

Type-Based Verification of Assembly Language for Compiler Debugging

Bor-Yuh Evan Chang Adam Chlipala
George C. Necula Robert R. Schneck

University of California, Berkeley

TLDI 2005
Long Beach, California
January 10, 2005

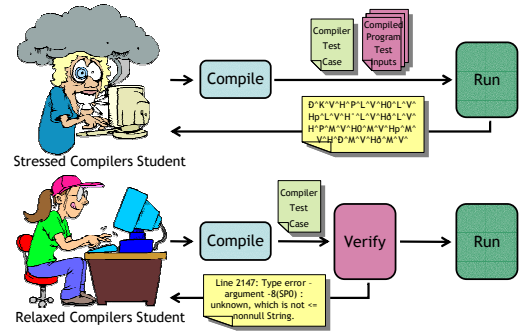
Problem and Motivation

- Type check **assembly code** produced by **existing** compilers for a Java-like language
- Validate that certifying compilation aids compiler debugging
 - using undergraduate compilers course as a test bed
- Introduce undergraduates to using compiler technology for software quality

Problem and Motivation

- Starting Point: Bytecode verification
 - Effective, simple, and very accessible
 - But does not work for assembly code
- Why assembly code?
 - Native code compilers are more complex
 - Also debug JITs
- Why existing compilers?
 - Force to make the verifier flexible
 - Better accessibility to students (i.e., require little or no annotations)
 - More test cases to validate “certified compilation for debugging”

Debugging Compilers with Verification



Challenges

```

class P {
  int f;
  int m() { ... }
}
class C extends P {
  int m() { ... }
}
...
P p = new P();
P c = new C();
c.m();
...
    
```

invokevirtual P.m()

```

branch (= r_c 0) L_abort
r_tmp := m[r_c + 8]
r_tmp := m[r_tmp + 12]
r_arg0 := r_c
r_ra := &L_ret
jump [r_tmp]
L_ret:
    
```

Challenges

```

class P {
  int f;
  int m() { ... }
}
class C extends P {
  int m() { ... }
}
...
P p = new P();
P c = new C();
c.m();
...
    
```

invokevirtual P.m()

```

branch (= r_c 0) L_abort (r_c : P, ...)
r_tmp := m[r_c + 8] (r_c : nonnull P, ...)
r_tmp := m[r_tmp + 12] (r_tmp : disp(P), ...)
r_arg0 := r_c (r_tmp : meth(P,12), ...)
r_ra := &L_ret (r_arg0 : P, ...)
jump [r_tmp]
L_ret: (r_rv : int, ...)
    
```

Challenges

```

class P {
  int f;
  int m() { ... }
}
class C extends P {
  int m() { ... }
}
...
P p = new P();
P c = new C();
c.m();
...

```

```

invokevirtual P.m()
branch (= r_c 0) L_abort
r_tmp := m[r_c + 8]
r_tmp := m[r_tmp + 12]
r_arg0 := r_p
r_ra := &L_ret
jump [r_tmp]
L_ret:

```

unsound

$\langle r_c : P, \dots \rangle$
 $\langle r_c : \text{nonnull } P, \dots \rangle$
 $\langle r_{tmp} : \text{disp } P, \dots \rangle$
 $\langle r_{tmp} : \text{methid } P, \dots \rangle$
 $\langle r_{rv} : \text{int}, \dots \rangle$

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 7

Challenges

```

class P {
  int f;
  int m() { ... }
}
class C extends P {
  int m() { ... }
}
...
P p = new P();
P c = new C();
int f = p.f;
p.m();
x = f + 1;

```

```

branch (= r_p 0) L_abort
r_tmp := r_p + 12
r_f := m[r_tmp]
branch (= r_p 0) L_abort
r_tmp := m[r_p + 8]
r_tmp := m[r_tmp + 12]
r_arg0 := r_p
r_ra := &L_ret
jump [r_tmp]
L_ret:
r_x := r_f + 1

```

reordering and optimization

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 8

Challenges

```

class P {
  int f;
  int m() { ... }
}
class C extends P {
  int m() { ... }
}
...
P p = new P();
P c = new C();
int f = p.f;
p.m();
x = f + 1;

```

```

branch (= r_p 0) L_abort
r_tmp := r_p + 12
r_f := m[r_tmp]
r_x := r_f + 1
branch (= r_p 0) L_abort
r_tmp := m[r_tmp - 4]
r_tmp := m[r_tmp + 12]
r_arg0 := r_p
r_ra := &L_ret
jump [r_tmp]
L_ret:

```

“funny” pointer arithmetic

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 9

Outline

- Overview
- Abstract State
 - Types
 - Join algorithm
- Verification Procedure
- Educational Experience
- Concluding Remarks

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 10

Overview

- To maximize accessibility (i.e., minimize annotations)
 - infer types at intermediate program points using abstract interpretation
 - willingness to have a limited amount of specialization to a compilation strategy
- To handle assembly code
 - use (simple) dependent types
- To be less sensitive to the compilation strategy
 - assign types to intermediate values lazily

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 11

Abstract State

abstract state $\langle \Sigma ; \Gamma \rangle$

value state $\Sigma ::= \mathbf{x}_0 = e_0, \mathbf{x}_1 = e_1, \dots, \mathbf{x}_{n-1} = e_{n-1}$

type state $\Gamma ::= \cdot \mid \Gamma, \alpha : \tau$

symbolic values α, β

expressions $e ::= n \mid \alpha \mid \&L \mid e_0 + e_1 \mid e_0 \cdot e_1 \mid \dots$

- Keep intermediate expressions in symbolic form
- **Symbolic values** abstract irrelevant expressions keeping only type information
- Intermediate expressions track many dependencies between registers

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 12

Why Symbolic Values?

- Used to track dependencies between registers

$$r_1 = r_2 \wedge r_1 = r_3 + 4$$

↓

$$r_1 = \alpha + 4, r_2 = \alpha + 4, r_3 = \alpha$$

- Value state form convenient for abstract interpretation

$$r_1 := r_2$$

$$\Sigma$$

$$\Sigma[r_1 \rightarrow \Sigma(r_2)]$$

1/10/2005

Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging

13

Values

values $v ::= n_0 \cdot \&L + n_1 \cdot \alpha + n_2 \mid \alpha = n \mid \alpha \neq n \mid \alpha < n \mid \dots$

- Define a normalization of expressions to values
 - $\Gamma \vdash e \Downarrow v \triangleright \Gamma'$
- Values give the expressions of interest
 - For our case, address computation and comparisons with constants
 - Would vary depending on the source language or compiler of interest
 - Not all equal expressions normalize to the same value, so value forms must be chosen carefully
- Registers are equal if they map to the same values (after normalization)

1/10/2005

Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging

14

Types

- Primitive Types and Object References

types $\tau ::= T \mid \text{word} \mid \perp \mid [\text{nonnull}] b \mid \tau \text{ ptr} \mid \dots$
 bounds $b ::= C$ bounded above by C
 | exactly C bounded above and below by C
 | classof(α) same class as α

- Types for Run-time Structures

types $\tau ::= \dots$
 | disp(α)
 | meth(α, n)
 | ...

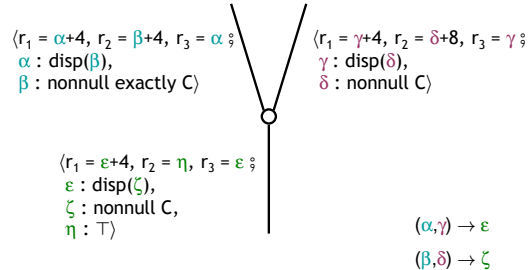
dispatch table of α
 Tracks that the value is the dispatch table of object α

1/10/2005

Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging

15

Join



1/10/2005

Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging

16

Join

- Joining of type state basically given by subtyping
 - subtyping lattice still fairly simple - dictated mostly by the class inheritance hierarchy
- Values are (fancy) labels for equivalence classes of registers
 - lattice of conjunctions of equalities on registers
 - first normalize expressions to values for join
 - each distinct pair of values in the input state gives an equivalence class in the result state

1/10/2005

Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging

17

Outline

- Overview
- Abstract State
 - Types
 - Join algorithm
- Verification Procedure
- Educational Experience
- Concluding Remarks

1/10/2005

Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging

18

Verification Example

```

class P {
  int f;
  int m() { ... }
}
class C extends P
{ ... }

```

branch ($= r_p 0$) L_{abort} ($r_p = \rho \S \rho : P$)
 $\langle \dots \S \rho : \text{nonnull } P \rangle$
 $r_{\text{tmp}} := r_p + 12$ ($r_{\text{tmp}} = \rho + 12 \S \dots$)
 $r_f := m[r_{\text{tmp}}]$ ($r_f = \phi \S \phi : \text{int}$)
 $r_x := r_f + 1$ ($r_x = \phi + 1 \S \phi : \text{int}$)
 $r_{\text{tmp}} := m[r_{\text{tmp}} - 4]$
 $r_{\text{tmp}} := m[r_{\text{tmp}} + 12]$
 $r_{\text{arg0}} := r_p$
 $r_{\text{ra}} := \&L_{\text{ret}}$
 jump $[r_{\text{tmp}}]$
 $L_{\text{ret}}:$

Typing Rule

$$\Gamma \vdash e \Downarrow \alpha + n \triangleright \Gamma'$$

$$\Gamma' \vdash \alpha : \text{nonnull } C \triangleright \Gamma''$$

(field at offset n of class C has type τ)

$$\Gamma \vdash e : \tau \text{ ptr} \triangleright \Gamma'$$

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 19

Verification Example

```

class P {
  int f;
  int m() { ... }
}
class C extends P
{ ... }

```

branch ($= r_p 0$) L_{abort} ($r_p = \rho \S \rho : P$)
 $\langle \dots \S \rho : \text{nonnull } P \rangle$
 $r_{\text{tmp}} := r_p + 12$ ($r_{\text{tmp}} = \rho + 12 \S \dots$)
 $r_f := m[r_{\text{tmp}}]$ ($r_f = \phi \S \phi : \text{int}$)
 $r_x := r_f + 1$ ($r_x = \phi + 1 \S \phi : \text{int}$)
 $r_{\text{tmp}} := m[r_{\text{tmp}} - 4]$ ($r_{\text{tmp}} = \theta \S \theta : \text{disp}(\rho)$)
 $r_{\text{tmp}} := m[r_{\text{tmp}} + 12]$ ($r_{\text{tmp}} = \sigma \S \sigma : \text{meth}(\rho, 12)$)
 $r_{\text{arg0}} := r_p$ ($r_{\text{arg0}} = \rho \S \dots$)
 $r_{\text{ra}} := \&L_{\text{ret}}$ (Calling $\text{meth}(\rho, 12)$ checks $r_{\text{arg0}} = \rho$)
 jump $[r_{\text{tmp}}]$ ($r_{\text{rv}} = \mu \S \mu : \text{int}$)
 $L_{\text{ret}}:$

Typing Rule

$$\Gamma \vdash e \Downarrow \alpha + 8 \triangleright \Gamma'$$

$$\Gamma' \vdash \alpha : \text{nonnull } C \triangleright \Gamma''$$

$$\Gamma \vdash e : \text{disp}(\alpha) \text{ ptr} \triangleright \Gamma''$$

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 20

Outline

- Overview
- Abstract State
 - Types
 - Join algorithm
- Verification Procedure
- Educational Experience
- Concluding Remarks

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 21

Comparing Compiler Development

Without Coolaid With Coolaid

Legend:

- test passed, type safe (observably correct)
- test passed, type safe but Coolaid failed (incompleteness)
- test passed, type error (hidden type error)
- test failed, type error (visible type error)
- test failed, type safe (semantic error)

• Good compilations increase: 53% to 75%

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 22

Comparing Compiler Development

Without Coolaid With Coolaid

Legend:

- test passed, type safe (observably correct)
- test passed, type safe but Coolaid failed (incompleteness)
- test passed, type error (hidden type error)
- test failed, type error (visible type error)
- test failed, type safe (semantic error)

• Type errors decrease: 44% to 19%
 - even 29% to 15% if only consider those that are visible in testing

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 23

Comparing Compiler Development

Without Coolaid With Coolaid

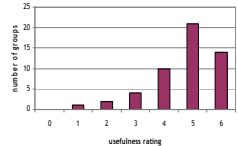
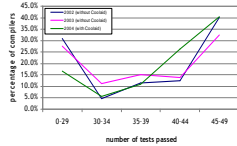
Legend:

- test passed, type safe (observably correct)
- test passed, type safe but Coolaid failed (incompleteness)
- test passed, type error (hidden type error)
- test failed, type error (visible type error)
- test failed, type safe (semantic error)

• False alarms: 6% to 1%
 - price for ease of use

1/10/2005 Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging 24

Student Perspective



- Traditional testing procedure
 - 49 test cases / inputs
- Mean scores:
 - 33 (67%) in 2002
 - 34 (69%) in 2003
 - 39 (79%) in 2004
- 0: "counter-productive"
- 6: "can't imagine being without it"
- A favorite comment:
"I would be totally lost without Coolaid. I learn best when I am using it hands-on I was able to really understand stack conventions and optimizations and to appreciate them."

1/10/2005

Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging

25

Conclusion

- Extended data-flow based type inference of intermediate languages to assembly code
 - Able to retrofit existing compilers
 - Minimal annotations for accessibility
- Provided experimental confirmation that certifying compilation helps early compiler debugging
- Introduced undergraduates to the idea of improving software quality with program analysis

1/10/2005

Chang et al.: Type-Based Verification of Assembly Language for Compiler Debugging

26