

Typed Compilation Against Non-Manifest Base Classes

Christopher League

Long Island University
christopher.league@liu.edu

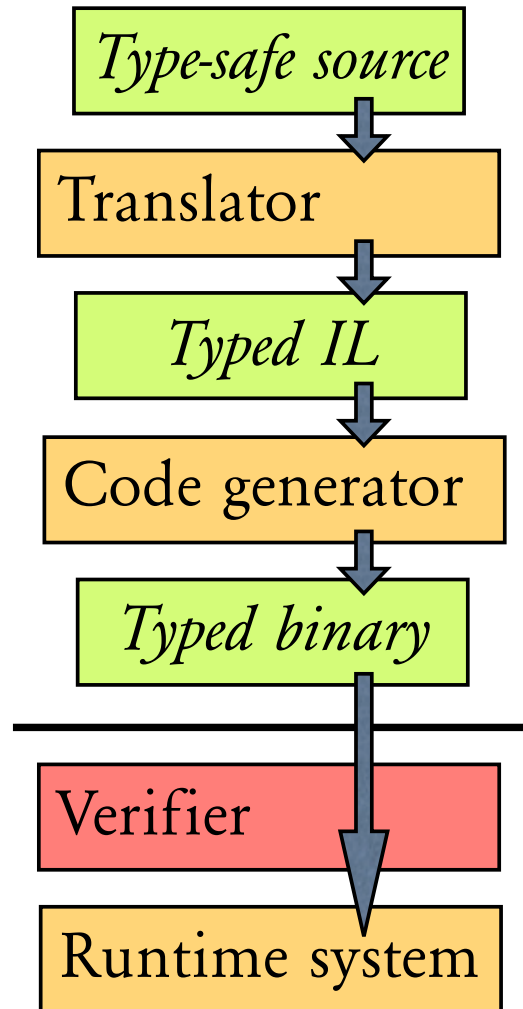
Stefan Monnier

Université de Montréal
monnier@iro.umontreal.ca

FTfJP workshop
26 July 2005

Goal: certifying compiler

...for Java-like languages



- Type safety is necessary for security
- Fairly well understood for core Java/C#

[Colby et al., PLDI '00]

[League et al., TOPLAS '02, CC '03]

- What about *next generation* object-oriented features?

A characteristic of next-gen OO

- Looser coupling between classes in a hierarchy
 - Classes as module parameters
 - First-class classes
 - Mixins, Traits
 - Java's binary compatibility

Q: How much information about a base class is needed to compile its derived class?

Non-manifest base classes

- Base class not available for inspection when derived class is compiled
- Implementations use *dictionary* to map method/field names to their locations in the object layout
 - Dictionary lookup may occur at link time or run time

Example in Objective Caml

```
module type CIRCLE =  
sig type spec  
  class circle : spec ->  
    object method center : float*float  
      method radius : float  
    end  
end  
end
```

Example in Objective Caml

```
module CircleBBox =  
  functor (C : CIRCLE) ->  
  struct  
    class bbox arg = object (self)  
      inherit C.circle arg  
      method bounds =  
        let (x,y) = self#center in  
        let r = self#radius in  
        ((x-r,y-r), (x+r,y+r))  
    end  
  end
```

Non-manifest
base class

Where are
these methods?

Granularity of IL

- Method dispatch is *atomic* at source level: `x.m(...);`
- In IL, it should be decomposed into:
 - null check (if applicable)
 - dictionary lookup
 - method dereference
 - function call (indirect)

A generic IL for OO languages

- *Links* [Fisher et al., ESOP '00]
 - Its primitives express a “wide range of class-based OO features, including various forms of method dispatch.”
 - But, it is not typed.

Our contribution

- A sound & decidable type system for *Links*
 - using the *Type-Safe Certified Binaries* framework, “for explicitly representing complex propositions and proofs in typed ILs and object code” [Shao et al., POPL '02]
 - type and proof system based on *Calculus of Inductive Constructions* [Coquand & Paulin-Mohring, '90]

Outline

1. Review of *Links* primitives
2. Introduction to *TSCB* framework
3. Our computation language: LITL
4. Encoding objects

Links syntax

$$\begin{aligned} e ::= & x \mid n \mid e_1 + e_2 \mid \lambda x. e \mid e_1 e_2 \\ & \mid \langle e_1, \dots, e_n \rangle \mid e_1 @ e_2 \\ & \mid e_1 @ e_2 \leftarrow e_3 \mid e; \langle e_1, \dots, e_n \rangle \\ & \mid \{l_1 = e_1, \dots, l_n = e_n\} \mid e \# l \end{aligned}$$

Links syntax

$$\begin{aligned} e ::= & x \mid n \mid e_1 + e_2 \mid \lambda x. e \mid e_1 e_2 \\ & \mid \langle e_1, \dots, e_n \rangle \mid e_1 @ e_2 \\ & \mid e_1 @ e_2 \leftarrow e_3 \mid e; \langle e_1, \dots, e_n \rangle \\ & \mid \{l_1 = e_1, \dots, l_n = e_n\} \mid e \# l \end{aligned}$$

- Untyped lambda calculus

Links syntax

$$e ::= x \mid n \mid e_1 + e_2 \mid \lambda x. e \mid e_1 e_2$$
$$\mid \langle e_1, \dots, e_n \rangle \mid e_1 @ e_2$$
$$\mid e_1 @ e_2 \leftarrow e_3 \mid e; \langle e_1, \dots, e_n \rangle$$
$$\mid \{l_1 = e_1, \dots, l_n = e_n\} \mid e \# l$$

- Natural numbers and addition
 - to represent *offsets* of fields & methods

Links syntax

$$\begin{aligned} e ::= & x \mid n \mid e_1 + e_2 \mid \lambda x. e \mid e_1 e_2 \\ & \mid \langle e_1, \dots, e_n \rangle \mid e_1 @ e_2 \\ & \mid e_1 @ e_2 \leftarrow e_3 \mid e; \langle e_1, \dots, e_n \rangle \\ & \mid \{l_1 = e_1, \dots, l_n = e_n\} \mid e \# l \end{aligned}$$

- Tuple construction and projection
 - tuples represent objects and virtual function tables.
 - projection is *cheap*.

Links syntax

$$\begin{aligned} e ::= & x \mid n \mid e_1 + e_2 \mid \lambda x. e \mid e_1 e_2 \\ & \mid \langle e_1, \dots, e_n \rangle \mid e_1 @ e_2 \\ & \mid e_1 @ e_2 \leftarrow e_3 \mid e; \langle e_1, \dots, e_n \rangle \\ & \mid \{l_1 = e_1, \dots, l_n = e_n\} \mid e \# l \end{aligned}$$

- Functional update and extension
 - needed for overriding and adding new methods

Links syntax

$$\begin{aligned} e ::= & x \mid n \mid e_1 + e_2 \mid \lambda x. e \mid e_1 e_2 \\ & \mid \langle e_1, \dots, e_n \rangle \mid e_1 @ e_2 \\ & \mid e_1 @ e_2 \leftarrow e_3 \mid e; \langle e_1, \dots, e_n \rangle \\ & \mid \{l_1 = e_1, \dots, l_n = e_n\} \mid e \# l \end{aligned}$$

- Dictionary construction and lookup
 - lookup involves search (can be expensive)
 - used to map method labels to offsets

Compiling method dispatch

If the vtable is at offset 0 of object x ,
and d maps method labels to vtable offsets,
then the method dispatch

$$x.m(\dots)$$

expands to

$$((x @ 0) @ (d \# m)) x \dots$$

Typing this seems hopeless!

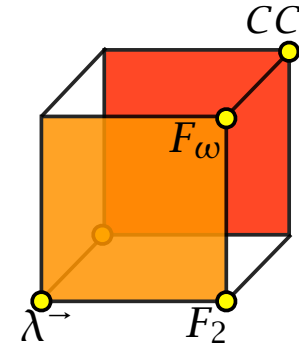
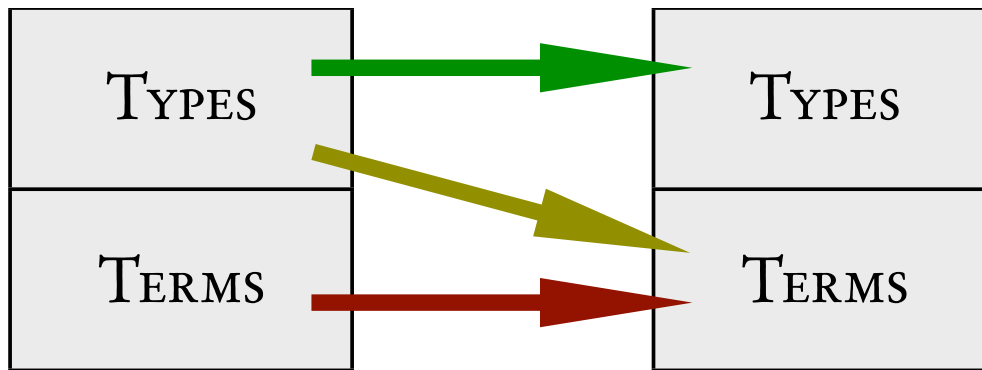
$$((x @ 0) @ (d \# m)) x \dots$$

- What we don't know *could* hurt us:
 - size/structure of vtable
 - offset returned by dictionary
 - is the index within bounds?
 - does it return a function pointer?
- Subtle connection between x and d
- Typed IL must capture these invariants

Outline

1. Review of *Links* primitives
2. Introduction to *TSCB* framework
3. Our computation language: LITL
4. Encoding objects

Dependencies in the λ cube



twice : $nat \rightarrow nat$

λ^-

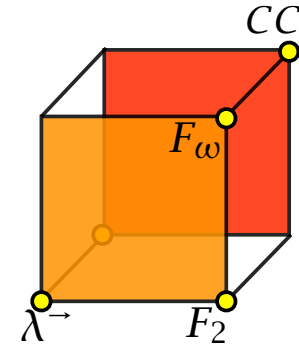
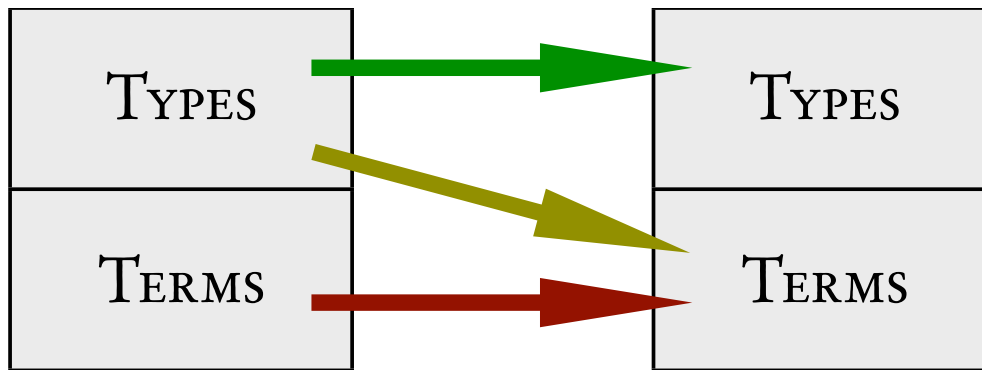
id : $\forall \alpha : \Omega. \alpha \rightarrow \alpha$

F_2

list : $\Omega \rightarrow \Omega$

F_ω

Traditional typed ILs stop here



$twice : nat \rightarrow nat$

λ^-

$id : \forall \alpha : \Omega. \alpha \rightarrow \alpha$

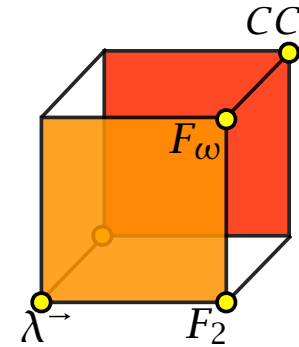
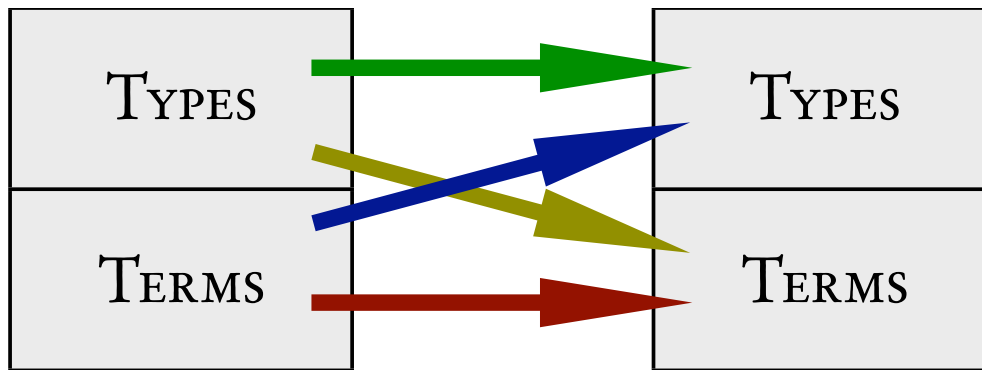
F_2

$list : \Omega \rightarrow \Omega$

F_ω

...and augment TERMS with *effects*
(like non-termination)

Dependencies in CC



twice : $nat \rightarrow nat$

λ^-

id : $\forall \alpha : \Omega. \alpha \rightarrow \alpha$

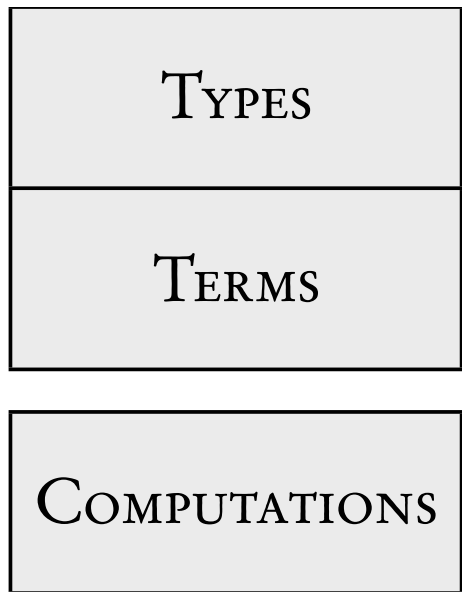
F_2

list : $\Omega \rightarrow \Omega$

array : $nat \rightarrow \Omega \rightarrow \Omega$

...now adding
effects is
dangerous!

Solution: add another layer



- The ‘types’ that govern computations are inductively defined TERMS in CIC.
 - We can use the *Coq Proof Assistant* as a type checker for CIC.
 - Sorts: SET, PROP, TYPE

Outline

1. Review of *Links* primitives
2. Introduction to *TSCB* framework
3. Our computation language: LITL*
4. Encoding objects

*LITL Is Typed Links

Preliminary definitions

Inductive $nat : SET \equiv$

| $O : nat$

| $S : nat \rightarrow nat.$

Parameter $sym : SET. (* \text{ to represent labels } *)$

Inductive $option (A : SET) : SET \equiv$

| $Some : A \rightarrow option A$

| $None : option A.$

Reasoning about things

Inductive $lt : nat \rightarrow nat \rightarrow \text{PROP} \equiv$

| $ltzs : \Pi n : nat. lt\ 0\ (S\ n)$

| $ltss : \Pi n\ m : nat. lt\ n\ m \rightarrow lt\ (S\ n)\ (S\ m).$

Inductive $eq\ (A : \text{TYPE})\ (x : A) : A \rightarrow \text{PROP} \equiv$

$refl_equal : eq\ A\ x\ x.$

Types for Links primitives

Inductive $Ty : SET \equiv$

| $arw : Ty \rightarrow Ty \rightarrow Ty$

| $snat : nat \rightarrow Ty$

| $tup : nat \rightarrow (nat \rightarrow Ty) \rightarrow Ty$

| $dict : (sym \rightarrow option Ty) \rightarrow Ty$

| $mu : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty$

| $all : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty$

| $ex : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty.$

Types for Links primitives

Inductive $Ty : SET \equiv$

| $arw : Ty \rightarrow Ty \rightarrow Ty$

| $snat : nat \rightarrow T$

| $tup : nat \rightarrow$

| $dict : (sym \rightarrow$

| $mu : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty$

| $all : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty$

| $ex : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty.$

The type of functions
(introduced by λ)

Types for Links primitives

Inductive $Ty : SET \equiv$

| $arw : Ty \rightarrow Ty \rightarrow Ty$

| $snat : nat \rightarrow Ty$

| $tup : nat \rightarrow (nat \rightarrow Ty) \rightarrow Ty$

| $dict : (sym \rightarrow Ty) \rightarrow Ty$

| $mu : \Pi k : SET$

| $all : \Pi k : SET$

| $ex : \Pi k : SET$

Singleton type for natural numbers:

$\vdash 0 : snat\ 0$

$\vdash 1 : snat\ (S\ 0)$

$\vdash 2 : snat\ (S\ (S\ 0))$

Types for Links primitives

Inductive $Ty : SET \equiv$

| $arw : Ty \rightarrow Ty \rightarrow Ty$

| $snat : nat \rightarrow Ty$

| $tup : nat \rightarrow (nat \rightarrow Ty) \rightarrow Ty$

| $dict : (sym \rightarrow nat \rightarrow Ty) \rightarrow Ty$

| $mu : \Pi k : SET$

| $all : \Pi k : SET$

| $ex : \Pi k : SET$

The type of tuples:

$\Delta; \Gamma \vdash e_1 : tup\ n\ f$

$\Delta; \Gamma \vdash e_2 : snat\ i$

$\Delta \vdash \sigma : lt\ i\ n$

$\Delta; \Gamma \vdash e_1 @ e_2 [\sigma] : f\ i$

Types for Links primitives

Inductive $Ty : SET \equiv$

| $arw : Ty \rightarrow Ty \rightarrow Ty$

| $snat : nat \rightarrow Ty$

| $tup : nat \rightarrow (nat \rightarrow Ty) \rightarrow Ty$

| $dict : (sym \rightarrow option Ty) \rightarrow Ty$

| $mu : \Pi k : SET (k \rightarrow Ty) \rightarrow Ty$

| $all : \Pi k : SET$

| $ex : \Pi k : SET$

The type of dictionaries:

$\Delta; \Gamma \vdash e : dict\ g$

$\Delta \vdash \sigma : eq\ (g\ l)\ (Some\ \tau)$

$\Delta; \Gamma \vdash e\#l[\sigma] : \tau$

Types for Links primitives

Inductive $Ty : SET \equiv$

| $arw : Ty \rightarrow Ty \rightarrow Ty$

| $snat : nat \rightarrow$

| $tup : nat \rightarrow$

| $dict : (sym \rightarrow option Ty) \rightarrow Ty$

| $mu : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty$

| $all : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty$

| $ex : \Pi k : SET. (k \rightarrow Ty) \rightarrow Ty.$

Recursive types, universal and existential quantifiers.

Type-annotated Links syntax

$e ::= x \mid n \mid e_1 + e_2 \mid \lambda x:\tau. e$

| $\Lambda \alpha:\sigma. e$ $\Delta \vdash \sigma : eq \tau_1 \tau_2$

| $\langle e_1, \dots, e_n \rangle$ $\Delta; \Gamma \vdash e : \tau_1$

| $e_1 @ e_2$ $\frac{\Delta; \Gamma \vdash e_1 : \tau_1 \quad \Delta; \Gamma \vdash e_2 : \tau_2}{\Delta; \Gamma \vdash e_1 @ e_2 : \tau_2}$

| $\{l_1 = e_1, \dots, l_n = e_n\}$ $\Delta; \Gamma \vdash e_i : \tau_i$

| $\text{cast} [\sigma] e$ | $[\tau_1, e \triangleright \tau_2]$

| $\text{open } e_1 \text{ as } [\alpha, x] \text{ in } e_2$

| $\text{fold } e \text{ as } \tau$ | $\text{unfold } e$

Outline

1. Review of *Links* primitives
2. Introduction to *TSCB* framework
3. Our computation language: LITL
4. Encoding objects

Recall: method dispatch in Links

$$(x @ (d \# m)) x \dots$$

Suppose that:

$$\Delta; \Gamma \vdash x : \text{tup } n f$$
$$\Delta; \Gamma \vdash d : \text{dict } g$$

These are unknown

The constraints on n, f, g are:

$$\Delta; \Gamma \vdash \sigma_1 : \text{eq } (g \ m) \ (\text{Some } (\text{snat } i))$$
$$\Delta; \Gamma \vdash \sigma_2 : \text{lt } i \ n$$
$$\Delta; \Gamma \vdash \sigma_3 : \forall \beta : \text{Ty}. \text{eq } (f \ \beta \ i) \ (\text{arw } \beta \ \tau)$$

These say “object has a method m at offset i returning type τ ”

Encapsulating the object rep.

Definition $Rep : \text{SET} \equiv$
 $(nat \times (Ty \rightarrow nat \rightarrow Ty) \times$
 $(sym \rightarrow option Ty)).$

Definition $size \equiv \lambda r : Rep.$
 $match\ r\ with\ (n, _, _) \Rightarrow n\ end.$

Definition $tupfn \equiv \lambda r : Rep.$
 $match\ r\ with\ (_, f, _) \Rightarrow f\ end.$

Definition $dictfn \equiv \lambda r : Rep.$
 $match\ r\ with\ (_, _, g) \Rightarrow g\ end.$

Expressing the constraints

Inductive *HasMethod*

$$\begin{aligned} & (r : \text{Rep}) (m : \text{sym}) (t : \text{Ty}) : \text{SET} \equiv \\ & \text{method} : \Pi i : \text{nat}. \text{lt } i \text{ (size } r) \rightarrow \\ & \text{eq (dictfn } r \text{ } m) (\text{Some (snat } i)) \rightarrow \\ & (\Pi \text{self}. \text{eq (tupfn } r \text{ self } i) (\text{arw self } t)) \rightarrow \\ & \text{HasMethod } r \text{ } m \text{ } t. \end{aligned}$$

The type of an object

$\exists r:Rep.$

$\exists p:HasMethod\ r\ m\ t.$

$\mu self:Ty.$

$tup\ (size\ r)\ (tupfn\ r\ self) \times dict\ (dictfn\ r)$

Conclusions

- *LITL* = sound, low-level typed IL with dictionaries & tuples
 - efficient target for various OO languages
 - even when layout of base class unknown
- *TSCB* is a powerful framework!
- A separate issue: encoding the subtype relationships of the source language
 - works in many cases; more work needed