(\xi, \tilde{\Gamma}), m) \\
 &\quad \text{where } m = \mathcal{D}_{\text{let}}(\tilde{\Gamma}) \text{ if } \xi \in \text{ftv}(\tilde{\tau}), \\
 &\quad \text{otherwise } m = \infty \\
 \tilde{\mathcal{R}}_{\text{let}}(\xi^m, \tilde{\Gamma}[x : \tilde{\sigma}]) &= \tilde{\mathcal{R}}_{\text{let}}(\xi^m, \tilde{\Gamma})
 \end{aligned}$$

#### Definition 11

$\mathcal{R}_{\text{let}}(\xi^m, \Gamma) = \tilde{\mathcal{R}}_{\text{let}}(\xi, \mathcal{E}_r(\Gamma))$

$\Gamma ::= \Gamma[\text{letd}] \mid \dots$	$\Gamma$ with let definien marks
$s = (c, \Gamma, d, k)$	states
Initial Machine State $s_0 = (e, \emptyset, 0, \bullet)$	
$(x, \Gamma, d, k) \mapsto_{\mathcal{L}} (\mathcal{I}(\Gamma(x)), \Gamma, d, k)$	[wl-var]
$(\lambda x.e, \Gamma, d, k) \mapsto_{\mathcal{L}} (e, \Gamma[x : \xi^d], d$	[wl- $\lambda$ -in]
$\lambda :: k)$	
$(\tau, \Gamma[x : \tau'], d, \lambda :: k) \mapsto_{\mathcal{L}} (\tau' \rightarrow \tau, \Gamma, d, k)$	[wl- $\lambda$ -out]
$(e_1 e_2, \Gamma, d, k) \mapsto_{\mathcal{L}} (e_1, \Gamma, d, (\square e_2) :: k)$	[wl-app-l]
$(\tau, \Gamma, d, (\square e_2) :: k) \mapsto_{\mathcal{L}} (e_2, \Gamma, d, (\tau \square) :: k)$	[wl-app-r]
$(\tau, \Gamma, d, (\tau' \square) :: k) \mapsto_{\mathcal{L}} (\theta \xi, \theta \Gamma, d, \theta k)$	[wl-app-out]
$\theta = \mathcal{U}(\tau', \tau \rightarrow \xi^\infty)$	
$(\text{let } x = e_1 \text{ in } e_2, \Gamma, d, k)$	[wl-let-def]
$\mapsto_{\mathcal{L}} (e_1, \Gamma[\text{letd}], d + 1, \text{let } x \text{ in } e_2 :: k)$	
$(\tau, \Gamma[\text{letd}], d, \text{let } x \text{ in } e_2 :: k)$	[wl-let-body]
$\mapsto_{\mathcal{L}} (e_2, \Gamma[x : \sigma], d - 1, \text{let } :: k)$	
$\sigma = \forall \bar{\alpha}. \{\bar{\alpha} / \mathcal{G}_d^{\text{let}}(\tau)\} \tau, \bar{\alpha} \text{ is fresh}$	
$(\tau, \Gamma[x : \sigma], d, \text{let } :: k) \mapsto_{\mathcal{L}} (\tau, \Gamma, d, k)$	[wl-let-out]
$\mathcal{G}_d^{\text{let}}(\tau) = \{\xi^m \in \text{ftv}(\tau) \mid m \geq d\}$	
Generic instantiation $\mathcal{I}$ , language, stack register $k$ , and frames $f$ remain the same as in the $\mathcal{W}_\lambda$ machine in Fig. 7.	

**Figure 14.** Algorithm  $\mathcal{W}$  Machine with let-Ranking ( $\mathcal{W}_L$ )

The univariable rank no longer indicates the exact  $\lambda$ -binding that has the first occurrence in the environment. Instead, the rank is the index of the nearest letd marker to the left of the  $\lambda$ -binding with the first occurrence. We can call this the *left bounding letd marker for the first occurrence*. Multiple, distinct univariables can share the same left bounding letd mark. These univariables have the same generalizability. As before, generic univariables are rank  $\infty$ . Univariables that are bound outside of any let definiens are rank 0. If the environment was  $[x : \xi][\text{letd}_1][y_1 : \phi_1][y_2 \phi_2][\text{letd}_2][z : \zeta]$ , then  $\xi, \phi_1, \phi_2$ , and  $\zeta$  would be ranked 0, 1, 1, and 2 respectively. The left bounding letd's for the  $\phi$ 's and  $\zeta$  are  $\text{letd}_1$  and  $\text{letd}_2$  respectively. Also, similar to the  $\mathcal{W}_\lambda$  rank function  $\tilde{\mathcal{R}}, \tilde{\mathcal{R}}_{\text{let}}$  still ignores let-bindings when computing the rank. As we have shown before, let-bindings have nothing to do with rank and generalizability.

Fig. 14 gives an abstract machine for a type inference algorithm that uses let-ranking. The [wl- $\lambda$ -in] and [wl- $\lambda$ -out] rules do change the depth counter. This is now the responsibility of the [wl-let-def] and [wl-let-body] rules that introduce and eliminate the let  $x$  in  $e$  definien binder frame on the stack respectively.

We can reuse the definition of unification from the  $\lambda$ -ranking system because it turns out that we can use  $\mathcal{L}$  limit substitution with no changes except for using the let-ranks in place of the  $\lambda$ -ranks. Since instantiation during unification still shifts around the first occurrence in exactly the same way it did before, the limit substitution is still necessary in those cases where a univariable of low rank is instantiated with a type with univariables of high rank. Because it is easier for univariables to share the same let-rank than to share  $\lambda$ -ranks, we do not expect the limit substitution to apply as often as it did with the  $\lambda$ -ranking.

The [wl-let-body] rule decrements the let-depth, pops the let  $x$  in  $e_2$  frame off the stack and the letd frame off the environment, and pushes a let onto the stack. When type checking the body of a let-expression, that let will introduce no additional polymorphism because it already generalized the types for the definien. Thus, the [wl-let-body] pops out from the definien's scope of generalization. This is why [wl-let-body] decrements the let depth. In contrast to  $\mathcal{W}_\lambda$ , the type environment is augmented but the depth is decremented. The let frame on the stack is a reminder to pop off this new let-binding on the environment.

As in  $\mathcal{W}_\lambda$ , we generalize univariables of rank greater than the depth at the point of generalization where we have popped out of the definien (and thus decremented the depth). Any univariables that we introduced inside the scope of the definien must have rank greater than or equal to the index of the letd mark we just popped off. Rémy's let-ranking scheme generalizes only those univariables that have rank exactly equal to the current depth [13]. This is due to a difference between how his and our systems rank generic univariables. His system ranks generic univariables with the current depth instead of  $\infty$ .

Many of the lemmas about the construction of  $\tilde{\mathcal{R}}$  hold for  $\tilde{\mathcal{R}}_{\text{let}}$ .

**Lemma 10 (Rank is  $\infty$  if not in  $\Gamma$ )**

If  $\xi \notin \text{ftv}(\tilde{\Gamma}_\lambda)$ , then  $\tilde{\mathcal{R}}_{\text{let}}(\xi, \tilde{\Gamma}) = \infty$ .

**Proof (sketch):** By straightforward induction on  $\tilde{\Gamma}$  ■

**Lemma 11**

If  $\xi^m$  occurs in a  $\lambda$ -binding in  $\Gamma$  that is left bounded by the  $n$ th letd mark and  $\xi^m$  occurs in no  $\lambda$ -binding to the left of that letd, then  $\tilde{\mathcal{R}}_{\text{let}}(\xi^m, \Gamma) = n$ .

**Proof (sketch):** By induction on  $\Gamma$  and lem. 10 ■

The properties also carry over with minimal changes.

**Definition 12 (Properties on  $\mathcal{W}_L$  Machine States)**

The following are the properties on valid  $\mathcal{W}_L$  machine states  $(c, \Gamma, d, k)$ :

**Prop<sub>1</sub>** For all  $\xi^m \in \text{ftv}(\Gamma)$ .  $\mathcal{R}(\xi^m, \Gamma) = m$

**Prop<sub>2</sub>** For all  $\xi^m \in \text{ftv}(c) \cup \text{ftv}(k)$  such that  $m \leq \mathcal{D}_{\text{let}}(\Gamma)$ ,  $\xi^m \in \text{ftv}(\Gamma)$

**Prop<sub>3</sub>**  $\mathcal{D}_{\text{let}}(\Gamma) = \#$  of let  $x$  in  $e$ 's (let definien frames) on  $k$

**Prop<sub>4</sub>** For all  $\xi^m \in \text{ftv}(k)$  such that  $m \leq r$  where  $r$  is the number of let  $x$  in  $e$  frames up to the frame with the occurrence of  $\xi^m$  on the stack  $k$ ,  $\xi^m \in \text{ftv}(\Gamma)$ .

**Lemma 12 (Properties Invariant Under Unification)**

If the def. 12 hold for machine state valid  $\mathcal{W}_L$  machine state  $(\tau, \Gamma, d, k)$  and  $\theta = \mathcal{U}(\tau', \tau \rightarrow \zeta^\infty)$  and  $\tau'$  occurs in  $k$  and  $\zeta$  is fresh, then those properties hold for  $(\theta c, \Gamma, d, k)$ .

**Proof (sketch):** This proof is very similar to that of Lemma 6. As the example in Fig. 15 shows, the limited univariable case also plays out as it did for  $\mathcal{W}_\lambda$ . The difference is that we look for left bounding letd marks for  $\lambda$ -bindings with first occurrences instead of the  $\lambda$ -bindings themselves. ■

Again, the properties hold trivially for the initial  $\mathcal{W}_L$  states, hence we only have to show that they hold under all the transition rules.

**Lemma 13 (Properties Are Invariant For  $\mathcal{W}_L$ )**

For all valid  $\mathcal{W}_L$  machine states  $s$ , if the def. 12 properties hold for  $s$  and  $s \mapsto_{\mathcal{L}} s'$ , then the properties must hold for  $s'$ .

**Proof (sketch):** This proof is similar to lem. 7. Lemma 12 provides the proof for the [wl-app-out] case. The two proofs differ in that some of the work for the [ $\lambda$ -out] case is shifted to the [wl-let-body] since now it does the rank decrementing. The intuition is that the environment and the stack are synchronized. Furthermore, letd and let  $x$  in  $e$  (let definien frames) partition the environment

$$\text{letd}_d \lambda \dots \lambda \xi^d \dots \text{letd}_m \lambda \dots \lambda \phi^m$$

letd mark #

(a) Instantiating  $\xi^d$  to  $\phi^m$

$$\text{letd}_d \lambda \dots \lambda \phi^m \dots \text{letd}_m \lambda \dots \lambda \phi^m$$

letd mark #

(b)  $\phi^m$  now has a bad rank. The first occurrence of  $\phi^m$  is left bounded by the  $d$ th let.

$$\text{letd}_d \lambda \dots \lambda \phi^d \dots \text{letd}_m \lambda \dots \lambda \phi^d$$

letd mark #

(c) Limited rank of  $\phi$  to  $d$

**Figure 15.** The limit substitution for unification instantiation works similar to  $\mathcal{W}_\lambda$

and stack respectively into the definien scopes that determine generalization. ■

Having shown that the properties are invariants under  $\mathcal{W}_L$ , we can proceed to show that the generalizability criteria are equal at the point of generalization.

**Lemma 14 (Equivalence of Generalizability Criteria)**

If  $(\tau, \Gamma[\text{letd}], d, \text{let } x \text{ in } e :: k) \mapsto_L (e, \Gamma[x : \sigma], d - 1, \text{let} :: k)$  where  $\sigma = \forall \bar{\alpha}. \{\bar{\alpha} / \mathcal{G}_{d-1}(\Gamma)\} \tau$ , then  $\mathcal{G}_{d-1}(\Gamma) = \mathcal{G}_\Gamma(\Gamma)$ .

**Proof (sketch):** This goes through similar to how the analogous  $\mathcal{W}_\lambda$  lemma goes through. ■

We can now wrap up by showing the  $\mathcal{W}_L$  to be equivalent to the classic  $\mathcal{W}$  machine.

**Theorem 2 (Equivalence of  $\mathcal{W}_L$  Machine to Classic  $\mathcal{W}$  Machine)**

$(e, \emptyset, 0, \bullet) \mapsto_L^* (\tau, \emptyset, 0, \bullet)$  iff  $(e, \emptyset, \bullet) \mapsto_{\mathcal{W}}^* (\tau, \emptyset, \bullet)$

**Proof (sketch):** Except for the [let-body] rule and let depths that are just carried around in the rest of the machine, the  $\mathcal{W}_L$  machine is identical to the classic  $\mathcal{W}$  machine. Lemma 14 shows that even the [wl-let-body] and [let-body] rules are equivalent. ■

## 6. The Value Restriction

The previous sections have considered only a pure object language. In this section, we introduce references into the object language, which complicates inference. Using Algorithm  $\mathcal{W}$  without modification on a language with references admits the possibility of unsound generalization. The case in point is given in example 1. Although `ref nil` has a seemingly unconstrained type  $\xi$  **list ref**,  $\xi$  cannot be generalized because there is only one `ref` cell, which can only have one monomorphic type.

**Example 1**

```
let x = ref nil
in (x := cons unit nil; hd(!x) 3)
```

To eliminate this unsoundness, we introduce the value restriction, a particularly simple solution to type checking imperative programs. The tradeoff is that some perfectly sound pure programs no

longer type check under the value restriction. In particular, this precludes returning polymorphic data structures [3]. The value restriction eliminates the need for the usually heavy type machinery necessary to prevent generalization of only the imperative univariabes. An early approach, Tofte’s imperative variables [17], revealed the internal use of state in types and thus did not provide implementation abstraction. Wright showed that the value restriction did not unduly limit expressiveness in practice [18]. Since then, the revised version of the Standard ML language (SML 97) adopted the value restriction in place of the earlier, more complex mechanisms for inference with imperative type variables, and production compilers such as SML/NJ have implemented it.

In Fig. 16, we introduce references and the standard imperative primitives into the language. All the notation is taken from Standard ML. We distinguish between primitive operators and types by **sans-serif** and **bold face** respectively. Thus, `ref` is the operator and `ref` is the type. We also add the conventional list type, list primitives (`cons` and `nil`), and unit to the language to make it more convenient to construct examples that use two easily distinguishable types.

Fig. 17 introduces an extension of the  $\mathcal{W}$  abstract machine that directly implements the value restriction. Wright’s value restriction distinguished value `let`’s and non-value `let`’s, a distinction reminiscent of Leroy’s call-by-name and call-by-value `let`’s [6] albeit in a much more restrictive form. Our abstract machine also distinguishes between the two syntactic forms of `let` by having two sets of `let`-rules, one for the value `let`’s and the other for the non-value `let`’s. The abstract machine treats value `let`’s the same as the pure language `let`-ranking ( $\mathcal{W}_L$ ) machine did. To distinguish non-value `let`’s, we add a new stack frame of the form `letn x in e`. The rules [letv-def] and [letn-def] push on a value definien frame `let x in e` and a non-value definien `letn x in e` frame respectively when entering the definien of the respective kind of `let`. After we get a type for the definien, we only generalize this type if we see a normal (value) `let x in e` frame on top of the stack as per rule [letv-body]. Otherwise, we see a `letn x in e` and therefore do not generalize the definien type when adding a new binding for this definien on the environment in the [letn-body] rule. In example 1, [letn-def] will push into the outer `let` and [letn-body] will add the appropriate binding to the environment without generalizing  $\xi$  in the type of  $x$ .

Note that non-value `let`-bindings may contain first occurrences of univariabes. In this sense, non-value `let`-bindings are the same as  $\lambda$ -bindings in the environment. In example 2, if the type of the definien  $x$  is  $\xi$  **list ref**, then  $\xi$  only appears in the binding introduced the non-value `let` in the environment. Consequently, we do not distinguish between these two kinds of bindings. Both  $\lambda$ ’s and non-value `let`’s may put bindings that restrict generalization into the environment. They are for all intents and purposes, the same kind of binding. When the abstract machine reaches the inner `let` in example 2, the `letn` binding in the environment contains the univariable  $\xi$  ( $[x : \xi \text{ list ref}]$ ), hence  $\xi$  would not be erroneously generalized for the type of the definien  $y$ ,  $\xi$  **list ref**. In order to pop off the monomorphic bindings introduced by [letn-body], we use `letn` frames. The `letn` frames serve the same role as  $\lambda$  frames except in the service of non-value `let`’s.

**Example 2**

```
let x = ref nil
in let y = x
   in (y := cons unit nil)(y := [\lambda x.x])
```

$e ::= x \mid \lambda x.t \mid e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2 \mid \text{ref} \mid !$	Terms
$\tau ::= \alpha \mid \xi \mid \tau \rightarrow \tau \mid \tau \text{ list} \mid \tau \text{ ref} \mid \text{unit}$	Types
$v ::= x \mid \lambda x.e \mid \text{ref} \mid ! \mid \text{unit} \mid \text{cons} \mid \text{nil}$	Values
$n ::= e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2 \mid t_1 ::= t_2$	Non-values

Figure 16. Reference ML

<b>Machine Stack Frames</b>	
$f ::= \text{let } x \text{ in } e$	Value let definien
$\mid \text{let}_n x \text{ in } e$	Non-value let definien
$\mid \text{let}_n$	Non-value let body
$\mid \dots$	Frames from classical $\mathcal{W}$ machine
<b>Value Restriction Let Rules</b>	
$(\text{let } x = v \text{ in } e, \Gamma, k) \mapsto_{\lambda V} (v, \Gamma, \text{let } x \text{ in } e :: k)$	[letv-def]
$(\tau, \Gamma, \text{let } x \text{ in } e :: k) \mapsto_{\lambda V} (e, \Gamma[x : \sigma], \text{let} :: k)$	[letv-body]
$\sigma = \forall \bar{\alpha}. \{\bar{\alpha} / \mathcal{G}_d(\tau)\} \tau$ where $\bar{\alpha}$ is fresh	
$(\text{let } x = n \text{ in } e', \Gamma, k) \mapsto_{\lambda V} (n, \Gamma, \text{let}_n x \text{ in } e' :: k)$	[letn-def]
$(\tau, \Gamma, \text{let}_n x \text{ in } e' :: k) \mapsto_{\lambda V} (e', \Gamma[x : \tau], \text{let}_n :: k)$	[letn-body]
$(\tau, \Gamma[x : \tau'], \text{let}_n :: k) \mapsto_{\lambda V} (\tau, \Gamma, k)$	[letn-out]
<b>Imperative Features</b>	
$(\text{ref } e, \Gamma, k) \mapsto_{\lambda V} (e, \Gamma, \text{ref} :: k)$	[ref-in]
$(\tau, \Gamma, \text{ref} :: k) \mapsto_{\lambda V} (\tau \text{ ref}, \Gamma, k)$	[ref-out]
$(!e, \Gamma, k) \mapsto_{\lambda V} (e, \Gamma, ! :: k)$	[deref-in]
$(\tau, \Gamma, ! :: k) \mapsto_{\lambda V} (\theta \xi, \theta \Gamma, \theta k)$	[deref-out]
$\xi$ is fresh $\theta = \mathcal{U}(\tau, \xi \text{ ref})$	
$(x := e, \Gamma, k) \mapsto_{\lambda V} (e, \Gamma, x := \square :: k)$	[set-in]
$(\tau, \Gamma, x := \square :: k) \mapsto_{\lambda V} (\text{unit}, \theta \Gamma, \theta k)$	[set-out]
$\theta = \mathcal{U}(\Gamma(x), \tau \text{ ref})$	
<b>List Rules</b>	
$(\text{nil}, \Gamma, k) \mapsto_{\lambda V} (\xi \text{ list}, \Gamma, k)$	[nil]
$(\text{cons}, \Gamma, k) \mapsto_{\lambda V} (\xi \text{ list} \rightarrow \xi \rightarrow \xi \text{ list}, \Gamma, k)$	[cons]
$\xi$ is fresh	

All the rules other than [let-def] and [let-body] carry over from the classical  $\mathcal{W}$  machine (Fig. 4).

Figure 17. Classical  $\mathcal{W}$  with References and Value Restriction

$e ::= x \mid \lambda x.e \mid e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2 \mid \text{ref} \mid !$	Terms
$\mid e_1 := e_2 \mid \text{cons} \mid \text{nil} \mid \text{unit}$	
$\tau ::= \alpha \mid \xi^d \mid \tau \rightarrow \tau \mid \tau \text{ list} \mid \tau \text{ ref} \mid \text{unit}$	Types
$v ::= x \mid \lambda x.e \mid \text{ref} \mid ! \mid \text{unit} \mid \text{cons} \mid \text{nil}$	Values
$n ::= e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2 \mid e_1 := e_2$	Non-values

Figure 18. Reference ML with Ranked Univariabes

## 7. Value Restriction with $\lambda$ -Ranking

We can also adapt the efficient ranking schemes to work with the value restriction. Because both  $\lambda$  and non-value let-introduced bindings may restrict generalization, our notion of nesting depth must count them both. Fortunately, because we can introduce monomorphic bindings into the environment for non-value let's, we can count them as  $\lambda$ -bindings and therefore reuse  $\mathcal{D}$  and  $\mathcal{R}$  from the  $\mathcal{W}_\lambda$  machine. The basic idea is to treat non-value let bodies as we do  $\lambda$  binding bodies in the foregoing machines and to distinguish non-value let definien so that we do not confuse them with value let definien whose types *do* generalize. However, there is one last complication in that non-value let definien may leak out univariabes that might be unsoundly generalized in the body. We deal with this situation by limiting the ranks of non-value definien types.

<b>Value Restriction Let Rules</b>	
$(\text{let } x = v \text{ in } e, \Gamma, d, k) \mapsto_{\lambda V} (v, \Gamma, d, \text{let } x \text{ in } e :: k)$	[ $\lambda v$ -letv-def]
$(\tau, \Gamma, d, \text{let } x \text{ in } e :: k) \mapsto_{\lambda V} (e, \Gamma[x : \sigma], d, \text{let} :: k)$	[ $\lambda v$ -letv-body]
$\sigma = \forall \bar{\alpha}. \{\bar{\alpha} / \mathcal{G}_d(\tau)\} \tau$ where $\bar{\alpha}$ is fresh	
$(\text{let } x = n \text{ in } e', \Gamma, d, k) \mapsto_{\lambda V} (n, \Gamma, d, \text{let}_n x \text{ in } e' :: k)$	[ $\lambda v$ -letn-def]
$(\tau, \Gamma, d, \text{let}_n x \text{ in } e' :: k) \mapsto_{\lambda V} (e', \Gamma[x : \tau'], d', \text{let}_n :: k')$	[ $\lambda v$ -letn-body]
$d' = d + 1 \quad \theta = \mathcal{L}_d(\tau) \quad \Gamma' = \theta(\Gamma) \quad \tau' = \theta(\tau)$	
$k' = \theta k$	
$(\tau, \Gamma[x : \tau'], d, \text{let}_n :: k) \mapsto_{\lambda V} (\tau, \Gamma, d - 1, k)$	[ $\lambda v$ -letn-out]
<b>Imperative Features</b>	
$(\text{ref } e, \Gamma, d, k) \mapsto_{\lambda V} (e, \Gamma, d, \text{ref} :: k)$	[ $\lambda v$ -ref-in]
$(\tau, \Gamma, d, \text{ref} :: k) \mapsto_{\lambda V} (\tau \text{ ref}, \Gamma, d, k)$	[ $\lambda v$ -ref-out]
$(!e, \Gamma, d, k) \mapsto_{\lambda V} (e, \Gamma, d, ! :: k)$	[ $\lambda v$ -deref-in]
$(\tau, \Gamma, d, ! :: k) \mapsto_{\lambda V} (\theta \xi^\infty, \theta \Gamma, d, \theta k)$	[ $\lambda v$ -deref-out]
$\xi^\infty$ is fresh $\theta = \mathcal{U}(\tau, \xi^\infty \text{ ref})$	
$(x := e, \Gamma, d, k) \mapsto_{\lambda V} (e, \Gamma, d, x := \square :: k)$	[ $\lambda v$ -set-in]
$(\tau, \Gamma, d, x := \square :: k) \mapsto_{\lambda V} (\text{unit}, \theta \Gamma, d, \theta k)$	[ $\lambda v$ -set-out]
$\theta = \mathcal{U}(\Gamma(x), \tau \text{ ref})$	
<b>List Rules</b>	
$(\text{nil}, \Gamma, d, k) \mapsto_{\lambda V} (\xi^\infty \text{ list}, \Gamma, d, k)$	[ $\lambda v$ -nil]
$(\text{cons}, \Gamma, d, k) \mapsto_{\lambda V} (\xi^\infty \text{ list} \rightarrow \xi^\infty \rightarrow \xi^\infty \text{ list}, \Gamma, d, k)$	[ $\lambda v$ -cons]
$\xi^\infty$ is fresh	

All the other rules, machine state, initial state, registers,  $\mathcal{G}_d$ , and  $\mathcal{I}$  carry over from  $\mathcal{W}_\lambda$ .

Figure 19.  $\lambda$ -Ranking  $\mathcal{W}$  with References and Value Restriction ( $\mathcal{W}_{\lambda V}$ )

$\lambda$ -ranking does nothing special with value let's. The rules look the same as in the classical  $\mathcal{W}$  with value restriction except now they carry around the  $\lambda$  depth. In contrast, there is a slight complication for non-value let's. If we leave high ranked univariabes in the non-value let definien alone, we admit the possibility that such contextual univariabes may be generalized at a later point. In example 2, although the we did not generalize the type of  $x$ ,  $\xi \text{ list ref}$ , at the outer let, a non-value let, it may erroneously generalize at the inner let, a value let. The [ $\lambda v$ -letn-body] avoids this unsound generalization by limiting the ranks in the non-value definien types. This measure ensures that no high rank univariable would be generalized in the body of the non-value let. For example 2, the [ $\lambda v$ -letn-body] rule would limit the rank of  $\xi$  to 0 where the type of  $x$  is  $\xi^1 \text{ list ref}$ , even though [ $\lambda v$ -letn-def] previously incremented the  $\lambda$ -depth to 1. When the machine reaches the inner let, the rank of  $\xi$  is too low to generalize. Using the intuition that non-value let  $x = e_1$  in  $e_2$  act like expressions of the form  $(\lambda x.e_2)e_1$ , [ $\lambda v$ -let-body] increments the depth when entering the body (and extending the environment with a monomorphic ( $\lambda$ -) binding). The [ $\lambda v$ -let-out] decrements the depth because it pops off that monomorphic binding from the environment when popping out from the body.

In the remaining rules, we simply add the depth register to the analogous classical  $\mathcal{W}$  with value restriction rules. In [ $\lambda v$ -deref-out] and [ $\lambda v$ -nil], we introduce fresh generic univariabes that should therefore be ranked  $\infty$ .

To account for the non-value let-introduced bindings in the environment, Prop<sub>3</sub> and Prop<sub>4</sub> must count both  $\lambda$  and let<sub>n</sub> frames on the stack because nesting depth now includes both. Def. 13 shows the Prop<sub>3</sub> and Prop<sub>4</sub> updated to reflect this. This is another reason why we need [ $\lambda v$ -letn-body] to introduce let<sub>n</sub> frames as distinct from regular let frames.

### Definition 13 (Properties for Value Restriction Machines)

- Prop<sub>1</sub>**  $\forall \xi^m \in \text{ftv}(\Gamma), \mathcal{R}(\xi^m, \Gamma) = m$   
**Prop<sub>2</sub>**  $\forall \xi^m \in \text{ftv}(c)$  such that  $m \leq \mathcal{D}(\Gamma), \xi^m \in \text{ftv}(\Gamma)$   
**Prop<sub>3</sub>**  $\mathcal{D}(\Gamma) = \#$  of  $\lambda$ 's and  $\text{let}_n$ 's on  $k$   
**Prop<sub>4</sub>**  $\forall \xi^m \in \text{ftv}(k)$  such that  $m \leq r$  where  $r$  is the number of  $\lambda$  and  $\text{let}_n$  frames up to the frame with the occurrence of  $\xi^m$  on the stack  $k, \xi^m \in \text{ftv}(\Gamma)$ .

We now have three cases where we use unification ([ $\lambda v$ -deref-out], [ $\lambda v$ -set-out], and the app-out rule), hence the properties must be invariant under all of them.

### Lemma 15 (Invariants under unification)

If the properties hold for  $s = (\tau, \Gamma, d, (\tau' \square) :: k), s = (\tau, \Gamma, d, ! :: k)$ , or  $s = (\tau, \Gamma, d, x := \square :: k) \in \mathcal{VM}S$  and  $\theta = \mathcal{U}(\tau', \tau \rightarrow \psi^\infty), \mathcal{U}(\tau, \psi^\infty)$ , or  $\mathcal{U}(\Gamma(x), \tau \text{ ref})$  respectively, then they must also hold for  $s' = (\theta \psi^\infty, \theta \Gamma, d, \theta k)$  for the first two cases or  $s' = (\text{unit}, \theta \Gamma, d, \theta k)$ .

**Proof (sketch):** The definition of  $\mathcal{U}$  has not changed at all. The other cases for  $s$  and  $\theta$  are simple adaptations of the proof for lem. 6. They are actually simpler than the general [app-out] case. ■

Next, we need to show that the properties are invariant under the new rules, especially the non-value let rules.

### Lemma 16 (Properties invariant over all $\mathcal{W}_{\lambda v}$ rules)

If the properties hold for  $\mathcal{W}_{\lambda v}$  state  $s \in \mathcal{VM}S$  and  $s \mapsto_{\lambda v} s'$ , then they hold for  $s'$ .

**Proof (sketch):** Again, the [app-out] case is handled by lem. 15. All rules that were present in  $\mathcal{W}_\lambda$  are handled in the same way. ■

At this point, we have everything we need to prove the analogs of Corollary 1, lemmas 8, and 9. The proofs are identical to those in  $\mathcal{W}_\lambda$ . We are now ready to connect  $\mathcal{W}_{\lambda v}$  to  $\mathcal{W}_v$ .

### Theorem 3 (Equivalence of $\mathcal{W}_{\lambda v}$ and $\mathcal{W}_v$ )

$(t, \emptyset, 0, \bullet) \mapsto_{\lambda v}^* (\tau, \emptyset, 0, \bullet)$  iff  $(t, \emptyset, \bullet) \mapsto_v^* (\tau, \emptyset, \bullet)$

**Proof (sketch):** This proof proceeds in the same way as thm. 1.  $\mathcal{W}_{\lambda v}$  and  $\mathcal{W}_v$  are identical except for the generalization criteria and innocuous carrying around of ranks. ■

## 8. Value Restriction with Let-Ranking

We can carry over much of the  $\mathcal{W}_l$  machine, so Fig. 20 only lists the new or changed rules in the  $\mathcal{W}$  with value restriction and let-ranking machine,  $\mathcal{W}_{lv}$ . The rules for let expressions are now also responsible for letd marking. However, it turns out that only the rules for value let's do the letd marking just as they did in  $\mathcal{W}_l$  and the rules for non-value let's are responsible for the definien type limitation as they were in  $\mathcal{W}_{\lambda v}$ . Thus, for value let's, the let depth is incremented only while inside the definien ([ $lv$ -letv-def]) and it is decremented once we start the body ([ $lv$ -letv-body]) just as it is in  $\mathcal{W}_l$ . Non-value let's, however, do not increment the let-depth because non-value let definien introduce no scope of generalization. A non-value let  $x = n$  in  $e$  is treated in the same way as would  $(\lambda x.e)n$ . The [ $lv$ -letn-def] rule enters the definien after pushes on a  $\text{let}_n x$  in  $e$  frame onto the stack. Once we get a type for the definien, the [ $lv$ -letn-body] pops off the  $\text{let}_n x$  in  $e$  in order to start the body under the environment extended with the limited but monomorphic definien type. We still need to limit the definien type ranks for the same reasons as in the previous section.

$(\text{let } x = v \text{ in } e, \Gamma, d, k)$	$\mapsto_{LV} (v, \Gamma[\text{letd}], d + 1, \text{let } x \text{ in } e :: k)$	[ $lv$ -letv-def]
$(\tau, \Gamma[\text{letd}], d, \text{let } x \text{ in } e :: k)$	$\mapsto_{LV} (e, \Gamma[x : \sigma], d - 1, \text{let} :: k)$	[ $lv$ -letv-body]
	$\sigma = \forall \bar{\alpha}. \{\bar{\alpha} / \mathcal{G}_{d-1}(\tau)\} \tau$ where $\bar{\alpha}$ is fresh	
$(\text{let } x = n \text{ in } e, \Gamma, d, k)$	$\mapsto_{LV} (n, \Gamma, d, \text{let}_n x \text{ in } e :: k)$	[ $lv$ -letn-def]
$(\tau, \Gamma, d, \text{let}_n x \text{ in } e :: k)$	$\mapsto_{LV} (e, (\theta \Gamma)[x : \theta(\tau)], d, \text{let}_n :: \theta k)$	[ $lv$ -letn-body]
	$\theta = \mathcal{L}_d(\tau)$	
$(\tau, \Gamma[x : \tau'], d, \text{let}_n :: k)$	$\mapsto_{LV} (\tau, \Gamma, d, k)$	[ $lv$ -letn-out]

All imperative features, list rules, and generic instantiation  $\mathcal{I}$  carry over from the  $\mathcal{W}_{\lambda v}$ , Fig. 19. All the remaining rules, machine state, and initial state carry over from  $\mathcal{W}_l$ , Fig. 14.

**Figure 20.** Let-Ranking  $\mathcal{W}$  with References and Value Restriction ( $\mathcal{W}_{lv}$ )

Non-value let's also do not add letd definien marks to the environment. Everything in the non-value let definien is no more generalizable than anything outside of that non-value let. Thus, the depth ranking here refers to the value let nesting depth only. Non-value let's are ignored just as  $\lambda$ 's are with respect to determining the depth. So, unlike in  $\mathcal{W}_{\lambda v}$ , [ $lv$ -letn-body] and [ $lv$ -letn-out] do not change the current depth. The [ $lv$ -letn-out] is relegated to the small role of popping off the monomorphic definien binding once we are done with the body.

Because non-value let's do not affect the let-ranking scheme, we can reuse the properties in def. 12 ( $\mathcal{W}_l$ ) and consequently all the lemmas for the  $\mathcal{W}_l$  machine. The new cases to the proofs are straightforward because everything, including the non-value let's can be encoded in terms of the  $\mathcal{W}_l$  language.

### Theorem 4 (Equivalence of $\mathcal{W}_{lv}$ and $\mathcal{W}_v$ )

$(e, \emptyset, 0, \bullet) \mapsto_{LV}^* (\tau, \emptyset, 0, \bullet)$  iff  $(e, \emptyset, \bullet) \mapsto_v^* (\tau, \emptyset, \bullet)$ .

**Proof (sketch):** This proof follows the form of the proof for  $\mathcal{W}_l$ . The new cases are straightforward. ■

## 9. Related Work

Most of the literature on Algorithm  $\mathcal{W}$  and the value restriction fall into two categories. The first category focuses on general frameworks for understanding  $\mathcal{W}$  and the value restriction. Rémy [13] formalizes a ranked variable system using ranks based on let-depth. Pottier and Rémy [12] offer a model implementation using the let-ranking technique. As mentioned earlier, Rémy's let-ranking differs somewhat from ours when it instantiates generic univariates to the current let-depth rather than rank  $\infty$ . He also proves a principal typing result. The work relies on a complex theory based on unificands, thus providing a very general framework that can be extended to other equational theories. It does not address the value restriction. It is not clear whether this development is equivalent to our  $\mathcal{W}_l$  abstract machine, but we conjecture that they are similar. In contrast to Rémy's work, this paper is concerned primarily about the engineering of inference algorithms using ranked type variables. This is reflected by the presentations of the machines that more closely follow implementations.

Pottier [11] and Skalka and Pottier [15] describe a variant of the Hindley-Milner type system (HM(X)) with the value restriction. They prove its soundness. Wright provides empirical evidence that the value restriction is practical [18]. His work relies on the pure Algorithm  $\mathcal{W}$  for inference after translating non-value let's away. Thus, implicitly, soundness for the value restriction is obtained by treating non-value let's as equivalent to corresponding beta-

redexes, where the bindings become  $\lambda$  bindings (e.g., let  $x = n_1$  in  $e_2$  translates to  $(\lambda x.e_2)n_1$ ). Pottier, Skalka and Pottier, and Wright do not address efficient algorithms for inference.

The second category focuses on expanding the set of well-typed programs by either refining the value restriction or by developing a type discipline for imperative type variables, neither of which are considered in this paper. Garrigue uses a subtyping relation to relax the value restriction and thereby expand the set of safe, well-typed programs while simultaneously preserving the implementation abstraction property of the value restriction [3]. His technique types some programs involving the handling of polymorphic data that would not type check under the ordinary value restriction. Tofte provides a soundness proof for the imperative type system of SML 90 [17]. The soundness of the value restriction is the limiting case of Tofte’s proof where all type variables are imperative type variables. There are also a number of other contributions that offer more powerful and complex alternatives to the value restriction [6, 7, 16]. These treatments do not address using type variable ranks to make generalization more efficient.

## 10. Future Work

Both the  $\mathcal{W}_{\lambda V}$  and  $\mathcal{W}_{LV}$  machines use a rank-limiting coercion on the univariates of the let-definien type to prevent unsound generalization of definien univariates in the body. However, for efficiency, we can avoid computing and applying this limit substitution if we introduced generic univariates with more precise initial ranks.

We believe that there is an alternative to this rank coercion approach. Another version of the value restriction machines can initialize generic type variables to the current depth rather than infinity. This setup would essentially pre-limit the ranks of generic univariates in a non-value let definien. In particular, we observe that the only requirement on the rank given to a generic univariate is that the rank must be greater than or equal to  $\mathcal{D}_{\text{let}}(\Gamma)$ . In fact,  $\mathcal{D}_{\text{let}}(\Gamma)$  is a perfectly sensible rank for generic univariates. The [var] rule can use a modified generic instantiation that assigns the current let depth as the rank for the new generic univariates. Similarly, the [nil] rule introduces a univariate ranked at the current depth. This technique may be related to Rémy’s variant of the let-ranking scheme. We plan to investigate these alternative machines and their connection to Rémy’s technique.

The abstract machine approach to formalizing inference can also be extended to other type systems and inference approaches. Local type inference offers partial inference for a variant of System  $F_{\leq}$  [10]. One alternative approach to inference for the Hindley-Milner type system is Lee and Yi’s formalization of the folklore top-down alternative to  $\mathcal{W}$ , algorithm  $\mathcal{M}$  [5]. The abstract machine approach may offer some further insight into these inference algorithms.

## 11. Conclusion

We provide a direct and accessible approach to formalizing the various versions of type inference with ranked variables. This particular formalism is also very close to implementation. Using  $\mathcal{W}_{LV}$  as a reference, we implemented a let-ranked machine for type checking programs under value restriction semantics. This implementation can be easily adapted to follow the other machines. This formalism also exposes some interesting properties of the type environment that we rely upon for use with the various ranking techniques.

Using this abstract machine type inference formalism, we have proved these systems correct with respect to the original Algorithm  $\mathcal{W}$  (for a pure version of ML), and with respect to a version with the value restriction for the more realistic impure ML. From this experience, we believe that the formulation of type inference

algorithms as abstract machines is an effective basis for reasoning about such systems and relating them to each other.

## References

- [1] Luís Damas. unpublished note, 1984.
- [2] Matthias Felleisen and Robert Hieb. A revised report on the syntactic theories of sequential control and state. *Theoretical Computer Science*, 103(2):235–271, 1992.
- [3] Jacques Garrigue. Relaxing the value restriction. *Functional and Logic Programming*, 2004.
- [4] George Kuan, David MacQueen, and Robert Bruce Findler. A rewriting semantics for type inference. In Rocco De Nicola, editor, *Programming Languages and Systems, 16th European Symposium on Programming, ESOP 2007*, volume 4421, March 2007.
- [5] Oukseh Lee and Kwangkeun Yi. Proofs about a folklore let-polymorphic type inference algorithm. *ACM Transactions on Programming Languages and Systems*, 20(4):707–723, July 1998.
- [6] Xavier Leroy. Polymorphism by name for references and continuations. In *POPL ’93: Proceedings of the 20th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 220–231, New York, NY, USA, 1993. ACM Press.
- [7] Xavier Leroy and Pierre Weis. Polymorphic type inference and assignment. In *POPL ’91: Proceedings of the 18th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 291–302, New York, NY, USA, 1991. ACM Press.
- [8] Harry G. Mairson. Deciding ml typability is complete for deterministic exponential time. In *POPL*, pages 382–401, 1990.
- [9] Robin Milner. A theory of type polymorphism in programming. *JCSS*, 17:348–375, 1978.
- [10] Benjamin C. Pierce and David N. Turner. Local type inference. In *POPL ’98: Proceedings of the 25th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 252–265, New York, NY, USA, 1998. ACM Press.
- [11] François Pottier. A semi-syntactic soundness proof for HM(X). Technical Report RR-4150, INRIA, 2001.
- [12] François Pottier and Didier Rémy. The essence of ML type inference. In Benjamin C. Pierce, editor, *Advanced Topics in Types and Programming Languages*, chapter 10, pages 389–489. MIT Press, 2005.
- [13] Didier Rémy. Extending ML type system with a sorted equational theory. Research Report 1766, Institut National de Recherche en Informatique et Automatique, Rocquencourt, BP 105, 78 153 Le Chesnay Cedex, France, 1992.
- [14] John C. Reynolds. Definitional interpreters for higher-order programming languages. In *ACM ’72: Proceedings of the ACM annual conference*, pages 717–740, New York, NY, USA, 1972. ACM Press.
- [15] Christian Skalka and François Pottier. Syntactic type soundness for HM(X). In *Proceedings of the Workshop on Types in Programming (TIP’02)*, volume 75 of *Electronic Notes in Theoretical Computer Science*, Dagstuhl, Germany, July 2002.
- [16] Geoffrey Smith and Dennis Volpano. Polymorphic typing of variables and references. *ACM Trans. Program. Lang. Syst.*, 18(3):254–267, 1996.
- [17] Mads Tofte. Type inference for polymorphic references. *Inf. Comput.*, 89(1):1–34, 1990.
- [18] Andrew K. Wright. Simple imperative polymorphism. *Lisp Symp. Comput.*, 8(4):343–355, 1995.