

Scalable Defect Detection

Manuvir Das, Zhe Yang, Daniel Wang
Center for Software Excellence
Microsoft Corporation

Part II High-Quality Scalable Checking *using modular path-sensitive analysis*

Zhe Yang
Center for Software Excellence
Microsoft

Secret Sauce for a Practical Checker

Keys to high-quality scalable checkers

- Scalability: *checking each function in isolation*
- Quality: *path sensitivity and defect prioritization*

Approach proven by our experience at Microsoft

- espX: buffer-overflow checker, widely deployed and used to get 20,000+ bugs found and fixed
- μ SpaCE: checker-building SDK, used by non-experts to build domain-rule-enforcing checkers

7/20/2007

Quality Checker via Path Sensitive Analysis

3

Scalability and Quality Overview

Scalability: Inter-Procedural Analysis?

Lecture 1 (by Manuvir): scalable inter-procedural analysis is possible, with

- Good Techniques: summarization, etc.
- Constraints on problems: finite automata, etc.

But

- **Intractable** for complex states (buffer overrun).
- **Mismatch** with the modular reasoning by devs.
 - “If an error is detected, who to blame”

7/20/2007

Quality Checker via Path Sensitive Analysis

5

Linear Scalability by Modular Analysis

If we can afford to analyze each function in isolation

- Scales up linearly in # of functions and scales out
 - Allows using complex states for accuracy
- But it's a big “if”.

- For example, is this function safe?

```
void f(int *buf, size_t n)
{ for (size_t i=0; i <= n; i++) buf[i] = 0; }
```

- Modular analysis requires specifications of the usage context (e.g., “*buf* has *n* elements”).

7/20/2007

Quality Checker via Path Sensitive Analysis

6

Assumption: specification possible

- “Did you say specifications?”
 - Isn’t it a pipe dream to design practical spec langs?
 - Who is going to add specs to millions of functions?
- This is the subject of Lecture 3 on SAL (by Dan)
- For now, assume functions come equipped with necessary specifications of contexts.
 - Say “void f(int *buf, size_t n)” →
“void f(int<n> *buf, size_t n)”
- So we can discuss modular checking in full detail.

7/20/2007

Quality Checker via Path Sensitive Analysis

7

Quality: The measures

- **Accurate**: fix rate (% bugs fixed), false positive rate (% of reported bugs deemed noise)
 - Dev’s perspective: frustrated with bogus issues.
- **Comprehensive**: validation rate (% of safe code), false negative rate (% of missed issues)
 - Exec’s perspective: measure of coverage/progress.
- **Clear** and **Actionable**: easy to understand the reported defects and take appropriate actions

7/20/2007

Quality Checker via Path Sensitive Analysis

8

Quality Measures: Historical Perspective at Microsoft

- Early years
 - Bugs found by static analysis met with excitement
 - **Accuracy** is the obvious tool quality for devs
- After a few years of worm-induced news
 - “how many bugs are left?”
 - Measure of coverage calls for **comprehensive** validation
- Use of symbolic abstraction improved coverage
 - “I can’t understand what this message is saying.”
 - “There are so many issues and so little time left.”
 - Messages need to be **clear** and **prioritized**

7/20/2007

Quality Checker via Path Sensitive Analysis

9

Achieving Quality

Clarity for developers to take action

- Using path-sensitive analysis instead of data flow analysis (since devs reason with paths)

Conflicting Goals: (**Accurate**) defect detection vs (**comprehensive**) validation?

- Both: use comprehensive validation as a basis, and then expose defects through prioritization

7/20/2007

Quality Checker via Path Sensitive Analysis

10

Detection

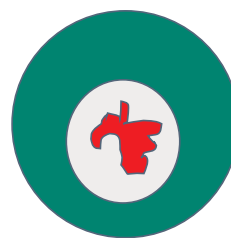


7/20/2007

Quality Checker via Path Sensitive Analysis

11

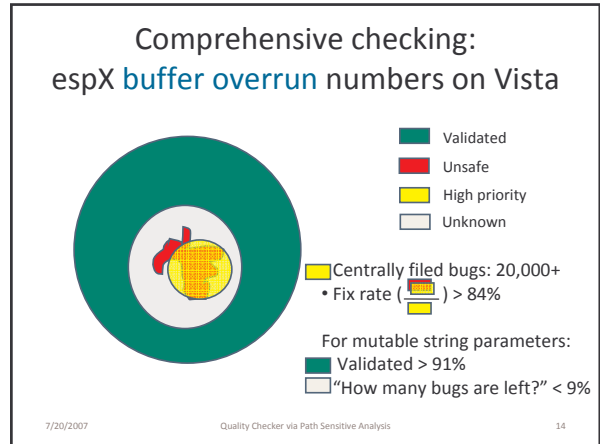
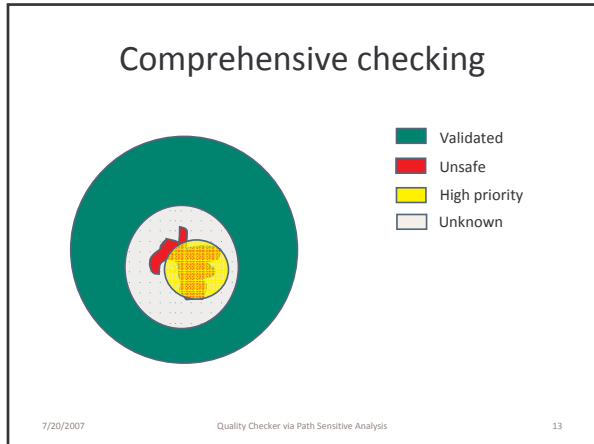
Validation



7/20/2007

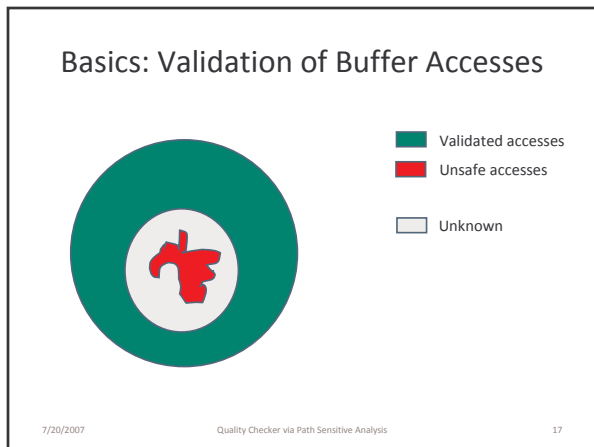
Quality Checker via Path Sensitive Analysis

12



espX: Buffer Overrun Checker

- ### Outline of this section of talk
- Basics of a buffer-overrun checker
 - Prior Art: merge-based dataflow analysis
 - Our Approach: path-sensitive analysis
 - Warning bucketing for prioritization
- 7/20/2007 Quality Checker via Path Sensitive Analysis 16



Example 1

```

BYTE<n> *Alloc(size_t n);
void FillRects(RECT<n> *r, size_t n);
void FillPoints(POINT<n> *p, size_t n);

void Fill(unsigned int r, unsigned int p)
{
  BYTE *buf = Alloc(r * sizeof(RECT) + p * sizeof(POINT));
  FillRects((RECT *)buf, r);
  buf += r * sizeof(RECT);
  FillPoints((POINT *)buf, p);
}
  
```

$16 \times r$ $8 \times p$
Rectangles Points

7/20/2007 Quality Checker via Path Sensitive Analysis 18

“Instrumenting” the Program

```

BYTE *buf = Alloc(r * sizeof(RECT) +
                 p * sizeof(POINT));
assume: offset(buf) = 0 □
bcap(buf) = 16 × r + 8 × p
assert: offset(buf) + 16 × r ≤ bcap(buf)
FillRects((RECT *)buf, r);

buf += r * sizeof(RECT);

assert: offset(buf) + 8 × p ≤ bcap(buf)
FillPoints((POINT *)buf, p);
    
```

7/20/2007

Quality Checker via Path Sensitive Analysis

19

Analysis of Example 1

```

BYTE *buf = Alloc(r * sizeof(RECT) +
                 p * sizeof(POINT));
assume: offset(buf) = 0 □
bcap(buf) = 16 × r + 8 × p
assert: offset(buf) + 16 × r ≤ bcap(buf)
FillRects((RECT *)buf, r);

buf += r * sizeof(RECT);

assert: offset(buf) + 8 × p ≤ bcap(buf)
FillPoints((POINT *)buf, p);
    
```

$\{r \geq 0; p \geq 0\}$

$\{bcap(buf) = 16 \times r + 8 \times p;$
 $offset(buf) = 0; r \geq 0; p \geq 0\}$
 $[offset(buf) + 16 \times r = 16 \times r$
 $\leq bcap(buf)]$ PASS

$\{bcap(buf) = 16 \times r + 8 \times p;$
 $offset(buf) = 16 \times r; r \geq 0; p \geq 0\}$
 $[offset(buf) + 8 \times p =$
 $16 \times r + 8 \times p \leq bcap(buf)]$ PASS

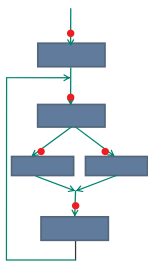
Need symbolic state tracking +
 linear integer theorem prover

7/20/2007

Quality Checker via Path Sensitive Analysis

20

Dataflow Analysis



Task: find invariants at CFG-nodes

Find a map $A: V \rightarrow Abs$
 stable under $T: E \times Abs \rightarrow Abs$

Abs : lattice of abstract values
 Stability condition:
 $A(v) = \bigcup \{T(e, A(u)) : e = (u, v) \in E\}$
 or $A = \tilde{T}(A)$ is a fixed-point of \tilde{T}

If T monotone, Abs complete, then
 least $A = \bigcup \{\tilde{T}^i(\perp) : i = 0, \dots\}$.
 This terminates if Abs finite in height.

Given flow graph (V, E)

Work-list algorithm used in practice.

7/20/2007

Quality Checker via Path Sensitive Analysis

Work-list algorithm used in practice.

Dataflow Analysis for Buffer Overruns [Dor et al:PLDI 2003]

To track the symbolic states, use a dataflow analysis

- Abs = Set of Linear Inequality Constraints
- T : suitably abstracted from concrete semantics
- But what about the join operator?
- I.e., how do you merge two sets of linear constraints into another set of linear constraints (that is implied by both of them)?
- Answer: Polyhedra (Cousot/Halbwachs:POPL78)

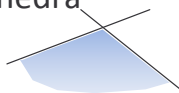
7/20/2007

Quality Checker via Path Sensitive Analysis

22

The Lattice of Polyhedra

- Geometric Interpretation:
 - One linear inequality gives a half-space
 - A set of linear inequalities is a (maybe-not-closed) convex polyhedron (n-dimensional polygon)
- Join-operator needs to find the smallest enclosing polyhedron (convex hull problem)
- Algorithm involves lots of linear programming
- infinite-height lattice: termination for loops?

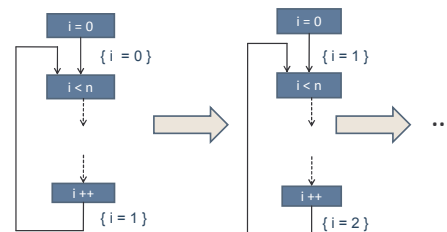


7/20/2007

Quality Checker via Path Sensitive Analysis

23

What about Loops ?

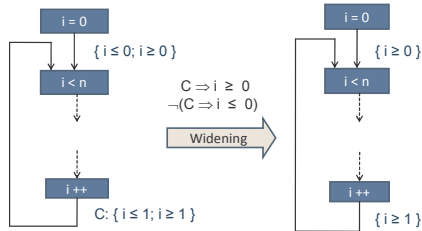


7/20/2007

Quality Checker via Path Sensitive Analysis

24

Loop Widening



7/20/2007

Quality Checker via Path Sensitive Analysis

25

Handling loops

- Use a loop widening algorithm to ensure that the analysis terminates
 - Widening operator: Weaken the constraints along a back edge of a loop in a way that ensures that finite number of such weakenings is sufficient
 - Mathematically, any chain formed by repeated application of the widening operator is finite.

7/20/2007

Quality Checker via Path Sensitive Analysis

26

Issues with Polyhedra

- Complexity (implementation & cost)
 - Several restricted version proposed and used:
 - Octagons (at most two variables; coefficients 1, -1)
 - Arbitrary predetermined shapes.
 - Inaccuracy
 - The convex hull won't be accurate closure
 - Real-numbered coefficients would appear
 - An approximation for Integer Linear Constraints
- Bigger issue: feedback to developer

7/20/2007

Quality Checker via Path Sensitive Analysis

27

Example 2

```

BYTE<size> *Alloc(size_t size);
void StringCopy(wchar_t*<en>*dest, const wchar_t {null-terminated} *src, size_t n);

void ProcessString(wchar_t *str)
{
    wchar_t buf[100];
    wchar_t *tmp = &buf;

    int len = wcslen(str) + 1;
    if (len > 100)
        tmp = (wchar_t *)Alloc(len);

    StringCopy(tmp, str, len);
    ...
}

```

Should be Alloc(len * sizeof(wchar_t))
Buffer overrun
{bcap(tmp) = 200; len ≤ 100} vs {bcap(tmp) = len; len > 100}

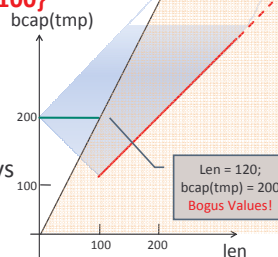
7/20/2007

Quality Checker via Path Sensitive Analysis

28

Example 2 with Polyhedra

- Merging **{bcap(tmp) = 200; len ≤ 100}** and **{bcap(tmp) = len; len > 100}**
- That is:
 - $bcap(tmp) \geq len$
 - $bcap(tmp) + len \geq 200$
 - $bcap(tmp) \geq len + 200$
- Obscure message to devs
 Need path-based analysis:
"overflow when len ≥ 100"

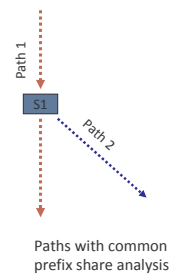


7/20/2007

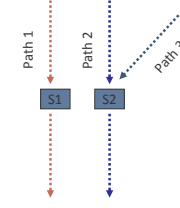
Quality Checker via Path Sensitive Analysis

29

Sharing in Path Sensitive Analysis



Paths with common prefix share analysis



Paths which reach the same program point with the same symbolic state share a suffix

7/20/2007

Quality Checker via Path Sensitive Analysis

30

Path-Sensitive Dataflow Analysis

- In its simplest form, path-sensitive analysis can be characterized as a dataflow analysis
- Find $A: V \rightarrow \text{Set}(\text{state})$ using $t: E \times \text{state} \rightarrow \text{state}$
- $T: \text{Set}(\text{state}) \rightarrow \text{Set}(\text{state})$ is the point-wise lifted version of t , using set-union.
- $\text{Set}(\text{state})$ is a complete lattice; T is monotone
- But $\text{Set}(\text{state})$ is infinite in height when the universe of states is infinite.

7/20/2007

Quality Checker via Path Sensitive Analysis

31

Widening in Path-Sensitive Analysis

- Issue: We share paths only when the symbolic states are the same at a node; but loops induce infinite number of states.
- “Widening? But what state to widen against?”
- Idea: at back edge, widen against the path itself
- Solution: extend the state to record the path-history of states at each loop entry node.
- $\text{WidenedState} = \text{LoopNestingLevel} \rightarrow \text{State}$
[s_1, s_2, s_3]: state is s_3 now, and was s_1 at loop level i
- Exercise: work out the detail

7/20/2007

Quality Checker via Path Sensitive Analysis

32

Fast Theorem Prover for Integer Linear Inequalities?

- Not asking for: constraint solver, completeness
- Observation 1: developers reasoning about linear constraints in a simple way; often: a proof is just a linear combination with small integer coefficients.
- Observation 2: difference constraint theorem prover is easy to construct. [CLR:alg-textbook]
- Exercise: figure out an algorithm.

7/20/2007

Quality Checker via Path Sensitive Analysis

33

Elements of the checker

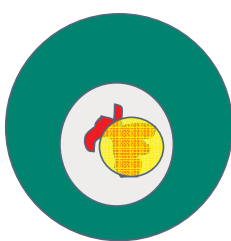
- Symbolic state tracking with linear inequalities
 - Path sensitive analysis
 - Path sensitive loop widening
- Fast linear integer theorem prover

7/20/2007

Quality Checker via Path Sensitive Analysis

34

Defect Bucketing/Prioritization



- Validated
- Unsafe
- High priority
- Unknown

7/20/2007

Quality Checker via Path Sensitive Analysis

35

Example 1 – provable error

```
if (CanProve(buffer index < buffer size))
    Validated Access
else
    if (CanProve(buffer index >= buffer size))
        Provable Error
    else
        Possible Error
```

- e.g. passing byte count instead of element count

```
wcscpy_s(buf, sizeof(buf), s); espX Warning 26000
```

7/20/2007

Quality Checker via Path Sensitive Analysis

36

Example 2 – incorrect validation

```
int glob[BUF_SIZE];
bool read(int i, int *val) {
    if (i > BUF_SIZE) // Off by one
        return false;
    assert: i < BUF_SIZE
    *val = glob[i];
    ...
}
```

espX Warning 26014:
Cannot prove: $i < BUF_SIZE$
Can prove: $i < BUF_SIZE + 1$

e.g. MS01-033(Code Red), MS04-036(NNTP), MS04-035(SMTP)

7/20/2007

Quality Checker via Path Sensitive Analysis

37

Example 3 – missing validation

```
void Transform(char<size> *dest,
               const char<null-terminated> *src,
               size_t size) {
    assert: strlen(src) + 1 <= size
    memcpy(dest, src, strlen(src) + 1);
}
```

espX Warning 26015:
Constraint set does not relate size and strlen(src)

e.g. MS03-026(Blaster), MS05-039 (Zotob)

7/20/2007

Quality Checker via Path Sensitive Analysis

38

Warning bucketing criteria

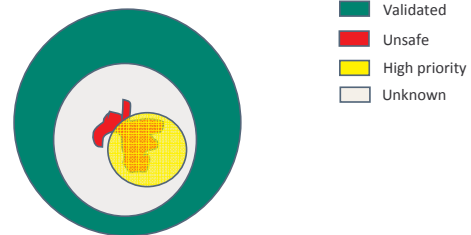
- Are heuristics based on observations of common coding mistakes
- Are semantic, not syntactic, in nature
 - Makes them robust
- Validated by security bulletin bug data and Watson crash data

7/20/2007

Quality Checker via Path Sensitive Analysis

39

Precision Improvements



7/20/2007

Quality Checker via Path Sensitive Analysis

40

Loop Invariant Inference

```
void StripSpaces(char<n> *dest, char *src, size_t n)
```

```
{
    while (*src != 0 && n > 1) {
        assume: n > 1
        if (*src != ' ') {
            assert: offset(dest) < n0
            *dest++ = *src;
            n--;
        }
        src++;
    }
    *dest = '\0';
}
```

← Need loop invariant
offset(dest) + n = n₀

espX deduces offset(dest) and n are *synchronized variables* in the loop

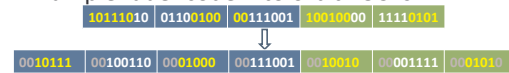
7/20/2007

Quality Checker via Path Sensitive Analysis

41

Combining Theorem Provers

- Example: uuencode into 6-bit ASCII's



```
void uuencode(BYTE<n> *src,
              BYTE<(n+2)/3*4> *dest, size_t n)
```

(Real spec added by a developer to real code)

- A second layer of theorem prover to uninterpreted operations, integer divisions, modular arithmetic, bitwise operations & etc.

7/20/2007

Quality Checker via Path Sensitive Analysis

42

espX Summary

- espX have made comprehensive defect detection a reality for buffers
 - Tens of thousands of bugs found and fixed
 - “How many bugs are left ?”
 - < 9% for mutable string buffers in Vista
 - Specifications also important (Details in Lecture 3)
- Achieved using
 - Modular Path-sensitive analysis
 - Careful warning bucketing and prioritization
 - Assortment of precision-refinement techniques

7/20/2007

Quality Checker via Path Sensitive Analysis

43

Devs want to build good checkers, too

- Developers who are domain experts often want to enforce certain domain-specific rules
- Encouraged by our work, they want to go static
- E.g.: Project Goldmine (Internationalization)

```
void IssueMessage()
{
    ::MessageBox(NULL,
        L"Failed to load file", ← Hard-coded strings
        MB_ERROR | MB_OK);      passed to user facing API
}
```

7/20/2007

Quality Checker via Path Sensitive Analysis

44

Developer-Generated Analyses: μSpaCE

(Or Better: You Checker)



How Do We Share Our Expertise?

- We understand static analysis well, but we can't solve problems for all domains
- MS solution: an SDK for domain experts to build path-sensitive dataflow analysis
- Challenge: intelligible explanation, i.e., without lattice, monotone function, join, etc.
- Our explanation: based on “Path Iteration”

7/20/2007

Quality Checker via Path Sensitive Analysis

46

Path Iteration

- Get set of paths; traverse them separately
- Simulation-style code:

For all paths p.
{
...
For all edges e in p.
...
}

- Limitation:
 - Cannot have full coverage
 - No sharing of analysis across paths

```
int f(int i, int n)
{
    if (i) = 1;
    else = 2;
    i++;
    while (i < n)
        i++;
    return i;
}
```

7/20/2007

Quality Checker via Path Sensitive Analysis

47

- Get set of paths; traverse them separately
- Explicit simulation state:

For all paths p.
{
Φ = Φ_{initial};
For all edges e in p.
Φ = t(Φ, e);
}

- Benefit of abstraction
 - Under-the-hood improvements

```
int f(int i, int n)
{
    if (i) = 1;
    else = 2;
    i++;
    while (i < n)
        i++;
    return i;
}
```

7/20/2007

Quality Checker via Path Sensitive Analysis

48

Path-Sensitive Analysis

- Define transfer function with explicit abstract state
- The μ SpaCE engine maintains
 - A state set pr. node
 - Reuses path computations
 - Covers state space 100%
- The fine print:
 - State domain needs to be finite
 - Or else widening operator needed

```

int f(int i, int n)
{
    if (i)
        i = 1;
    else
        i = 2;
    i = 3;
    while (i < n)
        i = 4;
    return i;
}
    
```

7/20/2007

Quality Checker via Path Sensitive Analysis

49

μ SpaCE SDK for Building Checks

- SDK: a concise core (with virtual transfer functions) + oracles (memory model, spec semantics, etc.)
- Multiple clients in one year
 - Goldmine (C/C++/.NET): intl. checker & meta data gen.
 - espC (C/C++): concurrency checker
 - iCatcher (.NET): cross-site scripting checker for ASP.NET
 - NullPtr (C/C++/.NET): spec-based null-ptr checker
- All these clients have found real bugs; they are getting deployed company wide

7/20/2007

Quality Checker via Path Sensitive Analysis

50

Summary

Keys to high-quality scalable checkers

- Scalability: *checking each function in isolation*
- Quality: *path sensitivity and defect prioritization*

Approach proven by our experience at Microsoft

- espX: buffer-overrun checker, widely deployed and used to get 20,000+ bugs found and fixed
- μ SpaCE: checker-building SDK, used by non-experts to build domain-rule-enforcing checkers

7/20/2007

Quality Checker via Path Sensitive Analysis

51

Exercises & Recommended Readings

for "High-Quality Scalable Checking using Modular Path-Sensitive Analysis"

Exercises:

1. **Path-Sensitive Loop Widening:** Work out the mathematical detail of path-sensitive loop widening. Prove that the algorithm terminates.
2. **Defect Prioritization:** Using the observation that linear combinations with small integer coefficients would suffice most of the time for developers' reasoning, construct a heuristic prover for linear inequality constraints based on graph search. A good starting point is the reduction of solving difference constraints to the shortest path problem (Cormen, et al. 2001). How would you add the missing power for modular arithmetic?
3. **Defect Prioritization:** A path-sensitive analysis uses widening to achieve coverage while ensuring termination and maintaining a reasonable cost. But when you want to show a potential issue to a developer, it is best to show them the specific information along some paths. For one defect found, there could be infinite number of paths leading to the defect, due to the presence of loops. How would you present the set of infinite paths to the developer in an understandable manner?

If you feel like you can send your solutions, comments, or questions to me (zhey.yang@microsoft.com)

Recommended Readings

- On Buffer Access Checking
 (Orr, Rubin and Sagie 2005). Acme-ecore framework for buffer access checking.
- (Cousot and Halbwachs 1978). Original paper on the lattice of polyhedra and its application to find linear constraints in a program.
- (Puckert, et al. 2000). Our work on buffer access checking. This article offers an overview of the workflow in use at Microsoft, the optimization language, the reference algorithm of the analyzer, as well as a high-level sketch of the path-sensitive checker.
- General Introduction to Program Analysis and Domain Theory
 (DeRemortel 1984). An efficient introduction to the techniques of program analysis.
- (Heintz, Norton and Hedin 1998). Standard textbook on program analysis.
- (Winkler 1995). Chapter 8 offers a good introduction to domain theory.
- Path-Sensitive Dataflow Analysis
 (Blenc, Narain and Sagie 1992). Some work on scalable inter-procedural analysis. The simplest form of path-sensitive analysis can be understood as graph reachability.

...

7/20/2007

Quality Checker via Path Sensitive Analysis

52

<http://www.microsoft.com/cse>
<http://research.microsoft.com/users/zhey.zhe.yang@microsoft.com>

© 2007 Microsoft Corporation. All rights reserved.
 This presentation is for informational purposes only.
 MICROSOFT MAKES NO WARRANTIES, EXPRESS OR IMPLIED, IN THIS SUMMARY.

7/20/2007

Quality Checker via Path Sensitive Analysis

53