



ASTs

ASTs are abstract

- They don't contain all information in the program
 E.g., spacing, comments, brackets, parentheses
- Any ambiguity has been resolved
 - E.g., a + b + c produces the same AST as (a + b) + c

Disadvantages of ASTs

- AST has many similar forms
- E.g., for, while, repeat...until
- ∎ E.g., if, ?:, switch
- Expressions in AST may be complex, nested • (42 * y) + (z > 5 ? 12 * z : z + 20)
- Want simpler representation for analysis
 ...at least, for dataflow analysis

Control-Flow Graph (CFG)

- · A directed graph where
 - Each node represents a statement
 - Edges represent control flow
- · Statements may be
 - Assignments x := y op z or x := op z
 - Copy statements x := y
 - Branches goto L or if x relop y goto L
 - ∎ etc.



Variations on CFGs

- We usually don't include declarations (e.g., int x;)
 - But there's usually something in the implementation
- May want a unique entry and exit node
 Won't matter for the examples we give
- May group statements into basic blocks
 - A sequence of instructions with no branches into or out of the block



CFG vs. AST

- CFGs are much simpler than ASTs
 - Fewer forms, less redundancy, only simple expressions
- But...AST is a more faithful representation
 - CFGs introduce temporaries
 - Lose block structure of program
- So for AST,
 - Easier to report error + other messages
 - Easier to explain to programmer
 - Easier to unparse to produce readable code

Data Flow Analysis

- A framework for proving facts about programs
- · Reasons about lots of little facts
- · Little or no interaction between facts
 - Works best on properties about how program computes
- Based on all paths through program
 Including infeasible paths

Available Expressions

- An expression e is available at program point p if
 - e is computed on every path to p, and
 - the value of e has not changed since the last time e is computed on p
- Optimization
 - If an expression is available, need not be recomputed
 - (At least, if it's still in a register somewhere)







Terminology

- A *joint point* is a program point where two branches meet
- Available expressions is a *forward must* problem
 - Forward = Data flow from in to out
 - Must = At join point, property must hold on all paths that are joined



Let s be a statement

- succ(s) = { immediate successor statements of s }
- pred(s) = { immediate predecessor statements of s}
- In(s) = program point just before executing s
- Out(s) = program point just after executing s
- $ln(s) = \bigcap_{s' \in pred(s)} Out(s')$ • $Out(s) = Gen(s) \cup (ln(s) - Kill(s))$
 - Note: These are also called *transfer functions*

Liveness Analysis

- A variable v is *live* at program point p if
 - ${\scriptstyle \bullet v}$ will be used on some execution path originating from p...
 - before v is overwritten
- Optimization
 - If a variable is not live, no need to keep it in a register
 - If variable is dead at assignment, can eliminate assignment

Data Flow Equations

- Available expressions is a forward must analysis
 - Data flow propagate in same dir as CFG edges
- Expr is available only if available on all paths
- Liveness is a backward may problem
- To know if variable live, need to look at future uses
 Variable is live if used on some path
- Out(s) = $\mathbf{U}_{s' \in \text{succ}(s)} \ln(s')$
- $ln(s) = Gen(s) \cup (Out(s) Kill(s))$













Partial Orders

- A partial order is a pair (P, \leq) such that
 - $\bullet \leq \subseteq P \times P$
 - \leq is reflexive: $x \leq x$
 - _ \leq is anti-symmetric: $x \leq y$ and $y \leq x \Rightarrow x = y$
 - _ < is transitive: $x \leq y$ and $y \leq z \Rightarrow x \leq z$

Lattices

- A partial order is a lattice if □ and □ are defined on any set:
 I is the *meet* or *greatest lower bound* operation:
 x □ y ≤ x and x □ y ≤ y
 if z ≤ x and z ≤ y, then z ≤ x □ y
- Lis the join or least upper bound operation:
 - $x \le x \sqcup y \text{ and } y \le x \sqcup y$ if $x \le z$ and $y \le z$, then $x \sqcup y \le z$







■temp := Gen(s) ∪ (In(s) - Kill(s)) ■if (temp != Out(s)) { -Out(s) := temp

-W := Ŵ ∪ succ(s) ■} •until W = ∅



Termination

- · We know the algorithm terminates because
 - The lattice has finite height
 - The operations to compute In and Out are monotonic
 - On every iteration, we remove a statement from the worklist and/or move down the lattice

Forward Data Flow, Again

•Out(s) = Top for all statements s •W := { all statements } (worklist) •repeat ■Take s from W ■temp := f_s(¬_{s' € pred(s)} Out(s')) (f_s monotonic *transfer fn*) ■if (temp != Out(s)) { -Out(s) := temp -W := W ∪ succ(s) ■} •until W = Ø

Lattices (P, ≤)

- Available expressions
 - P = sets of expressions
 - ∎ S1 ⊓ S2 = S1 ∩ S2
 - Top = set of all expressions
- Reaching Definitions
 - P = set of definitions (assignment statements)
 - ∎ S1 ⊓ S2 = S1 U S2
 - Top = empty set

Fixpoints

- ·We always start with Top
 - Every expression is available, no defns reach this point
 - Most optimistic assumption
 - Strongest possible hypothesis
 = true of fewest number of states
- Revise as we encounter contradictions
- Always move down in the lattice (with meet)
- Result: A greatest fixpoint

Lattices (P, ≤), cont'd

- Live variables
 - P = sets of variables
 - ∎S1 □ S2 = S1 ∪ S2
 - Top = empty set
- Very busy expressions
 - P = set of expressions
 - $\blacksquare \texttt{S1} \sqcap \texttt{S2} = \texttt{S1} \cap \texttt{S2}$
 - Top = set of all expressions



Termination Revisited• How many times can we apply this step:• $temp := f_s(\sqcap_{s' \in pred(s)} Out(s'))$ • if (temp != Out(s)) { ... }• Claim: Out(s) only shrinks• Proof: Out(s) starts out as top• So temp must be ≤ than Top after first step• Assume Out(s') shrinks for all predecessors s' of s• Then $\sqcap_{s' \in pred(s)} Out(s')$ shrinks• Since f monotonic, f ($\sqcap_{s' \in pred(s)} Out(s')$) shrinks

Termination Revisited (cont'd)

- A descending chain in a lattice is a sequence • $x0 \supseteq x1 \supseteq x2 \supseteq ...$
- The *height* of a lattice is the length of the longest descending chain in the lattice
- Then, dataflow must terminate in O(n k) time
 - \blacksquare **n** = # of statements in program
 - k = height of lattice
 - assumes meet operation takes O(1) time

Relationship to Section 2.4 of Book (NNH)

- MFP (Maximal Fixed Point) solution general iterative algorithm for monotone frameworks
 - always terminates
 - always computes the right solution

Least vs. Greatest Fixpoints

- Dataflow tradition: Start with Top, use meet
 - $\hfill\blacksquare$ To do this, we need a meet semilattice with top
- meet semilattice = meets defined for any set
- Computes greatest fixpoint
- Denotational semantics tradition: Start with Bottom, use join
- Computes least fixpoint

Distributive Data Flow Problems

• By monotonicity, we also have $f(x \sqcap y) \leq f(x) \sqcap f(y)$

• A function f is distributive if

 $f(x\sqcap y)=f(x)\sqcap f(y)$



Accuracy of Data Flow Analysis

- Ideally, we would like to compute the meet over all paths (MOP) solution:
 - Let f be the transfer function for statement s
 - If p is a path $\{s_1, ..., s_n\}$, let $f_n = f_n; ...; f_1$
 - Let path(s) be the set of paths from the entry to s

$\mathrm{MOP}(s) = \sqcap_{p \in \mathrm{path}(s)} f_p(\top)$

• If a data flow problem is distributive, then solving the data flow equations in the standard way yields the MOP solution, i.e., MFP = MOP

What Problems are Distributive?

- Analyses of *how* the program computes
 - Live variables
 - Available expressions
 - Reaching definitions
 - Very busy expressions
- All Gen/Kill problems are distributive



MOP vs MFP

- Computing MFP is always safe: MFP \sqsubseteq MOP
- When distributive: MOP = MFP
- When non-distributive: MOP may not be computable (decidable)
 - e.g., MOP for constant propagation (see Lemma 2.31 of NNH)

Practical Implementation

- Data flow facts = assertions that are true or false at a program point
- Represent set of facts as bit vector
 - Fact, represented by bit i
 - Intersection = bitwise and, union = bitwise or, etc
- "Only" a constant factor speedup
 But very useful in practice

Basic Blocks

- A basic block is a sequence of statements s.t.
 - No statement except the last in a branch
 - There are no branches to any statement in the block except the first
- · In practical data flow implementations,
 - Compute Gen/Kill for each basic block
 Compose transfer functions
 - Store only In/Out for each basic block
 - Typical basic block ~5 statements

Order Matters

- Assume forward data flow problem
 Let G = (V, E) be the CFG
 - Let k be the height of the lattice
- If G acyclic, visit in topological order
 Visit head before tail of edge

Running time O(|E|)

No matter what size the lattice

Order Matters — Cycles

- If G has cycles, visit in reverse postorder
 - Order from depth-first search
- Let Q = max # back edges on cycle-free path
 Nesting depth
 - Back edge is from node to ancestor on DFS tree
- Then if ∀x.f(x) ≤ x (sufficient, but not necessary)
 Running time is O((Q+1)|E|)
 - Note direction of req't depends on top vs. bottom

Flow-Sensitivity

- Data flow analysis is flow-sensitive
 - The order of statements is taken into account
 - I.e., we keep track of facts per program point
- Alternative: Flow-insensitive analysis
 - Analysis the same regardless of statement order
 - Standard example: types
 - /* x : int */ x := ... /* x : int */

Terminology Review

- Must vs. May
- (Not always followed in literature)
- Forwards vs. Backwards
- Flow-sensitive vs. Flow-insensitive
- Distributive vs. Non-distributive

Another Approach: Elimination

- Recall in practice, one transfer function per basic block
- Why not generalize this idea beyond a basic block?
 - "Collapse" larger constructs into smaller ones, combining data flow equations
 - Eventually program collapsed into a single node!
 - "Expand out" back to original constructs, rebuilding information

Lattices of Functions

- Let (P, ≤) be a lattice
- Let M be the set of monotonic functions on P
- Define $f \leq_{\epsilon} g$ if for all x, $f(x) \leq g(x)$
- Define the function $f \sqcap g$ as
 - $\blacksquare (f \sqcap g) (x) = f(x) \sqcap g(x)$
- Claim: (M, \leq_{f}) forms a lattice







Height of Function Lattice

- Assume underlying lattice (P, ≤) has finite height
 - What is height of lattice of monotonic functions?
 - Claim: finite
- Therefore, g(j) converges

Non-Reducible Flow Graphs

- Elimination methods usually only applied to *reducible* flow graphs
 - Ones that can be collapsed
 - Standard constructs yield only reducible flow graphs
- Unrestricted goto can yield non-reducible graphs

Comments

- Can also do backwards elimination
 - Not quite as nice (regions are usually single *entry* but often not single *exit*)
- For bit-vector problems, elimination efficient
 - Easy to compose functions, compute meet, etc.
- Elimination originally seemed like it might be faster than iteration
 - Not really the case

Data Flow Analysis and Functions

- What happens at a function call?
- Lots of proposed solutions in data flow analysis literature
- In practice, only analyze one procedure at a time
- Consequences
 - Call to function kills all data flow facts
 - May be able to improve depending on language, e.g., function call may not affect locals

More Terminology

- An analysis that models only a single function at a time is *intraprocedural*
- An analysis that takes multiple functions into account is *interprocedural*
- An analysis that takes the whole program into account is...guess?
- Note: *global* analysis means "more than one basic block," but still within a function

Data Flow Analysis and The Heap

• Data Flow is good at analyzing local variables

- But what about values stored in the heap?
- Not modeled in traditional data flow
- In practice: *x := e
 - Assume all data flow facts killed (!)
 - Or, assume write through x may affect any variable whose address has been taken
- In general, hard to analyze pointers

Data Flow Analysis and Optimization

- Moore's Law: Hardware advances double computing power every 18 months.
- Proebsting's Law: Compiler advances double computing power every 18 *years*.