Program Transformations as Abstract Interpretation MPRI — Cours 2.6 "Interprétation abstraite : application à la vérification et à l'analyse statique"

Xavier Rival

INRIA

Jan, 8th, 2019

Program transformations and static analysis

Previous lectures: focus on static analysis techniques, i.e.

- take one program as argument
- Compute some semantic properties of the program e.g., compute an over-approximation of the reachable states e.g., verify the absence of runtime errors

Today: we consider program transformations

- functions that compute a program from another program
- thus, we will consider not a single program but two
- different set of issues
 - abstract interpretation to reason about and verify the transformation
 - static analysis to enable the transformation

Compilation

- Transforms programs in high level languages (OCaml, C, Java) into assembly
- Verifies (e.g., types) and Optimizes

Source code:

```
int f( int a, int b ){
    int x0 = a - b;
    if( x0 > 0 )
        return x0 * (a + b);
    else return 0;
}
```

Compiled code:

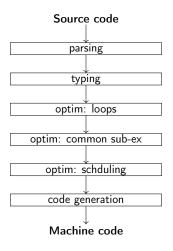
.file "foo.c" .text .globl f .type f, @function f: .LFBO: .cfi_startproc pushl %ebp .cfi_def_cfa_offset 8 .cfi_offset 5, -8 movl %esp, %ebp .cfi_def_cfa_register 5 subl %16, %esp movl 12(%ebp), %eax movl 8(%ebp), %edx movl %edx, %ecx subl %eax, %ecx movl %eax, ~ecx movl %eax, -4(%ebp) cmpl \$0, -4(%ebp) jle.L2 movl 12(%ebp), %eax movl 8(%ebp), %eax addl %edx, %eax imull -4(%ebp), %eax jmp L3 .L2: movl %0, %eax .L3: leave .cfi_restore 5 .cfi_def_cfa 4, 4 ret .cfi_endproc .LFEO: .size f, .-f .ident "GCC: (Gentoo 4.7.3-r1 p1.4, pie-0.5 .section .note.GNU-stack,"",@progbits

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019 3 / 96

Compilation phases



- Parsing: can be considered a static analysis
- Typing: static analysis
- Optimizations: enabled by static analysis
 - e.g., code removed if proved dead
 - e.g., expressions shared if common
- Code generation: by induction on syntax...

Slicing

Slice extraction

- \bullet a slice ${\cal S}$ is a syntactic subset of a program ${\cal P}$
- it is usually extracted following a criterion that describes an observation of the program that is under study
- there are many definitions of slicing criteria: a specific statement, a specific variable, the conjunction of both...

Applications:

• program understanding:

you are given a program, and need to understand how it works...

• program debugging:

a bug was identified, where x stores an unexpected value at line N...

• program maintenance:

a legacy code needs to be extended; what will intended changes do ?

Slicing

Example: slice to understand the value of a at line 5

1:	input(x);		1:	input(x);
2 :	<pre>input(y);</pre>		2 :	input(y);
3 :	a = 4 * x + 8;	\rightarrow	3 :	a = 4 * x + 8;
4:	b = 3 - 2 * y + a;		4 :	b = 3 - 2 * y + a;
5:	c = a + b;		5:	c = a + b;

Algorithm:

- compute dependences: usually, a dependence graph describes what x *immediately* depends on, at line N
- 2 extract a set of slice dependences from the slicing criterion
- S collect the corresponding statements and produce the slice

Effectively, 1 and 2 are a static analysis

Partial evaluation

Specialization and optimization of programs

- start from a very general program
- + possibly some assumptions on the input values
- compute a program that behaves similarly on those programs that satisfy the inputs
- partial evaluation of all program statements that can be, but may also involve unrolling of loop, duplication of functions...

Applications:

• practical:

design a software for several products, and specialize it for each product

theoretical: Futamura's projections
 compilation = specialization of an interpreter to a program

Partial evaluation

unfolding of the loop for a number of iterations
propagation of the value of b through the loop
simplification of conditions and removal of b

Questions about program transformations

Soundness:

- in what sense can we say a transformation is sound ?
- what properties should it preserves ? what properties should it modify ?
- how to semantically specify a transformation ?

Use of semantic information:

• transformations often need semantic properties of programs, to decide what code to generate...

e.g., for compiler optimizations, dependence information...

• in some cases the transformation itself may be potentially non terminating, and require a widening for convergence *e.g.*, partial evaluation

Example: semantics of C volatile variables

From the ANSI C'99 / C'11 standards

For every read from or write to a volatile variable that would be performed by a straightforward interpreter for C, exactly one load or store from/to the memory location allocated to the variable should be performed.

In other words:

- volatile variables should be assumed to be **modifiable** by the external world **at any time** (this is a worst case assumption)
- multiple accesses to a single volatile variable should never be optimized into a single read (this is a very strong constraint on the optimizers)

Do compilers follow this semantics ? NO...

Example: C compiler and volatile variables

Study by E. Eide and J. Regher, "Volatiles are Mis-compiled, and What to Do about it" (EMSOFT'2008)

- 13 compilers tested
- none of them is exempt of volatile bugs
- possible consequences:
 - incorrect computations
 - more serious crashes, such as system hangs
- one example on the next slide, more in the paper...

Since then, the **CompCert compiler** was tested free of volatile bugs using the same technique...

Example: C compiler and volatile variables

```
Compiler: LLVM GCC 2.2 (IA 32)
```

Buggy optimization:		
<pre>volatile int a; void foo(void){ int i; for(i = 0; i < 3; i + +){ a+ = 7; } }</pre>	leal movl leal movl addl	a,%eax 7(%eax),%ecx %ecx,a 14(%eax),%ecx %ecx,a \$21,%eax %eax,a

Only ONE load to a

- loop unrolled three times
- three stores (correct), but only one load (incorrect)

Main points of the lecture

Formalize soundness of program transformations:

- compare the semantics of two programs
- select the semantics to be compared by abstraction

Consider some verification techniques:

- invariant verification approach
- local equivalence proof...

These are partly inspired from static analysis techniques

Outline

Introduction to program transformations

2 Compilation correctness

- 3 Correctness of optimizing compilation
- 4 Application to the verification of compiled code
- 5 Application to certified compilation
- 6 Conclusion

Formalizing correctness: assumptions

Source language: C like imperative language

• very simplified: no procedure, library functions, etc

Assembly language: RISC style (similar to Power-PC)

- registers: differentiated dep. on types (floating-point, integers)
- memory access: direct, indirect, stack-based
- condition register:

Tests and branchings are **separate** operations Conditional branching: tests the value of the condition register

Compiler:

- the lecture is not about showing a compiler...
- we first assume no optimization and consider optimizations later

Transition systems

We assume a (source or compiled) program is a transition system $\mathcal{S} = (\mathbb{S}, \rightarrow, \mathbb{S}_{\mathcal{I}})$:

- $\mathbb{S}=\mathbb{L}\times\mathbb{M}$ is the set of states, where $\mathbb{M}=\mathbb{X}\to\mathbb{V}$
- $\bullet \ {\rightarrow} \subseteq \mathbb{S} \times \mathbb{S}$ is the transition relation
- $\mathbb{S}_{\mathcal{I}} \subseteq \mathbb{S}$ is the set of initial states

We consider their finite traces semantics:

- $\llbracket S \rrbracket = \{ \langle s_0, \dots, s_n \rangle \in \mathbb{S}^* \mid s_0 \in \mathbb{S}_{\mathcal{I}} \land \forall i, s_i \to s_{i+1} \}$
- it can be defined as a least fix-point: [S] = Ifp F

$$\begin{array}{rccc} F : & \mathcal{P}(\mathbb{S}^*) & \longrightarrow & \mathcal{P}(\mathbb{S}^*) \\ & X & \longmapsto & \{ \langle s_0 \rangle \mid s \in \mathbb{S}_{\mathcal{I}} \} \\ & & \cup \{ \langle s_0, \dots, s_n, s_{n+1} \rangle \\ & & & \mid \langle s_0, \dots, s_n \rangle \in X \land s_n \to s_{n+1} \} \end{array}$$

(exercise)

A very minimal imperative language

1 :=**I**-values | x (x $\in X$) e ::= expressions $\begin{array}{ccc} | & c & (c \in \mathbb{V}) \\ | & 1 & (l-value) \\ | & e \oplus e & (arith operation, comparison) \end{array}$ s ::= statements if(e){s} (condition) while(e){s} (loop)

Other extensions, not considered at this stage:

- functions
- collection of arithmetic data types, structures, unions, pointers
- compilation units...

A basic, PPC-like assembly language: principles

We now consider a (very simplified) assembly language

- machine integers: sequences of 32-bits (set: \mathbb{B}^{32})
- instructions are encoded over 32-bits (set: $\mathbb{I}_{\mathrm{MIPS}}$) and stored into the same space as data (*i.e.*, $\mathbb{I}_{\mathrm{MIPS}} \subseteq \mathbb{B}^{32}$)
- loads and store instructions, with relative addressing instructions
- conditional branching is indirect: comparison instruction sets condition register cr (comparison flag) conditional branching instruction reads cr and branches accordingly

Memory locations

- program counter pc (current instruction address)
- general purpose registers r_0, \ldots, r_{31}
- \bullet main memory (RAM) $\textbf{Addrs} \to \mathbb{B}^{32}$ where $\textbf{Addrs} \subseteq \mathbb{B}^{32}$
- condition register cr

Then: $\mathbb{X}^{c} = \{\mathbf{pc}, \mathbf{cr}, \mathbf{r}_{0}, \dots, \mathbf{r}_{31}\} \uplus \mathbf{Addrs}$

Compilation correctness

A basic, PPC-like assembly language: instruction set

Instruction encoded into 32-bits words:

Instruction set

$v, \mathit{dst}, o \in \mathbb{B}^{32}, cr \in \{ \mathrm{L}^{2} \}$	Γ, EQ, GT
i ::= ($\in \mathbb{I}_{\mathrm{MIPS}}$)	
li r _d , v	load $v \in \mathbb{B}^{32}$
$ $ add $\mathbf{r}_d, \mathbf{r}_{s0}, \mathbf{r}_{s1}$	addition
addi $\mathbf{r}_d, \mathbf{r}_{s0}, \mathbf{v}$	add. $ u \in \mathbb{V}' \subset \mathbb{B}^{32}$
$ $ sub $\mathbf{r}_d, \mathbf{r}_{s0}, \mathbf{r}_{s1}$	subtraction
$ $ cmp r_{s0}, r_{s1}	comparison
b dst	branch
$ $ blt $\langle cr \rangle$ dst	cond. branch
$ $ Id \mathbf{r}_d, o	absolute load
$ $ st \mathbf{r}_d, o	absolute store
$ $ Idx $\mathbf{r}_d, o, \mathbf{r}_x$	relative load (aka indeXed load)
$ $ stx $\mathbf{r}_d, o, \mathbf{r}_x$	relative store (aka indeXed store)

A basic, PPC-like assembly language: states

Definition: state

- A state is a tuple $s = (pc, \rho, cr, \mu)$ which comprises:
 - a program counter value $pc \in \mathbb{B}^{32}$
 - a function mapping each general purpose register to its value $\rho: \{0, \dots, 31\} \to \mathbb{B}^{32}$
 - a condition register value $cr \in \{LT, EQ, GT\}$
 - a function mapping each memory cell to its value μ : Addrs $\rightarrow \mathbb{B}^{32}$

Equivalently, we can also write s = (l, m), where

- the control state l is the current pc value
- the memory state *m* is the triple (ρ, cr, μ)

A basic, PPC-like assembly language: instruction set

We assume a state $s = (pc, (\rho, cr, \mu))$ and that $\mu(pc) = i$. Then:

• if $i = \mathbf{li} \mathbf{r}_d, v$, then:

$$s \rightarrow (pc + 4, (\rho[d \mapsto v], cr, \mu))$$

• if $i = \text{add } \mathbf{r}_d, \mathbf{r}_{s0}, \mathbf{r}_{s1}$, then:

$$\mathsf{s} o (\mathsf{pc}+\mathsf{4},(
ho[\mathsf{d}\mapsto (
ho(\mathsf{s0})+
ho(\mathsf{s1}))],\mathsf{cr},\mu))$$

• if $i = addi r_d, r_{s0}, v$, then:

$$s
ightarrow (
hoc+4, (
ho[d \mapsto (
ho(s0)+v)], cr, \mu))$$

• if $i = \operatorname{sub} \mathbf{r}_d, \mathbf{r}_{s0}, \mathbf{r}_{s1}$, then:

$$s \rightarrow (pc + 4, (\rho[d \mapsto (\rho(s0) - \rho(s1))], cr, \mu))$$

A basic, PPC-like assembly language: instruction set

We assume a state $s = (pc, (\rho, cr, \mu))$ and that $\mu(pc) = i$. Then:

• if $i = \operatorname{cmp} \mathbf{r}_{s0}, \mathbf{r}_{s1}$, then:

$$s \rightarrow \begin{cases} (\rho c + 4, (\rho, \mathrm{LT}, \mu)) & \text{if } \rho(\mathfrak{s0}) < \rho(\mathfrak{s1}) \\ (\rho c + 4, (\rho, \mathrm{EQ}, \mu)) & \text{if } \rho(\mathfrak{s0}) = \rho(\mathfrak{s1}) \\ (\rho c + 4, (\rho, \mathrm{GT}, \mu)) & \text{if } \rho(\mathfrak{s0}) > \rho(\mathfrak{s1}) \end{cases}$$

• if $i = \mathbf{blt} \langle cond \rangle dst$, then:

$$s
ightarrow \left\{ egin{array}{cc} (dst,(
ho,{f cr},\mu)) & {
m if} \ cr=cond\ (
hockstruck+4,(
ho,{f cr},\mu)) & {
m otherwise} \end{array}
ight.$$

• if $i = \mathbf{b} \, dst$, then:

 $s
ightarrow (\textit{dst}, (
ho, \textit{cr}, \mu))$

Xavier Rival (INRIA)

Program Transformations

A basic, PPC-like assembly language: instruction set

We assume a state $s = (pc, (\rho, cr, \mu))$ and that $\mu(pc) = i$. Then:

• if $i = \mathbf{Idx} \mathbf{r}_d, o, \mathbf{r}_x$, then: $s \rightarrow \begin{cases} (pc+4, (\rho[d \mapsto \mu(\rho(x)+o)], cr, \mu)) & \text{if } \mu(\rho(x)+o) \text{ is defined} \\ \Omega & \text{otherwise} \end{cases}$ • if $i = \mathbf{Id} \mathbf{r}_d, o$, then: $s \rightarrow \begin{cases} (pc+4, (\rho[d \mapsto \mu(o)], cr, \mu)) & \text{if } \mu(o) \text{ is defined} \\ \Omega & \text{otherwise} \end{cases}$ • if $i = \operatorname{stx} \mathbf{r}_d, o, \mathbf{r}_x$, then: $s \to \begin{cases} (pc+4, (\rho, cr, \mu[\rho(x)+o) \mapsto \rho(d)])) & \text{if } \mu(\rho(x)+o) \text{ is defined} \\ \Omega & \text{otherwise} \end{cases}$

• if $i = \text{Id } \mathbf{r}_d, o$, then effect can be deduced from the above cases

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019 23 / 96

Output of a non optimizing compiler

Assumptions and conventions:

- t is an array of integers initialized to $\mathtt{t}=\{\mathtt{0};\mathtt{1};\mathtt{4};-\mathtt{1}\}$
- \bullet i,x are integer variables
- $\bullet\,$ in the assembly, \underline{x} denotes the address of x

source code	compiled code
l_0^s i := i + 1;	<i>l</i> ₀ ^c ld r ₀ , <u>i</u>
	l_1^c addi $r_0, r_0, 1$
	l_2^c st $\mathbf{r}_0, \underline{\mathbf{i}}$
l_1^s $x := x + t[i];$	l_3^c ld $\mathbf{r}_0, \underline{\mathbf{x}}$
	l ₄ ^c ld r ₁ , <u>i</u>
	l_5^c ldx $\mathbf{r}_2, \underline{\mathbf{t}}, \mathbf{r}_1$
	l_6^c add $\mathbf{r}_0, \mathbf{r}_0, \mathbf{r}_2$
	l_7^c st $\mathbf{r}_0, \underline{\mathbf{x}}$
$l_2^s \ldots$	$l_8^c \dots$

Is it sound ? What property does it preserve ?

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019 24 / 96

A source level execution

$$\left\langle \begin{pmatrix} \mathrm{i}\mapsto 1; \\ \mathrm{x}\mapsto 1; \\ \boldsymbol{\ell}_0^s, \ t[0]\mapsto 0; \\ t[1]\mapsto 1; \\ t[2]\mapsto 4; \\ t[3]\mapsto -1; \end{pmatrix}, \begin{pmatrix} \mathrm{i}\mapsto 2; \\ \mathrm{x}\mapsto 1; \\ \mathbf{t}[0]\mapsto 0; \\ \boldsymbol{\ell}_1^s, \ t[0]\mapsto 0; \\ t[1]\mapsto 1; \\ t[2]\mapsto 4; \\ t[3]\mapsto -1; \end{pmatrix}, \begin{pmatrix} \mathrm{i}\mapsto 2; \\ \mathrm{x}\mapsto 5; \\ \mathbf{t}[0]\mapsto 0; \\ \mathbf{t}[1]\mapsto 0; \\ t[2]\mapsto 4; \\ t[3]\mapsto -1; \end{pmatrix}, \right\rangle$$

Correctness of compilation:

- we cannot find the same execution in the assembly: as memory locations are not the same at all
- thus, we expect a "similar" trace

Compilation correctness

Corresponding assembly level execution

6 ^c	ld r ₀ , <u>i</u>	l_4^c	ld r ₁ , <u>i</u>
l_1^c	$\textbf{addi } \textbf{r}_0, \textbf{r}_0, 1$	l_5^c	$Idx\; r_2, \underline{\mathtt{t}}, r_1$
l_2^c	st r ₀ , <u>i</u>	l_6^c	$\text{add } \textbf{r}_0, \textbf{r}_0, \textbf{r}_2$
l3 ^c	$Id \; r_0, \underline{x}$	l_7^c	$\text{st } r_0, \underline{x}$

We consider an assembly level trace starting from a similar state:

state s ^c _i	s ₀ ^c	s_1^c	s ₂ ^c	s ₃ ^c	s ₄ ^c	s ₅ ^c	s ₆ ^c	s ₇ ^c	<i>s</i> ₈ ^c
control state <i>pc</i> _i	l_0^c	l_1^c	l2 ^c	l3 ^c	l ₄ ^c	l ₅ ^c	l ₆ c	l7 ^c	l ₈ ^c
register state $\rho_i(0)$	45	1	2	2	1	1	1	5	5
register state $\rho_i(1)$	-5	-5	-5	-5	-5	2	2	2	2
register state $\rho_i(2)$	89	89	89	89	89	89	4	4	4
memory state $\mu_i(\underline{i})$	1	1	1	2	2	2	2	2	2
memory state $\mu_i(\underline{x})$	1	1	1	1	1	1	1	1	5
memory state $\mu_i(\underline{t} + 0)$	0	0	0	0	0	0	0	0	0
memory state $\mu_i(\underline{t}+1)$	1	1	1	1	1	1	1	1	1
memory state $\mu_i(\underline{t}+2)$	4	4	4	4	4	4	4	4	4
memory state $\mu_i(\underline{t}+3)$	-1	-1	-1	-1	-1	-1	-1	-1	-1

Compilation correctness

Source and assembly executions compared

state s;	<i>s</i> ^{<i>s</i>} ₀	s_1^s	<i>s</i> ^{<i>s</i>} ₂
control state l_i^s	l_0^s	l_1^s	l_2^s
memory state $m_i^s(i)$	1	2	2
memory state $m_i^s(x)$	1	1	5
memory state $m_i^s(t[0])$	0	0	0
memory state $m_i^s(t[1])$	1	1	1
memory state $m_i^s(t[2])$	4	4	4
memory state $m_i^s(t[3])$	-1	-1	-1

Much more information in the assembly trace:

- registers values
- more control states

state s _i ^c	s ₀ ^c	s_1^c	s ₂ ^c	s ₃ ^c	s ₄ ^c	s ₅ ^c	s ₆ ^c	s ₇ ^c	<i>s</i> ₈ ^c
control state <i>pc_i</i>	l_0^c	l_1^c	l_2^c	l3 ^c	l_4^c	l_5^c	l6 ^c	l7 ^c	l ₈ ^c
register state $\rho_i(0)$	45	1	2	2	1	1	1	5	5
register state $\rho_i(1)$	$^{-5}$	$^{-5}$	$^{-5}$	$^{-5}$	-5	2	2	2	2
register state $\rho_i(2)$	89	89	89	89	89	89	4	4	4
memory state $\mu_i(\underline{i})$	1	1	1	2	2	2	2	2	2
memory state $\mu_i(\underline{x})$	1	1	1	1	1	1	1	1	5
memory state $\mu_i(\underline{t} + 0)$	0	0	0	0	0	0	0	0	0
memory state $\mu_i(\underline{t}+1)$	1	1	1	1	1	1	1	1	1
memory state $\mu_i(\underline{t}+2)$	4	4	4	4	4	4	4	4	4
memory state $\mu_i(\underline{t} + 3)$	-1	-1	-1	-1	-1	-1	-1	-1	-1

An abstraction approach

state s ^s _i	s_0^s			s_1^s					<i>s</i> ^s ₂
control state l_i^s	l_0^s			l_1^s					l_2^s
memory state $m_i^s(i)$	1			2					2
memory state $m_i^s(\mathbf{x})$	1			1					5
memory state $m_i^s(t[0])$	0			0					0
memory state $m_i^s(t[1])$	1			1					1
memory state $m_i^s(t[2])$	4			4					4
memory state $m_i^s(t[3])$	-1			-1					$^{-1}$
state s _i ^c	s ₀ ^c	s_1^c	s ₂ ^c	s ₃ ^c	s ₄ ^c	s ₅ ^c	s ₆ ^c	s ₇ ^c	<i>s</i> ₈ ^c
control state pci	l ₀ ^c	l1 ^c	b	l3 ^c	l ₄ ^c	l5 ^c	l ₆ ^c	l7 ^C	l8 ^c
	0	1	- 2	5	4	5	0	-7	0
register state $\rho_i(0)$	45	1	2	2	1	1	1	5	5
register state $\rho_i(0)$ register state $\rho_i(1)$		_							
	45	1	2	2	1	1	1	5	5
register state $\rho_i(1)$	45 5	1 -5	2 -5	2 -5	1 -5	1 2	1 2	5	5 2
register state $\rho_i(1)$ register state $\rho_i(2)$	45 -5 89	1 -5 89	2 -5 89	2 -5 89	1 -5 89	1 2 89	1 2 4	5 2 4	5 2 4
register state $\rho_i(1)$ register state $\rho_i(2)$ memory state $\mu_i(\underline{i})$	45 -5 89 1	1 -5 89 1	2 -5 89 1	2 -5 89 2	1 -5 89 2	1 2 89 2	1 2 4 2	5 2 4 2	5 2 4 2
register state $\rho_i(1)$ register state $\rho_i(2)$ memory state $\mu_i(\underline{i})$ memory state $\mu_i(\underline{x})$	45 -5 89 1 1	1 -5 89 1 1	2 -5 89 1 1	2 -5 89 2 1	1 -5 89 2 1	1 2 89 2 1	1 2 4 2 1	5 2 4 2 1	5 2 4 2 5
register state $\rho_i(1)$ register state $\rho_i(2)$ memory state $\mu_i(\underline{i})$ memory state $\mu_i(\underline{x})$ memory state $\mu_i(\underline{x} + 0)$	45 -5 89 1 1 0	1 -5 89 1 1 0	2 -5 89 1 1 0	2 5 89 2 1 0	1 -5 89 2 1 0	1 2 89 2 1 0	1 2 4 2 1 0	5 2 4 2 1 0	5 2 4 2 5 0

• We can abstract away intermediate control states

An abstraction approach

state s ^s	<i>s</i> ^s 0	1		s_1^s					<i>s</i> ^s ₂
control state l_i^s	l ₀ s			l_1^s					l_2^s
memory state $m_i^s(i)$	1	1		2					2
memory state $m_i^s(\mathbf{x})$	1			1					5
memory state $m_i^s(t[0])$	0			0					0
memory state $m_i^s(t[1])$	1			1					1
memory state $m_i^s(t[2])$	4			4					4
memory state $m_i^s(t[3])$	-1			-1					$^{-1}$
state s _i ^c	s ₀ ^c	s_1^c	s ₂ ^c	s ₃ ^c	S4C	S_5^c	s ₆ c	s ₇ ^c	<i>s</i> ₈ ^c
control state pci	l ₀ ^c	l_1^c	l_2^c	l3 ^c	l_4^c	l ₅ ^c	l ₆ ^c	l_7^c	l8 ^c
register state $\rho_i(0)$	45	1	2	2	1	1	1	5	5
register state $\rho_i(1)$	-5	-5	-5	-5	-5	2	2	2	2
register state $\rho_i(2)$	89	89	89	89	89	89	4	4	4
(.)	1		- 11	0	0	0	2	2	2
memory state $\mu_i(\underline{i})$	1	1	1	2	2			~	~
memory state $\mu_i(\underline{i})$ memory state $\mu_i(\underline{x})$	1	1	1	2	2	1	1	1	5
,	-	1 1 0	1 1 0	_		1 0	2 1 0		_
memory state $\mu_i(\underline{x})$	1	-	1 1 0 1	1	1	2 1 0 1	1	1	5
memory state $\mu_i(\underline{x})$ memory state $\mu_i(\underline{t} + 0)$	1 0	-	1 1 0 1 4	1 0	1	2 1 0 1 4	1	1	5

• Intermediate control states abstracted; we can forget registers

An abstraction approach

state s ^s	<i>s</i> ₀ ^s			s_1^s					s ₂ s
control state l_i^s	l ₀ s			l_1^s					l_2^s
memory state $m_i^s(i)$	1			2					2
memory state $m_i^s(\mathbf{x})$	1			1					5
memory state $m_i^s(t[0])$	0			0					0
memory state $m_i^s(t[1])$	1			1					1
memory state $m_i^s(t[2])$	4			4					4
memory state $m_i^s(t[3])$	-1			-1					-1
state s _i ^c	s ₀ ^c	s_1^c	s ₂ ^c	s ₃ ^c	S4C	S_5^C	S6C	s ₇ ^c	s ₈ ^c
control state pci	l ₀ ^c	l_1^c	l_2^c	l3 ^c	l_4^c	l_5^c	l ₆ ^c	l_7^c	l8 ^c
register state $\rho_i(0)$	45	1	2	2	1	1	1	5	5
register state $\rho_i(1)$	-					~	~		0
· · · · · · · · · · · · · · · · · · ·	-5	-5	-5	-5	-5		2	2	2
register state $\rho_i(2)$	-5 89	-5 89	-5 89	—5 89	-5 89	2 89	2	2	2
• • • • • •		~	-5 89 1			2 89 2	2 4 2		
register state $\rho_i(2)$	89	~	-5 89 1 1	89	89	2 89 2 1	2 4 2 1	4	4
register state $\rho_i(2)$ memory state $\mu_i(\underline{i})$	89 1	~	-5 89 1 1 0	89 2	89 2	2 89 2 1 0	2	4	4
register state $\rho_i(2)$ memory state $\mu_i(\underline{i})$ memory state $\mu_i(\underline{x})$ memory state $\mu_i(\underline{t} + 0)$ memory state $\mu_i(\underline{t} + 1)$	89 1 1	89 1 1	1	89 2 1	89 2 1	2 89 2 1 0 1	2	4 2 1	4 2 5
register state $\mu_i(\underline{i})$ memory state $\mu_i(\underline{i})$ memory state $\mu_i(\underline{x})$ memory state $\mu_i(\underline{t} + 0)$	89 1 1 0	89 1 1	1	89 2 1 0	89 2 1	2 89 2 1 0 1 4	2	4 2 1	4 2 5 0

• Registers and intermediate control states removed We get very similar traces !

Syntactic relations

What we did remove:

- intermediate control states
- memory locations associated to registers

What we did preserve:

• control states in correspondence:

$$l_0^s \leftrightarrow l_0^c \qquad l_1^s \leftrightarrow l_3^c \qquad l_2^s \leftrightarrow l_8^c$$

• memory location in correspondence:

$$\begin{array}{cccc} \mathbf{i} \leftrightarrow \underline{\mathbf{i}} & \mathbf{x} \leftrightarrow \underline{\mathbf{x}} & \mathbf{i} \leftrightarrow \underline{\mathbf{i}} \\ \mathbf{t}[\mathbf{0}] \leftrightarrow \underline{\mathbf{t}} + \mathbf{0} & \mathbf{t}[\mathbf{1}] \leftrightarrow \underline{\mathbf{t}} + \mathbf{1} & \mathbf{t}[\mathbf{2}] \leftrightarrow \underline{\mathbf{t}} + \mathbf{2} \\ \mathbf{t}[\mathbf{3}] \leftrightarrow \underline{\mathbf{t}} + \mathbf{3} \end{array}$$

Intuitively, we did apply an abstraction (to a single trace)

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019 29 / 96

Syntactic relations

Definition

We define two syntactic mappings:

- Between control points: $\pi_{\mathbf{I}} : \mathbb{L}'_s \to \mathbb{L}'_c$ (where $\mathbb{L}'_i \subseteq \mathbb{L}_i$)
- Between memory locations: $\pi_{\mathbf{x}} : \mathbb{X}'_{\mathbf{s}} \to \mathbb{X}'_{\mathbf{c}}$ (where $\mathbb{X}'_{i} \subseteq \mathbb{X}_{i}$)

We consider only subsets X', \ldots of X, \ldots For instance:

- Some variables in the source code may be removed
- Registers in P_c may not correspond to variables of P_s
- One statement in P_s corresponds to several instructions in P_c

In practice, π_{I}, π_{x} are provided by the compiler:

- Linking information
- Line table

• Debugging information: Stabs, COFF...

Syntactic relations

Definition

We define two syntactic mappings:

- Between control points: $\pi_{\mathbf{I}} : \mathbb{L}'_s \to \mathbb{L}'_c$ (where $\mathbb{L}'_i \subseteq \mathbb{L}_i$)
- Between memory locations: $\pi_{\mathbf{x}} : \mathbb{X}'_{\mathbf{s}} \to \mathbb{X}'_{\mathbf{c}}$ (where $\mathbb{X}'_{i} \subseteq \mathbb{X}_{i}$)

For our example:

- Control points:
 - $\mathbb{L}'_{s} = \{l_{0}^{s}, l_{1}^{s}, l_{2}^{s}\}$ and $\mathbb{L}'_{c} = \{l_{0}^{c}, l_{3}^{c}, l_{8}^{c}\}$
 - $\blacktriangleright \ \pi_{\mathbf{I}}: l_0^s \mapsto l_0^c; \ l_1^s \mapsto l_3^c; \ l_2^s \mapsto l_8^c$
- Memory locations:

$$\begin{array}{l} \bullet \quad \mathbb{X}'_{s} = \{\mathbf{i}, \mathbf{x}, \mathbf{t}[0], \mathbf{t}[1], \mathbf{t}[2], \mathbf{t}[3]\} \text{ and } \mathbb{X}'_{c} = \{\underline{\mathbf{i}}, \underline{\mathbf{x}}, \underline{\mathbf{t}}, \underline{\mathbf{t}} + 1, \underline{\mathbf{t}} + 2, \underline{\mathbf{t}} + 3\} \\ \bullet \quad \pi_{\mathbf{x}} : \begin{cases} \mathbf{i} \quad \mapsto \quad \underline{i} \\ \mathbf{x} \quad \mapsto \quad \underline{x} \\ \mathbf{t}[n] \quad \mapsto \quad \underline{t} + n \end{cases}$$

State observational abstraction

We now formalize the process to project out irrelevant behaviors:

- in states
- in traces
- in the semantics

We consider the assembly level first:

Definition: state abstraction

We let the compiled code-level memory state abstraction $\Psi_c^{\mathbf{m}}$ be defined by:

$$\begin{array}{rcl} \Psi^{\mathbf{m}}_c: & (\mathbb{X}_c \to \mathbb{V}) & \longrightarrow & (\mathbb{X}'_c \to \mathbb{V}) \\ & m & \longmapsto & \lambda(x \in \mathbb{X}'_c) \cdot m(x) \end{array}$$

Similar definition at the source level...

(though no variable needs to be abstracted at this point, we will make use of that possibility further in this course)

Xavier Rival (INRIA)

Program Transformations

Compilation correctness

State observational abstraction: example

We recall that

$$\begin{aligned} \mathbb{X}'_s &= \{\mathtt{i}, \mathtt{x}, \mathtt{t}[0], \mathtt{t}[1], \mathtt{t}[2], \mathtt{t}[3] \} \\ \mathbb{X}'_c &= \{\underline{\mathtt{i}}, \underline{\mathtt{x}}, \underline{\mathtt{t}}, \underline{\mathtt{t}} + 1, \underline{\mathtt{t}} + 2, \underline{\mathtt{t}} + 3 \end{aligned}$$

Then $\Psi_c^{\mathbf{m}} : (pc, (\rho, \mathbf{cr}, \mu)) \longmapsto \mu$

So, in particular:

$$\Psi_{c}^{\mathbf{m}}:\left(\begin{array}{ccccc} pc & \mapsto & t_{0}^{c} \\ \rho: & 0 & \mapsto & 45 \\ 1 & \mapsto & -5 \\ 2 & \mapsto & 4 \\ \mu: & \underline{i} & \mapsto & 1 \\ & \underline{x} & \mapsto & 1 \\ & \underline{t}+0 & \mapsto & 0 \\ & \underline{t}+1 & \mapsto & 1 \\ & \underline{t}+2 & \mapsto & 4 \\ & \underline{t}+3 & \mapsto & -1 \end{array}\right) \mapsto \left(\begin{array}{cccc} \mu: & \underline{i} & \mapsto & 1 \\ & \underline{x} & \mapsto & 1 \\ & \underline{t}+0 & \mapsto & 0 \\ & \underline{t}+1 & \mapsto & 1 \\ & \underline{t}+2 & \mapsto & 4 \\ & \underline{t}+3 & \mapsto & -1 \end{array}\right)$$

Trace observational abstraction

We can now lift the same abstraction principle to traces:

Definition: trace abstraction

We let the compiled code-level trace abstraction Ψ_c^{tr} be defined by:

$$\begin{split} \Psi_{c}^{\mathsf{tr}} : & (\mathbb{L}_{c} \times (\mathbb{X}_{c} \to \mathbb{V}))^{\star} & \longrightarrow & (\mathbb{L}_{c}' \times (\mathbb{X}_{c}' \to \mathbb{V}))^{\star} \\ & \langle (l_{0}, m_{0}), \dots, (l_{n}, m_{n}) \rangle & \longmapsto & \langle (l_{k_{0}}, \Psi_{c}^{\mathsf{m}}(m_{k_{0}})), \dots, (l_{k_{m}}, \Psi_{c}^{\mathsf{m}}(m_{k_{m}})) \rangle \\ & \text{where:} \begin{cases} \{k_{0}, \dots, k_{m}\} = \{k \mid 0 \leq k \leq n \land l_{k} \in \mathbb{L}_{c}'\} \\ k_{0} < \dots < k_{m} \end{cases} \end{split}$$

Similar definition at the source level...

(though no control state / variable needs to be abstracted at this point, we will make use of that possibility further in this course)

Compilation correctness

Trace observational abstraction: example

control state pci	l ₀ ^c	l_1^c	l2 ^c	l3 ^c	l4 ^c	l_5^c	l ₆ ^c	l ₇ ^c	l ₈ c
register state $\rho_i(0)$	45	1	2	2	1	1	1	5	5
register state $\rho_i(1)$	-5	-5	-5	-5	-5	2	2	2	2
register state $\rho_i(2)$	89	89	89	89	89	89	4	4	4
memory state $\mu_i(\underline{i})$	1	1	1	2	2	2	2	2	2
memory state $\mu_i(\underline{x})$	1	1	1	1	1	1	1	1	5
memory state $\mu_i(\underline{t} + 0)$	0	0	0	0	0	0	0	0	0
memory state $\mu_i(\underline{t} + 1)$	1	1	1	1	1	1	1	1	1
memory state $\mu_i(\underline{t} + 2)$	4	4	4	4	4	4	4	4	4
memory state $\mu_i(\underline{t} + 3)$	-1	-1	-1	-1	-1	-1	-1	-1	-1

	control state <i>pc</i> _i	l ₀ ^c	l3 ^c	l ₈ c
	memory state $\mu_i(\underline{i})$	1	2	2
	memory state $\mu_i(\underline{x})$	1	1	5
\mapsto	memory state $\mu_i(\underline{t} + 0)$	0	0	0
	memory state $\mu_i(\underline{t} + 1)$	1	1	1
	memory state $\mu_i(\underline{t} + 2)$	4	4	4
	memory state $\mu_i(\underline{t} + 3)$	-1	-1	-1

 Ψ^{tr} :

Observable behaviors inclusions

Applying this systematically to all traces results into an abstraction:

Result: compiled code observational abstraction We let α_c^r be the compiled code observational abstraction:

$$\begin{array}{ccc} \alpha_{c}^{r} : & \mathcal{P}((\mathbb{L}_{c} \times (\mathbb{X}_{c} \to \mathbb{V}))^{\star}) & \longrightarrow & \mathcal{P}((\mathbb{L}_{c}^{r} \times (\mathbb{X}_{c}^{r} \to \mathbb{V}))^{\star}) \\ & \mathcal{E} & \longmapsto & \{\Psi_{c}^{tr}(\sigma) \mid \sigma \in \mathcal{E}\} \end{array}$$

It defines a Galois connection with an adjoint concretization γ_c^r :

$$(\mathcal{P}((\mathbb{L}_c imes (\mathbb{X}_c o \mathbb{V}))^{\star}), \subseteq) \xleftarrow{\gamma_c^r}{\alpha_c^r} (\mathcal{P}((\mathbb{L}_c' imes (\mathbb{X}_c' o \mathbb{V}))^{\star}), \subseteq)$$

 α^c_c is monotone and the concrete domain is a complete lattice; the concretization function follows and is defined by γ^c_c(ε') = ∪_ε{ε | α^c_c(ε) ⊆ ε'} = {σ | Ψ^{tr}(σ) ∈ ε'}
 The observational semantics is defined by: [[P_c]]_{obs} = α^c_c([[P_c]])

Xavier Rival (INRIA)

Correctness by semantic equivalence

- The same construction holds at the source level
- The resulting traces are very similar, up-to a basic renaming
- To define it, we assume the syntactic mappings π_{I}, π_{x} are bijective

Memory state renaming

We let the memory state renaming function be defined by:

$$\begin{array}{ccc} \pi_{\mathbf{m}} : & (\mathbb{X}'_{s} \to \mathbb{V}) & \longrightarrow & (\mathbb{X}'_{c} \to \mathbb{V}) \\ & m & \longmapsto & m \circ \pi_{\mathbf{x}}^{-1} \end{array}$$

Trace renaming

We let the trace renaming function be defined by:

$$\begin{aligned} \pi_{\mathbf{t}} : & \mathbb{L}'_{\mathbf{s}} \times (\mathbb{X}'_{\mathbf{s}} \to \mathbb{V}) & \longrightarrow & \mathbb{L}'_{\mathbf{c}} \times (\mathbb{X}'_{\mathbf{c}} \to \mathbb{V}) \\ & \langle (\ell_0, m_0), \dots, (\ell_n, m_n) \rangle & \longmapsto & \langle (\pi_{\mathbf{l}}(\ell_0), \pi_{\mathbf{m}}(m_0)), \dots, (\pi_{\mathbf{l}}(\ell_n), \pi_{\mathbf{m}}(m_n)) \rangle \end{aligned}$$

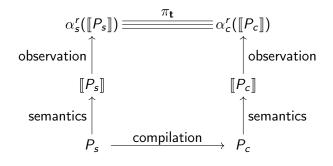
Correctness by semantic equivalence

We can now state the compilation correctness definition

Definition: compilation correctness

Compilation of P_s into P_c is correct with respect to π_I, π_x if and only if π_t establishes a bijection between $\alpha_s^r(\llbracket P_s \rrbracket)$ and $\alpha_c^r(\llbracket P_c \rrbracket)$.

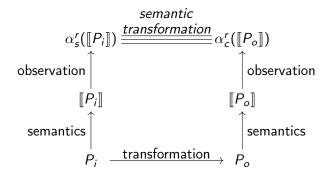
This definition can be illustrated by the diagram:



Correctness by semantic equivalence

This approach generalizes to other program transformations

This definition can be illustrated by the diagram:



Choice of another concrete semantics: consequences

New compilation correctness definition $\forall \rho \in \mathbb{M}, \ [\![P_c]\!]_{rel} \equiv [\![P_s]\!]_{rel} \text{ modulo } \pi_{I}, \pi_{x}$

- This new definition is much weaker:
 - Correctness assumes no relation about
 - intermediate control states
 - non terminating executions
 - More compilers are considered correct
 - Weaker relation between source and compiled programs This new definition really misses something, and impedes verification

Ways to circumvent the limitation:

- Include the whole trace into the final state! Back to the previous definition, hard to formalize, says nothing about ∞...
- **2** Better way: get it right first and choose the right semantics!

Choice of another concrete semantics

We have built our definition of compilation correctness upon **operational** (trace) semantics.

What if we abstracted into another observational semantics ?

Alternate choice: let us consider a more abstract semantics

For instance, relational semantics (equivalent to denotational semantics)

- Notation forinitial (resp. final) control states: l_{\vdash} (resp. l_{\dashv})
- Notation for non-termination written ∞ ;
- Observational semantics: relations between $\mathbb M$ and $\mathbb M\cup\{\infty\}$
- Observational abstraction defined by collecting for all traces:

$$\begin{array}{rcl} \langle (\ell_{\vdash}, \rho), \dots, (\ell_{\dashv}, \rho') \rangle & \mapsto & (\rho, \rho') \\ \sigma = \langle (\ell_{\vdash}, \rho), \dots \rangle & \mapsto & (\rho, \infty) \text{ if } \sigma \text{ infinite} \end{array}$$

• Denotational semantics defined by:

$$\llbracket P \rrbracket_{\mathrm{rel}} = \{ (\rho, \rho') \mid \ldots \} \uplus \{ (\rho, \infty) \mid \ldots \}$$

Xavier Rival (INRIA)

Outline

- Introduction to program transformations
- 2 Compilation correctness
- 3 Correctness of optimizing compilation
 - 4 Application to the verification of compiled code
- 5 Application to certified compilation
- 6 Conclusion

Optimizations

Until now we focused on non-optimizing compilation

In practice, compilers perform various optimizations

- Elimination: dead code, dead variables...
- Instruction scheduling: Instruction-Level-Parallelism...
- Global transformations: Propagation of common expressions...
- Structural transformations: Loop unrolling...

Consequences: π_{I} , π_{x} , \mathbb{L}'_{i} , \mathbb{X}'_{i} may not be defined

Framework extension:

- Redefine the "most precise observation preserved by compilation"
- Would be more difficult with bissimulations
- Next slides: consider a few optimizations...

Xavier Rival (INRIA)

Correctness of optimizing compilation

Dead-code elimination definition

Principle

Do not compile statements of the source program that provably never are executed

- This saves space as smaller executables get generated
- It also improves runtime as some tests may be removed (when they always produce the same result)

Example:

Dead-code elimination correctness

How to set up a formal definition of compilation, that considers dead-code elimination correct ?

- we have to abstract away all labels removed by the optimizations
- this is trivial: we should simply not include them in \mathbb{L}_{s}'
- thus, our previous definition of compilation correctness already accommodates dead-code elimination

Compilation correctness in presence of dead-code elimination
Same definition as before

Correctness of optimizing compilation

Dead-variable elimination definition

Principle

Discard entirely the variables that are never used anymore (the compiler may reuse cells of dead local variables as well)

- This obviously both saves space and improves runtime
- There is a caveat though: this may change the error semantics indeed, expressions may be optimized away, so a program that normally fails (*e.g.*, on a division by zero) may not fail after optimization

$$\begin{array}{l} x := y;\\ \text{while}(i < 10) \{\\ x := x + 1;\\ y := y - x - 1;\\ i := i + 1;\\ \}\\ \text{use}(x); \end{array}$$

- x read after the loop, but not y
- thus, y can be removed with no observable change
- the purple statement disappears
- but y does not disappear everywhere

Xavier Rival	(INRIA)
--------------	---------

Program Transformations

Dead-variable elimination correctness

How to set up a formal definition of compilation, that considers dead-variable elimination correct ?

- variables may need be removed at certain program points
- it is not possible to simply remove the dead variables from X_s altogether: in the example, this would not be correct, as y would be completely lost
- thus, $\pi_{\mathbf{x}}$ should be relational

Compilation correctness in presence of variable-code elimination Similar definition as before, but with $\pi_x : \mathbb{L}'_s \times \mathbb{X}'_s \to \mathbb{X}'_c$ instead.

Exercise: formalize the new definition, inspired from the previous one, and with $\pi_{\mathbf{x}} : \mathbb{L}'_s \times \mathbb{X}'_s \to \mathbb{X}'_c$ instead

Path modifying optimizations

Some optimization deeply modify the control flow paths:

- loop unrolling
- Ioop exchange
- loop tiling
- loop interchange
- flattening of conditions

Gains:

- more efficient code, due to fewer conditions (unrolling, tiling)
- enabling of other optimizations, *e.g.*, vectorization (tiling, interchange...)

In the next few slides, we consider the case of loop unrolling

Loop unrolling example

Assumption: a for loop run an even number of times (loop unrolling may also apply to loops run a non statically known number of times, but it is more complex in that case)

sour	ce code	optim	ized code
		l_0^o	i := 0;
l_0^s	i := 0;	l_1^o	while (i < 1000)
l_1^s	while(i < 1000)	l_2^o	x := x * y;
l_2^s	$\mathbf{x} := \mathbf{x} * \mathbf{y};$	l_3^o	$\mathtt{y}:=\mathtt{y}-\mathtt{1};$
l_3^s	y := y - 1;	l_4^o	x := x * y;
l_4^s	i := i + 1;	l_5^o	$\mathtt{y}:=\mathtt{y}-\mathtt{1};$
l_5^s	}	l_6^o	i := i + 2;
		l_7^o	}

Control state correspondence π_{I} is clearly broken:

$$\pi_{\mathbf{I}}: \left\{ \begin{array}{ccc} l_2^s & \leftrightarrow & l_2^o \\ l_2^s & \leftrightarrow & l_4^o \end{array} \right.$$

Xavier Rival (INRIA)

Loop unrolling source and assembly traces

We consider executions in the source and the optimized code, and only display control states at the assignment to x and the values of i, y:

• At the source code level:

control state	l_2^s	l_2^s	l_2^s	l_2^s
value of i	0	1	2	3
value of y	1200	1199	1198	1197

• At the compiled code level:

control state	l_2^o	l ₄ °	l_2^o	l_4^o
value of i	0	0	2	2
value of y	1200	1199	1198	1197

As expected:

- the correlation between the values of i and the other variables is lost
- the real correspondence is between values of other variables and iterations even-ness

Xavier Rival (INRIA)

Program Transformations

Loop unrolling observational abstractions

How to set up a formal definition of compilation, that accepts loop unrolling as correct ?

- the loop counter variable i should be excluded from $\mathbb{X}_s, \mathbb{X}_o$
- each control state in the source loop should be divided into a pair of labels, that carry an even-ness tab:

$$\begin{array}{cccc} \mathcal{L}_2^s & \mapsto & \mathcal{L}_2^{s,e}, \mathcal{L}_2^{s,o} \\ \mathcal{L}_3^s & \mapsto & \mathcal{L}_3^{s,e}, \mathcal{L}_3^{s,o} \\ \dots & \mapsto & \dots \end{array}$$

• the trace abstraction function Ψ_s^{tr} should map each loop body state into a state with a consistent iteration even-ness

This amounts to doing an even-ness based trace partitioning

Loop unrolling observational abstractions

We can consider the traces again:

source code	control state	l_2^s	l_2^s	l_2^s	l_2^s
	value of i	0	1	2	3
	value of y	1200	1199	1198	1197
source code, abstract	control state	$l_2^{s,e}$	$l_2^{s,o}$	$l_2^{s,e}$	$l_2^{s,o}$
	value of i	0	1	2	3
	value of y	1200	1199	1198	1197
optimized code	control state	l_2^o	l4 ^o	l_2^o	l_4^o
	value of i	0	0	2	2
	value of y	1200	1199	1198	1197

We observe the following control state correspondence:

Loop unrolling correctness

Then, the definition follows a very similar form as before:

Compilation correctness in presence of loop unrolling Similar definition as before, but with:

- trace partitioning α_s^r abstraction
- a mapping π_{I} that preserves even-ness

Instruction scheduling: instruction level parallelism

We now consider optimizations that modify the code **locally**, and take **instruction scheduling** as an example.

Instruction-level parallelism is a feature of modern processors:

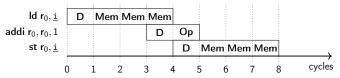
- one instruction = one or several cycles
 - memory typically slow: load, store take several cycles speed depends on the content of cache (hit/miss); can be 100 cycles!
 - arithmetic operations are usually faster
- Pipeline: run several instructions in parallel
- Some instructions cannot be evaluated in parallel due to dependences
- Scheduling: re-ordering of instructions so as to limit the number of *stall* cycles

Instruction level parallelism example

Assumptions:

- arith. instructions: 1 cycle instruction decoding, 1 cycle op.
- load/store instructions: 1 cycle instruction decoding, 3 cycle op.
- CPU: can have at the same time, one instruction in decoding, one in arithmetic stage, several doing memory read / write

We consider the code below:



Then, we observe a two cycles stall after the load

Consequence of this observation: instruction scheduling More efficient code is generated if there are more instructions between load/store instruction and uses of the values loaded/stored

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019

54 / 96

Instruction scheduling example

source code	non optimized code	optimized code
l_0^s i:=i+1;	<i>l</i> ₀ ^{<i>a</i>} ld r ₀ , <u>i</u>	ld r₀, <u>i</u>
	l_1^a addi $r_0, r_0, 1$	l_1^o ld $\mathbf{r}_1, \underline{\mathbf{x}}$
	l_2^a st $\mathbf{r}_0, \underline{\mathbf{i}}$	l_2^o addi $r_0, r_0, 1$
l_1^s $x := x + t[i];$	l_3^a ld r_1, \underline{x}	l_3^o ldx $\mathbf{r}_2, \mathbf{t}, \mathbf{r}_0$
	l_4^a ldx $\mathbf{r}_2, \underline{\mathbf{t}}, \mathbf{r}_0$	l_4^o st r_0, \underline{i}
	l_5^a add r_1, r_1, r_2	l_5^o add r_1, r_1, r_2
	l_6^a st r_1, \underline{x}	l_6^o st r_1, \underline{x}
$l_2^s \ldots$	l_7^a	$l_7^o \ldots$
Without optimization:	Without op	timization:

4	stall	cycles,	14	cycles	total
	ocum	<i>cyc.co</i> ,		0,0.00	cocar

l ₀	\leftrightarrow	l ₀
I_1^s	\leftrightarrow	ľŝ
I_2^s	\leftrightarrow	l_7^2

Without optimization: 2 stall cycles, 12 cycles total

$$egin{array}{cccc} I_0^s & \leftrightarrow & I_0^o \ I_1^s & \leftrightarrow & \ref{array}{cccc} I_1^s \ I_2^s & \leftrightarrow & I_7^o \end{array}$$

Instruction scheduling observational abstractions

Issues to fix our definition:

- Instructions execution order modified:
 - $l_1^a
 ightarrow l_2^a$ and $l_2^a
 ightarrow l_3^a$ are postponed
- Mapping π_{I} is broken:
 - The intermediate state l_1^s has no clear counterpart in the assembly
 - For i, it corresponds to l_5^o
 - For x, it corresponds to l_1^o
 - In general: this happens for all control points! (except for initial points, final points)

Thus, we need a relational mapping (π_{l}, π_{x}) , *i.e.*, a single function taking care of both variables and control states:

Relational syntactic mapping

A relational syntactic mapping is defined by an injective function

$$\pi_{\mathbb{X}\times\mathbb{X}}: (\mathbb{L}'_{s}\times\mathbb{X}'_{s}) \longrightarrow (\mathbb{L}_{c}\times\mathbb{X}_{c})$$

Xavier Rival (INRIA)

Instruction scheduling observational abstractions

Intuition

A source control state l^s corresponds to a **fictitious control state** where values of corresponding locations are gathered at different points in the execution of the optimized, compiled code

source code l_0^s i := i + 1;	optimized code f_0^o ld r_0, \underline{i}	We then have:
ℓ_1^s x := x + t[i];	$ \begin{array}{ll} \mathcal{L}_1^o & \text{Id } r_1, \underline{x} \\ \mathcal{L}_2^o & \text{addi } r_0, r_0, 1 \\ \mathcal{L}_3^o & \text{Idx } r_2, \underline{t}, r_0 \end{array} $	$ \begin{aligned} \pi_{\mathbb{X}\times\mathbb{X}} : & (l_0^s, \mathbf{i}) & \mapsto & (l_0^o, \underline{\mathbf{i}}) \\ & (l_0^s, \mathbf{x}) & \mapsto & (l_0^o, \underline{\mathbf{x}}) \\ & (l_1^s, \mathbf{i}) & \mapsto & (l_5^o, \underline{\mathbf{i}}) \end{aligned} $
l ₂ ^s	$\begin{array}{ll} \mathcal{l}_4^o & \text{st } \mathbf{r}_0, \underline{i} \\ \mathcal{l}_5^o & \text{add } \mathbf{r}_1, \mathbf{r}_1, \mathbf{r}_2 \\ \mathcal{l}_6^o & \text{st } \mathbf{r}_1, \underline{x} \\ \mathcal{l}_7^o & \dots \end{array}$	$\begin{array}{rccc} (l_1^s, \mathbf{x}) & \mapsto & (l_1^o, \underline{\mathbf{x}}) \\ (l_2^s, \mathbf{i}) & \mapsto & (l_7^o, \underline{\mathbf{i}}) \\ (l_2^s, \mathbf{x}) & \mapsto & (l_7^o, \underline{\mathbf{x}}) \end{array}$

Instruction scheduling correctness

The source level observational abstraction is unchanged.

Optimized level observational abstraction

Optimized code observational abstraction α_s^r abstracts traces into sequences of states observed at fictitious points

We now obtain:

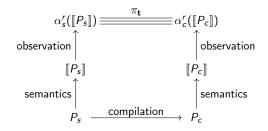
Compilation correctness in presence of instruction scheduling Similar definition as before, but with:

- optimized code observational abstraction α_s^r derived from $\pi_{\mathbb{X} \times \mathbb{X}}$
- semantic mapping π_t derived from $\pi_{\mathbb{X} \times \mathbb{X}}$

Compilation correctness

Definition: compilation correctness

Compilation of P_s into P_c is correct with respect to π_I, π_x (resp., $\pi_{\mathbb{X} \times \mathbb{X}}$) if and only if π_t establishes a bijection between $\alpha'_s(\llbracket P_s \rrbracket)$ and $\alpha'_c(\llbracket P_c \rrbracket)$.



Main idea: optimizations handled as standard compilation, but with more complex mappings, and observational abstractions

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019 59 / 96

On the formalization of program transformations

Methodology:

- Set up the standard semantics
- ② Define the observation preserved by the transformation
- Oerive the corresponding abstractions
- Stablish the correctness at the abstract level

Advantages of this approach:

- The framework can be extended (*e.g.*, with more complex abstractions)
- Abstract Interpretation theorems apply (e.g., fix-point transfers)

Other extensions:

- Define the transformation at the semantic level
- Derive an implementation of the transformation, from the definition

Outline

- Introduction to program transformations
- 2 Compilation correctness
- 3 Correctness of optimizing compilation
- 4 Application to the verification of compiled code
 - 5 Application to certified compilation
 - 6 Conclusion

Verifying compiled code

Kinds of properties:

- safety (no runtime errors, no overflows, no NaN...)
- security (no undesired information flow, in the sense of non-interference)

Two benefits:

- of course, verifying the generated code...
- but also, that the compiler does not turn a correct (already verified) program into an incorrect assembly one...

In the following, we consider safety properties and invariants

The invariant translation approach

Process

- **()** Analyze the source program P_s and compute an invariant \mathcal{I}_s
- **2** Translate \mathcal{I}_s into assembly level candidate invariant \mathcal{I}_t
- Perform an assembly level check of \mathcal{I}_t

Motivation:

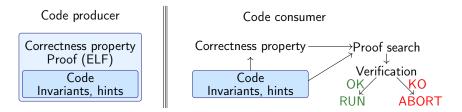
- inferring invariants is hard in general...
- and even more so at the assembly level

due to an important loss of structure at compile time (data-structures flattened, control flow more complex, additional steps to perform an arithmetic assignment –with separate load and store– or a test –with separate test and branching instructions)

Example 1: Proof Carrying Codes (PCC)

Principle:

- "Code producer": provides code and proof annotations in binaries (*i.e.*, proof of correctness),
- "Code consumer": checks the safety of the code
 - consistence of annotations: very quick proof search, from invariants
 - ② annotations \Rightarrow the safety property we wish to enforce



Context: execution of **non-trusted** code downloaded in the Internet *e.g.*, it could contain a **security bug** (information leak, buffer overflow)

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019

Example 2: TAL, compiled code certification by abstract interpretation

Typed and type safe assembly language:

- Java bytecode: interpreted (rather slow at runtime)
- TALx86: annotations for an assembly language closed to Intel 80x86
- Removing types \Rightarrow executable code
- A specific compiler translate source level types

Advantages:

- Ensure the safety of linkage thanks to types Linkage of object files usually *not* sound
- Improve the reliability of **optimizations** Constraint: they should *preserve* types!
- Compilation of type-safe versions of C (CCured, CClone)

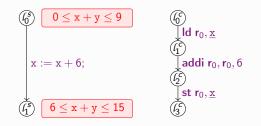
Certification of assembly code

Principle similar to PCC and TAL

but computation of invariants by abstract interpretation

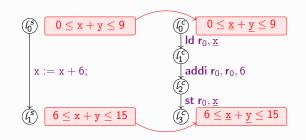
Xavier Rival (INRIA)

Assembly level verification of invariants



Start with invariants on the source code

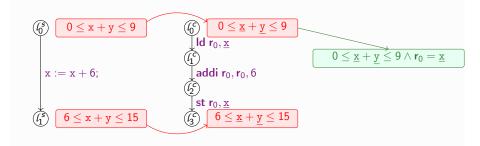
Assembly level verification of invariants



• Translates those invariants

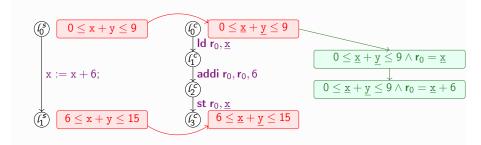
but not all control states are decorated

Assembly level verification of invariants



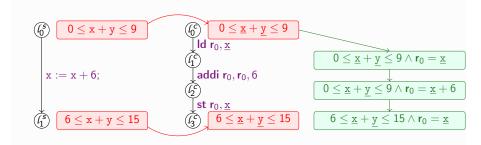
• Propagates the invariants and computes refined local invariants

Assembly level verification of invariants



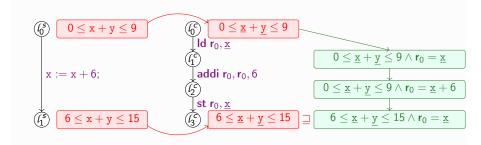
Propagates the invariants and computes refined local invariants

Assembly level verification of invariants



Propagates the invariants and computes refined local invariants

Assembly level verification of invariants



• Checks invariance at the end of the computation

Source static analysis: assumptions

• We assume an abstraction of sets of stores defined by an abstraction function for sets of stores

$$\alpha_{\mathsf{num}} : (\mathcal{P}(\mathbb{M}_s), \subseteq) \to (\mathbb{D}_{\mathsf{num}}^{\sharp}, \sqsubseteq)$$

• We derive an abstraction for sets of executions:

$$\begin{array}{rcl} \alpha_{i,s}: & \mathcal{P}(\mathbb{S}_{P}^{\star}) & \longrightarrow & \mathbb{L}_{s} \to \mathbb{D}_{\mathsf{num}}^{\sharp} \\ & X & \longmapsto & (\ell \in \mathbb{L}_{s}) \mapsto \alpha_{\mathsf{num}}(\{m \mid \langle \dots, (\ell, m), \dots \rangle \in X\}) \end{array}$$

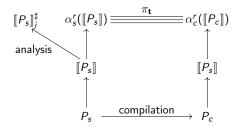
• We assume also a source code static analysis, that computes a sound over-approximation of the behaviors of the program:

$$\alpha_{i,s}(\llbracket P_s \rrbracket) \sqsubseteq \llbracket P_s \rrbracket_i^{\sharp}$$

Abstract invariant translation

Two abstractions have been defined:

- Abstraction for static analysis of Ps
- Abstraction for defining compilation correctness



Those abstractions are in general not comparable

Xavier Rival (INRIA)

Abstract invariant translation

We can derive **another abstraction**, more abstract than both α_s^r and $\alpha_{i,s}$:

- theoretical result: Galois-connections of a concrete domain form a lattice
- in practice, this common abstraction should abstract away all the elements that are not in L'_s, X'_s:

e.g., all dead variables, all unreachable control states...

e.g., in case of loop unrolling, it should perform the same trace partitioning

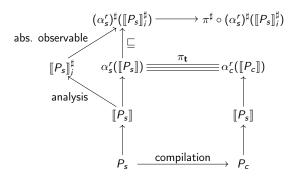
Moreover, $\pi_{\mathbf{l}}, \pi_{\mathbf{x}}$ induce a safe abstract invariant translation function $\pi^{\sharp} : (\mathbb{L}'_{s} \to \mathbb{D}^{\sharp}_{\mathsf{num}}) \to (\mathbb{L}'_{c} \to \mathbb{D}^{\sharp}_{\mathsf{num}})$

- for each pair of control points in correspondence in $\pi_{\rm I}$
- it maps numerical invariants among variables of *P_s* into numerical invariants among variables of *P_c*

Abstract invariant translation

Invariant translation process:

Apply π[#] to an abstract invariant [[P_s]][#]_i computed for P_s
 Result: a candidate invariant π[#]([[P_s]][#]_i) for P_c



Invariant translation: soundness

Soundness lemma

If:

- the compilation $P_s \rightarrow P_c$ is sound with respect to $\pi_{\mathbf{I}}, \pi_{\mathbf{x}}$;
- the analysis of P_s computes a sound $\llbracket P_s \rrbracket_i^{\sharp} \alpha_{i,s}(\llbracket P_s \rrbracket) \sqsubseteq \llbracket P_s \rrbracket_i^{\sharp}$

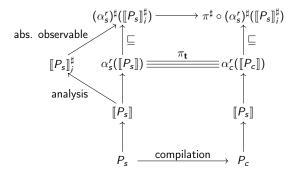
Then, $\pi^{\sharp}((\alpha_{s}^{r})^{\sharp}(\llbracket P_{s} \rrbracket_{i}^{\sharp}))$ is a sound approximation of $\llbracket P_{c} \rrbracket$:

$$\alpha_{i,r,c}(\llbracket P_c \rrbracket) \sqsubseteq \pi^{\sharp}((\alpha_s^r)^{\sharp}(\llbracket P_s \rrbracket_i^{\sharp}))$$

Consequence of the choice of another observational semantics for compilation correctness: If $\alpha_s^r(\llbracket P_s \rrbracket)$, $\alpha_c^r(\llbracket P_c \rrbracket)$ are weakened, then the invariants that can be translated are also weakened

Invariant translation: soundness

Proof summarized:



Assumptions are very strong:

compilation, analysis, translation need to be correct

We need an independent verification of translated invariants

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019

72 / 96

Independent verification of translated invariants

Principle of invariant checking: post-fixpoint checking

Theorem: invariant verification

Using a concretization function $\gamma,$

- The domain of the function F is a CPO,
- The *concrete* function *F* is continuous,
- $F \circ \gamma \subseteq \gamma \circ F^{\sharp}$,

```
• F^{\sharp}(x) \sqsubseteq x,
```

Then, Ifp $F \sqsubseteq \gamma(x)$

Proof left as exercise

• Only the verifier needs to be sound even if the assumptions of the translation soundness lemma are not met *i.e.*, we can have an incorrect compiler, translate an incorrect invariant, and still obtain and check a correct translated invariant !

Xavier Rival (INRIA)

Program Transformations

Independent verification of translated invariants

Principle of invariant checking: post-fixpoint checking

Theorem: invariant verification

Using a concretization function $\gamma,$

- The domain of the function F is a CPO,
- The *concrete* function *F* is continuous,
- $F \circ \gamma \subseteq \gamma \circ F^{\sharp}$,

•
$$F^{\sharp}(x) \sqsubseteq x$$
,

Then, **Ifp** $F \sqsubseteq \gamma(x)$

Invariant checking refines abstract predicates:

this phase also produces more precise abstract properties about:

- memory locations in $\mathbb{X}_c \setminus \mathbb{X}'_c$
- program points in $\mathbb{L}_c \setminus \mathbb{L}'_c$

```
In practice, every cycle of the compiled code control flow graph
Xavier Rival (INRIA) Program Transformations Jan, 8th, 2019
```

73 / 96

Invariant checking and difficulties

We consider the verification of invariants around a condition test Assumptions:

- $x \in [0, 12]$ at the entry point;
- we wish to verify the assert in the compiled code;
- we use a non relational abstract domain: intervals

Source code:

Compiled code:

$if(x \le 5)$ {	0	$Id \mathbf{r}_0, \underline{\mathbf{x}}$
$assert(\mathrm{x} \leq 5);$	4	li r ₁ ,5
	8	$cmp r_0, r_1$
}else{	12	$blt \langle GT \rangle \ell$ # (jump point)
	16	# true branch contents
}	l:	<pre># false branch contents</pre>

Invariant checking and difficulties

0: $\underline{x} \in [0, 12]$ $ld \mathbf{r}_0, \underline{x}$ 4: $li \mathbf{r}_1, 5$ 8: $cmp \mathbf{r}_0, \mathbf{r}_1$ 12: $blt \langle GT \rangle \ell$ # (jump point) 16:

Invariant checking and difficulties

```
0: \underline{x} \in [0, 12]

ld r_0, \underline{x}

4: \underline{x} \in [0, 12] \land r_0 \in [0, 12]

li r_1, 5

8:

cmp r_0, r_1

12:

blt \langle GT \rangle \ell # (jump point)

16:
```

Invariant checking and difficulties

16 :

Invariant checking and difficulties

- $\begin{array}{ll} 0: & \underline{\mathbf{x}} \in [0, 12] \\ & \textbf{Id} \ \textbf{r}_0, \mathbf{x} \end{array}$
- 8: $\underline{x} \in [0, 12] \land r_0 \in [0, 12] \land r_1 \in [5, 5]$ cmp r_0, r_1
- $\begin{array}{ll} 12: & \underline{x} \in [0,12] \wedge r_0 \in [0,12] \wedge r_1 \in [5,5] \wedge cr \in \{\mathrm{LT},\mathrm{EQ},\mathrm{GT}\} \\ & \mbox{blt}\langle \mathrm{GT} \rangle \ \ell & \mbox{ \# (jump point)} \end{array}$

16 :

Invariant checking and difficulties

- $\begin{array}{ll} 0: & \underline{\mathbf{x}} \in [0, 12] \\ & \text{Id } \mathbf{r}_0, \mathbf{x} \end{array}$
- 8: $\underline{x} \in [0, 12] \land r_0 \in [0, 12] \land r_1 \in [5, 5]$ cmp r_0, r_1
- $\textbf{16}: \quad \underline{x} \in [0,12] \land \textbf{r}_0 \in [0,12] \land \textbf{r}_1 \in [5,5] \land \textbf{cr} \in \{\mathrm{LT},\mathrm{EQ}\}$

Invariant checking and difficulties

- $\begin{array}{ll} 0: & \underline{x} \in [0,12] \\ & \text{Id } r_0, x \end{array}$
- 8: $\underline{x} \in [0, 12] \land r_0 \in [0, 12] \land r_1 \in [5, 5]$ cmp r_0, r_1
- $\textbf{16}: \quad \underline{x} \in [0,12] \land \textbf{r}_0 \in [0,12] \land \textbf{r}_1 \in [5,5] \land \textbf{cr} \in \{\mathrm{LT},\mathrm{EQ}\}$

The condition at the branch point is not precise

The range of x was not refined by the test:

- the test and branching are **independent** relations between test results and values need be tracked
- the test is made on a copy of x equalities between copies need be tracked by the verifier

Xavier Rival (INRIA)

Program Transformations

Refinement of the verifier

Relation between test and branching:

- \bullet each value in $\{LT, EQ, GT\}$ should be bound to the ranges of the other location
- this is obtained by a value partitioning, based on the value of cr:

$$\begin{array}{rcl} \gamma: & (\{\mathrm{LT}, \mathrm{EQ}, \mathrm{GT}\} \to \mathbb{D}_{\mathsf{num}}^{\sharp}) & \longrightarrow & \mathcal{P}(\mathbb{M}) \\ & \phi^{\sharp} & \longmapsto & \{m \mid m \in \gamma_{\mathsf{num}} \circ \phi^{\sharp} \circ m(\mathsf{cr})\} \end{array}$$

Equalities between copies, *e.g.*, of \underline{x} and \mathbf{r}_0 :

- an equality abstraction abstracts partitions of X_c
- \bullet replacement of $\mathbb{D}_{num}^{\sharp}$ with a reduced product of $\mathbb{D}_{num}^{\sharp}$ and an equality abstraction

Invariant checking: fixed

 $blt \langle GT \rangle \ell$ # (jump point)

16 :

In general, invariant checking is incomplete... It may require some refinement in the verifier

Xavier Rival (INRIA)

Program Transformations

Invariant checking: fixed

 $\begin{array}{lll} 0: & \underline{x} \in [0,12] \\ & \mbox{Id} \ \mbox{r}_0, \underline{x} \\ 4: & \underline{x} \in [0,12] \wedge \mbox{r}_0 \in [0,12] \wedge \underline{x} = \mbox{r}_0 \\ & \mbox{Ii} \ \mbox{r}_1, 5 \\ 8: \end{array}$

 $cmp r_0, r_1$

12 :

 $blt \langle GT \rangle \ \ell$ # (jump point)

16 :

In general, invariant checking is incomplete... It may require some refinement in the verifier

Xavier Rival (INRIA)

Program Transformations

Invariant checking: fixed

- $\begin{array}{ll} \mbox{li } r_1, 5 \\ 8: & \underline{x} \in [0, 12] \wedge r_0 \in [0, 12] \wedge r_1 \in [5, 5] \wedge \underline{x} = r_0 \\ \mbox{cmp } r_0, r_1 \end{array}$

12 :

 $blt \langle GT \rangle \ell$ # (jump point)

16 :

In general, invariant checking is incomplete... It may require some refinement in the verifier

Xavier Rival (INRIA)

Program Transformations

Invariant checking: fixed

16 :

In general, invariant checking is incomplete... It may require some refinement in the verifier

Xavier Rival (INRIA)

Program Transformations

Invariant checking: fixed

In general, invariant checking is incomplete... It may require some refinement in the verifier

Xavier Rival (INRIA)

Program Transformations

Outline

- Introduction to program transformations
- 2 Compilation correctness
- 3 Correctness of optimizing compilation
- 4 Application to the verification of compiled code
- 5 Application to certified compilation
 - 6 Conclusion

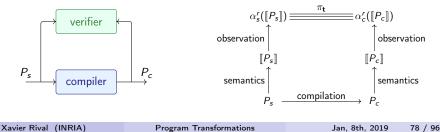
Verifying a compiler result

Principle: verify the semantic equivalence between source and compiled programs

Verification process: translation validation

- **()** Establish mappings $\pi_{\mathbf{I}}, \pi_{\mathbf{x}}$ between source and compiled programs
- Prove (with a specialized prover) the semantic equivalence of each basic block

Process:



Application to certified compilation

A technique based on fixpoint transfer

Foundation: fixpoint transfer

Theorem

Let $F_s : \mathcal{P}(\mathbb{S}_s^*) \to \mathcal{P}(\mathbb{S}_s^*)$ and $F_c : \mathcal{P}(\mathbb{S}_c^*) \to \mathcal{P}(\mathbb{S}_c^*)$ and $\pi_t : \mathbb{S}_s^* \to \mathbb{S}_c^*$ (complete for join), such that:

- F_s , F_c are monotone
- $\pi_{t}(\emptyset) = \emptyset$ (\emptyset least element);
- $\pi_{\mathbf{t}} \circ F_{s} = F_{c} \circ \pi_{\mathbf{t}}$

then both functions have a least fixpoint and:

 $\mathsf{lfp}\,F_c = \pi_{\mathsf{t}}(\mathsf{lfp}\,F_s)$

Proof: exercise

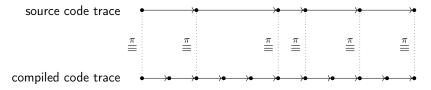
But the theorem does not apply directly: source and compiled executions are not correlated step-by-step

Xavier Rival (INRIA)

Application to certified compilation

A technique based on fixpoint transfer

Equivalence of source and assembly traces:



- standard semantics [[P_s]] and [[P_c]] are expressed as least fixpoints, but not directly correlated by π_x, π_l
- observational semantics α^r_s([[P_s]]) and α^r_c([[P_c]]) are directly correlated by not expressed as least fixpoint

We need fixpoint definitions for $\alpha_s^r(\llbracket P_s \rrbracket), \alpha_c^r(\llbracket P_c \rrbracket)$ (e.g., each basic block in the assembly code should be one computation step)

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019 80 / 96

Symbolic transfer functions: definition

A language to describe the effect of a basic block

- basic blocks usually contain series of assignment: we flatten sequences of assignments into parallel assignments
- a basic block may branch to several points (often two)
- no loop: each cycle in the compiled code control flow graph is associated to at least one control state in the source

Symbolic transfer functions

Symbolic transfer functions are defined by the grammar:

$$\begin{array}{lll} \delta(\in \mathbb{T}) & ::= & \Box & \text{no transition (dead branch, error)} \\ & & \left\lfloor \overrightarrow{x} \leftarrow \overrightarrow{e} \right\rfloor & \text{parallel assignment} \\ & & \left\lfloor c ? \delta_0 \mid \delta_1 \right\rfloor & \text{conditional} \end{array}$$

Intuitively, a symbolic transfer function is a store transformer

Symbolic transfer functions: semantics

Semantic domain:

- ullet \perp corresponds to the absence of behavior (error, blocking)
- $\llbracket \delta \rrbracket \in \mathbb{M} \to \mathbb{M} \cup \{\bot\}$

Denotational Semantics:

•
$$\llbracket \Box \rrbracket(\rho) = \bot$$

• $\llbracket \lfloor \vec{x} \leftarrow \vec{e} \rfloor \rrbracket(\rho) = \rho [\forall i, \llbracket x_i \rrbracket(\rho) \leftarrow \llbracket e_i \rrbracket(\rho)]$
if $\forall i, \llbracket x_i \rrbracket(\rho) \neq \text{error and } \forall i, \llbracket e_i \rrbracket(\rho) \neq \text{error}$
 $\llbracket \lfloor x \leftarrow e \rfloor \rrbracket(\rho) = \bot$ otherwise
• $\llbracket \lfloor e ? \delta_0 \mid \delta_1 \rfloor \rrbracket(\rho) = \begin{cases} \llbracket \delta_0 \rrbracket(\rho) & \text{if } \llbracket e \rrbracket(\rho) = \text{true} \\ \llbracket \delta_1 \rrbracket(\rho) & \text{if } \llbracket e \rrbracket(\rho) = \text{false} \\ \bot & \text{if } \llbracket e \rrbracket(\rho) = \text{error} \end{cases}$

Note: observe the identity is described by $\iota = \lfloor \cdot \leftarrow \cdot \rfloor$ (parallel assignment, with empty support)

Xavier Rival (INRIA)

Symbolic transfer functions: example

Encoding of a few instructions:

• "Addition" l_0 : addi r_0, r_1, v ; l_1 : ...:

$$\delta_{\ell_0,\ell_1} = \lfloor \mathbf{r}_0 \leftarrow \mathbf{r}_1 + \mathbf{v} \rfloor$$

• "Comparison" $\ell_0 : cmp \ r_0, r_1; \ \ell_1 : \ldots$:

$$\begin{split} \delta_{f_0,f_1} &= \lfloor \mathbf{r}_0 < \mathbf{r}_1 ?\\ & \lfloor \mathbf{cr} \leftarrow \mathrm{LT} \rfloor\\ & \mid \lfloor \mathbf{r}_0 = \mathbf{r}_1 ? \lfloor \mathbf{cr} \leftarrow \mathrm{EQ} \rfloor \mid \lfloor \mathbf{cr} \leftarrow \mathrm{GT} \rfloor \rfloor \rfloor \end{split}$$

• "Conditional branching" $l_0 : blt \langle LT \rangle \ l_1; \ l_2 : \ldots$:

$$\begin{split} \delta_{\textit{f}_0,\textit{f}_1} &= \lfloor \mathbf{cr} = \mathrm{LT} ? \iota \mid \Box \rfloor \\ \delta_{\textit{f}_0,\textit{f}_2} &= \lfloor \mathbf{cr} = \mathrm{LT} ? \Box \mid \iota \rfloor \end{split}$$

Symbolic transfer functions: example

Encoding of a few instructions:

• "Load" l_0 : ldx r_d , o, r_x ; l_1 :...:

$$\delta_{l_0,l_1} = \lfloor \mathbf{r}_d \leftarrow \mu(o + \mathbf{r}_x) \rfloor$$

• "Load" $l_0 : Id r_d, o; l_1 : ...:$

$$\delta_{\mathit{l}_{0},\mathit{l}_{1}} = \lfloor \mathbf{r}_{d} \leftarrow \mu(o) \rfloor$$

• "Store" $l_0 : stx r_d, o, r_x; l_1 : ...:$

$$\delta_{l_0,l_1} = \lfloor \mu(o + \mathbf{r}_x) \leftarrow \mathbf{r}_d \rfloor$$

The encoding of the source semantics is straightforward

Symbolic transfer functions: composition operation

Assumptions: memory locations are either equal or non-overlapping

Theorem

We can define a fully syntactic composition operation $\otimes:\mathbb{T}\times\mathbb{T}\to\mathbb{T}$ such that:

 $\llbracket \delta_0 \otimes \delta_1 \rrbracket \simeq \llbracket \delta_0 \rrbracket \circ \llbracket \delta_1 \rrbracket$

Full proof left as exercise; we consider a few cases:

•
$$\Box \otimes \delta = \Box$$

• $\delta \otimes \Box = \Box$
• $\delta \otimes \lfloor c ? \delta_0 \mid \delta_1 \rfloor = \lfloor c ? \delta \otimes \delta_0 \mid \delta \otimes \delta_1 \rfloor$
• $\lfloor x_0 \leftarrow e_0 \rfloor \otimes \lfloor x_1 \leftarrow e_1 \rfloor = \begin{cases} \lfloor x_0 \leftarrow e_0 [x_1 \leftarrow e_1] \rfloor & \text{if } x_0 = x_1 \\ x_1 \leftarrow e_1 \end{bmatrix} & \text{otherwise} \\ \text{when aliasing cannot be determined statically, use a symbolic} \\ \text{predicate is} alias(x, y), \text{ and return } \lfloor \text{is} alias(x, y) ? \delta_0 \mid \delta_1 \rfloor \end{cases}$

Xavier Rival (INRIA)

Symbolic transfer functions: composition operation

Example:

• no aliasing between
$$x, y, z$$

(*i.e.*, locations x, y, z are disjoint pairwise)
• $\delta_0 = \begin{bmatrix} x & \leftarrow & y+4 \\ y & \leftarrow & 3 \end{bmatrix}$
• $\delta_1 = \lfloor y \leftarrow z+1 \rfloor$
• here:

Then:

$$\delta_0 \otimes \delta_1 = \left[\begin{array}{ccc} x & \leftarrow & z+5 \\ y & \leftarrow & 3 \end{array} \right]$$

Note that y is overwritten, and the expression written into x takes into account that assignment

Translation validation with symbolic transfer functions

Application of symbolic transfer functions: Definition of a new program (labeled transition system) P'_c

Program Reduction

• States: L'_c

• \rightarrow is defined by a table of symbolic transfer functions: $(I, \rho) \rightarrow (I', \rho') \iff$ $\begin{cases} \exists l_0, \dots, l_n \in \mathbb{L}_c \setminus \mathbb{L}'_c, \\ \rho' = \llbracket \delta_{l_n, l'} \otimes \dots \otimes \delta_{l_i, l_{i+1}} \otimes \delta_{l_{i-1}, l_i} \otimes \dots \otimes \delta_{l, l_0} \rrbracket(\rho) \end{cases}$

Symbolic semantic abstraction

- Semantics: $\llbracket P'_c \rrbracket = \operatorname{lfp} F'_c$ where F'_c is derived from P'_c
- Soundness property: α^r_c([[P_c]]) = [[P^r_c]] = lfp F^r_c
 Proof: by induction on the length of the traces of P^r_c

Application to certified compilation

Translation validation: example (condition test)

Compiled code:

Source code:

- STF to the true branch: $\delta^{s} = \lfloor \mathbf{x} \leq 5 ? \iota \mid \Box \rfloor$

STF to
$$l$$
:

$$\delta_l^c = \lfloor \underline{\mathbf{x}} < 5 ?$$

$$\begin{bmatrix} \mathbf{r}_0 \leftarrow \mu(\underline{\mathbf{x}}) \\ \mathbf{r}_1 \leftarrow 5 \\ \mathbf{cr} \leftarrow \mathbf{LT} \end{bmatrix}$$

$$\| \dots \|$$
STF in P_c' :

$$\delta_l^c = \lfloor \underline{\mathbf{x}} < 5 ? \iota \mid \lfloor \underline{\mathbf{x}} = 5 ? \iota \mid \Box \rfloor \end{bmatrix}$$

Translation validation and optimization: instruction scheduling

source code	optimized code	Syntactic mappings:
$l_0^s \texttt{i} := \texttt{i} + 1;$	$\begin{array}{ll} \mathcal{I}_0^o & \text{Id } \mathbf{r}_0, \underline{i} \\ \mathcal{I}_1^o & \text{Id } \mathbf{r}_1, \underline{x} \\ \mathcal{I}_2^o & \text{addi } \mathbf{r}_0, \mathbf{r}_0, 1 \end{array}$	$\pi_{\mathbb{X}\times\mathbb{X}}: \begin{array}{ccc} (l_0^s, \mathbf{i}) & \mapsto & (l_0^o, \underline{i}) \\ (l_0^s, \mathbf{x}) & \mapsto & (l_0^o, \underline{x}) \\ \end{array}$
$l_1^s \mathtt{x} := \mathtt{x} + \mathtt{t}[\mathtt{i}];$	$ \begin{array}{c} \bar{\mathcal{L}}_{3}^{o} & \text{Idx } \mathbf{r}_{2}, \underline{\mathbf{t}}, \mathbf{r}_{0} \\ \mathcal{L}_{4}^{o} & \text{st } \mathbf{r}_{0}, \underline{\mathbf{i}} \\ \mathcal{L}_{5}^{o} & \text{add } \mathbf{r}_{1}, \mathbf{r}_{1}, \mathbf{r}_{2} \end{array} $	$egin{array}{cccc} (l_1^{s}, {f i}) &\mapsto & (l_5^{o}, {f \underline{i}}) \ (l_1^{s}, {f x}) &\mapsto & (l_1^{o}, {f \underline{x}}) \ (l_2^{s}, {f i}) &\mapsto & (l_7^{o}, {f \underline{i}}) \ (l_2^{s}, {f x}) &\mapsto & (l_7^{o}, {f \underline{x}}) \end{array}$
l_2^s	ℓ_6^o st $\mathbf{r}_1, \underline{\mathbf{x}}$ ℓ_7^o	Thus, $l_f^o = i \mathbb{Q} l_5^o$; $\mathbf{x} \mathbb{Q} l_1^o$

• Source level transfer functions:

$$\delta_{\ell_0^s,\ell_1^s} = \lfloor \mathtt{i} \leftarrow \mathtt{i} + 1 \rfloor \qquad \delta_{\ell_1^s,\ell_2^s} = \lfloor \mathtt{x} \leftarrow \mathtt{x} + \mathtt{t}[\mathtt{i}] \rfloor$$

• Optimized level transfer functions (registered not displayed):

 $\delta_{l_0^o,l_f^o} = \lfloor \mu(\mathtt{i}) \leftarrow \mu(\mathtt{i}) + 1 \rfloor \qquad \delta_{l_f^o,l_f^o} = \lfloor \mu(\underline{\mathtt{x}}) \leftarrow \mu(\underline{\mathtt{x}}) + \mu(\underline{\mathtt{t}} + \mu(\underline{\mathtt{i}})) \rfloor$

Translation validation and optimizations

Program reduction:

- produces a set of symbolic transfer functions that encode the transition relation of the program up-to observational abstraction
- abstracts the effect of optimizations
 as in the instruction scheduling example
 loop unrolling would result into unrolling at the source level
 (partitioning)

Translation validation:

based on a specialized prover, to establish equivalence of transfer functions

Outline

- Introduction to program transformations
- 2 Compilation correctness
- 3 Correctness of optimizing compilation
- 4 Application to the verification of compiled code
- 5 Application to certified compilation
- 6 Conclusion

Conclusion

Formalization of Compilation:

- At the concrete level: independent from analysis
- Very broad; works as well for
 - other architectures
 - optimizations (use of other abstractions)
- Algorithms for certified compilation described in the abstract interpretation frameworks:
 - Invariant translation
 - Invariant checking
 - Translation validation
 - Compiler formal certification

Symbolic transfer functions and use in static analysis and program transformations.

This approach applies to other program transformations

Xavier Rival (INRIA)

Program Transformations

Jan, 8th, 2019 92 / 96

Semantics

Program transformations: P. Cousot and R. Cousot.
 Systematic design of program transformation frameworks by abstract interpretation.

In Conference Record of the 29th Symposium on Principles of Programming Languages (POPL'02), pages 178–190, Portland, Oregon, January 2002.

• Relation between types and static analysis:

P. Cousot,

Types as Abstract Interpretations.

In POPL'97, pages 316-331, Paris, January 1997.

- Symbolic transfer functions:
 - C. Colby and P. Lee.
 - Trace-based program analysis.

In 23rd POPL, pages 195–207, St. Petersburg Beach, (Florida USA), 1996.

Bibliography: Certified Compilation

Proof Carrying Codes: G. C. Necula.
 Proof-Carrying Code.
 In 24th POPL, pages 106–119, 1997.

- Typed Assembly languages:
 G. Morrisett, D. Tarditi, P. Cheng, C. Stone, R. Harper, and P. Lee.
 The TIL/ML Compiler: Performance and Safety Through Types.
 In WCSSS, 1996.
- Abstract invariant translation (after compilation):
 X. Rival.

Abstract Interpretation-based Certification of Assembly Code. In 4th VMCAI, New York (USA), 2003.

Conclusion

Bibliography: Certified Compilation

- Translation validation: A. Pnueli, O. Strichman, and M. Siegel. Translation Validation for Synchronous Languages. In *ICALP'98*, pages 235–246. Springer-Verlag, 1998.
- Formal proof:

X. Leroy.

Formal certification of a compiler back-end, or: programming a compiler with a proof assistant.

In POPL'06, Charleston, january 2006.

• A generic frameork:

X. Rival.

Symbolic-Transfer Function-Based Approaches to Compilation Certification

In POPL'04, Venice, january 2004.

Xavier Rival (INRIA)

Program Transformations

Assignment: proofs

Read the paper *Systematic design of program transformation frameworks by abstract interpretation*, by Patrick Cousot and Radhia Cousot

Proofs based on fixpoint techniques:

- **(**) show the correctness of the invariant checking algorithm (slide 73)
- **2** show the correctness of the translation validation theorem (slide 79)