

Compiler Correctness

KC Sivaramakrishnan

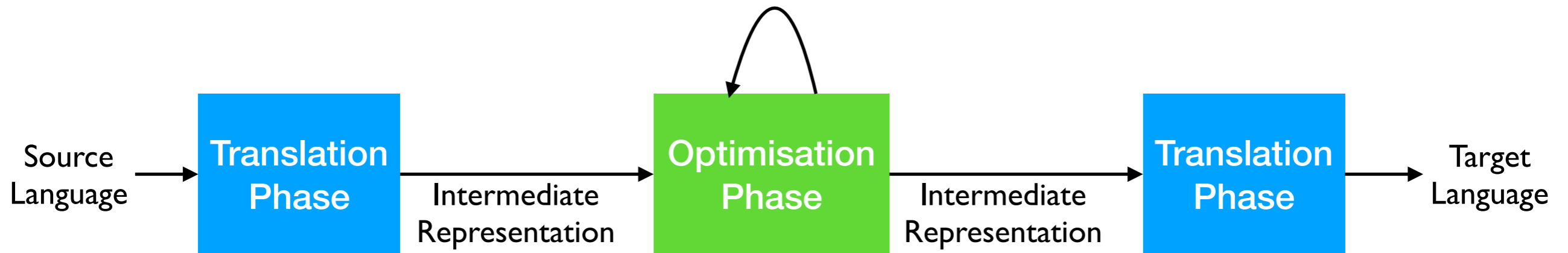
Spring 2021

IIT
MADRAS

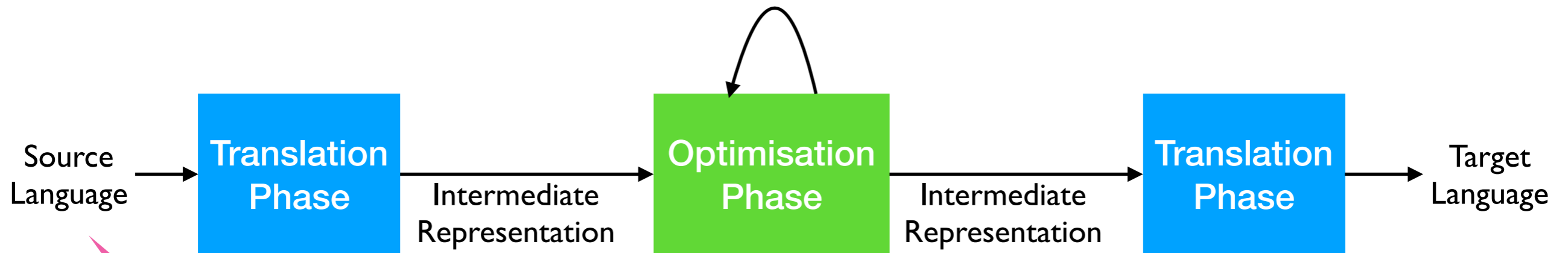


MADRAS IIT

Why prove compilers correct?

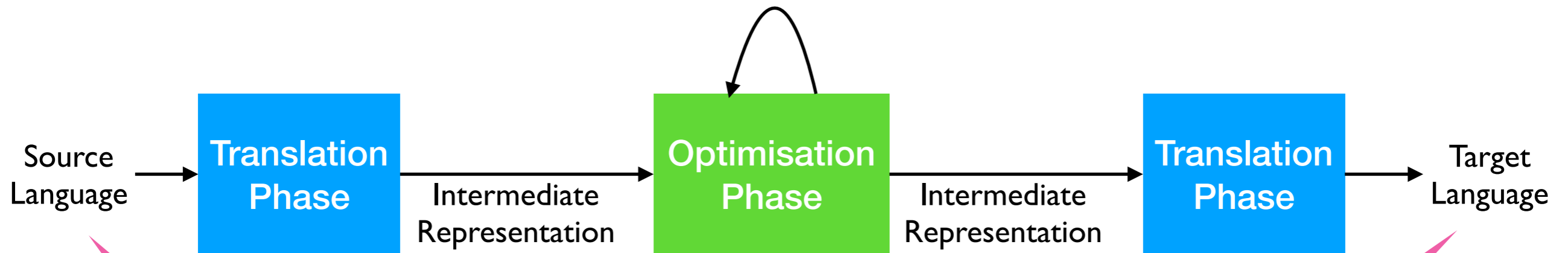


Why prove compilers correct?



I have some semantics for the source language

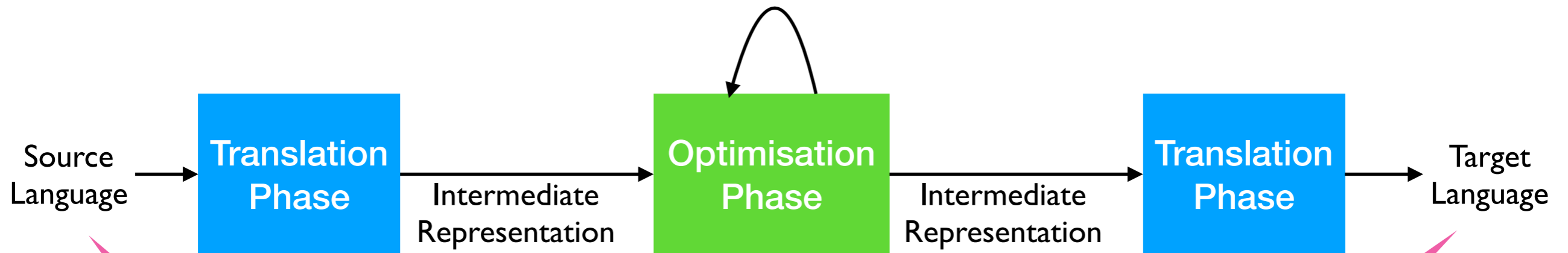
Why prove compilers correct?



I have some semantics for the source language

Does the target language have *equivalent* semantics?

Why prove compilers correct?



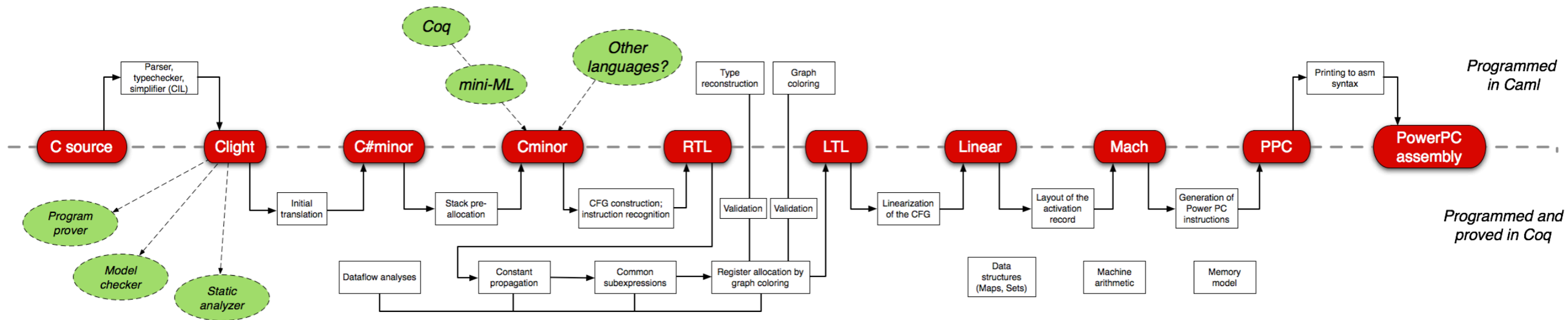
I have some semantics for the source language

Does the target language have *equivalent* semantics?

Idea: Prove each step of the translation and optimisation correct through *Simulation*

Compcert Verified C Compiler

<https://compcert.org/>



Language

Numbers	n	\in	\mathbb{N}
Variables	x	\in	Strings
Expressions	e	$::=$	$n \mid x \mid e + e \mid e - e \mid e \times e$
Commands	c	$::=$	$\text{skip} \mid x \leftarrow e \mid c; c \mid \text{if } e \text{ then } c \text{ else } c \mid \text{while } e \text{ do } c \mid \underline{\text{out}(e)}$

Language

Numbers $n \in \mathbb{N}$
Variables $x \in \text{Strings}$
Expressions $e ::= n \mid x \mid e + e \mid e - e \mid e \times e$
Commands $c ::= \text{skip} \mid x \leftarrow e \mid c; c \mid \text{if } e \text{ then } c \text{ else } c \mid \text{while } e \text{ do } c \mid \underline{\text{out}(e)}$

- An optimising compiler should preserve the behaviour of the program in terms of the output *traces* generated.

Language

Numbers	n	\in	\mathbb{N}
Variables	x	\in	Strings
Expressions	e	$::=$	$n \mid x \mid e + e \mid e - e \mid e \times e$
Commands	c	$::=$	$\text{skip} \mid x \leftarrow e \mid c; c \mid \text{if } e \text{ then } c \text{ else } c \mid \text{while } e \text{ do } c \mid \underline{\text{out}(e)}$

- An optimising compiler should preserve the behaviour of the program in terms of the output *traces* generated.
- Equivalence of traces is fundamental correctness property
 - ◆ Invariants — safety
 - ◆ Trace equivalence — liveness (not only for terminating programs)

Labelled transition semantics

$$\begin{array}{c}
 \hline
 (v, \text{out}(e)) \xrightarrow{0}^{\llbracket e \rrbracket v} (v, \text{skip}) \\
 \hline
 \hline
 (v, x \leftarrow e) \xrightarrow{0}^{\epsilon} (v[x \mapsto \llbracket e \rrbracket v], \text{skip}) \quad (v, \text{skip}; c_2) \xrightarrow{0}^{\epsilon} (v, c_2) \\
 \hline
 \begin{array}{cc}
 \frac{\llbracket e \rrbracket v \neq 0}{(v, \text{if } e \text{ then } c_1 \text{ else } c_2) \xrightarrow{0}^{\epsilon} (v, c_1)} & \frac{\llbracket e \rrbracket v = 0}{(v, \text{if } e \text{ then } c_1 \text{ else } c_2) \xrightarrow{0}^{\epsilon} (v, c_2)} \\
 \frac{\llbracket e \rrbracket v \neq 0}{(v, \text{while } e \text{ do } c_1) \xrightarrow{0}^{\epsilon} (v, c_1; \text{while } e \text{ do } c_1)} & \frac{\llbracket e \rrbracket v = 0}{(v, \text{while } e \text{ do } c_1) \xrightarrow{0}^{\epsilon} (v, \text{skip})}
 \end{array} \\
 \hline
 \frac{(v, c) \xrightarrow{0}^{\ell} (v', c')}{(v, C[c]) \xrightarrow{c}^{\ell} (v', C[c'])}
 \end{array}$$

Traces

- Finite sequences of outputs and termination events

$$\frac{}{\cdot \in \text{Tr}(s)} \quad \frac{}{\text{terminate} \in \text{Tr}((v, \text{skip}))} \quad \frac{s \xrightarrow{\epsilon}_c s' \quad t \in \text{Tr}(s')}{t \in \text{Tr}(s)} \quad \frac{s \xrightarrow{n}_c s' \quad t \in \text{Tr}(s')}{\text{out}(n) \bowtie t \in \text{Tr}(s)}$$

- A trace is allowed to end even if the program hasn't terminated

Traces

- Finite sequences of outputs and termination events

$$\frac{}{\cdot \in \text{Tr}(s)} \quad \frac{}{\text{terminate} \in \text{Tr}((v, \text{skip}))} \quad \frac{s \xrightarrow{\epsilon}_c s' \quad t \in \text{Tr}(s')}{t \in \text{Tr}(s)} \quad \frac{s \xrightarrow{n}_c s' \quad t \in \text{Tr}(s')}{\text{out}(n) \bowtie t \in \text{Tr}(s)}$$

- A trace is allowed to end even if the program hasn't terminated

DEFINITION 9.1 (Trace inclusion). For commands c_1 and c_2 , let $c_1 \leq c_2$ iff $\text{Tr}(c_1) \subseteq \text{Tr}(c_2)$.

DEFINITION 9.2 (Trace equivalence). For commands c_1 and c_2 , let $c_1 \simeq c_2$ iff $\text{Tr}(c_1) = \text{Tr}(c_2)$.

Traces

- Finite sequences of outputs and termination events

$$\frac{}{\cdot \in \text{Tr}(s)} \quad \frac{}{\text{terminate} \in \text{Tr}((v, \text{skip}))} \quad \frac{s \xrightarrow{\epsilon}_c s' \quad t \in \text{Tr}(s')}{t \in \text{Tr}(s)} \quad \frac{s \xrightarrow{n}_c s' \quad t \in \text{Tr}(s')}{\text{out}(n) \bowtie t \in \text{Tr}(s)}$$

- A trace is allowed to end even if the program hasn't terminated

DEFINITION 9.1 (Trace inclusion). For commands c_1 and c_2 , let $c_1 \preceq c_2$ iff $\text{Tr}(c_1) \subseteq \text{Tr}(c_2)$.

DEFINITION 9.2 (Trace equivalence). For commands c_1 and c_2 , let $c_1 \simeq c_2$ iff $\text{Tr}(c_1) = \text{Tr}(c_2)$.

$$\begin{aligned} & c_1 \simeq c_2 \\ \iff & \text{Tr}(c_1) = \text{Tr}(c_2) \\ \iff & \text{Tr}(c_1) \subseteq \text{Tr}(c_2) \wedge \text{Tr}(c_2) \subseteq \text{Tr}(c_1) \\ \iff & c_1 \preceq c_2 \wedge c_2 \preceq c_1 \end{aligned}$$

Constant-folding

- Optimisation is
 - ◆ Find all maximal program subexpressions that don't contain variables
 - ◆ Replace each subexpression with its known constant value
- The optimisation only changes variable free-expressions
 - ◆ The original and the optimised program match at each *step*

Basic Simulation Relation

DEFINITION 9.3 (Simulation relation). We say that binary relation R over states of our object language is a *simulation relation* iff:

- (1) Whenever $(v_1, \text{skip}) R (v_2, c_2)$, it follows that $c_2 = \text{skip}$.
- (2) Whenever $s_1 R s_2$ and $s_1 \xrightarrow{\ell}_c s'_1$, there exists s'_2 such that $s_2 \xrightarrow{\ell}_c s'_2$ and $s'_1 R s'_2$.

$$\begin{array}{ccc} s_1 & \xrightarrow{R} & s_2 \\ \downarrow \forall \xrightarrow{\ell}_c & & \downarrow \exists \xrightarrow{\ell}_c \\ s'_1 & \xleftarrow{R^{-1}} & s'_2 \end{array}$$

THEOREM 9.4. *If there exists a simulation R such that $s_1 R s_2$, then $s_1 \simeq s_2$.*

Simulation for Constant Folding

THEOREM 9.5. *For any v and c , $(v, c) \simeq (v, \text{cfold}_1(c))$.*

PROOF. By a simulation argument using this relation:

$$(v_1, c_1) R (v_2, c_2) = v_1 = v_2 \wedge c_2 = \text{cfold}_1(c_1)$$

Basic Simulation Relation

THEOREM 9.4. *If there exists a simulation R such that $s_1 R s_2$, then $s_1 \simeq s_2$.*

$$\begin{array}{ccc} s_1 & \xrightarrow{R} & s_2 \\ \downarrow \forall \xrightarrow{\ell}_c & & \downarrow \exists \xrightarrow{\ell}_c \\ s'_1 & \xleftarrow{R^{-1}} & s'_2 \end{array}$$

- We prove two trace inclusion directions separately
- Left-to-right $\text{Tr}(s_1) \subseteq \text{Tr}(s_2)$ proved by induction on traces on the left
- Right-to-left $\text{Tr}(s_2) \subseteq \text{Tr}(s_1)$ proved similarly
 - ♦ Depends on operational semantics being total and deterministic

Simulation with skipping

- Consider extension of constant folding to conditionals

$$\frac{\text{cfoldArith}(e_1) = 1}{\text{cfold}(\text{if } e_1 \text{ then } e_2 \text{ else } e_3) = \text{cfold}(e_2)}$$

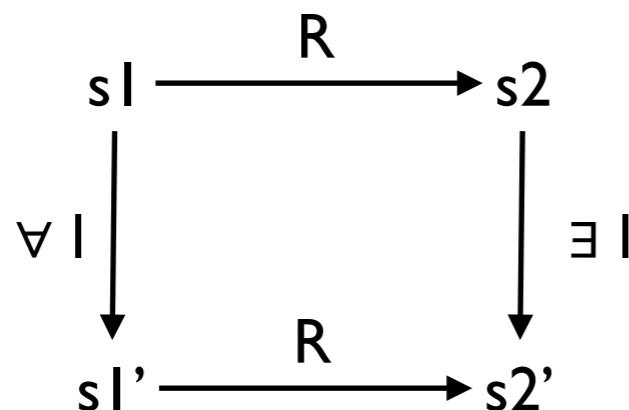
$$\frac{\text{cfoldArith}(e_1) = 0}{\text{cfold}(\text{if } e_1 \text{ then } e_2 \text{ else } e_3) = \text{cfold}(e_3)}$$

- We can no longer use our basic simulation
 - ✦ Steps are intentionally skipped in the optimised program

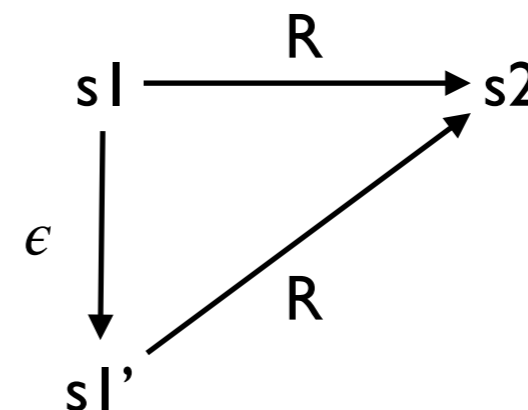
Simulation with Skipping: Faulty

DEFINITION 9.6 (Simulation relation with skipping (*faulty* version!)). We say that binary relation R over states of our object language is a *simulation relation with skipping* iff:

- (1) Whenever $(v_1, \text{skip}) R (v_2, c_2)$, it follows that $c_2 = \text{skip}$.
- (2) Whenever $s_1 R s_2$ and $s_1 \xrightarrow{\ell}_c s'_1$, then either:
 - (a) there exists s'_2 such that $s_2 \xrightarrow{\ell}_c s'_2$ and $s'_1 R s'_2$,
 - (b) or $\ell = \epsilon$ and $s'_1 R s_2$.



(a)



(b)

Faulty Optimisation: Advertising

$\text{withAds}(\text{while } 1 \text{ do Skip}) = \text{while } 1 \text{ do Out}(0)$
 $\text{withAds}(c) = c$

Faulty Optimisation: Advertising

$$\begin{aligned} \text{withAds}(\text{while } 1 \text{ do Skip}) &= \text{while } 1 \text{ do Out}(0) \\ \text{withAds}(c) &= c \end{aligned}$$

- We can use the candidate simulation relation

Faulty Optimisation: Advertising

$$\begin{aligned} \text{withAds}(\text{while } 1 \text{ do Skip}) &= \text{while } 1 \text{ do Out}(0) \\ \text{withAds}(c) &= c \end{aligned}$$

- We can use the candidate simulation relation

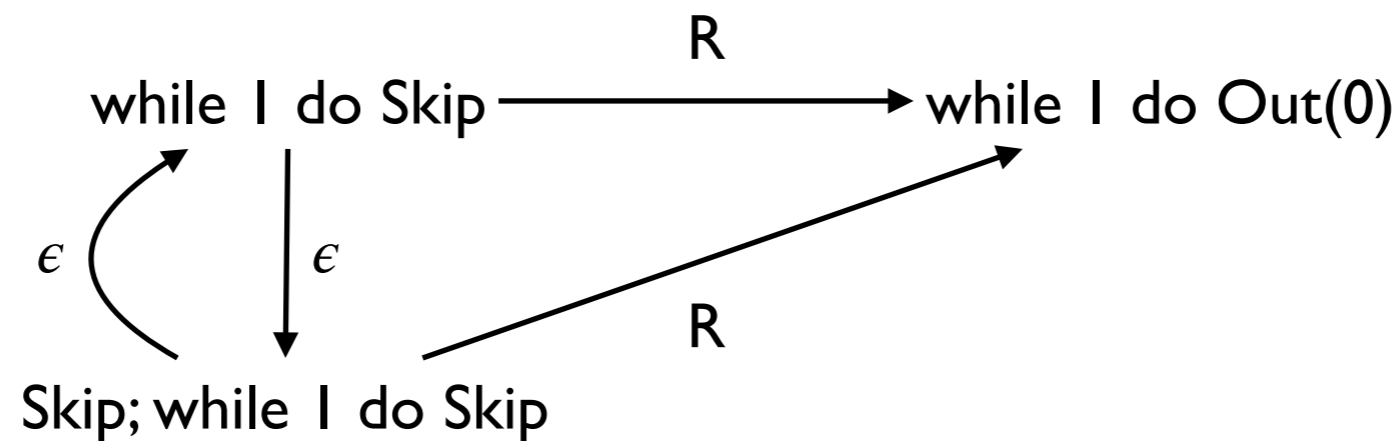
$$(v_1, c_1) R (v_2, c_2) = c_1 \in \{\text{while } 1 \text{ do skip}, (\text{skip}; \text{while } 1 \text{ do skip})\}$$

Faulty Optimisation: Advertising

$$\begin{aligned} \text{withAds}(\text{while } 1 \text{ do Skip}) &= \text{while } 1 \text{ do Out}(0) \\ \text{withAds}(c) &= c \end{aligned}$$

- We can use the candidate simulation relation

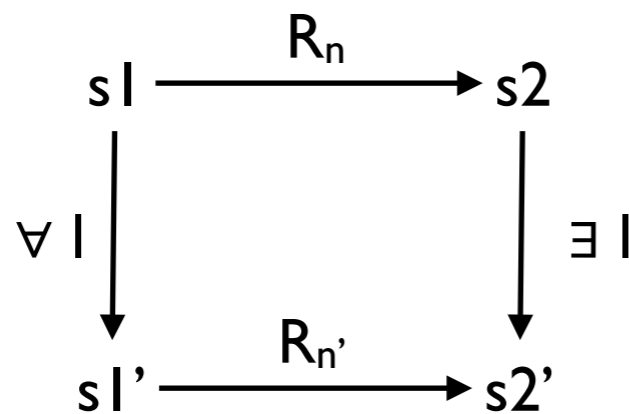
$$(v_1, c_1) R (v_2, c_2) = c_1 \in \{\text{while } 1 \text{ do skip}, (\text{skip}; \text{while } 1 \text{ do skip})\}$$



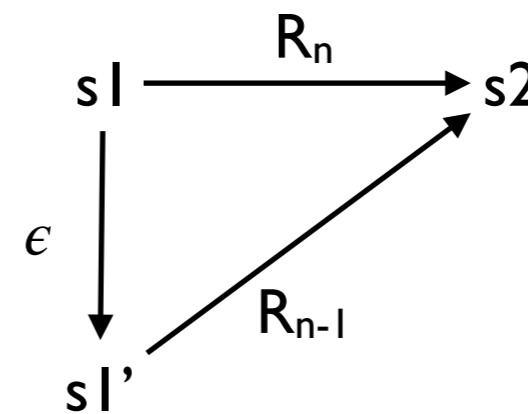
Simulation with skipping

DEFINITION 9.7 (Simulation relation with skipping). We say that an \mathbb{N} -indexed family of binary relations R_n over states of our object language is a *simulation relation with skipping* iff:

- (1) Whenever $(v_1, \text{skip}) R_n (v_2, c_2)$, it follows that $c_2 = \text{skip}$.
- (2) Whenever $s_1 R_n s_2$ and $s_1 \xrightarrow{\ell}_c s'_1$, then either:
 - (a) there exist n' and s'_2 such that $s_2 \xrightarrow{\ell}_c s'_2$ and $s'_1 R_{n'} s'_2$,
 - (b) or $n > 0$, $\ell = \epsilon$, and $s'_1 R_{n-1} s_2$.



(a)



(b)

Trace Equivalence for Simulation with Skipping

THEOREM 9.8. *If there exists a simulation with skipping R such that $s_1 R_n s_2$, then $s_1 \simeq s_2$.*

THEOREM 9.9. *For any v and c , $(v, c) \simeq (v, \text{cfold}_2(c))$.*

PROOF. By a simulation argument (with skipping) using this relation:

$$(v_1, c_1) R_n (v_2, c_2) = v_1 = v_2 \wedge c_2 = \text{cfold}_2(c_1) \wedge \text{countlfs}(c_1) < n$$