

CMSC 631 – Program Analysis and Understanding Spring 2006

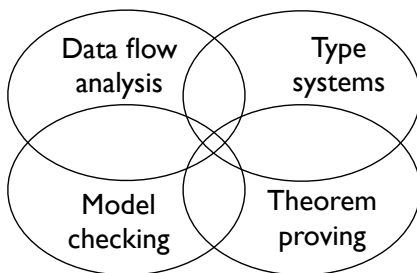
New Directions

Program Analysis and Understanding

- This semester, we've covered a lot of material about programs and programming languages
- Main areas of static program analysis:
 - Data flow analysis
 - Abstract interpretation
 - Type systems
 - Theorem Proving
 - Model checking
- Today: An assortment of things we didn't cover

2

The Space of Static Analyses



- Are these four really different?
 - What is the connection between them?

3

Type Systems and Theorem Proving

- Type systems are really dumb theorem provers
 - Only theorems are (basically) “*e* has type *t*,” but...
 - Type checking and inference are decidable and efficient (often)
 - Type systems also know what they can't prove
 - Can usually also prove “*e* has no type”
- Theorem proving systems are really smart
 - Can prove very ambitious theorems
 - Often undecidable, running time hard to predict
 - Hard to understand what works and what doesn't
 - Encoded in the decision procedure

4

Proof-Carrying Code (Necula and Lee)

- How do you know a program is safe to run?
 - It's pretty difficult to decide given just the code!
 - Coming up with the proof of correctness is hard
 - Idea: code comes with a safety proof
 - You just check the proof, which is often “easy”
- Applications:
 - Mobile/distributed code
 - Compilers (do you trust gcc?)

5

Types and PCC

- So what are these proofs?
 - In theory, any thing you like
 - In applications so far, type (and memory) safety!
- But...type systems are easy theorem provers
 - Except we need type safety proofs for executable code
 - Translation from high-level source obscures details
- Enter typed assembly language (Morrisett et al)
 - Bring types all the way through the compiler

6

Data Flow and Model Checking

- Schmidt, "Data Flow Analysis is Model Checking of Abstract Interpretations." POPL98.
 - State space: Program execution tree
 - Each conditional branch is a fork in the tree
- Consider very-busy expressions:

$$VBE(p) = Used(p) \cup (\text{notMod}(p) \cap (\bigcap_{p' \in \text{succ}(p)} VBE(p')))$$
- Reformatted as model checking the exec. space:

$$\text{isVBE}(e) = \nu Z. \text{isUsed}(e) \vee (\text{notMod}(e) \wedge \square Z)$$
 (here ν is the greatest fixpoint operator)

7

Model Checking and Theorem Proving

- Model checkers are fully automated theorem provers
 - Again, they prove "dumb" theorems
 - But somewhat smarter than type systems
 - E.g., they handle concurrency, complicated properties
 - But don't do a good job with complex structures
 - E.g., functions, data structures

8

Model Checking and Type Systems

- Naik and Palsberg, "A type system equivalent to a model checker"
 - Shows how to construct a type system that accepts exactly the set of programs that a model checker passes

9

Denotational Semantics

- What mathematical structures do programs represent?
 - How do we reason *compositionally* about programs?

State : Variables \rightarrow Values

$\llbracket \cdot \rrbracket$: Statement \rightarrow (State \rightarrow State)

$\llbracket \text{skip} \rrbracket$: $\lambda s. s$

$\llbracket x=e \rrbracket$: $\lambda s. s[\nu \setminus x]$ where $\nu = \llbracket e \rrbracket s$

$\llbracket \text{if } B \text{ then } C \text{ else } C' \rrbracket$: $\lambda s. \text{if}(\llbracket B \rrbracket s, \llbracket C \rrbracket s, \llbracket C' \rrbracket s)$ where
 if(true, ν, ν') = ν
 if(false, ν, ν') = ν'

10

Loops in Denotational Semantics

- Loops are tricky:
 - Want $\llbracket \text{while } B \text{ do } C \rrbracket$ to be defined in terms of B and C
 - $\llbracket \text{while } B \text{ do } C \rrbracket = \lambda s. s$ if $\llbracket B \rrbracket s = \text{false}$
 - $\llbracket \text{while } B \text{ do } C \rrbracket = \llbracket C; \text{while } B \text{ do } C \rrbracket$ if $\llbracket B \rrbracket s = \text{true}$
 - But that's not compositional reasoning!
 - `while` is defined in terms of itself
- Solution: Need to compute a fixpoint
 - Define *domains* on which minimal fixpoints exist

11

Complete Partial Orders

- A partial order (P, \sqsubseteq) is a set P and a reflexive, transitive, antisymmetric binary relation \sqsubseteq
- A partial order *has a bottom* if it has a least element \perp
- An ω -chain is infinite increasing sequence
 - $x_0 \sqsubseteq x_1 \sqsubseteq x_2 \sqsubseteq \dots$
- A partial order is *complete* (a "cpo") if every ω -chain has a least upper bound
 - Written $\sqcup \{x_i \mid i \in \omega\}$

(Following Abadi, CS263)

12

Continuous Functions

- Let P_1 and P_2 be two complete partial orders
- A function $f : P_1 \rightarrow P_2$ is *continuous* if
 - It is monotonic
 - $x \sqsubseteq y \Rightarrow f(x) \sqsubseteq f(y)$
 - For all ω -chains
 - $\sqcup_{i \in \omega} f(x_i) = f(\sqcup_{i \in \omega} x_i)$

13

Fixed-Point Theorem

- Let P be a cpo with bottom
- Let $f : P \rightarrow P$ be a continuous function
- Let $f^i(x) = f(f(\dots f(x)))$ (i times)
- Define $\text{fix}(f) = \sqcup_{i \in \omega} f^i(\perp)$
- Then $\text{fix}(f)$ is the *least fixed point* of f
 - $f(\text{fix}(f)) = \text{fix}(f)$
 - if $f(x) = x$ then $\text{fix}(f) \sqsubseteq x$

14

Proof: First step

- Claim: $\perp \sqsubseteq f(\perp) \sqsubseteq f(f(\perp)) \sqsubseteq \dots$
 - I.e., $f^i(\perp)$ forms an ω -chain
- Proof:
 - $\perp \sqsubseteq f(\perp)$ definition of \perp
 - $f(\perp) \sqsubseteq f(f(\perp))$ monotonicity
 - $f(f(\perp)) \sqsubseteq f(f(f(\perp)))$ monotonicity
 - ...

15

Proof: f is a Fixpoint

$$\begin{aligned}
 f(\text{fix}(f)) &= f(\sqcup_{i \in \omega} f^i(\perp)) && \text{by definition} \\
 &= \sqcup_{i \in \omega} (f(f^i(\perp))) && \text{by continuity} \\
 &= \sqcup_{i \in \omega} f^{i+1}(\perp) \\
 &= (\sqcup_{i \in \omega} f^{i+1}(\perp)) \sqcup \perp \\
 &= \sqcup_{i \in \omega} f^i(\perp) \\
 &= \text{fix}(f) && \text{by definition}
 \end{aligned}$$

16

Proof: f is the Least Fixed Point

- Let y be another fixed point, $f(y) = y$
- Then
 - $\perp \sqsubseteq y$ by definition of bottom
 - $f(\perp) \sqsubseteq f(y) = y$ by monotonicity
 - $f(f(\perp)) \sqsubseteq f(y) = y$ by monotonicity
 - $\dots f^i(\perp) \sqsubseteq y$ for all $i \in \omega$
 - $\sqcup (f^i(\perp)) \sqsubseteq y$
 - $\text{fix}(f) \sqsubseteq y$

17

A Useful CPO

- Let F be the set of functions $\text{State} \rightarrow (\text{State} \cup \perp)$
- Define $f \sqsubseteq g$ if $f(x) = g(x)$ or $f(x) = \perp$
- Then F is a CPO with bottom

18

Denotational Semantics of While

- Goal: $\llbracket \text{while } B \text{ do } C \rrbracket$ defined in terms of B and C
- Let $G = \lambda f. \lambda s. \text{if } \llbracket B \rrbracket(s) \text{ then } f(\llbracket C \rrbracket(s)) \text{ else } s$
 - G “unrolls” one iteration of the loop, using f for the recursive call
 - Notice $G : F \rightarrow F$ and G continuous
- Define $\llbracket \text{while } B \text{ do } C \rrbracket = \text{fix}(G)$
- Then $\text{fix}(G) = G(\text{fix}(G)) = \lambda s. \text{if } \llbracket B \rrbracket(s) \text{ then } \text{fix}(G)(\llbracket C \rrbracket(s)) \text{ else } s$
 - $\text{fix}(G)$ is the least function with this property

19

Denotational Semantics

- A very compelling theory
 - Composition reasoning very powerful
 - Requires a lot of math
 - Makes some proofs easier
- Today, operational semantics mostly used
 - A lot simpler to understand
 - Reduces to a lot of symbol pushing
 - But hard to reuse results

20

Language-Based Security

- Writing secure software is hard
 - Adversary is malicious: looking for bugs
 - Hard to test for security flaws
 - Often errors on non-covered paths
- Not many mechanisms in languages for security
 - Type and memory safety help (e.g., don't use C)
 - One exception: Stack inspection in Java
 - But what does it mean? What security can it achieve?

21

Secure Information Flow

- A popular notion of security: *non-interference*
 - Idea: Program is a function $H \times L \rightarrow H' \times L'$
 - H = high security, L = low security
 - High-security inputs should not leak to low-security outputs
 - Leaving L fixed and changing only H should not change L'
 - Is this a safety property? A liveness property?
 - What evidence shows this property is violated?

22

Enforcing Non-Interference

- Types distinguish high- and low-security data
 - Guarantee H never flows to L
 - Dual of *tainted/untainted* type qualifiers
- But wait! What about the following:
 $\text{if } (H) \text{ then } L := 1 \text{ else } L := 0$
 - No direct flow from H to L
 - This is a *covert channel*
 - Need to make PC high-security in this case

23

Enforcing Non-Interference (cont'd)

- But wait! What about multi-threaded code?
 $\text{if } (H) \text{ then } \langle \text{do a lot of work} \rangle \text{ else } \text{sleep}(100)$
 - Other process may observe schedule to find H
 - Need to make sides of conditional take equal time
- But wait! What if we're supposed to leak info?
 $\text{if } (\text{passwd matches}) \text{ then } \text{log-in} \text{ else } \text{fail}$
 - Need some way to *declassify* information
 - In fact, this is the key to making this all work
 - The jury is still out on whether any of this is practical

24

Multi-Threaded Programming

- Writing multi-threaded programs is hard
 - The scheduler can interleave threads unpredictably
 - Means that it's hard to understand all the possible behaviors of your program
 - And if you do find a bug, it's hard to reproduce
 - Multi-threaded programs trade off safety and liveness
 - A safe program won't have (harmful) race conditions
 - Typically use mutual exclusion locks to enforce
 - But locking forces threads to block
 - A live program will make progress
 - Reducing the amount of locking improves liveness

25

Mistakes Have Consequences

- A data race in the Therac-25, a radiation therapy machine, gave patients massive overdoses of radiation, killing at least five people
- A data race was partially responsible for the northeastern US blackout of August 14, 2004, one of the worst in North American history

26

Parallelism is Becoming the Norm

- Desktop machines with >1 CPU are not very expensive
- Chip-Level Multiprocessors are increasingly common
 - Intel, AMD, IBM, Sun are building/will build them
- An effort to keep up with Moore's law

27

How Do We Program These Things?

- One idea: *atomic sections*

```
atomic { s } /* execute s atomically */
```

 - Up to the programming language to decide how to implement this
- Ideas for implementing this
 - Hardware transactional memory
 - Software transactional memory
 - Typically both of these will be optimistic concurrency
 - Mutual exclusion
 - Inferred by some static analysis

28