

From INTERPRETATION TO COMPILATION

INTERPRETER FOR L : a program \mathcal{E} in a language L'
executing L -programs

$\mathcal{E}(P)$ = result of P

COMPILER FOR L IN L' : a program \mathcal{C} translating L -programs to
 L' -programs

Usually

L : source language

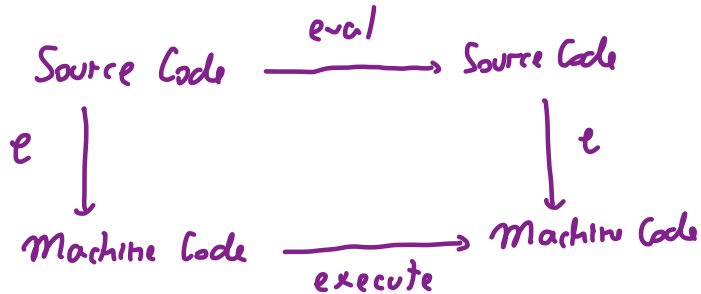
L' : machine language (executable by machine)

Focus on efficiency, energy consumption, code optimisation

INTENSIONAL PROGRAM BEHAVIOURS

What is compiler correctness?

↳ The generated code must meet the semantics of the source program



$$\text{execute}(e(p)) \stackrel{?}{=} e(\text{eval}(p))$$

⇒ Need formal semantics of both $\begin{cases} \text{source language} \\ \text{machine language} \end{cases}$

Why compiler correctness?

"We tested thirteen production-quality C-compilers and, for each, found situations in which the compiler generated incorrect code [...]"

Eide & Røgehaug Emsoft

"To improve the quality of C compilers, we created Csmith [...]
Every compiler we tested was found to crash and also to silently generate wrong code [...]"

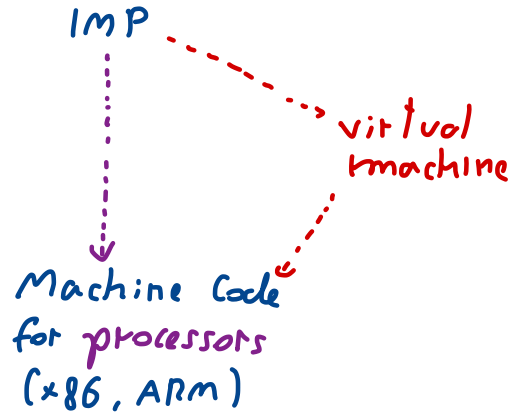
Yang, Chen, Eide & Røgehaug

PLDI 2011

Our goal : Formal verification of a (non-optimising)
compiler for Imp

Use techniques that scale to real world languages

How to compile IMP?



- similar to a real machine
 - program = sequences of instructions

⇒ no compound exp
no control structure

- close to the source lang.
(instructions reflect base operations in source lang.)

IMP virtual machine

4 components:

- Code C : list of instructions
- Code pointer pc : position of the instruction in C we are executing
- Store s : association variables-values
- Stack σ : list of integer values (to save intermediate results)

INSTRUCTIONS

Inductive *instr* : Type :=

- | Iconst (n: \mathbb{Z})
- | Ivar (x: ident)
- | Isetvar (x: ident)
- | Iadd
- | Iopp
- | Ibranch (d: \mathbb{Z})
- | Ibeq (d1: \mathbb{Z}) (d2: \mathbb{Z})
- | Ible (d1: \mathbb{Z}) (d2: \mathbb{Z})
- | Ihalt

- "push integer n"
- "push value of x"
- "pop integer and assign it to x"
- "pop two integers; push their sum"
- "pop one integer; push its opposite"
- "skip "forward" of d instructions"
- "pop two integers; skip d1 instructions if equal; skip d2 instructions if not"
- "stop"

NB. All instructions implement PC by 1; branching instructions of d+1

Ex.

$\langle \text{code}, \text{pc}, \text{store}, \text{stack} \rangle$ } configurations : state of the machine

$\langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 0, x \mapsto 12, [] \rangle$



Ex.

$\langle \text{code}, \text{pc}, \text{store}, \text{stack} \rangle$ } configurations : state of the machine

$\langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 0, x \mapsto 12, [] \rangle$



$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 1, x \mapsto 12, [12] \rangle$



Ex.

$\langle \text{code}, \text{pc}, \text{store}, \text{stack} \rangle$ } configurations : state of the machine

$\langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 0, x \mapsto 12, [] \rangle$



$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 1, x \mapsto 12, [12] \rangle$



$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 2, x \mapsto 12, [1, 12] \rangle$



Ex.

$\langle \text{code}, \text{pc}, \text{store}, \text{stack} \rangle$ } configurations : state of the machine

$\langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 0, x \mapsto 12, [] \rangle$

↑

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 1, x \mapsto 12, [12] \rangle$

↑

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 2, x \mapsto 12, [1, 12] \rangle$

↑

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 3, x \mapsto 12, [13] \rangle$

↑

Ex.

$\langle \text{code}, \text{pc}, \text{store}, \text{stack} \rangle$ } configurations : state of the machine

$\langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 0, x \mapsto 12, [] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 1, x \mapsto 12, [12] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 2, x \mapsto 12, [1, 12] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 3, x \mapsto 12, [13] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 4, x \mapsto 13, [] \rangle$

Ex.

$\langle \text{code}, \text{pc}, \text{store}, \text{stack} \rangle$ } configurations : state of the machine

$\langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 0, x \mapsto 12, [] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 1, x \mapsto 12, [12] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 2, x \mapsto 12, [1, 12] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 3, x \mapsto 12, [13] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 4, x \mapsto 13, [] \rangle$

$\rightarrow \langle \text{Ivar}(x); \text{Iconst}(1); \text{Iadd}; \text{Isetvar}(x); \text{Ibranch}(-5), 0, x \mapsto 13, [] \rangle$

We specify the vm using *operational semantics*

→ The pc -th element of C is $I(\text{const}(m))$

$$C[pc] = I(\text{const}(m))$$

$$\langle C, pc, s, \sigma \rangle \rightarrow \langle C, pc+1, s, m :: \sigma \rangle$$

$$C[pc] = I\text{Var}(x)$$

$$\langle C, pc, s, \sigma \rangle \rightarrow \langle C, pc+1, s, s(x) :: \sigma \rangle$$

$$C[pc] = I\text{SetVar}(x)$$

$$\langle C, pc, s, m :: \sigma \rangle \rightarrow \langle C, pc+1, s[x \mapsto m], \sigma \rangle$$

$$C[pc] = I_{add}$$

$$\langle C, pc, s, m :: m :: \sigma \rangle \rightarrow \langle C, pc+1, s, (m+m) :: \sigma \rangle$$

$$\vdots$$

$$C[pc] = I_{branch}(d)$$

$$\langle C, pc, s, \sigma \rangle \rightarrow \langle C, pc+1+d, s, \sigma \rangle$$

$$\vdots$$

How to encode op. sem in Coq?

Definition **code** := list instruction

Definition **store** := ident \rightarrow option \mathbb{Z}

Definition **stack** := list \mathbb{Z}

Definition **config** := $\mathbb{Z} * \text{store} * \text{stack}$ (\leadsto no code!)

Inductive **Transition** ($C : \text{code}$): $\underbrace{\text{config} \rightarrow \text{config} \rightarrow \text{Prop}}_{\text{A relation on configuration}} :=$

| rule_1const : $\forall pc \ s \ \sigma,$
 if $C[pc] = \text{Some } (Ivar \ n)$

\vdots

Behaviours of the VM

- Termination

$\langle C, pc, s, \sigma \rangle \rightarrow^* \langle C, pc', s', \sigma' \rangle$ and $C[pc'] = \text{Ihalt}$

Notation. $\langle C, pc, s, \sigma \rangle \Downarrow s'$

- Divergence

Infinitely many transitions from $\langle C, pc, s, \sigma \rangle$

Notation. $\langle C, pc, s, \sigma \rangle \Uparrow$

- Going wrong Otherwise

NB Machine programs can go wrong

E.g. $\langle Iadd; Ihalt, pc, s, \cdot \rangle$ (empty stack)

Digression : What is a program behaviour?

- Not a precisely defined notion...
- Semantics usually induce notions of behaviour
- In case of an SOS (Δ, \rightarrow) behaviour defined thus

① Specify a collection $\Delta_F \subseteq \Delta$ of final states

coherence. $\forall s \in \Delta_F. \neg \exists s' \in \Delta. s \rightarrow s'$

② Termination

$s \Downarrow$ iff $\exists s' \in \Delta_F. s \rightarrow^* s'$

Recall that, for $R \subseteq A \times A$, R^* is the reflexive and transitive closure of R .

$R^* = \text{smallest } \mathcal{K} \subseteq A \times A \text{ s.t.}$

- $R \subseteq \mathcal{K}$
- $\text{Id} \subseteq \mathcal{K}$
- $\mathcal{K}; \mathcal{K} \subseteq \mathcal{K}$

$$\left. \begin{array}{l} R \subseteq \mathcal{K} \\ \text{Id} \subseteq \mathcal{K} \\ \mathcal{K}; \mathcal{K} \subseteq \mathcal{K} \end{array} \right\} \Rightarrow \frac{R \subseteq \mathcal{Q} \quad \text{Id} \subseteq \mathcal{Q} \quad \mathcal{Q}; \mathcal{Q} \subseteq \mathcal{Q}}{R^* \subseteq \mathcal{Q}} \quad (*\text{-induction})$$

Constructively: $R^* = \bigcup_{n \geq 0} R^{(n)}$

where $R^{(0)} = \text{Id}$

$R^{(m+1)} = R; R^{(m)}$

$$\frac{\frac{a R a'}{a R^* a'} \quad \overline{a R^* a}}{\left(\frac{\overbrace{a R^* b} \quad b R^* c}{a R^* c} \right) a R b}$$

Termination

$$s \Downarrow : \text{ff} \quad \exists s'. s' \in \Delta_F \wedge s \rightarrow^* s'$$

Divergence

$$s \Uparrow : \text{ff} \quad \forall s'. s \rightarrow^* s' \Rightarrow \exists s''. s' \rightarrow s''$$

Goes Wrong

$$s \nabla : \text{ff} \quad \exists s'. \overbrace{s' \notin \Delta_F}^{\text{stack states}} \wedge s' \nrightarrow \wedge s \rightarrow^* s' \\ (\neg \exists s'', s' \rightarrow s'')$$

NB. These behaviours are exhaustive

In case of VM

$$\Delta = \{ \langle c, pc, s, \sigma \rangle \mid c \text{ code, } pc \in \mathbb{Z}, \dots \}$$

$$\Delta_F = \{ \langle c, pc, s, \sigma \rangle \mid \dots \quad c[pc] = \text{halt} \}$$

In case of IMP

$$\Delta = \{ (c, s) \mid c \in \text{Comm} \ \& \ s: \text{Var} \rightarrow \text{option}(\mathbb{Z}) \}$$

$$\Delta_F = \{ (\text{skip}, s) \mid s: \text{Var} \rightarrow \text{option}(\mathbb{Z}) \}$$

Prop. Programs do not go wrong:

$$\forall (c, s). \quad (c, s) \Downarrow \vee (c, s) \Uparrow$$

} safety properties

Compilation

We translate IMP programs into VM code

$\text{comp}_A: \text{AExp} \rightarrow \text{Code}$

$\text{comp}_A(x) = \text{Ivar}(x)$

$\text{comp}_A(n) = \text{Iconst}(n)$

$\text{comp}_A(e_1 + e_2) = \text{comp}_A(e_1) :: \text{comp}_A(e_2) :: \text{Iadd}$


reverse polish notation

$$\text{comp}_B : B\text{Exp} \times \mathbb{Z} \times \mathbb{Z} \rightarrow \text{Code}$$

$$\text{comp}_B (e_1 = e_2, d_0, d_1) = \text{comp}_A (e_1) :: \text{comp}_A (e_2) :: \text{Ieq} (d_0, d_1)$$

$$\text{comp}_B (e_1 \leq e_2, d_0, d_1) = \text{comp}_A (e_1) :: \text{comp}_A (e_2) :: \text{Ileq} (d_0, d_1)$$

$$\text{comp} : \text{Com} \rightarrow \text{Code}$$

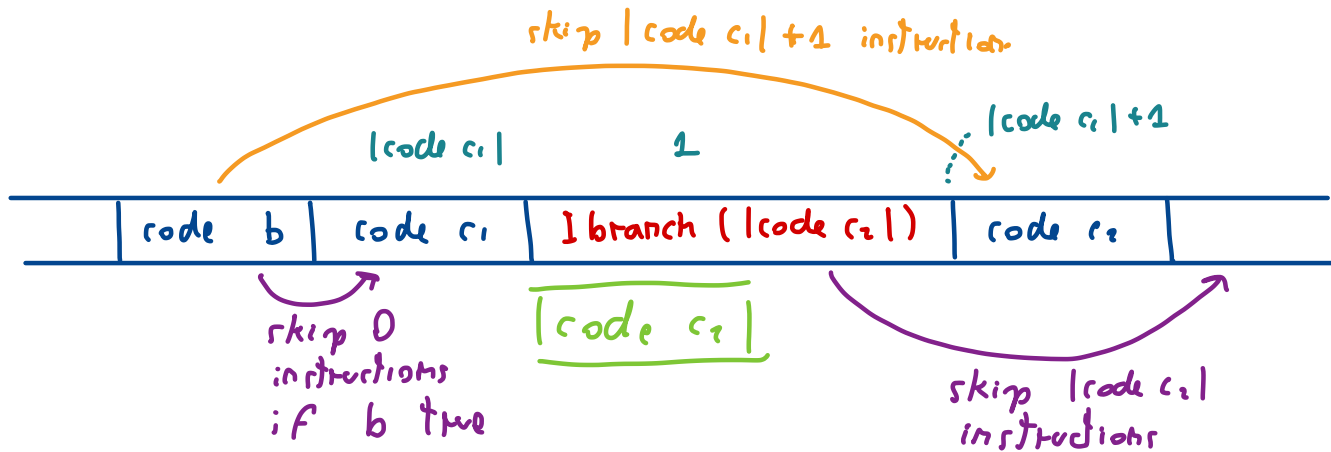
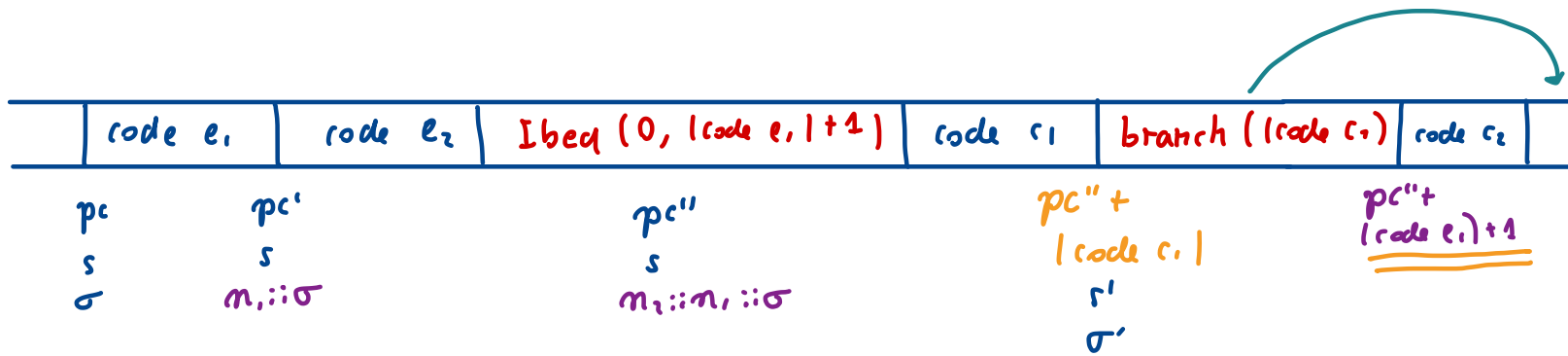
$$\text{comp} (\text{skip}) = []$$

$$\text{comp} (x := e) = \text{comp}(e) :: \text{Iset} \text{ var}(x)$$

$$\text{comp} (c_1 ; c_2) = \text{comp}(c_1) :: \text{comp}(c_2)$$

$$\begin{aligned} \text{comp} (\text{if } b \text{ then } c_1 \text{ else } c_2) = & \text{comp}_B (b, 0, |\text{comp}(c_1)| + 1) \\ & :: \text{comp}(c_1) :: \text{branch} (|\text{comp}(c_2)|) \\ & :: \text{comp}(c_2) \end{aligned}$$

?



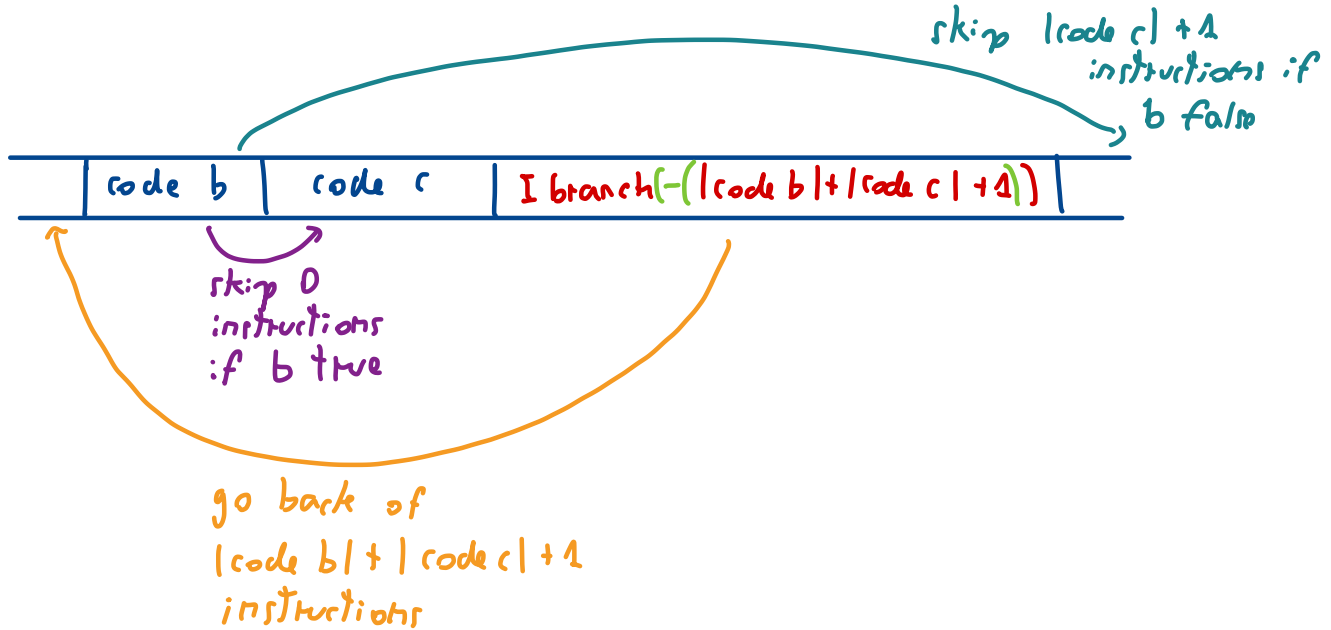
if b then c_1 else c_2

$|code\ c_1| = m_1$ $|code\ c_2| = m_2$

$\text{comp}(\text{while } b \text{ do } c) = \text{comp}(b, 0, |\text{code } c| + 1)$

$:: \text{comp}(c)$

$:: \text{Ibranch}(|\text{code } b| + |\text{code } c| + 1)$

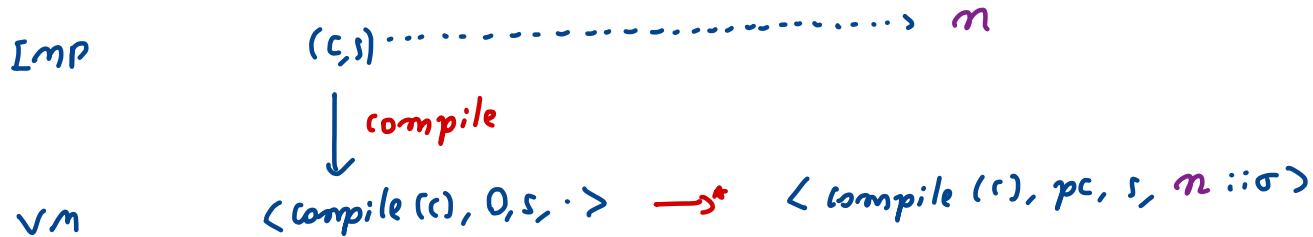


Finally compilation

$\text{compile}(c) = \text{comp}(c) :: \text{Ihoff}$

New way to execute IMP programs

- ① IMP operational semantics \leadsto interpreter
② VM operational semantics \leadsto compiler } Are they the same ??



Compiler Correctness

Intuition

$$\forall c \in \text{Com}, \forall s \in \text{Store}. \left(\langle c, s \rangle \Downarrow s' \implies \langle \text{compile}(c), 0, s, \cdot \rangle \Downarrow s' \right)$$

Q. . What about **divergence**?

. What if we require \iff ?

. What about **going wrong**?

Def. Given languages L_1, L_2 and programs $P_1 \in L_1, P_2 \in L_2$,
 let $\text{beh}(P_i)$ be the collection of **observable behaviours** of P_i .
 For us, possible behaviours are: $\Downarrow, \Uparrow, \bot$.

Then:

- $P_1 \approx P_2$ (**bisimulation**) if $\text{beh}(P_1) = \text{beh}(P_2)$
- $P_1 \leq P_2$ (**simulation**) if $\text{beh}(P_1) \subseteq \text{beh}(P_2)$
 - Forward simulation \leq
 - Backward simulation \geq
- $P_1 \sqsubseteq^f P_2$ (**correct-only forward sim**) if

$$\exists \notin \text{beh}(P_1) \Rightarrow \text{beh}(P_1) \subseteq \text{beh}(P_2)$$
- $P_1 \sqsupseteq^b P_2$ (**correct-only backward sim**) if

$$\exists \notin \text{beh}(P_1) \Rightarrow \text{beh}(P_1) \supseteq \text{beh}(P_2)$$

Which notion of correctness should we pick?

① **Bisimulation** is the **strongest** notion ($\approx \subseteq \lesssim \cap \gtrsim \cap \subseteq^f \cap \subseteq^b$)

\leadsto Usually too strong

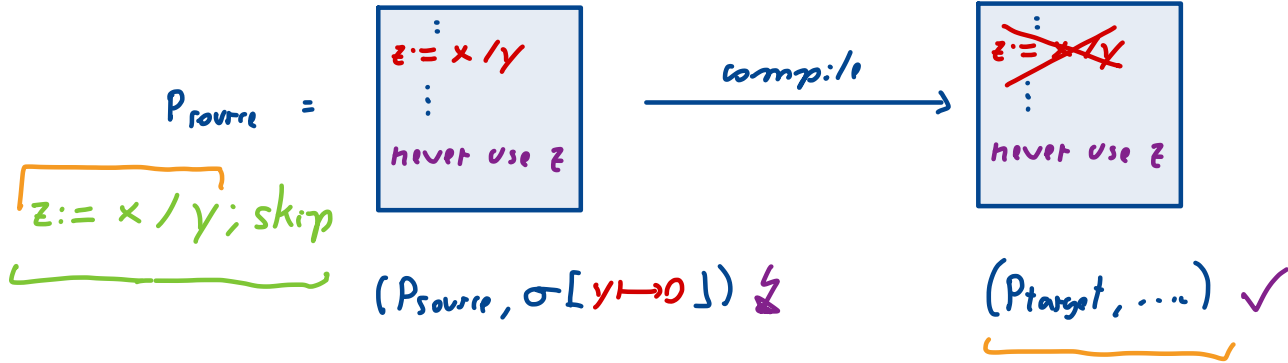
E.g. C has **nondeterministic** semantics (evaluation order of some expressions is not specified), but compilers choose an order while generating **deterministic** machine code

$$\text{beh}(P_{\text{source}}) \not\approx \text{beh}(P_{\text{target}})$$

$$P_{\text{source}} \gtrsim P_{\text{target}}$$

Lemma. If $(\mathcal{S}, \rightarrow)$ is such that \rightarrow is **deterministic** (i.e. $\forall s, s_1, s_2 \in \mathcal{S}. (s \rightarrow s_1 \wedge s \rightarrow s_2) \implies s_1 = s_2$), then $|\text{beh}(s)| = 1$, for any s

② Compilers optimise away 'going wrong' behaviours



\leadsto backward correct-only simulation

$$z \notin \text{beh}(P_{\text{source}}) \Rightarrow \text{beh}(P_{\text{source}}) \supseteq \text{beh}(P_{\text{target}})$$

$$P_{\text{source}} \supseteq^b P_{\text{target}}$$

$$\text{In particular: } P_{\text{source}} \not\sim \Rightarrow P_{\text{source}} \supseteq^b P_{\text{target}}$$

③ \exists^b seems the right notion for IMP and VM.

Working with \exists^b usually not easy.

What about \exists^f ?

(3.1) \exists^f easier than \exists^b

(3.2) \exists^f weaker than \exists^b : P_{target} can have more behaviours than P_{source}

but...

Lemma. If $(\mathcal{A}, \rightarrow)$ deterministic, then $\exists^f \subseteq \exists^b$

We can thus rely on \exists^f for IMP and VM

To prove correctness of compilation in VM, we thus have to prove **forward simulation**

$$\textcircled{2} \quad \forall c, s, s'. \left(\langle c, s \rangle \Downarrow s' \rightarrow \langle \text{compile}(c), 0, s, \cdot \rangle \Downarrow s' \right)$$

Sketch of the proof. Main techniques used

• induction on op. semantics or syntax

• case analysis

• inversion

syntactic shape

premises

IH
⋮

rule op. sem

conclusion rule
op. sem.

2.1 Arithmetic

Lemma. $\forall e_1, e_2 : \text{code}, \forall s, \forall \sigma.$

$$\langle e_1; \text{comp}(e); e_2, |e_1|, s, \sigma \rangle \rightarrow^*$$

$$\langle e_1; \text{comp}(e); e_2, |e_1| + |\text{comp}(e)|, s, \llbracket e \rrbracket(s) : \sigma \rangle$$

Corollary.

$$\forall e, s, \sigma. \langle \text{comp}(e), 0, s, \cdot \rangle \rightarrow^* \langle \text{comp}(e), |\text{comp}(e)|+1, s, \text{del}(s) \rangle \quad \because \text{denote } e \text{ } s$$

Proof. By induction on e

1.2 Boolean Expressions

Lemma $\forall e_1, e_2. \langle e_1; \text{comp}(b, d); e_2, |e_1|, s, \sigma \rangle$
 $\rightarrow^* \langle e_1; \text{comp}(b, d); e_2, pc, s, \sigma \rangle$

where $pc = \begin{cases} |e_1| + |\text{comp}(b)| & \text{if } \llbracket b \rrbracket s = \text{true} \\ |e_1| + |\text{comp}(b)| + \delta & \text{otherwise} \end{cases}$

Theorem

$\forall c, s, s', e_1, e_2, \sigma \dots$ implies

$$\langle c, s \rangle \Downarrow_{s'} \longrightarrow \langle e_1; \text{comp}(c); e_2, |e_1|, s, \sigma \rangle$$

$$\rightarrow^* \langle e_1; \text{comp}(c); e_2, |e_1| + |\text{comp}(c)|, s', \sigma \rangle$$

Proof, Induction on $\langle c, s \rangle \Downarrow_{s'}$

To prove Ξ^f , we still need to handle divergence...

How to do that?