NATIONAL
ICT AUSTRALIA

# Static Analysis

Ralf Huuck

**Australian Government**
**Department of Communications,
Information Technology and the Arts**
**Australian Research Council**

NICTA Members
ANU
UNSW
Department of State and Regional Development
BusinessACT
ACT GOVERNMENT

Victoria
The University of Sydney
Queensland Government
Griffith
The University of Queensland

NICTA Partners

---

NATIONAL
ICT AUSTRALIA

## Outline

- Introduction
- Basic Definitions
- Data Flow Analysis
- Abstract Interpretation
- Syntactical Model Checking
- Summary

---

NATIONAL
ICT AUSTRALIA

# Introduction

---

NATIONAL
ICT AUSTRALIA

## What is Static Analysis?

Static Analysis subsumes all methods which derive information from programs without actually executing them.

**Static analysis** is the term applied to the analysis of computer software that is performed without actually executing programs.
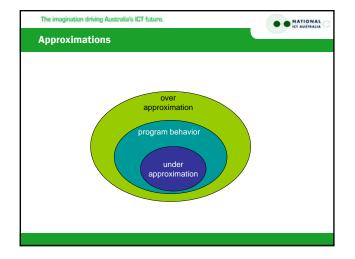
- wikipedia

## The Trouble with Software

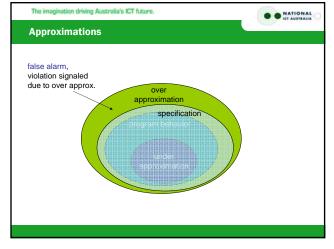The state space of programs is in theory infinite.
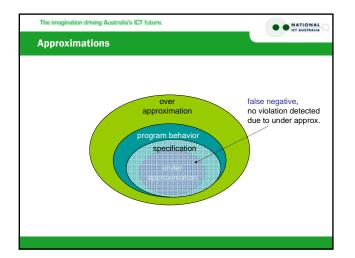
Computation depends on

- integers
- reals
- etc.

In reality it is finite, e.g., 32-bit representations, which is still practically infinite when exploring all combinations.
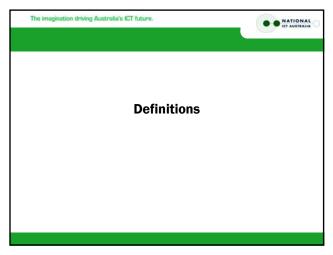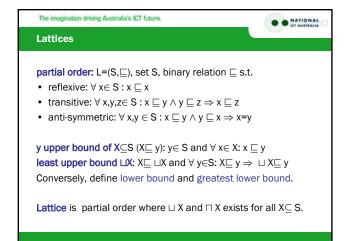
## Decidability

Properties on infinite state spaces are typically **undecidable**, i.e., there is no general algorithm to decide if they are true or false.

Alan Turing
1912-1954

**Rice's theorem:** Any nontrivial property about the language recognized by a Turing machine is undecidable.

But we can still attempt to give useful **approximate** solutions.
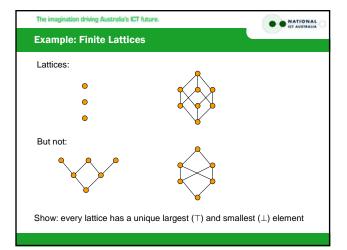
## Approximations



over approximation

program behavior

under approximation

## Approximations



false alarm, violation signaled due to over approx.

over approximation

specification

program behavior

under approximation

## Approximations



over approximation

false negative, no violation detected due to under approx.

program behavior
specification

under approximation

# Definitions

## Lattices

**partial order:** $L=(S,\sqsubseteq)$, set S, binary relation $\sqsubseteq$ s.t.

- reflexive: $\forall x \in S : x \sqsubseteq x$
- transitive: $\forall x,y,z \in S : x \sqsubseteq y \land y \sqsubseteq z \Rightarrow x \sqsubseteq z$
- anti-symmetric: $\forall x,y \in S : x \sqsubseteq y \land y \sqsubseteq x \Rightarrow x=y$

**y upper bound of $X \subseteq S$** ($X \sqsubseteq y$): $y \in S$ and $\forall x \in X: x \sqsubseteq y$
**least upper bound $\sqcup X$:** $X \sqsubseteq \sqcup X$ and $\forall y \in S: X \sqsubseteq y \Rightarrow \sqcup X \sqsubseteq y$
Conversely, define lower bound and greatest lower bound.

**Lattice** is partial order where $\sqcup X$ and $\sqcap X$ exists for all $X \subseteq S$.

## Example: Finite Lattices

Lattices:



But not:



Show: every lattice has a unique largest ($\top$) and smallest ($\bot$) element

## Monotone functions

f monotone: $\forall\, x,y \in S: x \sqsubseteq y \Rightarrow f(x) \sqsubseteq f(y)$

Note: compositions of monotone functions are monotone

Theorem: In lattice L with finite height every monotone function has a least fixed-point as:

$$fix(f) = \bigsqcup_{i \geq 0} f^i(\bot)$$

for which $f(fix(f))=fix(f)$.

---

## Closure

If $L_1,\ldots,L_n$ lattices of finite height, so is the product
$$L_1 \times \ldots \times L_n = \{(x_1,\ldots,x_n) \mid x_i \in L_i\}$$
where $\sqsubseteq$ is defined pointwise.

The height of the product is the sum of heights of its components.

Map: finite set A, lattice L with $|A|=|L|$ then with a pointwise order
$$A \longmapsto L = \{[a_1 \mapsto x_1,\ldots,a_n \mapsto a_n] \mid x_i \in L\}$$
is a lattice of height $|A|*height(L)$.

Other compositions: +, lift, flat, are also lattices of finite height again.

---

## Equation Systems (ES)

Equation System $F:L^n \rightarrow L^n$

$$x_1 = F_1 (x_1,\ldots,x_n)$$
$$\vdots$$
$$x_2 = F_2 (x_1,\ldots,x_n)$$
$$x_n = F_n (x_1,\ldots,x_n)$$

where
$x_i$ variables and
$F_i : L^n \rightarrow L$ collection of monotone functions

has a least fixed point for the function
$$F(x_1,\ldots,x_n)=(F_1(x_1,\ldots,x_n),\ldots,F_n(x_1,\ldots,x_n))$$

---

## Inequation System (IS)

Inequation System $F:L^n \rightarrow L^n$

$$x_1 \sqsubseteq F_1 (x_1,\ldots,x_n)$$
$$\vdots$$
$$x_2 \sqsubseteq F_2 (x_1,\ldots,x_n)$$
$$x_n \sqsubseteq F_n (x_1,\ldots,x_n)$$

We can show $x \sqsubseteq y \Leftrightarrow x=x \sqcap y$.
Thus, IS F equivalent to following ES:

$$x_1 = x_1 \sqcap F_1 (x_1,\ldots,x_n)$$
$$\vdots$$
$$x_2 = x_2 \sqcap F_2 (x_1,\ldots,x_n)$$
$$x_n = x_n \sqcap F_n (x_1,\ldots,x_n)$$

**Control Flow Graph (CFG)**

A control flow graph is a directed graph where nodes are program points and the edges represent the flow of control between these points.

---

**Example CFG**

Program

```
fun(n) {
    var f;
    f=1;
    while (n>0) {
        f=f*n;
        n=n-1
    }
    return f;
}
```

CFG



---

# Data Flow Analysis

---

**Terms**

Classical Data Flow Analysis is concerned about which data reaches which program point. The analysis is performed on program's CFG and expressed as an (in)equality system over a finite lattice.

Typical examples are:
- liveness
- available expressions
- very busy expressions
- reaching definitions

We will go trough each of them.

## Liveness

A variable is **live** before a program point if it will be read in the remaining program execution without being written first.

```
fun(n) {
    var f;
    f=1;
    while (n>0) {
        f=f*n;
        n=n-1
    }
    return f;
}
```

Which variables are live at which locations? Which are not?

---

## Liveness

Any analysis must find out the set of live variables for each location. We model the domain of live variables by the lattice $(2^{Vars}, \subseteq)$.

```
fun(n) {
    var f;
    f=1;
    while (n>0) {
        f=f*n;
        n=n-1
    }
    return f;
}
```

How does the lattice $(2^{Vars}, \subseteq)$ look like? Why is it a lattice?
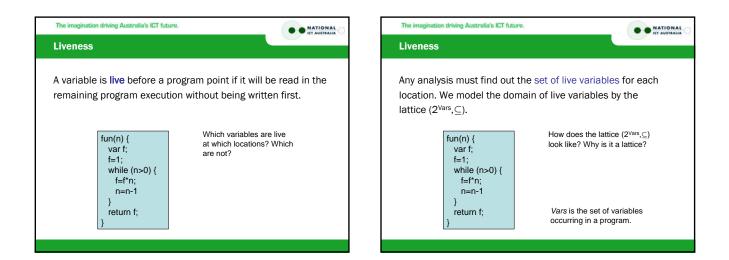
*Vars* is the set of variables occurring in a program.

---

## Questions
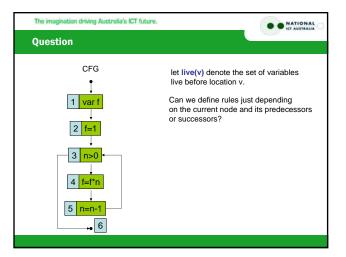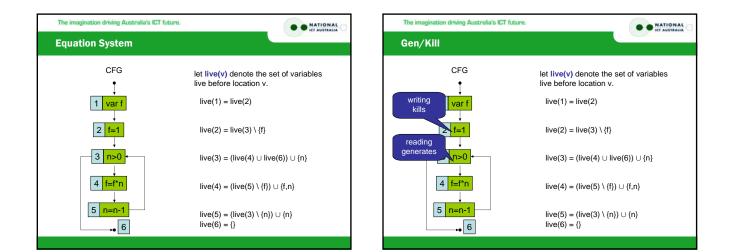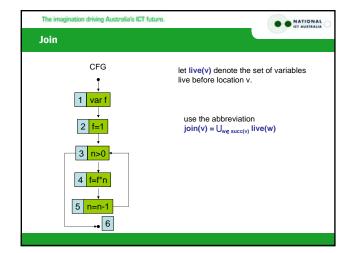
- Can we approximate (without executing the program) the set of live variables algorithmically?
- How does the set depend on our syntax?
- Can we define rules for each construct?
- How can rules lead to something that we can compute?
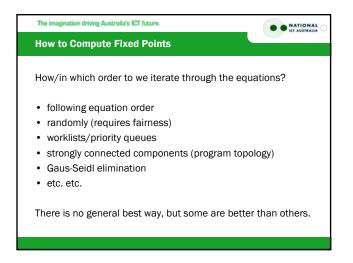
---

## Question

CFG

```
1  var f
2  f=1
3  n>0
4  f=f*n
5  n=n-1
6
```

let **live(v)** denote the set of variables live before location v.

Can we define rules just depending on the current node and its predecessors or successors?

6

## Equation System

CFG

1 | var f
2 | f=1
3 | n>0
4 | f=f*n
5 | n=n-1
6

let **live(v)** denote the set of variables live before location v.

live(1) = live(2)

live(2) = live(3) \ {f}

live(3) = (live(4) ∪ live(6)) ∪ {n}

live(4) = (live(5) \ {f}) ∪ {f,n}

live(5) = (live(3) \ {n}) ∪ {n}
live(6) = {}

## Gen/Kill

CFG

writing kills

reading generates

1 | var f
2 | f=1
3 | n>0
4 | f=f*n
5 | n=n-1
6

let **live(v)** denote the set of variables live before location v.

live(1) = live(2)

live(2) = live(3) \ {f}

live(3) = (live(4) ∪ live(6)) ∪ {n}

live(4) = (live(5) \ {f}) ∪ {f,n}

live(5) = (live(3) \ {n}) ∪ {n}
live(6) = {}

## Join

CFG

1 | var f
2 | f=1
3 | n>0
4 | f=f*n
5 | n=n-1
6

let **live(v)** denote the set of variables live before location v.

use the abbreviation
**join(v) = ∪$_{w \in succ(v)}$ live(w)**

## Join

CFG

1 | var f
2 | f=1
3 | n>0
4 | f=f*n
5 | n=n-1
6

let **live(v)** denote the set of variables live before location v.

live(1) = join(1) \ {f}

live(2) = join(2) \ {f}

live(3) = join(3) ∪ {n}

live(4) = (join(4) \ {f}) ∪ {f,n}

live(5) = (join(5) \ {n}) ∪ {n}
live(6) = {}

**Fixed Point**

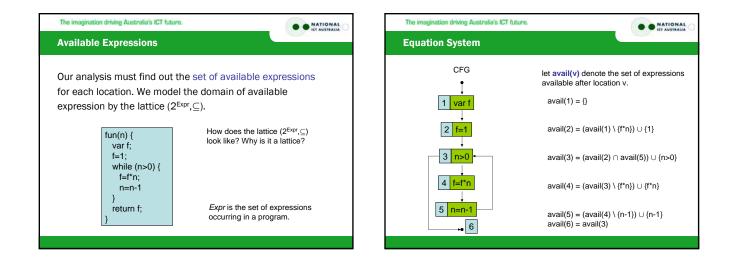let **live(v)** denote the set of variables live before location v.

The right hand-side of each equation is **monotone**, i.e., we can compute the fixed point of ES.

We are interested in the **least** fixed point.

Least fixed point: start with {} greatest start with Vars.

$live(1) = join(1) \setminus \{f\}$

$live(2) = join(2) \setminus \{f\}$

$live(3) = join(3) \cup \{n\}$

$live(4) = (join(4) \setminus \{f\}) \cup \{f,n\}$

$live(5) = (join(5) \setminus \{n\}) \cup \{n\}$
$live(6) = \{\}$

---

**Computing Fixed Point**

let **live(v)** denote the set of variables live before location v.

Computation on the board.

$live(1) = live(2)$

$live(2) = live(3) \setminus \{f\}$

$live(3) = (live(4) \cup live(6)) \cup \{n\}$

$live(4) = (live(5) \setminus \{f\}) \cup \{f,n\}$

$live(5) = (live(3) \setminus \{n\}) \cup \{n\}$
$live(6) = \{\}$

---

**Least Fixed Point Solution**

CFG

```
1  var f
2  f=1
3  n>0
4  f=f*n
5  n=n-1
6
```

let **live(v)** denote the set of variables live before location v.

$live(1) = \{n\}$

$live(2) = \{n\}$

$live(3) = \{f,n\}$

$live(4) = \{f,n\}$

$live(5) = \{f,n\}$
$live(6) = \{\}$

---

**How to Compute Fixed Points**

How/in which order to we iterate through the equations?

- following equation order
- randomly (requires fairness)
- worklists/priority queues
- strongly connected components (program topology)
- Gaus-Seidl elimination
- etc. etc.

There is no general best way, but some are better than others.

## Intermediate Summary

Seen so far:
- data flow analysis problem can be expressed in terms of fixed point over equations
- equations depend on syntax of program points and what is coming in/going out
- many ways to compute fixed point

We have seen Join as being union of successors and computation of least fixed point. This is not always so ...

---

## Available Expressions

An expression is **available** after a program point if its current value has been evaluated before and none of its variables are overwritten. (Good for optimizations)

```
fun(n) {
   var f;
   f=1;
   while (n>0) {
      f=f*n;
      n=n-1
   }
   return f;
}
```

Which expressions are available at which locations? Which are not?

---

## Available Expressions

Our analysis must find out the set of available expressions for each location. We model the domain of available expression by the lattice $(2^{Expr}, \subseteq)$.

```
fun(n) {
   var f;
   f=1;
   while (n>0) {
      f=f*n;
      n=n-1
   }
   return f;
}
```

How does the lattice $(2^{Expr}, \subseteq)$ look like? Why is it a lattice?

*Expr* is the set of expressions occurring in a program.

---

## Equation System

CFG

1  var f
2  f=1
3  n>0
4  f=f*n
5  n=n-1
6

let **avail(v)** denote the set of expressions available after location v.

$avail(1) = \{\}$

$avail(2) = (avail(1) \setminus \{f*n\}) \cup \{1\}$

$avail(3) = (avail(2) \cap avail(5)) \cup \{n>0\}$

$avail(4) = (avail(3) \setminus \{f*n\}) \cup \{f*n\}$

$avail(5) = (avail(4) \setminus \{n-1\}) \cup \{n-1\}$
$avail(6) = avail(3)$

## Join/Fixed Point

CFG

1 | var f
2 | f=1
3 | n>0
4 | f=f*n
5 | n=n-1
6 |

let **avail(v)** denote the set of variables available after location v.

We can use the abbreviation
**join(v) = ∩ _w∈ pred(v)_ avail(w)**
to see we have again a monotone framework.

This time we compute the **greatest fixed point**, as we like to have the maximum number of available expressions.

---

## Greatest Fixed Point Solution

CFG

1 | var f
2 | f=1
3 | n>0
4 | f=f*n
5 | n=n-1
6 |

let **avail(v)** denote the set of expressions available after location v.

avail(1) = {}

avail(2) = {1}

avail(3) = {n>0}

avail(4) = {n>0, f*n}

avail(5) = {f*n ,n-1}
avail(6) = {n>0}

---

## Very Busy Expressions

An expression is **very busy** before a program point if it definitely will be evaluated again before its value changes.

```
fun(n) {
    var f;
    f=1;
    while (n>0) {
       f=f*n;
       n=n-1
    }
    return f;
}
```

Which expressions are very busy at which locations? Which are not?

---

## Very Busy Expressions

Our analysis must find out the set of very busy expressions for each location. We model the domain of very busy expression by the lattice $(2^{Expr}, \subseteq)$.

```
fun(n) {
    var f;
    f=1;
    while (n>0) {
       f=f*n;
       n=n-1
    }
    return f;
}
```

How does the lattice $(2^{Expr}, \subseteq)$ look like? Why is it a lattice?

*Expr* is the set of expressions occurring in a program.

10

## Slide 1: Equation System

**Equation System**

CFG

1 var f
2 f=1
3 n>0
4 f=f*n
5 n=n-1
6

let **busy(v)** denote the set of expressions very busy before location v.

busy(1) = buys(2)

busy(2) = (busy(3) \ {f*n}) ∪ {1}

busy(3) = (busy(2) ∩ busy(5)) ∪ {n>0}

busy(4) = (busy(5) \ {f*n}) ∪ {f*n}

busy(5) = (busy(3) \ {n-1}) ∪ {n-1}
busy(6) = { }

## Slide 2: Join/Fixed Point

**Join/Fixed Point**

CFG

1 var f
2 f=1
3 n>0
4 f=f*n
5 n=n-1
6

let **busy(v)** denote the set of variables very busy before location v.

We can use the abbreviation
**join(v) = ∩ $_{w ∈ succ(v)}$ busy(w)**
to see we have again a monotone framework.

This time we compute the **least fixed point** again.

## Slide 3: Least Fixed Point Solution

**Least Fixed Point Solution**

CFG

1 var f
2 f=1
3 n>0
4 f=f*n
5 n=n-1
6

let **busy(v)** denote the set of expressions very busy before location v.

busy(1) = {n>0, 1}

busy(2) = {n>0, 1}

busy(3) = {n>0}

busy(4) = {n>0, n-1, f*n}

busy(5) = {n>0, n-1}
busy(6) = { }

## Slide 4: Reaching Definition

**Reaching Definition**

The effect of an assignment is **reaching** after a program point if it might have defined the current values of variables.

```
fun(n) {
    var f;
    f=1;
    while (n>0) {
        f=f*n;
        n=n-1
    }
    return f;
}
```

Which assignments are reaching at which locations? Which are not?

## Reaching Definitions

Our analysis must find out the set of reaching assignments for each location. We model the domain of reaching definitions by the lattice ($2^{Assign}, \subseteq$).

```
fun(n) {
    var f;
    f=1;
    while (n>0) {
        f=f*n;
        n=n-1
    }
    return f;
}
```

How does the lattice ($2^{Assign}, \subseteq$) look like? Why is it a lattice?

*Assign* is the set of assignments occurring in a program.

---

## Equation System

CFG

1 var f
2 f=1
3 n>0
4 f=f*n
5 n=n-1
6

let **reach(v)** denote the set of assignments may define variable values after location v.

reach(1) = {}

reach(2) = (reach(1) \ {f=1}) ∪ {f=1}

reach(3) = reach(2) ∪ reach(5)

reach(4) = (reach(3) \ {f=f*n,f=1}) ∪ {f=f*n}

reach(5) = (reach(4) \ {n=n-1}) ∪ {n=n-1}
reach(6) = reach(3)

---

## Join/Fixed Point

CFG

1 var f
2 f=1
3 n>0
4 f=f*n
5 n=n-1
6

let **reach(v)** denote the set of assignments may define variable values after location v.

We can use the abbreviation
**join(v) =⋃$_{w∈\,pred(v)}$ reach(w)**
to see we have again a monotone framework.

This time we compute the **least fixed point** again.

---

## Least Fixed Point Solution

CFG

1 var f
2 f=1
3 n>0
4 f=f*n
5 n=n-1
6

let **reach(v)** denote the set of assignments may define variable values after location v.

reach(1) = {}

reach(2) = {f=1}

reach(3) = {f=1, f=f*n, n=n-1}

reach(4) = {f=f*n, n=n-1}

reach(5) = {f=f*n, n=n-1}
reach(6) = {f=1, f=f*n, n=n-1}

## Summary Data Flow Analysis

- we have seen how to compute approximate solutions (we do not know if all the paths are executed!) to data flow problems
- Join can be intersection or union
- analysis has forward or backward nature
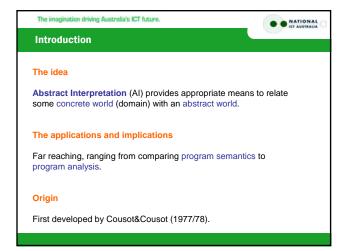- depending on the problem least or greatest fixed point

## Rules of Thumb

**forward analysis**: computes information about past behavior
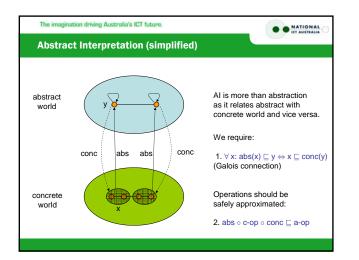
**backward analysis**: computes information about future behavior

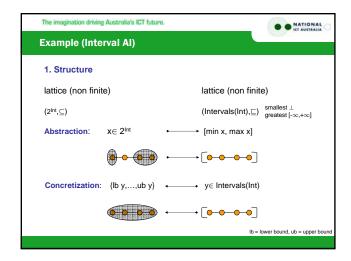**must analysis**: information that must be true (on all paths) and computes under-approximation

**may analysis**: information that may be true (on at least one paths) and computes over-approximation

# Abstract Interpretation

The Rough Guide

## Introduction

**The idea**

**Abstract Interpretation** (AI) provides appropriate means to relate some concrete world (domain) with an abstract world.

**The applications and implications**

Far reaching, ranging from comparing program semantics to program analysis.

**Origin**

First developed by Cousot&Cousot (1977/78).

## Abstract Interpretation (simplified)



AI is more than abstraction as it relates abstract with concrete world and vice versa.

We require:

1. $\forall x: abs(x) \sqsubseteq y \Leftrightarrow x \sqsubseteq conc(y)$
(Galois connection)

Operations should be safely approximated:

2. $abs \circ c\text{-}op \circ conc \sqsubseteq a\text{-}op$

---

## Example (Interval AI)

### 1. Structure

| Sets of integers | Intervals |
|---|---|
| {2} | [2,2] |
| {2,3,4} | [2,4] |
| {1,3,9} | [1,9] |
| $x \in 2^{Int}$ | smallest interval comprising x |

Intervals over approximate sets.

---

## Example (Interval AI)

### 1. Structure

lattice (non finite)                lattice (non finite)

$(2^{Int}, \subseteq)$         $(Intervals(Int), \sqsubseteq)$   smallest $\bot$
                                                   greatest $[-\infty,+\infty]$

**Abstraction**:     $x \in 2^{Int}$  ⟷  [min x, max x]

**Concretization**:   {lb y,…,ub y}  ⟷  $y \in Intervals(Int)$

lb = lower bound, ub = upper bound

---

## Example (Interval AI)

### 2. Operations

lattice (non finite)                lattice (non finite)

$(2^{Int}, \subseteq)$         $(Intervals(Int), \sqsubseteq)$

**Concrete world**:  {2,5,6} + {2,3} = {4,5,7,8,9}

**Abstract world**:  [2,6] + [2,3] = [4,9]

a) Introduce matching operator in abstract world for every operator in concrete world.

b) Check it satisfies safe approximation.

14

## Application

We like to compute all the possible values variables can take (collecting semantics).

```
var f, n;
n=3;
f=1;
while (n>0) {
  f=f*n;
  n=n-1
}
return f;
```

---

## Application

**before each location the variables might be as follows:**

1 var f,n
2 n=3
3 f=1
4 n>0
5 f=f*n
6 n=n-1
7

2: f: $\{-\infty, ...., +\infty\}$      n: $\{-\infty, ...., +\infty\}$

3: f: $\{-\infty, ...., +\infty\}$      n: $\{3\}$

4: f: $\{1\}$      n: $\{3\}$

5: f: $\{1,3,6\}$      n: $\{1,2,3\}$

6: f: $\{1,3,6\}$      n: $\{0,1,2,3\}$

7: f: $\{1,3,6\}$      n: $\{0\}$

---

## Application

1 var f,n
2 n=3
3 f=1
4 n>0
5 f=f*n
6 n=n-1
7

**Problems**:

1. We might need infinite space to store values.
2. We might not be able to compute them due to non-termination.

**Solution**:

1. can be overcome by using interval AI (over-approximation of values)

---

## Application

**before each location the variables might be as follows:**

1 var f,n
2 n=3
3 f=1
4 n>0
5 f=f*n
6 n=n-1
7

2: f: $[-\infty, +\infty]$    n: $[-\infty, +\infty]$

3: f: $[-\infty, +\infty]$    n: $[3,3]$

4: f: $[1,1]$    n: $[3,3]$

5: f: $[1,6]$    n: $[1,3]$

6: f: $[1,6]$    n: $[0,3]$

7: f: $[1,6]$    n: $[0,0]$

## Application

| 1 | var f,n |
| 2 | n=3 |
| 3 | f=1 |
| 4 | n>0 |
| 5 | f=f*n |
| 6 | n=n-1 |
| 7 | |

**Problems**:

1. We might need infinite space to store values.
2. We might not be able to compute them due to non-termination.

**Solution**:

1. can be overcome by using interval AI (over-approximation of all values)
2. not that easy: interval lattice has infinite width (not a problem) and infinite height (problem!).

## Acceleration

infinite height

AI:

We have relation between abstract and concrete domain.

We have relation between abstract and concrete operators.

In order to deal with infinite lattices (i.e. to compute a fixed point in finite time) we introduce an extra operator that can "jump" infinitely high.

Such an acceleration operator is called **widening** operator. Sometimes people speak of dynamic approximation operator.

## Application

| 1 | var f,n |
| 2 | n=3 |
| 3 | f=1 |
| 4 | n>0 |
| 5 | f=f*n |
| 6 | n=n-1 |
| 7 | |

put widening operator here

Idea:

1. if variable value increased compared to previous iteration jump to $+\infty$, if decreased to $-\infty$.

2. goto 2

Since our operations are monotone this is safe.

## Application

| 1 | var f,n |
| 2 | n=3 |
| 3 | f=1 |
| 4 | n>0 |
| 5 | f=f*n |
| 6 | n=n-1 |
| 7 | |

Result:

2: f: $[-\infty, +\infty]$   n: $[-\infty, +\infty]$

3: f: $[-\infty, +\infty]$   n: $[3,3]$

4: f: $[1,1]$   n: $[3,3]$

5: f: $[1,+\infty]$   n: $[-\infty,3]$

6: f: $[1, +\infty]$   n: $[-\infty,3]$

after first acceleration

## Slide 1

The imagination driving Australia's ICT future.

NATIONAL ICT AUSTRALIA

### Application

```
  1 | var f,n
  2 | n=3
  3 | f=1
  4 | n>0
  5 | f=f*n
  6 | n=n-1
  7 |
```

Result:

2: f: [-∞, +∞]    n: [-∞, +∞]

3: f: [-∞, +∞]    n: [3,3]

4: f: [1,1]    n: [3,3]

after second acceleration

approximation is very coarse

can narrow it down using condition as constraint

5: f: [-∞,+∞]    n: [-∞,+∞]

6: f: [-∞, +∞]    n: [-∞,+∞]

6: f: [-∞, +∞]    n: [-∞,+∞]

## Slide 2

The imagination driving Australia's ICT future.

NATIONAL ICT AUSTRALIA

### Application

```
  1 | var f,n
  2 | n=3
  3 | f=1
  4 | n>0
  5 | f=f*n
  6 | n=n-1
  7 |
```

Result:

2: f: [-∞, +∞]    n: [-∞, +∞]

3: f: [-∞, +∞]    n: [3,3]

4: f: [1,1]    n: [3,3]

We might loose a lot of information, but we are still able to tell that n>0.

5: f: [-∞,+∞]    n: [-∞,+∞]

6: f: [-∞, +∞]    n: [-∞,+∞]

6: f: [-∞, +∞]    n: [-\infty,0]

## Slide 3

The imagination driving Australia's ICT future.

NATIONAL ICT AUSTRALIA

### Summary

- AI relates abstract and concrete worlds (structure + operations)
- termination/safe approximation can be enforced by acceleration techniques

There are domains that capture more information, e.g., ployhedra.

AI is good for range approximation, i.e., array access, range check of operations, general buffer overflows.

## Slide 4

The imagination driving Australia's ICT future.

NATIONAL ICT AUSTRALIA

### Model Checking Syntax

Just some ideas …

17

The imagination driving Australia's ICT future.

NATIONAL ICT AUSTRALIA

## Introduction

Syntax gives us some information:
- when are variable is used
- when a variable is declared
- when a variable is modified
- etc.

Can we make use of it to find bugs in programs?

---

The imagination driving Australia's ICT future.

NATIONAL ICT AUSTRALIA

## Syntactical Information

1 var f,n     declaration_f   declaration_n

2 n=3     modified_n

3 f=1     modified_f

4 n>0     used_n

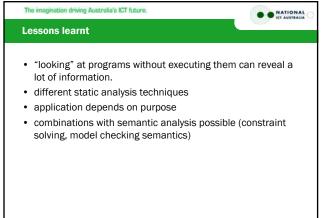5 f=f*n     used_n used_f modified_f

6 n=n-1     used_n modified_n

7

---

The imagination driving Australia's ICT future.

NATIONAL ICT AUSTRALIA

## Syntactical Information

CFG is a transition system

1 var f,n     declaration_f   declaration_n

2 n=3     modified_n

these are just like atomic propositions

3 f=1     modified_f

4 n>0     used_n

5 f=f*n     used_n used_f modified_f

6 n=n-1     used_n modified_n

7

---

The imagination driving Australia's ICT future.

NATIONAL ICT AUSTRALIA

## Kripke Structure

transition system   **+**   atomic propositions

Kripke Structure

So can we model check syntax?

18

## Model Checking Syntax



| Node | | Label |
|---|---|---|
| 1 | var f,n | declaration_f declaration_n |
| 2 | n=3 | modified_n |
| 3 | f=1 | modified_f |
| 4 | n>0 | used_n |
| 5 | f=f*n | used_n used_f modified_f |
| 6 | n=n-1 | used_n modified_n |
| 7 | | |

**Yes**, e.g.:

AG (delclaration$_f$ $\Rightarrow$ EF used$_f$)

AG (modified$_n$ $\Rightarrow$ EF used$_n$)

also if variables are initialized, certain protocols are respected, locks are released etc.

**But**: Abstraction is sometimes neither sound nor complete.

---

## Summary

- model checking syntax is good for finding bugs
- not so good for showing the absence of bugs/verification
- very efficient
- easy to use

---

## Summary

---

## Lessons learnt

- "looking" at programs without executing them can reveal a lot of information.
- different static analysis techniques
- application depends on purpose
- combinations with semantic analysis possible (constraint solving, model checking semantics)

NATIONAL
ICT AUSTRALIA

## Next Week

Model Checking Real-Time Systems