

Scaling Program Analysis

CS7575

Static program analysis is the analysis of computer software performed without actually executing programs.

Wikipedia

<http://surl.li/lboxd>

Shared notebook

With material from "Principles of Program Analysis" by Nielsen² & Hankin

1

Course Structure

Part 1: Lectures

Aim: cover chapters 1-4 of PoPA (dataflow, constraint-based, AI)

Part 2: Papers & Projects

Broad overview of applications of static analysis in software engineering, and compilers, including dynamic analysis, big code and machine learning

2

Workload

Option 1: TestOut + Talk

The midterm test and a paper prevention are your whole grade

Option 2: Repetition

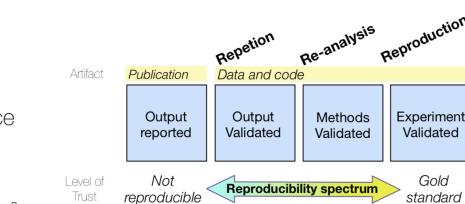
Pick a paper, re-run experiments, present and write a report

Option 3: Re-analysis

Option 2 + re-analyze the results

Option 4: Research

Work on something new of your choice



3

Program Analysis

Characterize possible executions behaviors of a program

Use results for:

Compiler optimizations

Program verification

Error detection

Vulnerability detection

Static Analysis

Compile-time characterization of all program executions by approximation

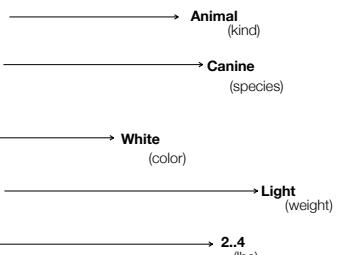
Dynamic Analysis

Run-time characterization of some program executions by observation

4

Abstraction

An abstraction is a property from some domain

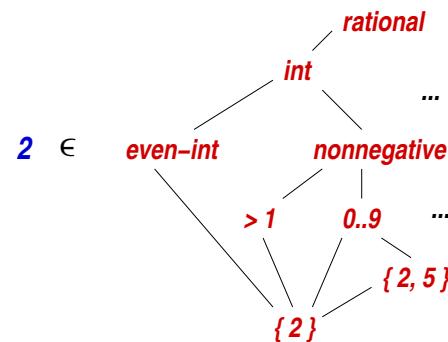


5

A static analysis is determined its choice of abstraction:
precise enough to capture the properties of interest, approximate enough to be computable

© David Schmidt

Value abstractions are common



All the properties listed on the right are abstractions of 2; the upwards lines denote \sqsubseteq , a loss of precision.

6

© David Schmidt

Abstractions can be relations

These Figures are from *Abstract Interpretation: Achievements and Perspectives* by Patrick Cousot, Proc. SSGRR 2000.

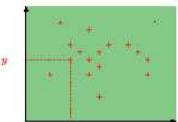


Fig. 2
SIGNS

$$\begin{cases} x \geq 0 \\ y \geq 0 \end{cases}$$

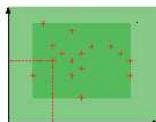


Fig. 3
INTERVALS

$$\begin{cases} x \in [3, 27] \\ y \in [4, 32] \end{cases}$$

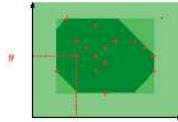


Fig. 4
OCTAGONS

$$\begin{cases} 3 \leq x \leq 27 \\ x + y \leq 88 \\ 4 \leq y \leq 32 \\ x - y \leq 61 \end{cases}$$

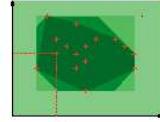


Fig. 5
POLYHEDRA

$$\begin{cases} 7x + 31y \leq 325 \\ 21x + 7y \geq 0 \end{cases}$$

© David Schmidt

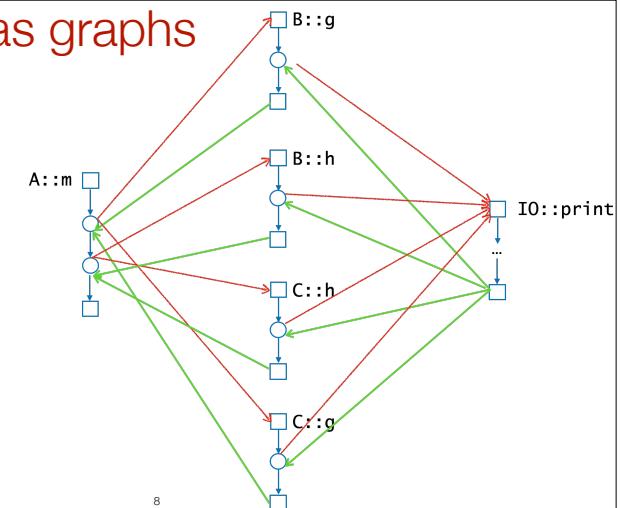
Abstractions as graphs

```

class A {
    m( x ) {
        x.g();
        x.h();
    }
}

class B {
    g( ) {
        IO.print("B::g");
    }
    h( ) {
        IO.print("B::h");
    }
}

class C {
    g( ) {
        IO.print("B::g");
    }
    h( ) {
        IO.print("B::h");
    }
}
  
```



8

Intermezzo Static Analysis at Facebook

contributed articles

Key lessons for designing static analysis tools deployed to find bugs in hundreds of millions of lines of code

ELEONOR DISTEFANO, MANUEL FÄHNDRICH, FRANCESCO LOGOZZO, AND PETER W. O'HEARN

Scaling Static Analyses at Facebook

STATIC ANALYSIS TOOLS are programs that examine, and attempt to draw conclusions about, the source of other programs without running them. At Facebook, we have developed a number of static analysis tools that employ reasoning techniques similar to those from program verification. The tools we describe here have been used to detect bugs in our codebase, and to the security of our services, they perform sometimes complex reasoning spanning multiple files and multiple components. In this paper, we present our rationale for developing static analysis tools, and the challenges that arise in scaling them to handle codebases as large as ours.

Content for Static Analysis at Facebook

Facebook has developed static analysis tools to prevent bugs that would otherwise be introduced into engineering workflows in a way that attempts to bring value while minimizing friction.

Infer targets our mobile apps as well as our backend C++ code, and Zoncolan targets our website. Infer has found over 100 thousand reported issues fixed by developers before code reaches production. Zoncolan targets the millions of lines of Java code that run on our website, and additionally

Key Insights

- Advanced static analysis techniques can scale to analyze code bases as large as Facebook's.
- Data analysis should enable a developer to quickly identify and understand reports (false positives, false negatives, and other errors).
- Programs that generate reports should be part of the developer's workflow.

9

It is important for a static analysis tool to be able to handle different types of code. For example, static analysis tools can handle the same different bugs can have different severity levels, depending on the context and the nature of the bug. For example, a bug in a critical component of a system might be more important than a bug in a less critical component. Additionally, the frequency of occurrence of a bug can also affect its importance. For example, if a certain type of bug occurs frequently, it may be more important than a bug of similar severity that occurs rarely.

We have several means to collect information about bugs. One way is to use static analysis tools, such as Infer or Zoncolan. These tools can analyze code and report bugs that are the same, or different, depending on the company's needs. For example, Infer can report various types of bugs, such as memory leaks, null pointer exceptions, and so on. Zoncolan can report bugs related to security, such as SQL injection, cross-site scripting, and so on.

Monte Zoncolan is a mountain in the Cimone Alps, Italy, with an elevation of 1,750 meters. It is one of the most demanding climbs in professional road bicycle racing.

10

Intermezzo Static Analysis at Facebook

Not all bugs are the same!

A memory leak on a seldom-used service not as important as a vulnerability granting access to unauthorized information. A crash, such as a Java NPE, that happens hourly is more important than the same bug happening yearly.

Facebook built two tools: **Infer** and **Zoncolan**

Infer

Target 10M lines of Java, Objective-C, C++ on Android/iOS for FB, Instagram, Messenger & WhatsApp, and backend C++ or Java

Origins: Peter O'Hearn ++ (UCL~Monoidics) separation logic

Zoncolan

Target 100M lines of Hack for Facebook's website

Origins: Fahndrich (EPFLBerkeley~MSR) & Logozzo (Polytechnique/R. Cousot~MSR)

Monte Zoncolan is a mountain in the Cimone Alps, Italy, with an elevation of 1,750 meters. It is one of the most demanding climbs in professional road bicycle racing.

Distefano, Fahndrich, Logozzo, O'Hearn: *Scaling Static Analyses at Facebook*, CACM'19, doi=3351434.3338112

10

Intermezzo Static Analysis at Facebook

Should **Infer** report errors, if so where?
What are the abstractions?

```
r = new R();
...
x = r.m;
x = r.m;
```

11

Intermezzo Static Analysis at Facebook

Should **Infer** report errors, if so where?
What are the abstractions?

```
@ThreadSafe
class R {
    int m;
    void pWOMT() {
        Utils.assertMainThread();
        synchronized(this){ m++; }
    }

    int uROMT() {
        Utils.assertMainThread();
        return m;
    }

    synchronized int pROMT() {
        Utils.assertMainThread();
        return m;
    }
}
```

```
synchronized void pWOMTB() {
    m+=2;
}

int uROMTB() {
    return m;
}
```

12

```

1 <?hh
2 class AddMemberToGroup extends FacebookEndpoint {
3     private function getIDs (): (string, int) {
4         // User input, untrusted
5         return tuple((string) $this->getRequest('member_id'),
6             (int) $this->getRequest('gid'));
7     }
8
9     public function render(): :xhp {
10        list($member_id, $group_id) = $this->getIDs();
11        return this->getConfirmationForm($group_id, $member_id);
12    }
13
14    public function getConfirmationForm
15        (int $group_id, string $member_id): :xhp {           (5, user - input)
16        $url = "https://facebook.com/groups/add_member/" .      ↓
17        $member_id;                                         (10, GenericGroupForm :: getIDs)
18
19        return
20        <form method="post" action={$url}>
21            <input name="gid" value={$group_id}/>
22            <input name="action" value="add"/>          (11, AddMemberToGroup :: getConfirmationForm). 2
23        </form>;                                         Should Zoncolan report errors, if so where?
24    }
25 }                                         via string concat
                                         (20, Form :: action)
                                         What are the abstractions?
```

13

Intermezzo Static Analysis at Facebook

Takeaways

What makes a static analysis “good”?

14

Intermezzo Static Analysis at Facebook

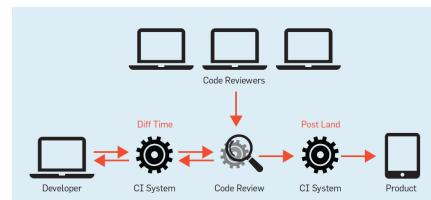
Takeaways:

Goodness = H5+H6+Scalable

Infer

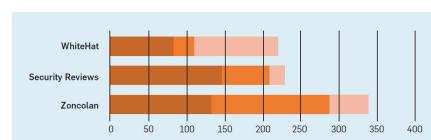
70% fix rate when integrated with CI and
0% fix rate otherwise

15 minutes / diff time (10M codebase)



Zoncolan

80% actioned rate and 11 “missed” bugs
100M LOC/30 minutes in batch mode



15

Intermezzo Static Analysis at Facebook

Takeaways

Do analyses need to scale across procedure boundaries?

| | Intra-procedural Infra-file | Inter-procedural Intra-file | Inter-procedural Inter-file |
|---------------------------|--------------------------------|--------------------------------|--------------------------------|
| Null Deref (Java) | 43 | 9 | 48 |
| Null Deref (ObjC) | 73 | 5 | 24 |
| Thread Safety | 36 | 12 | 53 |
| Allocates Memory | 0 | 2 | 98 |
| Bad Pointer Comparison | 100 | 0 | 0 |

too recent fixes in several bug categories

<https://engineering.fb.com/developer-tools/finding-inter-procedural-bugs-at-scale-with-infer-static-analyzer/>

16

Intermezzo Static Analysis at Facebook

Scalability

Infer is compositional; consider this bug trace

```
src/com/duckduckgo/mobile/android/activity/DuckDuckGo.java:867: error: NULL_DEREFERENCE  
object `feedObject` last assigned on line 866 could be null and is dereferenced by call  
to `feedItemSelected(...)` at line 867.
```

```
864.  
865.     public void feedItemSelected(String feedId) {  
866.         FeedObject feedObject = DDGApplication.getDB().selectFeedById(feedId);  
867.         feedItemSelected(feedObject);
```

<https://engineering.fb.com/developer-tools/finding-inter-procedural-bugs-at-scale-with-infer-static-analyzer/>

17

```
src/com/duckduckgo/mobile/android/db/DDgDB.java:483: start of procedure  
selectFeedById(...)
```

```
482.  
483.     public FeedObject selectFeedById(String id){  
484.         FeedObject out = null;  
485.         Cursor c = null;  
486.         try {  
487.             c = this.db.query(FEED_TABLE, null, "_id=?", new String[]{id}, null, null, null);  
488.             if (c.moveToFirst()) {  
489.                 out = getFeedObject(c);  
490.             }  
491.         } finally {  
492.             if(c!=null) {  
493.                 c.close();  
494.             }  
495.         }  
496.         return out;  
497.     }
```

18

Intermezzo Static Analysis at Facebook

```
src/com/duckduckgo/mobile/android/activity/DuckDuckGo.java:842: start of procedure  
feedItemSelected(...)
```

```
841.  
842.     public void feedItemSelected(FeedObject feedObject) {  
843.  
844.         DDGControlVar.currentFeedObject = feedObject;  
845.         DDGControlVar.mDuckDuckGoContainer.sessionType = SESSIONTYPE.SESSION_FEED;  
846.  
847.         String url = feedObject.getUrl();
```

19

Intermezzo Static Analysis at Facebook

PRE: this.db \mapsto - (db is non-null and points to a valid object)

POST1: return == null (we are returning a null feed object)

POST2: return \mapsto - (we are returning a valid, non-null feed object)

```
483.     public FeedObject selectFeedById(String id){  
484.         FeedObject out = null;  
485.         Cursor c = null;  
486.         try {  
487.             c = this.db.query(FEED_TABLE, null, "_id=?", new String[]{id}, null, null, null);  
488.             if (c.moveToFirst()) {  
489.                 out = getFeedObject(c);  
490.             }  
491.         } finally {  
492.             if(c!=null) {  
493.                 c.close();  
494.             }  
495.         }  
496.         return out;
```

20

Intermezzo Static Analysis at Facebook

PRE: `feedObject` $\mapsto -$ (`feedObject` is a valid, non-`null` object)

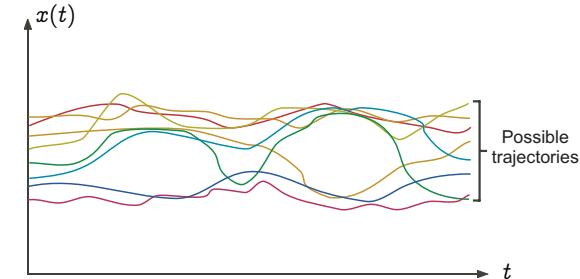
```
865. public void feedItemSelected(String feedId) {  
866.     FeedObject feedObject = DDGApplication.getDB().selectFeedById(feedId);  
867.     feedItemSelected(feedObject);
```

```
842. public void feedItemSelected(FeedObject feedObject) {  
843.  
844.     DDGControlVar.currentFeedObject = feedObject;  
845.     DDGControlVar.mDuckDuckGoContainer.sessionType = SESSIONTYPE.SESSION;  
846.  
847.     String url = feedObject.getUrl();
```

21

Principles of Program Analysis

The *concrete semantics* of a program is a mathematical model of the set of possible execution in possible environments



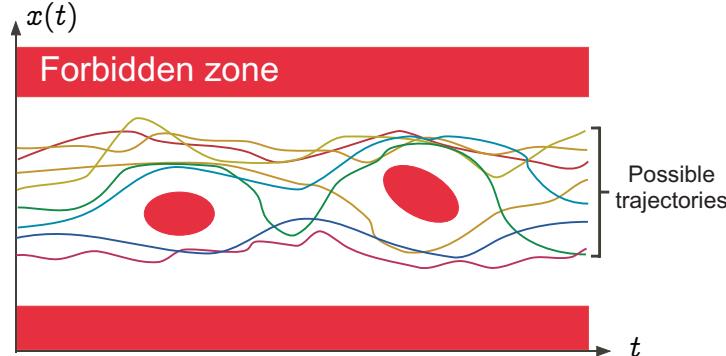
The concrete semantics of a program is *not computable*; all non trivial questions on the concrete semantics are *undecidable*

22

© Patrick Cousot

Safety Properties

The *safety properties* of a program express that no execution in any environment can reach an erroneous state



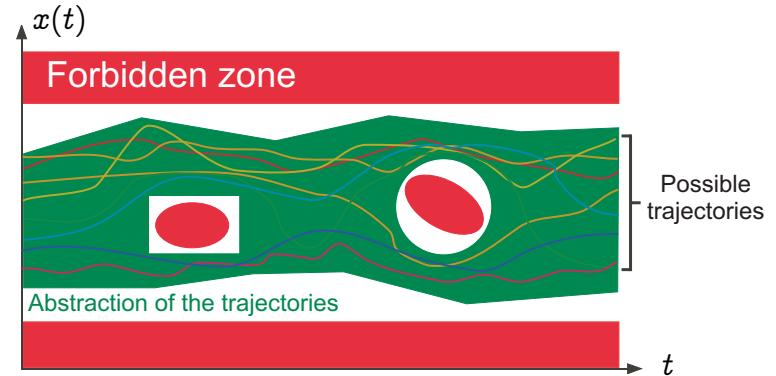
A safety proof consists in proving that the intersection of the concrete semantics and the forbidden zone is empty

23

© Patrick Cousot

Static Program Analysis

The abstract semantics is computed automatically from the program text according to pre-defined abstractions

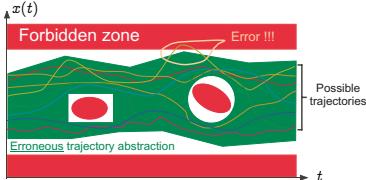


24

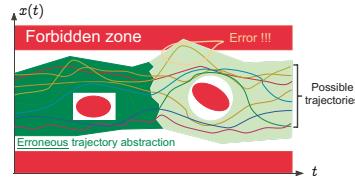
© Patrick Cousot

Errors

e.g.) unsound analysis



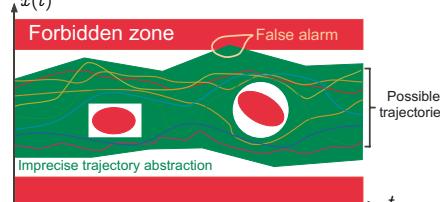
e.g.) bounded model checking



25

© Patrick Cousot

false positive



Static Analysis in Compilers

Consider the C-like program

```
i := 0;
while i <= n do
    j := 0;
    while j <= m do
        temp := Base(A) + i * (m+1) + j;
        Cont(temp) := Cont(Base(B) + i * (m+1) + j)
                     + Cont(Base(C) + i * (m+1) + j);
        j := j+1
    od;
    i := i+1
od
```

Nielsen, Nielsen, Hankin: Principles of Program Analysis, Chapter 1

26

© Nielsen^2, Hankin

Available Expressions analysis & Common Subexpression Elimination

```
i := 0;
while i <= n do      first computation
    j := 0;
    while j <= m do
        temp := Base(A) + i*(m+1) + j;
        Cont(temp) := Cont(Base(B) + i*(m+1) + j)
                     + Cont(Base(C) + i*(m+1) + j);
        j := j+1
    od;
    i := i+1
od
```

re-computations

```
t1 := i * (m+1) + j;
temp := Base(A) + t1;
Cont(temp) := Cont(Base(B)+t1)
             + Cont(Base(C)+t1);
```

Nielsen, Nielsen, Hankin: Principles of Program Analysis, Chapter 1

27

© Nielsen^2, Hankin

Detection of Loop Invariants & Invariant Code Motion

```
i := 0;
while i <= n do
    j := 0;
    loop invariant
    while j <= m do
        t1 := i * (m+1) + j;
        temp := Base(A) + t1;
        Cont(temp) := Cont(Base(B)+t1)
                     + Cont(Base(C)+t1);
        j := j+1
    od;
    i := i+1
od
```

Nielsen, Nielsen, Hankin: Principles of Program Analysis, Chapter 1

28

```
t2 := i * (m+1);
while j <= m do
    t1 := t2 + j;
    temp := ...
    Cont(temp) := ...
    j := ...
od
```

© Nielsen^2, Hankin

Detection of Induction Variables & Reduction of Strength

```
i := 0;
while i <= n do
  j := 0;
  t2 := i * (m+1);
  while j <= m do
    t1 := t2 + j;
    temp := Base(A) + t1;
    Cont(temp) := Cont(Base(B) + t1)
                 + Cont(Base(C) + t1);
    j := j+1
  od;
  i := i+1
od
```

induction variable

```
i := 0;
t3 := 0;
while i <= n do
  j := 0;
  t2 := t3;
  while j <= m do ... od
  i := i + 1;
  t3 := t3 + (m+1)
od
```

Nielsen, Nielsen, Hankin: Principles of Program Analysis. Chapter 1

29

© Nielsen^2, Hankin

Equivalent Expressions analysis & Copy Propagation

```
i := 0;
t3 := 0;
while i <= n do
  j := 0;
  t2 := t3;
  while j <= m do
    t1 := t2 + j;
    temp := Base(A) + t1;
    Cont(temp) := Cont(Base(B) + t1)
                 + Cont(Base(C) + t1);
    j := j+1
  od;
  i := i+1;
  t3 := t3 + (m+1)
```

Nielsen, Nielsen, Hankin: Principles of Program Analysis. Chapter 1

30

© Nielsen^2, Hankin

Live Variables analysis & Dead Code Elimination

```
i := 0;
t3 := 0;
while i <= n do
  j := 0;
  t2 := t3;
  while j <= m do
    t1 := t3 + j;
    temp := Base(A) + t1;
    Cont(temp) := Cont(Base(B) + t1)
                 + Cont(Base(C) + t1);
    j := j+1
  od;
  i := i+1;
  t3 := t3 + (m+1)
od
```

dead variable

```
i := 0;
t3 := 0;
while i <= n do
  j := 0;
  while j <= m do
    t1 := t3 + j;
    temp := Base(A) + t1;
    Cont(temp) := Cont(Base(B) + t1)
                 + Cont(Base(C) + t1);
    j := j+1
  od;
  i := i+1;
  t3 := t3 + (m+1)
```

Nielsen, Nielsen, Hankin: Principles of Program Analysis. Chapter 1

31

© Nielsen^2, Hankin

Summary

Analysis

Available expressions analysis

Detection of loop invariants

Detection of induction variables

Equivalent expression analysis

Live variables analysis

Transformation

Common subexpression elimination

Invariant code motion

Strength reduction

Copy propagation

Dead code elimination

Nielsen, Nielsen, Hankin: Principles of Program Analysis. Chapter 1

32

© Nielsen^2, Hankin

Approaches to Program Analysis

A family of techniques ...

- data flow analysis
- constraint based analysis
- abstract interpretation
- type and effect systems
- ...

... that differ in their focus:

- algorithmic methods
- semantic foundations
- language paradigms
 - imperative
 - object oriented
 - logical
 - functional
 - concurrent

Nielsen, Nielsen, Hankin: *Principles of Program Analysis*. Chapter 1

33

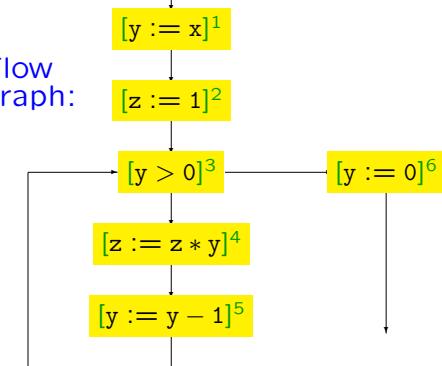
© Nielsen^2, Hankin

Data Flow Analysis

Program with labels for elementary blocks:

```
[y := x]1;
[z := 1]2;
while [y > 0]3 do
  [z := z * y]4;
  [y := y - 1]5
od;
[y := 0]6
```

Flow graph:

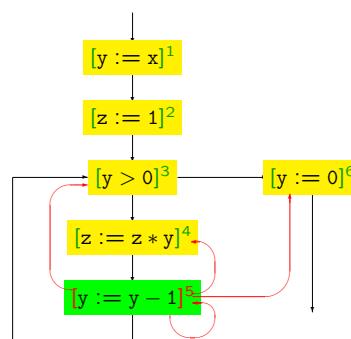


34

© Nielsen^2, Hankin

Reaching Definition

The assignment $[x := a]^\ell$ reaches ℓ' if there is an execution where x was last assigned at ℓ



35

© Nielsen^2, Hankin

Reaching Definition

```
{}((x,?), (y,?), (z,?))

[y := x]1;
[z := 1]2;
while [y > 0]3 do
  [z := z * y]4;
  [y := y - 1]5
od;
[y := 0]6
```

© Nielsen^2, Hankin

36

Reaching Definition

```

[y := x]1; ───────── {((x,?), (y,?), (z,?))}
[z := 1]2; ───────── {((x,?), (y, 1), (z,?))}
while [y > 0]3 do
    ───────── {((x,?), (y, 1), (z, 2))}
    [z := z * y]4;
    ───────── {((x,?), (y, 1), (z, 2))}
    [y := y - 1]5
    ───────── {((x,?), (y, 1), (z, 2))}
od;
[y := 0]6 ───────── {((x,?), (y, 1), (z, 2))}

```

37

© Nielsen^2, Hankin

Reaching Definition

```

[y := x]1; ───────── {((x,?), (y,?), (z,?))}
[z := 1]2; ───────── {((x,?), (y, 1), (z,?))}
while [y > 0]3 do
    ───────── {((x,?), (y, 1), (y, 5), (z, 2), (z, 4))}
    [z := z * y]4;
    ───────── {((x,?), (y, 1), (y, 5), (z, 2), (z, 4))}
    [y := y - 1]5
    ───────── {((x,?), (y, 1), (y, 5), (z, 2), (z, 4))}
od;
[y := 0]6 ───────── {((x,?), (y, 1), (y, 5), (z, 2), (z, 4))}

```

38

© Nielsen^2, Hankin

The “best” solution

```

[y := x]1; ───────── {((x,?), (y,?), (z,?))}
[z := 1]2; ───────── {((x,?), (y, 1), (z,?))}
while [y > 0]3 do
    ───────── {((x,?), (y, 1), (y, 5), (z, 2), (z, 4))}
    [z := z * y]4;
    ───────── {((x,?), (y, 1), (y, 5), (z, 4))}
    [y := y - 1]5
    ───────── {((x,?), (y, 5), (z, 4))}
od;
[y := 0]6 ───────── {((x,?), (y, 6), (z, 2), (z, 4))}

```

39

© Nielsen^2, Hankin

A solution

```

[y := x]1; ───────── {((x,?), (y,?), (z,?))}
[z := 1]2; ───────── {((x,?), (y, 1), (z,?))}
while [y > 0]3 do
    ───────── {((x,?), (y, 1), (y, 5), (z, 2), (z, 4))}
    [z := z * y]4;
    ───────── {((x,?), (y, 1), (y, 5), (z, 2), (z, 4))}
    [y := y - 1]5
    ───────── {((x,?), (y, 1), (y, 5), (z, 2), (z, 4))}
od;
[y := 0]6 ───────── {((x,?), (y, 6), (z, 2), (z, 4))}

```

40

© Nielsen^2, Hankin

An unsound solution

```

[y := x]1;           {(x,?), (y,?), (z,?)}
[  

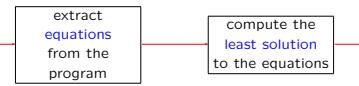
[z := 1]2;           {(x,?), (y, 1), (z,?)}
while [y > 0]3 do   {(x,?), (y, 1), (z, 2), (y, 5), (z, 4)}
    [z := z * y]4;   {(x,?), (y, 1), (y, 5), (z, 4)}
    [y := y - 1]5;   {(x,?), (y, 5), (z, 4)}
od;
[y := 0]6           {(x,?), (y, 6), (z, 2), (z, 4)}

```

41

© Nielsen^2, Hankin

Reaching Definitions



Analysis information:

- $RD_o(\ell)$: information available at the entry of block ℓ
- $RD_\bullet(\ell)$: information available at the exit of block ℓ

Two kinds of equations

$$\begin{array}{c} RD_o(\ell) \\ [x := a]^\ell \\ RD_\bullet(\ell) \end{array}$$

$$RD_o(\ell) \setminus \{(x, \ell') \mid \ell' \in \text{Lab}\} \cup \{(x, \ell)\} = RD_\bullet(\ell)$$

$$\begin{array}{c} [...]^{\ell_1} \\ RD_\bullet(\ell_1) \\ [...]^{\ell_2} \\ RD_o(\ell_2) \\ [...]^\ell \\ RD_\bullet(\ell_1) \cup RD_\bullet(\ell_2) = RD_o(\ell) \end{array}$$

42

© Nielsen^2, Hankin

Flow through assignments and tests

```

[y := x]1;           RD_\bullet(1) = RD_o(1) \ {(y, \ell) \mid \ell \in \text{Lab}} \cup \{(y, 1)\}
[z := 1]2;           RD_\bullet(2) = RD_o(2) \ {(z, \ell) \mid \ell \in \text{Lab}} \cup \{(z, 2)\}
while [y > 0]3 do   RD_\bullet(3) = RD_o(3)
    [z := z * y]4;   RD_\bullet(4) = RD_o(4) \ {(z, \ell) \mid \ell \in \text{Lab}} \cup \{(z, 4)\}
    [y := y - 1]5;   RD_\bullet(5) = RD_o(5) \ {(y, \ell) \mid \ell \in \text{Lab}} \cup \{(y, 5)\}
od;
[y := 0]6           RD_\bullet(6) = RD_o(6) \ {(y, \ell) \mid \ell \in \text{Lab}} \cup \{(y, 6)\}

```

 $\text{Lab} = \{1, 2, 3, 4, 5, 6\}$ 6 equations in
 $RD_o(1), \dots, RD_\bullet(6)$

43

© Nielsen^2, Hankin

Flow along the control

```

RD_o(1) = {(x,?), (y,?), (z,?)}
[y := x]1;           RD_o(2) = RD_\bullet(1)
[z := 1]2;           RD_o(3) = RD_\bullet(2) \cup RD_\bullet(5)
while [y > 0]3 do   RD_o(4) = RD_\bullet(3)
    [z := z * y]4;   RD_o(5) = RD_\bullet(4)
    [y := y - 1]5;   RD_o(6) = RD_\bullet(3)
od;
[y := 0]6           RD_o(6) = RD_\bullet(3)

```

 $\text{Lab} = \{1, 2, 3, 4, 5, 6\}$ 6 equations in
 $RD_o(1), \dots, RD_\bullet(6)$

44

© Nielsen^2, Hankin

Summary of equation system

$$\begin{aligned} RD_\bullet(1) &= RD_\circ(1) \setminus \{(y, \ell) \mid \ell \in \text{Lab}\} \cup \{(y, 1)\} \\ RD_\bullet(2) &= RD_\circ(2) \setminus \{(z, \ell) \mid \ell \in \text{Lab}\} \cup \{(z, 2)\} \\ RD_\bullet(3) &= RD_\circ(3) \\ RD_\bullet(4) &= RD_\circ(4) \setminus \{(z, \ell) \mid \ell \in \text{Lab}\} \cup \{(z, 4)\} \\ RD_\bullet(5) &= RD_\circ(5) \setminus \{(y, \ell) \mid \ell \in \text{Lab}\} \cup \{(y, 5)\} \\ RD_\bullet(6) &= RD_\circ(6) \setminus \{(y, \ell) \mid \ell \in \text{Lab}\} \cup \{(y, 6)\} \end{aligned}$$

$$\begin{aligned} RD_\circ(1) &= \{(x, ?), (y, ?), (z, ?)\} \\ RD_\circ(2) &= RD_\bullet(1) \\ RD_\circ(3) &= RD_\bullet(2) \cup RD_\bullet(5) \\ RD_\circ(4) &= RD_\bullet(3) \\ RD_\circ(5) &= RD_\bullet(4) \\ RD_\circ(6) &= RD_\bullet(3) \end{aligned}$$

45

© Nielsen^2, Hankin

- **12 sets:** $RD_\circ(1), \dots, RD_\bullet(6)$
all being subsets of $\text{Var} \times \text{Lab}$
- **12 equations:**
 $RD_j = F_j(RD_\circ(1), \dots, RD_\bullet(6))$
- **one function:**
 $F : \mathcal{P}(\text{Var} \times \text{Lab})^{12} \rightarrow \mathcal{P}(\text{Var} \times \text{Lab})^{12}$
- we want the **least fixed point** of F — this is the **best solution** to the equation system

A simple iterative algorithm

- **Initialisation**
 $RD_1 := \emptyset; \dots; RD_{12} := \emptyset;$

- **Iteration**
while $RD_j \neq F_j(RD_1, \dots, RD_{12})$ for some j
do
 $RD_j := F_j(RD_1, \dots, RD_{12})$

The algorithm terminates and computes the least fixed point of F .

46

© Nielsen^2, Hankin

The example equations

| RD _o | 1 | 2 | 3 | 4 | 5 | 6 |
|-----------------|---|---|---|---|---|---|
| 0 | ∅ | ∅ | ∅ | ∅ | ∅ | ∅ |

| RD _• | 1 | 2 | 3 | 4 | 5 | 6 |
|-----------------|---|---|---|---|---|---|
| 0 | ∅ | ∅ | ∅ | ∅ | ∅ | ∅ |

The equations:

$$\begin{aligned} RD_\bullet(1) &= RD_\circ(1) \setminus \{(y, \ell) \mid \dots\} \cup \{(y, 1)\} \\ RD_\bullet(2) &= RD_\circ(2) \setminus \{(z, \ell) \mid \dots\} \cup \{(z, 2)\} \\ RD_\bullet(3) &= RD_\circ(3) \\ RD_\bullet(4) &= RD_\circ(4) \setminus \{(z, \ell) \mid \dots\} \cup \{(z, 4)\} \\ RD_\bullet(5) &= RD_\circ(5) \setminus \{(y, \ell) \mid \dots\} \cup \{(y, 5)\} \\ RD_\bullet(6) &= RD_\circ(6) \setminus \{(y, \ell) \mid \dots\} \cup \{(y, 6)\} \\ \\ RD_\circ(1) &= \{(x, ?), (y, ?), (z, ?)\} \\ RD_\circ(2) &= RD_\bullet(1) \\ RD_\circ(3) &= RD_\bullet(2) \cup RD_\bullet(5) \\ RD_\circ(4) &= RD_\bullet(3) \\ RD_\circ(5) &= RD_\bullet(4) \\ RD_\circ(6) &= RD_\bullet(3) \end{aligned}$$

47

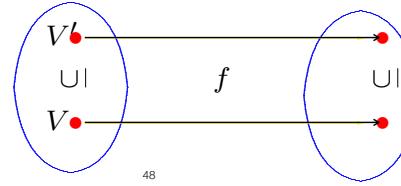
© Nielsen^2, Hankin

Why does it work? (1)

A function $f : \mathcal{P}(S) \rightarrow \mathcal{P}(S)$ is a **monotone function** if

$$V \subseteq V' \Rightarrow f(V) \subseteq f(V')$$

(the larger the argument – the larger the result)



48

© Nielsen^2, Hankin

Why does it work? (2)

A set L equipped with an ordering \subseteq satisfies the **Ascending Chain Condition** if all chains

$$V_0 \subseteq V_1 \subseteq V_2 \subseteq V_3 \subseteq \dots$$

stabilise, that is, if there exists some n such that $V_n = V_{n+1} = V_{n+2} = \dots$

If S is a **finite** set then $\mathcal{P}(S)$ equipped with the subset ordering \subseteq satisfies the Ascending Chain Condition — the chains cannot grow forever since each element is a subset of a finite set.

Fact

For a given program $\text{Var} \times \text{Lab}$ will be a finite set so $\mathcal{P}(\text{Var} \times \text{Lab})$ with the subset ordering satisfies the Ascending Chain Condition.

49

© Nielsen^2, Hankin

Correctness of the algorithm

Initialisation

$\text{RD}_1 := \emptyset; \dots; \text{RD}_{12} := \emptyset;$

Invariant: $\vec{\text{RD}} \subseteq F^n(\vec{\emptyset})$ since $\vec{\text{RD}} = \vec{\emptyset}$ is the least element

Iteration

while $\text{RD}_j \neq F_j(\text{RD}_1, \dots, \text{RD}_{12})$ for some j

do assume $\vec{\text{RD}}$ is $\vec{\text{RD}'}$ and $\vec{\text{RD}'} \subseteq F^n(\vec{\emptyset})$

$\text{RD}_j := F_j(\text{RD}_1, \dots, \text{RD}_{12})$

then $\vec{\text{RD}} \subseteq F(\vec{\text{RD}'}) \subseteq F^{n+1}(\vec{\emptyset}) = F^n(\vec{\emptyset})$ when $\text{lfp}(F) = F^n(\vec{\emptyset})$

If the algorithm terminates then it computes the least fixed point of F .

The algorithm terminates because $\text{RD}_j \subset F_j(\text{RD}_1, \dots, \text{RD}_{12})$ is only possible finitely many times since $\mathcal{P}(\text{Var} \times \text{Lab})^{12}$ satisfies the Ascending Chain Condition.

51

© Nielsen^2, Hankin

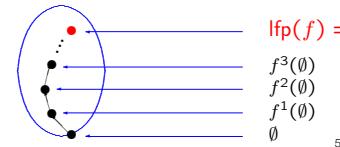
Why does it work? (3)

Let $f : \mathcal{P}(S) \rightarrow \mathcal{P}(S)$ be a **monotone function**. Then

$$\emptyset \subseteq f(\emptyset) \subseteq f^2(\emptyset) \subseteq f^3(\emptyset) \subseteq \dots$$

Assume that S is a finite set; then the **Ascending Chain Condition** is satisfied. This means that the chain cannot be growing infinitely so there exists n such that $f^n(\emptyset) = f^{n+1}(\emptyset) = \dots$

$f^n(\emptyset)$ is the **least fixed point** of f



$\text{lfp}(f) = f^n(\emptyset) = f^{n+1}(\emptyset) \text{ for some } n$

50

© Nielsen^2, Hankin

Intermezzo Static Analysis at Google

contributed articles



Two-billion line of code

Single Git repository

Bazel build system

C++, Java, Python, Go...

Programmers don't use static analysis because:

Not integrated into the developer's workflow

Not actionable.

Not trustworthy due to false positives

Not manifest in practice.

Too expensive to fix.

Warnings not understood.

Sadowski, Attardian, Eagle, Miller-Cushman, Jaspan. Lessons from Building Static Analysis Tools at Google. CACM'18. DOI: 10.1145/3188720

52

Intermezzo Static Analysis at Google

Lessons from FindBugs: a “failure” in three acts

- Act 1: **Bug dashboard** (2006) FindBugs as centralized tool, nightly runs over the entire Google Java codebase, produced database of bugs. Unused because dashboard outside developers’ workflow, and hard to find “new” bugs
- Act 2. **Filing bugs.** (2008-2009) Manually triage of issues found by nightly runs & filing bug reports for important ones. In May 2009, 100s of engineers did a “Fixit” week. Reviewed 3,954 warnings, only 16% (640) were fixed, despite that 44% (1,746) resulted in bug reports. Many bug not important enough and manual triage does not scale
- Act 3. **Code review integration.** (2010-2011) Run when changes sent for review, posting results as comments on code-review thread. The tooling attempted to show only new warnings but sometimes miscategorized issues as new. Discontinued when the code-review tool was replaced because: false positives caused loss of confidence



David Hovemeyer



Bill Pugh

C Sadowski, E Attandlian, A Eagle, L Miller-Cushon, C Jaspan. Lessons from Building Static Analysis Tools at Google. CACM'18. DOI: 10.1145/3188720

53

Intermezzo Static Analysis at Google

Things that worked...

- Thing #1: **Compiler workflow.** (1) Implement new compiler checks, (2) run updated compiler over entire Google codebase, (3) programmatically fix all problems in the codebase, (4) once the codebase was cleansed of an issue, the new diagnostic becomes a compiler error and offer an automatic patch. Survey show 57% of developers like the patches, only 2% unhappy.
- Thing #2. **Code review checks.** Integration into Google’s code review tool, Critique; static analysis results exposed using Tricorder, Google’s program-analysis platform. A code review check should:
Be understandable. Be actionable and easy to fix. The fix may require more time than a compiler check. Produce less than 10% effective false positives. Have the potential for significant impact on code quality.

C Sadowski, E Attandlian, A Eagle, L Miller-Cushon, C Jaspan. Lessons from Building Static Analysis Tools at Google. CACM'18. DOI: 10.1145/3188720

54

Intermezzo Static Analysis at Google

No interprocedural or whole-program analysis

The types of analyses deployed:

- **Style checkers:** Checkstyle, Pylint, and Golint...
- **Bug-finding tools that extend compilers:** Error Prone, ClangTidy, Clang Thread Safety Analysis, Gcov, and the Checker Framework..., abstract-syntax-tree pattern-match tools, type-based checks, and unused variable analysis
- **Analyzers that make calls to production services:** e.g. check if an employee mentioned in a comment is still at Google
- **Analyzers that examine properties of build outputs:** e.g. the size of binaries

C Sadowski, E Attandlian, A Eagle, L Miller-Cushon, C Jaspan. Lessons from Building Static Analysis Tools at Google. CACM'18. DOI: 10.1145/3188720

55

Constraint-Based Analysis

Example: Control Flow Analysis

```
let f = fn x => x 7
      g = fn y => y
      h = fn z => 3
in f g + f (g h)
      ↓      ↓
      g 7    f h
            ↓
            h 7
```

Aim: For each function application, which function abstractions may be applied?

| function applications | function abstractions that may be applied |
|-----------------------|---|
| x 7 | g, h |
| f g | f |
| g h | g |
| f (g h) | f |

56

© Nelson^2, Hankin

An application of control flow analysis

```
let f = fn x => x 7
  g = fn y => y
  h = fn z => 3
in f g + f (g h)
```

Partial evaluation of function call:

```
let f = fn x => case x of
  g: 7
  h: 3
  g = fn y => y
  h = fn z => 3
in f g + f (g h)
```

Aim: For each function application, which function abstractions may be applied?

| function applications | function abstractions that may be applied |
|-----------------------|---|
| x 7 | g, h |
| f g | f |
| g h | g |
| f (g h) | f |

57

© Nielsen^2, Hankin

The underlying analysis problem

```
let f = fn x => x 7
```

```
  g = fn y => y
```

```
  h = fn z => 3
```

```
in f g + f (g h)
```

Aim: for each **function application**, which function abstractions may be applied?

The analysis will compute:

- for each **subexpression**, which function abstractions may it denote?
 - e.g. **(g h)** may evaluate to **h** introduce abstract cache **C**
- for each **variable**, which function abstractions may it denote?
 - e.g. **x** may be **g** or **h** introduce abstract environment **R**

58

© Nielsen^2, Hankin

The best solution to the analysis problem

Add labels to subexpressions:

```
let f = fn x => (x1 72)3
  g = fn y => y4
  h = fn z => 35
in (f6 g7)8 + (f9 (g10 h11)12)13
```

| R | variable may be bound to |
|---|--|
| x | {fn y => y ⁴ , fn z => 3 ⁵ } |
| y | {fn z => 3 ⁵ } |
| z | Ø |
| f | {fn x => (x ¹ 7 ²) ³ } |
| g | {fn y => y ⁴ } |
| h | {fn z => 3 ⁵ } |

59

© Nielsen^2, Hankin

| C | subexpression may evaluate to |
|----|--|
| 1 | {fn y => y ⁴ , fn z => 3 ⁵ } |
| 2 | Ø |
| 3 | Ø |
| 4 | {fn z => 3 ⁵ } |
| 5 | Ø |
| 6 | {fn x => (x ¹ 7 ²) ³ } |
| 7 | {fn y => y ⁴ } |
| 8 | Ø |
| 9 | {fn x => (x ¹ 7 ²) ³ } |
| 10 | {fn y => y ⁴ } |
| 11 | {fn z => 3 ⁵ } |
| 12 | {fn z => 3 ⁵ } |
| 13 | Ø |

Constraint-Based Analysis

How to automate the analysis



Analysis information:

- **R(x)**: information available for the **variable x**
- **C(ℓ)**: information available at the **subexpression labelled ℓ**

Three kinds of constraints

- let-bound variables evaluate to their abstraction
- variables evaluate to their (abstract) values
- if a function abstraction is applied to an argument then
 - the argument is a possible value of the formal parameter
 - the value of the body of the abstraction is a possible value of the application

60

© Nielsen^2, Hankin

let-bound variables

let-bound variables evaluate to their abstractions

`let f = fn x => e` gives rise to the constraint $\{fn\ x\ =>\ e\} \subseteq R(f)$

$$\begin{array}{lcl} \text{let } f = fn\ x \Rightarrow (x^1\ 7^2)^3 & \longleftarrow & \{fn\ x \Rightarrow (x^1\ 7^2)^3\} \subseteq R(f) \\ g = fn\ y \Rightarrow y^4 & \longleftarrow & \{fn\ y \Rightarrow y^4\} \subseteq R(g) \\ \text{in } (f^5\ g^6)^7 & & \end{array}$$

61

© Nielsen^2, Hankin

Variables

Variables evaluate to their abstract values

x^ℓ gives rise to the constraint $R(x) \subseteq C(\ell)$

$$\begin{array}{lcl} \text{let } f = fn\ x \Rightarrow (x^1\ 7^2)^3 & \longleftarrow & R(x) \subseteq C(1) \\ g = fn\ y \Rightarrow y^4 & \longleftarrow & R(y) \subseteq C(4) \\ \text{in } (f^5\ g^6)^7 & \longleftarrow & \begin{cases} R(f) \subseteq C(5) \\ R(g) \subseteq C(6) \end{cases} \end{array}$$

62

© Nielsen^2, Hankin

Function application (1)

if a function abstraction is applied to an argument then

- the argument is a possible value of the formal parameter
- the value of the body of the abstraction is a possible value of the application

$$\begin{array}{lcl} \text{let } f = fn\ x \Rightarrow (x^1\ 7^2)^3 & \longleftarrow & \text{if } (fn\ y \Rightarrow y^4) \in C(1) \\ & & \text{then } C(2) \subseteq R(y) \text{ and } C(4) \subseteq C(3) \\ g = fn\ y \Rightarrow y^4 & & \text{if } (fn\ x \Rightarrow (x^1\ 7^2)^3) \in C(1) \\ \text{in } (f^5\ g^6)^7 & & \text{then } C(2) \subseteq R(x) \text{ and } C(3) \subseteq C(3) \\ & & \text{if } (fn\ y \Rightarrow y^4) \in C(5) \\ & & \text{then } C(6) \subseteq R(y) \text{ and } C(4) \subseteq C(7) \\ & & \text{if } (fn\ x \Rightarrow (x^1\ 7^2)^3) \in C(5) \\ & & \text{then } C(6) \subseteq R(x) \text{ and } C(3) \subseteq C(7) \end{array}$$

Conditional constraints

© Nielsen^2, Hankin

Summary of constraint system

$$\begin{aligned} & \{fn\ x \Rightarrow (x^1\ 7^2)^3\} \subseteq R(f) \\ & \{fn\ y \Rightarrow y^4\} \subseteq R(g) \\ & R(x) \subseteq C(1) \\ & R(y) \subseteq C(4) \\ & R(f) \subseteq C(5) \\ & R(g) \subseteq C(6) \\ & (fn\ y \Rightarrow y^4) \in C(1) \Rightarrow C(2) \subseteq R(y) \\ & (fn\ y \Rightarrow y^4) \in C(1) \Rightarrow C(4) \subseteq C(3) \\ & (fn\ x \Rightarrow (x^1\ 7^2)^3) \in C(1) \Rightarrow C(2) \subseteq R(x) \\ & (fn\ x \Rightarrow (x^1\ 7^2)^3) \in C(1) \Rightarrow C(3) \subseteq C(3) \\ & (fn\ y \Rightarrow y^4) \in C(5) \Rightarrow C(6) \subseteq R(y) \\ & (fn\ y \Rightarrow y^4) \in C(5) \Rightarrow C(4) \subseteq C(7) \\ & (fn\ x \Rightarrow (x^1\ 7^2)^3) \in C(5) \Rightarrow C(6) \subseteq R(x) \\ & (fn\ x \Rightarrow (x^1\ 7^2)^3) \in C(5) \Rightarrow C(3) \subseteq C(7) \end{aligned}$$

```
let f = fn x => (x^1 7^2)^3
g = fn y => y^4
in (f^5 g^6)^7
```

- 11 sets: $R(x), R(y), R(f), R(g), C(1), \dots, C(7)$; all being subsets of the set **Abstr** of function abstractions
- the constraints can be reformulated as a function: $F : \mathcal{P}(\text{Abstr})^{11} \rightarrow \mathcal{P}(\text{Abstr})^{11}$
- we want the least fixed point of F — this is the best solution to the constraint system

$\mathcal{P}(S)$ is the set of all subsets of the set S ; e.g. $\mathcal{P}(\{0, 1\}) = \{\emptyset, \{0\}, \{1\}, \{0, 1\}\}$.

© Nielsen^2, Hankin

Constraint-Based Analysis

The constraints define a function

$$\begin{aligned}
 & \{ \text{fn } x \Rightarrow (x^1 7^2)^3 \} \subseteq R(f) \\
 & \{ \text{fn } y \Rightarrow y^4 \} \subseteq R(g) \\
 & R(x) \subseteq C(1) \\
 & R(y) \subseteq C(4) \\
 & R(f) \subseteq C(5) \\
 & R(g) \subseteq C(6) \\
 & (\text{fn } y \Rightarrow y^4) \in C(1) \Rightarrow C(2) \subseteq R(y) \\
 & (\text{fn } y \Rightarrow y^4) \in C(1) \Rightarrow C(4) \subseteq C(3) \\
 & (\text{fn } x \Rightarrow (x^1 7^2)^3) \in C(1) \Rightarrow C(2) \subseteq R(x) \\
 & (\text{fn } x \Rightarrow (x^1 7^2)^3) \in C(1) \Rightarrow C(3) \subseteq C(3) \\
 & (\text{fn } y \Rightarrow y^4) \in C(5) \Rightarrow C(6) \subseteq R(y) \\
 & (\text{fn } y \Rightarrow y^4) \in C(5) \Rightarrow C(4) \subseteq C(7) \\
 & (\text{fn } x \Rightarrow (x^1 7^2)^3) \in C(5) \Rightarrow C(6) \subseteq R(x) \\
 & (\text{fn } x \Rightarrow (x^1 7^2)^3) \in C(5) \Rightarrow C(3) \subseteq C(7)
 \end{aligned}$$

65

© Nielsen^2, Hankin

$F : \mathcal{P}(\text{Abstr})^{11} \rightarrow \mathcal{P}(\text{Abstr})^{11}$
is defined by

$$\begin{aligned}
 F_{R(f)}(\dots, R_f, \dots) &= \emptyset \\
 F_{R_f}(R_x, \dots, C_1, \dots) &= C_1 \cup R_x \\
 F_{C(1)}(R_x, \dots, C_1, \dots) &= C_1 \cup R_x \\
 F_{R(y)}(\dots, R_y, \dots, C_1, C_2, \dots, C_5, C_6, \dots) &= \\
 R_y \cup \{a \in C_2 \mid \text{fn } y \Rightarrow y^4 \in C_1\} & \\
 \cup \{a \in C_6 \mid \text{fn } y \Rightarrow y^4 \in C_5\} &
 \end{aligned}$$

Constraint-Based Analysis

How to solve the constraints

- Initialisation

$$X_1 := \emptyset; \dots; X_{11} := \emptyset;$$

- Iteration

while $X_j \neq F_{X_j}(X_1, \dots, X_{11})$ for some j
do
 $X_j := F_{X_j}(X_1, \dots, X_{11})$

The algorithm terminates and computes the least fixed point of F

writing
 X_1, \dots, X_{11} for
 $R(x), \dots, R(g), C(1), \dots, C(7)$

In practice we propagate smaller contributions than F_{X_i} , e.g. a constraint at a time.

66

© Nielsen^2, Hankin

Constraint-Based Analysis

The example constraint system

| R | x | y | f | g |
|---|--------------------------------|---|--------------------------------|--------------------------------|
| 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | $\text{fn } x \Rightarrow .^3$ | 0 |
| 2 | 0 | 0 | $\text{fn } x \Rightarrow .^3$ | $\text{fn } y \Rightarrow .^4$ |
| 3 | 0 | 0 | $\text{fn } x \Rightarrow .^3$ | $\text{fn } y \Rightarrow .^4$ |
| 4 | 0 | 0 | $\text{fn } x \Rightarrow .^3$ | $\text{fn } y \Rightarrow .^4$ |
| 5 | $\text{fn } y \Rightarrow .^4$ | 0 | $\text{fn } x \Rightarrow .^3$ | $\text{fn } y \Rightarrow .^4$ |
| 6 | $\text{fn } y \Rightarrow .^4$ | 0 | $\text{fn } x \Rightarrow .^3$ | $\text{fn } y \Rightarrow .^4$ |

| C | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|--------------------------------|---|---|---|--------------------------------|--------------------------------|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | $\text{fn } x \Rightarrow .^3$ | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | $\text{fn } x \Rightarrow .^3$ | $\text{fn } y \Rightarrow .^4$ | 0 |
| 5 | 0 | 0 | 0 | 0 | $\text{fn } x \Rightarrow .^3$ | $\text{fn } y \Rightarrow .^4$ | 0 |
| 6 | $\text{fn } y \Rightarrow .^4$ | 0 | 0 | 0 | $\text{fn } x \Rightarrow .^3$ | $\text{fn } y \Rightarrow .^4$ | 0 |

67

The constraints:

$$\begin{aligned}
 & \{ \text{fn } x \Rightarrow .^3 \} \subseteq R(f) \quad (1) \\
 & \{ \text{fn } y \Rightarrow .^4 \} \subseteq R(g) \quad (2) \\
 & R(x) \subseteq C(1) \quad (6) \\
 & R(y) \subseteq C(4) \\
 & R(f) \subseteq C(5) \quad (3) \\
 & R(g) \subseteq C(6) \quad (4) \\
 & (\text{fn } y \Rightarrow .^4) \in C(1) \Rightarrow C(2) \subseteq R(y) \\
 & (\text{fn } y \Rightarrow .^4) \in C(1) \Rightarrow C(4) \subseteq C(3) \\
 & (\text{fn } x \Rightarrow .^3) \in C(1) \Rightarrow C(2) \subseteq R(x) \\
 & (\text{fn } x \Rightarrow .^3) \in C(1) \Rightarrow C(3) \subseteq C(3) \\
 & (\text{fn } y \Rightarrow .^4) \in C(5) \Rightarrow C(6) \subseteq R(y) \\
 & (\text{fn } y \Rightarrow .^4) \in C(5) \Rightarrow C(4) \subseteq C(7) \\
 & (\text{fn } x \Rightarrow .^3) \in C(5) \Rightarrow C(6) \subseteq R(x) \quad (5) \\
 & (\text{fn } x \Rightarrow .^3) \in C(5) \Rightarrow C(3) \subseteq C(7)
 \end{aligned}$$

© Nielsen^2, Hankin

Intermezzo In Defence of Soundness

An argument in favor of $2 + 2 = 5$

Analyses are expected to be sound, i.e. their result models all possible executions of the target program, or an over-approximation to stay tractable.

In practice, soundness is commonly eschewed: all realistic whole-program analysis tool purposely make unsound choices. Most/all published whole-program analyses omit conservative handling of language features of real languages.

Some specific features are best under-approximated by pretending that perfectly possible behaviors cannot happen. E.g. highly dynamic language constructs, such as eval, may be assumed to do nothing

We term such analyses as soundy. The concept of soundness attempts to capture the balance of over-approximated handling of most language features, yet deliberately under-approximated handling of a feature subset

Livshits, Sridharan, Smaragdakis, Lhoták, Amaral, Evan Chang, Guyer, Khedker, Moller, Vardoulakis. In Defense of Soundness: A Manifesto. CACM'15. DOI: 10.1145/2644805

68

Intermezzo In Defence of Soundness

An argument in favor of $2 + 2 = 5$

Some of the sources of unsoundness, sorted by language.

| Language | Examples of commonly ignored features | Consequences of not modeling these features |
|------------|--|--|
| C/C++ | setjmp/longjmp ignored effects of pointer arithmetic "manufactured" pointers | ignores arbitrary side effects to the program heap |
| Java/C# | Reflection JNI | can render much of the codebase invisible for analysis "invisible" code may create invisible side effects in programs |
| JavaScript | eval, dynamic code loading data flow through the DOM | missing execution missing data flow in program |

Read: "A True Positives Theorem for a Static Race Detector" POPL'19.

"We prove a True Positives Theorem stating that under certain assumptions, which reflect the way that product code uses concurrency, an idealized theoretical version of the analysis never reports a false positive.

The theorem was motivated by the desire to understand the observation that our tool was providing remarkably accurate signal to developers.

Technically, our result can be seen as saying that the analysis computes an under-approximation of an over-approximation, which is the reverse of the more usual situation in static analysis."

$2+2=4$ under some conditions.

69

Abstract Interpretation



- We have the analysis **old**: it has already been proved **correct** but it is **inefficient**, or maybe even uncomputable

- We want the analysis **new**: it has to be **correct** as well as **efficient!**

- Can we develop **new** from **old**?

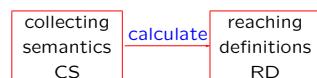
abstract interpretation !

70

© Nielsen^2, Hankin

Abstract Interpretation

Example: Collecting Semantics and Reaching Definitions



The **collecting semantics** CS

- collects the set of traces that can reach a given program point
- has an easy correctness proof
- is uncomputable

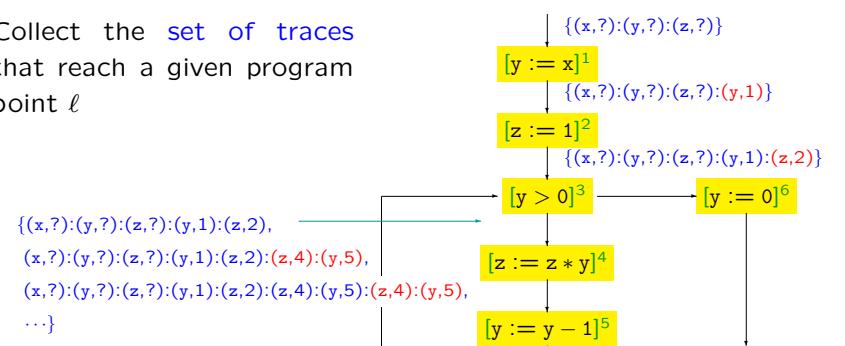
The **reaching definitions analysis** RD is as before

71

© Nielsen^2, Hankin

Example: Collecting Semantics

Collect the **set of traces** that reach a given program point ℓ



72

© Nielsen^2, Hankin

Abstract Interpretation

How to proceed

As before:

- extract a set of equations defining the possible sets of traces
- compute the least fixed point of the set of equations

And furthermore:

- prove the correctness: the set of traces computed by the analysis is a superset of the possible traces

73

© Nielsen^2, Hankin

Abstract Interpretation

Two kinds of equations

$$\begin{array}{c}
 \downarrow \text{CS}_o(\ell) \\
 [\underline{x} := a]^{\ell} \\
 \downarrow \text{CS}_o(\ell) \\
 \left\{ \text{trace} : (\underline{x}, \ell) \mid \text{trace} \in \text{CS}_o(\ell) \right\} = \text{CS}_o(\ell)
 \end{array}
 \quad \quad
 \begin{array}{c}
 \text{CS}_o(\ell_1) \quad \text{CS}_o(\ell_2) \\
 \swarrow \quad \searrow \\
 [...]^{\ell_1} \quad [...]^{\ell_2} \\
 \downarrow \text{CS}_o(\ell) \\
 [...]^{\ell}
 \end{array}
 \quad \quad
 \text{CS}_o(\ell_1) \cup \text{CS}_o(\ell_2) = \text{CS}_o(\ell)$$

74

© Nielsen^2, Hankin

Flow through assignments and tests

```

[y := x]1;           CS•(1) = {trace : (y, 1) | trace ∈ CSo(1)}
[z := 1]2;           CS•(2) = {trace : (z, 2) | trace ∈ CSo(2)}
while [y > 0]3 do
  CS•(3) = CSo(3)
  [z := z * y]4;
  CS•(4) = {trace : (z, 4) | trace ∈ CSo(4)}
  [y := y - 1]5
  CS•(5) = {trace : (y, 5) | trace ∈ CSo(5)}
od;
[y := 0]6             CS•(6) = {trace : (y, 6) | trace ∈ CSo(6)}

```

6 equations in
CS_o(1), ..., CS_o(6)

75

© Nielsen^2, Hankin

Flow along the control

```

[y := x]1;           CSo(1) = {(x, ?) : (y, ?) : (z, ?)}
[z := 1]2;           CSo(2) = CS•(1)
while [y > 0]3 do
  CSo(3) = CS•(2) ∪ CSo(5)
  [z := z * y]4;
  CSo(4) = CS•(3)
  [y := y - 1]5
  CSo(5) = CS•(4)
od;
[y := 0]6             CSo(6) = CS•(3)

```

6 equations in
CS_o(1), ..., CS_o(6)

76

© Nielsen^2, Hankin

Summary of Collecting Semantics

$$\begin{aligned} \text{CS}_\bullet(1) &= \{\text{trace} : (y, 1) \mid \text{trace} \in \text{CS}_\circ(1)\} \\ \text{CS}_\bullet(2) &= \{\text{trace} : (z, 2) \mid \text{trace} \in \text{CS}_\circ(2)\} \\ \text{CS}_\bullet(3) &= \text{CS}_\circ(3) \\ \text{CS}_\bullet(4) &= \{\text{trace} : (z, 4) \mid \text{trace} \in \text{CS}_\circ(4)\} \\ \text{CS}_\bullet(5) &= \{\text{trace} : (y, 5) \mid \text{trace} \in \text{CS}_\circ(5)\} \\ \text{CS}_\bullet(6) &= \{\text{trace} : (y, 6) \mid \text{trace} \in \text{CS}_\circ(6)\} \end{aligned}$$

$$\begin{aligned} \text{CS}_\circ(1) &= \{(x, ?) : (y, ?) : (z, ?)\} \\ \text{CS}_\circ(2) &= \text{CS}_\bullet(1) \\ \text{CS}_\circ(3) &= \text{CS}_\bullet(2) \cup \text{CS}_\bullet(5) \\ \text{CS}_\circ(4) &= \text{CS}_\bullet(3) \\ \text{CS}_\circ(5) &= \text{CS}_\bullet(4) \\ \text{CS}_\circ(6) &= \text{CS}_\bullet(3) \end{aligned}$$

- **12 sets:** $\text{CS}_\circ(1), \dots, \text{CS}_\bullet(6)$
all being subsets of Trace
- **12 equations:**
 $\text{CS}_j = G_j(\text{CS}_\circ(1), \dots, \text{CS}_\bullet(6))$
- **one function:**
 $G : \mathcal{P}(\text{Trace})^{12} \rightarrow \mathcal{P}(\text{Trace})^{12}$
- we want the **least fixed point** of G — but it is **uncomputable!**

77

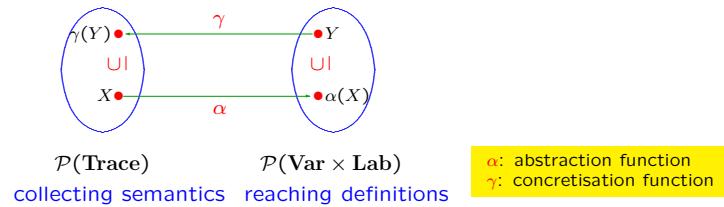
© Nielsen^2, Hankin

Example: Inducing an analysis

Galois Connections

A Galois connection between two sets is a pair of (α, γ) of functions between the sets satisfying

$$X \subseteq \gamma(Y) \Leftrightarrow \alpha(X) \subseteq Y$$



78

© Nielsen^2, Hankin

Semantically Reaching Definitions

For a single trace:

$$\begin{aligned} \text{trace: } & (x, ?) : (y, ?) : (z, ?) : (y, 1) : (z, 2) \\ \text{SRD(trace): } & \{(x, ?), (y, 1), (z, 2)\} \end{aligned}$$

For a set of traces:

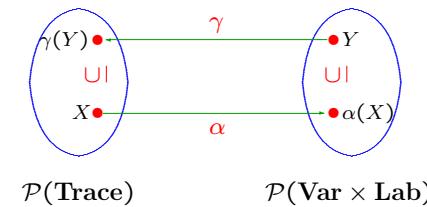
$$\begin{aligned} X \in \mathcal{P}(\text{Trace}): & \{(x, ?) : (y, ?) : (z, ?) : (y, 1) : (z, 2), \\ & (x, ?) : (y, ?) : (z, ?) : (y, 1) : (z, 2) : (z, 4) : (y, 5)\} \\ \text{SRD}(X): & \{(x, ?), (y, 1), (z, 2), (z, 4), (y, 5)\} \end{aligned}$$

79

© Nielsen^2, Hankin

Galois connection for Reaching Definitions analysis

$$\begin{aligned} \alpha(X) &= \text{SRD}(X) \\ \gamma(Y) &= \{\text{trace} \mid \text{SRD}(\text{trace}) \subseteq Y\} \end{aligned}$$



Galois connection:

$$X \subseteq \gamma(Y) \Leftrightarrow \alpha(X) \subseteq Y$$

© Nielsen^2, Hankin

Abstract Interpretation

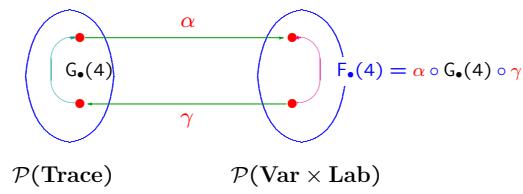
Inducing the Reaching Definitions analysis (1)

Known:

- $G_\bullet(4)$ defined on $\mathcal{P}(\text{Trace})$
- the Galois connection (α, γ)

Calculate:

- $F_\bullet(4)$ defined on $\mathcal{P}(\text{Var} \times \text{Lab})$
as $F_\bullet(4) = \alpha \circ G_\bullet(4) \circ \gamma$



$\mathcal{P}(\text{Trace})$

$\mathcal{P}(\text{Var} \times \text{Lab})$

81

© Nielsen^2, Hankin

Inducing the Reaching Definitions analysis (2)

$$\begin{aligned}
 \text{RD}_\bullet(4) &= F_\bullet(4)(\dots, \text{RD}_\circ(4), \dots) \\
 &= \alpha(G_\bullet(4)(\gamma(\dots, \text{RD}_\circ(4), \dots))) \text{ using } F_\bullet(4) = \alpha \circ G_\bullet(4) \circ \gamma \\
 &= \alpha(\{tr : (z, 4) \mid tr \in \gamma(\text{RD}_\circ(4))\}) \\
 &\quad \text{using } G_\bullet(4)(\dots, \text{CS}_\circ(4), \dots) = \{tr : (z, 4) \mid tr \in \text{CS}_\circ(4)\} \\
 &= \text{SRD}(\{tr : (z, 4) \mid tr \in \gamma(\text{RD}_\circ(4))\}) \quad \text{using } \alpha = \text{SRD} \\
 &= (\text{SRD}(\{tr \mid tr \in \gamma(\text{RD}_\circ(4))\}) \setminus \{(z, \ell) \mid \ell \in \text{Lab}\}) \cup \{(z, 4)\} \\
 &= (\alpha(\gamma(\text{RD}_\circ(4)) \setminus \{(z, \ell) \mid \ell \in \text{Lab}\}) \cup \{(z, 4)\}) \text{ using } \alpha = \text{SRD} \\
 &= (\text{RD}_\circ(4) \setminus \{(z, \ell) \mid \ell \in \text{Lab}\}) \cup \{(z, 4)\} \quad \text{using } \alpha \circ \gamma = id
 \end{aligned}$$

just as before!

82

© Nielsen^2, Hankin

Scaling Program Analysis

CS7575

- Lecture 2 — Data flow analysis; PoPA Chapter 2 (part 1)

83

Data flow analysis

Characterizes the *flow* of analysis facts through a program

Analyses differ in aspects such as:

- dependence on statement ordering wrt control flow operations
flow-sensitive vs. *flow-insensitive*
- direction of flow
forward vs. *backward*
- qualification over paths
may vs. *must*
- scope
interprocedural vs *intraprocedural*

84

© Thomas Knoll

Intraprocedural data flow analysis

Classical analyses:

- Available Expression
- Reaching Definition
- Very Busy Expression
- Live Variables

Derived analyses:

- Use-definition
- Definition-use

85

© Thomas Knoll

© Nielsen^2, Hankin

While language

Abstract syntax

$$\begin{aligned} a &::= x \mid n \mid a_1 \text{ op}_a a_2 \\ b &::= \text{true} \mid \text{false} \mid \text{not } b \mid b_1 \text{ op}_b b_2 \mid a_1 \text{ op}_r a_2 \\ S &::= [x := a]^\ell \mid [\text{skip}]^\ell \mid S_1; S_2 \mid \\ &\quad \text{if } [b]^\ell \text{ then } S_1 \text{ else } S_2 \mid \text{while } [b]^\ell \text{ do } S \end{aligned}$$

Example $[z := 1]^1; \text{while } [x > 0]^2 \text{ do } ([z := z * y]^3; [x := x - 1]^4)$

86

© Nielsen^2, Hankin

$[z := 1]^1; \text{while } [x > 0]^2 \text{ do } ([z := z * y]^3; [x := x - 1]^4)$

```

init(…) = 1
final(…) = {2}
labels(…) = {1, 2, 3, 4}
flow(…) = {(1, 2), (2, 3),
            (3, 4), (4, 2)}
flowRT(…)

```

87

© Nielsen^2, Hankin

$\text{init}(S)$ is the label of the first elementary block of S :

Helper functions

$\text{init} : \text{Stmt} \rightarrow \text{Lab}$

$$\begin{aligned} \text{init}([x := a]^\ell) &= \ell \\ \text{init}([\text{skip}]^\ell) &= \ell \\ \text{init}(S_1; S_2) &= \text{init}(S_1) \\ \text{init}(\text{if } [b]^\ell \text{ then } S_1 \text{ else } S_2) &= \ell \\ \text{init}(\text{while } [b]^\ell \text{ do } S) &= \ell \end{aligned}$$

$\text{final}(S)$ is the set of labels of the last elementary blocks of S :

$\text{final} : \text{Stmt} \rightarrow \mathcal{P}(\text{Lab})$

$$\begin{aligned} \text{final}([x := a]^\ell) &= \{\ell\} \\ \text{final}([\text{skip}]^\ell) &= \{\ell\} \\ \text{final}(S_1; S_2) &= \text{final}(S_2) \\ \text{final}(\text{if } [b]^\ell \text{ then } S_1 \text{ else } S_2) &= \text{final}(S_1) \cup \text{final}(S_2) \\ \text{final}(\text{while } [b]^\ell \text{ do } S) &= \{\ell\} \end{aligned}$$

© Nielsen^2, Hankin

Helper functions

$\text{flow}(S)$ and $\text{flow}^R(S)$ are representations of how control flows in S :

$$\text{flow}, \text{flow}^R : \text{Stmt} \rightarrow \mathcal{P}(\text{Lab} \times \text{Lab})$$

$$\begin{aligned}\text{flow}([x := a]^\ell) &= \emptyset \\ \text{flow}([\text{skip}]^\ell) &= \emptyset \\ \text{flow}(S_1; S_2) &= \text{flow}(S_1) \cup \text{flow}(S_2) \\ &\cup \{(\ell, \text{init}(S_2)) \mid \ell \in \text{final}(S_1)\} \\ \text{flow}(\text{if } [b]^\ell \text{ then } S_1 \text{ else } S_2) &= \text{flow}(S_1) \cup \text{flow}(S_2) \\ &\cup \{(\ell, \text{init}(S_1)), (\ell, \text{init}(S_2))\} \\ \text{flow}(\text{while } [b]^\ell \text{ do } S) &= \text{flow}(S) \cup \{(\ell, \text{init}(S))\} \\ &\cup \{(\ell', \ell) \mid \ell' \in \text{final}(S)\} \\ \text{flow}^R(S) &= \{(\ell, \ell') \mid (\ell', \ell) \in \text{flow}(S)\}\end{aligned}$$

89

© Nielsen^2, Hankin

Helper functions

A statement consists of a set of *elementary blocks*

$$\text{blocks} : \text{Stmt} \rightarrow \mathcal{P}(\text{Blocks})$$

$$\begin{aligned}\text{blocks}([x := a]^\ell) &= \{[x := a]^\ell\} \\ \text{blocks}([\text{skip}]^\ell) &= \{[\text{skip}]^\ell\} \\ \text{blocks}(S_1; S_2) &= \text{blocks}(S_1) \cup \text{blocks}(S_2) \\ \text{blocks}(\text{if } [b]^\ell \text{ then } S_1 \text{ else } S_2) &= \{[b]^\ell\} \cup \text{blocks}(S_1) \cup \text{blocks}(S_2) \\ \text{blocks}(\text{while } [b]^\ell \text{ do } S) &= \{[b]^\ell\} \cup \text{blocks}(S)\end{aligned}$$

A statement S is *label consistent* if and only if any two elementary statements $[S_1]^\ell$ and $[S_2]^\ell$ with the same label in S are equal: $S_1 = S_2$

A statement *where all labels are unique* is automatically label consistent

90

© Nielsen^2, Hankin

Available Expression Analysis

For each point, which expressions *must* be computed, and not later modified, on all paths.

Example:

$$[x := a+b]^1; [y := a*b]^2; \text{while } [y > a+b]^3 \text{ do } ([a := a+1]^4; [x := a+b]^5)$$

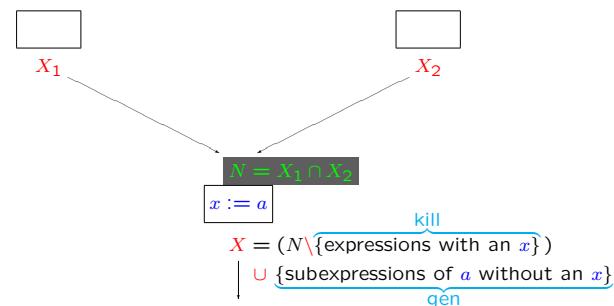
point of interest
↓

91

© Nielsen^2, Hankin

For each point, which expressions *must* be computed, and not later modified, on all paths.

$$[x := a+b]^1; [y := a*b]^2; \text{while } [y > a+b]^3 \text{ do } ([a := a+1]^4; [x := a+b]^5)$$



92

© Nielsen^2, Hankin

| | |
|---|--|
| $N = X_1 \cap X_2$ $x := a$ $X = (N \setminus \{\text{expressions with an } x\})$ $\quad \cup \{\text{subexpressions of } a \text{ without an } x\}$ $\quad \text{gen}$ | <p style="text-align: center;">kill and gen functions</p> $\begin{aligned} kill_{AE}([x := a]^\ell) &= \{a' \in AExp_\star \mid x \in FV(a')\} \\ kill_{AE}([\text{skip}]^\ell) &= \emptyset \\ kill_{AE}([b]^\ell) &= \emptyset \\ gen_{AE}([x := a]^\ell) &= \{a' \in AExp(a) \mid x \notin FV(a')\} \\ gen_{AE}([\text{skip}]^\ell) &= \emptyset \\ gen_{AE}([b]^\ell) &= AExp(b) \end{aligned}$ |
| data flow equations: $AE =$ | |
| $AE_{entry}(\ell) = \begin{cases} \emptyset & \text{if } \ell = init(S_\star) \\ \cap\{AE_{exit}(\ell') \mid (\ell', \ell) \in flow(S_\star)\} & \text{otherwise} \end{cases}$ | |
| $AE_{exit}(\ell) = (AE_{entry}(\ell) \setminus kill_{AE}(B^\ell)) \cup gen_{AE}(B^\ell)$ where $B^\ell \in blocks(S_\star)$ | |

93

© Nielsen^2, Hankin

| | | |
|---|---------------------|------------------|
| [x:=a+b]¹; [y:=a*b]²; while [y>a+b]³ do ([a:=a+1]⁴; [x:=a+b]⁵) | | |
| ℓ | $kill_{AE}(\ell)$ | $gen_{AE}(\ell)$ |
| 1 | \emptyset | $\{a+b\}$ |
| 2 | \emptyset | $\{a*b\}$ |
| 3 | \emptyset | $\{a+b\}$ |
| 4 | $\{a+b, a*b, a+1\}$ | \emptyset |
| 5 | \emptyset | $\{a+b\}$ |

| | | |
|--|--|--|
| kill and gen functions | | |
| $kill_{AE}([x := a]^\ell) = \{a' \in AExp_\star \mid x \in FV(a')\}$ | | |
| $kill_{AE}([\text{skip}]^\ell) = \emptyset$ | | |
| $kill_{AE}([b]^\ell) = \emptyset$ | | |
| $gen_{AE}([x := a]^\ell) = \{a' \in AExp(a) \mid x \notin FV(a')\}$ | | |
| $gen_{AE}([\text{skip}]^\ell) = \emptyset$ | | |
| $gen_{AE}([b]^\ell) = AExp(b)$ | | |

94

| | | |
|--|--|--|
| [x:=a+b]¹; [y:=a*b]²; while [y>a+b]³ do ([a:=a+1]⁴; [x:=a+b]⁵) | | |
| $AE_{entry}(1) = \emptyset$ | | |
| $AE_{entry}(2) = AE_{exit}(1)$ | | |
| $AE_{entry}(3) = AE_{exit}(2) \cap AE_{exit}(5)$ | | |
| $AE_{entry}(4) = AE_{exit}(3)$ | | |
| $AE_{entry}(5) = AE_{exit}(4)$ | | |
| $AE_{exit}(1) = AE_{entry}(1) \cup \{a+b\}$ | | |
| $AE_{exit}(2) = AE_{entry}(2) \cup \{a*b\}$ | | |
| $AE_{exit}(3) = AE_{entry}(3) \cup \{a+b\}$ | | |
| $AE_{exit}(4) = AE_{entry}(4) \setminus \{a+b, a*b, a+1\}$ | | |
| $AE_{exit}(5) = AE_{entry}(5) \cup \{a+b\}$ | | |
| data flow equations: $AE =$ | | |
| $AE_{entry}(\ell) = \begin{cases} \emptyset & \text{if } \ell = init(S_\star) \\ \cap\{AE_{exit}(\ell') \mid (\ell', \ell) \in flow(S_\star)\} & \text{otherwise} \end{cases}$ | | |
| $AE_{exit}(\ell) = (AE_{entry}(\ell) \setminus kill_{AE}(B^\ell)) \cup gen_{AE}(B^\ell)$ where $B^\ell \in blocks(S_\star)$ | | |

95

| | | |
|---|--------------------|-------------------|
| [x:=a+b]¹; [y:=a*b]²; while [y>a+b]³ do ([a:=a+1]⁴; [x:=a+b]⁵) | | |
| Largest solution: | | |
| ℓ | $AE_{entry}(\ell)$ | $AE_{exit}(\ell)$ |
| 1 | \emptyset | $\{a+b\}$ |
| 2 | $\{a+b\}$ | $\{a+b, a*b\}$ |
| 3 | $\{a+b\}$ | $\{a+b\}$ |
| 4 | $\{a+b\}$ | \emptyset |
| 5 | \emptyset | $\{a+b\}$ |

96

© Nielsen^2, Hankin

Available Expression Analysis

For each point, which expressions *must* be computed, and not later modified, on all paths.

Why the largest solution?

$$\begin{aligned} \text{AE}_{\text{entry}}(1) &= \emptyset \\ \text{AE}_{\text{entry}}(2) &= \text{AE}_{\text{exit}}(1) \\ \text{AE}_{\text{entry}}(3) &= \text{AE}_{\text{exit}}(2) \cap \text{AE}_{\text{exit}}(5) \\ \text{AE}_{\text{entry}}(4) &= \text{AE}_{\text{exit}}(3) \\ \text{AE}_{\text{entry}}(5) &= \text{AE}_{\text{exit}}(4) \\ \\ \text{AE}_{\text{exit}}(1) &= \text{AE}_{\text{entry}}(1) \cup \{a+b\} \\ \text{AE}_{\text{exit}}(2) &= \text{AE}_{\text{entry}}(2) \cup \{a*b\} \\ \text{AE}_{\text{exit}}(3) &= \text{AE}_{\text{entry}}(3) \cup \{a+b\} \\ \text{AE}_{\text{exit}}(4) &= \text{AE}_{\text{entry}}(4) \setminus \{a+b, a*b, a+1\} \\ \text{AE}_{\text{exit}}(5) &= \text{AE}_{\text{entry}}(5) \cup \{a+b\} \end{aligned}$$

97

© Nelson^2, Hankin

Available Expression Analysis

For each point, which expressions *must* be computed, and not later modified, on all paths.

Why the largest solution?

$$\begin{aligned} \text{AE}_{\text{entry}}(1) &= \emptyset \\ \text{AE}_{\text{entry}}(2) &= \{a+b\} \\ \text{AE}_{\text{entry}}(3) &= \{a+b, a*b\} \cap \text{AE}_{\text{exit}}(5) \\ \text{AE}_{\text{entry}}(4) &= \{a+b, a*b\} \cap \text{AE}_{\text{exit}}(5) \cup \{a+b\} \\ \text{AE}_{\text{entry}}(5) &= \text{AE}_{\text{exit}}(5) - \{a+b, a*b, a+1\} \\ \\ \text{AE}_{\text{exit}}(1) &= \{a+b\} \\ \text{AE}_{\text{exit}}(2) &= \{a+b, a*b\} \\ \text{AE}_{\text{exit}}(3) &= \{a+b, a*b\} \cap \text{AE}_{\text{exit}}(5) \cup \{a+b\} \\ \text{AE}_{\text{exit}}(4) &= \text{AE}_{\text{exit}}(5) - \{a+b, a*b, a+1\} \\ \text{AE}_{\text{exit}}(5) &= \text{AE}_{\text{exit}}(5) - \{a*b, a+1\} \cup \{a+b\} \\ \\ \text{AExp} &= \{a+b, a*b, a+1\} \end{aligned}$$

98

© Nelson^2, Hankin

Available Expression Analysis

For each point, which expressions *must* be computed, and not later modified, on all paths.

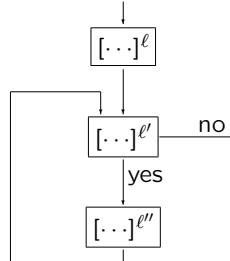
$[z := x+y]^{\ell}; \text{while } [\text{true}]^{\ell'} \text{ do } [\text{skip}]^{\ell''}$

Why the largest solution?

$$\begin{aligned} \text{AE}_{\text{entry}}(\ell) &= \emptyset \\ \text{AE}_{\text{entry}}(\ell') &= \text{AE}_{\text{exit}}(\ell) \cap \text{AE}_{\text{exit}}(\ell'') \\ \text{AE}_{\text{entry}}(\ell'') &= \text{AE}_{\text{exit}}(\ell') \\ \\ \text{AE}_{\text{exit}}(\ell) &= \text{AE}_{\text{entry}}(\ell) \cup \{x+y\} \\ \text{AE}_{\text{exit}}(\ell') &= \text{AE}_{\text{entry}}(\ell') \\ \text{AE}_{\text{exit}}(\ell'') &= \text{AE}_{\text{entry}}(\ell'') \end{aligned}$$

99

© Nelson^2, Hankin



Available Expression Analysis

For each point, which expressions *must* be computed, and not later modified, on all paths.

$[z := x+y]^{\ell}; \text{while } [\text{true}]^{\ell'} \text{ do } [\text{skip}]^{\ell''}$

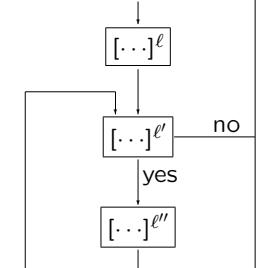
Why the largest solution?

$$\begin{aligned} \text{AE}_{\text{entry}}(\ell) &= \emptyset \\ \text{AE}_{\text{entry}}(\ell') &= \{x+y\} \cap \text{AE}_{\text{exit}}(\ell'') \\ \text{AE}_{\text{entry}}(\ell'') &= \{x+y\} \cap \text{AE}_{\text{exit}}(\ell'') \\ \\ \text{AE}_{\text{exit}}(\ell) &= \{x+y\} \\ \text{AE}_{\text{exit}}(\ell') &= \{x+y\} \cap \text{AE}_{\text{exit}}(\ell'') \\ \text{AE}_{\text{exit}}(\ell'') &= \{x+y\} \cap \text{AE}_{\text{exit}}(\ell'') \end{aligned}$$

100

© Nelson^2, Hankin

$\text{AExp} = \{x + y\}$



Retrospective: Gary Kildall

A number of techniques have evolved for the compile-time analysis of program structure in order to locate redundant computations, perform constant computations, and reduce the number of store-load sequences between memory and high-speed registers. Some of these techniques provide analysis of only straight-line sequences of instructions [5,6,9,14, 17,18,19,20,27,29,34,36,38,39,43,45,46], while others take the program branching structure into account [2,3,4,10,11,12,13,15,23,30,32,33,35]. The purpose here is to describe a single program flow analysis algorithm which extends all of these straight-line optimizing techniques to include branching structure. The algorithm is presented formally and is shown to be correct. Implementation of the flow analysis algorithm in a practical compiler is also discussed.

point of interest

[x:=5]¹; [y:=1]²; while [x>1]³ do ([y:=x*y]⁴; [x:=x-1]⁵)



G. Kildall. A unified approach to global program optimization. POPL73 doi=10.1145/512927.512945

Reaching Definition Analysis

For each point, which assignments *may* have been made, and not overwritten, on some paths.

Example:

point of interest

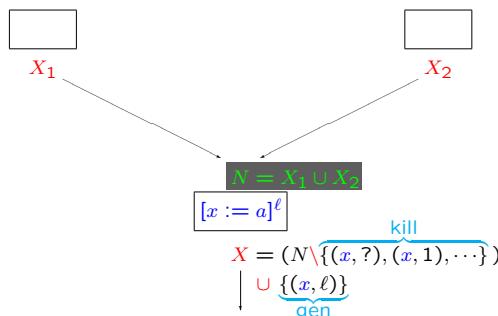
[x:=5]¹; [y:=1]²; while [x>1]³ do ([y:=x*y]⁴; [x:=x-1]⁵)

102

© Nielsen^2, Hankin

For each point, which assignments *may* have been made, and not overwritten, on some paths.

[x:=5]¹; [y:=1]²; while [x>1]³ do ([y:=x*y]⁴; [x:=x-1]⁵)



103

© Nielsen^2, Hankin

kill and gen functions

$$\begin{aligned} \text{kill}_{\text{RD}}([x := a]^\ell) &= \{(x, ?)\} \\ &\cup \{(x, \ell') \mid B^{\ell'} \text{ is an assignment to } x \text{ in } S_*\} \\ \text{kill}_{\text{RD}}([\text{skip}]^\ell) &= \emptyset \\ \text{kill}_{\text{RD}}([b]^\ell) &= \emptyset \\ \text{gen}_{\text{RD}}([x := a]^\ell) &= \{(x, \ell)\} \\ \text{gen}_{\text{RD}}([\text{skip}]^\ell) &= \emptyset \\ \text{gen}_{\text{RD}}([b]^\ell) &= \emptyset \end{aligned}$$

data flow equations: **RD=**

$$\begin{aligned} \text{RD}_{\text{entry}}(\ell) &= \begin{cases} \{(x, ?) \mid x \in FV(S_*)\} & \text{if } \ell = \text{init}(S_*) \\ \bigcup \{\text{RD}_{\text{exit}}(\ell') \mid (\ell', \ell) \in \text{flow}(S_*)\} & \text{otherwise} \end{cases} \\ \text{RD}_{\text{exit}}(\ell) &= (\text{RD}_{\text{entry}}(\ell) \setminus \text{kill}_{\text{RD}}(B^\ell)) \cup \text{gen}_{\text{RD}}(B^\ell) \\ &\text{where } B^\ell \in \text{blocks}(S_*) \end{aligned}$$

104

© Nielsen^2, Hankin

For each point, which assignments *may* have been made, and not overwritten, on some paths.

$[x:=5]^1; [y:=1]^2; \text{while } [x>1]^3 \text{ do } ([y:=\textcolor{red}{x*y}]^{\downarrow 4}; [x:=x-1]^5)$

$$\begin{aligned} RD_{entry}(1) &= \{(x, ?), (y, ?)\} \\ RD_{entry}(2) &= RD_{exit}(1) \\ RD_{entry}(3) &= RD_{exit}(2) \cup RD_{exit}(5) \\ RD_{entry}(4) &= RD_{exit}(3) \\ RD_{entry}(5) &= RD_{exit}(4) \end{aligned}$$

$$\begin{aligned} RD_{exit}(1) &= (RD_{entry}(1) \setminus \{(x, ?), (x, 1), (x, 5)\}) \cup \{(x, 1)\} \\ RD_{exit}(2) &= (RD_{entry}(2) \setminus \{(y, ?), (y, 2), (y, 4)\}) \cup \{(y, 2)\} \\ RD_{exit}(3) &= RD_{entry}(3) \\ RD_{exit}(4) &= (RD_{entry}(4) \setminus \{(y, ?), (y, 2), (y, 4)\}) \cup \{(y, 4)\} \\ RD_{exit}(5) &= (RD_{entry}(5) \setminus \{(x, ?), (x, 1), (x, 5)\}) \cup \{(x, 5)\} \end{aligned}$$

© Nielsen^2, Hankin

| ℓ | $kill_{RD}(\ell)$ | $gen_{RD}(\ell)$ |
|--------|------------------------------|------------------|
| 1 | $\{(x, ?), (x, 1), (x, 5)\}$ | $\{(x, 1)\}$ |
| 2 | $\{(y, ?), (y, 2), (y, 4)\}$ | $\{(y, 2)\}$ |
| 3 | \emptyset | \emptyset |
| 4 | $\{(y, ?), (y, 2), (y, 4)\}$ | $\{(y, 4)\}$ |
| 5 | $\{(x, ?), (x, 1), (x, 5)\}$ | $\{(x, 5)\}$ |

For each point, which assignments *may* have been made, and not overwritten, on some paths.

$[x:=5]^1; [y:=1]^2; \text{while } [x>1]^3 \text{ do } ([y:=\textcolor{red}{x*y}]^{\downarrow 4}; [x:=x-1]^5)$

Smallest solution:

| ℓ | $RD_{entry}(\ell)$ | $RD_{exit}(\ell)$ |
|--------|--------------------------------------|--------------------------------------|
| 1 | $\{(x, ?), (y, ?)\}$ | $\{(y, ?), (x, 1)\}$ |
| 2 | $\{(y, ?), (x, 1)\}$ | $\{(x, 1), (y, 2)\}$ |
| 3 | $\{(x, 1), (y, 2), (y, 4), (x, 5)\}$ | $\{(x, 1), (y, 2), (y, 4), (x, 5)\}$ |
| 4 | $\{(x, 1), (y, 2), (y, 4), (x, 5)\}$ | $\{(x, 1), (y, 4), (x, 5)\}$ |
| 5 | $\{(x, 1), (y, 4), (x, 5)\}$ | $\{(y, 4), (x, 5)\}$ |

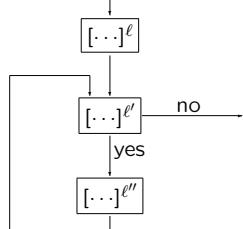
© Nielsen^2, Hankin

Why smallest solution?

$[z:=x+y]^{\ell}; \text{while } [\text{true}]^{\ell'} \text{ do } [\text{skip}]^{\ell''}$

Equations:

$$\begin{aligned} RD_{entry}(\ell) &= \{(x, ?), (y, ?), (z, ?)\} \\ RD_{entry}(\ell') &= RD_{exit}(\ell) \cup RD_{exit}(\ell'') \\ RD_{entry}(\ell'') &= RD_{exit}(\ell') \\ RD_{exit}(\ell) &= (RD_{entry}(\ell) \setminus \{(z, ?)\}) \cup \{(z, \ell)\} \\ RD_{exit}(\ell') &= RD_{entry}(\ell') \\ RD_{exit}(\ell'') &= RD_{entry}(\ell'') \end{aligned}$$



After some simplification: $RD_{entry}(\ell') = \{(x, ?), (y, ?), (z, \ell)\} \cup RD_{entry}(\ell'')$

Many solutions to this equation: any superset of $\{(x, ?), (y, ?), (z, \ell)\}$

107

© Nielsen^2, Hankin

Very Busy Expression Analysis

An expression is very busy at exit of a label, if, no matter what path is taken from the label, the expression is always used before any of the variables occurring in it are redefined.

For each point, which expression *must* be very busy at the exit from that point.

point of interest

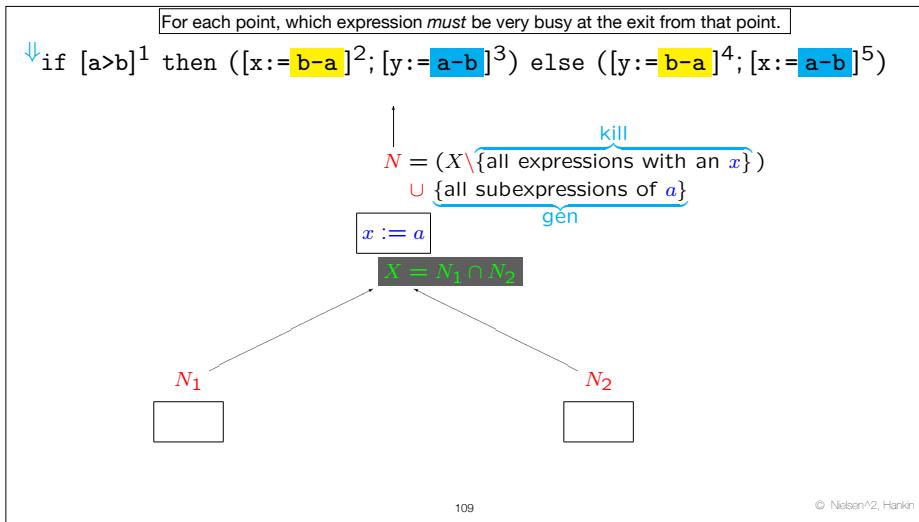
$\Downarrow \text{if } [a>b]^1 \text{ then } ([x:=\textcolor{blue}{b-a}]^2; [y:=\textcolor{red}{a-b}]^3) \text{ else } ([y:=\textcolor{blue}{b-a}]^4; [x:=\textcolor{red}{a-b}]^5)$

The analysis enables a transformation into

$[t1:=\textcolor{blue}{b-a}]^A; [t2:=\textcolor{red}{b-a}]^B;$
 $\text{if } [a>b]^1 \text{ then } ([x:=t1]^2; [y:=t2]^3) \text{ else } ([y:=t1]^4; [x:=t2]^5)$

108

© Nielsen^2, Hankin



109

© Nielsen^2, Hankin

kill and *gen* functions

$$\begin{array}{lcl} \text{kill}_{VB}([x := a]^\ell) & = & \{a' \in AExp_\star \mid x \in FV(a')\} \\ \text{kill}_{VB}([\text{skip}]^\ell) & = & \emptyset \\ \text{kill}_{VB}([b]^\ell) & = & \emptyset \\ \text{gen}_{VB}([x := a]^\ell) & = & AExp(a) \\ \text{gen}_{VB}([\text{skip}]^\ell) & = & \emptyset \\ \text{gen}_{VB}([b]^\ell) & = & AExp(b) \end{array}$$

data flow equations: $VB =$

$$VB_{exit}(\ell) = \begin{cases} \emptyset & \text{if } \ell \in final(S_\star) \\ \cap \{VB_{entry}(\ell') \mid (\ell', \ell) \in flow^R(S_\star)\} & \text{otherwise} \end{cases}$$

$$VB_{entry}(\ell) = (VB_{exit}(\ell) \setminus kill_{VB}(B^\ell)) \cup gen_{VB}(B^\ell)$$

where $B^\ell \in blocks(S_\star)$

110

© Nielsen^2, Hankin

For each point, which expression must be very busy at the exit from that point.

```
if [a>b]1 then ([x:=b-a]2; [y:=a-b]3) else ([y:=b-a]4; [x:=a-b]5)
```

$VB_{entry}(1) = VB_{exit}(1)$
 $VB_{entry}(2) = VB_{exit}(2) \cup \{b-a\}$
 $VB_{entry}(3) = \{a-b\}$
 $VB_{entry}(4) = VB_{exit}(4) \cup \{b-a\}$
 $VB_{entry}(5) = \{a-b\}$

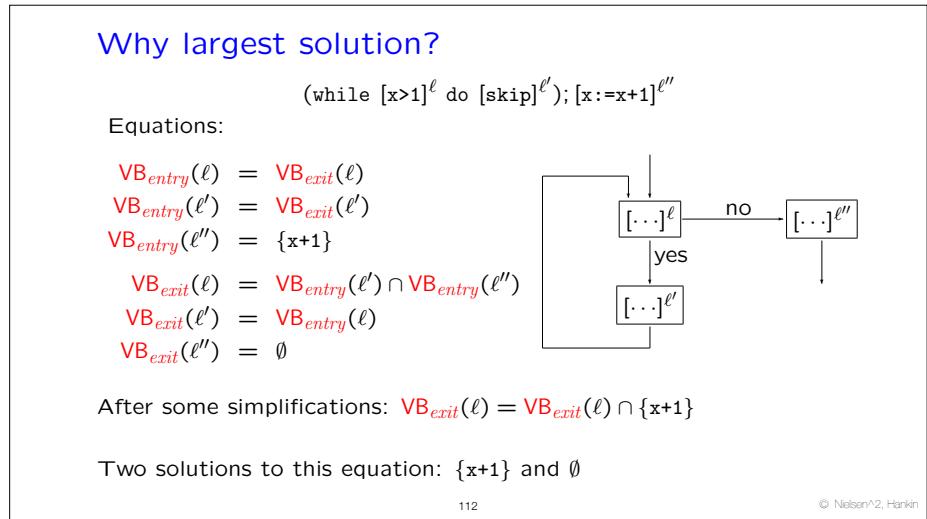
$VB_{exit}(1) = VB_{entry}(2) \cap VB_{entry}(4)$
 $VB_{exit}(2) = VB_{entry}(3)$
 $VB_{exit}(3) = \emptyset$
 $VB_{exit}(4) = VB_{entry}(5)$
 $VB_{exit}(5) = \emptyset$

Largest solution:

| ℓ | $VB_{entry}(\ell)$ | $VB_{exit}(\ell)$ |
|--------|--------------------|-------------------|
| 1 | $\{a-b, b-a\}$ | $\{a-b, b-a\}$ |
| 2 | $\{a-b, b-a\}$ | $\{a-b\}$ |
| 3 | $\{a-b\}$ | \emptyset |
| 4 | $\{a-b, b-a\}$ | $\{a-b\}$ |
| 5 | $\{a-b\}$ | \emptyset |

111

© Nielsen^2, Hankin



112

© Nielsen^2, Hankin

Live Variable Analysis

A variable is live at the exit from a label, if there is a path from the label to a use of the variable on which it is not overwritten.

For each point, which variables *may* be live at the exit from that point.

point of interest

\downarrow
 $[x := 2]^1; [y := 4]^2; [x := 1]^3; (\text{if } [y > x]^4 \text{ then } [z := y]^5 \text{ else } [z := y * y]^6); [x := z]^7$

The analysis enables a transformation into

$[y := 4]^2; [x := 1]^3; (\text{if } [y > x]^4 \text{ then } [z := y]^5 \text{ else } [z := y * y]^6); [x := z]^7$

113

© Nielsen^2, Hankin

Live Variable Analysis

A variable is live at the exit from a label, if there is a path from the label to a use of the variable on which it is not overwritten.

For each point, which variables *may* be live at the exit from that point.

$$N = (X \setminus \{x\}) \cup \{\text{all variables of } a\}$$

$x := a$

$X = N_1 \cup N_2$

N_1

N_2

114

© Nielsen^2, Hankin

For each point, which variables *may* be live at the exit from that point.

$$N = (X \setminus \{x\}) \cup \{\text{all variables of } a\}$$

$x := a$

$X = N_1 \cup N_2$

kill and gen functions

$$\begin{aligned} kill_{LV}([x := a]^\ell) &= \{x\} \\ kill_{LV}([\text{skip}]^\ell) &= \emptyset \\ kill_{LV}([b]^\ell) &= \emptyset \\ gen_{LV}([x := a]^\ell) &= FV(a) \\ gen_{LV}([\text{skip}]^\ell) &= \emptyset \\ gen_{LV}([b]^\ell) &= FV(b) \end{aligned}$$

data flow equations: $LV =$

$$LV_{exit}(\ell) = \begin{cases} \emptyset & \text{if } \ell \in \text{final}(S_*) \\ \bigcup \{LV_{entry}(\ell') \mid (\ell', \ell) \in \text{flow}^R(S_*)\} & \text{otherwise} \end{cases}$$

$$LV_{entry}(\ell) = (LV_{exit}(\ell) \setminus kill_{LV}(B^\ell)) \cup gen_{LV}(B^\ell)$$

where $B^\ell \in \text{blocks}(S_*)$

115

© Nielsen^2, Hankin

\downarrow
 $[x := 2]^1; [y := 4]^2; [x := 1]^3; (\text{if } [y > x]^4 \text{ then } [z := y]^5 \text{ else } [z := y * y]^6); [x := z]^7$

| ℓ | $kill_{LV}(\ell)$ | $gen_{LV}(\ell)$ |
|--------|-------------------|------------------|
| 1 | {x} | \emptyset |
| 2 | {y} | \emptyset |
| 3 | {x} | \emptyset |
| 4 | \emptyset | {x, y} |
| 5 | {z} | \emptyset |
| 6 | {z} | {y} |
| 7 | {x} | {z} |

Smallest solution:

| ℓ | $LV_{entry}(\ell)$ | $LV_{exit}(\ell)$ |
|--------|--------------------|-------------------|
| 1 | \emptyset | \emptyset |
| 2 | \emptyset | {y} |
| 3 | {y} | {x, y} |
| 4 | {x, y} | {y} |
| 5 | {y} | {z} |
| 6 | {y} | {z} |
| 7 | {z} | \emptyset |

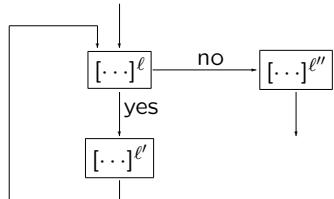
© Nielsen^2, Hankin

Why smallest solution?

(while [x>1] $^\ell$ do [skip] $^{\ell'}$); [x:=x+1] $^{\ell''}$

Equations:

$$\begin{aligned}\text{LV}_{\text{entry}}(\ell) &= \text{LV}_{\text{exit}}(\ell) \cup \{\text{x}\} \\ \text{LV}_{\text{entry}}(\ell') &= \text{LV}_{\text{exit}}(\ell') \\ \text{LV}_{\text{entry}}(\ell'') &= \{\text{x}\} \\ \text{LV}_{\text{exit}}(\ell) &= \text{LV}_{\text{entry}}(\ell') \cup \text{LV}_{\text{entry}}(\ell'') \\ \text{LV}_{\text{exit}}(\ell') &= \text{LV}_{\text{entry}}(\ell) \\ \text{LV}_{\text{exit}}(\ell'') &= \emptyset\end{aligned}$$



After some calculations: $\text{LV}_{\text{exit}}(\ell) = \text{LV}_{\text{exit}}(\ell) \cup \{x\}$

Many solutions to this equation: any superset of $\{x\}$

117

© Nielsen^2, Hankin

Retrospective: Fran Allen



Fran Allen
Turing Award (2006)

```

graph LR
    Data[Data] --> IV[Input Validation]
    IV --> DT[Data Transformation]
    IV --> DA[Data Analysis]
    DT --> OG[Output Generation]
    DA --> OG
    OG --> Report[Report]
  
```

Retrospective

1960-1965: Data Flow

1965-1970: Control Flow

1970-1975: Program Structure

1975-1980: Program Semantics

Data Flow

- 1960: Datalogic (D. H. Dill)
- 1965: PERT (Project Evaluation and Review Technique)
- 1970: DCE (Data Flow Computation)
- 1975: DCE (Data Flow Computation)
- 1980: DCE (Data Flow Computation)

Control Flow

- 1960: Datalogic (D. H. Dill)
- 1965: PERT (Project Evaluation and Review Technique)
- 1970: DCE (Data Flow Computation)
- 1975: DCE (Data Flow Computation)
- 1980: DCE (Data Flow Computation)

Program Structure

- 1960: Datalogic (D. H. Dill)
- 1965: PERT (Project Evaluation and Review Technique)
- 1970: DCE (Data Flow Computation)
- 1975: DCE (Data Flow Computation)
- 1980: DCE (Data Flow Computation)

Program Semantics

- 1960: Datalogic (D. H. Dill)
- 1965: PERT (Project Evaluation and Review Technique)
- 1970: DCE (Data Flow Computation)
- 1975: DCE (Data Flow Computation)
- 1980: DCE (Data Flow Computation)

Derived Information

- *Use-Definition chains* or *ud chains*:

each `use` of a variable is linked to all `assignments` that reach it

```
[x:=0]1; [x:=3]2; (if [z=x]3 then [z:=0]4 else [z:=x]5); [y:=x]6; [x:=y+z]7
```

- *Definition-Use chains* or *du chains*:

each `assignment` to a variable is linked to all `uses` of it

```
[x:=0]1; [x]:=32; (if [z=x]3 then [z:=0]4 else [z:=x]5); [y:=x]6; [x:=y+z]7
```

ud chains

$\text{ud} : \text{Var}_\star \times \text{Lab}_\star \rightarrow \mathcal{P}(\text{Lab}_\star)$

given by

$$\begin{aligned}\text{ud}(x, \ell') &= \{\ell \mid \text{def}(x, \ell) \wedge \exists \ell'': (\ell, \ell'') \in \text{flow}(S_\star) \wedge \text{clear}(x, \ell'', \ell')\} \\ &\cup \{? \mid \text{clear}(x, \text{init}(S_\star), \ell')\}\end{aligned}$$

where



- $\text{def}(x, \ell)$ means that the block ℓ assigns a value to x
- $\text{clear}(x, \ell, \ell')$ means that none of the blocks on a path from ℓ to ℓ' contains an assignments to x but that the block ℓ' uses x (in a test or on the right hand side of an assignment)

121

© Nielsen^2, Hankin

ud chains - an alternative definition

$\text{UD} : \text{Var}_\star \times \text{Lab}_\star \rightarrow \mathcal{P}(\text{Lab}_\star)$

is defined by:

$$\text{UD}(x, \ell) = \begin{cases} \{\ell' \mid (x, \ell') \in \text{RD}_{\text{entry}}(\ell)\} & \text{if } x \in \text{gen}_{\text{LV}}(B^\ell) \\ \emptyset & \text{otherwise} \end{cases}$$

One can show that:

$$\text{ud}(x, \ell) = \text{UD}(x, \ell)$$

122

© Nielsen^2, Hankin

du chains

$\text{du} : \text{Var}_\star \times \text{Lab}_\star \rightarrow \mathcal{P}(\text{Lab}_\star)$

given by

$$\text{du}(x, \ell) = \begin{cases} \{\ell' \mid \text{def}(x, \ell) \wedge \exists \ell'': (\ell, \ell'') \in \text{flow}(S_\star) \wedge \text{clear}(x, \ell'', \ell')\} \\ \quad \text{if } \ell \neq ? \\ \{\ell' \mid \text{clear}(x, \text{init}(S_\star), \ell')\} \\ \quad \text{if } \ell = ? \end{cases}$$



One can show that:

$$\text{du}(x, \ell) = \{\ell' \mid \ell \in \text{ud}(x, \ell')\}$$

123

© Nielsen^2, Hankin

Example:

$$[x := 0]^1; [x := 3]^2; (\text{if } [z = x]^3 \text{ then } [z := 0]^4 \text{ else } [z := x]^5); [y := x]^6; [x := y + z]^7$$

| $\text{ud}(x, \ell)$ | x | y | z | $\text{du}(x, \ell)$ | x | y | z |
|----------------------|-------------|-------------|-------------|----------------------|---------------|-------------|-------------|
| 1 | \emptyset | \emptyset | \emptyset | 1 | \emptyset | \emptyset | \emptyset |
| 2 | \emptyset | \emptyset | \emptyset | 2 | $\{3, 5, 6\}$ | \emptyset | \emptyset |
| 3 | $\{2\}$ | \emptyset | $\{?\}$ | 3 | \emptyset | \emptyset | \emptyset |
| 4 | \emptyset | \emptyset | \emptyset | 4 | \emptyset | \emptyset | $\{7\}$ |
| 5 | $\{2\}$ | \emptyset | \emptyset | 5 | \emptyset | \emptyset | $\{7\}$ |
| 6 | $\{2\}$ | \emptyset | \emptyset | 6 | \emptyset | $\{7\}$ | \emptyset |
| 7 | \emptyset | $\{6\}$ | $\{4, 5\}$ | 7 | \emptyset | \emptyset | $\{3\}$ |
| | | | | ? | \emptyset | \emptyset | $\{3\}$ |

124

© Nielsen^2, Hankin

Theoretical Properties

- Structural operational semantics
- Correctness of Live Variable Analysis

125

© Nielsen^2, Hankin

A *state* is a mapping from variables to integers:

$$\sigma \in \text{State} = \text{Var} \rightarrow \mathbb{Z}$$

The semantics of arithmetic and boolean expressions

$$\mathcal{A} : \text{AExp} \rightarrow (\text{State} \rightarrow \mathbb{Z}) \quad (\text{no errors allowed})$$

$$\mathcal{B} : \text{BExp} \rightarrow (\text{State} \rightarrow \text{T}) \quad (\text{no errors allowed})$$

The *transitions* of the semantics are of the form

$$\langle S, \sigma \rangle \rightarrow \sigma' \quad \text{and} \quad \langle S, \sigma \rangle \rightarrow \langle S', \sigma' \rangle$$

$$\langle [x := a]^\ell, \sigma \rangle \rightarrow \sigma[x \mapsto \mathcal{A}[a]\sigma]$$

$$\langle [\text{skip}]^\ell, \sigma \rangle \rightarrow \sigma$$

$$\frac{\langle S_1, \sigma \rangle \rightarrow \langle S'_1, \sigma' \rangle}{\langle S_1; S_2, \sigma \rangle \rightarrow \langle S'_1; S_2, \sigma' \rangle}$$

$$\frac{\langle S_1, \sigma \rangle \rightarrow \sigma'}{\langle S_1; S_2, \sigma \rangle \rightarrow \langle S_2, \sigma' \rangle}$$

$$\langle \text{if } [b]^\ell \text{ then } S_1 \text{ else } S_2, \sigma \rangle \rightarrow \langle S_1, \sigma \rangle \quad \text{if } \mathcal{B}[b]\sigma = \text{true}$$

$$\langle \text{if } [b]^\ell \text{ then } S_1 \text{ else } S_2, \sigma \rangle \rightarrow \langle S_2, \sigma \rangle \quad \text{if } \mathcal{B}[b]\sigma = \text{false}$$

$$\langle \text{while } [b]^\ell \text{ do } S, \sigma \rangle \rightarrow \langle (S; \text{while } [b]^\ell \text{ do } S), \sigma \rangle \quad \text{if } \mathcal{B}[b]\sigma = \text{true}$$

$$\langle \text{while } [b]^\ell \text{ do } S, \sigma \rangle \rightarrow \sigma \quad \text{if } \mathcal{B}[b]\sigma = \text{false}$$

126

© Nielsen^2, Hankin

Example:

$$\begin{aligned} & \langle [y:=x]^1; [z:=1]^2; \text{while } [y>1]^3 \text{ do } ([z:=z*y]^4; [y:=y-1]^5); [y:=0]^6, \sigma_{300} \rangle \\ & \rightarrow \langle [z:=1]^2; \text{while } [y>1]^3 \text{ do } ([z:=z*y]^4; [y:=y-1]^5); [y:=0]^6, \sigma_{330} \rangle \end{aligned}$$

© Nielsen^2, Hankin

Equations and Constraints

Equation system $\text{LV}=(S_\star)$:

$$\text{LV}_{\text{exit}}(\ell) = \begin{cases} \emptyset & \text{if } \ell \in \text{final}(S_\star) \\ \cup \{\text{LV}_{\text{entry}}(\ell') \mid (\ell', \ell) \in \text{flow}^R(S_\star)\} & \text{otherwise} \end{cases}$$

$$\text{LV}_{\text{entry}}(\ell) = (\text{LV}_{\text{exit}}(\ell) \setminus \text{kill}_{\text{LV}}(B^\ell)) \cup \text{gen}_{\text{LV}}(B^\ell) \\ \text{where } B^\ell \in \text{blocks}(S_\star)$$

Constraint system $\text{LV}^\subseteq(S_\star)$:

$$\text{LV}_{\text{exit}}(\ell) \supseteq \begin{cases} \emptyset & \text{if } \ell \in \text{final}(S_\star) \\ \cup \{\text{LV}_{\text{entry}}(\ell') \mid (\ell', \ell) \in \text{flow}^R(S_\star)\} & \text{otherwise} \end{cases}$$

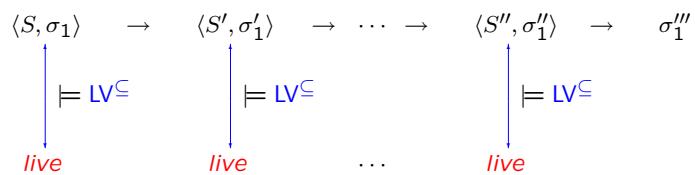
$$\text{LV}_{\text{entry}}(\ell) \supseteq (\text{LV}_{\text{exit}}(\ell) \setminus \text{kill}_{\text{LV}}(B^\ell)) \cup \text{gen}_{\text{LV}}(B^\ell) \\ \text{where } B^\ell \in \text{blocks}(S_\star)$$

128

© Nielsen^2, Hankin

Lemma

A solution *live* to the constraint system is preserved during computation



Proof: requires a lot of machinery — see the book.

Correctness Relation

$$\sigma_1 \sim_V \sigma_2$$

means that for all practical purposes the two states σ_1 and σ_2 are equal: only the values of the live variables of V matters and here the two states are equal.

Example:

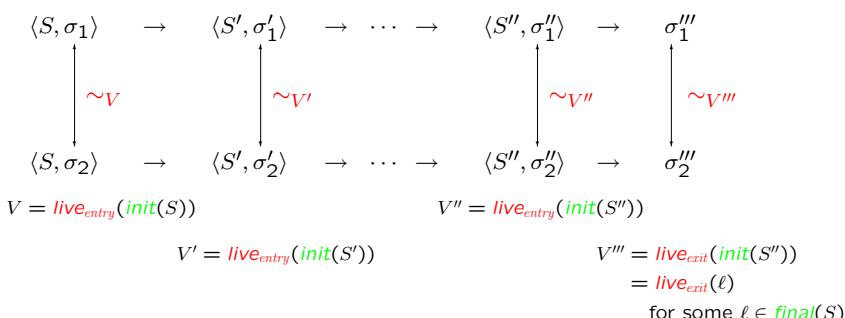
Consider the statement $[x := y + z]$.

Let $V_1 = \{y, z\}$. Then $\sigma_1 \sim_{V_1} \sigma_2$ means $\sigma_1(y) = \sigma_2(y) \wedge \sigma_1(z) = \sigma_2(z)$

Let $V_2 = \{x\}$. Then $\sigma_1 \sim_{V_2} \sigma_2$ means $\sigma_1(x) = \sigma_2(x)$

Correctness Theorem

The relation “ \sim ” is *invariant* under computation: the live variables for the initial configuration remain live throughout the computation.



Retrospective: Ullman



Jeffrey D. Ullman

Acta Informatica 7, 305–317 (1977)
© by Springer-Verlag 1977

Monotone Data Flow Analysis Framework

John B. Kam and Jeffrey D. Ullman

Received March 24, 1

Summary. We consider a generalization of Kildall's lattice theoretic approach to data flow analysis, which we call *monotone data flow analysis frameworks*. Many flow analysis problems which appear in practice meet the monotonicity condition but not Kildall's condition called *distributivity*. We show that the maximal fixed point solution exists for every instance of every monotone framework, and that it can be obtained by Kildall's algorithm. However, whenever the framework is monotone but not distributive, there are instances in which the desired solution—the “meet over all paths solution” —differs from the maximal fixed point. Finally, we show the nonexistence of an algorithm to compute the meet over all paths solution for monotone frameworks.

In this paper, we present a general class of gilpin's called *monotone data analysis frameworks* introduced by a new concept called *monotone solvability*. The paper also shows that for monotone frameworks, the MOP solution for an individual program does not necessarily coincide with the maximum fixed point solution to the corresponding set of nonlinear equations. Several methods for approximating the class of frameworks can be mentioned. We conclude the paper by showing that there exists an algorithm which, given an arbitrary monotone type framework, will compute the MOP for each program.

Monotone Frameworks

- Monotone and Distributive Frameworks
- Instances of Frameworks
- Constant Propagation Analysis

133

© Nielsen^2, Hankin

The Overall Pattern

Each of the four classical analyses take the form

$$\text{Analysis}_\circ(\ell) = \begin{cases} \ell & \text{if } \ell \in E \\ \sqcup \{\text{Analysis}_\bullet(\ell') \mid (\ell', \ell) \in F\} & \text{otherwise} \end{cases}$$

$$\text{Analysis}_\bullet(\ell) = f_\ell(\text{Analysis}_\circ(\ell))$$

where

- \sqcup is \cap or \cup (and \sqcup is \cup or \cap),
- F is either $\text{flow}(S_*)$ or $\text{flow}^R(S_*)$,
- E is $\{\text{init}(S_*)\}$ or $\text{final}(S_*)$,
- \circ specifies the initial or final analysis information, and
- f_ℓ is the transfer function associated with $B^\ell \in \text{blocks}(S_*)$.

134

© Nielsen^2, Hankin

The Principle: forward versus backward

- The *forward analyses* have F to be $\text{flow}(S_*)$ and then Analysis_\circ concerns entry conditions and Analysis_\bullet concerns exit conditions; the equation system presupposes that S_* has isolated entries.
- The *backward analyses* have F to be $\text{flow}^R(S_*)$ and then Analysis_\circ concerns exit conditions and Analysis_\bullet concerns entry conditions; the equation system presupposes that S_* has isolated exits.

135

© Nielsen^2, Hankin

The Principle: union versus intersection

- When \sqcup is \cap we require the *greatest sets* that solve the equations and we are able to detect properties satisfied by *all execution paths* reaching (or leaving) the entry (or exit) of a label; the analysis is called a *must*-analysis.
- When \sqcup is \cup we require the *smallest sets* that solve the equations and we are able to detect properties satisfied by *at least one execution path* to (or from) the entry (or exit) of a label; the analysis is called a *may*-analysis.

136

© Nielsen^2, Hankin

Property Spaces

The *property space*, L , is used to represent the data flow information, and the *combination operator*, $\sqcup: \mathcal{P}(L) \rightarrow L$, is used to combine information from different paths.

- L is a *complete lattice*, that is, a partially ordered set, (L, \sqsubseteq) , such that each subset, Y , has a least upper bound, $\sqcup Y$.
- L satisfies the *Ascending Chain Condition*; that is, each ascending chain eventually stabilises (meaning that if $(l_n)_n$ is such that $l_1 \sqsubseteq l_2 \sqsubseteq l_3 \sqsubseteq \dots$, then there exists n such that $l_n = l_{n+1} = \dots$).

137

© Nielsen^2, Hankin

Example: Reaching Definitions

- $L = \mathcal{P}(\text{Var}_* \times \text{Lab}_*)$ is partially ordered by subset inclusion so \sqsubseteq is \subseteq
- the least upper bound operation \sqcup is \cup and the least element \perp is \emptyset
- L satisfies the Ascending Chain Condition because $\text{Var}_* \times \text{Lab}_*$ is finite (unlike $\text{Var} \times \text{Lab}$)

138

© Nielsen^2, Hankin

Example: Available Expressions

- $L = \mathcal{P}(\text{AExp}_*)$ is partially ordered by superset inclusion so \sqsubseteq is \supseteq
- the least upper bound operation \sqcup is \cap and the least element \perp is AExp_*
- L satisfies the Ascending Chain Condition because AExp_* is finite (unlike AExp)

139

© Nielsen^2, Hankin

Transfer Functions

The set of transfer functions, \mathcal{F} , is a set of *monotone functions* over L , meaning that

$$l \sqsubseteq l' \text{ implies } f_\ell(l) \sqsubseteq f_\ell(l')$$

and furthermore they fulfil the following conditions:

- \mathcal{F} contains *all* the transfer functions $f_\ell : L \rightarrow L$ in question (for $\ell \in \text{Lab}_*$)
- \mathcal{F} contains the *identity function*
- \mathcal{F} is *closed under composition* of functions

140

© Nielsen^2, Hankin

Frameworks

A *Monotone Framework* consists of:

- a complete lattice, L , that satisfies the Ascending Chain Condition; we write \sqcup for the least upper bound operator
- a set \mathcal{F} of monotone functions from L to L that contains the identity function and that is closed under function composition

A *Distributive Framework* is a Monotone Framework where additionally all functions f in \mathcal{F} are required to be *distributive*:

$$f(l_1 \sqcup l_2) = f(l_1) \sqcup f(l_2)$$

141

© Nielsen^2, Hankin

Instances

An *instance* of a Framework consists of:

- the complete lattice, L , of the framework
- the space of functions, \mathcal{F} , of the framework
- a finite flow, F (typically $\text{flow}(S_*)$ or $\text{flow}^R(S_*)$)
- a finite set of *extremal labels*, E (typically $\{\text{init}(S_*)\}$ or $\text{final}(S_*)$)
- an *extremal value*, $\iota \in L$, for the extremal labels
- a mapping, f_* , from the labels Lab_* to transfer functions in \mathcal{F}

142

© Nielsen^2, Hankin

Equations of the Instance:

$$\text{Analysis}_{\circ}(\ell) = \bigsqcup \{\text{Analysis}_{\bullet}(\ell') \mid (\ell', \ell) \in F\} \sqcup \iota_E^\ell$$

where $\iota_E^\ell = \begin{cases} \iota & \text{if } \ell \in E \\ \perp & \text{if } \ell \notin E \end{cases}$

$$\text{Analysis}_{\bullet}(\ell) = f_\ell(\text{Analysis}_{\circ}(\ell))$$

Constraints of the Instance:

$$\text{Analysis}_{\circ}(\ell) \sqsupseteq \bigsqcup \{\text{Analysis}_{\bullet}(\ell') \mid (\ell', \ell) \in F\} \sqcup \iota_E^\ell$$

where $\iota_E^\ell = \begin{cases} \iota & \text{if } \ell \in E \\ \perp & \text{if } \ell \notin E \end{cases}$

$$\text{Analysis}_{\bullet}(\ell) \sqsupseteq f_\ell(\text{Analysis}_{\circ}(\ell))$$

143

© Nielsen^2, Hankin

The Examples Revisited

| | Available Expressions | Reaching Definitions | Very Busy Expressions | Live Variables |
|---------------|---|---|------------------------------|-----------------------------|
| L | $\mathcal{P}(\text{AExp}_*)$ | $\mathcal{P}(\text{Var}_* \times \text{Lab}_*)$ | $\mathcal{P}(\text{AExp}_*)$ | $\mathcal{P}(\text{Var}_*)$ |
| \sqsubseteq | \supseteq | \subseteq | \supseteq | \subseteq |
| \sqcup | \cap | \cup | \cap | \cup |
| \perp | AExp_* | \emptyset | AExp_* | \emptyset |
| ι | \emptyset | $\{(x, ?) \mid x \in \text{FV}(S_*)\}$ | \emptyset | \emptyset |
| E | $\{\text{init}(S_*)\}$ | $\{\text{init}(S_*)\}$ | $\text{final}(S_*)$ | $\text{final}(S_*)$ |
| F | $\text{flow}(S_*)$ | $\text{flow}(S_*)$ | $\text{flow}^R(S_*)$ | $\text{flow}^R(S_*)$ |
| \mathcal{F} | $\{f : L \rightarrow L \mid \exists l_k, l_g : f(l) = (l \setminus l_k) \cup l_g\}$ | | | |
| f_ℓ | $f_\ell(l) = (l \setminus \text{kill}(B^\ell)) \cup \text{gen}(B^\ell)$ where $B^\ell \in \text{blocks}(S_*)$ | | | |

144

© Nielsen^2, Hankin

Bit Vector Frameworks

A *Bit Vector Framework* has

- $L = \mathcal{P}(D)$ for D finite
- $\mathcal{F} = \{f \mid \exists l_k, l_g : f(l) = (l \setminus l_k) \cup l_g\}$

Examples:

- Available Expressions
- Live Variables
- Reaching Definitions
- Very Busy Expressions

145

© Nielsen^2, Hankin

The Constant Propagation Framework

An example of a Monotone Framework that is **not** a Distributive Framework

The aim of the *Constant Propagation Analysis* is to determine

For each program point, whether or not a variable has a constant value whenever execution reaches that point.

Example:

`[x:=6]^1; [y:=3]^2; while [x > y]^3 do ([x:=x - 1]^4; [z:=y * y]^6)`

The analysis enables a transformation into

`[x:=6]^1; [y:=3]^2; while [x > 3]^3 do ([x:=x - 1]^4; [z:=9]^6)`

147

© Nielsen^2, Hankin

Lemma: Bit Vector Frameworks are always Distributive Frameworks

Proof

$$\begin{aligned} f(l_1 \sqcup l_2) &= \begin{cases} f(l_1 \cup l_2) \\ f(l_1 \cap l_2) \end{cases} &= \begin{cases} ((l_1 \cup l_2) \setminus l_k) \cup l_g \\ ((l_1 \cap l_2) \setminus l_k) \cup l_g \end{cases} \\ &= \begin{cases} ((l_1 \setminus l_k) \cup (l_2 \setminus l_k)) \cup l_g \\ ((l_1 \setminus l_k) \cap (l_2 \setminus l_k)) \cup l_g \end{cases} &= \begin{cases} ((l_1 \setminus l_k) \cup l_g) \cup ((l_2 \setminus l_k) \cup l_g) \\ ((l_1 \setminus l_k) \cap l_g) \cap ((l_2 \setminus l_k) \cap l_g) \end{cases} \\ &= \begin{cases} f(l_1) \cup f(l_2) \\ f(l_1) \cap f(l_2) \end{cases} &= f(l_1) \sqcup f(l_2) \end{aligned}$$

- $id(l) = (l \setminus \emptyset) \cup \emptyset$
- $f_2(f_1(l)) = (((l \setminus l_k^1) \cup l_g^1) \setminus l_k^2) \cup l_g^2 = (l \setminus (l_k^1 \cup l_k^2)) \cup ((l_g^1 \setminus l_k^2) \cup l_g^2)$
- monotonicity follows from distributivity
- $\mathcal{P}(D)$ satisfies the Ascending Chain Condition because D is finite

146

© Nielsen^2, Hankin

Elements of L

$$\widehat{\text{State}}_{\text{CP}} = ((\text{Var}_* \rightarrow \mathbf{Z}^\top) \perp, \sqsubseteq)$$

Idea:

- \perp is the least element: no information is available
- $\hat{\sigma} \in \text{Var}_* \rightarrow \mathbf{Z}^\top$ specifies for each variable whether it is constant:
 - $\hat{\sigma}(x) \in \mathbf{Z}$: x is constant and the value is $\hat{\sigma}(x)$
 - $\hat{\sigma}(x) = \top$: x might not be constant

148

© Nielsen^2, Hankin

Partial Ordering on L

The partial ordering \sqsubseteq on $(\text{Var}_* \rightarrow \mathbf{Z}^{\top})_{\perp}$ is defined by

$$\begin{aligned}\forall \hat{\sigma} \in (\text{Var}_* \rightarrow \mathbf{Z}^{\top})_{\perp}: \quad \perp &\sqsubseteq \hat{\sigma} \\ \forall \hat{\sigma}_1, \hat{\sigma}_2 \in \text{Var}_* \rightarrow \mathbf{Z}^{\top}: \quad \hat{\sigma}_1 &\sqsubseteq \hat{\sigma}_2 \quad \text{iff} \quad \forall x: \hat{\sigma}_1(x) \sqsubseteq \hat{\sigma}_2(x)\end{aligned}$$

where $\mathbf{Z}^{\top} = \mathbf{Z} \cup \{\top\}$ is partially ordered as follows:

$$\begin{aligned}\forall z \in \mathbf{Z}^{\top}: z &\sqsubseteq \top \\ \forall z_1, z_2 \in \mathbf{Z}: (z_1 &\sqsubseteq z_2) \Leftrightarrow (z_1 = z_2)\end{aligned}$$

149

© Nielsen^2, Hankin

Transfer Functions in \mathcal{F}

$$\mathcal{F}_{\text{CP}} = \{f \mid f \text{ is a monotone function on } \widehat{\text{State}}_{\text{CP}}\}$$

Lemma

Constant Propagation as defined by $\widehat{\text{State}}_{\text{CP}}$ and \mathcal{F}_{CP} is a Monotone Framework

150

© Nielsen^2, Hankin

Instances

Constant Propagation is a forward analysis, so for the program S_* :

- the flow, F , is $\text{flow}(S_*)$,
- the extremal labels, E , is $\{\text{init}(S_*)\}$,
- the extremal value, ι_{CP} , is $\lambda x. \top$, and
- the mapping, f^{CP} , of labels to transfer functions is as shown next

151

© Nielsen^2, Hankin

Constant Propagation Analysis

$$\mathcal{A}_{\text{CP}} : \text{AExp} \rightarrow (\widehat{\text{State}}_{\text{CP}} \rightarrow \mathbf{Z}^{\top})$$

$$\mathcal{A}_{\text{CP}}[x]\hat{\sigma} = \begin{cases} \perp & \text{if } \hat{\sigma} = \perp \\ \hat{\sigma}(x) & \text{otherwise} \end{cases}$$

$$\mathcal{A}_{\text{CP}}[n]\hat{\sigma} = \begin{cases} \perp & \text{if } \hat{\sigma} = \perp \\ n & \text{otherwise} \end{cases}$$

$$\mathcal{A}_{\text{CP}}[a_1 \text{ op}_a a_2]\hat{\sigma} = \mathcal{A}_{\text{CP}}[a_1]\hat{\sigma} \text{ op}_a \mathcal{A}_{\text{CP}}[a_2]\hat{\sigma}$$

transfer functions: f_{ℓ}^{CP}

$$[x := a]^{\ell}: f_{\ell}^{\text{CP}}(\hat{\sigma}) = \begin{cases} \perp & \text{if } \hat{\sigma} = \perp \\ \hat{\sigma}[x \mapsto \mathcal{A}_{\text{CP}}[a]\hat{\sigma}] & \text{otherwise} \end{cases}$$

$$[\text{skip}]^{\ell}: f_{\ell}^{\text{CP}}(\hat{\sigma}) = \hat{\sigma}$$

$$[b]^{\ell}: f_{\ell}^{\text{CP}}(\hat{\sigma}) = \hat{\sigma}$$

© Nielsen^2, Hankin

152

Lemma

Constant Propagation is **not** a Distributive Framework

Proof

Consider the transfer function f_ℓ^{CP} for $[y:=x*x]^\ell$

Let $\hat{\sigma}_1$ and $\hat{\sigma}_2$ be such that $\hat{\sigma}_1(x) = 1$ and $\hat{\sigma}_2(x) = -1$

Then $\hat{\sigma}_1 \sqcup \hat{\sigma}_2$ maps x to \top — $f_\ell^{\text{CP}}(\hat{\sigma}_1 \sqcup \hat{\sigma}_2)$ maps y to \top

Both $f_\ell^{\text{CP}}(\hat{\sigma}_1)$ and $f_\ell^{\text{CP}}(\hat{\sigma}_2)$ map y to 1 — $f_\ell^{\text{CP}}(\hat{\sigma}_1) \sqcup f_\ell^{\text{CP}}(\hat{\sigma}_2)$ maps y to 1

153

© Nielsen^2, Hankin

Equation Solving

- The MFP solution — “Maximum” (actually least) Fixed Point
 - Worklist algorithm for Monotone Frameworks
- The MOP solution — “Meet” (actually join) Over all Paths

154

© Nielsen^2, Hankin

The MFP Solution

- Idea: iterate until stabilisation.

Worklist Algorithm

Input: An instance $(L, \mathcal{F}, F, E, \iota, f)$ of a Monotone Framework

Output: The MFP Solution: MFP_o, MFP_\bullet

Data structures:

- **Analysis:** the current analysis result for block entries (or exits)
- The worklist **W**: a list of pairs (ℓ, ℓ') indicating that the current analysis result has changed at the entry (or exit) to the block ℓ and hence the entry (or exit) information must be recomputed for ℓ'

155

© Nielsen^2, Hankin

Worklist Algorithm

- Step 1 **Initialisation (of W and Analysis)**
 $W := \text{nil};$
for all (ℓ, ℓ') in F do $W := \text{cons}((\ell, \ell'), W);$
for all ℓ in F or E do
 if $\ell \in E$ then $\text{Analysis}[\ell] := \iota$ else $\text{Analysis}[\ell] := \perp_L;$
- Step 2 **Iteration (updating W and Analysis)**
while $W \neq \text{nil}$ do
 $\ell := \text{fst}(\text{head}(W)); \ell' = \text{snd}(\text{head}(W)); W := \text{tail}(W);$
 if $f_\ell(\text{Analysis}[\ell]) \not\subseteq \text{Analysis}[\ell']$ then
 $\text{Analysis}[\ell'] := \text{Analysis}[\ell'] \sqcup f_\ell(\text{Analysis}[\ell]);$
 for all ℓ'' with (ℓ', ℓ'') in F do $W := \text{cons}((\ell', \ell''), W);$
- Step 3 **Presenting the result (MFP_o and MFP_\bullet)**
for all ℓ in F or E do
 $MFP_o(\ell) := \text{Analysis}[\ell];$
 $MFP_\bullet(\ell) := f_\ell(\text{Analysis}[\ell])$

156

© Nielsen^2, Hankin

Correctness

The worklist algorithm always terminates and it computes the least (or MFP) solution to the instance given as input.

Complexity

Suppose that E and F contain at most $b \geq 1$ distinct labels, that F contains at most $e \geq b$ pairs, and that L has finite height at most $h \geq 1$.

Count as basic operations the applications of f_ℓ , applications of \sqcup , or updates of Analysis.

Then there will be at most $O(e \cdot h)$ basic operations.

Example: Reaching Definitions (assuming unique labels):

$O(b^2)$ where b is size of program: $O(h) = O(b)$ and $O(e) = O(b)$.

157

© Nielsen^2, Hankin

The MOP Solution

– Idea: propagate analysis information along [paths](#).

Paths

The paths up to [but not including](#) ℓ :

$$\text{path}_o(\ell) = \{[\ell_1, \dots, \ell_{n-1}] \mid n \geq 1 \wedge \forall i < n : (\ell_i, \ell_{i+1}) \in F \wedge \ell_n = \ell \wedge \ell_1 \in E\}$$

The paths up to [and including](#) ℓ :

$$\text{path}_\bullet(\ell) = \{[\ell_1, \dots, \ell_n] \mid n \geq 1 \wedge \forall i < n : (\ell_i, \ell_{i+1}) \in F \wedge \ell_n = \ell \wedge \ell_1 \in E\}$$

Transfer functions for a path $\vec{\ell} = [\ell_1, \dots, \ell_n]$:

$$f_{\vec{\ell}} = f_{\ell_n} \circ \dots \circ f_{\ell_1} \circ id$$

158

© Nielsen^2, Hankin

The MOP Solution

The solution up to [but not including](#) ℓ :

$$MOP_o(\ell) = \bigsqcup \{f_{\vec{\ell}}(i) \mid \vec{\ell} \in \text{path}_o(\ell)\}$$

The solution up to [and including](#) ℓ :

$$MOP_\bullet(\ell) = \bigsqcup \{f_{\vec{\ell}}(i) \mid \vec{\ell} \in \text{path}_\bullet(\ell)\}$$

Precision of the MOP versus MFP solutions

The MFP solution safely approximates the MOP solution: $MFP \sqsupseteq MOP$ (“because” $f(x \sqcup y) \sqsupseteq f(x) \sqcup f(y)$ when f is monotone).

For Distributive Frameworks the MFP and MOP solutions are equal: $MFP = MOP$ (“because” $f(x \sqcup y) = f(x) \sqcup f(y)$ when f is distributive).

159

© Nielsen^2, Hankin

Lemma

Consider the MFP and MOP solutions to an instance $(L, \mathcal{F}, F, B, \iota, f.)$ of a Monotone Framework; then:

$$MFP_o \sqsupseteq MOP_o \text{ and } MFP_\bullet \sqsupseteq MOP_\bullet$$

If the framework is distributive and if $\text{path}_o(\ell) \neq \emptyset$ for all ℓ in E and F then:

$$MFP_o = MOP_o \text{ and } MFP_\bullet = MOP_\bullet$$

160

© Nielsen^2, Hankin

Decidability of MOP and MFP

The MFP solution is always computable (meaning that it is decidable) because of the **Ascending Chain Condition**.

The MOP solution is often uncomputable (meaning that it is undecidable): the existence of a general algorithm for the MOP solution would imply the decidability of the *Modified Post Correspondence Problem*, which is known to be undecidable.

161

© Nielsen^2, Hankin

Lemma

The MOP solution for Constant Propagation is undecidable.

Proof: Let u_1, \dots, u_n and v_1, \dots, v_n be strings over the alphabet $\{1, \dots, 9\}$; let $|u|$ denote the length of u ; let $\llbracket u \rrbracket$ be the natural number denoted.

The Modified Post Correspondence Problem is to determine whether or not $u_{i_1} \cdots u_{i_m} = v_{i_1} \cdots v_{i_n}$ for some sequence i_1, \dots, i_m with $i_1 = 1$.

```
x:=\llbracket u_1 \rrbracket; y:=\llbracket v_1 \rrbracket;
while [...] do
    if [...] then x:=x * 10^{|u_1|} + \llbracket u_1 \rrbracket; y:=y * 10^{|v_1|} + \llbracket v_1 \rrbracket else
        :
        if [...] then x:=x * 10^{|u_n|} + \llbracket u_n \rrbracket; y:=y * 10^{|v_n|} + \llbracket v_n \rrbracket else skip
    [z:=abs((x-y)*(x-y)))^\ell
```

Then $MOP_\bullet(\ell)$ will map z to 1 if and only if the Modified Post Correspondence Problem has no solution. This is undecidable.

162

© Nielsen^2, Hankin

Scaling Program Analysis

CS7575

- Lecture 3 — Data flow analysis; PoPA Chapter 2 (part 2) - Interprocedural + Shape

163

Interprocedural Analysis

- The problem
- MVP: “Meet” over Valid Paths
- Making context explicit
- Context based on call-strings
- Context based on assumption sets

```
proc add(val x, y, res z) is
    if (x is Integer) then
        z := x + y;
    else if (x is String) then
        z := concatenate(x,y);
    end

    add(40, 2, a);
    b := a+1;
    add("4", "2", s);
    println(s)
```

164

© Nielsen^2, Hankin

Syntax for procedures

Programs: $P_\star = \text{begin } D_\star \text{ end}$

Declarations: $D ::= D; D \mid \text{proc } p(\text{val } x; \text{res } y) \text{ is }^{\ell_n} S \text{ end }^{\ell_x}$

Statements: $S ::= \dots \mid [\text{call } p(a, z)]_{\ell_r}^{\ell_c}$

Example:

```
begin proc fib(val z, u; res v) is1
    if [z<3]2 then [v:=u+1]3
        else ([call fib(z-1,u,v)]4; [call fib(z-2,v,v)]6)
    end8;
    [call fib(x,0,y)]910
```

165

© Nielsen^2, Hankin

Flow graphs for procedure calls

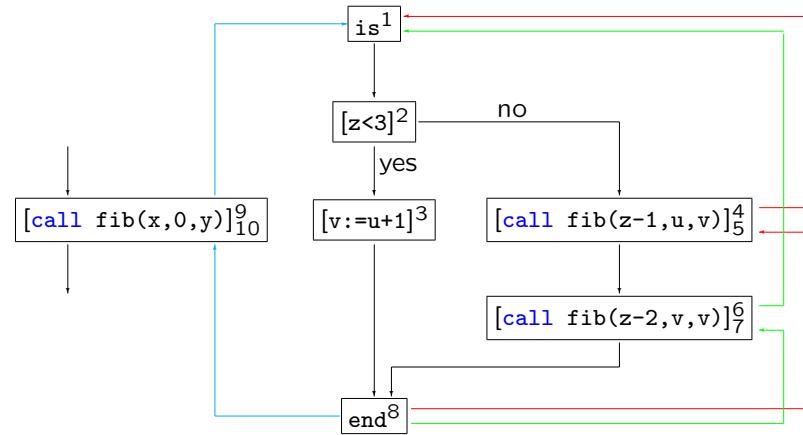
$\text{init}([\text{call } p(a, z)]_{\ell_r}^{\ell_c}) = \ell_c$
 $\text{final}([\text{call } p(a, z)]_{\ell_r}^{\ell_c}) = \{\ell_r\}$
 $\text{blocks}([\text{call } p(a, z)]_{\ell_r}^{\ell_c}) = \{[\text{call } p(a, z)]_{\ell_r}^{\ell_c}\}$
 $\text{labels}([\text{call } p(a, z)]_{\ell_r}^{\ell_c}) = \{\ell_c, \ell_r\}$
 $\text{flow}([\text{call } p(a, z)]_{\ell_r}^{\ell_c}) = \{(\ell_c; \ell_n), (\ell_x; \ell_r)\}$
 if $\text{proc } p(\text{val } x; \text{res } y) \text{ is }^{\ell_n} S \text{ end }^{\ell_x}$ is in D_\star

- $(\ell_c; \ell_n)$ is the flow corresponding to *calling* a procedure at ℓ_c and entering the procedure body at ℓ_n , and
- $(\ell_x; \ell_r)$ is the flow corresponding to exiting a procedure body at ℓ_x and *returning* to the call at ℓ_r .

167

© Nielsen^2, Hankin

proc fib(val z, u; res v)



Flow graphs for procedure declarations

For each procedure declaration $\text{proc } p(\text{val } x; \text{res } y) \text{ is }^{\ell_n} S \text{ end }^{\ell_x}$ of D_\star :

$\text{init}(p) = \ell_n$
 $\text{final}(p) = \{\ell_x\}$
 $\text{blocks}(p) = \{\text{is }^{\ell_n}, \text{end }^{\ell_x}\} \cup \text{blocks}(S)$
 $\text{labels}(p) = \{\ell_n, \ell_x\} \cup \text{labels}(S)$
 $\text{flow}(p) = \{(\ell_n, \text{init}(S))\} \cup \text{flow}(S) \cup \{(\ell, \ell_x) \mid \ell \in \text{final}(S)\}$

168

© Nielsen^2, Hankin

Flow graphs for programs

For the program $P_\star = \text{begin } D_\star \text{ end}$:

$$\begin{aligned} \text{init}_\star &= \text{init}(S_\star) \\ \text{final}_\star &= \text{final}(S_\star) \\ \text{blocks}_\star &= \bigcup \{ \text{blocks}(p) \mid \text{proc } p(\text{val } x; \text{res } y) \text{ is } \ell_n S \text{ end } \ell_x \text{ is in } D_\star \} \\ &\quad \cup \text{blocks}(S_\star) \\ \text{labels}_\star &= \bigcup \{ \text{labels}(p) \mid \text{proc } p(\text{val } x; \text{res } y) \text{ is } \ell_n S \text{ end } \ell_x \text{ is in } D_\star \} \\ &\quad \cup \text{labels}(S_\star) \\ \text{flow}_\star &= \bigcup \{ \text{flow}(p) \mid \text{proc } p(\text{val } x; \text{res } y) \text{ is } \ell_n S \text{ end } \ell_x \text{ is in } D_\star \} \\ &\quad \cup \text{flow}(S_\star) \\ \text{interflow}_\star &= \{ (\ell_c, \ell_n, \ell_x, \ell_r) \mid \text{proc } p(\text{val } x; \text{res } y) \text{ is } \ell_n S \text{ end } \ell_x \text{ is in } D_\star \\ &\quad \text{and } [\text{call } p(a, z)]_{\ell_r}^{\ell_c} \text{ is in } S_\star \} \end{aligned}$$

169

© Nielsen^2, Hankin

Example:

```
begin proc fib(val z, u; res v) is1
    if [z<3]2 then [v:=u+1]3
    else ([call fib(z-1,u,v)]4; [call fib(z-2,v,v)]5)
end8;
[call fib(x,0,y)]910
```

We have

$$\begin{aligned} \text{flow}_\star &= \{(1,2), (2,3), (3,8), \\ &\quad (2,4), (4;1), (8;5), (5,6), (6;1), (8;7), (7,8), \\ &\quad (9;1), (8;10)\} \end{aligned}$$

$$\text{interflow}_\star = \{(9, 1, 8, 10), (4, 1, 8, 5), (6, 1, 8, 7)\}$$

and $\text{init}_\star = 9$ and $\text{final}_\star = \{10\}$.

170

© Nielsen^2, Hankin

Intermezzo

PEPM 1991: New Haven, Connecticut, USA 

Demand-Driven Static Backward Slicing of Unstructured Programs 

Home > Conferences and Workshops > PEPM

Charles Consel, Olivier Danvy: Husni Khanfar, Björn Lisper, Abu Naser Masud: Husni Khanfar, Björn Lisper, Saad Mubeen: Sebastian Altmeier, Björn Lisper, Claire Maiza, Jan Reineke, Christine Rochange: Niklas Holst: 

2015

[c56]  Charles Consel, Olivier Danvy: Husni Khanfar, Björn Lisper, Abu Naser Masud:                                         <img alt="document icon" data-bbox="7015 6

MVP: “Meet” over Valid Paths

Complete Paths

We need to match procedure entries and exits:

A *complete path* from ℓ_1 to ℓ_2 in P_* has proper nesting of procedure entries and exits; and a procedure returns to the point where it was called:

$$\begin{aligned} CP_{\ell_1, \ell_2} &\longrightarrow \ell_1 && \text{whenever } \ell_1 = \ell_2 \\ CP_{\ell_1, \ell_3} &\longrightarrow \ell_1, CP_{\ell_2, \ell_3} && \text{whenever } (\ell_1, \ell_2) \in \text{flow}_* \\ CP_{\ell_c, \ell} &\longrightarrow \ell_c, CP_{\ell_n, \ell_x}, CP_{\ell_r, \ell} && \text{whenever } P_* \text{ contains [call } p(a, z)]_{\ell_r}^{\ell_c} \\ &&& \text{and proc } p(\text{val } x; \text{res } y) \text{ is } \ell_n S \text{ end }^{\ell_x} \end{aligned}$$

More generally: whenever $(\ell_c, \ell_n, \ell_x, \ell_r)$ is an element of interflow_* (or interflow_*^R for backward analyses); see the book.

173

© Nielsen^2, Hankin

Valid Paths

A *valid path* starts at the entry node init_* of P_* , all the procedure exits match the procedure entries but some procedures might be entered but not yet exited:

$$\begin{aligned} VP_* &\longrightarrow VP_{\text{init}_*, \ell} && \text{whenever } \ell \in \text{Lab}_* \\ VP_{\ell_1, \ell_2} &\longrightarrow \ell_1 && \text{whenever } \ell_1 = \ell_2 \\ VP_{\ell_1, \ell_3} &\longrightarrow \ell_1, VP_{\ell_2, \ell_3} && \text{whenever } (\ell_1, \ell_2) \in \text{flow}_* \\ VP_{\ell_c, \ell} &\longrightarrow \ell_c, CP_{\ell_n, \ell_x}, VP_{\ell_r, \ell} && \text{whenever } P_* \text{ contains [call } p(a, z)]_{\ell_r}^{\ell_c} \\ &&& \text{and proc } p(\text{val } x; \text{res } y) \text{ is } \ell_n S \text{ end }^{\ell_x} \\ VP_{\ell_c, \ell} &\longrightarrow \ell_c, VP_{\ell_n, \ell} && \text{whenever } P_* \text{ contains [call } p(a, z)]_{\ell_r}^{\ell_c} \\ &&& \text{and proc } p(\text{val } x; \text{res } y) \text{ is } \ell_n S \text{ end }^{\ell_x} \end{aligned}$$

174

© Nielsen^2, Hankin

The MVP solution

$$MVP_o(\ell) = \bigsqcup \{f_\ell(\iota) \mid \vec{\ell} \in vpath_o(\ell)\}$$

where

$$vpath_o(\ell) = \{[\ell_1, \dots, \ell_{n-1}] \mid n \geq 1 \wedge \ell_n = \ell \wedge [\ell_1, \dots, \ell_n] \text{ is a valid path}\}$$

The MVP solution may be undecidable for lattices satisfying the Ascending Chain Condition, just as was the case for the MOP solution.

175

© Nielsen^2, Hankin

The MVP solution

$$MVP_o(\ell) = \bigsqcup \{f_\ell(\iota) \mid \vec{\ell} \in vpath_o(\ell)\}$$

where

$$vpath_o(\ell) = \{[\ell_1, \dots, \ell_{n-1}] \mid n \geq 1 \wedge \ell_n = \ell \wedge [\ell_1, \dots, \ell_n] \text{ is a valid path}\}$$

The MVP solution may be undecidable for lattices satisfying the Ascending Chain Condition, just as was the case for the MOP solution.

176

© Nielsen^2, Hankin

QUESTIONS?

Making Context Explicit

Starting point: an instance $(L, \mathcal{F}, F, E, \iota, f.)$ of a Monotone Framework

- the analysis is **forwards**, i.e. $F = \text{flow}_*$ and $E = \{\text{init}_*\}$;
- the complete lattice is a powerset, i.e. $L = \mathcal{P}(D)$;
- the transfer functions in \mathcal{F} are completely additive; and
- each f_ℓ is given by $f_\ell(Y) = \bigcup\{\phi_\ell(d) \mid d \in Y\}$ where $\phi_\ell : D \rightarrow \mathcal{P}(D)$.

177

© Nielsen^2, Hankin

An embellished monotone framework

- $L' = \mathcal{P}(\Delta \times D)$;
- the transfer functions in \mathcal{F}' are completely additive; and
- each f'_ℓ is given by $f'_\ell(Z) = \bigcup\{\{\delta\} \times \phi_\ell(d) \mid (\delta, d) \in Z\}$.

Ignoring procedures, the data flow equations will take the form:

$$\begin{aligned} A_\bullet(\ell) &= f'_\ell(A_\circ(\ell)) \\ &\quad \text{for all labels that do not label a procedure call} \\ A_\circ(\ell) &= \bigsqcup\{A_\bullet(\ell') \mid (\ell', \ell) \in F \text{ or } (\ell'; \ell) \in F\} \sqcup \iota_E'' \\ &\quad \text{for all labels (including those that label procedure calls)} \end{aligned}$$

178

© Nielsen^2, Hankin

Example:

Detection of Signs Analysis as a Monotone Framework:

$(L_{\text{sign}}, \mathcal{F}_{\text{sign}}, F, E, \iota_{\text{sign}}, f^{\text{sign}})$ where $\text{Sign} = \{-, 0, +\}$ and

$$L_{\text{sign}} = \mathcal{P}(\text{Var}_* \rightarrow \text{Sign})$$

The transfer function f_ℓ^{sign} associated with the assignment $[x := a]^\ell$ is

$$f_\ell^{\text{sign}}(Y) = \bigcup\{\phi_\ell^{\text{sign}}(\sigma^{\text{sign}}) \mid \sigma^{\text{sign}} \in Y\}$$

where $Y \subseteq \text{Var}_* \rightarrow \text{Sign}$ and

$$\phi_\ell^{\text{sign}}(\sigma^{\text{sign}}) = \{\sigma^{\text{sign}}[x \mapsto s] \mid s \in \mathcal{A}_{\text{sign}}[a](\sigma^{\text{sign}})\}$$

179

© Nielsen^2, Hankin

Example (cont.):

Detection of Signs Analysis as an embellished monotone framework

$$L'_{\text{sign}} = \mathcal{P}(\Delta \times (\text{Var}_* \rightarrow \text{Sign}))$$

The transfer function associated with $[x := a]^\ell$ will now be:

$$f_\ell^{\text{sign}'}(Z) = \bigcup\{\{\delta\} \times \phi_\ell^{\text{sign}}(\sigma^{\text{sign}}) \mid (\delta, \sigma^{\text{sign}}) \in Z\}$$

180

© Nielsen^2, Hankin

Transfer functions for procedure declarations

Procedure declarations

$\text{proc } p(\text{val } x; \text{res } y) \text{ is}^{\ell_n} S \text{ end}^{\ell_x}$

have two transfer functions, one for entry and one for exit:

$f_{\ell_n}, f_{\ell_x} : \mathcal{P}(\Delta \times D) \rightarrow \mathcal{P}(\Delta \times D)$

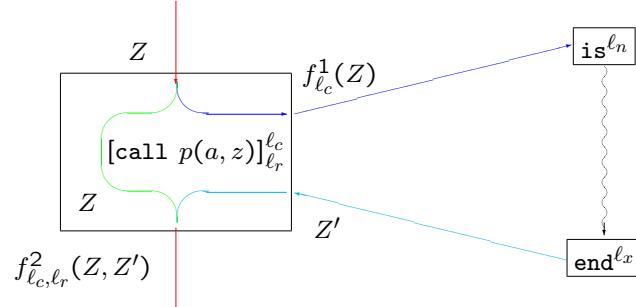
For simplicity we take both to be the identity function (thus incorporating procedure entry as part of procedure call, and procedure exit as part of procedure return).

181

© Nielsen^2, Hankin

Procedure calls and returns

$\text{proc } p(\text{val } x; \text{res } y)$



182

© Nielsen^2, Hankin

Transfer functions for procedure calls

Procedure calls $[\text{call } p(a, z)]_{\ell_c}^{\ell_c}$ have two transfer functions:

For the *procedure call*

$f_{\ell_c}^1 : \mathcal{P}(\Delta \times D) \rightarrow \mathcal{P}(\Delta \times D)$

and it is used in the equation:

$$A_{\bullet}(\ell_c) = f_{\ell_c}^1(A_{\circ}(\ell_c)) \text{ for all procedure calls } [\text{call } p(a, z)]_{\ell_c}^{\ell_c}$$

For the *procedure return*

$f_{\ell_c, \ell_r}^2 : \mathcal{P}(\Delta \times D) \times \mathcal{P}(\Delta \times D) \rightarrow \mathcal{P}(\Delta \times D)$

and it is used in the equation:

$$A_{\bullet}(\ell_r) = f_{\ell_c, \ell_r}^2(A_{\circ}(\ell_c), A_{\circ}(\ell_r)) \text{ for all procedure calls } [\text{call } p(a, z)]_{\ell_c}^{\ell_c}$$

(Note that $A_{\circ}(\ell_r)$ will equal $A_{\bullet}(\ell_x)$ for the relevant procedure exit.)

© Nielsen^2, Hankin

Different Kinds of Context

- **Call Strings** — contexts based on control

– Call strings of unbounded length

– Call strings of bounded length (k)

- **Assumption Sets** — contexts based on data

– Large assumption sets ($k = 1$)

184

© Nielsen^2, Hankin

Call Strings of Unbounded Length

$$\Delta = \text{Lab}^*$$

Transfer functions for procedure call

$$f_{\ell_c}^1(Z) = \bigcup \{ \{\delta'\} \times \phi_{\ell_c}^1(d) \mid (\delta, d) \in Z \wedge \delta' = [\delta, \ell_c] \}$$

$$f_{\ell_c, \ell_r}^2(Z, Z') = \bigcup \{ \{\delta\} \times \phi_{\ell_c, \ell_r}^2(d, d') \mid (\delta, d) \in Z \wedge (\delta', d') \in Z' \wedge \delta' = [\delta, \ell_c] \}$$

185

© Nielsen^2, Hankin

Example:

Recalling the statements:

$$\text{proc } p(\text{val } x; \text{res } y) \text{ is } S \text{ end } \ell_x \quad [\text{call } p(a, z)]_{\ell_r}^{\ell_c}$$

Detection of Signs Analysis:

$$\phi_{\ell_c}^{\text{sign}1}(\sigma^{\text{sign}}) = \{ \sigma^{\text{sign}} \underbrace{[x \mapsto s][y \mapsto s']}_{\text{initialise formals}} \mid s \in \mathcal{A}_{\text{sign}}[[a]](\sigma^{\text{sign}}), s' \in \{-, 0, +\} \}$$

$$\phi_{\ell_c, \ell_r}^{\text{sign}2}(\sigma_1^{\text{sign}}, \sigma_2^{\text{sign}}) = \{ \sigma_2^{\text{sign}} \underbrace{[x \mapsto \sigma_1^{\text{sign}}(x)][y \mapsto \sigma_1^{\text{sign}}(y)]}_{\text{restore formals}} \underbrace{[z \mapsto \sigma_2^{\text{sign}}(y)]}_{\text{return result}} \mid \dots \}$$

© Nielsen^2, Hankin

Example:

Recalling the statements:

$$\text{proc } p(\text{val } x; \text{res } y) \text{ is } S \text{ end } \ell_x \quad [\text{call } p(a, z)]_{\ell_r}^{\ell_c}$$

Detection of Signs Analysis:

$$\phi_{\ell_c}^{\text{sign}1}(\sigma^{\text{sign}}) = \{ \sigma^{\text{sign}} \underbrace{[x \mapsto s][y \mapsto s']}_{\text{initialise formals}} \mid s \in \mathcal{A}_{\text{sign}}[[a]](\sigma^{\text{sign}}), s' \in \{-, 0, +\} \}$$

$$\phi_{\ell_c, \ell_r}^{\text{sign}2}(\sigma_1^{\text{sign}}, \sigma_2^{\text{sign}}) = \{ \sigma_2^{\text{sign}} \underbrace{[x \mapsto \sigma_1^{\text{sign}}(x)][y \mapsto \sigma_1^{\text{sign}}(y)]}_{\text{restore formals}} \underbrace{[z \mapsto \sigma_2^{\text{sign}}(y)]}_{\text{return result}} \mid \dots \}$$

© Nielsen^2, Hankin

Call Strings of Bounded Length

$$\Delta = \text{Lab}^{\leq k}$$

Transfer functions for procedure call

$$f_{\ell_c}^1(Z) = \bigcup \{ \{\delta'\} \times \phi_{\ell_c}^1(d) \mid (\delta, d) \in Z \wedge \delta' = [\delta, \ell_c] \}$$

$$f_{\ell_c, \ell_r}^2(Z, Z') = \bigcup \{ \{\delta\} \times \phi_{\ell_c, \ell_r}^2(d, d') \mid (\delta, d) \in Z \wedge (\delta', d') \in Z' \wedge \delta' = [\delta, \ell_c] \}$$

188

© Nielsen^2, Hankin

QUESTIONS?

A special case: call strings of length $k = 0$

$$\Delta = \{\Lambda\}$$

Note: this is equivalent to having no context information!

Specialising the transfer functions:

$$f_{\ell_c}^1(Y) = \bigcup \{\phi_{\ell_c}^1(d) \mid d \in Y\}$$

$$f_{\ell_c, \ell_r}^2(Y, Y') = \bigcup \{\phi_{\ell_c, \ell_r}^2(d, d') \mid d \in Y \wedge d' \in Y'\}$$

(We use that $\mathcal{P}(\Delta \times D)$ isomorphic to $\mathcal{P}(D)$.)

189

© Nielsen^2, Hankin

A special case: call strings of length $k = 1$

$$\Delta = \text{Lab} \cup \{\Lambda\}$$

Specialising the transfer functions:

$$f_{\ell_c}^1(Z) = \bigcup \{\{\ell_c\} \times \phi_{\ell_c}^1(d) \mid (\delta, d) \in Z\}$$

$$f_{\ell_c, \ell_r}^2(Z, Z') = \bigcup \{\{\delta\} \times \phi_{\ell_c, \ell_r}^2(d, d') \mid (\delta, d) \in Z \wedge (\ell_c, d') \in Z'\}$$

190

© Nielsen^2, Hankin

Large Assumption Sets ($k = 1$)

$$\Delta = \mathcal{P}(D)$$

Transfer functions for procedure call

$$f_{\ell_c}^1(Z) = \bigcup \{\{\delta'\} \times \phi_{\ell_c}^1(d) \mid (\delta, d) \in Z \wedge \delta' = \{d'' \mid (\delta, d'') \in Z\}\}$$

$$f_{\ell_c, \ell_r}^2(Z, Z') = \bigcup \{\{\delta\} \times \phi_{\ell_c, \ell_r}^2(d, d') \mid (\delta, d) \in Z \wedge (\delta', d') \in Z' \wedge \delta' = \{d'' \mid (\delta, d'') \in Z\}\}$$

191

© Nielsen^2, Hankin

Intermezzo: PEPM'91

$$\begin{aligned} \text{Ans} &= \text{LAB} \rightarrow \mathcal{P}(\text{LAB} + \text{PRIM}) \\ \text{D} &= \mathcal{P}(\text{D})^* \rightarrow \text{VEnv} \rightarrow \text{LAB} \rightarrow \text{Ans} \\ \text{VEnv} &= \text{VAR} \rightarrow \mathcal{P}(\text{D}) \end{aligned}$$

$$\begin{aligned} \mathcal{P}R\ell &= f \left\langle \left\{ \lambda a v_x e_x \ell_c, [e_c \mapsto \{\ell_{\text{op}}\}] \right\} \right\rangle \sqcup \text{c}_\text{top} \\ \text{where } \{f\} &= \mathcal{A} \ell \{ \} \\ \mathcal{A}[\text{c}_x] e_x &= e v \\ \mathcal{A}[\text{b}_x] e_x &= \emptyset \\ \mathcal{A}[\text{p}_x] e_x &= \{ \text{p} p \} \\ \mathcal{A}[\ell : \lambda (v_1 \dots v_n) c] e_x &= \begin{cases} \lambda a v_x \ell_c, [e_x \mapsto \{f\}] \sqcup \text{length}(av) = n \longrightarrow \mathcal{C} c e' \\ \text{otherwise } [] \end{cases} \\ \text{where } e' &= e \sqcup [v_i \mapsto av[i]] \end{aligned}$$

$$\begin{aligned} \mathcal{C}[c : (f a_1 \dots a_n) e] &= \bigcup \{f'(a'_1 \dots a'_n) e' c \mid f' \in \mathcal{A} f, e' \\ \text{where } a'_i &= \mathcal{A} a_i, e' \\ \mathcal{C}[\text{f} : (\text{letrec } ((f_1 t_1) \dots) c)] e &= \mathcal{C} c e' \\ \text{where } e' &= e \sqcup [f_i \mapsto \mathcal{A} t_i, e] \end{aligned}$$

$$\begin{aligned} \mathcal{P}[\star] &= \lambda (a b c) e \ell_c, [\ell_c \mapsto \{+\}] \sqcup \text{bad argument} \longrightarrow [] \\ \text{otherwise } &\bigcup \{c'(\emptyset) e i c_{t_c, t_c} \mid c' \in c\} \\ \mathcal{P}[\text{if}] &= \lambda (p c a) e \ell_c, [\ell_c \mapsto \{\text{if}\}] \sqcup \begin{pmatrix} \text{bad argument} \longrightarrow [] \\ \text{otherwise } \bigcup \{c'(\emptyset) e i c_{t_c, t_c} \mid c' \in c\} \\ \sqcup \bigcup \{a'(\emptyset) e i c_{t_c, t_c} \mid a' \in a\} \end{pmatrix} \end{aligned}$$

Figure 6: OCFA

Shivers. The Semantics of Scheme Control-Flow Analysis, PEPM'91 <https://doi.org/10.1145/115865.115884>

190

Shape Analysis

Goal: to obtain a finite representation of the shape of the heap of a language with pointers.

The analysis result can be used for

- detection of pointer aliasing
- detection of sharing between structures
- software development tools
 - detection of errors like dereferences of nil-pointers
- program verification
 - reverse transforms a non-cyclic list to a non-cyclic list

193

© Nielsen^2, Hankin

Syntax of the pointer language

$$\begin{aligned} a &::= p \mid n \mid a_1 \text{ op}_a a_2 \mid \text{nil} \\ p &::= x \mid x.\text{sel} \\ b &::= \text{true} \mid \text{false} \mid \text{not } b \mid b_1 \text{ op}_b b_2 \mid a_1 \text{ op}_r a_2 \mid \text{op}_p p \\ S &::= [p:=a]^\ell \mid [\text{skip}]^\ell \mid S_1; S_2 \mid \\ &\quad \text{if } [b]^\ell \text{ then } S_1 \text{ else } S_2 \mid \text{while } [b]^\ell \text{ do } S \mid \\ &\quad [\text{malloc } p]^\ell \end{aligned}$$

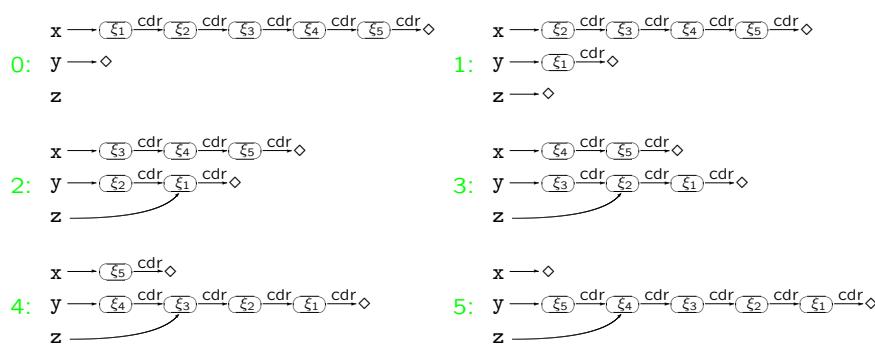
Example

```
[y:=nil]1;
while [not is-nil(x)]2 do
  ([z:=y]3; [y:=x]4; [x:=x.cdr]5; [y.cdr:=z]6);
[z:=nil]7
```

194

© Nielsen^2, Hankin

Reversal of a list



195

© Nielsen^2, Hankin

Intermezzo: meaning of car/cdr

The IBM 704 had 36-bit words and a 15-bit address space. ... The first 15-bit field was the operand address and the second held a decrement. ... Indexing was a subtractive process on the 704, hence the value to be loaded into an index register was called a "decrement". The 704 hardware had special instructions for accessing the address and decrement fields in a word. As a result it was efficient to use those two fields to store within a single word the two pointers needed for a list.

Thus, "CAR" is "Contents of the Address part of the Register". The term "register" in this context refers to "memory location".

Precursors to Lisp included functions:

- car ("contents of the address part of register number"),
 - cdr ("contents of the decrement part of register number"),
- each of which took a machine address as an argument, loaded the corresponding word from memory, and extracted the appropriate bits.

Wikipedia
H. Gelernter, J. R. Hansen, C. L. Gerberich. "A Fortran-Compiled List-Processing Language". J. ACM 1960 doi:10.1145/321021.321022

A Fortran-Compiled List-Processing Language*

H. GELERTNER, J. R. HANSEN, AND C. L. GERBERICH
International Business Machines Corp., Yorktown Heights, N.Y.

Abstract. A compiled computer language for the manipulation of symbolic expressions organized in storage as Newell-Shaw-Simon lists has been developed as a tool to move more common list-processing programming to the simulation of a geometry theorem-proving machine on the IBM 704 high-speed electronic digital computer. Statements in the language are similar to standard Fortran statements, but include additional statements and functions appended to the standard Fortran library. The algebraic structure of certain statements in this language corresponds directly to the structure of a list in memory pointers to the head and tail of the list. The language is designed to be used with a variable stepsize. The many programming advantages resulting from the use of Fortran, and in particular, the ability to use the facilities of the compiler to generate efficient programs, and the flexibility offered by an NSS list organization of storage make the language particularly useful where, as in the case of our theorem-proving program, intermediate data of unpredictable form, complexity, and length may be generated.

1. Introduction

Until recently, digital computer design has been strongly oriented toward increased speed and facility in the arithmetic manipulation of numbers, for it is in this mode of operation that most calculations are performed. With greater appreciation of the ultimate capacity of such machines, however, and with increased understanding of the techniques of list manipulation, many computers now have been developed which deal largely with symbolic data. These are purely symbolic and processes that are logicistic rather than arithmetic. One such effort is the simulation of a geometry theorem-proving machine being investigated by Newell, Shaw, and Simon at the University of Michigan (see below). The many programming advantages resulting from the use of Fortran, and in particular, the ability offered by an NSS list organization of storage make the language particularly useful where, as in the case of our theorem-proving program, intermediate data of unpredictable form, complexity, and length may be generated.

To simulate the geometry theorem-proving machine, the authors have developed a massive simulation program with a characteristic feature in common with many other such symbol-manipulating routines, and in particular, with those intended to carry out some algorithmic processes by the use of lists of items. The lists of items in these programs are generally unpredictable in their form, complexity, and length. Arbitrary lists of information may or may not contain as data an arbitrary number of items or symbols. To simulate beforehand to each possible list a block of storage sufficient to contain symbolic material may be difficult, if not impossible, since one cannot exhaust all available fast-access storage as well as prescribe rigidly the organization of information in the list. A program failure caused by some list exceeding its allocated space will leave a remainder of storage in almost empty could be expected as a not uncommon occurrence.

Faced with this problem, Newell, Shaw, and Simon, in programming a heuristic theorem-proving system for the propositional calculus, simulated by program-

* Received April, 1959. Part of the material contained herein was presented at the meeting of the Association, September 1-3, 1959.

87

Structural Operational Semantics

A configurations consists of

- a state $\sigma \in \text{State} = \text{Var}_\star \rightarrow (\mathbf{Z} + \text{Loc} + \{\diamond\})$
mapping variables to values, locations (in the heap) or the nil-value
- a heap $\mathcal{H} \in \text{Heap} = (\text{Loc} \times \text{Sel}) \rightarrow_{\text{fin}} (\mathbf{Z} + \text{Loc} + \{\diamond\})$
mapping pairs of locations and selectors to values, locations in the heap or the nil-value

197

© Nielsen^2, Hankin

Pointer expressions

$$\wp : \text{PExp} \rightarrow (\text{State} \times \text{Heap}) \rightarrow_{\text{fin}} (\mathbf{Z} + \{\diamond\} + \text{Loc})$$

is defined by

$$\begin{aligned}\wp[x](\sigma, \mathcal{H}) &= \sigma(x) \\ \wp[x.\text{sel}](\sigma, \mathcal{H}) &= \begin{cases} \mathcal{H}(\sigma(x), \text{sel}) & \text{if } \sigma(x) \in \text{Loc} \text{ and } \mathcal{H} \text{ is defined on } (\sigma(x), \text{sel}) \\ \text{undefined} & \text{otherwise} \end{cases}\end{aligned}$$

Arithmetic and boolean expressions

$$\mathcal{A} : \text{AExp} \rightarrow (\text{State} \times \text{Heap}) \rightarrow_{\text{fin}} (\mathbf{Z} + \text{Loc} + \{\diamond\})$$

$$\mathcal{B} : \text{BExp} \rightarrow (\text{State} \times \text{Heap}) \rightarrow_{\text{fin}} \mathbf{T}$$

198

© Nielsen^2, Hankin

Statements

Clauses for assignments:

$$\langle [x := a]^\ell, \sigma, \mathcal{H} \rangle \rightarrow \langle \sigma[x \mapsto \mathcal{A}[a](\sigma, \mathcal{H})], \mathcal{H} \rangle$$

if $\mathcal{A}[a](\sigma, \mathcal{H})$ is defined

$$\langle [x.\text{sel} := a]^\ell, \sigma, \mathcal{H} \rangle \rightarrow \langle \sigma, \mathcal{H}[(\sigma(x), \text{sel}) \mapsto \mathcal{A}[a](\sigma, \mathcal{H})] \rangle$$

if $\sigma(x) \in \text{Loc}$ and $\mathcal{A}[a](\sigma, \mathcal{H})$ is defined

Clauses for malloc:

$$\langle [\text{malloc } x]^\ell, \sigma, \mathcal{H} \rangle \rightarrow \langle \sigma[x \mapsto \xi], \mathcal{H} \rangle$$

where ξ does not occur in σ or \mathcal{H}

$$\langle [\text{malloc } (x.\text{sel})]^\ell, \sigma, \mathcal{H} \rangle \rightarrow \langle \sigma, \mathcal{H}[(\sigma(x), \text{sel}) \mapsto \xi] \rangle$$

where ξ does not occur in σ or \mathcal{H} and $\sigma(x) \in \text{Loc}$

199

© Nielsen^2, Hankin

Shape graphs

The analysis will operate on *shape graphs* (S, H, is) consisting of

- an abstract state, S ,
- an abstract heap, H , and
- sharing information, is , for the abstract locations.

The nodes of the shape graphs are *abstract locations*:

$$\text{ALoc} = \{n_X \mid X \subseteq \text{Var}_\star\}$$

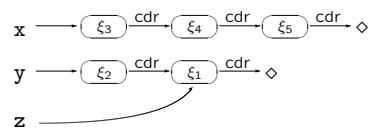
Note: there will only be *finitely* many abstract locations

200

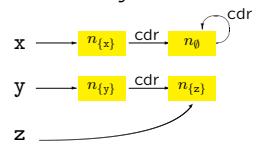
© Nielsen^2, Hankin

Example

In the semantics:



In the analysis:



Abstract Locations

The abstract location n_X represents the location $\sigma(x)$ if $x \in X$

The abstract location n_\emptyset is called the *abstract summary location*: n_\emptyset represents all the locations that cannot be reached directly from the state without consulting the heap

Invariant 1 If two abstract locations n_X and n_Y occur in the same shape graph then either $X = Y$ or $X \cap Y = \emptyset$

201

© Nielsen^2, Hankin

Abstract states and heaps

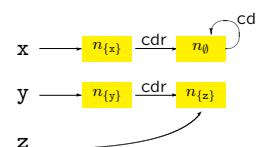
$$S \in AState = \mathcal{P}(\text{Var}_* \times ALoc)$$

abstract states

$$H \in AHeap = \mathcal{P}(ALoc \times Sel \times ALoc)$$

abstract heap

Invariant 2 If x is mapped to n_X by the abstract state S then $x \in X$

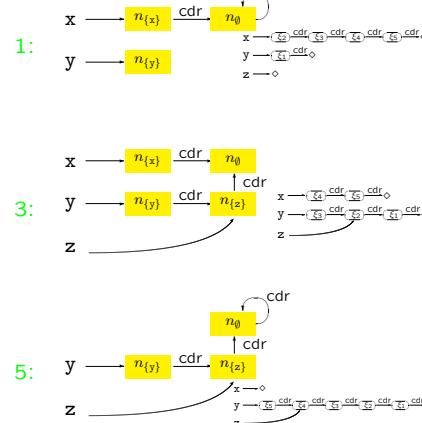
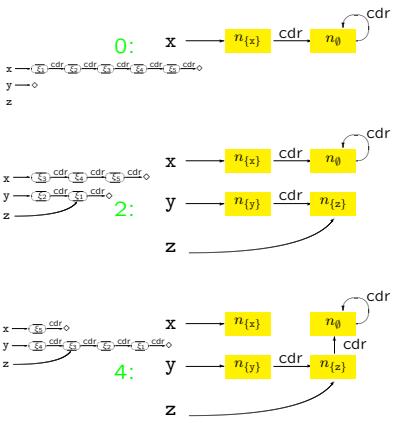


Invariant 3 Whenever (n_V, sel, n_W) and $(n_{V'}, sel', n_{W'})$ are in the abstract heap H then either $V = \emptyset$ or $W = W'$

202

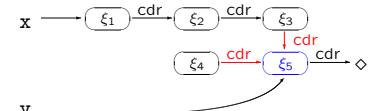
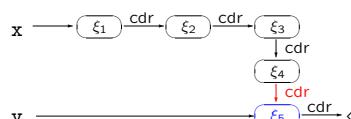
© Nielsen^2, Hankin

Reversal of a list

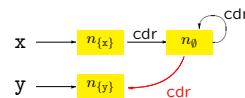


203

Sharing in the heap



Give rise to the same shape graph: is: the abstract locations that *might* be shared due to pointers in the heap:

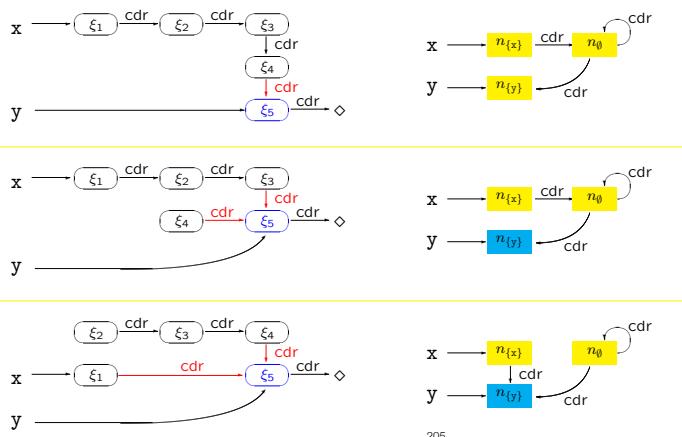


n_X is included in is if it might represent a location that is the target of more than one pointer in the heap

204

© Nielsen^2, Hankin

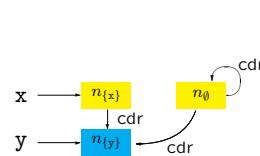
Examples: sharing in the heap



© Nielsen^2, Hankin

Sharing information

The **implicit** sharing information of the abstract heap must be consistent with the **explicit** sharing information:



Invariant 4 If $n_X \in \text{is}$ then either

- $(n_\emptyset, \text{sel}, n_X)$ is in the abstract heap for some sel , or
- there are two distinct triples (n_V, sel_1, n_X) and (n_W, sel_2, n_X) in the abstract heap

Invariant 5 Whenever there are two distinct triples (n_V, sel_1, n_X) and (n_W, sel_2, n_X) in the abstract heap and $X \neq \emptyset$ then $n_X \in \text{is}$

206

© Nielsen^2, Hankin

The complete lattice of shape graphs

A *shape graph* is a triple (S, H, is) where

$$\begin{aligned} S &\in \text{AState} = \mathcal{P}(\text{Var}_* \times \text{ALoc}) \\ H &\in \text{AHeap} = \mathcal{P}(\text{ALoc} \times \text{Sel} \times \text{ALoc}) \\ \text{is} &\in \text{IsShared} = \mathcal{P}(\text{ALoc}) \end{aligned}$$

and $\text{ALoc} = \{n_Z \mid Z \subseteq \text{Var}_*\}$.

A shape graph (S, H, is) is *compatible* if it fulfills the five invariants.

The analysis computes over *sets of compatible shape graphs*

$$\text{SG} = \{(S, H, \text{is}) \mid (S, H, \text{is}) \text{ is compatible}\}$$

207

© Nielsen^2, Hankin

The analysis

An instance of a *forward* Monotone Framework with the complete lattice of interest being $\mathcal{P}(\text{SG})$

A *may analysis*: each of the sets of shape graphs computed by the analysis may contain shape graphs that cannot really arise

Aspects of a *must analysis*: each of the individual shape graphs (in a set of shape graphs computed by the analysis) will be the best possible description of some (σ, \mathcal{H})

208

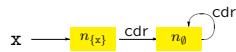
© Nielsen^2, Hankin

The analysis

Equations:

$$\begin{aligned} \text{Shape}_c(\ell) &= \begin{cases} \iota & \text{if } \ell = \text{init}(S_*) \\ \cup\{\text{Shape}_c(\ell') \mid (\ell', \ell) \in \text{flow}(S_*)\} & \text{otherwise} \end{cases} \\ \text{Shape}_\bullet(\ell) &= f_\ell^{\text{SA}}(\text{Shape}_c(\ell)) \end{aligned}$$

Example: The extremal value ι for the list reversal program

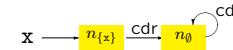


- x points to a non-cyclic list with at least three elements

209

© Nielsen^2, Hankin

$\text{Shape}_\bullet(1)$ for $[y := \text{nil}]^1$

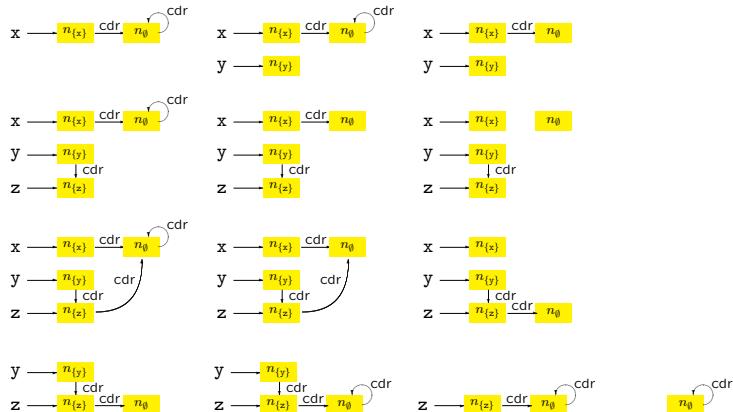


Note: we do not record nil-values in the analysis

210

© Nielsen^2, Hankin

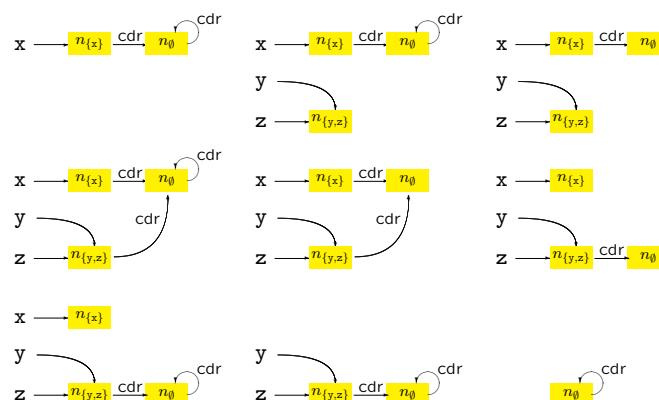
$\text{Shape}_\bullet(2)$ for $[\text{not is-nil}(x)]^2$



211

© Nielsen^2, Hankin

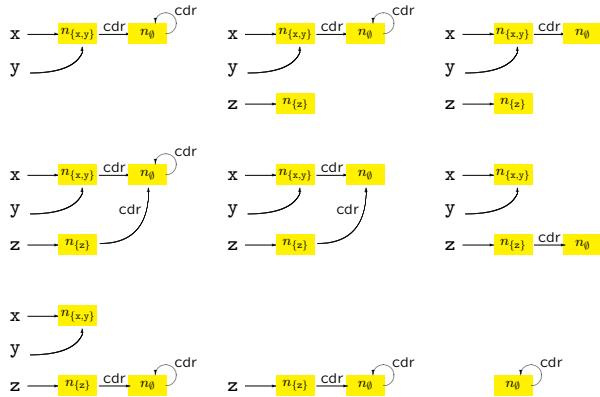
$\text{Shape}_\bullet(3)$ for $[z := y]^3$



212

© Nielsen^2, Hankin

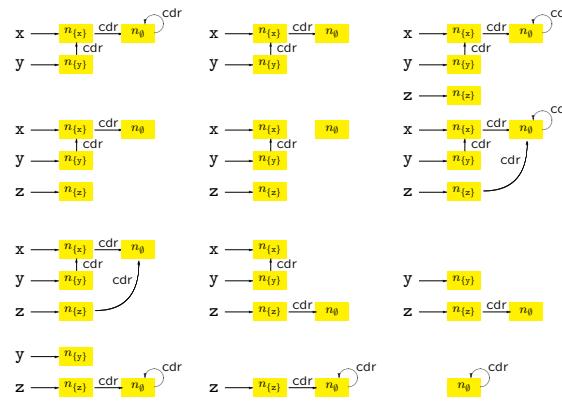
Shape_•(4) for $[y:=x]^4$



213

© Nielsen^2, Hankin

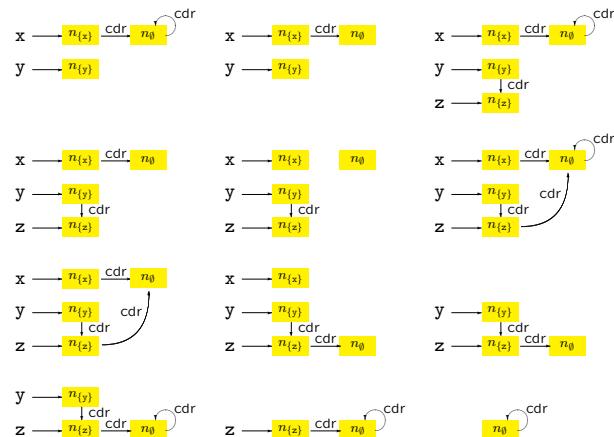
Shape_•(5) for $[x:=x.cdr]^5$



214

© Nielsen^2, Hankin

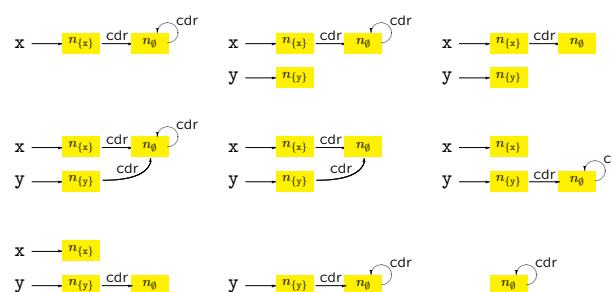
Shape_•(6) for $[y.cdr:=z]^6$



215

© Nielsen^2, Hankin

Shape_•(7) for $[z:=nil]^7$



– upon termination y points to a non-circular list

– a more precise analysis taking tests into account will know that x is nil upon termination

216

© Nielsen^2, Hankin

Transfer functions

$$f_\ell^{\text{SA}} : \mathcal{P}(\text{SG}) \rightarrow \mathcal{P}(\text{SG})$$

has the form:

$$f_\ell^{\text{SA}}(\text{SG}) = \bigcup \{ \phi_\ell^{\text{SA}}((\mathbf{S}, \mathbf{H}, \mathbf{is})) \mid (\mathbf{S}, \mathbf{H}, \mathbf{is}) \in \text{SG} \}$$

where

$$\phi_\ell^{\text{SA}} : \text{SG} \rightarrow \mathcal{P}(\text{SG})$$

specifies how a *single* shape graph (in $\text{Shape}_\circ(\ell)$) may be transformed into a *set* of shape graphs (in $\text{Shape}_\bullet(\ell)$) by the elementary block.

217

© Nielsen^2, Hankin

Transfer function for $[b]^\ell$ and $[\text{skip}]^\ell$

We are only interested in the shape of the heap – and it is not changed by these elementary blocks:

$$\phi_\ell^{\text{SA}}((\mathbf{S}, \mathbf{H}, \mathbf{is})) = \{(\mathbf{S}, \mathbf{H}, \mathbf{is})\}$$

218

© Nielsen^2, Hankin

Transfer function for $[x := a]^\ell$

— where a is of the form n , $a_1 \text{ op}_a a_2$ or nil

$$\phi_\ell^{\text{SA}}((\mathbf{S}, \mathbf{H}, \mathbf{is})) = \{\text{kill}_x((\mathbf{S}, \mathbf{H}, \mathbf{is}))\}$$

where $\text{kill}_x((\mathbf{S}, \mathbf{H}, \mathbf{is})) = (\mathbf{S}', \mathbf{H}', \mathbf{is}')$ is

$$\mathbf{S}' = \{(z, k_x(n_Z)) \mid (z, n_Z) \in \mathbf{S} \wedge z \neq x\}$$

$$\mathbf{H}' = \{(k_x(n_V), \text{sel}, k_x(n_W)) \mid (n_V, \text{sel}, n_W) \in \mathbf{H}\}$$

$$\mathbf{is}' = \{k_x(n_X) \mid n_X \in \mathbf{is}\}$$

and

$$k_x(n_Z) = n_{Z \setminus \{x\}}$$

Idea: all abstract locations are renamed to not having x in their name set

219

© Nielsen^2, Hankin

Transfer function for $[x := a]^\ell$

— where a is of the form n , $a_1 \text{ op}_a a_2$ or nil

$$\phi_\ell^{\text{SA}}((\mathbf{S}, \mathbf{H}, \mathbf{is})) = \{\text{kill}_x((\mathbf{S}, \mathbf{H}, \mathbf{is}))\}$$

where $\text{kill}_x((\mathbf{S}, \mathbf{H}, \mathbf{is})) = (\mathbf{S}', \mathbf{H}', \mathbf{is}')$ is

$$\mathbf{S}' = \{(z, k_x(n_Z)) \mid (z, n_Z) \in \mathbf{S} \wedge z \neq x\}$$

$$\mathbf{H}' = \{(k_x(n_V), \text{sel}, k_x(n_W)) \mid (n_V, \text{sel}, n_W) \in \mathbf{H}\}$$

$$\mathbf{is}' = \{k_x(n_X) \mid n_X \in \mathbf{is}\}$$

and

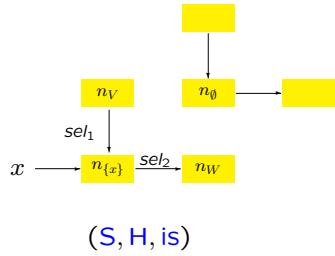
$$k_x(n_Z) = n_{Z \setminus \{x\}}$$

Idea: all abstract locations are renamed to not having x in their name set

220

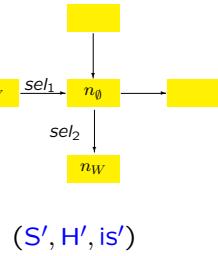
© Nielsen^2, Hankin

The effect of $[x := \text{nil}]^\ell$



221

© Nielsen^2, Hankin



Transfer function for $[x := y]^\ell$ when $x \neq y$

$$\phi_\ell^{\text{SA}}((S, H, is)) = \{(S'', H'', is'')\}$$

where $(S', H', is') = \text{kill}_x((S, H, is))$ and

$$S'' = \{(z, g_x^y(n_Z)) \mid (z, n_Z) \in S'\} \\ \cup \{(x, g_x^y(n_Y)) \mid (y, n_Y) \in S' \wedge y' = y\}$$

$$H'' = \{(g_x^y(n_V), sel, g_x^y(n_W)) \mid (n_V, sel, n_W) \in H'\}$$

$$is'' = \{g_x^y(n_Z) \mid n_Z \in is'\}$$

and

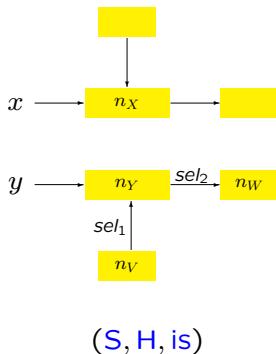
$$g_x^y(n_Z) = \begin{cases} n_{Z \cup \{x\}} & \text{if } y \in Z \\ n_Z & \text{otherwise} \end{cases}$$

Idea: all abstract locations are renamed to also have x in their name set if they already have y

222

© Nielsen^2, Hankin

The effect of $[x := y]^\ell$ when $x \neq y$



223

© Nielsen^2, Hankin

 (S'', H'', is'')

Transfer function for $[x := y.sel]^\ell$ when $x \neq y$

Remove the old binding for x :

strong nullification

$$(S', H', is') = \text{kill}_x((S, H, is))$$

Establish the new binding for x :

1. There is no abstract location n_Y such that $(y, n_Y) \in S'$ – or there is an abstract location n_Y such that $(y, n_Y) \in S'$ but no n_Z such that $(n_Y, sel, n_Z) \in H'$
2. There is an abstract location n_Y such that $(y, n_Y) \in S'$ and there is an abstract location $n_U \neq n_\emptyset$ such that $(n_Y, sel, n_U) \in H'$
3. There is an abstract location n_Y such that $(y, n_Y) \in S'$ and $(n_Y, sel, n_\emptyset) \in H'$

224

© Nielsen^2, Hankin

Case 1 for $[x:=y.sel]^\ell$

Assume there is no abstract location n_Y such that $(y, n_Y) \in S'$

$$\phi_\ell^{\text{SA}}((S, H, \text{is})) = \{(S', H', \text{is}')\}$$

OBS: dereference of a nil-pointer

Assume there is an abstract location n_Y such that $(y, n_Y) \in S'$ but there is no abstract location n such that $(n_Y, sel, n) \in H'$

$$\phi_\ell^{\text{SA}}((S, H, \text{is})) = \{(S', H', \text{is}')\}$$

OBS: dereference of a non-existing sel-field

225

© Nielsen^2, Hankin

Case 2 for $[x:=y.sel]^\ell$

Assume there is an abstract location n_Y such that $(y, n_Y) \in S'$ and there is an abstract location $n_U \neq n_\emptyset$ such that $(n_Y, sel, n_U) \in H'$.

The abstract location n_U will be renamed to include the variable x using the function:

$$h_x^U(n_Z) = \begin{cases} n_{U \cup \{x\}} & \text{if } Z = U \\ n_Z & \text{otherwise} \end{cases}$$

We take

$$\phi_\ell^{\text{SA}}((S, H, \text{is})) = \{(S'', H'', \text{is}'')\}$$

where $(S', H', \text{is}') = \text{kill}_x((S, H, \text{is}))$ and

$$S'' = \{(z, h_x^U(n_Z)) \mid (z, n_Z) \in S'\} \cup \{(x, h_x^U(n_U))\}$$

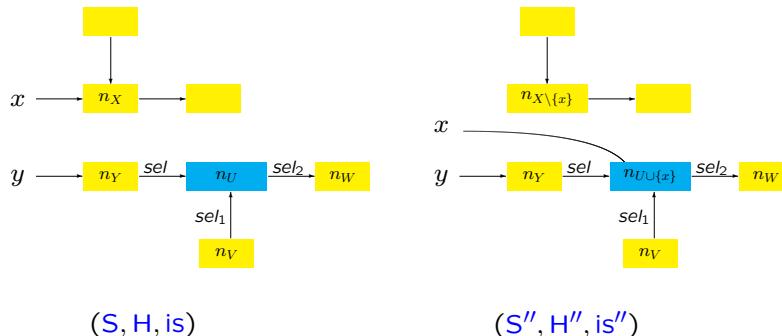
$$H'' = \{(h_x^U(n_V), sel', h_x^U(n_W)) \mid (n_V, sel', n_W) \in H'\}$$

$$\text{is}'' = \{h_x^U(n_Z) \mid n_Z \in \text{is}'\}$$

226

© Nielsen^2, Hankin

The effect of $[x:=y.sel]^\ell$ in Case 2



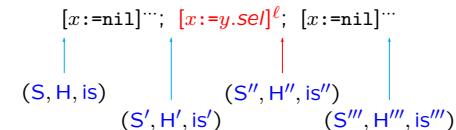
227

© Nielsen^2, Hankin

Case 3 for $[x:=y.sel]^\ell$ (1)

Assume that there is an abstract location n_Y such that $(y, n_Y) \in S'$ and furthermore $(n_Y, sel, n_\emptyset) \in H'$.

We have to *materialise* a new abstract location $n_{\{x\}}$ from n_\emptyset .



Idea:

$$(S', H', \text{is}') = (S''', H''', \text{is}''') = \text{kill}_x((S'', H'', \text{is}''))$$

228

© Nielsen^2, Hankin

Case 3 for $[x := y.\text{sel}]^\ell$ (2)

Transfer function:

$$\begin{aligned}\phi_\ell^{\text{SA}}((S, H, \text{is})) &= \{(S'', H'', \text{is}'') \mid (S'', H'', \text{is}'') \text{ is compatible} \wedge \\ &\quad \text{kill}_x((S'', H'', \text{is}'')) = (S', H', \text{is}') \wedge \\ &\quad (x, n_{\{x\}}) \in S'' \wedge (n_Y, \text{sel}, n_{\{x\}}) \in H''\}\end{aligned}$$

where $(S', H', \text{is}') = \text{kill}_x((S, H, \text{is}))$.

229

© Nielsen^2, Hankin

Case 3 for $[x := y.\text{sel}]^\ell$ (2)

Transfer function:

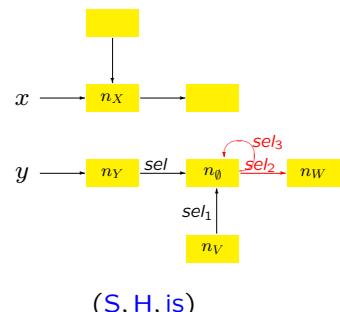
$$\begin{aligned}\phi_\ell^{\text{SA}}((S, H, \text{is})) &= \{(S'', H'', \text{is}'') \mid (S'', H'', \text{is}'') \text{ is compatible} \wedge \\ &\quad \text{kill}_x((S'', H'', \text{is}'')) = (S', H', \text{is}') \wedge \\ &\quad (x, n_{\{x\}}) \in S'' \wedge (n_Y, \text{sel}, n_{\{x\}}) \in H''\}\end{aligned}$$

where $(S', H', \text{is}') = \text{kill}_x((S, H, \text{is}))$.

230

© Nielsen^2, Hankin

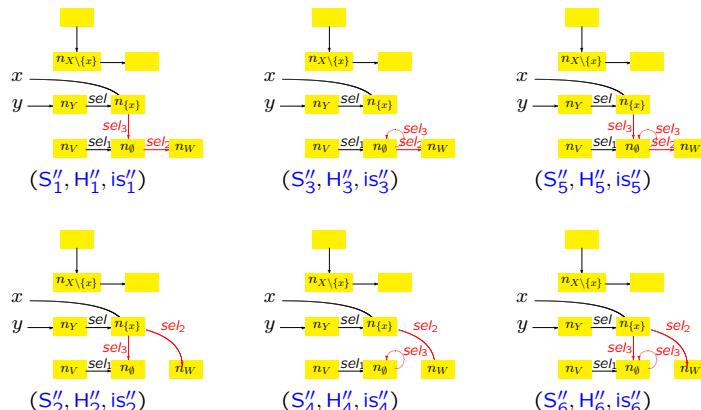
The effect of $[x := y.\text{sel}]^\ell$ in Case 3 (1)



231

© Nielsen^2, Hankin

The effect of $[x := y.\text{sel}]^\ell$ in Case 3 (2)



232

© Nielsen^2, Hankin

QUESTIONS?

Transfer function for $[x.\text{sel}:=a]^\ell$

— where a is of the form n , $a_1 \ op_a\ a_2$ or nil .

If there is no n_X such that $(x, n_X) \in S$ then f_ℓ^{SA} is the identity.

If there is n_X such that $(x, n_X) \in S$ but that there is no n_U such that $(n_X, sel, n_U) \in H$ then f_ℓ^{SA} is the identity.

If there are abstract locations n_X and n_U such that $(x, n_X) \in \mathbf{S}$ and $(n_X, sel, n_U) \in \mathbf{H}$ then

$$\phi_\ell^{\text{SA}}((\mathbf{S}, \mathbf{H}, \text{is})) = \{\textcolor{red}{kill}_{x.\text{sel}}((\mathbf{S}, \mathbf{H}, \text{is}))\}$$

where $\text{kill}_{x.\text{sel}}((S, H, \text{is})) = (S', H', \text{is}')$ is given by

$$s' = s$$

$$\texttt{H}' = \{(n_V, sel', n_W) \mid (n_V, sel', n_W) \in \texttt{H} \wedge \neg(X = V \wedge sel = sel')\}$$

$$\text{is}' = \begin{cases} \text{is} \setminus \{n_U\} & \text{if } n_U \in \text{is} \wedge \#\text{into}(n_U, H') \leq 1 \wedge \neg \exists(n_\emptyset, \text{sel}', n_U) \in H' \\ \text{is} & \text{otherwise} \end{cases}$$

233

© Nielsen^2, Hankin

Transfer function for $\text{[malloc } x\text{]}^\ell$

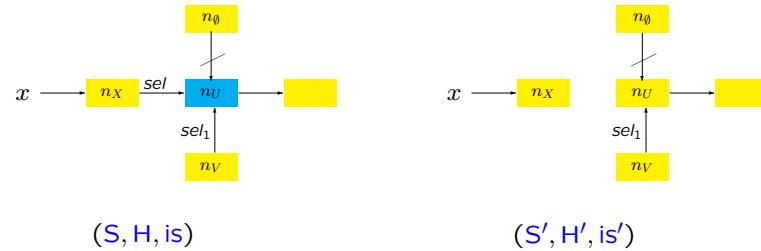
$$\phi_\ell^{\text{SA}}((\mathbf{S}, \mathbf{H}, \mathbf{is})) = \{(\mathbf{S}' \cup \{(x, n_{\{x\}})\}, \mathbf{H}', \mathbf{is}')\}$$

where $(S', H', \text{is}') = \text{kill}_x(S, H, \text{is})$.

235

© Nielsen^2, Hankin

The effect of $[x.\text{sel} := \text{nil}]^\ell$ when $\#\text{into}(n_U, \mathcal{H}') \leq 1$



234

© Nielsen^2, Hankin

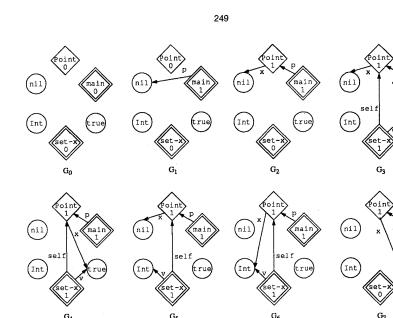


Fig. 6: Abstract Object Graphs. Creation counts are displayed under node names.

Compile-Time Analysis of Object-Oriented Programs

Jan Vitek, R. Nigel Horspool and James S. Uhl
Department of Computer Science
University of Victoria *

Abstract. Generation of efficient code for object-oriented programs requires knowledge of object lifetimes and method bindings. For object-oriented languages that have automatic storage management and dynamic look-up of methods, the compiler must maintain such knowledge by performing static analysis of the program. We present a method for determining the lifetime and potential classes of each object in an object-oriented program as well as a safe approximation of their lifetimes. These results are obtained using abstract domains that approximate memory configurations and interprocedure call patterns of the program. We present several alternatives for these abstract domains that permit a trade-off between accuracy and complexity of the overall analysis.

1 Introduction

The object-oriented approach to programming has become an accepted programming paradigm, joining other paradigms such as imperative, functional and relational programming. This new paradigm is normally associated with the concepts of *class*, *method* and *inheritance*. Different object-oriented languages implement these concepts with varying degrees of dynamic behaviour.

instances of a class are created dynamically which, when no longer referenced, are collected by the garbage collector. This is the mechanism used in Smalltalk and in most other dynamic implementations. In general this means that a run-time search through the inheritance hierarchy finds the appropriate method as required for every message send. Finally, SMALLTALK is dynamically typed; there are no type declarations and methods are type-checked at run-time.

In an ideal world, an object-oriented language would be as dynamic as SMALL-TALK and as efficient as C++. The language would provide dynamic features, but

a compiler would analyze the code and determine whether or not these features are used. The user would then only pay a run-time penalty when, and where, the language is used in a truly dynamic manner.

* P.O. Box 3055, Victoria, BC, Canada V8W 3P6. {jvitek, nigeh, juhl}@csr.uvic.ca

Scaling Program Analysis

CS7575

- Lecture 4 – Constraint Based Analysis - PoPA Chapter 3

237

© Nielsen^2, Hankin

```
let f = fn x => x 1;
      g = fn y => y+2;
      h = fn z => z+3
in (f g) + (f h)
```

The aim of Control Flow Analysis:

For each function application, which functions may be applied?

Control Flow Analysis computes the interprocedural flow relation used when formulating interprocedural Data Flow Analysis.

239

© Nielsen^2, Hankin

The Dynamic Dispatch Problem

```
: [call p(p1,1,v)] $_{\ell_r}^{\ell_c^1}$ 
:
[call p(p2,2,v)] $_{\ell_r}^{\ell_c^2}$ 
:
proc p(procval q, val x, res y) is $^{\ell_n}$ 
:
[call q (x,y)] $_{\ell_r}^{\ell_c^p}$ 
:
end $^{\ell_x}$ 
```

which procedure is called?

These problems arise for:

- imperative languages with procedures as parameters
- object oriented languages
- functional languages

238

© Nielsen^2, Hankin

Syntax of the Fun Language

Syntactic categories:

| | | |
|------------|-------|-----------------------------------|
| $e \in$ | Exp | expressions (or labelled terms) |
| $t \in$ | Term | terms (or unlabelled expressions) |
| $f, x \in$ | Var | variables |
| $c \in$ | Const | constants |
| $op \in$ | Op | binary operators |
| $\ell \in$ | Lab | labels |

Syntax:

```
 $e ::= t^\ell$ 
 $t ::= c \mid x \mid \text{fn } x \Rightarrow e_0 \mid \text{fun } f \ x \Rightarrow e_0 \mid e_1 \ e_2$ 
       $\mid \text{if } e_0 \text{ then } e_1 \text{ else } e_2 \mid \text{let } x = e_1 \text{ in } e_2 \mid e_1 \ op \ e_2$ 
```

(Labels correspond to program points or nodes in the parse tree.) © Nielsen^2, Hankin

$((\text{fn } x \Rightarrow x^1)^2 (\text{fn } y \Rightarrow y^3)^4)^5$

```
(let f = (fn x => (x1 x2)3)4;
in (let g = (fn y => y5)6;
     in (let h = (fn z => z7)8
          in ((f9 g10)11 + (f12 h13)14)15)16)17)18
```

```
(let g = (fun f x => (f1 (fn y => y2)3)4)5
in (g6 (fn z => z7)8)9)10
```

241

© Nielsen^2, Hankin

Abstract 0-CFA Analysis

- Abstract domains
- Specification of the analysis
- Well-definedness of the analysis

242

© Nielsen^2, Hankin

Towards defining the Abstract Domains

The *result* of a 0-CFA analysis is a pair $(\hat{\mathcal{C}}, \hat{\rho})$:

- $\hat{\mathcal{C}}$ is the *abstract cache* associating abstract values with each labelled program point
- $\hat{\rho}$ is the *abstract environment* associating abstract values with each variable

243

© Nielsen^2, Hankin

$((\text{fn } x \Rightarrow x^1)^2 (\text{fn } y \Rightarrow y^3)^4)^5$

Three guesses of a 0-CFA analysis result:

| | $(\hat{\mathcal{C}}_e, \hat{\rho}_e)$ | $(\hat{\mathcal{C}}'_e, \hat{\rho}'_e)$ | $(\hat{\mathcal{C}}''_e, \hat{\rho}''_e)$ |
|---|---------------------------------------|---|--|
| 1 | {fn y => y ³ } | {fn y => y ³ } | {fn x => x ¹ , fn y => y ³ } |
| 2 | {fn x => x ¹ } | {fn x => x ¹ } | {fn x => x ¹ , fn y => y ³ } |
| 3 | \emptyset | \emptyset | {fn x => x ¹ , fn y => y ³ } |
| 4 | {fn y => y ³ } | {fn y => y ³ } | {fn x => x ¹ , fn y => y ³ } |
| 5 | {fn y => y ³ } | {fn y => y ³ } | {fn x => x ¹ , fn y => y ³ } |
| x | {fn y => y ³ } | \emptyset | {fn x => x ¹ , fn y => y ³ } |
| y | \emptyset | \emptyset | {fn x => x ¹ , fn y => y ³ } |

244

© Nielsen^2, Hankin

```
(let g = (fun f x => (f1 (fn y => y2)3)4)5
  in (g6 (fn z => z7)8)9)10
```

Abbreviations:

$$\begin{aligned} f &= \text{fun } f \ x \ => (f^1 \ (\text{fn } y \ => y^2)^3)^4 \\ \text{id}_y &= \text{fn } y \ => y^2 \\ \text{id}_z &= \text{fn } z \ => z^7 \end{aligned}$$

One guess of a 0-CFA analysis result:

$$\begin{array}{lll} \hat{C}_{lp}(1) = \{f\} & \hat{C}_{lp}(6) = \{f\} & \hat{\rho}_{lp}(f) = \{f\} \\ \hat{C}_{lp}(2) = \emptyset & \hat{C}_{lp}(7) = \emptyset & \hat{\rho}_{lp}(g) = \{f\} \\ \hat{C}_{lp}(3) = \{\text{id}_y\} & \hat{C}_{lp}(8) = \{\text{id}_z\} & \hat{\rho}_{lp}(x) = \{\text{id}_y, \text{id}_z\} \\ \hat{C}_{lp}(4) = \emptyset & \hat{C}_{lp}(9) = \emptyset & \hat{\rho}_{lp}(y) = \emptyset \\ \hat{C}_{lp}(5) = \{f\} & \hat{C}_{lp}(10) = \emptyset & \hat{\rho}_{lp}(z) = \emptyset \end{array}$$

© Nielsen^2, Hankin

Abstract Domains

Formally:

$$\begin{aligned} \hat{v} \in \widehat{\text{Val}} &= \mathcal{P}(\text{Term}) \quad \text{abstract values} \\ \hat{\rho} \in \widehat{\text{Env}} &= \text{Var} \rightarrow \widehat{\text{Val}} \quad \text{abstract environments} \\ \hat{C} \in \widehat{\text{Cache}} &= \text{Lab} \rightarrow \widehat{\text{Val}} \quad \text{abstract caches} \end{aligned}$$

An abstract value \hat{v} is a set of terms of the forms

- $\text{fn } x \ => e_0$
- $\text{fun } f \ x \ => e_0$

246

© Nielsen^2, Hankin

Specification of the 0-CFA

When is a proposed guess $(\hat{C}, \hat{\rho})$ of an analysis results an *acceptable 0-CFA analysis* for the program?

Different approaches:

- **abstract specification**
- syntax-directed and constraint-based specifications
- algorithms for computing the *best* result

247

© Nielsen^2, Hankin

Specification of the Abstract 0-CFA

$(\hat{C}, \hat{\rho}) \models e$ means that $(\hat{C}, \hat{\rho})$ is an *acceptable Control Flow Analysis* of the expression e

The relation \models has functionality:

$$\models : (\widehat{\text{Cache}} \times \widehat{\text{Env}} \times \text{Exp}) \rightarrow \{\text{true}, \text{false}\}$$

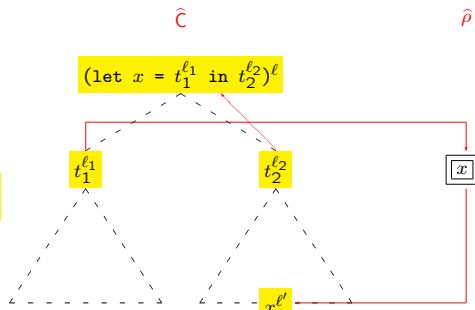
248

© Nielsen^2, Hankin

$(\hat{C}, \hat{\rho}) \models c^\ell$ always

$(\hat{C}, \hat{\rho}) \models x^\ell \text{ iff } \hat{\rho}(x) \subseteq \hat{C}(\ell)$

$(\hat{C}, \hat{\rho}) \models (\text{let } x = t_1^{\ell_1} \text{ in } t_2^{\ell_2})^\ell$
 $\text{iff } (\hat{C}, \hat{\rho}) \models t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models t_2^{\ell_2} \wedge$
 $\hat{C}(\ell_1) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_2) \subseteq \hat{C}(\ell)$



249

© Nielsen^2, Hankin

$(\hat{C}, \hat{\rho}) \models (\text{if } t_0^{\ell_0} \text{ then } t_1^{\ell_1} \text{ else } t_2^{\ell_2})^\ell$

$\text{iff } (\hat{C}, \hat{\rho}) \models t_0^{\ell_0} \wedge$

$(\hat{C}, \hat{\rho}) \models t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models t_2^{\ell_2} \wedge$

$\hat{C}(\ell_1) \subseteq \hat{C}(\ell) \wedge \hat{C}(\ell_2) \subseteq \hat{C}(\ell)$

$(\hat{C}, \hat{\rho}) \models (t_1^{\ell_1} op t_2^{\ell_2})^\ell$

$\text{iff } (\hat{C}, \hat{\rho}) \models t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models t_2^{\ell_2}$

250

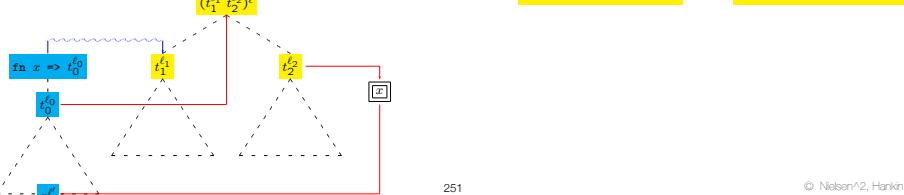
© Nielsen^2, Hankin

$(\hat{C}, \hat{\rho}) \models (\text{fn } x \Rightarrow t_0^{\ell_0})^\ell \text{ iff } \{\text{fn } x \Rightarrow t_0^{\ell_0}\} \subseteq \hat{C}(\ell)$

$(\hat{C}, \hat{\rho}) \models (t_1^{\ell_1} t_2^{\ell_2})^\ell$
 $\text{iff } (\hat{C}, \hat{\rho}) \models t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models t_2^{\ell_2} \wedge$

$(\forall (\text{fn } x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1) : (\hat{C}, \hat{\rho}) \models t_0^{\ell_0} \wedge$

$\hat{C}(\ell_2) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell))$



251

© Nielsen^2, Hankin

$(\hat{C}, \hat{\rho}) \models (\text{fun } f x \Rightarrow e_0)^\ell \text{ iff } \{\text{fun } f x \Rightarrow e_0\} \subseteq \hat{C}(\ell)$

$(\hat{C}, \hat{\rho}) \models (t_1^{\ell_1} t_2^{\ell_2})^\ell$

$\text{iff } (\hat{C}, \hat{\rho}) \models t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models t_2^{\ell_2} \wedge$

$(\forall (\text{fn } x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1) : (\hat{C}, \hat{\rho}) \models t_0^{\ell_0} \wedge$

$\hat{C}(\ell_2) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell)) \wedge$

$(\forall (\text{fun } f x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1) : (\hat{C}, \hat{\rho}) \models t_0^{\ell_0} \wedge$

$\hat{C}(\ell_2) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell) \wedge$

$\{\text{fun } f x \Rightarrow t_0^{\ell_0}\} \subseteq \hat{\rho}(f)\)$

252

© Nielsen^2, Hankin

Two guesses for $((\text{fn } x \Rightarrow x^1)^2 (\text{fn } y \Rightarrow y^3)^4)^5$

| | $(\hat{C}_e, \hat{\rho}_e)$ | $(\hat{C}'_e, \hat{\rho}'_e)$ |
|---|------------------------------------|------------------------------------|
| 1 | $\{\text{fn } y \Rightarrow y^3\}$ | $\{\text{fn } y \Rightarrow y^3\}$ |
| 2 | $\{\text{fn } x \Rightarrow x^1\}$ | $\{\text{fn } x \Rightarrow x^1\}$ |
| 3 | \emptyset | \emptyset |
| 4 | $\{\text{fn } y \Rightarrow y^3\}$ | $\{\text{fn } y \Rightarrow y^3\}$ |
| 5 | $\{\text{fn } y \Rightarrow y^3\}$ | $\{\text{fn } y \Rightarrow y^3\}$ |
| x | $\{\text{fn } y \Rightarrow y^3\}$ | \emptyset |
| y | \emptyset | \emptyset |

Checking the guesses:

$$(\hat{C}_e, \hat{\rho}_e) \models ((\text{fn } x \Rightarrow x^1)^2 (\text{fn } y \Rightarrow y^3)^4)^5$$

$$(\hat{C}'_e, \hat{\rho}'_e) \not\models ((\text{fn } x \Rightarrow x^1)^2 (\text{fn } y \Rightarrow y^3)^4)^5$$

© Nielsen^2, Hankin

Well-definedness of the Abstract 0-CFA

Difficulty: The clause for function application is *not* of a form that allows us to define $(\hat{C}, \hat{\rho}) \models e$ by Structural Induction in the expression e

$$\begin{aligned} & (\hat{C}, \hat{\rho}) \models (t_1^{\ell_1} t_2^{\ell_2})^\ell \\ \text{iff } & (\hat{C}, \hat{\rho}) \models t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models t_2^{\ell_2} \wedge \\ & (\forall (\text{fn } x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1) : (\hat{C}, \hat{\rho}) \models t_0^{\ell_0} \wedge \\ & \hat{C}(\ell_2) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell)) \end{aligned}$$

Solution: The relation \models is defined by **coinduction**, that is, as the **greatest fixed point** of a functional.

254

© Nielsen^2, Hankin

Theoretical Properties:

- structural operational semantics
- semantic correctness
- the existence of least solutions

255

© Nielsen^2, Hankin

Choice of Semantics

- operational or denotational semantics?
 - an operational semantics more easily models intensional properties
- small-step or big-step operational semantics?
 - a small-step semantics allows us to reason about looping programs
- operational semantics based on environments or substitutions?
 - an environment based semantics preserves the identity of functions

256

© Nielsen^2, Hankin

Configurations and Transitions

Semantic categories:

$$v \in \text{Val} \quad \text{values}$$

$$\rho \in \text{Env} \quad \text{environments}$$

defined by:

$$v ::= c \mid \text{close } t \text{ in } \rho \quad \text{closures}$$

$$\rho ::= [] \mid \rho[x \mapsto v]$$

Transitions have the form

$$\rho \vdash e_1 \rightarrow e_2$$

meaning that *one step* of computation of the expression e_1 in the environment ρ will transform it into e_2 .

257

© Nielsen^2, Hankin

Transitions

$$\rho \vdash x^\ell \rightarrow v^\ell \text{ if } x \in \text{dom}(\rho) \text{ and } v = \rho(x)$$

$$\rho \vdash (\text{fn } x \Rightarrow e_0)^\ell \rightarrow (\text{close } (\text{fn } x \Rightarrow e_0) \text{ in } \rho_0)^\ell \\ \text{where } \rho_0 = \rho \mid FV(\text{fn } x \Rightarrow e_0)$$

$$\rho \vdash (\text{fun } f x \Rightarrow e_0)^\ell \rightarrow (\text{close } (\text{fun } f x \Rightarrow e_0) \text{ in } \rho_0)^\ell \\ \text{where } \rho_0 = \rho \mid FV(\text{fun } f x \Rightarrow e_0)$$

static scope!

258

© Nielsen^2, Hankin

Intermediate Expressions and Terms

$$ie \in \text{IExp} \quad \text{intermediate expressions}$$

$$it \in \text{ITerm} \quad \text{intermediate terms}$$

extending the syntax:

$$\begin{aligned} ie &::= it^\ell \\ it &::= c \mid x \mid \text{fn } x \Rightarrow e_0 \mid \text{fun } f x \Rightarrow e_0 \mid ie_1 ie_2 \\ &\mid \text{if } ie_0 \text{ then } e_1 \text{ else } e_2 \mid \text{let } x = ie_1 \text{ in } e_2 \mid ie_1 \text{ op } ie_2 \\ &\mid \text{close } t \text{ in } \rho \mid \text{bind } \rho \text{ in } ie \end{aligned}$$

The correct form of transitions

$$\rho \vdash ie_1 \rightarrow ie_2$$

© Nielsen^2, Hankin

Transitions

$$\frac{\rho \vdash ie_1 \rightarrow ie'_1}{\rho \vdash (ie_1 ie_2)^\ell \rightarrow (ie'_1 ie_2)^\ell}$$

$$\frac{\rho \vdash ie_2 \rightarrow ie'_2}{\rho \vdash (v_1^{\ell_1} ie_2)^\ell \rightarrow (v_1^{\ell_1} ie'_2)^\ell}$$

$$\begin{aligned} \rho \vdash ((\text{close } (\text{fn } x \Rightarrow e_1) \text{ in } \rho_1)^{\ell_1} v_2^{\ell_2})^\ell &\rightarrow (\text{bind } \rho_1[x \mapsto v_2] \text{ in } e_1)^\ell \\ \rho \vdash ((\text{close } (\text{fun } f x \Rightarrow e_1) \text{ in } \rho_1)^{\ell_1} v_2^{\ell_2})^\ell &\rightarrow (\text{bind } \rho_2[x \mapsto v_2] \text{ in } e_1)^\ell \\ \text{where } \rho_2 = \rho_1[f \mapsto \text{close } (\text{fun } f x \Rightarrow e_1) \text{ in } \rho_1] \end{aligned}$$

$$\frac{\rho_1 \vdash ie_1 \rightarrow ie'_1}{\rho \vdash (\text{bind } \rho_1 \text{ in } ie_1)^\ell \rightarrow (\text{bind } \rho_1 \text{ in } ie'_1)^\ell}$$

$$\rho \vdash (\text{bind } \rho_1 \text{ in } v_1^{\ell_1})^\ell \rightarrow v_1^\ell \quad \text{the outermost label remains the same}$$

```

[ ] ⊢ ((fn x => x1)2 (fn y => y3)4)5
→ ((close (fn x => x1) in [ ])2 (fn y => y3)4)5
→ ((close (fn x => x1) in [ ])2 (close (fn y => y3) in [ ])4)5
→ (bind [x ↦ (close (fn y => y3) in [ ])] in x1)5
→ (bind [x ↦ (close (fn y => y3) in [ ])] in
      (close (fn y => y3) in [ ])1)5
→ (close (fn y => y3) in [ ])5

```

261

© Nielsen^2, Hankin

Transitions

$$\begin{array}{c}
 \frac{\rho \vdash ie_0 \rightarrow ie'_0}{\rho \vdash (\text{if } ie_0 \text{ then } e_1 \text{ else } e_2)^\ell \rightarrow (\text{if } ie'_0 \text{ then } e_1 \text{ else } e_2)^\ell} \\
 \frac{\rho \vdash (\text{if true}^{\ell_0} \text{ then } t_1^{\ell_1} \text{ else } t_2^{\ell_2})^\ell \rightarrow t_1^\ell}{\rho \vdash (\text{if false}^{\ell_0} \text{ then } t_1^{\ell_1} \text{ else } t_2^{\ell_2})^\ell \rightarrow t_2^\ell} \\
 \\
 \frac{\rho \vdash ie_1 \rightarrow ie'_1}{\rho \vdash (\text{let } x = ie_1 \text{ in } e_2)^\ell \rightarrow (\text{let } x = ie'_1 \text{ in } e_2)^\ell} \\
 \frac{\rho \vdash (\text{let } x = v^{\ell_1} \text{ in } e_2)^\ell \rightarrow (\text{bind } [x \mapsto v] \text{ in } e_2)^\ell}{\rho \vdash ie_1 \rightarrow ie'_1} \\
 \\
 \frac{\rho \vdash ie_1 \rightarrow ie'_1}{\rho \vdash (ie_1 \text{ op } ie_2)^\ell \rightarrow (ie'_1 \text{ op } ie_2)^\ell} \quad \frac{\rho \vdash ie_2 \rightarrow ie'_2}{\rho \vdash (v_1^{\ell_1} \text{ op } ie_2)^\ell \rightarrow (v_1^{\ell_1} \text{ op } ie'_2)^\ell} \\
 \frac{\rho \vdash (v_1^{\ell_1} \text{ op } v_2^{\ell_2})^\ell \rightarrow v^\ell \quad \text{if } v = v_1 \text{ op } v_2}{\rho \vdash ie_1 \rightarrow ie'_1} \quad \frac{\rho \vdash ie_2 \rightarrow ie'_2}{\rho \vdash (v_1^{\ell_1} \text{ op } v_2^{\ell_2})^\ell \rightarrow v^\ell \quad \text{if } v = v_1 \text{ op } v_2}
 \end{array}$$

© Nielsen^2, Hankin

```

[ ] ⊢ (let g = (fun f x => (f1 (fn y => y2)3)4)5
       in (g6 (fn z => z7)8)9)10
→ (let g = f5 in (g6 (fn z => z7)8)9)10
→ (bind [g ↦ f] in (g6 (fn z => z7)8)9)10
→ (bind [g ↦ f] in (f6 (fn z => z7)8)9)10
→ (bind [g ↦ f] in (f6 idz8)9)10
→ (bind [g ↦ f] in (bind [f ↦ f][x ↦ idz] in (f1 (fn y => y2)3)4)9)10
→* (bind [g ↦ f] in (bind [f ↦ f][x ↦ idz] in
      (bind [f ↦ f][x ↦ idy] in (f1 (fn y => y2)3)4)4)9)10
→* ...

```

Abbreviations:

$$\begin{aligned}
 f &= \text{close (fun f x => (f¹ (fn y => y²)³)⁴) in []} \\
 \text{id}_y &= \text{close (fn y => y²) in []} \\
 \text{id}_z &= \text{close (fn z => z⁷) in []}
 \end{aligned}$$

263

© Nielsen^2, Hankin

Semantic Correctness

A *subject reduction result*: an acceptable result of the analysis remains acceptable under evaluation

Analysis of intermediate expressions

$$\begin{array}{ll}
 (\hat{C}, \hat{\rho}) \models (\text{bind } \rho \text{ in } it_0^{\ell_0})^\ell & \\
 \text{iff} & (\hat{C}, \hat{\rho}) \models it_0^{\ell_0} \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell) \wedge \rho \mathcal{R} \hat{\rho} \\
 \\
 (\hat{C}, \hat{\rho}) \models (\text{close } t_0 \text{ in } \rho)^\ell & \\
 \text{iff} & \{t_0\} \subseteq \hat{C}(\ell) \wedge \rho \mathcal{R} \hat{\rho}
 \end{array}$$

264

© Nielsen^2, Hankin

Correctness Relation

The global abstract environment, $\hat{\rho}$ models all the local environments of the semantics

Correctness relation

$$\mathcal{R} : (\text{Env} \times \widehat{\text{Env}}) \rightarrow \{\text{true}, \text{false}\}$$

We demand that $\rho \mathcal{R} \hat{\rho}$ for all local environments, ρ , occurring in the intermediate expressions

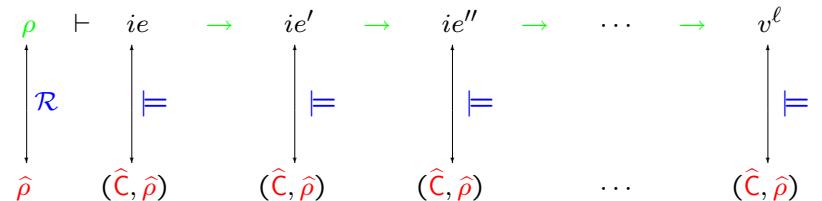
Define

$$\rho \mathcal{R} \hat{\rho} \quad \text{iff} \quad \forall x \in \text{dom}(\rho) \subseteq \text{dom}(\hat{\rho}) \quad \forall t_x \quad \forall \rho_x : \\ (\rho(x) = \text{close } t_x \text{ in } \rho_x) \Rightarrow (t_x \in \hat{\rho}(x) \wedge \rho_x \mathcal{R} \hat{\rho})$$

(Well-defined by induction in the size of ρ)

© Nielsen^2, Hankin

Correctness Result



266

© Nielsen^2, Hankin

Formal details of Correctness Result

Theorem:

If $\rho \mathcal{R} \hat{\rho}$ and $\rho \vdash ie \rightarrow ie'$ then $(\hat{C}, \hat{\rho}) \models ie$ implies $(\hat{C}, \hat{\rho}) \models ie'$.

Intuitively:

If there is a possible evaluation of the program such that the function at a call point evaluates to some abstraction, then this abstraction has to be in the set of possible abstractions computed by the analysis.

Observe: the theorem expresses that all acceptable analysis results remain acceptable under evaluation!

Thus we do not rely on the existence of a least or “best” solution.

© Nielsen^2, Hankin

Semantics:

$$[] \vdash ((\text{fn } x \Rightarrow x^1)^2 (\text{fn } y \Rightarrow y^3)^4)^5 \xrightarrow{*} (\text{close } (\text{fn } y \Rightarrow y^3) \text{ in } [])^5$$

| | $(\hat{C}_e, \hat{\rho}_e)$ |
|---|------------------------------------|
| 1 | $\{\text{fn } y \Rightarrow y^3\}$ |
| 2 | $\{\text{fn } x \Rightarrow x^1\}$ |
| 3 | \emptyset |
| 4 | $\{\text{fn } y \Rightarrow y^3\}$ |
| 5 | $\{\text{fn } y \Rightarrow y^3\}$ |
| x | $\{\text{fn } y \Rightarrow y^3\}$ |
| y | \emptyset |

Analysis:

$$(\hat{C}_e, \hat{\rho}_e) \models ((\text{fn } x \Rightarrow x^1)^2 (\text{fn } y \Rightarrow y^3)^4)^5$$

Correctness relation:

$$[] \mathcal{R} \hat{\rho}_e$$

Correctness theorem: $(\hat{C}_e, \hat{\rho}_e) \models (\text{close } (\text{fn } y \Rightarrow y^3) \text{ in } [])^5$

268

© Nielsen^2, Hankin



The History of T

Olin Shivers

269

This brings us to the summer of 1984. The mission was to build the world's most highly-optimizing Scheme compiler. We wanted to compete with C and Fortran. ... I had passed the previous semester at CMU becoming an expert on data-flow analysis, a topic on which I completely grooved.

All hot compilers do DFA. It is necessary for all the really cool optimisations, like loop-invariant hoisting, global register allocation, global common subexpression elimination, copy propagation, induction-variable elimination. I knew that no Scheme or Lisp compiler had ever provided these hot optimisations. I was burning to make it happen. ... So

when we divided up the compiler, I told everyone else to back off and loudly claimed DFA for my own. Fine, everyone said.

Lamping and I spent the rest of the summer failing. Taking trips to the Stanford library to look up papers. Hashing things out on white boards. Staring into space. Writing little bits of experimental code. Failing. Finding out **why** no one had ever provided DFA optimization for Scheme. In short, the fundamental item the classical data-flow analysis algorithms need to operate is not available in a Scheme program. It was really depressing. I was making more money than I'd ever made in my life (\$600/wk). ... It was **not** a good summer. I slunk back to CMU with my tail between my legs, having contributed not one line of code.

About three years after the summer, I **finally** figured out how to do data-flow analysis for Scheme, which ended a long, pretty unhappy period in my life.

I officially switched from being an AI student to being a PL student...

Olin Shivers. *The History of T*. <http://www.paulgraham.com/thist.html>

4.1 Handling Lambda

Defs can be given with a recursive definition: In a call $c:(f\dots a_i\dots)$, $\forall l \in \text{Defs}(f)$ of the form $(\lambda \dots v_i \dots) \dots$,

$$\text{Defs}(a_i) \subset \text{Defs}(v_i) \quad [\text{LAM}]$$

I.e. if lambda l can be called from call site c , then l 's i th variable can evaluate to any of the lambdas that c 's i th argument can evaluate to.

Olin Shivers. *Control Flow Analysis in Scheme*. PLDI 1988

271

Control Flow Analysis in Scheme

Olin Shivers
Carnegie Mellon University
<http://www.cse.cmu.edu/~shivers/>

Abstract

Traditional flow analysis techniques, such as the ones typically employed by optimizing Forth compilers, do not work for Scheme. This paper describes a new flow analysis technique — control flow analysis — which is applicable to Scheme. The basic idea is to use the control flow information gathered by control flow analysis to assist in performing a traditional flow analysis problem: induction variable management.

The techniques presented in this paper are backed up by working examples. They are also compared to the flow analysis used to select compilers, such as Common Lisp and ML.

1 The Task

Flow analysis is a traditional optimizing compiler technique for determining useful information about a program at compile time. Doing so can lead to significant improvements in the performance of a program. A flow analysis problem is a question of the form:

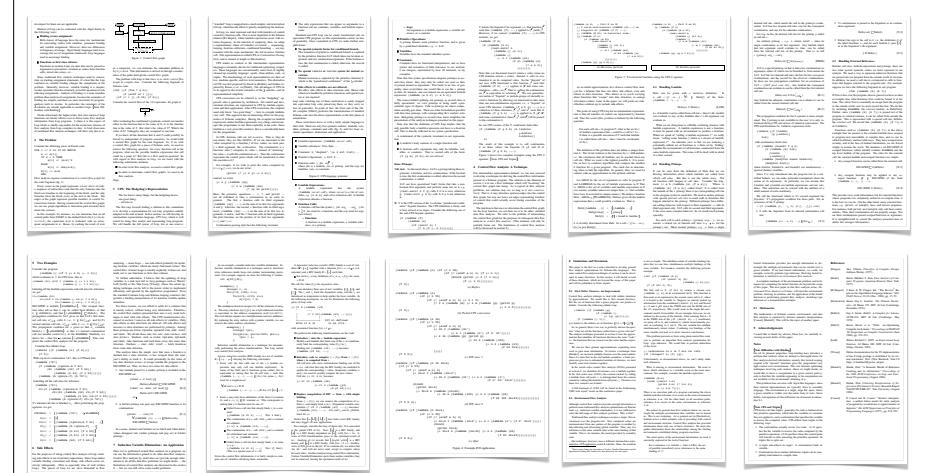
• Argument analysis: What is the range of values that a given reference to an integer variable is constrained to be while? For example, what is the range of i in $\text{for } i \text{ from } 1 \text{ to } n \text{ do } \dots$ to do array bounds checking at compile time?

• Loop invariant detection: Do all possible prior assignments to a given variable reference be excluded in computing loops?

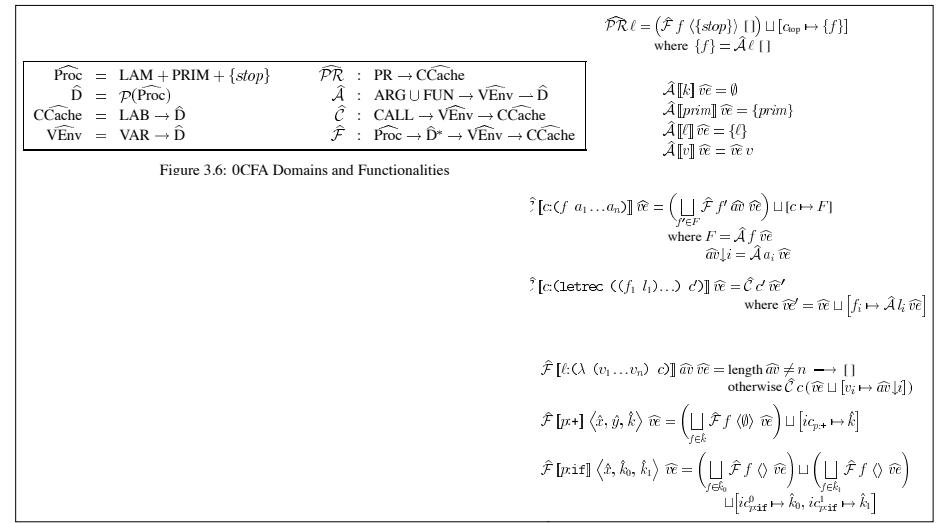
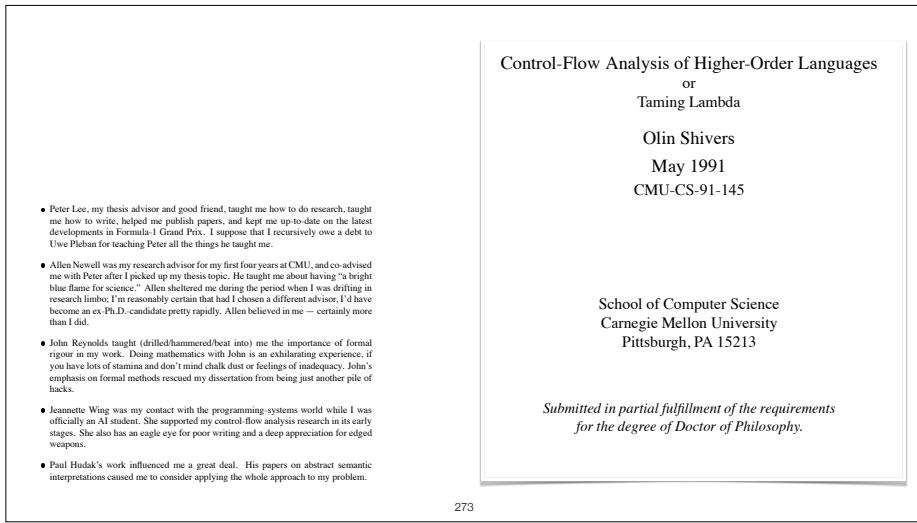
Over the last thirty years, standard techniques have been developed for solving these kinds of problems in various imperative and Algol-like languages (e.g., Pascal, C, Ada, Basic, and, finally

To appear at the ACM SIGPLAN '88 Conference on Programming Languages Design and Implementation, Atlanta, Georgia, June 22-24, 1988.

So it is clear that flow analysis has much potential with respect to compiling Lisp programs. Unfortunately, this potential is not realized because the traditional flow analysis techniques for the Algol family of languages that the traditional techniques



272



Syntax Directed 0-CFA Analysis

Reformulate the abstract specification:

- (i) Syntax directed specification
- (ii) Constructing a finite set of constraints
- (iii) Compute the least solution of the set of constraints

Common Phenomenon

A specification \models_A is reformulated into a specification \models_B ensuring that

$$(\widehat{\mathcal{C}}, \widehat{\rho}) \models_A e_\star \Leftarrow (\widehat{\mathcal{C}}, \widehat{\rho}) \models_B e_\star$$

so that " \models_B " is a *safe approximation* to " \models_A " and hence the best (i.e. least) solution to " $\models_B e_\star$ " will also be a solution to " $\models_A e_\star$ ".

If additionally

$$(\widehat{\mathcal{C}}, \widehat{\rho}) \models_A e_\star \Rightarrow (\widehat{\mathcal{C}}, \widehat{\rho}) \models_B e_\star$$

then we can be assured that *no solutions are lost* and hence the best (i.e. least) solution to " $\models_B e_\star$ " will also be the best (i.e. least) solution to " $\models_A e_\star$ ".

Syntax Directed Specification (1)

$$(\hat{C}, \hat{\rho}) \models_s (\text{fn } x \Rightarrow e_0)^\ell \quad \text{iff} \quad \{\text{fn } x \Rightarrow e_0\} \subseteq \hat{C}(\ell) \wedge \\ (\hat{C}, \hat{\rho}) \models_s e_0$$

$$(\hat{C}, \hat{\rho}) \models_s (\text{fun } f \ x \Rightarrow e_0)^\ell \quad \text{iff} \quad \{\text{fun } f \ x \Rightarrow e_0\} \subseteq \hat{C}(\ell) \wedge \\ (\hat{C}, \hat{\rho}) \models_s e_0 \wedge \{\text{fun } f \ x \Rightarrow e_0\} \subseteq \hat{\rho}(f)$$

$$(\hat{C}, \hat{\rho}) \models_s (t_1^{\ell_1} t_2^{\ell_2})^\ell \quad \text{iff} \quad (\hat{C}, \hat{\rho}) \models_s t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models_s t_2^{\ell_2} \wedge \\ (\forall (\text{fn } x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1) : \hat{C}(\ell_2) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell)) \wedge \\ (\forall (\text{fun } f \ x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1) : \hat{C}(\ell_2) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell))$$

277

© Nielsen^2, Hankin

$$(\hat{C}, \hat{\rho}) \models_s c^\ell \text{ always}$$

$$(\hat{C}, \hat{\rho}) \models_s x^\ell \quad \text{iff} \quad \hat{\rho}(x) \subseteq \hat{C}(\ell)$$

$$(\hat{C}, \hat{\rho}) \models_s (\text{if } t_0^{\ell_0} \text{ then } t_1^{\ell_1} \text{ else } t_2^{\ell_2})^\ell \quad \text{iff} \quad (\hat{C}, \hat{\rho}) \models_s t_0^{\ell_0} \wedge \\ (\hat{C}, \hat{\rho}) \models_s t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models_s t_2^{\ell_2} \wedge \\ \hat{C}(\ell_1) \subseteq \hat{C}(\ell) \wedge \hat{C}(\ell_2) \subseteq \hat{C}(\ell)$$

$$(\hat{C}, \hat{\rho}) \models_s (\text{let } x = t_1^{\ell_1} \text{ in } t_2^{\ell_2})^\ell \quad \text{iff} \quad (\hat{C}, \hat{\rho}) \models_s t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models_s t_2^{\ell_2} \wedge \\ \hat{C}(\ell_1) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_2) \subseteq \hat{C}(\ell)$$

$$(\hat{C}, \hat{\rho}) \models_s (t_1^{\ell_1} \ op \ t_2^{\ell_2})^\ell \quad \text{iff} \quad (\hat{C}, \hat{\rho}) \models_s t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models_s t_2^{\ell_2}$$

278

© Nielsen^2, Hankin

Constraint Based 0-CFA Analysis

$C_\star[[e_\star]]$ is a set of constraints of the form

$$\text{lhs} \subseteq \text{rhs}$$

$$\{t\} \subseteq \text{rhs}' \Rightarrow \text{lhs} \subseteq \text{rhs}$$

where

$$\text{rhs} ::= C(\ell) \mid r(x)$$

$$\text{lhs} ::= C(\ell) \mid r(x) \mid \{t\}$$

and all occurrences of t are of the form $\text{fn } x \Rightarrow e_0$ or $\text{fun } f \ x \Rightarrow e_0$

279

© Nielsen^2, Hankin

Constraint Based Control Flow Analysis (1)

$$C_\star[[\text{fn } x \Rightarrow e_0]^\ell] = \{\{ \text{fn } x \Rightarrow e_0 \} \subseteq C(\ell)\} \cup C_\star[[e_0]]$$

$$C_\star[[\text{fun } f \ x \Rightarrow e_0]^\ell] = \{\{ \text{fun } f \ x \Rightarrow e_0 \} \subseteq C(\ell)\} \cup C_\star[[e_0]] \cup \{\{ \text{fun } f \ x \Rightarrow e_0 \} \subseteq r(f)\}$$

$$C_\star[(t_1^{\ell_1} t_2^{\ell_2})^\ell] = C_\star[[t_1^{\ell_1}]] \cup C_\star[[t_2^{\ell_2}]] \cup \{\{t\} \subseteq C(\ell_1) \Rightarrow C(\ell_2) \subseteq r(x) \mid t = (\text{fn } x \Rightarrow t_0^{\ell_0}) \in \text{Term}_\star\} \cup \{\{t\} \subseteq C(\ell_1) \Rightarrow C(\ell_0) \subseteq C(\ell) \mid t = (\text{fn } x \Rightarrow t_0^{\ell_0}) \in \text{Term}_\star\} \cup \{\{t\} \subseteq C(\ell_1) \Rightarrow C(\ell_2) \subseteq r(x) \mid t = (\text{fun } f \ x \Rightarrow t_0^{\ell_0}) \in \text{Term}_\star\} \cup \{\{t\} \subseteq C(\ell_1) \Rightarrow C(\ell_0) \subseteq C(\ell) \mid t = (\text{fun } f \ x \Rightarrow t_0^{\ell_0}) \in \text{Term}_\star\}$$

(Eager rather than lazy unfolding – easy but costly.)

280

© Nielsen^2, Hankin

Constraint Based Control Flow Analysis (2)

$$\mathcal{C}_\star[\![c^\ell]\!] = \emptyset$$

$$\mathcal{C}_\star[\![x^\ell]\!] = \{r(x) \subseteq C(\ell)\}$$

$$\begin{aligned} \mathcal{C}_\star[\!(\text{if } t_0^{\ell_0} \text{ then } t_1^{\ell_1} \text{ else } t_2^{\ell_2})^\ell]\!] &= \mathcal{C}_\star[\![t_0^{\ell_0}]\!] \cup \mathcal{C}_\star[\![t_1^{\ell_1}]\!] \cup \mathcal{C}_\star[\![t_2^{\ell_2}]\!] \\ &\cup \{C(\ell_1) \subseteq C(\ell)\} \\ &\cup \{C(\ell_2) \subseteq C(\ell)\} \end{aligned}$$

$$\begin{aligned} \mathcal{C}_\star[\!(\text{let } x = t_1^{\ell_1} \text{ in } t_2^{\ell_2})^\ell]\!] &= \mathcal{C}_\star[\![t_1^{\ell_1}]\!] \cup \mathcal{C}_\star[\![t_2^{\ell_2}]\!] \\ &\cup \{C(\ell_1) \subseteq r(x)\} \cup \{C(\ell_2) \subseteq C(\ell)\} \end{aligned}$$

$$\mathcal{C}_\star[\!(t_1^{\ell_1} \text{ op } t_2^{\ell_2})^\ell]\!] = \mathcal{C}_\star[\![t_1^{\ell_1}]\!] \cup \mathcal{C}_\star[\![t_2^{\ell_2}]\!]$$

281

© Nielsen^2, Hankin

$$\begin{aligned} \mathcal{C}_\star[\!((\text{fn } x \Rightarrow x^1)^2 (\text{fn } y \Rightarrow y^3)^4)^5]\!] &= \\ &\{ \{ \text{fn } x \Rightarrow x^1 \} \subseteq C(2), \\ &r(x) \subseteq C(1), \\ &\{ \text{fn } y \Rightarrow y^3 \} \subseteq C(4), \\ &r(y) \subseteq C(3), \\ &\{ \text{fn } x \Rightarrow x^1 \} \subseteq C(2) \Rightarrow C(4) \subseteq r(x), \\ &\{ \text{fn } x \Rightarrow x^1 \} \subseteq C(2) \Rightarrow C(1) \subseteq C(5), \\ &\{ \text{fn } y \Rightarrow y^3 \} \subseteq C(2) \Rightarrow C(4) \subseteq r(y), \\ &\{ \text{fn } y \Rightarrow y^3 \} \subseteq C(2) \Rightarrow C(3) \subseteq C(5) \} \end{aligned}$$

282

© Nielsen^2, Hankin

Preservation of Solutions

Translating syntactic entities to sets of terms:

$$\begin{aligned} (\hat{C}, \hat{\rho})[\!C(\ell)]\! &= \hat{C}(\ell) \\ (\hat{C}, \hat{\rho})[\!r(x)]\! &= \hat{\rho}(x) \\ (\hat{C}, \hat{\rho})[\!\{t\}\!] &= \{t\} \end{aligned}$$

Satisfaction relation for constraints: $(\hat{C}, \hat{\rho}) \models_c (lhs \subseteq rhs)$

$$\begin{aligned} (\hat{C}, \hat{\rho}) \models_c (lhs \subseteq rhs) \\ \text{iff } & (\hat{C}, \hat{\rho})[\!lhs]\! \subseteq (\hat{C}, \hat{\rho})[\!rhs]\! \\ (\hat{C}, \hat{\rho}) \models_c (\{t\} \subseteq rhs' \Rightarrow lhs \subseteq rhs) \\ \text{iff } & (\{t\} \subseteq (\hat{C}, \hat{\rho})[\!rhs']\! \wedge (\hat{C}, \hat{\rho})[\!lhs]\! \subseteq (\hat{C}, \hat{\rho})[\!rhs]\!) \\ & \vee (\{t\} \not\subseteq (\hat{C}, \hat{\rho})[\!rhs']\!) \end{aligned}$$

Proposition: $(\hat{C}, \hat{\rho}) \models_s e_\star$ if and only if $(\hat{C}, \hat{\rho}) \models_c \mathcal{C}_\star[\!e_\star]\!$.

283

© Nielsen^2, Hankin

Solving the Constraints (1)

Input: a set of constraints $\mathcal{C}_\star[\!e_\star]\!$

Output: the least solution $(\hat{C}, \hat{\rho})$ to the constraints

Data structures: a graph with one node for each $C(\ell)$ and $r(x)$ (where $\ell \in \text{Lab}_\star$ and $x \in \text{Var}_\star$) and zero, one or two edges for each constraint in $\mathcal{C}_\star[\!e_\star]\!$

- **W:** the worklist of the nodes whose outgoing edges should be traversed
- **D:** an array that for each node gives an element of $\widehat{\text{Val}}_\star$
- **E:** an array that for each node gives a list of constraints influenced (and outgoing edges)

Auxiliary procedure:

procedure add(q, d) is if $\neg (d \subseteq D[q])$ then $D[q] := D[q] \cup d$;
 $W := \text{cons}(q, W)$;

284

© Nielsen^2, Hankin

Solving the Constraints (2)

Step 1 Initialisation

```
W := nil;
for q in Nodes do D[q] := ∅; E[q] := nil;
```

Step 2 Building the graph

```
for cc in  $C_*[[e_*]]$  do
  case cc of {t} ⊆ p: add(p, {t});
  p1 ⊆ p2: E[p1] := cons(cc, E[p1]);
  {t} ⊆ p ⇒ p1 ⊆ p2: E[p1] := cons(cc, E[p1]);
  E[p] := cons(cc, E[p]);
```

Step 3 Iteration

```
while W ≠ nil do
  q := head(W); W := tail(W);
  for cc in E[q] do
    case cc of p1 ⊆ p2: add(p2, D[p1]);
    {t} ⊆ p ⇒ p1 ⊆ p2: if t ∈ D[p] then add(p2, D[p1]);
```

Step 4 Recording the solution

```
for ℓ in Lab* do  $\hat{C}(\ell) := D[C(\ell)]$ ; for x in Var* do  $\hat{\rho}(x) := D[r(x)]$ ;
```

285

© Nielsen^2, Hankin

Example:

Initialisation of data structures

| p | $D[p]$ | $E[p]$ |
|--------|-------------|--|
| $C(1)$ | \emptyset | $[id_x \subseteq C(2) \Rightarrow C(1) \subseteq C(5)]$ |
| $C(2)$ | id_x | $[id_y \subseteq C(2) \Rightarrow C(3) \subseteq C(5), id_y \subseteq C(2) \Rightarrow C(4) \subseteq r(y), id_x \subseteq C(2) \Rightarrow C(1) \subseteq C(5), id_x \subseteq C(2) \Rightarrow C(4) \subseteq r(x)]$ |
| $C(3)$ | \emptyset | $[id_y \subseteq C(2) \Rightarrow C(3) \subseteq C(5)]$ |
| $C(4)$ | id_y | $[id_y \subseteq C(2) \Rightarrow C(4) \subseteq r(y), id_x \subseteq C(2) \Rightarrow C(4) \subseteq r(x)]$ |
| $C(5)$ | \emptyset | $[]$ |
| $r(x)$ | \emptyset | $[r(x) \subseteq C(1)]$ |
| $r(y)$ | \emptyset | $[r(y) \subseteq C(3)]$ |

286

© Nielsen^2, Hankin

Example:

Iteration steps

| W | $[C(4), C(2)]$ | $[r(x), C(2)]$ | $[C(1), C(2)]$ | $[C(5), C(2)]$ | $[C(2)]$ | $[]$ |
|--------|----------------|----------------|----------------|----------------|-------------|-------------|
| p | $D[p]$ | $D[p]$ | $D[p]$ | $D[p]$ | $D[p]$ | $D[p]$ |
| $C(1)$ | \emptyset | \emptyset | id_y | id_y | id_y | id_y |
| $C(2)$ | id_x | id_x | id_x | id_x | id_x | id_x |
| $C(3)$ | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset |
| $C(4)$ | id_y | id_y | id_y | id_y | id_y | id_y |
| $C(5)$ | \emptyset | \emptyset | \emptyset | id_y | id_y | id_y |
| $r(x)$ | \emptyset | id_y | id_y | id_y | id_y | id_y |
| $r(y)$ | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset |

287

© Nielsen^2, Hankin

Correctness:

Given input $C_*[[e_*]]$ the worklist algorithm terminates and the result $(\hat{C}, \hat{\rho})$ produced by the algorithm satisfies

$$(\hat{C}, \hat{\rho}) = \bigcap \{(\hat{C}', \hat{\rho}') \mid (\hat{C}', \hat{\rho}') \models_c C_*[[e_*]]\}$$

and hence it is the least solution to $C_*[[e_*]]$.

Complexity:

The algorithm takes at most $O(n^3)$ steps if the original expression e_* has size n .

288

© Nielsen^2, Hankin

Adding Data Flow Analysis

Idea: extend the set $\widehat{\text{Val}}$ to contain other abstract values than just abstractions

- powerset (possibly finite)
- complete lattice (possibly satisfying Ascending Chain Condition)

289

© Nielsen^2, Hankin

Abstract Values as Powersets

Let Data be a set of *abstract data values* (i.e. abstract properties of booleans and integers)

$$\hat{v} \in \widehat{\text{Val}}_d = \mathcal{P}(\text{Term} \cup \text{Data}) \quad \text{abstract values}$$

For each constant $c \in \text{Const}$ we need an element $d_c \in \text{Data}$

For each operator $op \in \text{Op}$ we need a total function

$$\hat{o}p : \widehat{\text{Val}}_d \times \widehat{\text{Val}}_d \rightarrow \widehat{\text{Val}}_d$$

typically

$$\hat{v}_1 \hat{o}p \hat{v}_2 = \bigcup \{d_{op}(d_1, d_2) \mid d_1 \in \hat{v}_1 \cap \text{Data}, d_2 \in \hat{v}_2 \cap \text{Data}\}$$

for some $d_{op} : \text{Data} \times \text{Data} \rightarrow \mathcal{P}(\text{Data})$

290

© Nielsen^2, Hankin

Example: Detection of Signs Analysis

$$\text{Data}_{\text{sign}} = \{\text{tt}, \text{ff}, -, 0, +\}$$

$$d_{\text{true}} = \text{tt}$$

$$d_7 = +$$

$\widehat{+}$ is defined from

| d_+ | tt | ff | - | 0 | + |
|-------|-------------|-------------|---------------|-------------|---------------|
| tt | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset |
| ff | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset |
| - | \emptyset | \emptyset | $\{-\}$ | $\{-\}$ | $\{-, 0, +\}$ |
| 0 | \emptyset | \emptyset | $\{-\}$ | $\{0\}$ | $\{+\}$ |
| + | \emptyset | \emptyset | $\{-, 0, +\}$ | $\{+\}$ | $\{+\}$ |

291

© Nielsen^2, Hankin

Abstract Values as Powersets (1)

$$(\hat{C}, \hat{p}) \models_d (\text{fn } x \Rightarrow e_0)^\ell \quad \text{iff} \quad \{\text{fn } x \Rightarrow e_0\} \subseteq \hat{C}(\ell)$$

$$(\hat{C}, \hat{p}) \models_d (\text{fun } f x \Rightarrow e_0)^\ell \quad \text{iff} \quad \{\text{fun } f x \Rightarrow e_0\} \subseteq \hat{C}(\ell)$$

$$\begin{aligned} (\hat{C}, \hat{p}) \models_d (t_1^{\ell_1} t_2^{\ell_2})^\ell &\quad \text{iff} \quad (\hat{C}, \hat{p}) \models_d t_1^{\ell_1} \wedge (\hat{C}, \hat{p}) \models_d t_2^{\ell_2} \wedge \\ &\quad (\forall (\text{fn } x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1) : \\ &\quad \quad (\hat{C}, \hat{p}) \models_d t_0^{\ell_0} \wedge \\ &\quad \quad \hat{C}(\ell_2) \subseteq \hat{p}(x) \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell)) \wedge \\ &\quad (\forall (\text{fun } f x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1) : \\ &\quad \quad (\hat{C}, \hat{p}) \models_d t_0^{\ell_0} \wedge \\ &\quad \quad \hat{C}(\ell_2) \subseteq \hat{p}(x) \wedge \hat{C}(\ell_0) \subseteq \hat{C}(\ell) \wedge \\ &\quad \quad \{\text{fun } f x \Rightarrow t_0^{\ell_0}\} \subseteq \hat{p}(f)) \end{aligned}$$

292

© Nielsen^2, Hankin

$$(\hat{C}, \hat{\rho}) \models_d c^\ell \quad \text{iff} \quad \{dc\} \subseteq \hat{C}(\ell)$$

$$(\hat{C}, \hat{\rho}) \models_d x^\ell \quad \text{iff} \quad \hat{\rho}(x) \subseteq \hat{C}(\ell)$$

$$(\hat{C}, \hat{\rho}) \models_d (\text{if } t_0^{\ell_0} \text{ then } t_1^{\ell_1} \text{ else } t_2^{\ell_2})^\ell \quad \text{iff} \quad (\hat{C}, \hat{\rho}) \models_d t_0^{\ell_0} \wedge$$

$$(d_{\text{true}} \in \hat{C}(\ell_0) \Rightarrow ((\hat{C}, \hat{\rho}) \models_d t_1^{\ell_1} \wedge \hat{C}(\ell_1) \subseteq \hat{C}(\ell))) \wedge$$

$$(d_{\text{false}} \in \hat{C}(\ell_0) \Rightarrow ((\hat{C}, \hat{\rho}) \models_d t_2^{\ell_2} \wedge \hat{C}(\ell_2) \subseteq \hat{C}(\ell)))$$

$$(\hat{C}, \hat{\rho}) \models_d (\text{let } x = t_1^{\ell_1} \text{ in } t_2^{\ell_2})^\ell \quad \text{iff} \quad (\hat{C}, \hat{\rho}) \models_d t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models_d t_2^{\ell_2} \wedge \hat{C}(\ell_1) \subseteq \hat{\rho}(x) \wedge \hat{C}(\ell_2) \subseteq \hat{C}(\ell)$$

$$(\hat{C}, \hat{\rho}) \models_d (t_1^{\ell_1} \text{ op } t_2^{\ell_2})^\ell \quad \text{iff} \quad (\hat{C}, \hat{\rho}) \models_d t_1^{\ell_1} \wedge (\hat{C}, \hat{\rho}) \models_d t_2^{\ell_2} \wedge \hat{C}(\ell_1) \widehat{\text{op}} \hat{C}(\ell_2) \subseteq \hat{C}(\ell)$$

293

© Nielsen^2, Hankin

Example:

```
(let f = (fn x => (if (x^1 > 0^2)^3 then (fn y => y^4)^5
                      else (fn z => 25^6)^7)^8)^9
    in ((f^10 3^11)^12 0^13)^14)^15
```

A pure 0-CFA analysis will not be able to discover that the `else`-branch of the conditional will never be executed.

When we combine the analysis with a Detection of Signs Analysis then the analysis can determine that only `fn y => y^4` is a possible abstraction at label 12.

294

© Nielsen^2, Hankin

| | $(\hat{C}, \hat{\rho})$ | $(\hat{C}, \hat{\rho})$ |
|----|--|---------------------------------|
| 1 | \emptyset | $\{+\}$ |
| 2 | \emptyset | $\{0\}$ |
| 3 | \emptyset | $\{\text{tt}\}$ |
| 4 | \emptyset | $\{0\}$ |
| 5 | $\{\text{fn } y => y^4\}$ | $\{\text{fn } y => y^4\}$ |
| 6 | \emptyset | \emptyset |
| 7 | $\{\text{fn } z => 25^6\}$ | \emptyset |
| 8 | $\{\text{fn } y => y^4, \text{ fn } z => 25^6\}$ | $\{\text{fn } y => y^4\}$ |
| 9 | $\{\text{fn } x => (\dots)^8\}$ | $\{\text{fn } x => (\dots)^8\}$ |
| 10 | $\{\text{fn } x => (\dots)^8\}$ | $\{\text{fn } x => (\dots)^8\}$ |
| 11 | \emptyset | $\{+\}$ |
| 12 | $\{\text{fn } y => y^4, \text{ fn } z => 25^6\}$ | $\{\text{fn } y => y^4\}$ |
| 13 | \emptyset | $\{0\}$ |
| 14 | \emptyset | $\{0\}$ |
| 15 | \emptyset | $\{0\}$ |
| f | $\{\text{fn } x => (\dots)^8\}$ | $\{\text{fn } x => (\dots)^8\}$ |
| x | \emptyset | $\{+\}$ |
| y | \emptyset | $\{0\}$ |
| z | \emptyset | \emptyset |

© Nielsen^2, Hankin

Abstract Values as Complete Lattices

A *monotone structure* consists of:

- a complete lattice L , and
- a set \mathcal{F} of monotone functions of $L \times L \rightarrow L$.

An *instance* of a monotone structure consists of the structure (L, \mathcal{F}) and

- a mapping ι . from the constants $c \in \text{Const}$ to values in L , and
- a mapping f . from the binary operators $op \in \text{Op}$ to functions of \mathcal{F} .

296

© Nielsen^2, Hankin

Abstract Domains

For the Control Flow Analysis:

$$\begin{aligned}\hat{v} \in \widehat{\text{Val}} &= \mathcal{P}(\text{Term}) \quad \text{abstract values} \\ \hat{\rho} \in \widehat{\text{Env}} &= \text{Var} \rightarrow \widehat{\text{Val}} \quad \text{abstract environments} \\ \hat{C} \in \widehat{\text{Cache}} &= \text{Lab} \rightarrow \widehat{\text{Val}} \quad \text{abstract caches}\end{aligned}$$

For the Data Flow Analysis:

$$\begin{aligned}\hat{d} \in \widehat{\text{Data}} &= L \quad \text{abstract data values} \\ \hat{\delta} \in \widehat{\text{DEnv}} &= \text{Var} \rightarrow \widehat{\text{Data}} \quad \text{abstract data environments} \\ \hat{D} \in \widehat{\text{DCache}} &= \text{Lab} \rightarrow \widehat{\text{Data}} \quad \text{abstract data caches}\end{aligned}$$

297

© Nielsen^2, Hankin

Abstract Values as Complete Lattices (1)

$$\begin{aligned}(\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D (\text{fn } x \Rightarrow e_0)^\ell &\quad \text{iff} \quad \{\text{fn } x \Rightarrow e_0\} \subseteq \hat{C}(\ell) \\ (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D (\text{fun } f x \Rightarrow e_0)^\ell &\quad \text{iff} \quad \{\text{fun } f x \Rightarrow e_0\} \subseteq \hat{C}(\ell) \\ (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D (t_1^{\ell_1} t_2^{\ell_2})^\ell &\quad \text{iff} \quad (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_1^{\ell_1} \wedge (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_2^{\ell_2} \wedge \\ &\quad (\forall (\text{fn } x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1)) : (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_0^{\ell_0} \wedge \\ &\quad \hat{C}(\ell_2) \subseteq \hat{\rho}(x) \wedge \hat{D}(\ell_2) \sqsubseteq \hat{\delta}(x) \wedge \\ &\quad \hat{C}(\ell_0) \subseteq \hat{C}(\ell) \wedge \hat{D}(\ell_0) \sqsubseteq \hat{D}(\ell) \wedge \\ &\quad (\forall (\text{fun } f x \Rightarrow t_0^{\ell_0}) \in \hat{C}(\ell_1)) : (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_0^{\ell_0} \wedge \\ &\quad \hat{C}(\ell_2) \subseteq \hat{\rho}(x) \wedge \hat{D}(\ell_2) \sqsubseteq \hat{\delta}(x) \wedge \\ &\quad \hat{C}(\ell_0) \subseteq \hat{C}(\ell) \wedge \hat{D}(\ell_0) \sqsubseteq \hat{D}(\ell) \wedge \\ &\quad \{\text{fun } f x \Rightarrow t_0^{\ell_0}\} \subseteq \hat{\rho}(f)\end{aligned}$$

298

© Nielsen^2, Hankin

Abstract Values as Complete Lattices (2)

$$\begin{aligned}(\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D c^\ell &\quad \text{iff} \quad \iota_C \sqsubseteq \hat{D}(\ell) \\ (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D x^\ell &\quad \text{iff} \quad \hat{\rho}(x) \subseteq \hat{C}(\ell) \wedge \hat{\delta}(x) \sqsubseteq \hat{D}(\ell) \\ (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D (\text{if } t_0^{\ell_0} \text{ then } t_1^{\ell_1} \text{ else } t_2^{\ell_2})^\ell &\quad \text{iff} \quad (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_0^{\ell_0} \wedge \\ &\quad (\iota_{\text{true}} \sqsubseteq \hat{D}(\ell_0) \Rightarrow (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_1^{\ell_1} \wedge \\ &\quad \hat{C}(\ell_1) \subseteq \hat{C}(\ell) \wedge \hat{D}(\ell_1) \sqsubseteq \hat{D}(\ell)) \wedge \\ &\quad (\iota_{\text{false}} \sqsubseteq \hat{D}(\ell_0) \Rightarrow (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_2^{\ell_2} \wedge \\ &\quad \hat{C}(\ell_2) \subseteq \hat{C}(\ell) \wedge \hat{D}(\ell_2) \sqsubseteq \hat{D}(\ell))\end{aligned}$$

299

© Nielsen^2, Hankin

Abstract Values as Complete Lattices (3)

$$\begin{aligned}(\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D (\text{let } x = t_1^{\ell_1} \text{ in } t_2^{\ell_2})^\ell &\quad \text{iff} \quad (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_1^{\ell_1} \wedge \\ &\quad (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_2^{\ell_2} \wedge \\ &\quad \hat{C}(\ell_1) \subseteq \hat{\rho}(x) \wedge \hat{D}(\ell_1) \sqsubseteq \hat{\delta}(x) \wedge \hat{C}(\ell_2) \subseteq \hat{C}(\ell) \wedge \hat{D}(\ell_2) \sqsubseteq \hat{D}(\ell) \\ (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D (t_1^{\ell_1} op t_2^{\ell_2})^\ell &\quad \text{iff} \quad (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_1^{\ell_1} \wedge (\hat{C}, \hat{D}, \hat{\rho}, \hat{\delta}) \models_D t_2^{\ell_2} \wedge \\ &\quad f_{op}(\hat{D}(\ell_1), \hat{D}(\ell_2)) \sqsubseteq \hat{D}(\ell)\end{aligned}$$

300

© Nielsen^2, Hankin

| | $(\bar{C}, \bar{\rho})$ | $(\bar{C}, \bar{\rho})$ | $(\bar{C}, \bar{\rho})$ | $(\bar{D}, \bar{\delta})$ |
|----|---|--|--|---------------------------|
| 1 | \emptyset | $\{+\}$ | \emptyset | $\{+\}$ |
| 2 | \emptyset | $\{0\}$ | \emptyset | $\{0\}$ |
| 3 | \emptyset | $\{\text{tt}\}$ | \emptyset | $\{\text{tt}\}$ |
| 4 | \emptyset | $\{0\}$ | \emptyset | $\{0\}$ |
| 5 | $\{\text{fn } y \Rightarrow y^4\}$ | $\{\text{fn } y \Rightarrow y^4\}$ | $\{\text{fn } y \Rightarrow y^4\}$ | \emptyset |
| 6 | \emptyset | \emptyset | \emptyset | \emptyset |
| 7 | $\{\text{fn } z \Rightarrow 25^6\}$ | \emptyset | \emptyset | \emptyset |
| 8 | $\{\text{fn } y \Rightarrow y^4, \text{fn } z \Rightarrow 25^6\}$ | $\{\text{fn } y \Rightarrow y^4\}$ | $\{\text{fn } y \Rightarrow y^4\}$ | \emptyset |
| 9 | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | \emptyset |
| 10 | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | \emptyset |
| 11 | \emptyset | $\{+\}$ | \emptyset | $\{+\}$ |
| 12 | $\{\text{fn } y \Rightarrow y^4, \text{fn } z \Rightarrow 25^6\}$ | $\{\text{fn } y \Rightarrow y^4\}$ | $\{\text{fn } y \Rightarrow y^4\}$ | \emptyset |
| 13 | \emptyset | $\{0\}$ | \emptyset | $\{0\}$ |
| 14 | \emptyset | $\{0\}$ | \emptyset | $\{0\}$ |
| 15 | \emptyset | $\{0\}$ | \emptyset | $\{0\}$ |
| f | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | $\{\text{fn } x \Rightarrow (\dots)^8\}$ | \emptyset |
| x | \emptyset | $\{+\}$ | \emptyset | $\{+\}$ |
| y | \emptyset | $\{0\}$ | \emptyset | $\{0\}$ |
| z | \emptyset | \emptyset | \emptyset | \emptyset |

301

© Nielsen^2, Hankin

Adding Context Information

Mono-variant analysis: does not distinguish the various instances of variables and program points from one another. (Compare with context-insensitive interprocedural analysis.) **0-CFA** is a typical example.

Poly-variant analysis: distinguishes between the various instances of variables and program points. (Compare with context-sensitive interprocedural analysis.)

302

© Nielsen^2, Hankin

A purely syntactic solution:

Expand

```
(let f = (fn x => x) in ((f f) (fn y => y)))
```

into

```
let f1 = (fn x1 => x1)
  in let f2 = (fn x2 => x2) in (f1 f2) (fn y => y)
```

and analyse the expanded expression.

The 0-CFA analysis is now able to deduce that the overall expression will evaluate to $\text{fn } y \Rightarrow y$ only.

303

© Nielsen^2, Hankin

A purely semantic solution: Uniform k -CFA

Idea: extend the set $\widehat{\text{Val}}$ to include context information

In a (uniform) **k -CFA** a context δ records the last k dynamic call points; hence contexts will be sequences of labels of length at most k and they will be updated whenever a function application is analysed. (Compare call strings of length at most k .)

304

© Nielsen^2, Hankin

Abstract Domains

$$\begin{aligned}
 \delta &\in \Delta = \text{Lab}^{\leq k} && \text{context information} \\
 ce &\in \text{CEnv} = \text{Var} \rightarrow \Delta && \text{context environments} \\
 \hat{v} &\in \widehat{\text{Val}} = \mathcal{P}(\text{Term} \times \text{CEnv}) && \text{abstract values} \\
 \hat{\rho} &\in \widehat{\text{Env}} = (\text{Var} \times \Delta) \rightarrow \widehat{\text{Val}} && \text{abstract environments} \\
 \hat{c} &\in \widehat{\text{Cache}} = (\text{Lab} \times \Delta) \rightarrow \widehat{\text{Val}} && \text{abstract caches}
 \end{aligned}$$

(Uniform because Δ used both for $\widehat{\text{Env}}$ and $\widehat{\text{Cache}}$.)

305

© Nielsen^2, Hankin

Acceptability Relation

$$(\hat{c}, \hat{\rho}) \models_{\delta}^{ce} e$$

where

- ce is the current context environment – will be changed when new bindings are made
- δ is the current context – will be changed when functions are called

Idea: The formula expresses that $(\hat{c}, \hat{\rho})$ is an acceptable analysis of e in the *context* specified by ce and δ .

306

© Nielsen^2, Hankin

Control Flow Analysis with Context (1)

$$\begin{aligned}
 (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} (\text{fn } x \Rightarrow e_0)^{\ell} &\quad \text{iff} \quad \{(\text{fn } x \Rightarrow e_0, ce)\} \subseteq \hat{c}(\ell, \delta) \\
 (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} (\text{fun } f \ x \Rightarrow e_0)^{\ell} &\quad \text{iff} \quad \{(\text{fun } f \ x \Rightarrow e_0, ce)\} \subseteq \hat{c}(\ell, \delta) \\
 (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} (t_1^{\ell_1} t_2^{\ell_2})^{\ell} &\quad \text{iff} \quad (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} t_1^{\ell_1} \wedge (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} t_2^{\ell_2} \wedge \\
 &\quad (\forall (\text{fn } x \Rightarrow t_0^{\ell_0}, ce_0) \in \hat{c}(\ell_1, \delta) : \\
 &\quad (\hat{c}, \hat{\rho}) \models_{\delta_0}^{ce_0} t_0^{\ell_0} \wedge \hat{c}(\ell_2, \delta) \subseteq \hat{\rho}(x, \delta_0) \wedge \hat{c}(\ell_0, \delta_0) \subseteq \hat{c}(\ell, \delta) \wedge \\
 &\quad \text{where } \delta_0 = [\delta, \ell]_k \text{ and } ce'_0 = ce_0[x \mapsto \delta_0]) \wedge \\
 &\quad (\forall (\text{fun } f \ x \Rightarrow t_0^{\ell_0}, ce_0) \in \hat{c}(\ell_1, \delta) : \\
 &\quad (\hat{c}, \hat{\rho}) \models_{\delta_0}^{ce_0} t_0^{\ell_0} \wedge \hat{c}(\ell_2, \delta) \subseteq \hat{\rho}(x, \delta_0) \wedge \hat{c}(\ell_0, \delta_0) \subseteq \hat{c}(\ell, \delta) \wedge \\
 &\quad \{(\text{fun } f \ x \Rightarrow t_0^{\ell_0}, ce)\} \subseteq \hat{\rho}(f, \delta_0) \\
 &\quad \text{where } \delta_0 = [\delta, \ell]_k \text{ and } ce'_0 = ce_0[f \mapsto \delta_0, x \mapsto \delta_0])
 \end{aligned}$$

© Nielsen^2, Hankin

Control Flow Analysis with Context (2)

$$\begin{aligned}
 (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} c^{\ell} &\quad \text{always} \\
 (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} x^{\ell} &\quad \text{iff} \quad \hat{\rho}(x, ce(x)) \subseteq \hat{c}(\ell, \delta) \\
 (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} (\text{if } t_0^{\ell_0} \text{ then } t_1^{\ell_1} \text{ else } t_2^{\ell_2})^{\ell} &\quad \text{iff} \quad (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} t_0^{\ell_0} \wedge (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} t_1^{\ell_1} \wedge (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} t_2^{\ell_2} \wedge \\
 &\quad \hat{c}(\ell_1, \delta) \subseteq \hat{c}(\ell, \delta) \wedge \hat{c}(\ell_2, \delta) \subseteq \hat{c}(\ell, \delta) \\
 (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} (\text{let } x = t_1^{\ell_1} \text{ in } t_2^{\ell_2})^{\ell} &\quad \text{iff} \quad (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} t_1^{\ell_1} \wedge (\hat{c}, \hat{\rho}) \models_{\delta}^{ce'} t_2^{\ell_2} \wedge \\
 &\quad \hat{c}(\ell_1, \delta) \subseteq \hat{\rho}(x, \delta) \wedge \hat{c}(\ell_2, \delta) \subseteq \hat{c}(\ell, \delta) \\
 &\quad \text{where } ce' = ce[x \mapsto \delta] \\
 (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} (t_1^{\ell_1} \text{ op } t_2^{\ell_2})^{\ell} &\quad \text{iff} \quad (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} t_1^{\ell_1} \wedge (\hat{c}, \hat{\rho}) \models_{\delta}^{ce} t_2^{\ell_2}
 \end{aligned}$$

© Nielsen^2, Hankin

Example:

$(\text{let } f = (\text{fn } x \Rightarrow x^1)^2 \text{ in } ((f^3 f^4)^5 (f \text{ fn } y \Rightarrow y^6)^7)^8)^9$

Contexts of interest for uniform 1-CFA:

- A: the initial context
- 5: the context when the application point labelled 5 has been passed
- 8: the context when the application point labelled 8 has been passed

Context environments of interest for uniform 1-CFA:

- $ce_0 = []$ the initial (empty) context environment
- $ce_1 = ce_0[f \mapsto \Lambda]$ the context environment for the analysis of the body of the let-construct
- $ce_2 = ce_0[x \mapsto 5]$ the context environment used for the analysis of the body of f initiated at the application point 5
- $ce_3 = ce_0[x \mapsto 8]$ the context environment used for the analysis of the body of f initiated at the application point 8.

309

© Nielsen^2, Hankin

Example: Let us take \hat{C}_{id}' and $\hat{\rho}_{id}'$ to be:

$$\begin{aligned}\hat{C}_{id}'(1, 5) &= \{(f \text{ fn } x \Rightarrow x^1, ce_0)\} & \hat{C}_{id}'(1, 8) &= \{(f \text{ fn } y \Rightarrow y^6, ce_0)\} \\ \hat{C}_{id}'(2, \Lambda) &= \{(f \text{ fn } x \Rightarrow x^1, ce_0)\} & \hat{C}_{id}'(3, \Lambda) &= \{(f \text{ fn } x \Rightarrow x^1, ce_0)\} \\ \hat{C}_{id}'(4, \Lambda) &= \{(f \text{ fn } x \Rightarrow x^1, ce_0)\} & \hat{C}_{id}'(5, \Lambda) &= \{(f \text{ fn } x \Rightarrow x^1, ce_0)\} \\ \hat{C}_{id}'(7, \Lambda) &= \{(f \text{ fn } y \Rightarrow y^6, ce_0)\} & \hat{C}_{id}'(8, \Lambda) &= \{(f \text{ fn } y \Rightarrow y^6, ce_0)\} \\ \hat{C}_{id}'(9, \Lambda) &= \{(f \text{ fn } y \Rightarrow y^6, ce_0)\} & \\ \hat{\rho}_{id}'(f, \Lambda) &= \{(f \text{ fn } x \Rightarrow x^1, ce_0)\} & \\ \hat{\rho}_{id}'(x, 5) &= \{(f \text{ fn } x \Rightarrow x^1, ce_0)\} & \hat{\rho}_{id}'(x, 8) &= \{(f \text{ fn } y \Rightarrow y^6, ce_0)\}\end{aligned}$$

This is an acceptable analysis result:

$$(\hat{C}_{id}', \hat{\rho}_{id}') \models_{\Lambda}^{ce_0} (\text{let } f = (f \text{ fn } x \Rightarrow x^1)^2 \text{ in } ((f^3 f^4)^5 (f \text{ fn } y \Rightarrow y^6)^7)^8)^9$$

310

© Nielsen^2, Hankin

Complexity

Uniform k -CFA has exponential worst case complexity even when $k = 1$

Assume that the expression has size n and that it has p different variables. Then Δ has $O(n)$ elements and hence there will be $O(p \cdot n)$ different pairs (x, δ) and $O(n^2)$ different pairs (ℓ, δ) . This means that $(\hat{C}, \hat{\rho})$ can be seen as an $O(n^2)$ tuple of values from $\widehat{\text{Val}}$. Since $\widehat{\text{Val}}$ itself is a powerset of pairs of the form (t, ce) and there are $O(n \cdot n^p)$ such pairs it follows that $\widehat{\text{Val}}$ has height $O(n \cdot n^p)$. Since $O(p) = O(n)$ we have the exponential worst case complexity.

0-CFA analysis has polynomial worst case complexity

It corresponds to letting Δ be a singleton. Repeating the above calculations we can see $(\hat{C}, \hat{\rho})$ as an $O(p + n)$ tuple of values from $\widehat{\text{Val}}$, and $\widehat{\text{Val}}$ will be a lattice of height $O(n)$.

© Nielsen^2, Hankin

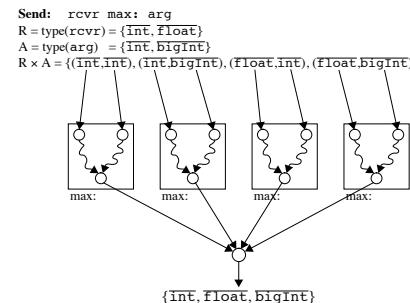


Fig. 6. CPA forms the cartesian product of the receiver type and argument types of the send it is analyzing. Each member of the product is propagated to a template reserved exclusively for that receiver and argument combination. The result type of the send is obtained by collecting the contributions from each template.

Ole Agesen, The Cartesian Product Algorithm: Simple and Precise Type Inference of Parametric Polymorphism ECOOP'95

Published in ECOOP'95 Conference Proceedings, Århus, Denmark, August 1995 (Springer-Verlag LNCS)

The Cartesian Product Algorithm

Simple and Precise Type Inference of Parametric Polymorphism

Ole Agesen
Computer Science Department, Stanford University
Stanford, CA 94304-5445
ole@cs.stanford.edu

Abstract. Concrete types and abstract types are different and serve different purposes. Concrete types, the focus of this paper, are essential to support compilation, application delivery, and debugging. Abstract types, on the other hand, are used to support polymorphism and inheritance. Inheritance relations because their presence limits polymorphism acceptability. This leaves us with type inference. Unfortunately, while polymorphism denotes the use of type inference, it has also been misinterpreted as type inference.

We review previous type inference algorithms that analyze code with parametric polymorphism and inheritance. We then propose a new algorithm that is more precise than previous algorithms. Our previous algorithms and theirs directly with inheritance, rather than relying on a preprocessor to expand away. Last, we compare our algorithm with the Self system since late 1993. We present measures to document its performance and compare it against several previous algorithms.

Keywords: concrete types, abstract types, type inference, polymorphism, inheritance, Self.

1 Introduction

1.1 Concrete versus Abstract Type

Concrete types (a.k.a. implementation types or representation types) are sets of classes. For example, an expression like $x : \text{int}$ means that the variable x has type int . Abstract types (a.k.a. interface types or supertypes) are sets of classes. For example, an expression like $x : \text{List}[\text{int}]$ means that the variable x has concrete type $\text{List}[\text{int}]$.

In contrast, abstract types (a.k.a. interface types or principal types) capture abstract properties of objects. They may not distinguish between different implementations of the same abstract data type: both list and array may have the same abstract type $\text{List}[\text{int}]$. Concrete types provide detailed information since they distinguish over different implementations.

Concrete and abstract types are extremes in a spectrum of type systems. Class-types, as found in BETA, C++, and Eiffel, are neither fully abstract (it is impossible to express the type of any object that has a push and pop method) nor fully concrete (it is impossible to implement an object with class-type Stack , does not reveal the specific subclass of which it is an instance).

Concrete types are the focus of this paper. Previous publications have argued that they directly support solving several important problems in object-oriented programming environments:

- **Compilation.** Concrete types support important optimizations such as elimination of dynamic dispatch and method inlining [17, 21].

- **Execution.** Concrete types support efficient compact applications by enabling an extractor to sift through the code and extract the parts that are relevant to the application.

- **Debugging.** Concrete types help a debugger understand programs by allowing a browser to track control flow (and thereby data flow) through dynamically dispatched message sends [3].



Dave Ungar obtained a PhD from UC Berkeley in 1985 with David Patterson on *The Design and Evaluation of a High-Performance Smalltalk System*; it won the 1986 ACM Dissertation Award. He was an assistant professor at Stanford. In 1991, he joined Sun and became a distinguished engineer. In 2006 he was recognized as a Distinguished Engineer by the ACM and in 2010 a Fellow. Ungar's 1984 paper, *Generation Scavenging: A Non-disruptive High Performance Storage Reclamation Algorithm* won a Retrospective SIGSOFT Impact Paper Award. He was awarded the Dahl-Nygaard Senior Prize in 2009.



Urs Hözle is a Swiss software engineer and technology executive. He is the senior vice president of technical infrastructure and Google Fellow at Google. As Google's eighth employee and its first VP of Engineering, he has shaped much of Google's development processes and infrastructure.



Craig Chambers has been at Google since 2007.^[1] Prior to this, he was a Professor at the University of Washington. He received his B.S. from MIT in 1986 and a Ph.D. from Stanford in 1992. He is known for Self, which introduced prototypes as an alternative to classes, and code-splitting, a compilation technique that generates separate code paths for fast and general cases to speed execution of dynamically-typed programs.



Ole Agesen joined VMware in 1999 and was appointed as a VMware Fellow in 2012. During his 12-year tenure as a principal engineer, he has earned 52 issued patents, 25 of which are assigned to VMware and made major contributions to VMware vSphere. Agesen worked for Sun Microsystems as a technical lead for the company's first server-class Java Virtual Machine, ExactVM. He holds a master's from Aarhus in Denmark and a Ph.D. in from Stanford.



Jeff Dean joined Google in '99, he heads Artificial Intelligence. The projects he worked:
 • Spanner, a scalable, multi-version, globally distributed, and replicated database
 • BigTable, a large-scale semi-structured storage system^[2]
 • MapReduce, a system for large-scale data processing applications^[3]
 • LevelDB, an open-source on-disk key-value store
 • DistBelief, a machine-learning system for deep neural nets refactored into TensorFlow



Scaling Program Analysis

CS7575

• Lecture 5 – Abstract Interpretation - PoPA Chapter 5

314

© Nielsen^2, Hankin

ABSTRACT INTERPRETATION : A UNIFIED LATTICE MODEL FOR STATIC ANALYSIS OF PROGRAMS BY CONSTRUCTION OR APPROXIMATION OF FIXPOINTS

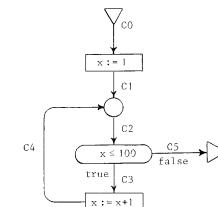
Patrick Cousot* and Radhia Cousot**

Laboratoire d'Informatique, U.S.M.G., BP. 53
38041 Grenoble cedex, France



Acknowledgments

We wish to thank M. Sintzoff for stimulating discussions. We were very lucky to have F. Blanc do the typing for us.



Let us note $[a, b]$ where $a \leq b$ the predicate $a \leq x \leq b$. The system of equations corresponding to the example is :

- (0) $C_0 = [-\infty, 1]$
- (1) $C_1 = [1, 1]$
- (2) $C_2 = C_1 \vee C_4$
- (3) $C_3 = C_2 \cap [-\infty, 100]$
- (4) $C_4 = C_3 + [1, 1]$
- (5) $C_5 = C_2 \cap [101, +\infty]$

Assignment statements are treated using an interval arithmetic (e.g. $[i, j] + [k, l] = [i+k, j+l]$) naturally extended to include the case of the empty interval. Similarly tests are treated using an "interval logic". Since there exist infinite Kleene's sequences (e.g. $[,] < [0, 0] < [0, 1] < \dots < [0, +\infty]$ for the program $x := 0 ; \text{while true do } x := x + 1$), we must use an approximation sequence. Hence the results will be somewhat inaccurate but runtime subscript tests may be inserted in the absence of certainty.

Let us define the widening ∇ of intervals by :

$$\begin{aligned} [-,] &\text{ is the null element of } \nabla \\ [-i, j] \nabla [k, l] &= \begin{cases} \text{if } k < i \text{ then } -\infty \text{ else } i \text{ fi,} \\ \frac{i+k}{2} \text{ if } k > j \text{ then } +\infty \text{ else } j \text{ fi} \end{cases} \end{aligned}$$

∇ satisfies the requirements of 9.1.3. According to 9.1.3.4 the system of equations is modified by :

$$(2) \quad C_2 = C_2 \nabla (C_1 \cup C_4)$$

The corresponding approximation sequence is :

$$\begin{aligned} C_1 &= [-,] \text{ for } i \in [0, 5] \\ * \quad C_1 &= [1, 1] \\ C_2 &= C_2 \nabla (C_1 \cup C_4) \\ &= [-,] \nabla ([1, 1] \cup [-,]) \\ &= [-,] \nabla [-,] \\ C_3 &= C_2 \cap [-\infty, 100] \\ &= [1, 1] \cap [-\infty, 100] \\ &= [1, 1] \\ C_4 &= C_3 + [1, 1] \\ &= [2, 2] \\ C_2 &= C_2 \nabla (C_1 \cup C_4) \\ &= [-,] \nabla ([1, 1] \cup [2, 2]) \\ &= [-,] \nabla [1, 2] \\ * \quad C_2 &= [1, +\infty] \\ C_3 &= C_2 \cap [-\infty, 100] \\ &= [1, +\infty] \cap [-\infty, 100] \end{aligned}$$

* $C_3 = [1, 100]$
 $C_4 = C_3 + [1, 1]$
 $= [1, 100] + [1, 1]$
 $* \quad C_4 = [2, 101]$
 Note $[1, 100] \cap C_4 = [1, 101] \leq C_2 = [1, +\infty]$
 stop on that path.
 $C_5 = C_2 \cap [101, +\infty]$
 $= [1, +\infty] \cap [101, +\infty]$
 $* \quad C_5 = [101, +\infty]$
 exit, stop.

The final context on each arc is marked by a star
 *. Note that the results are approximate ones,
 e.g. C_5 .

In this example the widening is a very rough operation which introduces a great loss of information. However it can be seen in the trace that tests behave like filters. Furthermore, for PASCAL like languages, one can first use the bounds given in the declaration of x before widening to infinite limits.

Early years (1972-73): formal semantics

Patrick Cousot: works on the **operational semantics** of programming languages and the **derivation of implementations from the formal definition**.

Patrick Cousot. *Définition interprétative et implantation de langages de programmation*. Thèse de Docteur Ingénieur en Informatique, Université Joseph Fourier, Grenoble, France, 14 Décembre 1974 (submitted in 1973 but defended after finishing military service with J.D. Ichbiah at CII).

Static analysis of the formal definition and transformation to get the implementation by “pre-evaluation” (similar to the more recent “partial evaluation”)

317

Before starting (1972-73): formal syntax

- Radhia Rezig: works on **precedence parsing** (R.W. Floyd, N. Wirth and H. Weber, etc.) for Algol 68
 - ➡ Pre-processing (by **static analysis and transformation**) of the grammar before building the bottom-up parser
- Patrick Cousot: works on **context-free grammar parsing** (J. Earley and F. De Remer)
 - ➡ Pre-processing (by **static analysis and transformation**) of the grammar before building the top-down parser

-
- Radhia Rezig. Application de la méthode de précédence totale à l'analyse d'Algol 68, Master thesis, Université Joseph Fourier, Grenoble, France, September 1972.
 - Patrick Cousot. Un analyseur syntaxique pour grammaires hors contexte ascendant sélectif et général. In Congrès AFCET 72, Brochure 1, pages 106-130, Grenoble, France, 6-9 November 1972.

1972

Intervals →

pas le niveau de "compréhension" des programmeurs. Les langages actuels ne sont pas faits pour l'optimisation. Entre autres, il y a certains faits sur un programme qui sont connus du programmeur et qui ne sont pas explicités dans le programme. On pourrait y remédier en incluant des assertions, tout comme on insère des déclarations de type pour les variables.

Exemple :

```
(1) - pour i de 0 à 10 faire a[i] := i ; fin ;
(2) - pour i de 11 à 10000 faire a[i] := 0 ; fin ;
(3) - a[(a[j] + 1) * a[j + 1]] := j ;
(4) - si a[j * j + 2 * j + 1] ≠ a[j] aller à étiquette ;
```

Pour un tel programme, il est important de savoir que $1 \leq j < 99$ (à charge éventuellement au système de le déduire à partir d'autres assertions), parce qu'en peut alors remplacer (4) par (4') :
(4') si $j < 10$ aller à étiquette ;

Assertions →

Cette insertion d'assertions peut donc servir de guide à une analyse automatique des programmes essentielle pour l'optimisation (mais également pour la mise au point, la documentation automatique, la décompilation, l'adaptation à un changement d'environnement d'exécution...).

Dans tous les exemples que nous avons pris, (équivalence de définitions de données, équivalence de définition d'opérateurs) nous avons conduit cette analyse sémantique à la main.

La possibilité de son automation, nous semble conditionner les progrès dans le domaine de l'optimisation de l'implantation automatisée d'un langage étant donnée sa définition, aussi bien que dans celui de l'optimisation des programmes [41].

Static analysis →

1974: The origins

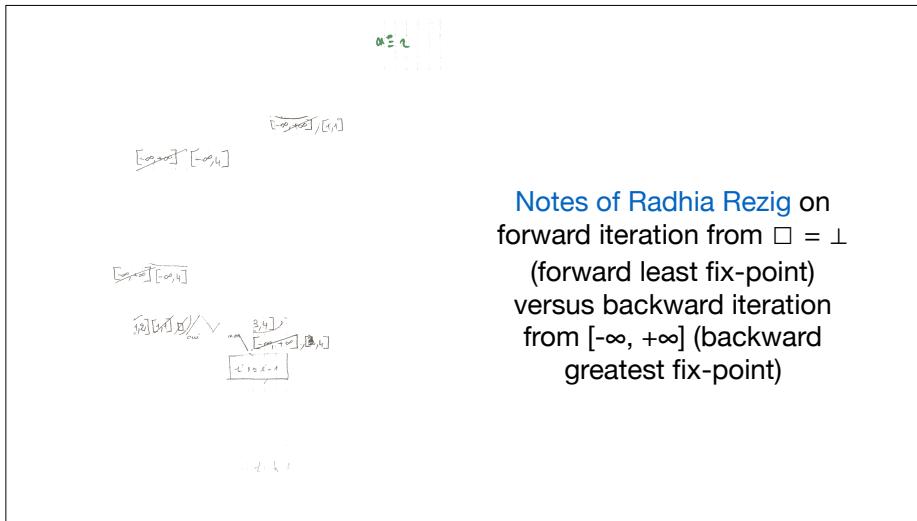
Radhia Rezig shows her **interval analysis** ideas to **Patrick Cousot**

- ➡ Patrick very critical on going backwards from $[-\infty, +\infty]$ and claims that going forward would be much better
- ➡ Patrick also very skeptical on forward termination for loops



Radhia comes back with the idea of extrapolating bounds to $\pm\infty$ for the forward analysis

We discover **widening = induction in the abstract** and that the idea is very general!!

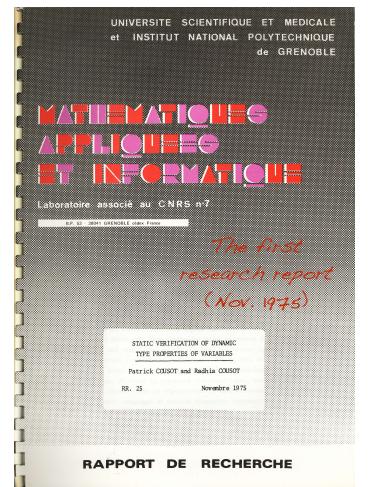


Notes of Radhia Rezig on
forward iteration from $\square = \perp$
(forward least fix-point)
versus backward iteration
from $[-\infty, +\infty]$ (backward
greatest fix-point)

The first abstract
interpreter with
widening
(as of 23 Sep. 1975)

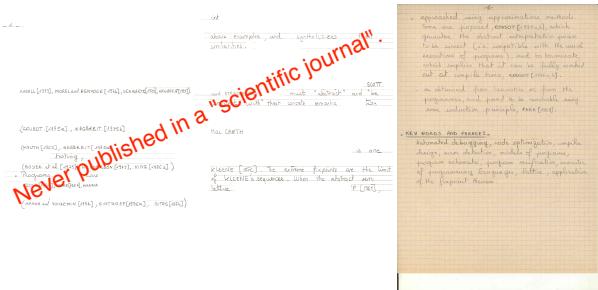
```

131 1) procedure interpretacion_dominio (argph =  $\{A = \{x_1, x_2, x_3, x_4\},$ 
132 1)2)  $B = \{y_1, y_2, y_3, y_4\},$   $C = \{z_1, z_2, z_3, z_4\}$ )
133 2) given  $charac$  as  $f(x)$  means  $f(x) = f(x)$  +  $0 \neq ref$ 
134 3) given  $charac$  as  $f(x) = g(x)$  means  $f(x) = g(x)$  +  $0 \neq ref$ 
135 4) macro  $(charac \ b \ ref) \equiv$  2) else
136 5) various functions  $f(x)$  are  $f(x)$  +  $0 \neq ref$  indications on  $x \in$ 
137 6) variables  $x, y, z$  are  $x = y$  +  $0 \neq ref$ 
138 7) check a domain  $D$  admits  $x$  +  $0 \neq ref$ 
139 8) macro  $(domain \ D) \equiv$   $D = \{x_1, x_2, x_3, x_4\}$ 
140 9) AS
141 10) used macro  $e =$ 
142 11) calculus context  $(arg\_expr, context, found, result)$ 
143 12) used macro  $t =$ 
144 13) calculus context  $(arg\_expr, context, found, result)$ 
145 14) calculus context  $(arg\_expr, result, val, found, result)$ 
146 15) calculus context  $(arg\_expr, result, val, found, result)$ ,  $V_1$ 
147 16) calculus context  $(arg\_expr, result, val, found, result)$ ,  $V_2$ 
148 17) used macro  $(P, Q) =$ 
149 18) used macro  $e =$ 
150 19) used macro  $t =$ 
151 20) your choose how to function  $f(x)$ 
152
153 21) context sort  $\equiv$   $\{1 \leq i \leq n \mid$  predicacion  $\alpha_i(x_i)$  and
154 22)  $\exists_i \neg \alpha_i(x_i)$  sort  $\equiv$   $\{1 \leq i \leq n \mid \alpha_i(x_i)\}$  then
155 23) sort  $\equiv$   $\{1 \leq i \leq n \mid$ 
156 24) calculus context  $(arg\_expr, result)$ 
157 25) context sort  $\equiv$   $\{1 \leq i \leq n \mid$ 
158 26) sort  $\equiv$   $\{1 \leq i \leq n \mid$ 
159 27) calculus context  $(arg\_expr, result)$ 
160 28) context local  $\equiv$   $\{1 \leq i \leq n \mid \alpha_i(x_i)\} \subseteq$  context sort  $\equiv$ 
161 29)  $\{1 \leq i \leq n \mid$ 
162 30) refinable  $\equiv$ 
163 31) givendom context  $\sigma(x_i, cont, contexts)$ 
164 32) dom  $\equiv$ 
165 33) cont  $\equiv$   $\{1 \leq i \leq n \mid$  local  $\alpha_i(x_i)$ 
166 34) context local  $\equiv$   $\{1 \leq i \leq n \mid$ 
167 35) choose  $\equiv$   $\{1 \leq i \leq n \mid$ 
168 36) final  $\equiv$ 
169 37) fin  $\equiv$ 
170 38) fin
```



On submitting to POPL1977

For 4th POPL'77, Patrick and Radhia submitted on August 12, 1976 hard copies of a two-hands written manuscript of 100 pages.



0 INTRODUCTION

A program denotes (complicated)

is to use that denotation

juice

An intuitive example of abstract interpretation, is the rule of signs which we borrow from SINTZOFF [1912] The Text :

55 * 17

KEY WORDS AND PHRASES

of programming language

325

$[-\infty, +\infty]$

$[0, \infty]$

$[-\infty, \underline{c}]$

$[-\infty, -3]$

$[\underline{c}, +\infty]$

$[3, +\infty]$

326

CONCLUDING REMARKS

days

decs
program spec

A Mundane Approach to Semantic Correctness

Semantics:

$$p \vdash v_1 \rightsquigarrow v_2$$

where $v_1, v_2 \in V$.

Note: \rightsquigarrow might be deterministic.

Program analysis:

$$p \vdash l_1 \triangleright l_2$$

where $l_1, l_2 \in L$.

Note: \triangleright should be deterministic:

$$f_p(l_1) = l_2.$$

What is the relationship between the semantics and the analysis?

Restrict attention to analyses where properties directly describe sets of values i.e. "*first-order*" analyses (rather than "*second-order*" analyses).

328

© Nelson^2, Hankin

Example: Data Flow Analysis

Structural Operational Semantics:

Values: $V = \text{State}$

Transitions:

$$S_\star \vdash \sigma_1 \rightsquigarrow \sigma_2 \quad \text{iff} \quad \langle S_\star, \sigma_1 \rangle \rightarrow^* \sigma_2$$

Constant Propagation Analysis:

Properties: $L = \widehat{\text{State}}_{\text{CP}} = (\text{Var}_\star \rightarrow \mathbf{Z}^\top)_\perp$

Transitions:

$$S_\star \vdash \hat{\sigma}_1 \triangleright \hat{\sigma}_2$$

iff

$$\begin{aligned} \hat{\sigma}_1 &= \iota \\ \hat{\sigma}_2 &= \sqcup \{ \text{CP}_\bullet(\ell) \mid \ell \in \text{final}(S_\star) \} \\ (\text{CP}_\circ, \text{CP}_\bullet) &\models \text{CP}^=(S_\star) \end{aligned}$$

329

© Nielsen^2, Hankin

Example: Control Flow Analysis

Structural Operational Semantics:

Values: $V = \text{Val}$

Transitions:

$$e_\star \vdash v_1 \rightsquigarrow v_2$$

iff

$$[] \vdash (e_\star v_1^{\ell_1})^{\ell_2} \rightarrow^* v_2^{\ell_2}$$

$$e_\star \vdash (\hat{p}_1, \hat{v}_1) \triangleright (\hat{p}_2, \hat{v}_2)$$

iff

$$\hat{c}(\ell_1) = \hat{v}_1$$

$$\hat{c}(\ell_2) = \hat{v}_2$$

$$\hat{p}_1 = \hat{p}_2 = \hat{p}$$

$$(\hat{c}, \hat{p}) \models (e_\star c^{\ell_1})^{\ell_2}$$

for some place holder constant c

330

© Nielsen^2, Hankin

Correctness Relations

$$R : V \times L \rightarrow \{\text{true}, \text{false}\}$$

Idea: $v R l$ means that the value v is described by the property l .

Correctness criterion: R is preserved under computation:

$$\begin{array}{llll} p \vdash v_1 & \rightsquigarrow & v_2 \\ \vdots & & \vdots \\ R & \Rightarrow & R \\ \vdots & & \vdots \\ p \vdash l_1 & \triangleright & l_2 \end{array}$$

logical relation:
 $(p \vdash \cdot \rightsquigarrow \cdot) \quad (R \Rightarrow R) \quad (p \vdash \cdot \triangleright \cdot)$

331

© Nielsen^2, Hankin

Admissible Correctness Relations

$$v R l_1 \wedge l_1 \sqsubseteq l_2 \Rightarrow v R l_2$$

$$(\forall l \in L' \subseteq L : v R l) \Rightarrow v R (\sqcap L')$$

Two consequences:

$$v R \top$$

$$v R l_1 \wedge v R l_2 \Rightarrow v R (l_1 \sqcap l_2)$$

Assumption: (L, \sqsubseteq) is a complete lattice.

332

© Nielsen^2, Hankin

Example: Data Flow Analysis

Correctness relation

$$R_{CP} : \text{State} \times \widehat{\text{State}}_{CP} \rightarrow \{\text{true}, \text{false}\}$$

is defined by

$$\sigma R_{CP} \hat{\sigma} \text{ iff } \forall x \in FV(S_*) : (\hat{\sigma}(x) = \top \vee \sigma(x) = \hat{\sigma}(x))$$

333

© Nielsen^2, Hankin

Example: Control Flow Analysis

Correctness relation

$$R_{CFA} : \text{Val} \times (\widehat{\text{Env}} \times \widehat{\text{Val}}) \rightarrow \{\text{true}, \text{false}\}$$

is defined by

$$v R_{CFA} (\hat{\rho}, \hat{v}) \text{ iff } v V(\hat{\rho}, \hat{v})$$

where V is given by:

$$v V(\hat{\rho}, \hat{v}) \text{ iff } \begin{cases} \text{true} & \text{if } v = c \\ t \in \hat{v} \wedge \forall x \in \text{dom}(\rho) : \rho(x) V(\hat{\rho}, \hat{\rho}(x)) & \text{if } v = \text{close } t \text{ in } \rho \end{cases}$$

334

© Nielsen^2, Hankin

Representation Functions

$$\beta : V \rightarrow L$$

Idea: β maps a value to the *best* property describing it.

Correctness criterion:

$$\begin{array}{ccccc} p \vdash v_1 & \xrightarrow{\sim} & v_2 & & \\ \beta \downarrow & \Rightarrow & \downarrow \beta & & \\ \sqcap & & \sqcap & & \\ p \vdash l_1 & \xrightarrow{\triangleright} & l_2 & & \end{array}$$

335

© Nielsen^2, Hankin

Equivalence of Correctness Criteria

Given a representation function β we define a correctness relation R_β by

$$v R_\beta l \text{ iff } \beta(v) \sqsubseteq l$$

Given a correctness relation R we define a representation function β_R by

$$\beta_R(v) = \sqcap \{l \mid v R l\}$$

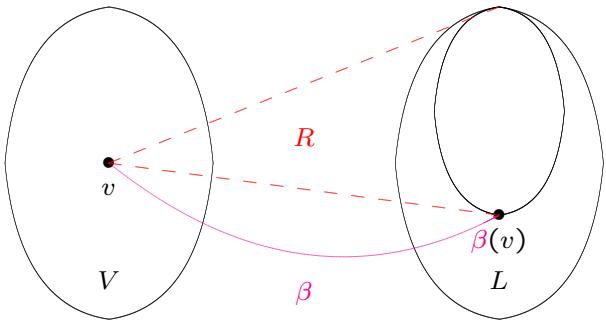
Lemma:

- (i) Given $\beta : V \rightarrow L$, then the relation $R_\beta : V \times L \rightarrow \{\text{true}, \text{false}\}$ is an admissible correctness relation such that $\beta_{R_\beta} = \beta$.
- (ii) Given an admissible correctness relation $R : V \times L \rightarrow \{\text{true}, \text{false}\}$, then β_R is well-defined and $R_{\beta_R} = R$.

336

© Nielsen^2, Hankin

Equivalence of Criteria: R is generated by β



337

© Nielsen^2, Hankin

Example: Data Flow Analysis

Representation function

$$\beta_{\text{CP}} : \text{State} \rightarrow \widehat{\text{State}}_{\text{CP}}$$

is defined by

$$\beta_{\text{CP}}(\sigma) = \lambda x. \sigma(x)$$

R_{CP} is generated by β_{CP} :

$$\sigma \ R_{\text{CP}} \ \widehat{\sigma} \quad \text{iff} \quad \beta_{\text{CP}}(\sigma) \sqsubseteq_{\text{CP}} \widehat{\sigma}$$

© Nielsen^2, Hankin

Example: Control Flow Analysis

Representation function

$$\beta_{\text{CFA}} : \text{Val} \rightarrow \widehat{\text{Env}} \times \widehat{\text{Val}}$$

is defined by

$$\beta_{\text{CFA}}(v) = \begin{cases} (\lambda x. \emptyset, \emptyset) & \text{if } v = c \\ (\beta_{\text{CFA}}^E(\rho), \{t\}) & \text{if } v = \text{close } t \text{ in } \rho \end{cases}$$

$$\begin{aligned} \beta_{\text{CFA}}^E(\rho)(x) &= \bigcup \{ \hat{\rho}_y(x) \mid \beta_{\text{CFA}}(\rho(y)) = (\hat{\rho}_y, \hat{v}_y) \text{ and } y \in \text{dom}(\rho) \} \\ &\cup \begin{cases} \{\hat{v}_x\} & \text{if } x \in \text{dom}(\rho) \text{ and } \beta_{\text{CFA}}(\rho(x)) = (\hat{\rho}_x, \hat{v}_x) \\ \emptyset & \text{otherwise} \end{cases} \end{aligned}$$

R_{CFA} is generated by β_{CFA} :

$$v \ R_{\text{CFA}} \ (\hat{\rho}, \hat{v}) \quad \text{iff} \quad \beta_{\text{CFA}}(v) \sqsubseteq_{\text{CFA}} (\hat{\rho}, \hat{v})$$

339

© Nielsen^2, Hankin

A Modest Generalisation

Semantics:

$$p \vdash v_1 \rightsquigarrow v_2$$

where $v_1 \in V_1, v_2 \in V_2$

Program analysis:

$$p \vdash l_1 \triangleright l_2$$

where $l_1 \in L_1, l_2 \in L_2$

$$\begin{array}{ccccc} p & \vdash & v_1 & \rightsquigarrow & v_2 \\ & & \vdots & & \vdots \\ & & R_1 & \Rightarrow & R_2 \\ & & \vdots & & \vdots \\ p & \vdash & l_1 & \triangleright & l_2 \end{array}$$

logical relation:

$$(p \vdash \cdot \rightsquigarrow \cdot) \ (R_1 \Rightarrow R_2) \ (p \vdash \cdot \triangleright \cdot)$$

340

© Nielsen^2, Hankin

Higher-Order Formulation

Assume that

- R_1 is an admissible correctness relation for V_1 and L_1
that is generated by the representation function $\beta_1 : V_1 \rightarrow L_1$
- R_2 is an admissible correctness relation for V_2 and L_2
that is generated by the representation function $\beta_2 : V_2 \rightarrow L_2$

Then the relation $R_1 \rightarrow R_2$ is an admissible correctness relation for $V_1 \rightarrow V_2$ and $L_1 \rightarrow L_2$

that is generated by the representation function $\beta_1 \rightarrow \beta_2$ defined by

$$(\beta_1 \rightarrow \beta_2)(\sim) = \lambda l_1. \bigsqcup \{ \beta_2(v_2) \mid \beta_1(v_1) \sqsubseteq l_1 \wedge v_1 \sim v_2 \}$$

341

© Nielsen^2, Hankin

Example:

Semantics:

$$\text{plus} \vdash (z_1, z_2) \rightsquigarrow z_1 + z_2$$

where $z_1, z_2 \in \mathbb{Z}$

Program analysis:

$$\text{plus} \vdash ZZ \triangleright \{z_1 + z_2 \mid (z_1, z_2) \in ZZ\}$$

where $ZZ \subseteq \mathbb{Z} \times \mathbb{Z}$

| | Correctness relations | Representation functions |
|----------|---|---|
| result | $R_{\mathbb{Z}}$ | $\beta_{\mathbb{Z}}(z) = \{z\}$ |
| argument | $R_{\mathbb{Z} \times \mathbb{Z}}$ | $\beta_{\mathbb{Z} \times \mathbb{Z}}(z_1, z_2) = \{(z_1, z_2)\}$ |
| plus | $(\text{plus} \vdash \cdot \rightsquigarrow \cdot)$ $(R_{\mathbb{Z} \times \mathbb{Z}} \rightarrow R_{\mathbb{Z}})$ $(\text{plus} \vdash \cdot \triangleright \cdot)$ | $(\beta_{\mathbb{Z} \times \mathbb{Z}} \rightarrow \beta_{\mathbb{Z}})(\text{plus} \vdash \cdot \rightsquigarrow \cdot)$ $\sqsubseteq (\text{plus} \vdash \cdot \triangleright \cdot)$ |

342

© Nielsen^2, Hankin

Approximation of Fixed Points

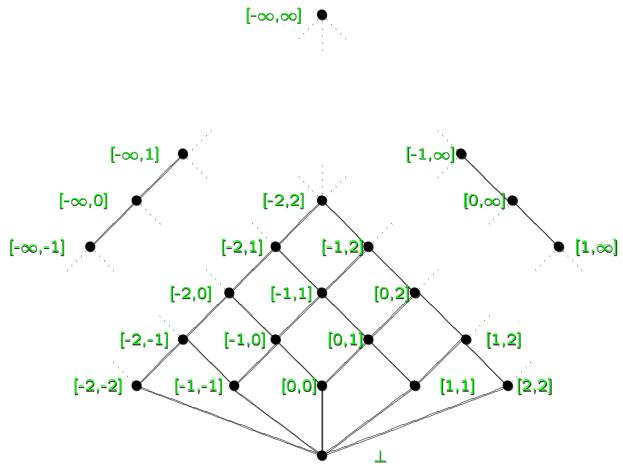
- Fixed points
- Widening
- Narrowing

Example: lattice of intervals for *Array Bound Analysis*

343

© Nielsen^2, Hankin

The complete lattice **Interval** = (**Interval**, \sqsubseteq)



© Nielsen^2, Hankin

Fixed points

Let $f : L \rightarrow L$ be a *monotone function* on a complete lattice $L = (L, \sqsubseteq, \sqcup, \sqcap, \perp, \top)$.

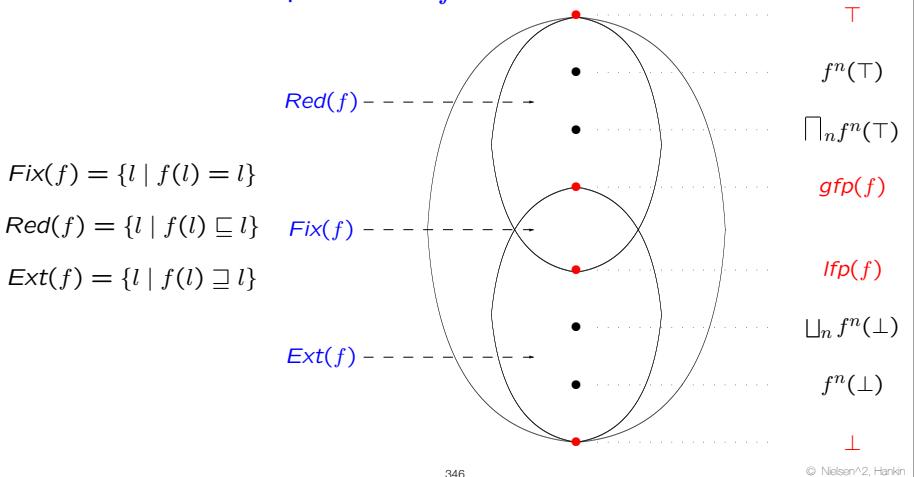
$$\begin{array}{ll} l \text{ is a } \text{fixed point} & \text{iff } f(l) = l \quad \text{Fix}(f) = \{l \mid f(l) = l\} \\ f \text{ is } \text{reductive} \text{ at } l & \text{iff } f(l) \sqsubseteq l \quad \text{Red}(f) = \{l \mid f(l) \sqsubseteq l\} \\ f \text{ is } \text{extensive} \text{ at } l & \text{iff } f(l) \sqsupseteq l \quad \text{Ext}(f) = \{l \mid f(l) \sqsupseteq l\} \end{array}$$

Tarski's Theorem ensures that

$$\begin{aligned} \text{lfp}(f) &= \sqcap \text{Fix}(f) = \sqcap \text{Red}(f) \in \text{Fix}(f) \subseteq \text{Red}(f) \\ \text{gfp}(f) &= \sqcup \text{Fix}(f) = \sqcup \text{Ext}(f) \in \text{Fix}(f) \subseteq \text{Ext}(f) \end{aligned}$$

© Nielsen^2, Hankin

Fixed points of f



Widening Operators

Problem: We cannot guarantee that $(f^n(\perp))_n$ eventually stabilises nor that its least upper bound necessarily equals $\text{lfp}(f)$.

Idea: We replace $(f^n(\perp))_n$ by a new sequence $(f_{\triangleright}^n)_n$ that is known to eventually stabilise and to do so with a value that is a safe (upper) approximation of the least fixed point.

The new sequence is parameterised on the widening operator \triangleright : an upper bound operator satisfying a finiteness condition.

347

© Nielsen^2, Hankin

Upper bound operators

$\triangleright : L \times L \rightarrow L$ is an *upper bound operator* iff

$$l_1 \sqsubseteq l_1 \triangleright l_2 \sqsupseteq l_2$$

for all $l_1, l_2 \in L$.

Let $(l_n)_n$ be a sequence of elements of L . Define the sequence $(l_{\triangleright}^n)_n$ by:

$$l_{\triangleright}^n = \begin{cases} l_n & \text{if } n = 0 \\ l_{\triangleright}^{n-1} \triangleright l_n & \text{if } n > 0 \end{cases}$$

Fact: If $(l_n)_n$ is a sequence and \triangleright is an upper bound operator then $(l_{\triangleright}^n)_n$ is an ascending chain; furthermore $l_{\triangleright}^n \sqsupseteq \sqcup \{l_0, l_1, \dots, l_n\}$ for all n .

348

© Nielsen^2, Hankin

Example:

Let int be an arbitrary but fixed element of **Interval**.

An upper bound operator:

$$\text{int}_1 \sqcup^{\text{int}} \text{int}_2 = \begin{cases} \text{int}_1 \sqcup \text{int}_2 & \text{if } \text{int}_1 \sqsubseteq \text{int} \vee \text{int}_2 \sqsubseteq \text{int}_1 \\ [-\infty, \infty] & \text{otherwise} \end{cases}$$

Example: $[1, 2] \sqcup^{[0,2]} [2, 3] = [1, 3]$ and $[2, 3] \sqcup^{[0,2]} [1, 2] = [-\infty, \infty]$.

Transformation of: $[0, 0], [1, 1], [2, 2], [3, 3], [4, 4], [5, 5], \dots$

If $\text{int} = [0, \infty]$: $[0, 0], [0, 1], [0, 2], [0, 3], [0, 4], [0, 5], \dots$

If $\text{int} = [0, 2]$: $[0, 0], [0, 1], [0, 2], [0, 3], [-\infty, \infty], [-\infty, \infty], \dots$

349

© Nielsen^2, Hankin

Widening operators

An operator $\nabla : L \times L \rightarrow L$ is a *widening operator* iff

- it is an upper bound operator, and
- for all ascending chains $(l_n)_n$ the ascending chain $(l_n^\nabla)_n$ eventually stabilises.

350

© Nielsen^2, Hankin

Widening operators

Given a monotone function $f : L \rightarrow L$ and a widening operator ∇ define the sequence $(f_\nabla^n)_n$ by

$$f_\nabla^n = \begin{cases} \perp & \text{if } n = 0 \\ f_\nabla^{n-1} \nabla f(f_\nabla^{n-1}) & \text{if } n > 0 \wedge f(f_\nabla^{n-1}) \sqsubseteq f_\nabla^{n-1} \\ f_\nabla^{n-1} & \text{otherwise} \end{cases}$$

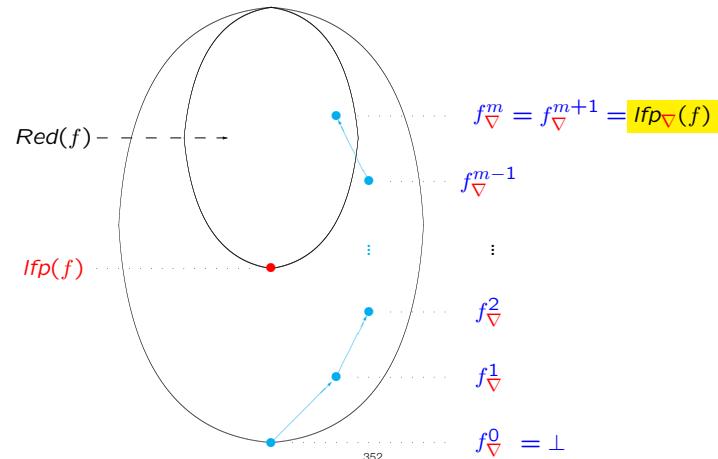
One can show that:

- $(f_\nabla^n)_n$ is an ascending chain that eventually stabilises
- it happens when $f(f_\nabla^m) \sqsubseteq f_\nabla^m$ for some value of m
- Tarski's Theorem then gives $f_\nabla^m \sqsupseteq \text{Ifp}(f)$

351

$$\boxed{\text{Ifp}_\nabla(f) = f_\nabla^m}$$

The widening operator ∇ applied to f



© Nielsen^2, Hankin

Example:

Let K be a *finite* set of integers, e.g. the set of integers explicitly mentioned in a given program.

We shall define a widening operator ∇ based on K .

Idea: $[z_1, z_2] \nabla [z_3, z_4]$ is

$$[LB(z_1, z_3) , UB(z_2, z_4)]$$

where

- $LB(z_1, z_3) \in \{z_1\} \cup K \cup \{-\infty\}$ is the best possible lower bound, and
- $UB(z_2, z_4) \in \{z_2\} \cup K \cup \{\infty\}$ is the best possible upper bound.

The effect: a change in any of the bounds of the interval $[z_1, z_2]$ can only take place finitely many times – corresponding to the cardinality of K

353

© Nielsen^2, Hankin

Example (cont.) — formalisation:

Let $z_i \in \mathbf{Z}' = \mathbf{Z} \cup \{-\infty, \infty\}$ and write:

$$LB_K(z_1, z_3) = \begin{cases} z_1 & \text{if } z_1 \leq z_3 \\ k & \text{if } z_3 < z_1 \wedge k = \max\{k \in K \mid k \leq z_3\} \\ -\infty & \text{if } z_3 < z_1 \wedge \forall k \in K : z_3 < k \end{cases}$$

$$UB_K(z_2, z_4) = \begin{cases} z_2 & \text{if } z_4 \leq z_2 \\ k & \text{if } z_2 < z_4 \wedge k = \min\{k \in K \mid z_4 \leq k\} \\ \infty & \text{if } z_2 < z_4 \wedge \forall k \in K : k < z_4 \end{cases}$$

$$int_1 \nabla int_2 = \begin{cases} \perp & \text{if } int_1 = int_2 = \perp \\ [LB_K(\inf(int_1), \inf(int_2)) , UB_K(\sup(int_1), \sup(int_2))] & \text{otherwise} \end{cases}$$

354

© Nielsen^2, Hankin

Example (cont.):

Consider the ascending chain $(int_n)_n$

$$[0, 1], [0, 2], [0, 3], [0, 4], [0, 5], [0, 6], [0, 7], \dots$$

and assume that $K = \{3, 5\}$.

Then $(int_n^\nabla)_n$ is the chain

$$[0, 1], [0, 3], [0, 3], [0, 5], [0, 5], [0, \infty], [0, \infty], \dots$$

which eventually stabilises.

355

© Nielsen^2, Hankin

Narrowing Operators

Status: Widening gives us an upper approximation $Ifp_\nabla(f)$ of the least fixed point of f .

Observation: $f(Ifp_\nabla(f)) \sqsubseteq Ifp_\nabla(f)$ so the approximation can be improved by considering the iterative sequence $(f^n(Ifp_\nabla(f)))_n$.

It will satisfy $f^n(Ifp_\nabla(f)) \sqsupseteq Ifp(f)$ for all n so we can stop at an arbitrary point.

The notion of **narrowing** is one way of encapsulating a termination criterion for the sequence.

356

© Nielsen^2, Hankin

Narrowing

An operator $\Delta : L \times L \rightarrow L$ is a *narrowing operator* iff

- $l_2 \sqsubseteq l_1 \Rightarrow l_2 \sqsubseteq (l_1 \Delta l_2) \sqsubseteq l_1$ for all $l_1, l_2 \in L$, and
- for all descending chains $(l_n)_n$ the sequence $(l_n^\Delta)_n$ eventually stabilises.

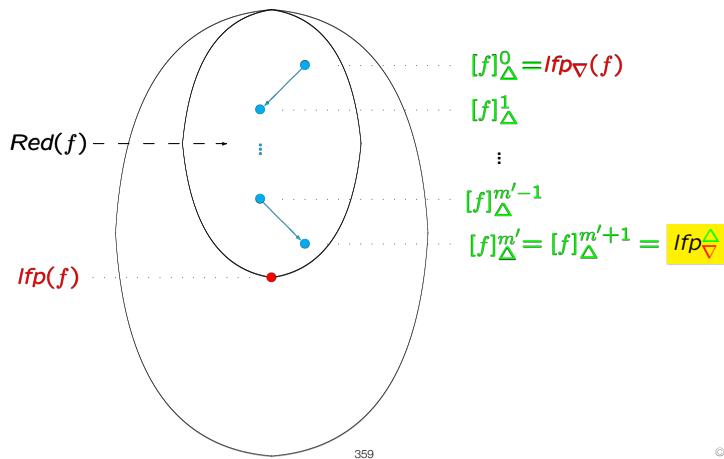
Recall: The sequence $(l_n^\Delta)_n$ is defined by:

$$l_n^\Delta = \begin{cases} l_n & \text{if } n = 0 \\ l_{n-1}^\Delta \Delta l_n & \text{if } n > 0 \end{cases}$$

357

© Nelson^2, Hankin

The narrowing operator Δ applied to f



© Nelson^2, Hankin

Narrowing

We construct the sequence $([f]_\Delta^n)_n$

$$[f]_\Delta^n = \begin{cases} Ifp_Δ(f) & \text{if } n = 0 \\ [f]_\Delta^{n-1} \Delta f([f]_\Delta^{n-1}) & \text{if } n > 0 \end{cases}$$

One can show that:

- $([f]_\Delta^n)_n$ is a descending chain where all elements satisfy $Ifp(f) \sqsubseteq [f]_\Delta^n$
- the chain eventually stabilises so $[f]_\Delta^{m'} = [f]_\Delta^{m'+1}$ for some value m'

$$Ifp_Δ(f) = [f]_\Delta^{m'}$$

358

Example:

The complete lattice (**Interval**, \sqsubseteq) has two kinds of infinite descending chains:

- those with elements of the form $[-\infty, z]$, $z \in \mathbb{Z}$
- those with elements of the form $[z, \infty]$, $z \in \mathbb{Z}$

Idea: Given some fixed non-negative number N
the narrowing operator Δ_N will force an infinite descending chain

$$[z_1, \infty], [z_2, \infty], [z_3, \infty], \dots$$

(where $z_1 < z_2 < z_3 < \dots$) to stabilise when $z_i > N$

Similarly, for a descending chain with elements of the form $[-\infty, z_i]$ the narrowing operator will force it to stabilise when $z_i < -N$

© Nelson^2, Hankin

Example (cont.) — formalisation:

Define $\Delta = \Delta_N$ by

$$int_1 \Delta int_2 = \begin{cases} \perp & \text{if } int_1 = \perp \vee int_2 = \perp \\ [z_1, z_2] & \text{otherwise} \end{cases}$$

where

$$z_1 = \begin{cases} \inf(int_1) & \text{if } N < \inf(int_2) \wedge \sup(int_2) = \infty \\ \inf(int_2) & \text{otherwise} \end{cases}$$

$$z_2 = \begin{cases} \sup(int_1) & \text{if } \inf(int_2) = -\infty \wedge \sup(int_2) < -N \\ \sup(int_2) & \text{otherwise} \end{cases}$$

361

© Nielsen^2, Hankin

Example (cont.):

Consider the infinite descending chain $([n, \infty])_n$

$$[0, \infty], [1, \infty], [2, \infty], [3, \infty], [4, \infty], [5, \infty], \dots$$

and assume that $N = 3$.

Then the narrowing operator Δ_N will give the sequence $([n, \infty]^{\Delta})_n$

$$[0, \infty], [1, \infty], [2, \infty], [3, \infty], [3, \infty], [3, \infty], \dots$$

362

© Nielsen^2, Hankin

Galois Connections

- Galois connections and adjunctions
- Extraction functions
- Galois insertions
- Reduction operators

363

© Nielsen^2, Hankin

Galois connections

$$\begin{array}{ccc} L & \xrightleftharpoons[\alpha]{\gamma} & M \end{array}$$

α : abstraction function
 γ : concretisation function

is a Galois connection if and only if

α and γ are monotone functions

that satisfy

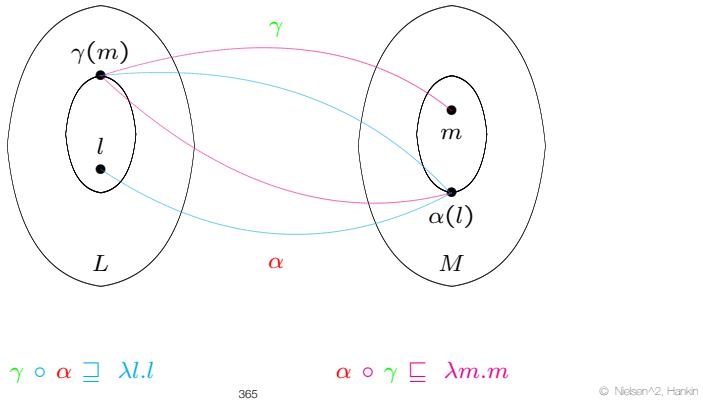
$$\gamma \circ \alpha \sqsupseteq \lambda l.l$$

$$\alpha \circ \gamma \sqsubseteq \lambda m.m$$

364

© Nielsen^2, Hankin

Galois connections



Example:

Galois connection

$(\mathcal{P}(\mathbf{Z}), \alpha_{\mathbf{ZI}}, \gamma_{\mathbf{ZI}}, \mathbf{Interval})$

with concretisation function

$$\gamma_{\mathbf{ZI}}(int) = \{z \in \mathbf{Z} \mid \inf(int) \leq z \leq \sup(int)\}$$

and abstraction function

$$\alpha_{\mathbf{ZI}}(Z) = \begin{cases} \perp & \text{if } Z = \emptyset \\ [\inf'(Z), \sup'(Z)] & \text{otherwise} \end{cases}$$

Examples:

$$\gamma_{\mathbf{ZI}}([0, 3]) = \{0, 1, 2, 3\}$$

$$\gamma_{\mathbf{ZI}}([0, \infty)) = \{z \in \mathbf{Z} \mid z \geq 0\}$$

$$\alpha_{\mathbf{ZI}}(\{0, 1, 3\}) = [0, 3]$$

$$\alpha_{\mathbf{ZI}}(\{2 * z \mid z > 0\}) = [2, \infty]$$

366

© Nielsen^2, Hankin

Galois connections from representation functions

A representation function $\beta : V \rightarrow L$ gives rise to a Galois connection

$$(\mathcal{P}(V), \alpha, \gamma, L)$$

where

$$\alpha(V') = \bigcup \{\beta(v) \mid v \in V'\}$$

$$\gamma(l) = \{v \in V \mid \beta(v) \sqsubseteq l\}$$

for $V' \subseteq V$ and $l \in L$.

367

© Nielsen^2, Hankin

Galois connections from extraction functions

An *extraction function*

$$\eta : V \rightarrow D$$

maps the values of V to their best descriptions in D .

It gives rise to a representation function $\beta_\eta : V \rightarrow \mathcal{P}(D)$ (corresponding to $L = (\mathcal{P}(D), \subseteq)$) defined by

$$\beta_\eta(v) = \{\eta(v)\}$$

The associated Galois connection is

$$(\mathcal{P}(V), \alpha_\eta, \gamma_\eta, \mathcal{P}(D))$$

where

$$\alpha_\eta(V') = \bigcup \{\beta_\eta(v) \mid v \in V'\} = \{\eta(v) \mid v \in V'\}$$

$$\gamma_\eta(D') = \{v \in V \mid \beta_\eta(v) \subseteq D'\} = \{v \mid \eta(v) \in D'\}$$

368

© Nielsen^2, Hankin

Example:

Extraction function

$$\text{sign} : \mathbb{Z} \rightarrow \text{Sign}$$

specified by

$$\text{sign}(z) = \begin{cases} - & \text{if } z < 0 \\ 0 & \text{if } z = 0 \\ + & \text{if } z > 0 \end{cases}$$

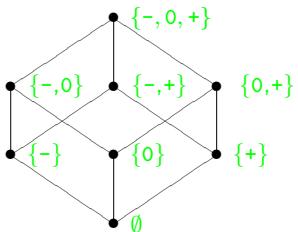
Galois connection

$$(\mathcal{P}(\mathbb{Z}), \alpha_{\text{sign}}, \gamma_{\text{sign}}, \mathcal{P}(\text{Sign}))$$

with

$$\alpha_{\text{sign}}(Z) = \{\text{sign}(z) \mid z \in Z\}$$

$$\gamma_{\text{sign}}(S) = \{z \in \mathbb{Z} \mid \text{sign}(z) \in S\}$$



© Nielsen^2, Hankin

Properties of Galois Connections

Lemma: If (L, α, γ, M) is a Galois connection then:

- α uniquely determines γ by $\gamma(m) = \sqcup\{l \mid \alpha(l) \sqsubseteq m\}$
- γ uniquely determines α by $\alpha(l) = \sqcap\{m \mid l \sqsubseteq \gamma(m)\}$
- α is completely additive and γ is completely multiplicative

In particular $\alpha(\perp) = \perp$ and $\gamma(\top) = \top$.

Lemma:

- If $\alpha : L \rightarrow M$ is completely additive then there exists (an upper adjoint) $\gamma : M \rightarrow L$ such that (L, α, γ, M) is a Galois connection.
- If $\gamma : M \rightarrow L$ is completely multiplicative then there exists (a lower adjoint) $\alpha : L \rightarrow M$ such that (L, α, γ, M) is a Galois connection.

Fact: If (L, α, γ, M) is a Galois connection then

- $\alpha \circ \gamma \circ \alpha = \alpha$ and $\gamma \circ \alpha \circ \gamma = \gamma$

370

© Nielsen^2, Hankin

Example:

Define $\gamma_{\text{IS}} : \mathcal{P}(\text{Sign}) \rightarrow \text{Interval}$ by:

$$\begin{array}{ll} \gamma_{\text{IS}}(\{-, 0, +\}) = [-\infty, \infty] & \gamma_{\text{IS}}(\{-, 0\}) = [-\infty, 0] \\ \gamma_{\text{IS}}(\{-, +\}) = [-\infty, \infty] & \gamma_{\text{IS}}(\{0, +\}) = [0, \infty] \\ \gamma_{\text{IS}}(\{-\}) = [-\infty, -1] & \gamma_{\text{IS}}(\{0\}) = [0, 0] \\ \gamma_{\text{IS}}(\{+\}) = [1, \infty] & \gamma_{\text{IS}}(\emptyset) = \perp \end{array}$$

Does there exist an abstraction function

$$\alpha_{\text{IS}} : \text{Interval} \rightarrow \mathcal{P}(\text{Sign})$$

such that $(\text{Interval}, \alpha_{\text{IS}}, \gamma_{\text{IS}}, \mathcal{P}(\text{Sign}))$ is a Galois connection?

371

© Nielsen^2, Hankin

Example (cont.):

Is γ_{IS} completely multiplicative?

- if yes: then there exists a Galois connection
- if no: then there cannot exist a Galois connection

Lemma: If L and M are complete lattices and M is finite then $\gamma : M \rightarrow L$ is completely multiplicative if and only if the following hold:

- $\gamma : M \rightarrow L$ is monotone,
- $\gamma(\top) = \top$, and
- $\gamma(m_1 \sqcap m_2) = \gamma(m_1) \sqcap \gamma(m_2)$ whenever $m_1 \not\sqsubseteq m_2 \wedge m_2 \not\sqsubseteq m_1$

We calculate

$$\begin{aligned} \gamma_{\text{IS}}(\{-, 0\} \cap \{-, +\}) &= \gamma_{\text{IS}}(\{-\}) = [-\infty, -1] \\ \gamma_{\text{IS}}(\{-, 0\}) \sqcap \gamma_{\text{IS}}(\{-, +\}) &= [-\infty, 0] \sqcap [-\infty, \infty] = [-\infty, 0] \end{aligned}$$

showing that there is no Galois connection involving γ_{IS} .

© Nielsen^2, Hankin

Galois Connections are the Right Concept

We use the mundane approach to correctness to demonstrate this for:

- Admissible correctness relations
- Representation functions

373

© Nielsen^2, Hankin

The mundane approach: correctness relations

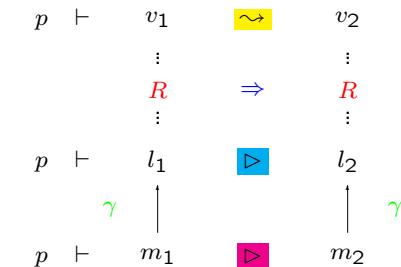
Assume

- $R : V \times L \rightarrow \{\text{true, false}\}$ is an admissible correctness relation
- (L, α, γ, M) is a Galois connection

Then $S : V \times M \rightarrow \{\text{true, false}\}$ defined by

$$v S m \quad \text{iff} \quad v R (\gamma(m))$$

is an admissible correctness relation between V and M



374

© Nielsen^2, Hankin

The mundane approach: representation functions

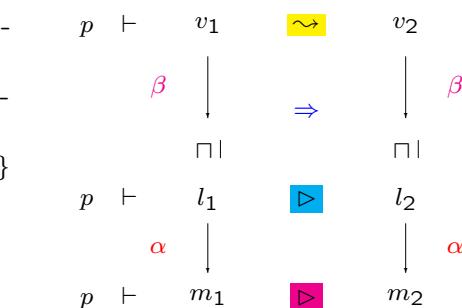
Assume

- $R : V \times L \rightarrow \{\text{true, false}\}$ is generated by $\beta : V \rightarrow L$
- (L, α, γ, M) is a Galois connection

Then $S : V \times M \rightarrow \{\text{true, false}\}$ defined by

$$v S m \quad \text{iff} \quad v R (\gamma(m))$$

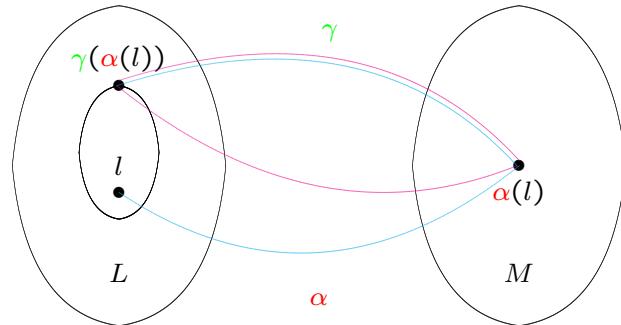
is generated by $\alpha \circ \beta : V \rightarrow M$



375

© Nielsen^2, Hankin

Galois Insertions



Monotone functions satisfying:

$$\gamma \circ \alpha \sqsubseteq \lambda l.l \quad \alpha \circ \gamma = \lambda m.m$$

376

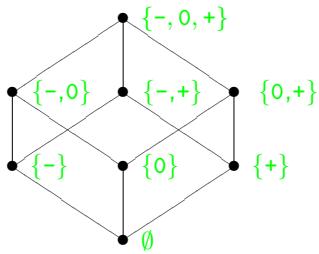
© Nielsen^2, Hankin

Example (1):

$(\mathcal{P}(\mathbb{Z}), \alpha_{\text{sign}}, \gamma_{\text{sign}}, \mathcal{P}(\text{Sign}))$

where $\text{sign} : \mathbb{Z} \rightarrow \text{Sign}$ is specified by:

$$\text{sign}(z) = \begin{cases} - & \text{if } z < 0 \\ 0 & \text{if } z = 0 \\ + & \text{if } z > 0 \end{cases}$$



Is it a Galois insertion?

377

© Nielsen^2, Hankin

Example (2):

$(\mathcal{P}(\mathbb{Z}), \alpha_{\text{signparity}}, \gamma_{\text{signparity}}, \mathcal{P}(\text{Sign} \times \text{Parity}))$

where $\text{Sign} = \{-, 0, +\}$ and $\text{Parity} = \{\text{odd, even}\}$

and $\text{signparity} : \mathbb{Z} \rightarrow \text{Sign} \times \text{Parity}$:

$$\text{signparity}(z) = \begin{cases} (\text{sign}(z), \text{odd}) & \text{if } z \text{ is odd} \\ (\text{sign}(z), \text{even}) & \text{if } z \text{ is even} \end{cases}$$

Is it a Galois insertion?

378

© Nielsen^2, Hankin

Example (1) reconsidered:

$(\mathcal{P}(\mathbb{Z}), \alpha_{\text{sign}}, \gamma_{\text{sign}}, \mathcal{P}(\text{Sign}))$

$$\text{sign}(z) = \begin{cases} - & \text{if } z < 0 \\ 0 & \text{if } z = 0 \\ + & \text{if } z > 0 \end{cases}$$

is a Galois insertion because sign is surjective.

Example (2) reconsidered:

$(\mathcal{P}(\mathbb{Z}), \alpha_{\text{signparity}}, \gamma_{\text{signparity}}, \mathcal{P}(\text{Sign} \times \text{Parity}))$

$$\text{signparity}(z) = \begin{cases} (\text{sign}(z), \text{odd}) & \text{if } z \text{ is odd} \\ (\text{sign}(z), \text{even}) & \text{if } z \text{ is even} \end{cases}$$

is not a Galois insertion because signparity is not surjective.

379

© Nielsen^2, Hankin

Properties of Galois Insertions

Lemma: For a Galois connection (L, α, γ, M) the following claims are equivalent:

- (i) (L, α, γ, M) is a Galois insertion;
- (ii) α is surjective: $\forall m \in M : \exists l \in L : \alpha(l) = m$;
- (iii) γ is injective: $\forall m_1, m_2 \in M : \gamma(m_1) = \gamma(m_2) \Rightarrow m_1 = m_2$; and
- (iv) γ is an order-similarity: $\forall m_1, m_2 \in M : \gamma(m_1) \sqsubseteq \gamma(m_2) \Leftrightarrow m_1 \sqsubseteq m_2$.

Corollary: A Galois connection specified by an extraction function $\eta : V \rightarrow D$ is a Galois insertion if and only if η is surjective.

380

© Nielsen^2, Hankin

Reduction Operators

Given a Galois connection (L, α, γ, M) it is **always** possible to obtain a Galois insertion by enforcing that the concretisation function γ is injective.

Idea: remove the superfluous elements from M using a *reduction operator*

$$\varsigma : M \rightarrow M$$

defined from the Galois connection.

Proposition: Let (L, α, γ, M) be a Galois connection and define the reduction operator $\varsigma : M \rightarrow M$ by

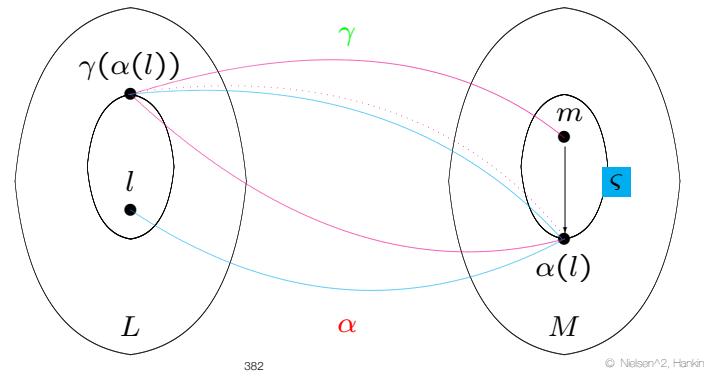
$$\varsigma(m) = \sqcap \{m' \mid \gamma(m) = \gamma(m')\}$$

Then $\varsigma[M] = (\{\varsigma(m) \mid m \in M\}, \sqsubseteq_M)$ is a complete lattice and $(L, \alpha, \gamma, \varsigma[M])$ is a Galois insertion.

381

© Nielsen^2, Hankin

The reduction operator $\varsigma : M \rightarrow M$



382

© Nielsen^2, Hankin

Reduction operators from extraction functions

Assume that the Galois connection $(\mathcal{P}(V), \alpha_\eta, \gamma_\eta, \mathcal{P}(D))$ is given by an extraction function $\eta : V \rightarrow D$.

Then the reduction operator ς_η is given by

$$\varsigma_\eta(D') = D' \cap \eta[V]$$

where $\eta[V] = \{d \in D \mid \exists v \in V : \eta(v) = d\}$.

Since $\varsigma_\eta[\mathcal{P}(D)]$ is isomorphic to $\mathcal{P}(\eta[V])$ the resulting Galois insertion is isomorphic to

$$(\mathcal{P}(V), \alpha_\eta, \gamma_\eta, \mathcal{P}(\eta[V]))$$

383

© Nielsen^2, Hankin

Systematic Design of Galois Connections

The “functional composition” (or “sequential composition”) of two Galois connections is also a Galois connection:

$$L_0 \xrightarrow[\alpha_1]{\gamma_1} L_1 \xrightarrow[\alpha_2]{\gamma_2} L_2 \xrightarrow[\alpha_3]{\gamma_3} \dots \xrightarrow[\alpha_k]{\gamma_k} L_k$$

A catalogue of techniques for combining Galois connections:

- independent attribute method
- relational method
- direct product
- direct tensor product
- reduced product
- reduced tensor product
- total function space
- monotone function space

384

© Nielsen^2, Hankin

Running Example: Array Bound Analysis

Approximation of the difference in magnitude between two numbers (typically the index and the bound):

- a Galois connection for approximating pairs (z_1, z_2) of integers by their difference $|z_1| - |z_2|$
- a Galois connection for approximating integers using a finite lattice $\{<-1, -1, 0, +1, >+1\}$
- a Galois connection for their functional composition

385

© Nielsen^2, Hankin

Example: Difference in Magnitude

$$(\mathcal{P}(\mathbf{Z} \times \mathbf{Z}), \alpha_{\text{diff}}, \gamma_{\text{diff}}, \mathcal{P}(\mathbf{Z}))$$

where the extraction function $\text{diff} : \mathbf{Z} \times \mathbf{Z} \rightarrow \mathbf{Z}$ calculates the difference in magnitude:

$$\text{diff}(z_1, z_2) = |z_1| - |z_2|$$

The abstraction and concretisation functions are

$$\begin{aligned}\alpha_{\text{diff}}(ZZ) &= \{|z_1| - |z_2| \mid (z_1, z_2) \in ZZ\} \\ \gamma_{\text{diff}}(Z) &= \{(z_1, z_2) \mid |z_1| - |z_2| \in Z\}\end{aligned}$$

for $ZZ \subseteq \mathbf{Z} \times \mathbf{Z}$ and $Z \subseteq \mathbf{Z}$.

386

© Nielsen^2, Hankin

Example: Finite Approximation

$$(\mathcal{P}(\mathbf{Z}), \alpha_{\text{range}}, \gamma_{\text{range}}, \mathcal{P}(\mathbf{Range}))$$

where $\mathbf{Range} = \{<-1, -1, 0, +1, >+1\}$

and the extraction function $\text{range} : \mathbf{Z} \rightarrow \mathbf{Range}$ is

$$\text{range}(z) = \begin{cases} <-1 & \text{if } z < -1 \\ -1 & \text{if } z = -1 \\ 0 & \text{if } z = 0 \\ +1 & \text{if } z = 1 \\ >+1 & \text{if } z > 1 \end{cases}$$

The abstraction and concretisation functions are

$$\begin{aligned}\alpha_{\text{range}}(Z) &= \{\text{range}(z) \mid z \in Z\} \\ \gamma_{\text{range}}(R) &= \{z \mid \text{range}(z) \in R\}\end{aligned}$$

for $Z \subseteq \mathbf{Z}$ and $R \subseteq \mathbf{Range}$.

387

© Nielsen^2, Hankin

Example: Functional Composition

$$(\mathcal{P}(\mathbf{Z} \times \mathbf{Z}), \alpha_{\mathbf{R}}, \gamma_{\mathbf{R}}, \mathcal{P}(\mathbf{Range}))$$

where

$$\begin{aligned}\alpha_{\mathbf{R}} &= \alpha_{\text{range}} \circ \alpha_{\text{diff}} \\ \gamma_{\mathbf{R}} &= \gamma_{\text{diff}} \circ \gamma_{\text{range}}\end{aligned}$$

The explicit formulae for the abstraction and concretisation functions

$$\begin{aligned}\alpha_{\mathbf{R}}(ZZ) &= \{\text{range}(|z_1| - |z_2|) \mid (z_1, z_2) \in ZZ\} \\ \gamma_{\mathbf{R}}(R) &= \{(z_1, z_2) \mid \text{range}(|z_1| - |z_2|) \in R\}\end{aligned}$$

correspond to the extraction function $\text{range} \circ \text{diff}$.

388

© Nielsen^2, Hankin

Approximation of Pairs

Independent Attribute Method

Let $(L_1, \alpha_1, \gamma_1, M_1)$ and $(L_2, \alpha_2, \gamma_2, M_2)$ be Galois connections.

The *independent attribute method* gives a Galois connection

$$(L_1 \times L_2, \alpha, \gamma, M_1 \times M_2)$$

where

$$\alpha(l_1, l_2) = (\alpha_1(l_1), \alpha_2(l_2))$$

$$\gamma(m_1, m_2) = (\gamma_1(m_1), \gamma_2(m_2))$$

389

© Nielsen^2, Hankin

Example: Detection of Signs Analysis

Given

$$(\mathcal{P}(\mathbf{Z}), \alpha_{\text{sign}}, \gamma_{\text{sign}}, \mathcal{P}(\text{Sign}))$$

using the extraction function `sign`.

The independent attribute method gives

$$(\mathcal{P}(\mathbf{Z}) \times \mathcal{P}(\mathbf{Z}), \alpha_{\text{ss}}, \gamma_{\text{ss}}, \mathcal{P}(\text{Sign}) \times \mathcal{P}(\text{Sign}))$$

where

$$\alpha_{\text{ss}}(Z_1, Z_2) = (\{\text{sign}(z) \mid z \in Z_1\}, \{\text{sign}(z) \mid z \in Z_2\})$$

$$\gamma_{\text{ss}}(S_1, S_2) = (\{z \mid \text{sign}(z) \in S_1\}, \{z \mid \text{sign}(z) \in S_2\})$$

390

© Nielsen^2, Hankin

Motivating the Relational Method

The independent attribute method often leads to imprecision!

Semantics: The expression $(x, -x)$ may have a value in

$$\{(z, -z) \mid z \in \mathbf{Z}\}$$

Analysis: When we use $\mathcal{P}(\mathbf{Z}) \times \mathcal{P}(\mathbf{Z})$ to represent sets of pairs of integers we cannot do better than representing $\{(z, -z) \mid z \in \mathbf{Z}\}$ by

$$(\mathbf{Z}, \mathbf{Z})$$

Hence the best property describing it will be

$$\alpha_{\text{ss}}(\mathbf{Z}, \mathbf{Z}) = (\{-, 0, +\}, \{-, 0, +\})$$

391

© Nielsen^2, Hankin

Relational Method

Let $(\mathcal{P}(V_1), \alpha_1, \gamma_1, \mathcal{P}(D_1))$ and $(\mathcal{P}(V_2), \alpha_2, \gamma_2, \mathcal{P}(D_2))$ be Galois connections.

The *relational method* will give rise to the Galois connection

$$(\mathcal{P}(V_1 \times V_2), \alpha, \gamma, \mathcal{P}(D_1 \times D_2))$$

where

$$\alpha(VV) = \bigcup \{\alpha_1(\{v_1\}) \times \alpha_2(\{v_2\}) \mid (v_1, v_2) \in VV\}$$

$$\gamma(DD) = \{(v_1, v_2) \mid \alpha_1(\{v_1\}) \times \alpha_2(\{v_2\}) \subseteq DD\}$$

Generalisation to arbitrary complete lattices: use *tensor products*.

392

© Nielsen^2, Hankin

Relational Method from Extraction Functions

Assume that the Galois connections $(\mathcal{P}(V_i), \alpha_i, \gamma_i, \mathcal{P}(D_i))$ are given by extraction functions $\eta_i : V_i \rightarrow D_i$ as in

$$\alpha_i(V'_i) = \{\eta_i(v_i) \mid v_i \in V'_i\}$$

$$\gamma_i(D'_i) = \{v_i \mid \eta_i(v_i) \in D'_i\}$$

Then the Galois connection $(\mathcal{P}(V_1 \times V_2), \alpha, \gamma, \mathcal{P}(D_1 \times D_2))$ has

$$\alpha(VV) = \{(\eta_1(v_1), \eta_2(v_2)) \mid (v_1, v_2) \in VV\}$$

$$\gamma(DD) = \{(v_1, v_2) \mid (\eta_1(v_1), \eta_2(v_2)) \in DD\}$$

which also can be obtained directly from the extraction function

$\eta : V_1 \times V_2 \rightarrow D_1 \times D_2$ defined by

$$\eta(v_1, v_2) = (\eta_1(v_1), \eta_2(v_2))$$

393

© Nielsen^2, Hankin

Example: Detection of Signs Analysis

Using the relational method we get a Galois connection

$$(\mathcal{P}(\mathbf{Z} \times \mathbf{Z}), \alpha_{SS'}, \gamma_{SS'}, \mathcal{P}(\mathbf{Sign} \times \mathbf{Sign}))$$

where

$$\alpha_{SS'}(ZZ) = \{(\text{sign}(z_1), \text{sign}(z_2)) \mid (z_1, z_2) \in ZZ\}$$

$$\gamma_{SS'}(SS) = \{(z_1, z_2) \mid (\text{sign}(z_1), \text{sign}(z_2)) \in SS\}$$

corresponding to an extraction function $\text{twosigns} : \mathbf{Z} \times \mathbf{Z} \rightarrow \mathbf{Sign} \times \mathbf{Sign}$ defined by

$$\text{twosigns}(z_1, z_2) = (\text{sign}(z_1), \text{sign}(z_2))$$

394

© Nielsen^2, Hankin

Advantages of the Relational Method

Semantics: The expression $(x, -x)$ may have a value in

$$\{(z, -z) \mid z \in \mathbf{Z}\}$$

In the present setting $\{(z, -z) \mid z \in \mathbf{Z}\}$ is an element of $\mathcal{P}(\mathbf{Z} \times \mathbf{Z})$.

Analysis: The best “relational” property describing it is

$$\alpha_{SS'}(\{(z, -z) \mid z \in \mathbf{Z}\}) = \{(-,+), (0,0), (+,-)\}$$

whereas the best “independent attribute” property was

$$\alpha_{SS}(\mathbf{Z}, \mathbf{Z}) = \{(-, 0, +), (-, 0, +)\}$$

395

© Nielsen^2, Hankin

Function Spaces

Total Function Space

Let (L, α, γ, M) be a Galois connection and let S be a set.

The Galois connection for the total function space

$$(S \rightarrow L, \alpha', \gamma', S \rightarrow M)$$

is defined by

$$\alpha'(f) = \alpha \circ f$$

$$\gamma'(g) = \gamma \circ g$$

396

© Nielsen^2, Hankin

Monotone Function Space

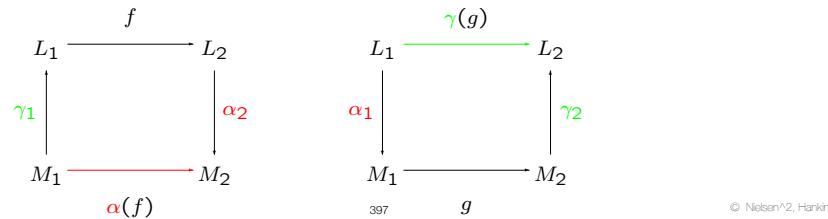
Let $(L_1, \alpha_1, \gamma_1, M_1)$ and $(L_2, \alpha_2, \gamma_2, M_2)$ be Galois connections.

The Galois connection for the *monotone function space*

$$(L_1 \rightarrow L_2, \alpha, \gamma, M_1 \rightarrow M_2)$$

is defined by

$$\alpha(f) = \alpha_2 \circ f \circ \gamma_1 \quad \gamma(g) = \gamma_2 \circ g \circ \alpha_1$$



Performing Analyses Simultaneously

Direct Product

Let $(L, \alpha_1, \gamma_1, M_1)$ and $(L, \alpha_2, \gamma_2, M_2)$ be Galois connections.

The *direct product* is the Galois connection

$$(L, \alpha, \gamma, M_1 \times M_2)$$

defined by

$$\alpha(l) = (\alpha_1(l), \alpha_2(l))$$

$$\gamma(m_1, m_2) = \gamma_1(m_1) \sqcap \gamma_2(m_2)$$

398

© Nielsen^2, Hankin

Example:

Combining the *detection of signs analysis* for pairs of integers with the *analysis of difference in magnitude*.

We get the Galois connection

$$(\mathcal{P}(\mathbf{Z} \times \mathbf{Z}), \alpha_{\text{SSR}}, \gamma_{\text{SSR}}, \mathcal{P}(\text{Sign} \times \text{Sign}) \times \mathcal{P}(\text{Range}))$$

where

$$\begin{aligned} \alpha_{\text{SSR}}(ZZ) &= (\{\text{sign}(z_1), \text{sign}(z_2)\} \mid (z_1, z_2) \in ZZ, \\ &\quad \{\text{range}(|z_1| - |z_2|) \mid (z_1, z_2) \in ZZ\}) \end{aligned}$$

$$\begin{aligned} \gamma_{\text{SSR}}(SS, R) &= \{(z_1, z_2) \mid (\text{sign}(z_1), \text{sign}(z_2)) \in SS\} \\ &\quad \sqcap \{(z_1, z_2) \mid \text{range}(|z_1| - |z_2|) \in R\} \end{aligned}$$

399

© Nielsen^2, Hankin

Motivating the Direct Tensor Product

The expression $(x, 3*x)$ may have a value in

$$\{(z, 3*z) \mid z \in \mathbf{Z}\}$$

which is described by

$$\alpha_{\text{SSR}}(\{(z, 3*z) \mid z \in \mathbf{Z}\}) = (\{(-,-), (0,0), (+,+)\}, \{0, <-1\})$$

But

- any pair described by $(0,0)$ will have a difference in magnitude described by 0
- any pair described by $(-, -)$ or $(+, +)$ will have a difference in magnitude described by <-1

and the analysis cannot express this.

400

© Nielsen^2, Hankin

Direct Tensor Product

Let $(\mathcal{P}(V), \alpha_1, \gamma_1, \mathcal{P}(D_1))$ and $(\mathcal{P}(V), \alpha_2, \gamma_2, \mathcal{P}(D_2))$ be Galois connections.

The *direct tensor product* is the Galois connection

$$(\mathcal{P}(V), \alpha, \gamma, \mathcal{P}(D_1 \times D_2))$$

defined by

$$\alpha(V') = \bigcup \{\alpha_1(\{v\}) \times \alpha_2(\{v\}) \mid v \in V'\}$$

$$\gamma(DD) = \{v \mid \alpha_1(\{v\}) \times \alpha_2(\{v\}) \subseteq DD\}$$

401

© Nielsen^2, Hankin

Direct Tensor Product from Extraction Functions

Assume that the Galois connections $(\mathcal{P}(V), \alpha_i, \gamma_i, \mathcal{P}(D_i))$ are given by *extraction functions* $\eta_i : V \rightarrow D_i$ as in

$$\alpha_i(V') = \{\eta_i(v) \mid v \in V'\}$$

$$\gamma_i(D'_i) = \{v \mid \eta_i(v) \in D'_i\}$$

The Galois connection $(\mathcal{P}(V), \alpha, \gamma, \mathcal{P}(D_1 \times D_2))$ has

$$\alpha(V') = \{(\eta_1(v), \eta_2(v)) \mid v \in V'\}$$

$$\gamma(DD) = \{v \mid (\eta_1(v), \eta_2(v)) \in DD\}$$

corresponding to the extraction function $\eta : V \rightarrow D_1 \times D_2$ defined by

$$\eta(v) = (\eta_1(v), \eta_2(v))$$

402

© Nielsen^2, Hankin

Example:

Using the direct tensor product to combine the *detection of signs analysis* for pairs of integers with the *analysis of difference in magnitude*.

$$(\mathcal{P}(\mathbb{Z} \times \mathbb{Z}), \alpha_{SSR'}, \gamma_{SSR'}, \mathcal{P}(\text{Sign} \times \text{Sign} \times \text{Range}))$$

is given by

$$\alpha_{SSR'}(ZZ) = \{(\text{sign}(z_1), \text{sign}(z_2), \text{range}(|z_1| - |z_2|)) \mid (z_1, z_2) \in ZZ\}$$

$$\gamma_{SSR'}(SSR) = \{(z_1, z_2) \mid (\text{sign}(z_1), \text{sign}(z_2), \text{range}(|z_1| - |z_2|)) \in SSR\}$$

corresponding to *twosignsrangle* : $\mathbb{Z} \times \mathbb{Z} \rightarrow \text{Sign} \times \text{Sign} \times \text{Range}$ given by

$$\text{twosignsrangle}(z_1, z_2) = (\text{sign}(z_1), \text{sign}(z_2), \text{range}(|z_1| - |z_2|))$$

403

© Nielsen^2, Hankin

Advantages of the Direct Tensor Product

The expression $(x, 3*x)$ may have a value in $\{(z, 3*z) \mid z \in \mathbb{Z}\}$ which in the direct tensor product can be described by

$$\alpha_{SSR'}(\{(z, 3*z) \mid z \in \mathbb{Z}\}) = \{(-, -, <-1), (0, 0, 0), (+, +, <-1)\}$$

compared to the direct product that gave

$$\alpha_{SSR}(\{(z, 3*z) \mid z \in \mathbb{Z}\}) = \{((-,-),(0,0),(+,+)), \{0,<-1\}\}$$

Note that the Galois connection is *not* a Galois insertion because

$$\gamma_{SSR'}(\emptyset) = \emptyset = \gamma_{SSR'}(\{(0, 0, <-1)\})$$

so $\gamma_{SSR'}$ is not injective and hence we do not have a Galois insertion.

404

© Nielsen^2, Hankin

From Direct to Reduced

Reduced Product

Let $(L, \alpha_1, \gamma_1, M_1)$ and $(L, \alpha_2, \gamma_2, M_2)$ be Galois connections.

The *reduced product* is the Galois *insertion*

$$(L, \alpha, \gamma, \varsigma[M_1 \times M_2])$$

defined by

$$\alpha(l) = (\alpha_1(l), \alpha_2(l))$$

$$\gamma(m_1, m_2) = \gamma_1(m_1) \sqcap \gamma_2(m_2)$$

$$\varsigma(m_1, m_2) = \sqcap\{(m'_1, m'_2) \mid \gamma_1(m_1) \sqcap \gamma_2(m_2) = \gamma_1(m'_1) \sqcap \gamma_2(m'_2)\}$$

405

© Nielsen^2, Hankin

Reduced Tensor Product

Let $(\mathcal{P}(V), \alpha_1, \gamma_1, \mathcal{P}(D_1))$ and $(\mathcal{P}(V), \alpha_2, \gamma_2, \mathcal{P}(D_2))$ be Galois connection.

The *reduced tensor product* is the Galois *insertion*

$$(\mathcal{P}(V), \alpha, \gamma, \varsigma[\mathcal{P}(D_1 \times D_2)])$$

defined by

$$\alpha(V') = \bigcup\{\alpha_1(\{v\}) \times \alpha_2(\{v\}) \mid v \in V'\}$$

$$\gamma(DD) = \{v \mid \alpha_1(\{v\}) \times \alpha_2(\{v\}) \subseteq DD\}$$

$$\varsigma(DD) = \bigcap\{DD' \mid \gamma(DD) = \gamma(DD')\}$$

406

© Nielsen^2, Hankin

Example: Array Bounds Analysis

The superfluous elements of $\mathcal{P}(\text{Sign} \times \text{Sign} \times \text{Range})$ will be removed when we use a reduced tensor product:

The reduction operator $\varsigma_{SSR'}$ amounts to

$$\varsigma_{SSR'}(SSR) = \bigcap\{SSR' \mid \gamma_{SSR'}(SSR) = \gamma_{SSR'}(SSR')\}$$

where $SSR, SSR' \subseteq \text{Sign} \times \text{Sign} \times \text{Range}$.

The singleton sets constructed from the following 16 elements

$$\begin{aligned} &(-, 0, <-1), (-, 0, -1), (-, 0, 0), \\ &(0, -, 0), (0, -, +1), (0, -, >+1), \\ &(0, 0, <-1), (0, 0, -1), (0, 0, +1), (0, 0, >+1), \\ &(0, +, 0), (0, +, +1), (0, +, >+1), \\ &(+, 0, <-1), (+, 0, -1), (+, 0, 0) \end{aligned}$$

will be mapped to the empty set (as they are useless).

407

© Nielsen^2, Hankin

Example (cont.): Array Bounds Analysis

The remaining 29 elements of $\text{Sign} \times \text{Sign} \times \text{Range}$ are

$$\begin{aligned} &(-, -, <-1), (-, -, -1), (-, -, 0), (-, -, +1), (-, -, >+1), \\ &(-, 0, +1), (-, 0, >+1), \\ &(-, +, <-1), (-, +, -1), (-, +, 0), (-, +, +1), (-, +, >+1), \\ &(0, -, <-1), (0, -, -1), (0, 0, 0), (0, +, <-1), (0, +, -1), \\ &(+, -, <-1), (+, -, -1), (+, -, 0), (+, -, +1), (+, -, >+1), \\ &(+, 0, +1), (+, 0, >+1), \\ &(+, +, <-1), (+, +, -1), (+, +, 0), (+, +, +1), (+, +, >+1) \end{aligned}$$

and they describe disjoint subsets of $\mathbb{Z} \times \mathbb{Z}$.

Any collection of properties can be described in 4 bytes.

408

© Nielsen^2, Hankin

Summary

The Array Bound Analysis has been designed from three simple Galois connections specified by extraction functions:

- (i) an analysis approximating integers by their sign,
- (ii) an analysis approximating pairs of integers by their difference in magnitude, and
- (iii) an analysis approximating integers by their closeness to 0, 1 and -1.

These analyses have been combined using:

- (iv) the relational product of analysis (i) with itself,
- (v) the functional composition of analyses (ii) and (iii), and
- (vi) the reduced tensor product of analyses (iv) and (v).

409

© Nielsen^2, Hankin

Induced Operations

Given: Galois connections $(L_i, \alpha_i, \gamma_i, M_i)$ so that M_i is more approximate than (i.e. is coarser than) L_i .

Aim: Replace an existing analysis over L_i with an analysis making use of the coarser structure of M_i .

Methods:

- **Inducing along the abstraction function:** move the computations from L_i to M_i .
- Application to Data Flow Analysis.
- **Inducing along the concretisation function:** move a widening from M_i to L_i .

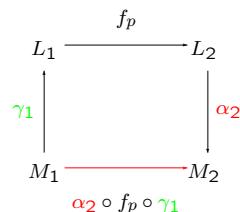
410

© Nielsen^2, Hankin

Inducing along the Abstraction Function

Given Galois connections $(L_i, \alpha_i, \gamma_i, M_i)$ so that M_i is more approximate than L_i .

Replace an existing analysis $f_p : L_1 \rightarrow L_2$ with a new and more approximate analysis $g_p : M_1 \rightarrow M_2$: take $g_p = \alpha_2 \circ f_p \circ \gamma_1$.



The analysis $\alpha_2 \circ f_p \circ \gamma_1$ is induced from f_p and the Galois connections.

411

© Nielsen^2, Hankin

Example:

A very precise analysis for plus based on $\mathcal{P}(\mathbf{Z})$ and $\mathcal{P}(\mathbf{Z} \times \mathbf{Z})$:

$$f_{\text{plus}}(ZZ) = \{z_1 + z_2 \mid (z_1, z_2) \in ZZ\}$$

Two Galois connections

$$\begin{aligned} & (\mathcal{P}(\mathbf{Z}), \alpha_{\text{sign}}, \gamma_{\text{sign}}, \mathcal{P}(\mathbf{Sign})) \\ & (\mathcal{P}(\mathbf{Z} \times \mathbf{Z}), \alpha_{\text{SS}'}, \gamma_{\text{SS}'}, \mathcal{P}(\mathbf{Sign} \times \mathbf{Sign})) \end{aligned}$$

An approximate analysis for plus based on $\mathcal{P}(\mathbf{Sign})$ and $\mathcal{P}(\mathbf{Sign} \times \mathbf{Sign})$:

$$g_{\text{plus}} = \alpha_{\text{sign}} \circ f_{\text{plus}} \circ \gamma_{\text{SS}'}$$

412

© Nielsen^2, Hankin

Example (cont.):

We calculate

$$\begin{aligned}
 g_{\text{plus}}(SS) &= \alpha_{\text{sign}}(f_{\text{plus}}(\gamma_{SS}(SS))) \\
 &= \alpha_{\text{sign}}(f_{\text{plus}}(\{(z_1, z_2) \in \mathbf{Z} \times \mathbf{Z} \mid (\text{sign}(z_1), \text{sign}(z_2)) \in SS\})) \\
 &= \alpha_{\text{sign}}(\{z_1 + z_2 \mid z_1, z_2 \in \mathbf{Z}, (\text{sign}(z_1), \text{sign}(z_2)) \in SS\}) \\
 &= \{\text{sign}(z_1 + z_2) \mid z_1, z_2 \in \mathbf{Z}, (\text{sign}(z_1), \text{sign}(z_2)) \in SS\} \\
 &= \bigcup \{s_1 \oplus s_2 \mid (s_1, s_2) \in SS\}
 \end{aligned}$$

where $\oplus : \text{Sign} \times \text{Sign} \rightarrow \mathcal{P}(\text{Sign})$ is the “addition” operator on signs (so e.g. $+ \oplus + = \{+\}$ and $+ \oplus - = \{-, 0, +\}$).

413

© Nielsen^2, Hankin

The Mundane Correctness of f_p carries over to g_p

The correctness relation R_i for V_i and L_i :

$$R_i : V_i \times L_i \rightarrow \{\text{true}, \text{false}\} \text{ is generated by } \beta_i : V_i \rightarrow L_i$$

Correctness of f_p means

$$(p \vdash \cdot \rightsquigarrow \cdot) (R_1 \rightarrow R_2) f_p$$

(with $R_1 \rightarrow R_2$ being generated by $\beta_1 \rightarrow \beta_2$).

The correctness relation S_i for V_i and M_i :

$$S_i : V_i \times M_i \rightarrow \{\text{true}, \text{false}\} \text{ is generated by } \alpha_i \circ \beta_i : V_i \rightarrow M_i$$

One can prove that

$$(p \vdash \cdot \rightsquigarrow \cdot) (R_1 \rightarrow R_2) f_p \wedge \alpha_2 \circ f_p \circ \gamma_1 \sqsubseteq g_p$$

$$\Rightarrow (p \vdash \cdot \rightsquigarrow \cdot) (S_1 \rightarrow S_2) g_p$$

with $S_1 \rightarrow S_2$ being generated by $(\alpha_1 \circ \beta_1) \rightarrow (\alpha_2 \circ \beta_2)$.

414

© Nielsen^2, Hankin

Fixed Points in the Induced Analysis

Let $f_p = \text{lfp}(F)$ for a monotone function $F : (L_1 \rightarrow L_2) \rightarrow (L_1 \rightarrow L_2)$.

The Galois connections $(L_i, \alpha_i, \gamma_i, M_i)$ give rise to a Galois connection $(L_1 \rightarrow L_2, \alpha, \gamma, M_1 \rightarrow M_2)$.

Take $g_p = \text{lfp}(G)$ where $G : (M_1 \rightarrow M_2) \rightarrow (M_1 \rightarrow M_2)$ is an “upper approximation” to F : we demand that $\alpha \circ F \circ \gamma \sqsubseteq G$.

Then for all $m \in M_1 \rightarrow M_2$:

$$G(m) \sqsubseteq m \Rightarrow F(\gamma(m)) \sqsubseteq \gamma(m)$$

and $\text{lfp}(F) \sqsubseteq \gamma(\text{lfp}(G))$ and $\alpha(\text{lfp}(F)) \sqsubseteq \text{lfp}(G)$

415

© Nielsen^2, Hankin

Application to Data Flow Analysis

A *generalised Monotone Framework* consists of:

- the property space: a complete lattice $L = (L, \sqsubseteq)$;
- the set \mathcal{F} of monotone functions from L to L .

An *instance* \mathbf{A} of a generalised Monotone Framework consists of:

- a finite flow, $F \subseteq \text{Lab} \times \text{Lab}$;
- a finite set of extremal labels, $E \subseteq \text{Lab}$;
- an extremal value, $\iota \in L$; and
- a mapping f from the labels Lab of F and E to monotone transfer functions from L to L .

416

© Nielsen^2, Hankin

Application to Data Flow Analysis

Let (L, α, γ, M) be a Galois connection.

Consider an instance B of the generalised Monotone Framework M that satisfies

- the mapping g from the labels Lab of F and E to monotone transfer functions of $M \rightarrow M$ satisfies $g_\ell \sqsubseteq \alpha \circ f_\ell \circ \gamma$ for all ℓ ; and
- the extremal value \jmath satisfies $\gamma(\jmath) = \iota$;

and otherwise B is as A .

One can show that a solution to the B -constraints gives rise to a solution to the A -constraints:

$$(B_o, B_\bullet) \models B^\exists \text{ implies } (\gamma \circ B_o, \gamma \circ B_\bullet) \models A^\exists$$

© Nielsen^2, Hankin

The Mundane Approach to Semantic Correctness

Here $F = \text{flow}(S_*)$ and $E = \{\text{init}(S_*)\}$.

Correctness of every solution to A^\exists amounts to:

Assume $(A_o, A_\bullet) \models A^\exists$ and $\langle S_*, \sigma_1 \rangle \rightarrow^* \sigma_2$.

Then $\beta(\sigma_1) \sqsubseteq \iota$ implies $\beta(\sigma_2) \sqsubseteq \sqcup\{A_\bullet(\ell) \mid \ell \in \text{final}(S_*)\}$.

where $\beta : \text{State} \rightarrow L$.

One can then prove the correctness result for B :

Assume $(B_o, B_\bullet) \models B^\exists$ and $\langle S_*, \sigma_1 \rangle \rightarrow^* \sigma_2$.

Then $(\alpha \circ \beta)(\sigma_1) \sqsubseteq \jmath$ implies $(\alpha \circ \beta)(\sigma_2) \sqsubseteq \sqcup\{B_\bullet(\ell) \mid \ell \in \text{final}(S_*)\}$.

418

© Nielsen^2, Hankin

Sets of States Analysis

Generalised Monotone Framework over $(\mathcal{P}(\text{State}), \sqsubseteq)$.

Instance SS for S_* :

- the flow F is $\text{flow}(S_*)$;
- the set E of extremal labels is $\{\text{init}(S_*)\}$;
- the extremal value ι is State ; and
- the transfer functions are given by f^{SS} :

$$[x := a]^\ell : f_\ell^{\text{SS}}(\Sigma) = \{\sigma[x \mapsto A[\![a]\!] \sigma] \mid \sigma \in \Sigma\}$$

$$[\text{skip}]^\ell : f_\ell^{\text{SS}}(\Sigma) = \Sigma$$

$$[b]^\ell : f_\ell^{\text{SS}}(\Sigma) = \Sigma$$

where $\Sigma \subseteq \text{State}$.

Correctness: Assume $(SS_o, SS_\bullet) \models SS^\exists$ and $\langle S_*, \sigma_1 \rangle \rightarrow^* \sigma_2$. Then $\sigma_1 \in \text{State}$ implies $\sigma_2 \in \sqcup\{SS_\bullet(\ell) \mid \ell \in \text{final}(S_*)\}$.

419

© Nielsen^2, Hankin

Constant Propagation Analysis

Generalised Monotone Framework over $\widehat{\text{State}}_{\text{CP}} = ((\text{Var} \rightarrow \mathbf{Z}^\top)_\perp, \sqsubseteq)$.

Instance CP for S_* :

- the flow F is $\text{flow}(S_*)$;
- the set E of extremal labels is $\{\text{init}(S_*)\}$;
- the extremal value ι is $\lambda x. \top$; and
- the transfer functions are given by the mapping f^{CP} :

$$[x := a]^\ell : f_\ell^{\text{CP}}(\hat{\sigma}) = \begin{cases} \perp & \text{if } \hat{\sigma} = \perp \\ \hat{\sigma}[x \mapsto A_{\text{CP}}[\![a]\!] \hat{\sigma}] & \text{otherwise} \end{cases}$$

$$[\text{skip}]^\ell : f_\ell^{\text{CP}}(\hat{\sigma}) = \hat{\sigma}$$

$$[b]^\ell : f_\ell^{\text{CP}}(\hat{\sigma}) = \hat{\sigma}$$

420

© Nielsen^2, Hankin

Galois Connection

The representation function $\beta_{CP} : \text{State} \rightarrow \widehat{\text{State}}_{CP}$ is defined by

$$\beta_{CP}(\sigma) = \sigma$$

This gives rise to a Galois connection

$$(\mathcal{P}(\text{State}), \alpha_{CP}, \gamma_{CP}, \widehat{\text{State}}_{CP})$$

where $\alpha_{CP}(\Sigma) = \sqcup\{\beta_{CP}(\sigma) \mid \sigma \in \Sigma\}$ and $\gamma_{CP}(\hat{\sigma}) = \{\sigma \mid \beta_{CP}(\sigma) \sqsubseteq \hat{\sigma}\}$.

One can show that for all labels ℓ

$$f_\ell^{CP} \sqsupseteq \alpha_{CP} \circ f_\ell^{SS} \circ \gamma_{CP} \quad \text{as well as } \gamma_{CP}(\lambda x. \top) = \text{State}$$

It follows that $[CP]$ is an upper approximation to the analysis induced from $[SS]$ and the Galois connection; therefore it is correct.

421

© Nielsen^2, Hankin

Inducing along the Concretisation Function

Given an upper bound operator

$$\nabla_M : M \times M \rightarrow M$$

and a Galois connection (L, α, γ, M) .

Define an upper bound operator

$$\nabla_L : L \times L \rightarrow L$$

by

$$l_1 \nabla_L l_2 = \gamma(\alpha(l_1) \nabla_M \alpha(l_2))$$

It defines a widening operator if one of the following conditions holds:

- (i) M satisfies the Ascending Chain Condition, or
- (ii) (L, α, γ, M) is a Galois insertion and $\nabla_M : M \times M \rightarrow M$ is a widening.

422

© Nielsen^2, Hankin

Precision of the Induced Widening Operator

Lemma: Let (L, α, γ, M) be a Galois insertion such that $\gamma(\perp_M) = \perp_L$ and let $\nabla_M : M \times M \rightarrow M$ be a widening operator.

Then the widening operator $\nabla_L : L \times L \rightarrow L$ defined by

$$l_1 \nabla_L l_2 = \gamma(\alpha(l_1) \nabla_M \alpha(l_2))$$

satisfies

$$\text{Ifp}_{\nabla_L}(f) = \gamma(\text{Ifp}_{\nabla_M}(\alpha \circ f \circ \gamma))$$

for all monotone functions $f : L \rightarrow L$.

423

© Nielsen^2, Hankin

Precision of the Induced Widening Operator

Corollary: Let M be of finite height, let (L, α, γ, M) be a Galois insertion (such that $\gamma(\perp_M) = \perp_L$), and let ∇_M equal the least upper bound operator \sqcup_M .

Then the above lemma shows that $\text{Ifp}_{\nabla_L}(f) = \gamma(\text{Ifp}(\alpha \circ f \circ \gamma))$.

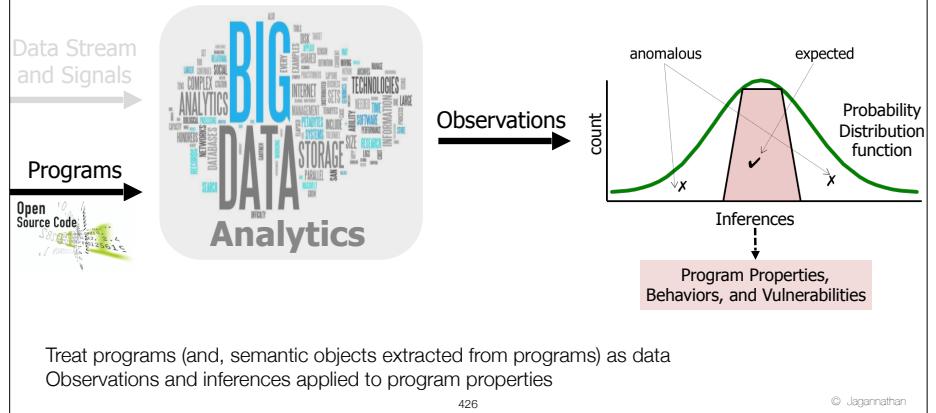
This means that $\text{Ifp}_{\nabla_L}(f)$ equals the result we would have obtained if we decided to work with $\alpha \circ f \circ \gamma : M \rightarrow M$ instead of the given $f : L \rightarrow L$; furthermore the number of iterations needed turn out to be the same. However, for all other operations the increased precision of L is available.

424

© Nielsen^2, Hankin

Analyzing Big Code

Apply principles of Big Data Analytics to a large corpus



426

© Jagannathan

Redundancies in the corpus exposed as dense components (enclaves) in the mined network

- Nodes represent properties facts, claims, and evidence
- Edges connect related properties

Anomalous properties have small number of connections

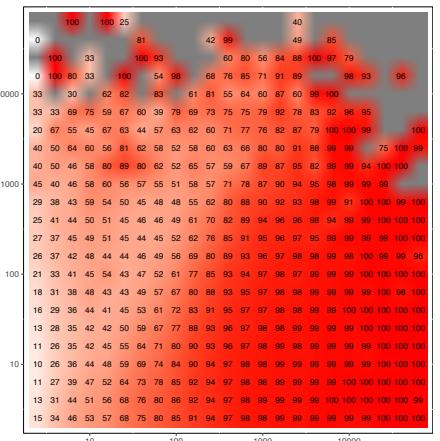
Likely invariants have large number of connections

<https://muse-portal.net/publications>

© Jagannathan

Duplication

| | | JavaScript |
|--------|-------------------------|-------------|
| Counts | # projects (total) | 4,479,173 |
| | # projects (non-fork) | 2,011,875 |
| | # projects (downloaded) | 1,778,679 |
| | # projects (analyzed) | 1,755,618 |
| | # files (analyzed) | 261,676,091 |



JavaScript

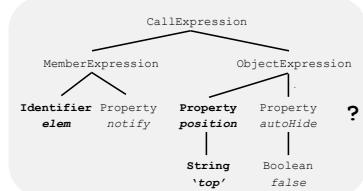
Lopes, Maj, Vitek, et al. DéjàVu: A Map of Code Duplicates on GitHub. OOPSLA'17

Probabilistic models of code

Program

```
elem.notify({
  position: 'top',
  autoHide: false,
?
});
```

AST



Probabilistic Model

$$P(\text{element} \mid S)$$

S is a learned semantic analyzer

Can predict entire expressions at once, precisely

Raychev, Bielik, Vechev. Probabilistic Model for Code with Decision Trees. ACM OOPSLA 2016

429

© Vechev

Language Translation

JavaScript

```
function addListener(element, event, handler) {
  if (!element.addEventListener) {
    element.attachEvent('on' + event, handler);
  } else {
    element.addEventListener(event, handler, false);
  }
}

function removeListener(element, event, handler) {
  if (!element.removeEventListener) {
    element.detachEvent('on' + event, handler);
  } else {
    element.removeEventListener(event, handler, false);
  }
}
```

CoffeeScript

```
listener = (el, event, handler) ->
  if (!el.addEventListener)
    el.attachEvent('on' + event, handler)
  else
    el.addEventListener(event, handler, false)

removeListener = (el, event, handler) ->
  if (!el.removeEventListener)
    el.detachEvent('on' + event, handler)
  else
    el.removeEventListener(event, handler, false)
```

Sounds easy, but very hard problem:
0 margin for error, different than NLP
(where similar intent fine)

Easy to get data, learning **hard to get scale**, with 0 mistakes!

Deep Learning **does not work well**

430

© Vechev

Statistical Deobfuscation

Major challenge: fast inference with precision.
Various works sacrifice one, making tool less usable in practice

Analyzer again learns **right few features** to use.
Fast to extract, yet precise.

```
function chunkData(str: string,
step: number) {
  var colNames = [];
  var len = str.length;
  var i = 0;
  for (; i < len; i += step)
    if (i + step < len)
      colNames.push(str.substring(i, i + step));
    else
      colNames.push(str.substring(i, len));
  return colNames;
}
```

JS NICE

JSNice.org

Raychev, Vechev. Predicting Program Properties from "Big Code". POPL 2015

© Vechev

Learn from thousands of commits

Example Fixes

cinnabar/msgcodec

Fix byte->int expansion bug. Apply 0xFF mask when returning int representation of byte data.

```
@Override
public int read() throws IOException {
-   return data[pos++];
+
+   return 0xFF & data[pos++];
}
```

Explanations from Others

Fix a bug when the byte read is -1.

Predict bugs elsewhere Make explainable predictions

This issue was fixed by 66 projects. Below we summarize their changes.

Fix byte->int expansion bug. Apply 0xFF mask when returning int representation of byte data.

```
@Override
public synchronized int read() throws IOException {
  final byte[] b = new byte[1];
  if (doRead(b, 0, 1) == -1) {
    return -1;
  }
  return b[0];
}
```

Semantic analysis over large corpus of code and learning using **interpretable** model

Key for displaying **actionable** suggestions to users

Enable learning of new static analysis checks [PLDI'18]

© Vechev

Code Vectors

Symbolic trace

```
int example() {
    buf = alloc(12);
    if (buf != 0) {
        bar(buf);
        free(buf);
        return 0;
    } else {
        return -ENOMEM;
    }
}
```

Deixe a paramétrica codificação de programas em sequências de palavras que (i) pode ser ajustada para capturar diferentes escolhas de representação na escala do espectro, principalmente sintática e (ii) é amigável a técnicas de aprendizado de vetores de palavras, e (iii) pode ser obtida eficientemente de programas.

Figure 1: An example procedure

```
call alloc(12);           call alloc(12);
assume alloc(12) != 0;    assume alloc(12) == 0;
call bar(alloc(12));
call free(alloc(12));
return 0;
```

(a) Trace 1

(b) Trace 2

Abstracted Symbolic Traces

```
lock(&obj->lock);
foo = alloc(12);
if (foo != 0) {
    obj->baz =
        bar(foo);

} else {
    unlock(
        &obj->lock);
    return -ENOMEM;
}

}
```

Há muitas palavras em uma trilha simbólica.

Local variáveis, funções definidas pelo usuário, literais e condições de caminho combinam para criar um vocabulário extraordinariamente grande.

```
Start
Called(lock)
Called(alloca)
RetEq(alloca, 0)
ParamShare(unlock, lock)
Called(unlock)
Error
RetError(-ENOMEM)
End
```