

# Slicing Modern Program Structures While Eliminating Irrelevant Loops

Torben Amtoft

Kansas State University

December 2009

Introduction

Standard Definitions

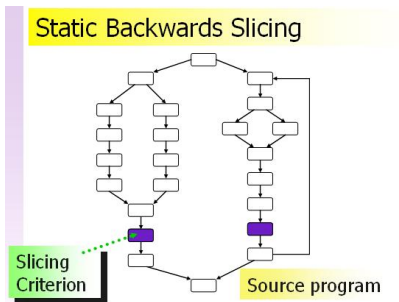
Meeting Challenge 1

Meeting Challenge 2

Technical Results

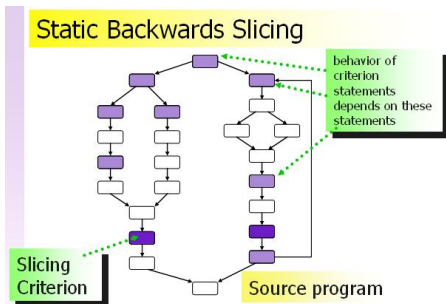
Conclusion

# What is Slicing?



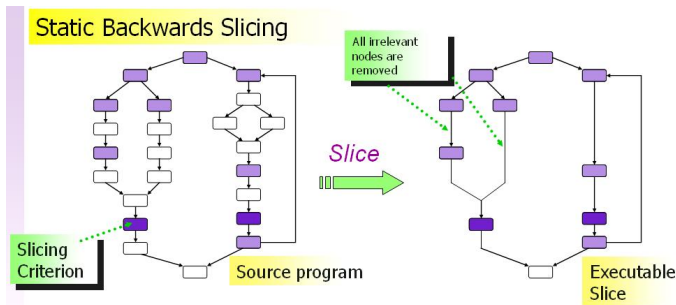
Pick one or more program points of interest  
(called the **slicing criterion**)

# What is Slicing?



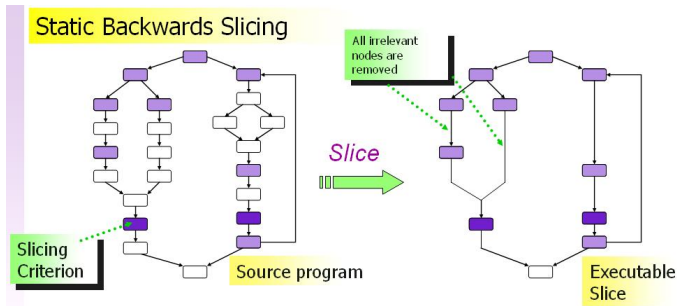
Walk backwards to find nodes that  $C$  depend on (called the *slice set*)

# What is Slicing?



Remove **irrelevant** nodes  
(or replace them by dummy nodes)

# What is Slicing?



Applications include

- ▶ compiler optimizations
- ▶ debugging
- ▶ model checking
- ▶ protocol understanding

# Challenge

A definition of **control dependency** that

1. is applicable to **arbitrary CFGs**

# Challenge

A definition of **control dependency** that

1. is applicable to **arbitrary CFGs**, including those that
  - ▶ are irreducible (as for state charts)

A definition of **control dependency** that

1. is applicable to **arbitrary CFGs**, including those that
  - ▶ are irreducible (as for state charts)
  - ▶ have **no end node** (as for reactive systems)



# Challenge

A definition of **control dependency** that

1. is applicable to **arbitrary CFGs**, including those that
  - ▶ are irreducible (as for state charts)
  - ▶ have **no end node** (as for reactive systems)
2. produces slices that are likely to be **manageable**

# Challenge

A definition of **control dependency** that

1. is applicable to **arbitrary CFGs**, including those that
  - ▶ are irreducible (as for state charts)
  - ▶ have **no end node** (as for reactive systems)
2. produces slices that are likely to be **manageable**

Our **proposed solution** meets these challenges and is correct in the sense that

- ▶ the observable behavior of source program is **prefix** of observable behavior of sliced program.

# Challenge

A definition of **control dependency** that

1. is applicable to **arbitrary CFGs**, including those that
  - ▶ are irreducible (as for state charts)
  - ▶ have **no end node** (as for reactive systems)
2. produces slices that are likely to be **manageable**

Our **proposed solution** meets these challenges and is correct in the sense that

- ▶ the observable behavior of source program is **prefix** of observable behavior of sliced program.
- ▶ this is expressed using **simulation**

A definition of **control dependency** that

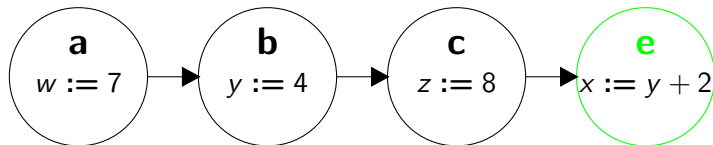
1. is applicable to **arbitrary CFGs**, including those that
  - ▶ are irreducible (as for state charts)
  - ▶ have **no end node** (as for reactive systems)
2. produces slices that are likely to be **manageable**

Our **proposed solution** meets these challenges and is correct in the sense that

- ▶ the observable behavior of source program is **prefix** of observable behavior of sliced program.
- ▶ this is expressed using **simulation**
- ▶ correctness proof easily fits conference page limits

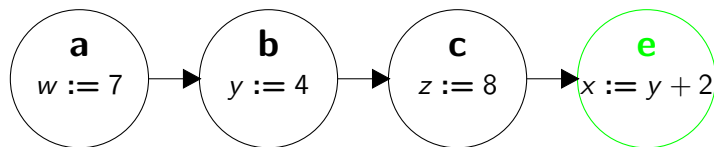
# Data Dependence

Source: final value of  $x = 6$



# Data Dependence

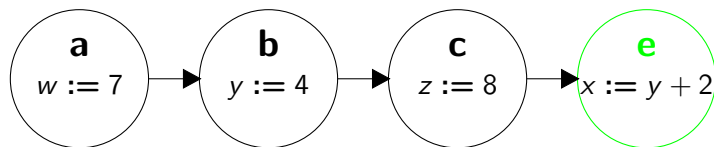
Source: final value of  $x = 6$



$e$  is **data dependent** on  $b$ , written  $b \xrightarrow{dd} e$ , since  $y$  is  
▶ used in  $e$ , defined in  $b$ , not defined in between

# Data Dependence

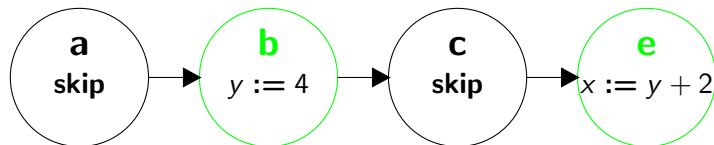
Source: final value of  $x = 6$



$e$  is **data dependent** on  $b$ , written  $b \xrightarrow{dd} e$ , since  $y$  is

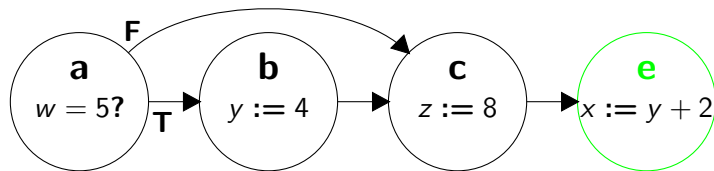
- ▶ used in  $e$ , defined in  $b$ , not defined in between

Slice: final value of  $x = 6$



# Control Dependence (CD)

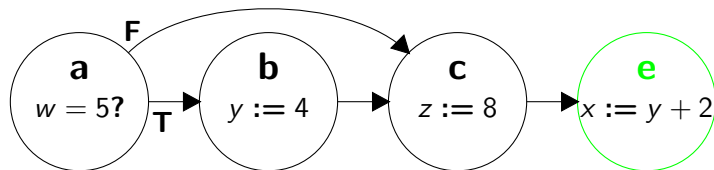
Source: final value of  $x = 6$  if  $w = 5$ , else  $2$





# Control Dependence (CD)

Source: final value of  $x = 6$  if  $w = 5$ , else  $2$

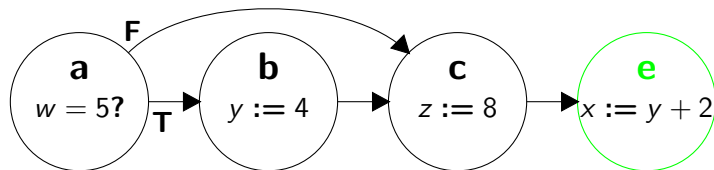


$b \xrightarrow{dd} e$  and  $b$  is **control dependent** on  $a$ ,  $a \xrightarrow{cd} b$ , since

- ▶  $a$  is not strictly **postdominated** by  $b$ : path  $[a..e] \not\equiv b$
- ▶ in  $]a..b[$ , all nodes postdominated by  $b$  (vacuously)

# Control Dependence (CD)

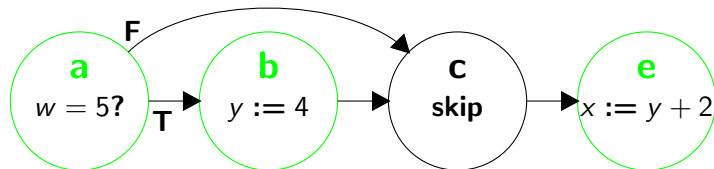
Source: final value of  $x = 6$  if  $w = 5$ , else  $2$



$b \stackrel{dd}{\rightarrow} e$  and  $b$  is **control dependent** on  $a$ ,  $a \stackrel{cd}{\rightarrow} b$ , since

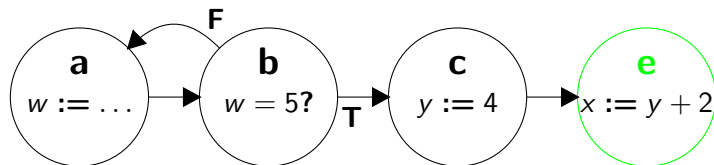
- ▶  $a$  is not strictly **postdominated** by  $b$ : path  $[a..e] \not\supseteq b$
- ▶ in  $]a..b[$ , all nodes postdominated by  $b$  (vacuously)

Slice: final value of  $x = 6$  if  $w = 5$ , else  $2$



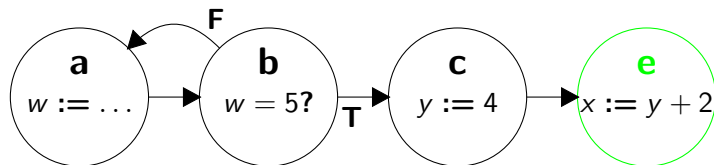
# Weak Control Dependence (WCD)

Source: final value of  $x = 6$  or  $\perp$



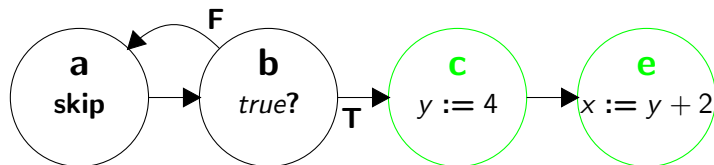
# Weak Control Dependence (WCD)

Source: final value of  $x = 6$  or  $\perp$



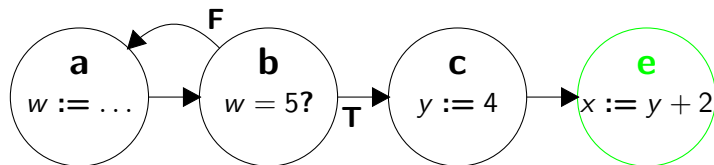
$c \xrightarrow{dd} e$  but  $b \not\xrightarrow{cd} c$ , since  $c$  strictly postdominates  $b$

Slice: final value of  $x = 6$



# Weak Control Dependence (WCD)

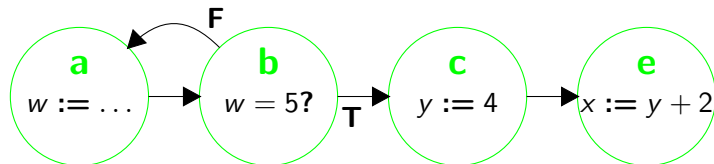
Source: final value of  $x = 6$  or  $\perp$



$c \xrightarrow{dd} e$  but  $b \not\xrightarrow{cd} c$ , since  $c$  strictly postdominates  $b$

But  $c$  is **weakly control dependent** on  $b$ ,  $b \xrightarrow{wcd} c$  (though no  $[b..e]$  bypasses  $c$ ,  $b$  may cause indefinite delay of  $c$ )

Slice: final value of  $x = 6$  or  $\perp$



# Removing End Node Assumption

- ▶ The definitions so far assume the existence of a (unique) **end node**.
- ▶ recall **Challenge 1**: extend the definition of control dependency to be applicable even to CFGs that **may not** have an end node

# Removing End Node Assumption

- ▶ The definitions so far assume the existence of a (unique) **end node**.
- ▶ recall **Challenge 1**: extend the definition of control dependency to be applicable even to CFGs that **may not** have an end node

Ranganath et al [ESOP'05 & TOPLAS'07] proposed a new definition, saying that  $h$  is control dependent on  $a$  iff

- ▶ from one of  $a$ 's successors,  $h$  **cannot** be avoided forever: all **maximal** paths go through  $h$ .
- ▶ from another of  $a$ 's successors,  $h$  **may** be avoided forever: there exists a maximal path not going through  $h$ .

# Removing End Node Assumption

- ▶ The definitions so far assume the existence of a (unique) **end node**.
- ▶ recall **Challenge 1**: extend the definition of control dependency to be applicable even to CFGs that **may not** have an end node

Ranganath et al [ESOP'05 & TOPLAS'07] proposed a new definition, saying that  $h$  is control dependent on  $a$  iff

- ▶ from one of  $a$ 's successors,  $h$  **cannot** be avoided forever: all **maximal** paths go through  $h$ .
- ▶ from another of  $a$ 's successors,  $h$  **may** be avoided forever: there exists a maximal path not going through  $h$ .

This is called **NTSCD** (non-termination **sensitive** CD)

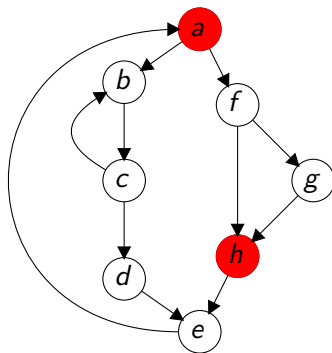
- ▶ **Conservative Extension**: in the **special case** where the CFG has an end node, NTSCD is **equivalent** to Podgurski & Clarke's weak control dependence.



# Illustrating NTSCD

Torben Amtoft

*h* NTSCD on *a*



Introduction

Standard Definitions

Meeting Challenge 1

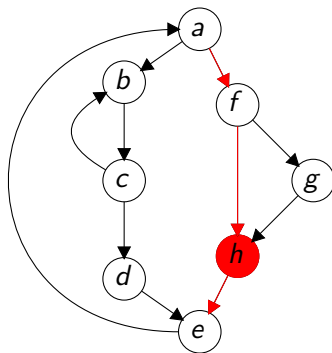
Meeting Challenge 2

Technical Results

Conclusion

# Illustrating NTSCD

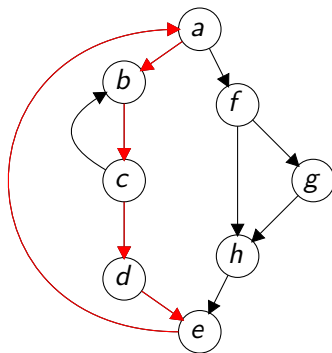
$h$  NTSCD on  $a$



If  $a$  selects  $f$  then  $h$  cannot be avoided

# Illustrating NTSCD

$h$  NTSCD on  $a$

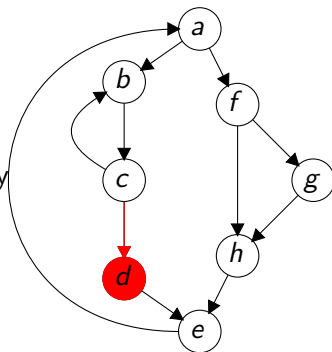


If  $a$  selects  $b$  then  $h$  can be avoided forever

# Illustrating NTSCD

## Termination is Preserved

$d$  NTSD on  $c$   
so inner loop **not** sliced away

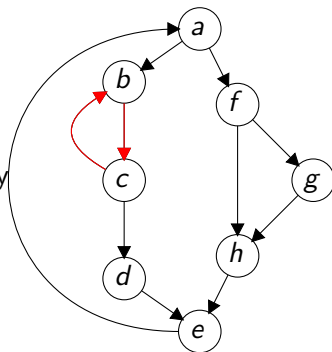


If  $c$  selects  $d$  then  $d$  cannot be avoided

# Illustrating NTSCD

## Termination is Preserved

$d$  NTSD on  $c$   
so inner loop **not** sliced away



If  $c$  selects  $b$  then  $d$  can be avoided forever

We aim at (weak) **bisimulation**:

- ▶ if a node can do an **observable** action in the source program, it can do so in the sliced program
- ▶ if a node can do an observable action in the sliced program, it can do so in the source program
- ▶ Also, for stores to be bisimilar, they must agree on “relevant” variables.

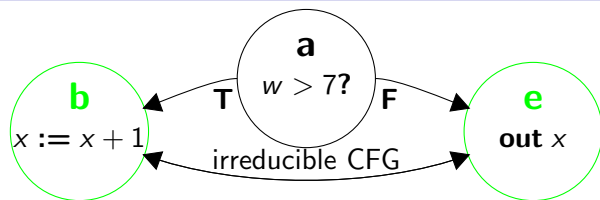
We aim at (weak) **bisimulation**:

- ▶ if a node can do an **observable** action in the source program, it can do so in the sliced program
- ▶ if a node can do an observable action in the sliced program, it can do so in the source program
- ▶ Also, for stores to be bisimilar, they must agree on “relevant” variables.

If the CFG is **reducible**, we can prove that slicing based on NTSCD gives rise to a bisimulation. Recall that a CFG is reducible if its edges can be split into

- ▶ forward edges which form a DAG;
- ▶ back edges where the target dominates the source

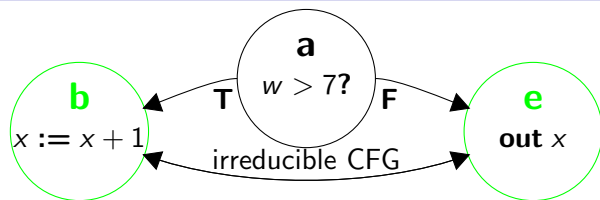
# Irreducible Graphs



Observable behavior:  $1, 2, 3, \dots$  if  $w > 7$ , else  $0, 1, 2, \dots$



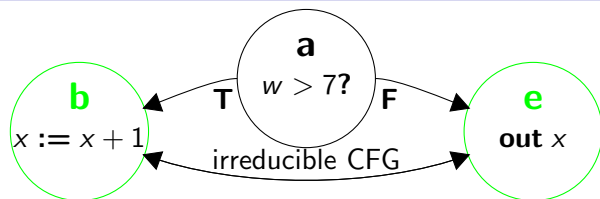
# Irreducible Graphs



Observable behavior:  $1, 2, 3, \dots$  if  $w > 7$ , else  $0, 1, 2, \dots$

- ▶ Thus  $a$  should be in slice set, but hard to see how  $a \xrightarrow{?cd} b$  or  $a \xrightarrow{?cd} e$  since all paths from  $a$  hit  $b$  and  $e$ .

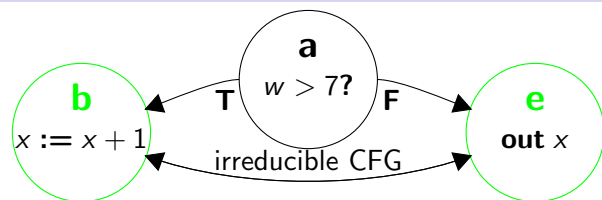
# Irreducible Graphs



Observable behavior:  $1, 2, 3, \dots$  if  $w > 7$ , else  $0, 1, 2, \dots$

- ▶ Thus  $a$  should be in slice set, but hard to see how  $a \xrightarrow{?cd} b$  or  $a \xrightarrow{?cd} e$  since all paths from  $a$  hit  $b$  and  $e$ .
- ▶ Instead we need a **ternary** relation, expressing that the **order** of  $b, c$  depends on  $a$ .

# Irreducible Graphs



Observable behavior:  $1, 2, 3, \dots$  if  $w > 7$ , else  $0, 1, 2, \dots$

- ▶ Thus  $a$  should be in slice set, but hard to see how  $a \xrightarrow{?cd} b$  or  $a \xrightarrow{?cd} e$  since all paths from  $a$  hit  $b$  and  $e$ .
- ▶ Instead we need a **ternary** relation, expressing that the **order** of  $b, c$  depends on  $a$ .

**Theorem** For an **arbitrary** CFG (with/without end nodes, reducible or not), the sliced program is **bisimilar** to the source program, provided the slice set is **closed** under

- ▶ data dependency
- ▶ NTSCD
- ▶ a certain kind of “**order dependency**”

The approach by Ranganath et al

- ▶ provides first theoretical foundation for slicing of CFGs that may have zero end nodes
- ▶ generalizes **weak** control dependence, and hence provides for termination-preserving slicing

The approach by Ranganath et al

- ▶ provides first theoretical foundation for slicing of CFGs that may have zero end nodes
- ▶ generalizes **weak** control dependence, and hence provides for termination-preserving slicing

Slices are typically larger, as they must include all nodes that influence guards of potential loops

- ▶ great, if you are slicing to preserve liveness properties for model checking
- ▶ not so great, if slicing for program understanding

The approach by Ranganath et al

- ▶ provides first theoretical foundation for slicing of CFGs that may have zero end nodes
- ▶ generalizes **weak** control dependence, and hence provides for termination-preserving slicing

Slices are typically larger, as they must include all nodes that influence guards of potential loops

- ▶ great, if you are slicing to preserve liveness properties for model checking
- ▶ not so great, if slicing for program understanding

We may like to discard nodes that may influence loops if they do not impact **observable** computation

- ▶ thus we should generalize  $\xrightarrow{cd}$  rather than  $\xrightarrow{wcd}$
- ▶ as was done by something called NTICD, but which was not helpful to establish the correctness of slicing

# Weak Order Dependence (WOD)

Recall [Challenge 2](#): find a dependence relation that

- ▶ slices away irrelevant loops
- ▶ does not assume the existence of an end node

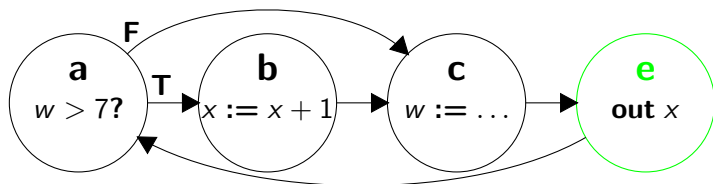






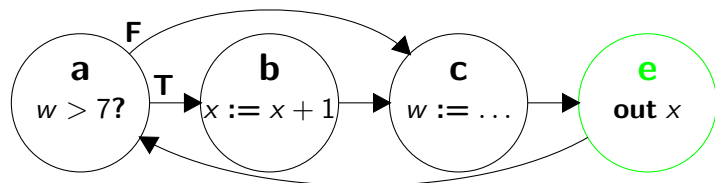
# Example: WOD for Conditionals

Source: observable behavior of the form  
1,1,2,2,2,3,4,4,...

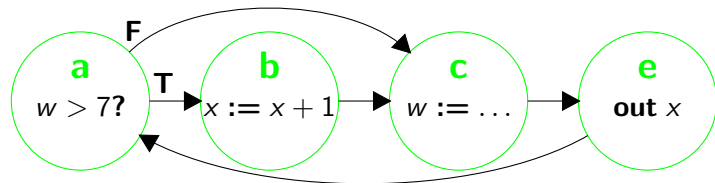


# Example: WOD for Conditionals

Source: observable behavior of the form  
1,1,2,2,2,3,4,4,...

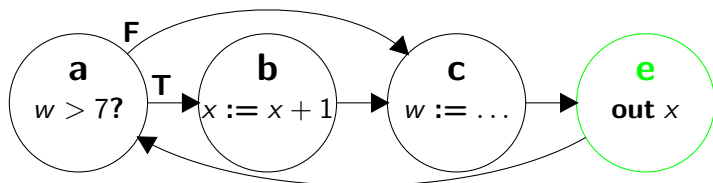


$b \xrightarrow{dd} e$  and  $c \xrightarrow{dd} a \xrightarrow{wod} b, e$



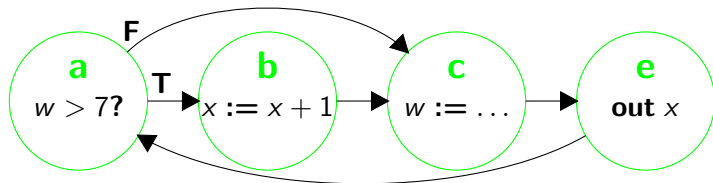
# Example: WOD for Conditionals

Source: observable behavior of the form  
 $1, 1, 2, 2, 2, 3, 4, 4, \dots$



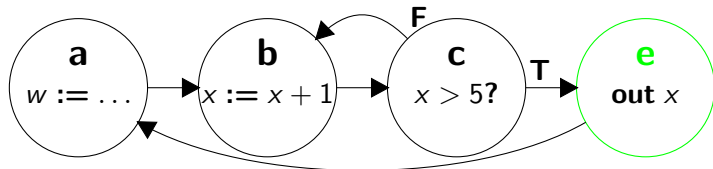
$b \xrightarrow{dd} e$  and  $c \xrightarrow{dd} a \xrightarrow{wod} b, e$  since

- ▶ path  $[a..b] \not\rightarrow e$  and path  $[a..e] \not\rightarrow b$
- ▶  $a$  has successor  $b$  with  $[b..b]$  and all  $[b..e]$  contain  $b$



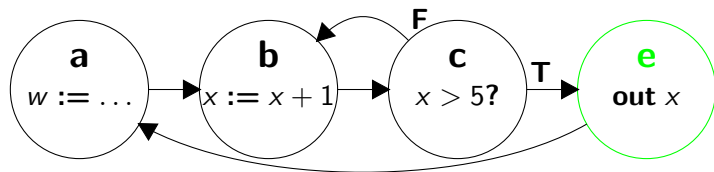
# Example: WOD for Relevant Loops

Source: observable behavior 6,7,8,9,10,...



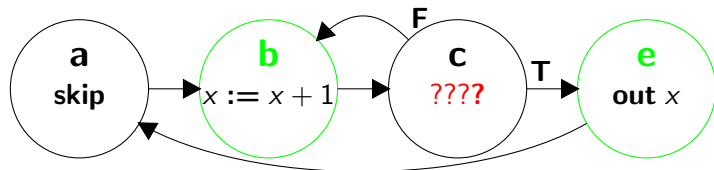
# Example: WOD for Relevant Loops

Source: observable behavior 6,7,8,9,10,...



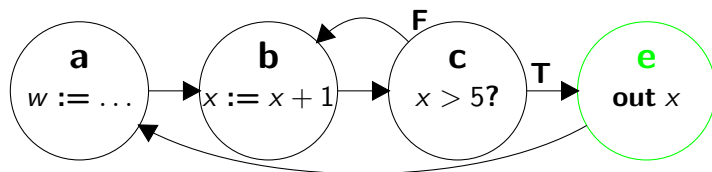
Slice: observable behavior if guard **wrongly** sliced away:

- ▶ 1,2,3,4,5,... if replaced by *true*
- ▶  $\epsilon$  if replaced by *false*



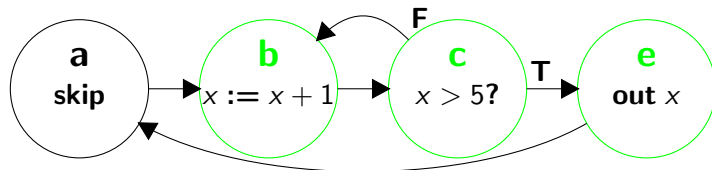
# Example: WOD for Relevant Loops

Source: observable behavior 6,7,8,9,10,...



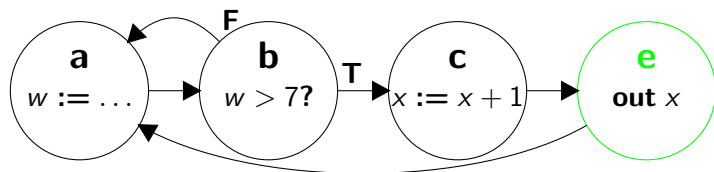
$c \xrightarrow{wod} b, e$  since

- ▶ path  $[c..b] \not\equiv e$  and path  $[c..e] \not\equiv b$
- ▶  $c$  has successor  $b$  with  $[b..b]$  and all  $[b..e]$  contain  $b$



# Example: WOD for Irrelevant Loops

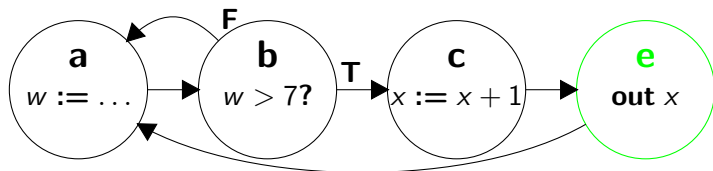
Source: observable behavior **1,2,3,4,...** or **prefix** thereof





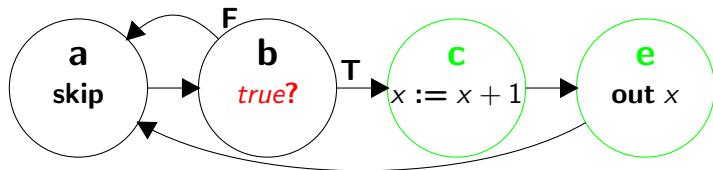
# Example: WOD for Irrelevant Loops

Source: observable behavior **1,2,3,4,...** or **prefix** thereof



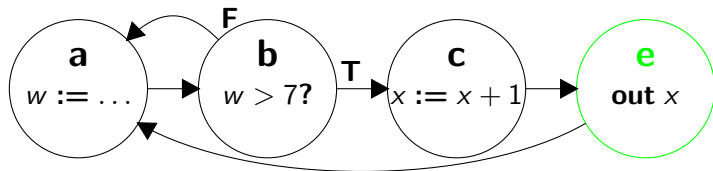
New guard chosen to get **closer** to **next observable**.

Slice: observable behavior **1,2,3,4,...**



# Example: WOD for Irrelevant Loops

Source: observable behavior  $1,2,3,4,\dots$  or prefix thereof

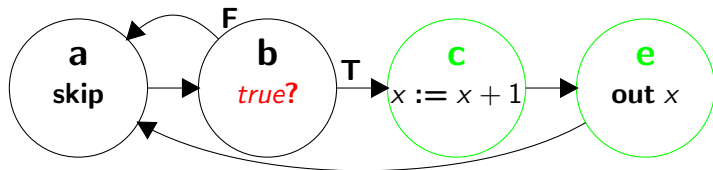


$b \stackrel{wod}{\not\rightarrow} c, e$  since

- ▶ no path  $[b..e] \not\equiv c$

New guard chosen to get closer to next observable.

Slice: observable behavior  $1,2,3,4,\dots$



# Relating WOD to CD

**Conservative extension** For a CFG with an end node which is part of slicing criterion, the **closures** of  $\xrightarrow{wod}$  and  $\xrightarrow{cd}$  **coincide**.

# Relating WOD to CD

**Conservative extension** For a CFG with an end node which is part of slicing criterion, the **closures** of  $\xrightarrow{wod}$  and  $\xrightarrow{cd}$  **coincide**.

**Alternative characterization** Current work with Mark Harman et al from King's College captures the intuition that we need a ternary relation because of one the three nodes must play the role as “local pseudo exit node”:

# Relating WOD to CD

**Conservative extension** For a CFG with an end node which is part of slicing criterion, the **closures** of  $\xrightarrow{wod}$  and  $\xrightarrow{cd}$  **coincide**.

**Alternative characterization** Current work with Mark Harman et al from King's College captures the intuition that we need a ternary relation because of one the three nodes must play the role as “local pseudo exit node”:

- ▶ Assume a CFG where all nodes are **reachable from each other** (as is often the case)
- ▶ for each node  $x$ , if we remove all outgoing edges from  $x$ , we get a graph  $G_{\rightarrow x}$  having  $x$  as end node.

# Relating WOD to CD

**Conservative extension** For a CFG with an end node which is part of slicing criterion, the **closures** of  $\xrightarrow{wod}$  and  $\xrightarrow{cd}$  **coincide**.

**Alternative characterization** Current work with Mark Harman et al from King's College captures the intuition that we need a ternary relation because of one the three nodes must play the role as “local pseudo exit node”:

- ▶ Assume a CFG where all nodes are **reachable from each other** (as is often the case)
- ▶ for each node  $x$ , if we remove all outgoing edges from  $x$ , we get a graph  $G_{\rightarrow x}$  having  $x$  as end node.
- ▶ If  $a \xrightarrow{wod} b, c$  in  $G$  then either
  - ▶  $a \xrightarrow{cd} b$  in  $G_{\rightarrow c}$ , or
  - ▶  $a \xrightarrow{cd} c$  in  $G_{\rightarrow b}$ .

# Relating WOD to CD

**Conservative extension** For a CFG with an end node which is part of slicing criterion, the **closures** of  $\xrightarrow{wod}$  and  $\xrightarrow{cd}$  **coincide**.

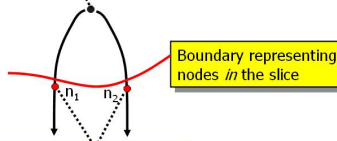
**Alternative characterization** Current work with Mark Harman et al from King's College captures the intuition that we need a ternary relation because of one the three nodes must play the role as "local pseudo exit node":

- ▶ Assume a CFG where all nodes are **reachable from each other** (as is often the case)
- ▶ for each node  $x$ , if we remove all outgoing edges from  $x$ , we get a graph  $G_{\rightarrow x}$  having  $x$  as end node.
- ▶ If  $a \xrightarrow{wod} b, c$  in  $G$  then either
  - ▶  $a \xrightarrow{cd} b$  in  $G_{\rightarrow c}$ , or
  - ▶  $a \xrightarrow{cd} c$  in  $G_{\rightarrow b}$ .
- ▶ if  $a \xrightarrow{cd} c$  in  $G_{\rightarrow b}$  with  $a \neq c$  then  $a \xrightarrow{wod} b, c$  in  $G$ .

# Unique Next Observable

- ▶ For the correctness proof, we shall assume the set of observables is **closed under  $\xrightarrow{dd}$  and  $\xrightarrow{wod}$** .
- ▶ Key advantage of being closed under  $\xrightarrow{wod}$ :  
**at most one next observable.**

A conditional that  
is not in the slice...

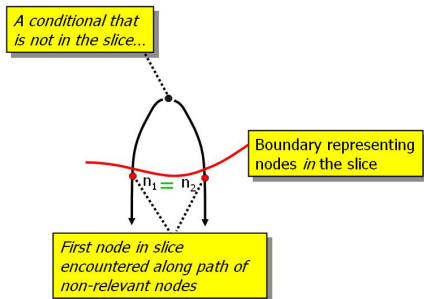


First node in slice  
encountered along path of  
non-relevant nodes



# Unique Next Observable

- ▶ For the correctness proof, we shall assume the set of observables is **closed under  $\xrightarrow{dd}$  and  $\xrightarrow{wod}$** .
- ▶ Key advantage of being closed under  $\xrightarrow{wod}$ :  
**at most one next observable.**



# Unique Next Observable

- ▶ For the correctness proof, we shall assume the set of observables is **closed under  $\xrightarrow{dd}$  and  $\xrightarrow{wod}$** .
- ▶ Key advantage of being closed under  $\xrightarrow{wod}$ : **at most one next observable**.

Recent classification [Sebastian Danicic]: **all** slices in the literature computes one of the two below:

- ▶ the least set where all nodes have **at most one** next observable.

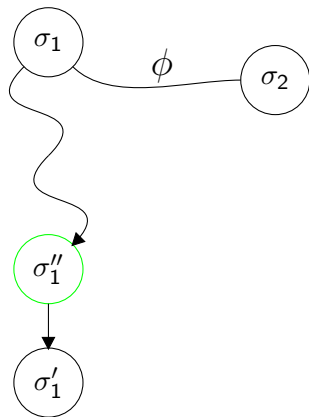
This is **weak commitment**, as computed by  $\xrightarrow{cd}$ , and for general CFGs by  $\xrightarrow{wod}$ .

- ▶ the least set where all nodes have at most one next observable, **and** satisfy that all nodes either
  - ▶ eventually will reach an observable, or
  - ▶ always avoids an observable

This is **strong commitment**, as computed by  $\xrightarrow{wcd}$ , and for general CFGs by  $\xrightarrow{ntscd} + \text{„dod“}$

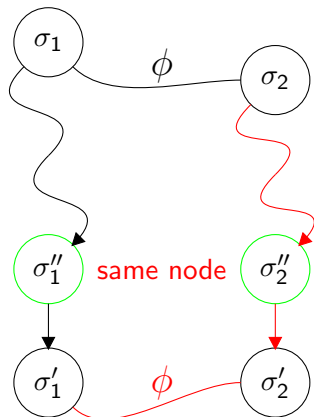
# Weak Simulation

- ▶ Slicing Correctness expressed using **weak simulation**
- ▶ A **weak simulation** relates **source states** to **slice states**
- ▶ A **program state** is an **node** and a **store**



# Weak Simulation

- ▶ Slicing Correctness expressed using **weak simulation**
- ▶ A **weak simulation** relates **source states** to **slice states**
- ▶ A **program state** is an **node** and a **store**



# Relevant Variables

A variable  $x$  is **relevant** at  $n$  if

- ▶  $x$  is **used** at an **observable** node  $q$
- ▶ there exists a path  $[n..q[$  along which  $x$  is **not defined**

A variable  $x$  is **relevant** at  $n$  if

- ▶  $x$  is **used** at an **observable** node  $q$
- ▶ there exists a path  $[n..q[$  along which  $x$  is **not defined**

Facts:

- ▶ this definition is **independent** of whether one considers the source or the slice
- ▶ if two nodes have the **same next observable**, they have the **same relevant variables**

# Defining a Concrete Weak Simulation

We define  $(n_1, \sigma_1) R (n_2, \sigma_2)$  if

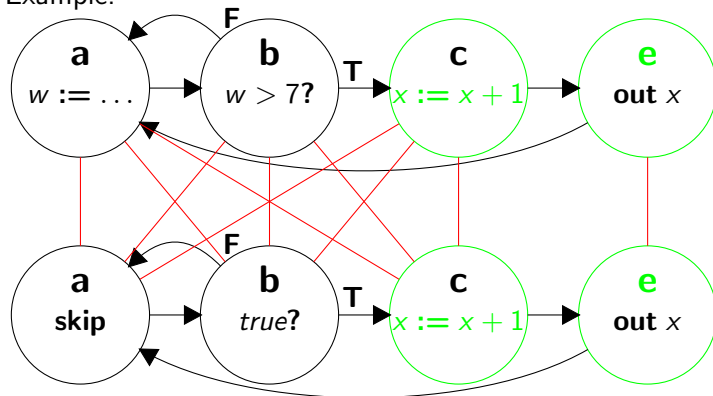
- ▶  $n_1$  and  $n_2$  have the **same next observable**
- ▶  $\sigma_1$  and  $\sigma_2$  **agree on the relevant variables**

# Defining a Concrete Weak Simulation

We define  $(n_1, \sigma_1) R (n_2, \sigma_2)$  if

- ▶  $n_1$  and  $n_2$  have the **same next observable**
- ▶  $\sigma_1$  and  $\sigma_2$  **agree on the relevant variables**

Example:





## Theorem

- ▶  $R$  is a weak simulation.

Thus any observable action by the source can be matched by the slice, but **not** necessarily vice versa.

## Ball & Horwitz, 1993

- ▶ assumes unique end node
- ▶ correctness expressed using pointwise history
- ▶ one history is allowed to be prefix of other

## Hatcliff et al, HOSC 2000

- ▶ assumes unique end node
- ▶ used for model checking, implying
  - ▶ sequencing of observables must be preserved
  - ▶ liveness properties must be preserved
- ▶ detailed correctness proof, assuming termination of both source and slice

## Hatcliff et al, SAS 1999

- ▶ considers multi-threading
- ▶ proposes bisimulation as correctness property
- ▶ does not work out the details

Foundation for slicing which

- ▶ handles CFGs that arise in modeling reactive systems
- ▶ produces manageable slices
- ▶ allows crisp correctness proof

Foundation for slicing which

- ▶ handles CFGs that arise in modeling reactive systems
- ▶ produces manageable slices
- ▶ allows crisp correctness proof

**Open question:** Is (bi)simulation a strictly stronger correctness criterion than “pointwise” (bi)simulation?

If so, we might modify  $\xrightarrow{wod}$  to allow for independent assignments to be swapped.

- ▶ Interprocedural analysis
- ▶ **Multi**-threading
- ▶ Connect to information flow analysis  
(Hammer & Krinke & Snelting)