Lecture 26
CIS 341: COMPILERS

#### Announcements

- HW 7: Optimization & Experiments
  - Due: April 29th
- Final Exam:
  - Thursday, May 7<sup>th</sup>
  - 9:00AM
  - Moore 216
- Visitor: Yaron Minsky of Jane St. Capital
  - Monday, April 27th
  - Lunch: noon 1:15 (Raisler Lounge) sign-up sheet on Piazza
  - Talk: 2:00 3:00 (Raisler Lounge)
     From Theory into Practice: the story of Incremental

# **COMPILER VERIFICATION**

Zdancewic CIS 341: Compilers

# **Compiler Correctness?**

• We have to relate the source and target language semantics across the compilation function  $\mathbb{C}[-]$ : source  $\rightarrow$  target.

```
cmd / st _{S} \mapsto^{*} SKIP / st'
iff
C[cmd] / C[st] _{T} \mapsto^{*} C[st']
```

- Is this enough?
- What if cmd goes into an infinite loop?

# **Comparing Behaviors**

- Consider two programs P1 and P2 possibly in different languages.
  e.g. P1 is an Oat program, P2 is its compilation to LL
- The semantics of the languages associate to each program a set of observable behaviors:

 $\mathfrak{B}(P)$  and  $\mathfrak{B}(P')$ 

• Note:  $|\mathfrak{B}(P)| = 1$  if P is deterministic, > 1 otherwise

# What is Observable?

• For C-like languages:

```
observable behavior ::=

| terminates(st) (i.e. observe the final state)

| diverges

| goeswrong
```

• For pure functional languages:

```
observable behavior ::=

| terminates(v) (i.e. observe the final value)

| diverges

| goeswrong
```

# What about I/O?

• Add a *trace* of input-output events performed:

	t	::= []	e :: t
coind.	Т	::= []	e :: T

(finite traces) (finite and infinite traces)

```
observable behavior ::=
| terminates(t, st)
| diverges(T)
| goeswrong(t)
```

terminates(t, st) (end in state st after trace t) diverges(T) (loop, producing trace T)

## **Examples**

- P1: print(1); / st ⇒ terminates(out(1)::[],st)
- P2:
   print(1); print(2); / st
   ⇒

terminates(out(1)::out(2)::[],st)

 P3: WHILE true DO print(1) END / st
 ⇒ diverges(out(1)::out(1):...)

• So  $\mathfrak{B}(P1) \neq \mathfrak{B}(P2) \neq \mathfrak{B}(P3)$ 

# **Bisimulation**

• Two programs P1 and P2 are bisimilar whenever:

 $\mathfrak{B}(P1) = \mathfrak{B}(P2)$ 

• The two programs are completely indistinguishable.

• But... this is often too strong in practice.

# **Compilation Reduces Nondeterminism**

- Some languages (like C) have underspecified behaviors:
  - Example: order of evaluation of expressions f() + g()
- Concurrent programs often permit nondeterminism
  - Classic optimizations can reduce this nondeterminism
  - Example:
    - a := x + 1; b := x + 1 || x := x+1

VS.

a := x + 1; b := a || x := x+1

- LLVM explicitly allows nondeterminism:
  - undef values (not part of LLVM lite)
  - see the discussion later

# **Backward Simulation**

• Program P2 can exhibit fewer behaviors than P1:

 $\mathfrak{B}(P1) \supseteq \mathfrak{B}(P2)$ 

- All of the behaviors of P2 are permitted by P1, though some of them may have been eliminated.
- Also called *refinement*.

# What about goeswrong?

• Compilers often translate away bad behaviors.

x := 1/y; x := 42 vs. (divide by 0 error)

x := 42 (always terminates)

- Justifications:
  - Compiled program does not "go wrong" because the program type checks or is otherwise formally verified
  - Or just "garbage in/garbage out"

# **Safe Backwards Simulation**

• Only require the compiled program's behaviors to agree if the source program could not go wrong:

goeswrong(t)  $\notin \mathfrak{B}(P1) \Rightarrow \mathfrak{B}(P1) \supseteq \mathfrak{B}(P2)$ 

 Idea: let S be the *functional specification* of the program: A set of behaviors not containing goeswrong(t).

- A program P satisfies the spec if  $\mathfrak{B}(P) \subseteq S$ 

• Lemma: If P2 is a safe backwards simulation of P1 and P1 satisfies the spec, then P2 does too.

# **Building Backward Simulations**



Idea: The event trace along a (target) sequence of steps originating from a compiled program must correspond to some source sequence. Tricky parts:

- Must consider all possible target steps
- If the compiler uses many target steps for once source step, we have invent some way of relating the intermediate states to the source.
- the compilation function goes the wrong way to help!

## **Safe Forwards Simulation**

• Source program's behaviors are a subset of the target's:

goeswrong(t)  $\notin \mathfrak{B}(P1) \Rightarrow \mathfrak{B}(P1) \subseteq \mathfrak{B}(P2)$ 

- P2 captures all the good behaviors of P1, but could exhibit more (possibly bad) behaviors.
- But: Forward simulation is significantly easier to prove:
  - Only need to show the existence of a compatible target trace.

## **Determinism!**

- Lemma: If P2 is deterministic then forward simulation implies backward simulation.
- Proof:  $\emptyset \subset \mathfrak{B}(P1) \subseteq \mathfrak{B}(P2) = \{b\}$  so  $\mathfrak{B}(P1) = \{b\}$ .
- Corollary: safe forward simulation implies safe backward simulation if P2 is deterministic.

# **Forward Simulations**



Idea: Show that every transition in the source program:is simulated by some sequence of transitions in the target

- while preserving a relation ~ between the states

# **Lock-step Forward Simulation**



A single source-program step is simulated by a single target step.

(Solid = assumptions, Dashed = must be shown)

# "Plus"-step Forward Simulation



A single source-program step is simulated by *one or more* target steps. (But only finitely many!)

(Solid = assumptions, Dashed = must be shown)

# **Optional Forward Simulation**



A single source-program step is simulated by zero steps in the target.

# **Problem with "Infinite Stuttering"**



An infinite sequence of source transitions can be "simulated" by 0 transitions in the target!

(This simulation doesn't preserve nontermination.)

#### **Solution: Disallow such "trivial" simulations**



Equip the source language with a measure  $|\sigma|$  and require that  $|\sigma_2| < |\sigma_1|$ .

The measure can't decrease indefinitely, so the target program must either take a step or the source must terminate.

The target diverges if the source program does.

# **Is Backward Simulation Hopeless?**

- Suppose the source & target languages are the same.
  - So they share the same definition of program state.
- Further suppose that the steps are very "small".
  - Abstract machine (i.e. no "complex" instructions).
- Further suppose that "compilation" is only a very minor change.
  - add or remove a single instruction
  - substitute a value for a variable
- Then: backward simulation is more achievable
  - it's easier to invent the "decompilation" function because the "compilation" function is close to trivial
- Happily: This is the situation for many LLVM optimizations

## **Lock-Step Backward Simulation**



o is either an "observable event" or a "silent event" o ::= e |  $\epsilon$ 

Example use: proving variable substitution correct.

# **Right-Option Backward Simulation**



- Either:
  - the source and target are in lock-step simulation.
  - Or
  - the source takes a silent transition to a smaller state

Example use: removing an instruction in the target.

# **Left-Option Backward Simulation**



- Either:
  - the source and target are in lock-step simulation.
  - Or
  - the target takes a silent transition to a smaller state

Example use: adding an instruction to the target.

Verifying optimizations at the LLVM level of abstraction.

# **EXAMPLE: VELLVM**

Zdancewic CIS 341: Compilers

# **Step 1: Define LLVM IR Semantics**

- Essentially: define an interpreter for LLVM IR code
- But: more complex than the LLVMlite we use in class
  - Aggregate / Structured data
  - Undefined behaviors
  - Nondeterminism
- So: can't be just an interpreter
  - Semantics is given by a relation

# **Other Parts of the LLVM IR**

```
op ::= %uid | constant | undef
bop ::= add | sub | mul | shl | ...
cmpop ::= eq | ne | slt | sle | ...
insn ::=
   %uid = alloca ty
   %uid = load ty op1
   store ty op1, op2
   %uid = getelementptr ty op1 ...
   %uid = call rt fun(...args...)
   •••
phi ::=
 \Phi [op1;lbl1]...[opn;lbln]
terminator ::=
   ret %ty op
   br op label %lbl1, label %lbl2
  br label %1b1
```

Operands Operations Comparison

Stack Allocation Load Store Address Calculation Function Calls

# **Sources of Undefined Behavior**

#### Target-dependent Results

• Uninitialized variables:

%v = add i32 %x, undef

• Uninitialized memory:

%ptr = alloca i32
%v = load (i32\*) %ptr

• Ill-typed memory usage

Nondeterminism

#### Fatal Errors

- Out-of-bounds accesses
- Access dangling pointers
- Free invalid pointers
- Invalid indirect calls

#### Stuck States

## **Sources of Undefined Behavior**

#### Target-dependent Results

• Uninitialized variables:

%v = add i32 %x, undef

• Uninitialized memory:

%ptr = alloca i32
%v = load (i32\*) %ptr

• Ill-typed memory usage

Nondeterminism

Defined by a predicate on the program configuration.

A program configuration is *stuck* if there is no transition it can make.



## LLVM's memory model

%ST = type {i10,[10 x i8\*]}

High-level Representation

i10
i8*

• Manipulate structured types.

```
%val = load %ST* %ptr
...
store %ST* %ptr, %new
```

# LLVM's memory model

	%ST =	= <b>type</b> { Low-level	[i10	,[1	10 x	i8*]}
High-level	R	epresentatio	n	N	lanipula	ate struct
Representatio	on	b(10, 136)	0		iamp an	
i10		b(10, 2)	1		%val :	= load
:0*		uninit	2		•••	
10*		uninit	3		store	%ST* %
18*		ptr(Blk32,0,0)	4	Se	emantic	s is give
18*		ptr(Blk32,0,1)	5	b	vte-oriented l	nted low
i8*		ptr(Blk32,0,2)	6	~	- naddi	ng & alig
i8*		ptr(Blk32,0,3)	7	_		
i8*		$ptr(B k_{32} \otimes 0)$	8	-	- physic	cal subtyp
i8*		ptr(Blk32,8,1)	9			
i8*		ptr(DIK32,0,1)	1			
i8*		ри(ыкз2,6,2)	0			
i8*		ptr(Blk32,8,3)	1 1			
			1 2			

tructured types.

```
oad %ST* %ptr
[* %ptr, %new
```

- given in terms of low-level memory.
  - alignment
  - ubtyping

# **Adapting CompCert's Memory Model**



- Data lives in blocks
- Represent pointers abstractly
   block + offset
- Deallocate by invalidating blocks
- Allocate by creating new blocks
  - infinite memory available

# **Dynamic Physical Subtyping**

[Nita, et al. *POPL* '08]



#### undef

• What is the value of **%y** after running the following?

- One plausible answer: 0
- Not LLVM's semantics!

(LLVM is more liberal to permit more aggressive optimizations)

#### undef

• Partially defined values are interpreted *nondeterministically* as sets of possible values:

[%x] = {a or b | a∈[i8 undef], b ∈[1]} = {1,3,5,...,255} [%y] = {a xor b | a∈[%x], b∈[%x]} = {0,2,4,...,254}

# **Nondeterministic Branches**



# **LLVM<sub>ND</sub> Operational Semantics**

• Define a transition relation:

$$f \vdash \sigma_1 \mapsto \sigma_2$$

- f is the program
- $\sigma$  is the program state: pc, locals( $\delta$ ), stack, heap
- Nondeterministic
  - $\delta$  maps local %uids to sets.
  - Step relation is nondeterministic
- Mostly straightforward (given the heap model)
  - Another wrinkle: phi-nodes executed atomically

#### **Need for Atomic Phi-node Updates**

blk:

%x = phi i32 [ %z, %blk ], [ 0, %pred ]
%z = phi i32 [ %x, %blk ], [ 1, %pred ]
%b = icmp leq %x %z
br %b %blk %succ

# **Operational Semantics**

	Small Step	Big Step
Nondeterministic	LLVM <sub>ND</sub>	
Deterministic		

# **Deterministic Refinement**

	Small Step	Big Step
Nondeterministic	LLVM <sub>ND</sub>	
Deterministic		

Instantiate 'undef' with default value (0 or null)  $\Rightarrow$  deterministic.

# **Big-step Deterministic Refinements**

	Small	Step	Big Step
Nondeterministic		LLVM <sub>ND</sub>	
Deterministic	LLVM <sub>Interp</sub>	$\approx$ LLVM <sub>D</sub>	

Bisimulation up to "observable events":

• external function calls

# **Big-step Deterministic Refinements**

	Small Step		Big Step
Nondeterministic		LLVM <sub>ND</sub>	
Deterministic	LLVM <sub>Interp</sub> ≈		$\gtrsim$ LLVM <sup>*</sup> <sub>DFn</sub> $\gtrsim$ LLVM <sup>*</sup> <sub>DB</sub>

Simulation up to "observable events":

- useful for encapsulating behavior of function calls
- large step evaluation of basic blocks

[Tristan, et al. *POPL '08*, Tristan, et al. *PLDI '09*]