

A Type-Preserving Compiler Infrastructure

Christopher League
Yale University
17 May 2002

Advisor

Zhong Shao

Committee

Kim Bruce
Arvind Krishnamurthy
Carsten Schürmann

Acknowledgment

Valery Trifonov

Thesis

- A strongly-typed compiler intermediate language can safely and efficiently accommodate very different programming languages

Christopher League

2

Mobile code, pervasive networks



- Wireless handheld computers
- Remotely programmable devices
- Browser applets



- Widely distributed computation

Christopher League

3

Security is critical



- We might not completely **trust** the programs we receive and run
- Must ensure they does not **misbehave**:
 - crash the device
 - exhaust resources
 - interfere with other programs/data
- Correctness is hard—focus on **safety**

Christopher League

4

Security toolbox: digital signature



- Confirms **identity** of producer, not **safety** of code
- I might not trust Microsoft, but would still run the code—assuming it is harmless

Christopher League

5

Security toolbox: reference monitor



impractical for monitoring fine-grained properties

- Code runs in a **sandbox**
 - interactions with outside world **mediated** by the monitor
- Hardware mechanisms
 - expensive context switches
 - not available on all devices
- Software rewriting
 - frequent dynamic checks

Christopher League

6

Security toolbox: language features

- Array bounds checking
- Garbage collection
- Exceptions
- Encapsulation, access control
- **Type systems**

Christopher League

7

Many vulnerabilities are type errors

Date	System	Kind
1/24/02	AOL ICQ	remotely exploitable buffer overflow
1/14/02	Solaris CDE	buffer overflow vulnerability
12/20/01	MS u-PNP	buffer overflow vulnerability
12/12/01	SysV 'login'	remotely exploitable buffer overflow
11/29/01	WU ftpd	format string vulnerability; free() on unallocated pointer
11/21/01	HP-UX lpd	remotely exploitable buffer overflow
10/25/01	Oraclegi AS	remotely exploitable buffer overflow
10/5/01	CDE ToolTalk	format string vulnerability

[CERT advisories]

Christopher League

8

The need for typed machine code

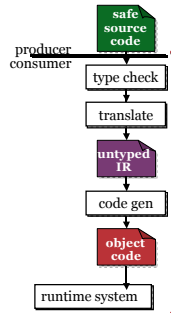
safe source code
(ML, Haskell, Java, C#)

- Is it enough to program in type-safe languages?

Christopher League

9

The need for typed machine code



- Is it enough to program in type-safe languages? **No.**

1. Microsoft is unlikely to ship **source code**
2. Still must **trust compiler**

[Thompson 1984]

- **Trusted Computing Base**
 - system modules responsible for **supporting security policy**

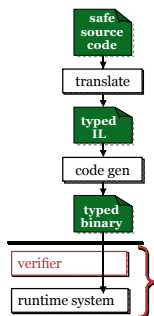
Christopher League

10

Vision: high assurance systems with minimal TCB

□ 1,000 LOC vs. 30,000+
(We are not there yet!)

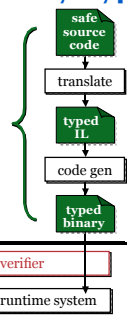
“Foundational proof-carrying code”
[Appel et al. 2001]



Christopher League

11

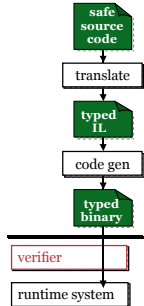
Key: type-preserving compiler



Christopher League

12

Until now: functional languages or "safe C"

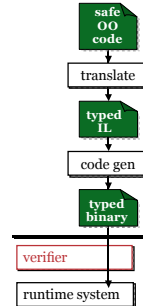


- TIL/ML [Tarditi, Morrisett et al. 1996]
- FLINT/ML [Shao 1997]
- Touchstone [Necula & Lee 1998]
- Popcorn
- Scheme-- [Morrisett et al. 1999]

Christopher League

13

Real world: object-oriented

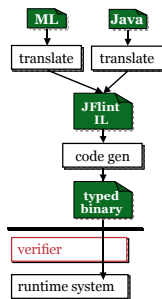


- Widely used (Java, C++)
- New, complex features
 - dynamic binding
 - dynamic casts
 - inheritance
 - interfaces
 - name-based types

Christopher League

14

Contributions



- Single typed IL that supports Java and ML
 - Not partial to either language
- Novel techniques for encoding most features of Java
 - Same efficiency as untyped
 - Preserves safety guarantees
- New high-level IL for Java
 - JVM primitives with functional flow
 - Easier to verify and optimize

Christopher League

15

Outline

- Mobile code security
- The need for typed machine code
- Dynamic dispatch security hole
- How our approach is different
- Safe and efficient object encoding
- Functional Java byte code
- A prototype compiler for Java and ML

Christopher League

16

Dynamic binding: essential to OO

- Inheritance without polymorphism is possible, but certainly **not very useful**.
- One can declare derived types, but the actual operation being called is **always known at compile time**.

[Booch 1994]

Christopher League

17

Efficient dynamic dispatch

```

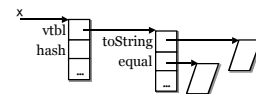
public static void
example ( Object x, Object y )
{ x.toString();

```

```

// compiles to:
(null check)
let r1 = x.vtbl ; // method suite
let r2 = r1.toString ; // method pointer
r2( x ) // "self application"
}

```



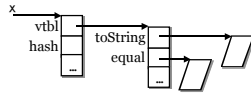
Christopher League

18

Must ensure the call is safe

```
public static void
example ( Object x, Object y )
{ x.toString();

  // compiles to:
  (null check)
  let r1 = x.vtbl;
  let r2 = r1.toString();
  r2(x); // this is sound
  r2(y)  // this is not!
}
```



Any unsoundness can be exploited

Christopher League

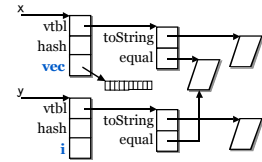
19

Cast arbitrary integer to pointer

```
class Ref extends Object
{ public byte[] vec;
  public String toString()
  { vec[13] = 42;
    return "Ha!";
  }
}
```

```
class Int extends Object
{ public int i;
}
```

```
example( new Ref(...),
         new Int(...));
```



```
public static void
example ( Object x, Object y )
{ let r1 = x.vtbl;
  let r2 = r1.toString();
  r2(y) // wrong!
}
```

Christopher League

20

This is a major security hole

in Cedilla Systems' PCC [Colby et al. 2000]

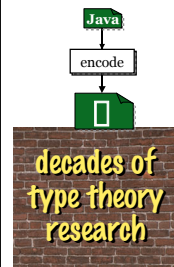
Special J

- Java compiler; generates annotated x86
- Types retain Java abstractions
 - (object_ty C)
 - (method_ty C SIG)
- Faulty axiom; undetected for 3+ years

Christopher League

21

Our approach is different



- Start with a strong foundation
 - “off the shelf” type theory
 - not specific to Java (or ML)
 - simple soundness proof
- Design complex encodings of Java features that
 - preserve type safety
 - maintain efficiency

Christopher League

23

Others have encoded objects in type theory

[Cardelli 84] [Cook et al. 89]
 [Pierce & Turner 94] [Bruce 94]
 [Hofmann & Pierce 95]
 [Abadi, Cardelli, Viswanathan 96]



- These are **models** of OOP
- Rather **inefficient**
 - extra indirections
 - extra function calls
- Assumed **subsumption**
 - a Circle is also a Shape

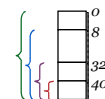
Christopher League

24

Type theory toolbox: rows

[Rémy 1993]

- Intuition: the **suffix** or **tail** of a record type, starting at a given offset



Christopher League

25

Type theory toolbox: recursion

- Intuition: μ notation for recursive definitions
 - list = { data: int, next: list }
- Replace recursive ref. with a **type variable**:
 - list = $\mu \square. \{ \text{data: int, next: } \square \}$
- let $x = \text{fold } y \text{ as list}$
- let $y : \{ \text{data: int, next: list } \} = \text{unfold } x$

Christopher League

26

Type theory toolbox: existentials

[Mitchell & Plotkin 1988]

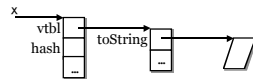
- Intuition: **hide** a type from outsiders
 - enforces abstract data types
- let $x1 : \square \square. \{ z : \square, f : \square \rightarrow \text{string} \}$
 = $\square \square = \text{int}, \{ z=42, f=\text{int2string} \}$
 : $\{ z : \square, f : \square \rightarrow \text{string} \} \square$
- open $x1$ as $\square \square. y : \{ z : \square, f : \square \rightarrow \text{string} \} \square$
 in $y.f (y.z)$

Christopher League

27

The type of Object

```
ObjTy[Object] =
  []fs::R8, ms::R4.
  μ self. ObjRcd[Object] fs ms self
```



```
ObjRcd[Object] fs ms self =
  { vtbl : { toString : self → string ;
            ms } ;
    hash : unsigned ;
    fs }
```

- Hide the **differences** between sub- and super-class.

Christopher League

28

Provably safe method invocation...

```
example (x, y : ObjTy[Object]) = rx
open x as []fx, mx, x1 : μ self. ObjRcd[Object] fx mx self [];
let x2 : ObjRcd[Object] fx mx rx = unfold x1 ;
let r1 : { toString : rx → string ; mx } = x2.vtbl ;
let r2 : rx → string = r1.toString ;
r2(x1)
```

```
ObjTy[Object] =
  []fs::R8, ms::R4.
  μ self. ObjRcd[Object] fs ms self
ObjRcd[Object] fs ms self =
  { vtbl : { toString : self → string ;
            ms } ;
    hash : unsigned ;
    fs }
```

Christopher League

29

...without sacrificing efficiency

```
example (x, y : ObjTy[Object]) = rx
open x as []fx, mx, x1 : μ self. ObjRcd[Object] fx mx self [];
let x2 : ObjRcd[Object] fx mx rx = unfold x1 ;
let r1 : { toString : rx → string ; mx } = x2.vtbl ;
let r2 : rx → string = r1.toString ;
r2(x1)
```

```
ObjTy[Object] =
  []fs::R8, ms::R4.
  μ self. ObjRcd[Object] fs ms self
ObjRcd[Object] fs ms self =
  { vtbl : { toString : self → string ;
            ms } ;
    hash : unsigned ;
    fs }
```

Christopher League

30

Techniques extend to most of Java

- classes
 - inheritance
 - dynamic dispatch
 - dynamic cast
 - mutual recursion
 - interfaces
 - constructors
 - super calls
 - subroutines
 - exceptions
 - privacy
- Chapters 3–5 contain:
 - Formal definition of source language (Featherweight Java)
 - Formal definition of intermediate language (JFlint)
 - Proofs that JFlint is **sound** and **decidable**
 - Type-directed translation
 - Proof that well-typed inputs yield well-typed outputs

Christopher League

31

Act III

- Mobile code security
- The need for typed machine code
- Dynamic dispatch security hole
- How our approach is different
- Safe and efficient object encoding
- Functional Java byte code
- A prototype compiler for Java and ML

Christopher League

32

System building

- A prototype **compiler** for Java and ML
- Many **practical** problems must be solved
 - Efficient implementation of IL
 - Large **semantic gap** between Java and JFlint

Christopher League

33

Where to start?

- Java
 - ↓ javac
 - x.println(y);
- JVMML byte code
 - 3 aload_0 # this
 - 4 getfield PrintStream C::x
 - 7 dload 2
 - 9 invokevirtual
 - void PrintStream::println(double)
- Many details are **not explicit** in Java source
- Java byte code has **untyped** local vars & **implicit** data flow

Christopher League

34

Two sets of concerns

1. data & control flow, type inference
2. expanding Java primitives

- JVMML byte code → JFlint
- 3 aload_0 # this
- 4 getfield PrintStream C::x
- 7 dload 2
- 9 invokevirtual
- void PrintStream::println(double)

Christopher League

35

A new IL to bridge the gap

- High-level Java primitives, types
- JVMML → \square^{JVM} → JFlint
- Functional control and data flow

Christopher League

36

A better Java byte code

- JVMML → \square^{JVM} → JFlint
- Fully explicit
 - Supports all of JVMML, yet is
 - Easier to verify and optimize
- Nastiest parts of JVM become **tractable**
 - Object initialization [Freund & Mitchell 1999]
 - Subroutines [Stata & Abadi 1998]
- Verification is just simple type checking
 - < 260 lines of ML code

[Chapter 7]

Christopher League

37

Example: Factorial

```
public int fact(int n)
{ int x;
  for (x = 1; n > 0; n--)
    x = x * n;
  return x;
}
```

```
public int fact (int n)
{  iconst_1
   istore_1 # x=1
   goto T
  L iload_1
   iload_0
   imul
   istore_1 # x=x * n
   iinc 0 -1 # n--
  T iload_0
   ifgt L # n > 0
   iload_1
   ireturn
}
```

Christopher League

38

Example: Factorial

```
public int fact(int n) =
letrec L = [(i : I, x : I) .
  let y = x * i;
  let j = i - 1;
  T(j, y)
and T = [(k : I, z : I) .
  if n > 0 then L(k, z)
  else return z
in T(n, 1)
```

```
public int fact (int n)
{  iconst_1
   istore_1 # x=1
   goto T
  L iload_1
   iload_0
   imul
   istore_1 # x=x * n
   iinc 0 -1 # n--
  T iload_0
   ifgt L # n > 0
   iload_1
   ireturn
}
```

Christopher League

39

Subroutines are tricky

- **jsr offset** — push return address, jump
- **ret var** — return to address in *var*
- Used to implement ‘finally’ blocks


```
try { A }
catch (Error e) { B; throw e }
finally { C }
```
- Can achieve complicated control flow

Christopher League

40

Continuation-passing style

[Steele 1978] [Kranz et al. 1986]

- The **higher-order** answer to flexible control flow
 - Represent return address as a **function**

Christopher League

41

1. ret need not obey stack discipline

```
void main(String[] args)
{ try {
  finally
  { while(true)
    { try {
      finally { break; }
    }
  }
}
```

```
M jsr A
  goto R
A astore_1 # return addr
L jsr B
  goto L
B pop # return addr
  ret 1
R return
```

Christopher League

42

1. ret need not obey stack discipline

```
letrec M = [( ) . A(R)
and A = [(k1 : () -> V) .
  letrec L = [( ) . B(L)
and B = [(k2 : () -> V) .
  ki()
  in L()
and R = [( ) . return
in M()
```

```
M jsr A
  goto R
A astore_1 # return addr
L jsr B
  goto L
B pop # return addr
  ret 1
R return
```

Christopher League

43

2. Subroutine might update local var

```

letrec M = [] . S(43, P)
and S = [(i: I, k: (I)→V) .
  let j = i - 1;
  k(j)
and P = [(i: I) .
  invoke printInt (i);
  return
in M()

M ldc 43
  istore_1 # i=43
  jsr S
  goto P
S astore_2 # return addr
  iinc 0 -1 # i--
  ret 2
P iload_1 # print i
  invoke printInt
  return

```

Christopher League

44

3. Polymorphic over untouched vars

```

letrec S = [(k: ()→V) . k()] in
letrec M = [] .
  let x1 = 3.14;
  letrec I = [] .
    invoke printFloat (x1);
    let x2 = 42;
    letrec R = [] .
      invoke printInt (x2);
      return
    in S(R)
  in S(I)
in M()

M ldc 3.14
  fstore_1 # x=3.14
  jsr S
I fload_1
  invoke printFloat
  ldc 42
  istore_1 # x=42
  jsr S
R iload_1
  invoke printInt
  return
S astore_2 # return addr
  ret 2

```

Christopher League

45

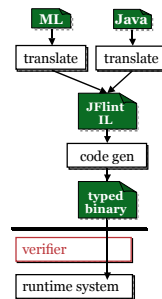
System overview

- Based on SML/NJ compiler 110.30
 - added a new type-preserving Java front end
 - interactively loads and runs Java classes
- Same back end and runtime system
- Front end \square 8k LOC; JFlint checker \square 1k LOC
- Runs CaffeineMark 3.0 (12 classes, \square 100k)
- Compile time just 60% longer than gcj
- Does **not** load native code!

Christopher League

46

Synergy with ML front end



- The type system is **impartial**

JFlint	Java	ML
\square	inheritance	parametric poly.
\square	object enc.	closures
μ	rec. classes	rec. datatypes
tags	dynamic cast	exceptions
rows	object enc.	—
records	vtable, objects	records, tuples
functions	methods	functions

Christopher League

47

Summary

- Mobile code security is **critical**
- High-assurance systems need **minimal** TCB
- Type-preserving compilers are the key
- With care, they scale to **real** languages
 - use type theory as the foundation;
 - focus on practical encodings
- A single **typed IL** can safely and efficiently support different kinds of source languages

Christopher League

48