

# Untangling Typechecking of Intersections and Unions

Joshua Dunfield\*

School of Computer Science, McGill University  
Montréal, Canada  
joshua@mpi-sws.org

Intersection and union types denote conjunctions and disjunctions of properties. Using bidirectional typechecking, intersection types are relatively straightforward, but union types present challenges. For union types, we can case-analyze a subterm of union type when it appears in evaluation position (replacing the subterm with a variable, and checking that term twice under appropriate assumptions). This technique preserves soundness in a call-by-value semantics. Sadly, there are so many choices of subterms that a direct implementation is not practical. But carefully transforming programs into let-normal form drastically reduces the number of choices. The key results are soundness and completeness: a typing derivation (in the system with too many subterm choices) exists for a program if and only if a derivation exists for the let-normalized program.

## 1 Introduction

To check programs in advanced type systems, it can be useful to split the traditional typing judgment  $e : A$  into two forms,  $e \uparrow A$  read “ $e$  synthesizes type  $A$ ” and  $e \downarrow A$  read “ $e$  checks against type  $A$ ”, and requiring that the user write annotations on redexes. This *bidirectional typechecking* (Pierce and Turner 1998) is decidable for many interesting features, including intersection and union types without syntactic markers. *Tridirectional* typechecking (Dunfield and Pfenning 2004; Dunfield 2007b) is essentially bidirectional, but union types are eliminated with the aid of a *tridirectional rule* that uses an evaluation context  $\mathcal{E}$ :

$$\frac{\Gamma; \Delta_1 \vdash e' \uparrow A \quad \Gamma; \Delta_2, \bar{x}:A \vdash \mathcal{E}[\bar{x}] \downarrow C}{\Gamma; \Delta_1, \Delta_2 \vdash \mathcal{E}[e'] \downarrow C} \text{direct}\mathbb{L}$$

In this rule,  $\Gamma$  is an ordinary variable context and  $\Delta_1, \Delta_2$  is the concatenation of *linear contexts*; linear variables  $\bar{x}$  in  $\Delta$ s essentially stand for subterms (occurrences) in the subject.  $\text{direct}\mathbb{L}$  gives  $e'$  a (linear) name  $\bar{x}$ , so that a left rule, which decomposes types in the context  $\Delta$ , can eliminate union types appearing in  $A$ . Instead of a direct union elimination rule like  $\vee E$ , we use  $\text{direct}\mathbb{L}$  together with a left rule  $\vee\mathbb{L}$ .

$$\frac{\Gamma \vdash e' \uparrow A \vee B \quad \Gamma, x:A \vdash \mathcal{E}[x] \downarrow C \quad \Gamma, y:B \vdash \mathcal{E}[y] \downarrow C}{\Gamma \vdash \mathcal{E}[e'] \downarrow C} [\vee E] \quad \frac{\Gamma; \Delta, \bar{x}:A \vdash e \downarrow C \quad \Gamma; \Delta, \bar{x}:B \vdash e \downarrow C}{\Gamma; \Delta, \bar{x}:A \vee B \vdash e \downarrow C} \vee\mathbb{L}$$

While evaluation contexts are defined syntactically, this rule is not syntax-directed in the usual sense: many terms have more than one decomposition into some  $\mathcal{E}[e']$  where the subterm  $e'$  can synthesize a type. Under a left-to-right (functions first, arguments second) call-by-value semantics, even  $f x$  has three decompositions  $\mathcal{E} = []$ ,  $\mathcal{E} = [] x$ ,  $\mathcal{E} = f []$ , so a straightforward implementation of a system with  $\text{direct}\mathbb{L}$  would require far too much backtracking. Compounded with backtracking due to intersection and union types (e.g., if  $f : (A_1 \rightarrow A_2) \wedge (B_1 \rightarrow B_2)$  we may have to try both  $f : A_1 \rightarrow A_2$  and  $f : B_1 \rightarrow B_2$ ), such an implementation would be hopelessly impractical.

---

\*Current address: MPI-SWS, Gottlieb-Daimler-Str. 49, 67663 Kaiserslautern, Germany.

This paper reformulates tridirectional typechecking (summarized in Section 2) to work on terms in a particular *let-normal form*, in which steps of computation are sequenced and intermediate computations are named. The let-normal transformation (Section 3) drastically constrains the decomposition by sequencing terms, forcing typechecking to proceed left to right (with an interesting exception). The results stated in Section 4 guarantee that the let-normal version of a program  $e$  is well typed under the let-normal version of the type system if and only if  $e$  is well typed under the tridirectional system.

The let-normal transformation itself is not complicated, though the motivation for my particular formulation is somewhat involved. The details of the transformation may be of interest to designers of advanced type systems, whether their need for a sequential form arises from typechecking itself (as in this case) or from issues related to compilation.

Unfortunately, the proofs (especially the proof of completeness) are very involved; I couldn't even fit all the statements of lemmas in this paper, much less sketch their proofs. I hope only to convey a shadow of the argument's structure.

This paper distills part of my dissertation (Dunfield 2007b, chapter 5). To simplify presentation, I omit tuples, datasort refinements, indexed types (along with universal and existential quantification, guarded types, and asserting types), and a greatest type  $\top$ .

## 2 Tridirectional Typechecking

We have functions, products, intersections, unions, and an empty type  $\perp$ . We'll use a unit type and other base types like `int` in examples. In the terms  $e$ , we have variables  $x$  (which are values) bound by  $\lambda x. e$ , variables  $u$  (not values) bound by  $\text{fix } u. e$ , and call-by-value application  $e_1 e_2$ . Note the lack of syntactic markers for intersections or unions. As usual,  $\mathcal{E}[e']$  is the evaluation context  $\mathcal{E}$  with its hole replaced by  $e'$ . To replace  $x$  with  $e_1$ , we write  $[e_1/x]e_2$ : “ $e_1$  for  $x$  in  $e_2$ ”.

Types  $A, B, C, D ::= A \rightarrow B \mid A \wedge B \mid A \vee B \mid \perp$   
 Terms  $e ::= x \mid u \mid \lambda x. e \mid e_1 e_2 \mid \text{fix } u. e$   
 Values  $v ::= x \mid \lambda x. e$   
 Evaluation contexts  $\mathcal{E} ::= [] \mid \mathcal{E} e \mid v \mathcal{E}$

Small-step reduction rules

$\mathcal{E}[(\lambda x. e) v] \mapsto \mathcal{E}[[v/x]e]$   
 $\mathcal{E}[\text{fix } u. e] \mapsto \mathcal{E}[[\text{fix } u. e] / u] e]$

We'll start by looking at the “left tridirectional” (in this paper, called just “tridirectional”) system. This system was presented in Dunfield and Pfenning (2004) and Dunfield (2007b, chapter 4); space allows only a cursory description.

The subtyping judgment (Figure 1) is  $A \leq B$ . Transitivity is admissible.  $\wedge$  does not distribute across  $\rightarrow$ , for reasons explained by Davies and Pfenning (2000).

Figure 2 gives the typing rules. The judgment  $\Gamma; \Delta \vdash e \uparrow A$  is read “ $e$  synthesizes type  $A$ ”, and  $\Gamma; \Delta \vdash e \Downarrow A$  is read “ $e$  checks against  $A$ ”. When synthesizing,  $A$  is output; when checking,  $A$  is input.

Contexts  $\Gamma ::= \cdot \mid \Gamma, x:A$   
 Linear contexts  $\Delta ::= \cdot \mid \Delta, \bar{x}:A$

Contexts  $\Gamma$  have regular variable declarations. Linear contexts  $\Delta ::= \cdot \mid \Delta, \bar{x}:A$  have linear variables. If  $\Gamma; \Delta \vdash e \dots$  is derivable, then “ $\Delta \Vdash e$  ok”, read “ $e$  OK under  $\Delta$ ”: each  $\bar{x}$  declared in  $\Delta$  appears exactly once in  $e$ , and  $e$  contains no other linear variables. Rules that decompose the subject, such as  $\rightarrow E$  decomposing  $e_1 e_2$  into  $e_1$  and  $e_2$ , likewise decompose  $\Delta$ .

Most rules follow a formula devised in Dunfield and Pfenning (2004): introduction rules, such as  $\rightarrow I$ , check; elimination rules, such as  $\rightarrow E$ , synthesize. Introduction forms like  $\lambda x. e$  thus construct

*synthesizing terms*, while elimination forms like  $e_1 e_2$  are *checked terms*. Some rules fall outside this classification. The *assumption* rules  $\text{var}$ ,  $\overline{\text{var}}$  and  $\text{fixvar}$  synthesize (an assumption  $x:A$  can be read  $x \uparrow A$ ). The *subsumption* rule  $\text{sub}$  allows a term that synthesizes  $A$  to check against a type  $B$ , provided  $A$  is a subtype of  $B$ . The rule  $\text{ctx-anno}$  permits *contextual typing annotations*; for example, in  $(\text{succ } x : (x:\text{odd} \vdash \text{even}, x:\text{even} \vdash \text{odd}))$ , the annotated term  $\text{succ } x$  is checked against  $\text{even}$  if  $x:\text{odd} \in \Gamma$ , and against  $\text{odd}$  if  $x:\text{even} \in \Gamma$ . The premise  $(\Gamma_0 \vdash A) \lesssim (\Gamma \vdash A)$  is derivable if the assumptions in  $\Gamma$  support the assumptions in  $\Gamma_0$ . For details, see Dunfield and Pfenning (2004).

Finally, we have *left rules*  $\wedge\mathbb{L}_1$ ,  $\wedge\mathbb{L}_2$ ,  $\vee\mathbb{L}$ ,  $\perp\mathbb{L}$  which act on linear assumptions  $\bar{x}:A$  where  $A$  is of intersection, union, or empty type. These act as elimination rules—for  $\vee$  and  $\perp$ , they are the *only* elimination rules.  $\wedge\mathbb{L}_1$  and  $\wedge\mathbb{L}_2$  are not useful alone (the ordinary eliminations  $\wedge\mathbb{E}_1$  and  $\wedge\mathbb{E}_2$  would do) but are needed to expose a nested  $\vee$  for  $\vee\mathbb{L}$ , or a  $\perp$  for  $\perp\mathbb{L}$ .

The backtracking required to choose between  $\wedge\mathbb{E}_1$  and  $\wedge\mathbb{E}_2$ , or between  $\vee\mathbb{I}_1$  and  $\vee\mathbb{I}_2$ , or between the related subtyping rules, as well as the need to check a single term more than once ( $\wedge\mathbb{I}$ ,  $\vee\mathbb{L}$ ) suggests that typechecking is exponential. In fact, Reynolds (1996, pp. 67–68) proved that for a closely related system, typechecking is PSPACE-hard. We can’t make typechecking polynomial, but “untangling”  $\text{direct}\mathbb{L}$  will remove one additional source of complexity.

$$\boxed{A \leq B} \quad \frac{B_1 \leq A_1 \quad A_2 \leq B_2}{A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2} \rightarrow \leq \quad \frac{A_1 \leq B}{A_1 \wedge A_2 \leq B} \wedge\mathbb{L}_1 \leq \quad \frac{A_2 \leq B}{A_1 \wedge A_2 \leq B} \wedge\mathbb{L}_2 \leq \quad \frac{A \leq B_1 \quad A \leq B_2}{A \leq B_1 \wedge B_2} \wedge\mathbb{R} \leq$$

$$\frac{}{\perp \leq A} \perp\mathbb{L} \leq \quad \frac{A_1 \leq B \quad A_2 \leq B}{A_1 \vee A_2 \leq B} \vee\mathbb{L} \leq \quad \frac{A \leq B_1}{A \leq B_1 \vee B_2} \vee\mathbb{R}_1 \leq \quad \frac{A \leq B_2}{A \leq B_1 \vee B_2} \vee\mathbb{R}_2 \leq$$

Figure 1: Subtyping

$$\frac{\Gamma(x) = A}{\Gamma; \cdot \vdash x \uparrow A} \text{var} \quad \frac{}{\Gamma; \bar{x}:A \vdash \bar{x} \uparrow A} \overline{\text{var}} \quad \frac{\Gamma, x:A; \cdot \vdash e \Downarrow B}{\Gamma; \cdot \vdash \lambda x. e \Downarrow A \rightarrow B} \rightarrow\mathbb{I} \quad \frac{\Gamma; \Delta_1 \vdash e_1 \uparrow A \rightarrow B \quad \Gamma; \Delta_2 \vdash e_2 \Downarrow A}{\Gamma; \Delta_1, \Delta_2 \vdash e_1 e_2 \uparrow B} \rightarrow\mathbb{E}$$

$$\frac{\Gamma; \Delta \vdash e \uparrow A \quad A \leq B}{\Gamma; \Delta \vdash e \Downarrow B} \text{sub} \quad \frac{\Gamma(u) = A}{\Gamma; \cdot \vdash u \uparrow A} \text{fixvar} \quad \frac{\Gamma, u:A; \cdot \vdash e \Downarrow A}{\Gamma; \cdot \vdash \text{fix } u. e \Downarrow A} \text{fix}$$

$$\frac{\Gamma \vdash e \text{ ok} \quad \Delta, \bar{x}:\perp \vdash e \text{ ok}}{\Gamma; \Delta, \bar{x}:\perp \vdash e \Downarrow C} \perp\mathbb{L} \quad \frac{(\Gamma_0 \vdash A) \lesssim (\Gamma \vdash A) \quad \Gamma; \Delta \vdash e \Downarrow A}{\Gamma; \Delta \vdash (e : \dots, (\Gamma_0 \vdash A), \dots) \uparrow A} \text{ctx-anno}$$

$$\frac{\Gamma; \Delta, \bar{x}:A \vdash e \Downarrow C}{\Gamma; \Delta, \bar{x}:A \wedge B \vdash e \Downarrow C} \wedge\mathbb{L}_1 \quad \frac{\Gamma; \Delta \vdash v \Downarrow A \quad \Gamma; \Delta \vdash v \Downarrow B}{\Gamma; \Delta \vdash v \Downarrow A \wedge B} \wedge\mathbb{I} \quad \frac{\Gamma; \Delta \vdash e \uparrow A \wedge B}{\Gamma; \Delta \vdash e \uparrow A} \wedge\mathbb{E}_1$$

$$\frac{\Gamma; \Delta, \bar{x}:B \vdash e \Downarrow C}{\Gamma; \Delta, \bar{x}:A \wedge B \vdash e \Downarrow C} \wedge\mathbb{L}_2 \quad \frac{\Gamma; \Delta \vdash e \uparrow A \wedge B}{\Gamma; \Delta \vdash e \uparrow B} \wedge\mathbb{E}_2$$

$$\frac{\Gamma; \Delta, \bar{x}:A \vdash e \Downarrow C \quad \Gamma; \Delta, \bar{x}:B \vdash e \Downarrow C}{\Gamma; \Delta, \bar{x}:A \vee B \vdash e \Downarrow C} \vee\mathbb{L} \quad \frac{\Gamma; \Delta \vdash e \Downarrow A}{\Gamma; \Delta \vdash e \Downarrow A \vee B} \vee\mathbb{I}_1 \quad \frac{\Gamma; \Delta \vdash e \Downarrow B}{\Gamma; \Delta \vdash e \Downarrow A \vee B} \vee\mathbb{I}_2$$

$$\frac{\Gamma; \Delta_1 \vdash e' \uparrow A \quad \Gamma; \Delta_2, \bar{x}:A \vdash \mathcal{E}[\bar{x}] \Downarrow C}{\Gamma; \Delta_1, \Delta_2 \vdash \mathcal{E}[e'] \Downarrow C} \text{direct}\mathbb{L} \quad \text{where } e' \text{ is not a linear variable}$$

Figure 2: The left tridirectional system

## 2.1 Tridirectional typechecking and evaluation contexts

Rule  $\text{direct}\mathbb{L}$ 's use of an evaluation context might give the impression that typechecking simply proceeds in the order in which terms are actually evaluated. However, this is not the case. The subject of  $\text{direct}\mathbb{L}$  is  $\mathcal{E}[e']$  where  $e'$  synthesizes a type, so certainly  $e'$  must be in *an* evaluation position, but there may be several such positions. Even a term as simple as  $f(x\ y)$  has 5 subterms in evaluation position, each corresponding to a different evaluation context  $\mathcal{E}$ :

$$\begin{array}{llll} \mathcal{E} = [](x\ y) & \text{and} & e' = f; & \mathcal{E} = f(x\ []) & \text{and} & e' = y; & \mathcal{E} = f\ [] & \text{and} & e' = (x\ y); \\ \mathcal{E} = f([]\ y) & \text{and} & e' = x; & & & & \mathcal{E} = [] & \text{and} & e' = f(x\ y). \end{array}$$

In fact, we may need to repeatedly apply  $\text{direct}\mathbb{L}$  to the same subject term with different choices of  $\mathcal{E}$ ! For example, we might use  $\mathcal{E} = [](x\ y)$  to name an  $f$  of union type, introducing  $\bar{f}:A \vee B$  into the context; then, case-analyze  $A \vee B$  with  $\vee E$ ; finally, choose  $\mathcal{E} = \bar{f}([]\ y)$  to name  $x$  (also of union type). Thus we are faced not with a choice over decompositions, but over many *sequences* of decompositions.

Typechecking cannot go strictly left to right. Given an ML-like `int` option type, containing `None` and some integer `Some(n)`, assume `None : none` and `Some(n) : some`. Then, if `map f (Some(n))` returns `Some(f n)` and `map f None` is `none`, then `map : (int → int) → ((some → some) ∧ (none → none))`. Similarly, a function filtering out negative integers could have type `filter : int → (some ∨ none)`.

Consider the term `(map f) (filter n)`. The term `(map f)` synthesizes  $(\text{some} \rightarrow \text{some}) \wedge (\text{none} \rightarrow \text{none})$ . This is an intersection type—we'll abbreviate it as  $(s \rightarrow s) \wedge (n \rightarrow n)$ —and the intersection must be eliminated so that rule  $\rightarrow E$  can be applied to `(map f) (filter x)`. However, we cannot commit to one part of the intersection yet, because we must first case-analyze the union type of the subterm `(filter x)`. We need to “jump over” `(map f)` to type `(filter x)`, so apply  $\text{direct}\mathbb{L}$  with evaluation context  $\mathcal{E} = [](\text{filter } x)$ , giving `(map f)` the name  $\bar{x}$ ; second, apply  $\text{direct}\mathbb{L}$  with context  $\mathcal{E} = \bar{x} []$ , synthesizing `some ∨ none` for `(filter x)`. Rule  $\vee\mathbb{L}$  splits on  $\bar{y}:\text{some} \vee \text{none}$ ; in its left subderivation  $\mathcal{D}_{\text{some}}$ , we have  $\bar{y}:\text{some}$ , so  $\wedge E_1$  on  $\bar{x} \uparrow (\text{some} \rightarrow \text{some}) \wedge (\text{none} \rightarrow \text{none})$  gives  $\bar{x} \uparrow \text{some} \rightarrow \text{some}$ , while its right subderivation  $\mathcal{D}_{\text{none}}$  has  $\bar{y}:\text{none}$ , so  $\wedge E_2$  gives  $\bar{x} \uparrow \text{none} \rightarrow \text{none}$ . Writing  $\Delta$  for  $\bar{x}:(s \rightarrow s) \wedge (n \rightarrow n)$ , the derivation is

$$\frac{\frac{\vdash \text{map } f \uparrow (s \rightarrow s) \wedge (n \rightarrow n)}{\vdash (\text{map } f) (\text{filter } x) \Downarrow C} \text{direct}\mathbb{L} \quad \frac{\frac{\Delta \vdash \text{filter } x \uparrow \text{some} \vee \text{none} \quad \frac{\mathcal{D}_{\text{some}} \quad \mathcal{D}_{\text{none}}}{\Delta, \bar{y}:(\text{some} \vee \text{none}) \vdash \bar{x} \bar{y} \Downarrow C} \vee\mathbb{L}}{\Delta \vdash \bar{x} (\text{filter } x) \Downarrow C} \text{direct}\mathbb{L}}{\vdash (\text{map } f) (\text{filter } x) \Downarrow C} \text{direct}\mathbb{L}}$$

where  $C$  is `some ∨ none`, and the derivations  $\mathcal{D}_{\text{some}}$  and  $\mathcal{D}_{\text{none}}$  are

$$\frac{\dots \vdash \bar{x} \uparrow (\text{some} \rightarrow \text{some}) \wedge (\text{none} \rightarrow \text{none})}{\dots \vdash \bar{x} \uparrow \text{some} \rightarrow \text{some}} \wedge E_1 \quad \text{and} \quad \frac{\dots \vdash \bar{x} \uparrow (\text{some} \rightarrow \text{some}) \wedge (\text{none} \rightarrow \text{none})}{\dots \vdash \bar{x} \uparrow \text{none} \rightarrow \text{none}} \wedge E_2$$

$$\begin{array}{c} \vdots \\ \Delta, \bar{y}:\text{some} \vdash \bar{x} \bar{y} \Downarrow C \end{array} \quad \begin{array}{c} \vdots \\ \Delta, \bar{y}:\text{none} \vdash \bar{x} \bar{y} \Downarrow C \end{array}$$

On a purely theoretical level, the tridirectional system is acceptable, but the nondeterminism is excessive. Xi approached (very nearly) the same problem by transforming the program so the term of  $\vee$  type appears before the term of  $\wedge$  type. (Actually, Xi had index-level quantifiers  $\Sigma$  and  $\Pi$  instead of  $\vee$  and  $\wedge$ , but these are analogous.) A standard let-normal translation  $|e|$  (Xi 1998, p. 86), where  $|e_1\ e_2| = \text{let } x_1 = |e_1| \text{ in let } x_2 = |e_2| \text{ in } x_1 x_2$  suffices for the examples above. (In Xi's system, existential variables are unpacked where a term of existential type is let-bound: an existential variable  $b'$  is unpacked at the binding of  $x_2$ , which appears before the application  $x_1 x_2$  at which the universal variable  $a$  must be instantiated.) Unfortunately, the translation interacts unpleasantly with bidirectionality: terms

such as  $\text{map } (\lambda x. e)$ , in which  $(\lambda x. e)$  must be checked, no longer typecheck because the  $\lambda$  becomes the right hand side of a  $\text{let}$ , in  $\text{let } x_1 = \text{map in let } x_2 = \lambda x. e \text{ in } x_1 x_2$  and  $\text{let}$ -bound expressions must synthesize a type, but  $\lambda x. e$  does not. Typechecking becomes incomplete in the sense that some programs that were well typed are not well typed after translation.

Xi ameliorated this incompleteness by treating  $e_1 v_2$  as a special case (Xi 1998, p. 139):  $|e_1 v_2| = \text{let } x_1 = |e_1| \text{ in } x_1 v_2$ . Now  $v_2$  (which is  $\lambda x. e$  in the above example) is in a checking position. This is adequate for non-synthesizing values, but terms such as  $\text{map}$  (case  $z$  of ...), where a non-synthesizing non-value is in checking position, remain untypable. It is not clear why Xi did not also have special cases for  $\text{case}$  and other non-synthesizing non-values, e.g.  $|e_1 (\text{case } e \text{ of } ms)| = \text{let } x_1 = |e_1| \text{ in } x_1 |\text{case } e \text{ of } ms|$ . Xi's translation is also incomplete for terms like  $f$  (case  $x$  of  $ms$ ). Suppose  $x$  synthesizes a union that must be analyzed to select the appropriate part of an intersection in the type of  $f$ . Since  $x$ 's scope—and thus the scope of its union—is entirely within the  $\text{let}$  created for the case, typechecking fails.

$$\begin{aligned} |f(\text{case } x \text{ of } ms)| &= \text{let } f_1 = f \text{ in let } x_0 = |\text{case } x \text{ of } ms| \text{ in } f_1 x_0 \\ &= \text{let } f_1 = f \text{ in let } x_0 = (\text{let } x_1 = x \text{ in case } x_1 \text{ of } |ms|) \text{ in } f_1 x_0 \end{aligned}$$

It could be argued that the cases in which Xi's translation fails are rare in practice. However, that may only increase confusion when such a case is encountered. I follow Xi's general approach of sequentializing the program before typechecking, but no programs are lost in my translation.

Do we need all the freedom that  $\text{direct}\mathbb{L}$  provides? No. At the very least, if we do not need to name a subterm, naming it anyway does no harm. But naming all the subterms only slightly reduces the nondeterminism. Clearly, a strategy of in-order traversal is sound (we can choose to apply  $\text{direct}\mathbb{L}$  from left to right if we like). It is tempting to think it is complete. In fact, it holds for many programs, but fails for a certain class of annotated terms. We will explain why as we present the general mechanism for enforcing a strategy of left-to-right traversal except for certain annotated terms.

### 3 Let-Normal Typechecking

We'll briefly mention previous work on let-normal form, then explain the ideas behind the variant here, including why we need a *principal synthesis of values* property. Because the most universal form of principality does not hold for a few terms, we introduce *slack bindings*.

Traditional let-normal or A-normal transformations (Moggi 1988; Flanagan et al. 1993) (1) explicitly sequence the computation, and (2) name the result of each intermediate computation. (Continuation-passing style (CPS) (Reynolds 1993) also (3) introduces named continuations. Thus let-normal form is also known as *two-thirds CPS*.) Many compilers for functional languages use some kind of let-normal form to facilitate optimizations; see, for instance, Tarditi et al. (1996), Reppy (2001), Chlipala et al. (2005), and Peyton Jones et al. (2006).

Our let-normal form will sequentialize the computation, but it does not only name intermediate *computations*, but *values* as well. In our let-normal type system,  $\text{direct}\mathbb{L}$  is replaced by a rule  $\text{let}$  that can *only* be applied to  $\text{let}$ ; see Figure 3.  $\mathcal{Q}$  is a special evaluation context, discussed below.

This  $\text{let}$  is a syntactic marker with no computational character. In contrast to let-normal translations for compilation purposes, there is no evaluation step (reduction) corresponding to a  $\text{let}$ . I won't even give a dynamic semantics for terms with lets. It would be easy; it's simply not useful here. If we insist on knowing what a let-normal term  $e$  means, we can use a standard call-by-value operational semantics over the term's reverse translation.

Instead of making explicit the order of computation, our let-normal form makes explicit the order of *typechecking*—or rather, the order in which  $\text{direct}\mathbb{L}$  names subterms in evaluation position. Thus, to be

$$\frac{\Gamma; \Delta_1 \vdash e' \uparrow A \quad \Gamma; \Delta_2, \bar{x}:A \vdash \mathcal{Q}[\bar{x}] \downarrow C}{\Gamma; \Delta_1, \Delta_2 \vdash \text{let } \bar{x}=e' \text{ in } \mathcal{Q}[\bar{x}] \downarrow C} \text{let}$$

$$\frac{\Gamma; \Delta, \sim\bar{x}=v \vdash \mathcal{Q}[\bar{x}] \downarrow C}{\Gamma; \Delta \vdash \text{let } \sim\bar{x}=v \text{ in } \mathcal{Q}[\bar{x}] \downarrow C} \text{let}\sim \quad \frac{\Gamma; \Delta_1 \vdash v \uparrow A \quad \Gamma; \Delta_2, \bar{x}:A \vdash e \downarrow C}{\Gamma; \Delta_1, \Delta_2, \sim\bar{x}=v \vdash e \downarrow C} \sim\text{var}$$

... plus all rules in Figure 2, except  $\text{direct}\mathbb{L}$

**Figure 3:** The *let-normal type system* for terms containing  $\text{let } \bar{x}$  bindings

complete with respect to the tridirectional system, the transformation must create a  $\text{let}$  for every subterm in synthesizing form: if an (untranslated) program contains a subterm  $e'$  in synthesizing form, it might be possible to name  $e'$  with  $\text{direct}\mathbb{L}$ , so the  $\text{let-normal}$  translation must bind  $e'$ . Otherwise, a chance to apply  $\forall\mathbb{L}$  is lost. Even variables  $x$  must be named, since they synthesize a type and so can be named in  $\text{direct}\mathbb{L}$ . This models an “aggressive” strategy of applying  $\text{direct}\mathbb{L}$ . On the other hand, checked terms like  $\lambda x.e$  can’t synthesize, so we won’t name them.

Another consequence of the  $\text{let-normal}$  form following typing, not evaluation, is that  $\text{let } \bar{x}=v_1 \text{ in } v_2$  is considered a value—after all, the original term  $[v_1/\bar{x}]v_2$  was a value, and we transformed a value into a non-value we could not apply value-restricted typing rules such as  $\wedge\text{I}$ , leading to incompleteness.

We define the translation by a judgment  $e \hookrightarrow L + e'$ , read “ $e$  translates to a sequence of  $\text{let}$ -bindings  $L$  with body  $e'$ ”. For example, the translation of  $f(x y)$ , which names every synthesizing subterm, is

$$\text{let } \bar{f}=f \text{ in let } \bar{x}=x \text{ in let } \bar{y}=y \text{ in let } \bar{z}=\bar{x} \bar{y} \text{ in let } \bar{a}=\bar{f} \bar{z} \text{ in } \bar{a}$$

This is expressed by the judgment  $f(x y) \hookrightarrow \bar{f}=f, \bar{x}=x, \bar{y}=y, \bar{z}=\bar{x} \bar{y}, \bar{a}=\bar{f} \bar{z} + \bar{a}$ . Figure 4 has the definition. Note that  $L + e'$  is not a term;  $+$  is punctuation. We write  $L$  in  $e'$  as shorthand: read  $e \hookrightarrow L + e'$  as “ $e \hookrightarrow L$  in  $e'$ ”. The divergent notations come from the multiple decompositions of a term into a pair of bindings and a “body”. For example,  $\text{let } \bar{x}_1=e_1 \text{ in let } \bar{x}_2=e_2 \text{ in } e_3$  can be written three ways: **(1)**  $\cdot$  in  $\text{let } \bar{x}_1=e_1 \text{ in let } \bar{x}_2=e_2 \text{ in } e_3$ , **(2)**  $(\bar{x}_1=e_1)$  in  $\text{let } \bar{x}_2=e_2 \text{ in } e_3$ , or **(3)**  $(\bar{x}_1=e_1, \bar{x}_2=e_2)$  in  $e_3$ . The last decomposition is *maximal*: it has the maximum number of bindings (and the smallest ‘body’), which is the case when the body isn’t a  $\text{let}$ . If  $e \hookrightarrow L + e'$  then  $L$  in  $e'$  is maximal.

Again, to model a complete strategy of  $\text{direct}\mathbb{L}$ -application, in  $e \hookrightarrow L + e'$  we need  $L$  to bind all the synthesizing subterms that could be in evaluation position (after applying  $\text{direct}\mathbb{L}$  zero or more times).

We syntactically partition terms into *pre-* and *anti-*values. A *pre-value*  $\check{e}$  is a value, such as  $x$ , or a term that can “become” a value via  $\text{direct}\mathbb{L}$ , such as  $x y$  which “becomes” the value  $\bar{z}$  in the derivation. (The háček  $\check{\cdot}$  above the  $e$  is shaped like a ‘v’ for ‘value’.) An *anti-value*  $\hat{e}$ , such as  $\text{fix } u.e$  (or case  $e$  of  $ms$ ) is not a value and cannot become a value.

$\text{direct}\mathbb{L}$  can replace any synthesizing subterm with a linear variable, so the pre-values must include both the values and the synthesizing forms. This leads to the following grammar for pre-values, with values  $x$ ,  $\bar{x}$ , and  $\lambda x.e$  and synthesizing forms  $(e : As)$ ,  $e_1 e_2$ ,  $u$ . (In the full system, the prevalues also include checking forms that can become values if all their subterms can, such as  $(e_1, e_2)$ .)

$$\begin{array}{l} \text{Pre-values } \check{e} ::= x \mid \bar{x} \mid (e : As) \mid \lambda x.e \mid e_1 e_2 \mid u \\ \text{Anti-values } \hat{e} ::= \text{fix } u.e \end{array}$$

The distinction matters for terms with sequences of immediate subterms such that at least two subterms in the sequence may be in evaluation position. Only application  $e_1 e_2$  has this property (and in the full system, pairs  $(e_1, e_2)$ ).  $\lambda x.e$  and  $\text{fix } u.e$  have no subterms in evaluation position at all.

$$\boxed{e \hookrightarrow L + e'} \text{ read “}e \text{ translates to bindings } L \text{ with result } e'\text{”}$$

$$\frac{}{x \hookrightarrow (\bar{\mathbf{x}}=x) + \bar{\mathbf{x}}} \quad \frac{e \hookrightarrow L + e'}{\lambda x. e \hookrightarrow \cdot + \lambda x. (L \text{ in } e')}$$

$$\frac{}{u \hookrightarrow (\bar{\mathbf{x}}=u) + \bar{\mathbf{x}}} \quad \frac{e \hookrightarrow L + e'}{\text{fix } u. e \hookrightarrow \cdot + \text{fix } u. (L \text{ in } e')}$$

$$\frac{\widehat{e}_1 \hookrightarrow L_1 + e'_1 \quad e_2 \hookrightarrow L_2 + e'_2}{\widehat{e}_1 e_2 \hookrightarrow L_1, \bar{\mathbf{x}}=e'_1(L_2 \text{ in } e'_2) + \bar{\mathbf{x}}} \quad \frac{\check{e}_1 \hookrightarrow L_1 + e'_1 \quad e_2 \hookrightarrow L_2 + e'_2}{\check{e}_1 e_2 \hookrightarrow L_1, L_2, \bar{\mathbf{x}}=e'_1 e'_2 + \bar{\mathbf{x}}}$$

$$\frac{e \hookrightarrow L + e' \quad e \text{ not a value}}{(e : As) \hookrightarrow L, \bar{\mathbf{x}}=(e' : As) + \bar{\mathbf{x}}} \quad \frac{v \hookrightarrow L + e'}{(v : As) \hookrightarrow L, \sim\bar{\mathbf{x}}=(e' : As) + \bar{\mathbf{x}}} \quad \frac{}{\bar{\mathbf{x}} \hookrightarrow \cdot + \bar{\mathbf{x}}}$$

**Figure 4:** The let-normal transformation

A telling example is  $(\text{fix } u. e) (\omega x)$  where  $\omega : \dots \rightarrow \perp$ . In the tridirectional system, this term has no synthesizing subterms in evaluation position. In particular,  $\omega x$  is not in evaluation position, so however we translate the term, we must not bind  $\omega x$  outside the outer application; if we did, we would add  $\bar{\mathbf{z}}:\perp$  to the context and could apply rule  $\perp\mathbb{L}$  to declare the outer application well typed while ignoring  $e!$  If  $e$  is ill-typed, this is actually unsound. On the other hand, in the term  $(f g) (\omega x)$  the left tridirectional system *can* bind  $\omega x$  before checking the pair, by applying  $\text{direct}\mathbb{L}$  with  $\mathcal{E} = []$   $(\omega x)$  (synthesizing a type for  $f g$ , ensuring soundness) to yield a subject  $\bar{\mathbf{x}}(\omega x)$  in which  $\omega x$  is in evaluation position.

The difference is that  $\text{fix } u. e$  is an anti-value, while  $f g$  is a pre-value. Therefore, given an application  $e_1 e_2$ , if  $e_1$  is some anti-value  $\widehat{e}_1$ , the translation places the bindings for subterms of  $e_2$  (e.g.  $\bar{\mathbf{z}}=\omega x$  above) *inside* the argument part. On the other hand, if  $e_1$  is a pre-value  $\check{e}_1$ , the translation puts the bindings for subterms of  $e_2$  *outside* the application. See the shaded rules in Figure 4.

*Elongated evaluation contexts*  $\mathcal{Q}$ , unlike ordinary evaluation contexts  $\mathcal{E}$ , can skip over pre-values.  $\mathcal{Q}$  is a sort of transitive closure of  $\mathcal{E}$ : if, by repeatedly replacing pre-values in evaluation position with values, some subterm is then in evaluation position, that subterm is in elongated evaluation position. In a sequence of  $\text{direct}\mathbb{L}$ -applications, subterms in evaluation position are replaced with linear variables, which are values. For example,  $z$  is not in evaluation position in  $(x y) z$ , but applying  $\text{direct}\mathbb{L}$  with  $\mathcal{E} = [] z$  yields a subderivation with subject  $\bar{\mathbf{x}} z$ , in which  $z$  is in evaluation position. A  $\mathcal{Q}$  is thus a path that can skip pre-values: if every intervening subterm is a pre-value (equivalently, if there is no intervening anti-value), the hole is in elongated evaluation position. The grammar for let-normal terms ensures that the body  $e_2$  of  $\text{let } \bar{\mathbf{x}}=e_1 \text{ in } e_2$  must have the form  $\mathcal{Q}[\bar{\mathbf{x}}]$ .

$$\begin{array}{ll}
\text{Elongated} & \mathcal{Q} ::= [] \mid \mathcal{Q}e \mid \check{e}\mathcal{Q} \mid (\mathcal{Q} : As) \\
\text{evaluation contexts} & \mid \text{let } \bar{\mathbf{x}}=\mathcal{Q} \text{ in } e \mid \text{let } \bar{\mathbf{x}}=\check{e} \text{ in } \mathcal{Q} \mid \text{let } \sim\bar{\mathbf{x}}=\mathcal{Q} \text{ in } e \mid \text{let } \sim\bar{\mathbf{x}}=v \text{ in } \mathcal{Q} \\
\text{Terms} & e ::= \dots \mid \text{let } \bar{\mathbf{x}}=e_1 \text{ in } \mathcal{Q}[\bar{\mathbf{x}}] \mid \text{let } \sim\bar{\mathbf{x}}=v_1 \text{ in } \mathcal{Q}[\bar{\mathbf{x}}] \\
\text{Values} & v ::= x \mid \lambda x. e \mid \bar{\mathbf{x}} \mid \text{let } \bar{\mathbf{x}}=v_1 \text{ in } v_2 \mid \text{let } \sim\bar{\mathbf{x}}=v_1 \text{ in } v_2 \\
\text{Eval. contexts} & \mathcal{E} ::= \dots \mid \text{let } \bar{\mathbf{x}}=\mathcal{E} \text{ in } e \mid \text{let } \bar{\mathbf{x}}=v \text{ in } \mathcal{E} \mid \text{let } \sim\bar{\mathbf{x}}=\mathcal{E} \text{ in } e \mid \text{let } \sim\bar{\mathbf{x}}=v \text{ in } \mathcal{E} \\
\text{Sequences of bindings} & L ::= \cdot \mid L, (\bar{\mathbf{x}}=e) \mid L, (\sim\bar{\mathbf{x}}=v)
\end{array}$$

### 3.1 Principal synthesis of values

A key step in completeness is the movement of let-bindings outward. To prove this preserves typing, we show that principal types (Hindley 1969) exist in certain cases. Consider the judgment  $x : (A_1 \rightarrow B) \wedge (A_2 \rightarrow B), y : A_1 \vee A_2; \cdot \vdash x y \Downarrow B$ . To derive this in the left tridirectional system, we need  $\text{direct}\mathbb{L}$  with  $\mathcal{E} = x []$  to name  $y$  as a new linear variable  $\bar{y} : A_1 \vee A_2$ . Then we use  $\vee\mathbb{L}$ ; we must now derive

$$x : (A_1 \rightarrow B) \wedge (A_2 \rightarrow B), \dots; \bar{y} : A_1 \vdash x \bar{y} \Downarrow B \quad \text{and} \quad x : (A_1 \rightarrow B) \wedge (A_2 \rightarrow B), \dots; \bar{y} : A_2 \vdash x \bar{y} \Downarrow B$$

Here, the scope of  $\bar{y}$  is  $x \bar{y}$ , and we synthesize a type for  $x$  *twice*, once in each branch:

$$\frac{\frac{\dots; \cdot \vdash x \uparrow A_1 \rightarrow B \quad \vdots}{\dots; \bar{y} : A_1 \vdash x \bar{y} \Downarrow B} \rightarrow E \quad \frac{\dots; \cdot \vdash x \uparrow A_2 \rightarrow B \quad \vdots}{\dots; \bar{y} : A_2 \vdash x \bar{y} \Downarrow B} \rightarrow E}{\dots; \bar{y} : A_1 \vee A_2 \vdash x \bar{y} \Downarrow B} \vee\mathbb{L}}{\dots, y : A_1 \vee A_2; \cdot \vdash y \uparrow A_1 \vee A_2 \quad \dots; \bar{y} : A_1 \vee A_2 \vdash x \bar{y} \Downarrow B} \text{direct}\mathbb{L}}{\frac{}{x : (A_1 \rightarrow B) \wedge (A_2 \rightarrow B), y : A_1 \vee A_2; \cdot \vdash x y \Downarrow B} \text{direct}\mathbb{L}}$$

However, when checking the translated term  $\text{let } \bar{x} = x \text{ in let } \bar{y} = y \text{ in let } \bar{z} = \bar{x} \bar{y} \text{ in } \bar{z}$  against  $B$ , we need to first name  $x$  as  $\bar{x}$ , then  $y$  as  $\bar{y}$ , then use  $\vee\mathbb{L}$  to decompose the union  $\bar{y} : A_1 \vee A_2$  with subject  $\text{let } \bar{z} = \bar{x} \bar{y} \text{ in } \bar{z}$ .

$$\frac{\dots; \cdot \vdash x \uparrow (A_1 \rightarrow B) \wedge (A_2 \rightarrow B) \quad \dots; \bar{x} : (A_1 \rightarrow B) \wedge (A_2 \rightarrow B) \vdash \text{let } \bar{y} = y \text{ in let } \bar{z} = \bar{x} \bar{y} \text{ in } \bar{z} \Downarrow B}{x : (A_1 \rightarrow B) \wedge (A_2 \rightarrow B), y : A_1 \vee A_2; \cdot \vdash \text{let } \bar{x} = x \text{ in let } \bar{y} = y \text{ in let } \bar{z} = \bar{x} \bar{y} \text{ in } \bar{z} \Downarrow B} \text{let}$$

But we only get one chance (highlighted above) to synthesize a type for  $x$ , so we must take care when using  $\text{let}$  to name  $x$ ; if we choose to synthesize  $x \uparrow A_1 \rightarrow B$  in  $\text{let}$ , we can't derive

$$\bar{x} : A_1 \rightarrow B, \bar{y} : A_2 \vdash (\text{let } \bar{z} = \bar{x} \bar{y} \text{ in } \bar{z}) \Downarrow B$$

but if we choose to synthesize  $x \uparrow A_2 \rightarrow B$  we can't get

$$\bar{x} : A_2 \rightarrow B, \bar{y} : A_1 \vdash (\text{let } \bar{z} = \bar{x} \bar{y} \text{ in } \bar{z}) \Downarrow B$$

The only choice that works is  $\Gamma(x)$ , which is  $(A_1 \rightarrow B) \wedge (A_2 \rightarrow B)$ , since given  $\bar{x} \uparrow (A_1 \rightarrow B) \wedge (A_2 \rightarrow B)$  we can synthesize  $\bar{x} \uparrow A_1 \rightarrow B$  and  $\bar{x} \uparrow A_2 \rightarrow B$  using  $\wedge E_1$  and  $\wedge E_2$ , respectively.

In the above situation,  $e' = x$  is a variable, so there is a best type  $C$ —namely  $\Gamma(x)$ —such that if  $x \uparrow C_1$  and  $x \uparrow C_2$  then  $x \uparrow C$ , from which follows (by rules  $\wedge E_{1,2}$  in the example above)  $x \uparrow C_1$  and  $x \uparrow C_2$ . We'll say that  $x$  has the property of *principal synthesis*. Which terms have this property? Variables do: the best type for some  $x$  is  $\Gamma(x)$ . On the other hand, it does not hold for many non-values:  $f x \uparrow A_1$  and  $f x \uparrow A_2$  do *not* imply  $f x \uparrow A_1 \wedge A_2$ , since the intersection introduction rule  $\wedge I$  is (1) restricted to values and (2) in the checking direction. Fortunately, we don't need it for non-values: Consider  $(e_1 e_2) y$ . Since  $(e_1 e_2)$  is not a value,  $y$  is not in evaluation position in  $(e_1 e_2) y$ , so even in the tridirectional system, to name  $y$  we must first name  $(e_1 e_2)$ . Here, the  $\text{let}$ -normal system is no more restrictive. Moreover, some values, such as pairs, are checking forms and *never* synthesize, so they do not have the principal synthesis property. But neither system binds values in checking form to linear variables.

Now, do all *values in synthesizing form* have the principal synthesis property? The only values in synthesizing form are ordinary variables  $x$ , linear variables  $\bar{x}$ , and annotated values  $(v : As)$ . For  $x$  or  $\bar{x}$  the principal type is simply  $\Gamma(x)$  or  $\Delta(\bar{x})$ . Unfortunately, principal types do not always exist for terms of the form  $(v : As)$ . For example,  $((\lambda x. x) : (\vdash \text{unit} \rightarrow \text{unit}), (\vdash \text{bool} \rightarrow \text{bool}))$  can synthesize  $\text{unit} \rightarrow \text{unit}$ , and it can synthesize  $\text{bool} \rightarrow \text{bool}$ , but it can't synthesize their intersection, so it has no principal type.



### 3.2 Slack bindings

Rather than restrict the form of annotations, we use a different kind of binding for  $(v : As)$ —a *slack binding*  $\sim\bar{x}=v$  where  $v$ 's type is synthesized not at its binding site, but at any point up to its use (rules  $\sim\text{var}$  and  $\text{let}\sim$  in Figure 3). Wherever  $\bar{x}$  is in scope, we can try rule  $\sim\text{var}$  to synthesize a type  $A$  for  $v$  and replace  $\sim\bar{x}=v$  with an ordinary linear variable typing  $\bar{x}:A$ . For example,  $((\lambda x.e) : (\vdash \text{int} \rightarrow \text{int})) y$  is translated to  $\text{let } \sim\bar{x} = ((\lambda x.e') : (\vdash \text{int} \rightarrow \text{int})) \text{ in } \text{let } \bar{y} = y \text{ in } \bar{x} \bar{y}$ . Now, we have several chances to use  $\sim\text{var}$  to synthesize the type of  $\bar{x}$ : just before checking  $\text{let } \bar{y} = y \text{ in } \bar{x} \bar{y}$ , or when checking  $\bar{x} \bar{y}$ . This is just like choosing when to apply  $\text{direct}\mathbb{L}$  in the tridirectional system. If all our bindings were slack we would have put ourselves in motion to no purpose, but we'll use slack bindings for  $(v : As)$  *only*. My experiments suggest that slack bindings are rare in practice (Dunfield 2007b, p. 187), and are certainly less problematic than the backtracking from intersections and unions themselves ( $\wedge E_{1,2}$ , etc.).

## 4 Results

The two major results are *soundness*: if the let-normal translation of a program is well typed in the let-normal type system, the original program is well typed in the left tridirectional system—and *completeness*: if a program is well typed in the left tridirectional type system, its translation is well typed in the let-normal type system. Once these are shown, it follows from Dunfield and Pfenning (2004) that the let-normal system is sound and complete with respect to a system (Dunfield and Pfenning 2003) for which preservation and progress hold under a call-by-value semantics.

At its heart, the let-normal system merely enforces a particular pattern of linear variable introductions (via  $\text{let}$ , instead of  $\text{direct}\mathbb{L}$ ). So it is no surprise that soundness holds. The proof is syntactic, but not too involved; see Dunfield (2007b, pp. 132–134).

**Corollary (Let-Normal Soundness).**

*If  $e \hookrightarrow L + e'$  and  $;\cdot \vdash L$  in  $e' \Downarrow C$  (let-normal system) then  $;\cdot \vdash e \Downarrow C$  (tridirectional system).*

However, completeness—that the let-normal system is not *strictly* weaker than the tridirectional system—is involved. What follows is the roughest sketch of the proof found in Dunfield (2007b, pp. 135–165). We want to show that given a well-typed term  $e$ , the let-normal translation  $L$  in  $e'$  where  $e \hookrightarrow L + e'$  is well-typed. To be precise, given a derivation  $\mathcal{D}$  deriving  $\Gamma; \Delta \vdash e \Downarrow C$  in the left tridirectional system, we must construct a derivation  $\Gamma; \Delta \vdash L$  in  $e' \Downarrow C$  in the let-normal system, where  $e \hookrightarrow L + e'$ . My attempts to prove this by straightforward induction on the derivation failed: thanks to  $\text{direct}\mathbb{L}$ , the relationship between  $e$  and  $\mathcal{D}$  is complex. Nor is  $L + e'$  compositional in  $e$ : for a given subterm of  $e$  there may not be a corresponding subterm of  $L + e'$ , because translation can insert bindings inside the translated subterm.

Instead, the completeness proof proceeds as follows:

1. *Mark  $e$  with lets* wherever  $\text{direct}\mathbb{L}$  is used in  $\mathcal{D}$ . However, if  $\wedge I$  or another subject-duplicating rule is used, the subderivations need not apply  $\text{direct}\mathbb{L}$  in the same way, resulting in distinct terms to which  $\wedge I$  cannot be applied. So we use step 2 inductively to obtain typing derivations for the canonical version of the subterm (the  $L + e'$  from  $e \hookrightarrow L + e'$ ), to which  $\wedge I$  can be applied.

This step centres on a lemma which produces a term with a let-system typing derivation. This term might not be canonical. For example, if the original tridirectional derivation for  $\lambda x.x$  didn't use  $\text{direct}\mathbb{L}$  at all, no let bindings are created, unlike the canonical let-normal term  $\lambda x.\text{let } \bar{x} = x \text{ in } \bar{x}$ .

2. *Transform* the marked term into the canonical  $L + e'$  in small steps, adding or moving one let at a time. Each small step preserves typing. We'll define a syntactic measure  $\mu$  that quantifies how different a term is from  $L + e'$ ; each let-manipulating step reduces the measure, bringing the term closer to  $L + e'$ . When the measure is all zeroes, the term is  $L + e'$ .

The *measure* of  $e'$  is

$$\mu(e') = \langle \text{unbound}_{\uparrow}(e'), \text{brittle}(e'), \text{prickly}(e'), \text{transposed}(e') \rangle$$

where:

- $\text{unbound}_{\uparrow}(e')$  is the number of subterms of  $e'$  in synthesizing form (that is, variables  $x$  and  $u$ , annotated terms  $(e : As)$ , and applications  $e_1 e_2$ ) that are not let-bound. The translation  $\hookrightarrow$  has let-bindings for all such terms, so an  $e'$  that does not bind such terms is quite far from being in canonical let-normal form.
- $\text{brittle}(e')$  is the number of let-bindings in  $e'$  of the form  $\text{let } \bar{x} = (v_1 : As) \text{ in } e_2$ . To correspond to the translation  $\hookrightarrow$ , we need to change such let-bindings to slack bindings  $\text{let } \sim \bar{x} = (v_1 : As) \text{ in } e_2$ . These terms are “brittle” because they need to be slackened.
- $\text{prickly}(e')$  is the number of let-bindings in  $e'$  that are not properly collected together at a *root*. A *root* is somewhere that the canonical translation  $\hookrightarrow$  may place a sequence of let-bindings. In the proof, we start by reducing the number of unbound synthesizing forms by inserting lets nearby, but some of these are too deep inside the term. For example, given a term  $c y$ , we first put a binding around the  $y$ , giving  $c(\text{let } \bar{y} = y \text{ in } \bar{y})$ . (To simplify the example,  $c$  is some constant or primitive operation that is never let-bound.) Then we bind the application, giving  $\text{let } \bar{a} = c(\text{let } \bar{y} = y \text{ in } \bar{y}) \text{ in } \bar{a}$ . But the canonical translation would be  $\text{let } \bar{y} = y \text{ in let } \bar{a} = c \bar{y} \text{ in } \bar{a}$ . Thus, a prickly binding needs to be lifted outward until it is in some sequence of let-bindings at the outside of the body of a  $\lambda$  or  $\text{fix}$ , or at the outside of the entire term  $e'$ .
- $\text{transposed}(e')$  is the number of *transposed variable pairs* in  $e'$ . If there are no prickly bindings, there may still be bindings that are out of order. For a term  $x y$ , the original derivation might have used  $\text{direct}\mathbb{L}$  first on  $y$  (with  $\mathcal{E} = x []$ ) then on  $x$  (with  $\mathcal{E} = [] \bar{y}$ ). In this case, Step 1 above would produce  $\text{let } \bar{y} = y \text{ in let } \bar{x} = x \text{ in } \bar{x} \bar{y}$ . Supposing this application is the body of some  $\lambda$ , these bindings are not prickly, but don't correspond to what  $\hookrightarrow$  would produce. Variables (and their bindings) are transposed if they are not used in the same order they were bound. Thus,  $\bar{x}$  and  $\bar{y}$  are transposed in  $\text{let } \bar{y} = y \text{ in let } \bar{x} = x \text{ in } \bar{x} \bar{y}$ , because  $\bar{y}$  is bound before  $\bar{x}$  but  $\bar{x}$  appears to the left of  $\bar{y}$  in the body of the let.

We interpret the quadruples lexicographically. Likewise, the proof of completeness relies on type preservation lemmas for each part of the quadruple: adding a let-binding preserves typing, changing a regular let-binding to a slack let-binding preserves typing, lifting a let-binding to a root preserves typing, and reordering the bindings of transposed variables preserves typing.

**Theorem** (Let-Normal Completeness).

If  $\cdot; \vdash e \Downarrow C$  (*tridirectional system*) and  $e \hookrightarrow L + e^*$  then  $\cdot; \vdash L$  in  $e^* \Downarrow C$  (*let-normal system*).

## 5 Related Work

The effects of transformation to continuation passing style on the precision of program analyses such as 0-CFA have been studied for some time (Sabry and Felleisen 1994). The effect depends on the specific

details of the CPS transform and the analysis done (Damian and Danvy 2001; Palsberg and Wand 2003). The “analysis” in this work is the process of bidirectional checking/synthesis. My soundness and completeness results show that my let-normal transformation does not affect the analysis. It is not clear if this means anything for more traditional let-normal transformations and compiler analyses.

## 6 Conclusion

Transforming programs into a let-normal form removes a major impediment to implementing tridirectional typechecking. The system is sound *and complete* with respect to a type assignment system for intersections and unions (Dunfield and Pfenning 2003), in contrast to systems (Xi 1998) in which completeness is lost. The tridirectional rule *can* be turned into something practical. A chain of soundness results (Dunfield 2007b, p. 165) guarantees that if we run a program  $e$  whose let-normal translation typechecks in the system in this paper, it will not go wrong.

Despite “untangling”  $\text{direct}\mathbb{L}$ , typechecking is still very time-consuming in the worst cases, thanks to checking terms several times in  $\wedge I$  and backtracking in  $\wedge E_{1,2}$ , etc. As implementing (an extended version of) this system shows (Dunfield 2007a), bad cases do occur in practice!

Parametric polymorphism is absent, but I have extended the tridirectional system and the let-normal implementation (Dunfield 2009), and the soundness and completeness results should still hold.

The major flaw of this work is its completeness proof, which uses purely syntactic methods, is complicated, and has not been mechanized. Ideally, it would be mechanized and/or proved more simply.

**Acknowledgments** Many thanks to Frank Pfenning for countless discussions about this research. Thanks also to the ITRS reviewers. Most of the work was done at Carnegie Mellon University with the support of the US National Science Foundation.

## References

- Adam Chlipala, Leaf Petersen, and Robert Harper. Strict bidirectional type checking. In *Workshop on Types in Language Design and Implementation (TLDI '05)*, pages 71–78, 2005.
- Daniel Damian and Olivier Danvy. Syntactic accidents in program analysis: on the impact of the CPS transformation. Technical Report BRICS-RS-01-54, University of Aarhus, 2001.
- Rowan Davies and Frank Pfenning. Intersection types and computational effects. In *ICFP*, pages 198–208, 2000.
- Joshua Dunfield. Refined typechecking with Stardust. In *Programming Languages meets Programming Verification (PLPV '07)*, 2007a.
- Joshua Dunfield. Greedy bidirectional polymorphism. In *ML Workshop (ML '09)*, 2009.
- Joshua Dunfield. *A Unified System of Type Refinements*. PhD thesis, Carnegie Mellon University, 2007b. CMU-CS-07-129.
- Joshua Dunfield and Frank Pfenning. Type assignment for intersections and unions in call-by-value languages. In *Found. Software Science and Computation Structures*, pages 250–266, 2003.
- Joshua Dunfield and Frank Pfenning. Tridirectional typechecking. In *POPL*, pages 281–292, 2004.

- Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. The essence of compiling with continuations. In *Programming Language Design and Implementation*, pages 237–247, 1993.
- R. Hindley. The principal type-scheme of an object in combinatory logic. *Trans. Am. Math. Soc.*, 146: 29–60, 1969.
- Eugenio Moggi. Computational lambda-calculus and monads. Technical Report ECS-LFCS-88-66, University of Edinburgh, 1988.
- Jens Palsberg and Mitchell Wand. CPS transformation of flow information. *J. Functional Programming*, 13(5):905–923, 2003.
- Simon Peyton Jones and the GHC developers. Glasgow Haskell Compiler Commentary. <http://hackage.haskell.org/trac/ghc/wiki/Commentary/Compiler/HscPipe>, 2006.
- Benjamin C. Pierce and David N. Turner. Local type inference. In *POPL*, pages 252–265, 1998. Full version in *ACM Trans. Programming Languages and Systems*, 22(1):1–44, 2000.
- John Reppy. Local CPS conversion in a direct-style compiler. In *ACM Workshop on Continuations (CW '01)*, pages 13–22, 2001.
- John C. Reynolds. The discoveries of continuations. *LISP and Symbolic Computation*, 6(3–4):233–247, 1993.
- John C. Reynolds. Design of the programming language Forsythe. Technical Report CMU-CS-96-146, Carnegie Mellon University, 1996.
- Amr Sabry and Matthias Felleisen. Is continuation-passing useful for data flow analysis? In *Programming Language Design and Implementation*, pages 1–12, 1994.
- D. Tarditi, G. Morrisett, P. Cheng, C. Stone, R. Harper, and P. Lee. TIL: A type-directed optimizing compiler for ML. In *Programming Language Design and Implementation*, pages 181–192, 1996.
- Hongwei Xi. *Dependent Types in Practical Programming*. PhD thesis, Carnegie Mellon University, 1998.