

Decompilation into Logic — Improved

Magnus O. Myreen, Michael J. C. Gordon
Computer Laboratory, University of Cambridge, UK

Konrad Slind
Rockwell Collins, USA

FMCAD'12 Cambridge

Machine Code

This talk is about machine-code verification.

```
0: E3A00000  
4: E3510000  
8: 12800001  
12: 15911000  
16: 1AFFFFFB
```

Verification of Machine Code

Challenges:

machine code



code

Verification of Machine Code

Challenges:

machine code

code

correctness

$\{P\}$ code $\{Q\}$

Verification of Machine Code

Challenges:

machine code

code

ARM/x86/PowerPC model
(1000...10,000 lines each)

correctness

{P} code {Q}

Verification of Machine Code

Challenges:

machine code

code

ARM/x86/PowerPC model
(1000...10,000 lines each)

correctness

{P} code {Q}

Contribution: tools/methods which

- expose as little as possible of the big models to the user
- makes non-automatic proofs independent of the models

Decompilation into Logic

Example: given some (hard-to-read) ARM machine code

```
0: E3A00000      mov r0, #0
4: E3510000      L: cmp r1, #0
8: 12800001      addne r0, r0, #1
12: 15911000      ldrne r1, [r1]
16: 1AFFFFFB      bne L
```

Decompilation into Logic

Example: given some (hard-to-read) ARM machine code

```
0: E3A00000      mov r0, #0
4: E3510000      L: cmp r1, #0
8: 12800001      addne r0, r0, #1
12: 15911000      ldrne r1, [r1]
16: 1AFFFFFB      bne L
```

our decompiler produces a readable function in HOL:

$$f(r_0, r_1, m) = \text{let } r_0 = 0 \text{ in } g(r_0, r_1, m)$$
$$g(r_0, r_1, m) = \text{if } r_1 = 0 \text{ then } (r_0, r_1, m) \text{ else}$$
$$\quad \text{let } r_0 = r_0 + 1 \text{ in}$$
$$\quad \text{let } r_1 = m(r_1) \text{ in}$$
$$\quad g(r_0, r_1, m)$$

Certificate Theorems

Decompiler **automatically** proves a certificate theorem, which states that **f** describes the effect of the ARM code:

$$f_{pre}(r_0, r_1, m) \Rightarrow$$

$$\{ (R0, R1, M) \text{ is } (r_0, r_1, m) * \text{PC } p * S \}$$

$$p : \text{E3A00000 E3510000 12800001 15911000 1AFFFFF B}$$

$$\{ (R0, R1, M) \text{ is } f(r_0, r_1, m) * \text{PC } (p + 20) * S \}$$

Read informally:

if initially reg 0, reg 1 and memory described by (r_0, r_1, m) , then the code **terminates** with reg 0, reg 1, memory as $f(r_0, r_1, m)$

Preconditions

Precondition f_{pre} keeps track of side conditions:

$$f_{pre}(r_0, r_1, m) = \text{let } r_0 = 0 \text{ in } g_{pre}(r_0, r_1, m)$$

$$g_{pre}(r_0, r_1, m) = \text{if } r_1 = 0 \text{ then } \mathit{true} \text{ else}$$
$$\quad \text{let } r_0 = r_0 + 1 \text{ in}$$
$$\quad \text{let } \mathit{cond} = r_1 \in \text{domain } m \wedge \mathit{aligned}(r_1) \text{ in}$$
$$\quad \text{let } r_1 = m(r_1) \text{ in}$$
$$\quad g_{pre}(r_0, r_1, m) \wedge \mathit{cond}$$

Decompilation into Logic

Strengths:

- ☑ **separates** definition of ISA model from program verification (program verification is done based on decompiler output)
- ☑ can implement **proof-producing program** synthesis from HOL based on decompilation (translation validation)
- ☑ has been shown to scale to **large verification projects**:
 - ▶ functional correctness of garbage collectors,
 - ▶ Lisp implementations (ARM, x86, PowerPC), and
 - ▶ the seL4 microkernel (approx. 12,000 lines of ARM)

Performance Ignored...

Weaknesses:

- **unnecessarily complicated** (for historical reasons, uses ideas from separation logic that are irrelevant to decompilation)
- sometimes **very slow** (never optimised for speed, separation logic composition rule slow to execute in LCF-style prover)
 - ▶ decompilation of the seL4 microkernel (approx. 12,000 lines of ARM) takes 8 hours (for gcc -O2 output)
- **not applicable** to code with general-purpose **code pointers**, e.g. jump to code pointer.

HOL: fully-expansive LCF-style prover

Proofs are performed in HOL4 — a **fully expansive** theorem prover

HOL4 theorem prover
(incl. decompiler)

HOL4 kernel

All proofs expand at runtime into primitive inferences in the HOL4 kernel.

The kernel implements the axioms and inference rules of higher-order logic.

HOL: fully-expansive LCF-style prover

Proofs are performed in HOL4 — a **fully expansive** theorem prover

HOL4 theorem prover
(incl. decompiler)

HOL4 kernel

All proofs expand at runtime into primitive inferences in the HOL4 kernel.

The kernel implements the axioms and inference rules of higher-order logic.

Example: proving $!0+1=1!$ using the simplifier requires 85 primitive inferences (0.0003 seconds). No hope of producing a very fast tool...

This Talk:

Improving Decompilation

Weaknesses of old approach:

- unnecessarily complicated
- sometimes very slow
- cannot handle code pointers

This Talk:

Improving Decompilation

Weaknesses of old approach:

- unnecessarily complicated
- sometimes very slow
- cannot handle code pointers

Contribution:

- ☑ simpler Hoare logic
- ☑ revised approach for better speed
- ☑ support for code pointers

New Hoare triple

Decompiler performs proofs in terms Hoare triples:

$$\{ pre \} code \{ post \}.$$

New Hoare triple

Decompiler performs proofs in terms Hoare triples:

$$\{ pre \} code \{ post \}.$$

Semantics parametrised by *assert* and *next*:

$$\forall state\ c. \text{assert}(code \cup c, pre)\ state \implies \\ \exists n. \text{assert}(code \cup c, post)\ (next^n(state))$$

New Hoare triple

Decompiler performs proofs in terms Hoare triples:

$$\{ pre \} code \{ post \}.$$

Semantics parametrised by *assert* and *next*:

$$\forall state c. \text{assert} (code \cup c, pre) state \implies \\ \exists n. \text{assert} (code \cup c, post) (next^n(state))$$

We instantiate these for each architecture (ARM, x86, PowerPC), e.g.

$$\text{arm_assert} (code, pc, r_0, r_1, \dots, cond) state = \\ (cond \implies code \text{ is in memory of } state \text{ and} \\ \text{the PC of } state \text{ is } pc \text{ and } \dots)$$

New Hoare triple

Decompiler performs proofs in terms Hoare triples:

$$\{ pre \} code \{ post \}.$$

Semantics parametrised by *assert* and *next*:

$$\forall state\ c. \text{assert} (code \cup c, pre)\ state \implies \\ \exists n. \text{assert} (code \cup c, post)\ (next^n(state))$$

We instantiate these for

pre/post are tuples, allows fast composition

$$\text{arm_assert} (code, pc, r_0, r_1, \dots, cond)\ state = \\ (cond \implies \text{code is in memory of } state \text{ and} \\ \text{the PC of } state \text{ is } pc \text{ and } \dots)$$

New Hoare triple

Decompiler performs proofs in terms Hoare triples:

$$\{ pre \} code \{ post \}.$$

Semantics parametrised by *assert* and *next*:

$$\forall state c. \text{assert} (code \cup c, pre) state \implies \\ \exists n. \text{assert} (code \cup c, post) (next^n(state))$$

We instantiate these for

pre/post are tuples, allows fast composition

$$\text{arm_assert} (code, pc, r_0, r_1, \dots, cond) state = \\ (cond \implies \text{code is in memory of } state \text{ and} \\ \text{the PC of } state \text{ is } pc \text{ and } \dots)$$

side-condition is accumulated in 'postcondition'

New Hoare triple

Decompiler performs proofs in terms Hoare triples:

$$\{pre\} code \{post\}$$

Semantics parametrised by *assert* and *next*:

$$\forall state c. \text{assert} (code \cup c, pre) state \implies$$

program counter value is explicitly part of pre/post

We instantiate these for $pre/post$ are tuples, allows fast composition

$$\text{arm_assert} (code, pc, r_0, r_1, \dots, cond) state =$$

$(cond \implies code \text{ is in memory of } state \text{ and}$
 $\text{the PC of } state \text{ is } pc \text{ and } \dots)$

side-condition is accumulated in 'postcondition'

Function extraction

Decompilation algorithm:

- Step 1:** evaluate underlying ISA model
(prove Hoare triples for each instruction)
- Step 2:** construct CFG and find ‘decompilation rounds’
(usually one round per loop)
- Step 3:** for each round, compose a Hoare triple theorem:

$$\begin{array}{l} \{pre[v_0 \dots v_n]\} \\ \text{code} \\ \{\text{let } (v'_0 \dots v'_n) = f(v_0 \dots v_n) \text{ in } post[v'_0 \dots v'_n]\} \end{array}$$

if the code contains a loop, apply a loop rule (next slide...)

Function extraction

Decompilation algorithm:

- Step 1:** evaluate underlying ISA model
(prove Hoare triples for each instruction)
- Step 2:** construct CFG and find ‘decompilation rounds’
(usually one round per loop)
- Step 3:** for each round, compose a Hoare triple theorem:

$$\begin{array}{l} \{pre[v_0 \dots v_n]\} \\ code \\ \{\text{let } (v'_0 \dots v'_n) = f(v_0 \dots v_n) \text{ in } post[v'_0 \dots v'_n]\} \end{array}$$

if the code

collects both update and side-conditions

Loop Rule

If there are loops in the code then apply:

$$\begin{aligned} & (\forall x. \{pre\ x\} \text{code} \{\text{if } g\ x \text{ then } pre\ (f\ x) \text{ else } post\ (d\ x)\}) \\ & \implies \\ & (\forall x. \{pre\ x\} \text{code} \{post\ (\text{tailrec } g\ f\ d\ x)\}) \end{aligned}$$

where tailrec is a function format that satisfies:

$$\text{tailrec } g\ f\ d\ x = \text{if } g\ x \text{ then tailrec } g\ f\ d\ (f\ x) \text{ else } d\ x$$

Definition of tailrec is in the paper.

Example

Assembly code:

```
L0: ldr r1, [r2, r3] ; load mem[r2+r3] into r1
L1: add r0, r1      ; add r1 to r0
L2: subs r3, #4    ; decrement r3 by 4
L3: bne L0         ; goto L0 if r3 ≠ 0
L4:
```

Example

Assembly code:

```
L0: ldr r1, [r2, r3] ; load mem[r2+r3] into r1
L1: add r0, r1      ; add r1 to r0
L2: subs r3, #4    ; decrement r3 by 4
L3: bne L0         ; goto L0 if r3 ≠ 0
L4:
```

Extracted function:

```
sum(cond, r0, r1, r2, r3, m) =
  let cond = cond ∧ valid_address (r2 + r3) m in
  let r1 = m(r2 + r3) in
  let r0 = r0 + r1 in
  let r3 = r3 - 4 in
  if r3 = 0 then (cond, r0, r1, r2, r3, m)
  else sum(cond, r0, r1, r2, r3, m)
```

Example

Assembly code:

```
L0: ldr r1, [r2, r3] ; load mem[r2+r3] into r1
L1: add r0, r1      ; add r1 to r0
L2: subs r3, #4    ; decrement r3 by 4
L3: bne L0         ; goto L0 if r3 ≠ 0
L4:
```

Extracted function:

```
sum(cond, r0, r1, r2, r3, m) =
  let cond = cond ∧ valid_address (r2 + r3) m in
  let r1 = m(r2 + r3) in
  let r0 = r0 + r1 in
  let r3 = r3 - 4 in
    if r3 = 0 then (cond, r0, r1, r2, r3, m)
    else sum(cond, r0, r1, r2, r3, m)
```

side-conditions part of function

Example

Assembly code:

```
L0: ldr r1, [r2, r3] ; load mem[r2+r3] into r1
L1: add r0, r1      ; add r1 to r0
L2: subs r3, #4    ; decrement r3 by 4
L3: bne L0         ; goto L0 if r3 ≠ 0
L4:
```

Extracted function:

```
sum(cond, r0, r1, r2, r3, m) =
  let cond = cond ∧ valid_address (r2 + r3) m in
  let r1 = m(r2 + r3) in
  let r0 = r0 + r1 in
  let r3 = r3 - 4 in
  if r3 = 0 then (cond, r0, r1, r2, r3, m)
  else sum(cond, r0, r1, r2, r3, m)
```

side-conditions part of function

Certificate theorem:

```
(sum(c, r0, r1, r2, r3, m) = (true, r'0, r'1, r'2, r'3, m')) ⇒
{ ARM state holds (r0, r1, r2, r3, m) }
E7921003 E0800001 E2533004 1AFFFFF8
{ ARM state holds (r'0, r'1, r'2, r'3, m')} }
```

Example

Assembly code:

```
L0: ldr r1, [r2, r3] ; load mem[r2+r3] into r1
L1: add r0, r1      ; add r1 to r0
L2: subs r3, #4    ; decrement r3 by 4
L3: bne L0         ; goto L0 if r3 ≠ 0
L4:
```

Extracted function:

```
sum(cond, r0, r1, r2, r3, m) =
  let cond = cond ∧ valid_address (r2 + r3) m in
  let r1 = m(r2 + r3) in
  let r0 = r0 + r1 in
  let r3 = r3 - 4 in
    if r3 = 0 then (cond, r0, r1, r2, r3, m)
    else sum(cond, r0, r1, r2, r3, m)
```

side-conditions part of function

side-conditions must be true

Certificate theorem:

```
(sum(c, r0, r1, r2, r3, m) = (true, r'0, r'1, r'2, r'3, m')) ⇒
{ ARM state holds (r0, r1, r2, r3, m) }
E7921003 E0800001 E2533004 1AFFFFFEB
{ ARM state holds (r'0, r'1, r'2, r'3, m') }
```

Code pointers

Assembly code:

```
L0: ldr r4, [r5, r6] ; load mem[r5+r6] into r4
L1: blx r4          ; call code-pointer r4
L2: subs r6, #4     ; decrement r6 by 4
L3: bne L0          ; goto L0 if r6 ≠ 0
L4:
```

Code pointers

Assembly code:

call to code pointer

```
L0: ldr r4, [r5, r6] ; load mem[r5+r6] into r4
L1: blx r4           ; call code-pointer r4
L2: subs r6, #4      ; decrement r6 by 4
L3: bne L0           ; goto L0 if r6 ≠ 0
L4:
```

Code pointers

Assembly code:

call to code pointer

```
L0: ldr r4, [r5, r6] ; load mem[r5+r6] into r4
L1: blx r4          ; call code-pointer r4
L2: subs r6, #4    ; decrement r6 by 4
L3: bne L0         ; goto L0 if r6 ≠ 0
L4:
```

Extracted function:

```
calls(cond, pc, r4, r5, r6, r14, m) =
  if pc = L0 then
    let cond = cond ∧ valid_address (r5 + r6) m in
    let r4 = m(r5 + r6) in
    let cond = cond ∧ word_aligned_address r4 in
    let (pc, r14) = (r4, L2) in
      (cond, pc, r4, r5, r6, r14, m)
  else if pc = L2 then
    let r6 = r6 - 4 in
      if r6 = 0 then (cond, L4, r4, r5, r6, r14, m)
      else (cond, L0, r4, r5, r6, r14, m)
  else (cond, pc, r4, r5, r6, r14, m)
```

Code pointers

Assembly code:

```
L0: ldr r4, [r5, r6] ; load mem[r5+r6] into r4
L1: blx r4          ; call code-pointer r4
L2: subs r6, #4     ; decrement r6 by 4
L3: bne L0         ; goto L0 if r6 ≠ 0
L4:
```

call to code pointer

PC is an input

Extracted function:

```
calls(cond, pc, r4, r5, r6, r14, m) =
  if pc = L0 then
    let cond = cond ∧ valid_address (r5 + r6) m in
    let r4 = m(r5 + r6) in
    let cond = cond ∧ word_aligned_address r4 in
    let (pc, r14) = (r4, L2) in
      (cond, pc, r4, r5, r6, r14, m)
  else if pc = L2 then
    let r6 = r6 - 4 in
      if r6 = 0 then (cond, L4, r4, r5, r6, r14, m)
      else (cond, L0, r4, r5, r6, r14, m)
  else (cond, pc, r4, r5, r6, r14, m)
```

Code pointers

Assembly code:

```
L0: ldr r4, [r5, r6] ; load mem[r5+r6] into r4
L1: blx r4          ; call code-pointer r4
L2: subs r6, #4     ; decrement r6 by 4
L3: bne L0          ; goto L0 if r6 ≠ 0
L4:
```

call to code pointer

PC is an input

Extracted function:

```
calls(cond, pc, r4, r5, r6, r14, m) =
  if pc = L0 then
    let cond = cond ∧ valid_address (r5 + r6) m in
    let r4 = m(r5 + r6) in
    let cond = cond ∧ word_aligned_address r4 in
    let (pc, r14) = (r4, L2) in
      (cond, pc, r4, r5, r6, r14, m)
  else if pc = L2 then
    let r6 = r6 - 4 in
      if r6 = 0 then (cond, L4, r4, r5, r6, r14, m)
      else (cond, L0, r4, r5, r6, r14, m)
  else (cond, pc, r4, r5, r6, r14, m)
```

test for value of PC

Code pointers

Assembly code:

```
L0: ldr r4, [r5, r6] ; load mem[r5+r6] into r4
L1: blx r4          ; call code-pointer r4
L2: subs r6, #4     ; decrement r6 by 4
L3: bne L0          ; goto L0 if r6 ≠ 0
L4:
```

call to code pointer

PC is an input

Extracted function:

```
calls(cond, pc, r4, r5, r6, r14, m) =
  if pc = L0 then
```

test for value of PC

```
    let cond = cond ∧ valid_address (r5 + r6) m in
    let r4 = m(r5 + r6) in
    let cond = cond ∧ word_aligned_address r4 in
    let (pc, r14) = (r4, L2) in
```

PC updated

```
      (cond, pc, r4, r5, r6, r14, m)
    else if pc = L2 then
      let r6 = r6 - 4 in
        if r6 = 0 then (cond, L4, r4, r5, r6, r14, m)
        else (cond, L0, r4, r5, r6, r14, m)
```

```
    else (cond, pc, r4, r5, r6, r14, m)
```

Code pointers

Assembly code:

```
L0: ldr r4, [r5, r6] ; load mem[r5+r6] into r4
L1: blx r4          ; call code-pointer r4
L2: subs r6, #4    ; decrement r6 by 4
L3: bne L0         ; goto L0 if r6 ≠ 0
L4:
```

call to code pointer

PC is an input

Extracted function:

```
calls(cond, pc, r4, r5, r6, r14, m) =
  if pc = L0 then
```

test for value of PC

```
  let cond = cond ∧ valid_address (r5 + r6) m in
  let r4 = m(r5 + r6) in
  let cond = cond ∧ word_aligned_address r4 in
  let (pc, r14) = (r4, L2) in
```

```
(cond, pc, r4, r5, r6, r14, m)
```

Certificate theorem:

```
(calls(c, pc, r4, r5, r6, r14, m) = (true, pc', r4', ...)) ⇒
{ ARM state holds (r4, r5, r6, r14, m) and PC is pc }
E7954006 E12FFF34 E2566004 1AFFFFF8B
{ ARM state holds (r4', r5', r6', r14', m') and PC is pc' }
```

Performance Numbers

Comparison between new and old approach.

Cost given in seconds (s) and HOL inference rules (i).

ARM machine code	instr.	time/cost of old	time/cost of new	reduction	model eval.
sum of array (Sec. I-A)	4	2.5 s (73,039 i)	0.3 s (4,019 i)	86 % (94 %)	7.8 s (1.5 Mi)
copying garbage collector [10]	89	50 s (1,526,281 i)	6.0 s (53,301 i)	88 % (97 %)	173 s (40 Mi)
1024-bit multiword addition	224	70 s (1,029,685 i)	1.2 s (10,802 i)	98 % (99 %)	37 s (8.9 Mi)
256-bit Skein hash function	1,352	5,366 s (21,432,926 i)	56 s (1,842,642 i)	99 % (91 %)	500 s (105 Mi)

Cost of evaluating the ISA model is separate as this cost is independent of decompilation approach.

Performance Numbers

Comparison between new and old approach.

Cost given in seconds (s) and HOL inference rules (i).

ARM machine code	instr.	time/cost of old	time/cost of new	reduction	model eval.
sum of array (Sec. I-A)	4	2.5 s (73,039 i)	0.3 s (4,019 i)	86 % (94 %)	7.8 s (1.5 Mi)
copying garbage collector [10]	89	50 s (1,526,281 i)	6.0 s (53,301 i)	88 % (97 %)	173 s (40 Mi)
1024-bit multiword addition	224	70 s (1,029,685 i)	1.2 s (10,802 i)	98 % (99 %)	37 s (8.9 Mi)
256-bit Skein hash function	1,352	5,366 s (21,432,926 i)	56 s (1,842,642 i)	99 % (91 %)	500 s (105 Mi)

Cost of evaluating the ISA model is separate as this cost is independent of decompilation approach.

Performance Numbers

Comparison between new and old approach.

Cost given in seconds (s) and HOL inference rules (i).

ARM machine code	instr.	time/cost of old	time/cost of new	reduction	model eval.
sum of array (Sec. I-A)	4	2.5 s (73,039 i)	0.3 s (4,019 i)	86 % (94 %)	7.8 s (1.5 Mi)
copying garbage collector [10]	89	50 s (1,526,281 i)	6.0 s (53,301 i)	88 % (97 %)	173 s (40 Mi)
1024-bit multiword addition	224	70 s (1,029,685 i)	1.2 s (10,802 i)	98 % (99 %)	37 s (8.9 Mi)
256-bit Skein hash function	1,352	5,366 s (21,432,926 i)	56 s (1,842,642 i)	99 % (91 %)	500 s (105 Mi)

Cost of evaluating the ISA model is separate as this cost is independent of decompilation approach.

More efficient ISA models...

Better, faster ISA models are in the pipeline...

ITP'12

Directions in ISA Specification

Anthony Fox

Computer Laboratory, University of Cambridge, UK

Abstract. This rough diamond presents a new domain-specific language (DSL) for producing detailed models of Instruction Set Architectures, such as ARM and x86. The language's design and methodology is discussed and we propose future plans for this work. Feedback is sought from the wider theorem proving community in helping establish future directions for this project. A parser and interpreter for the DSL has been developed in Standard ML, with an ARMv7 model used as a case study.

This paper describes recent work on developing a domain-specific language (DSL) for Instruction Set Architecture (ISA) specification. Various theorem

Summary

Decompilation:

- extracts** functions from machine code (ARM, x86, PowerPC)
- proves** that the extracted functions are faithful to the code
- useful in proof of **full functional correctness** (e.g. seL4, Lisp)

Summary

Decompilation:

- ☑ **extracts** functions from machine code (ARM, x86, PowerPC)
- ☑ **proves** that the extracted functions are faithful to the code
- ☑ useful in proof of **full functional correctness** (e.g. seL4, Lisp)

Improvements in this paper:

- ☑ simplified Hoare logic for **easier mechanisation/automation**
- ☑ significantly **improved performance**
- ☑ now more widely applicable (support for **code pointers**)

Summary

Decompilation:

- ☑ **extracts** functions from machine code (ARM, x86, PowerPC)
- ☑ **proves** that the extracted functions are faithful to the code
- ☑ useful in proof of **full functional correctness** (e.g. seL4, Lisp)

Improvements in this paper:

- ☑ simplified Hoare logic for **easier mechanisation/automation**
- ☑ significantly **improved performance**
- ☑ now more widely applicable (support for **code pointers**)

Questions?