Lecture 22
CIS 341: COMPILERS

Announcements

- HW 6: Dataflow Analysis and Optimizations
 - Available later today(?)
 - Due Next Thursday, April 16
- HW 7: Optimization & Experiments
 - Available next week
 - Due: April 29th

QUALITY OF DATAFLOW ANALYSIS SOLUTIONS

Best Possible Solution

- Suppose we have a control-flow graph.
- If there is a path p₁ starting from the root node (entry point of the function) traversing the nodes n₀, n₁, n₂, ... n_k
- The best possible information along the path p_1 is: $l_{p1} = F_{nk}(...F_{n2}(F_{n1}(F_{n0}(T)))...)$
- Best solution at the output is some $\mathfrak{l} \subseteq \mathfrak{l}_p$ for *all* paths p.
- Meet-over-paths (MOP) solution:

 $\prod_{p \in paths_{to[n]}} l_p$



What about quality of iterative solution?

- Does the iterative solution: $out[n] = F_n(\prod_{n' \in pred[n]} out[n'])$ compute the MOP solution?
- MOP Solution: $\prod_{p \in paths_{to[n]}} l_p$
- Answer: Yes, *if* the flow functions *distribute* over \square
 - Distributive means: $\prod_{i} F_{n}(\mathbf{l}_{i}) = F_{n}(\prod_{i} \mathbf{l}_{i})$
 - Proof is a bit tricky & beyond the scope of this class. (Difficulty: loops in the control flow graph might mean there are infinitely many paths...)
- Not all analyses give MOP solution
 - They are more conservative.

Reaching Definitions is MOP

- $F_n[x] = gen[n] \cup (x kill[n])$
- Does F_n distribute over meet $\square = U$?
- $F_n[x \sqcap y]$
 - = gen[n] U ((x U y) kill[n])
 - $= gen[n] \cup ((x kill[n]) \cup (y kill[n]))$
 - = (gen[n] U(x kill[n])) U (gen[n] U(y kill[n]))
 - $= F_n[x] \cup F_n[y]$
 - $= F_n[x] \prod F_n[y]$
- Therefore: Reaching Definitions with iterative analysis always terminates with the MOP (i.e. best) solution.

"Classic" Constant Propagation

- Constant propagation can be formulated as a dataflow analysis.
- Idea: propagate and fold integer constants in one pass:

$$x = 1;$$
 $x = 1;$
 $y = 5 + x;$ $y = 6;$
 $z = y * y;$ $z = 36;$

- Information about a single variable:
 - Variable is never defined.
 - Variable has a single, constant value.
 - Variable is assigned multiple values.

Domains for Constant Propagation

• We can make a constant propagation lattice \mathcal{L} for *one variable* like this:



- To accommodate multiple variables, we take the product lattice, with one element per variable.
 - Assuming there are three variables, x, y, and z, the elements of the product lattice are of the form (l_x, l_y, l_z) .
 - Alternatively, think of the product domain as a context that maps variable names to their "*abstract interpretations*"
- What are "meet" and "join" in this product lattice?
- What is the height of the product lattice?

Flow Functions

- Consider the node x = y op z ۲
- $F(\boldsymbol{\varrho}_{x'} \boldsymbol{\varrho}_{y'} \boldsymbol{\varrho}_{z}) = ?$
- $F(\boldsymbol{\ell}_{x'}, \boldsymbol{\ell}_{y'}, \top) = (\top, \boldsymbol{\ell}_{y'}, \top)$
- $F(l_x, T, l_z) = (T, T, l_z)$ If either input might have multiple values the result of the operation might too."

- F(𝒫_{x'} ⊥, 𝒫_z) = (⊥, ⊥, 𝒫_z)
 F(𝒫_{x'} 𝒫_{y'} ⊥) = (⊥, 𝒫_{y'} ⊥)
 "If either input is undefined the result of the operation is the result of the operation is too."
- $F(l_{x'}, i, j) = (i \text{ op } j, i, j) \int$ "If the inputs are known constants, calculate the output statically."
- Flow functions for the other nodes are easy... •
- Monotonic? ۲
- Distributes over meets? ۲

Iterative Solution



MOP Solution \neq **Iterative Solution**



Why not compute MOP Solution?

- If MOP is better than the iterative analysis, why not compute it instead?
 - ANS: exponentially many paths (even in graph without loops)
- O(n) nodes
- O(n) edges
- O(2ⁿ) paths*
 - At each branch there is a choice of 2 directions

* Incidentally, a similar idea can be used to force ML / Haskell type inference to need to construct a type that is exponentially big in the size of the program!

12

Dataflow Analysis: Summary

- Many dataflow analyses fit into a common framework.
- Key idea: *Iterative solution* of a system of equations over a *lattice* of constraints.
 - Iteration terminates if flow functions are monotonic.
 - Solution is equivalent to meet-over-paths answer if the flow functions distribute over meet (□).
- Dataflow analyses as presented work for an "imperative" intermediate representation.
 - The values of temporary variables are updated ("mutated") during evaluation.
 - Such mutation complicates calculations
 - SSA = "Single Static Assignment" eliminates this problem, by introducing more temporaries – each one assigned to only once.
 - Next up: Converting to SSA, finding loops and dominators in CFGs

LOOPS AND DOMINATORS

Zdancewic CIS 341: Compilers

Loops in Control-flow Graphs

- Taking into account loops is important for optimizations.
 - The 90/10 rule applies, so optimizing loop bodies is important
- Should we apply loop optimizations at the AST level or at a lower representation?
 - Loop optimizations benefit from other IR-level optimizations and vice-versa, so it is good to interleave them.
- Loops may be hard to recognize at the quadruple / LLVM IR level.
 - Many kinds of loops: while, do/while, for, continue, goto...
- Problem: *How do we identify loops in the control-flow graph?*

Definition of a Loop

- A *loop* is a set of nodes in the control flow graph.
 - One distinguished entry point called the *header*
- Every node is reachable from the header & the header is reachable from every node.
 - A loop is a strongly connected component
- No edges enter the loop except to the header
- Nodes with outgoing edges are called loop exit nodes



Nested Loops

- Control-flow graphs may contain many loops
- Loops may contain other loops:





The control tree depicts the nesting structure of the program.

Control-flow Analysis

- Goal: Identify the loops and nesting structure of the CFG.
- Control flow analysis is based on the idea of *dominators*:
- Node A *dominates* node B if the only way to reach B from the start node is through node A.
- An edge in the graph is a *back edge* if the target node dominates the source node.
- A loop contains at least one back edge.



Dominator Trees

- Domination is transitive:
 - if A dominates B and B dominates C then A dominates C
- Domination is anti-symmetric:
 - if A dominates B and B dominates A then A = B
- Every flow graph has a dominator tree
 - The Hasse diagram of the dominates relation



Dominator Dataflow Analysis

- We can define Dom[n] as a forward dataflow analysis.
 - Using the framework we saw earlier: Dom[n] = out[n] where:
- "A node B is dominated by another node A if A dominates *all* of the predecessors of B."
 - in[n] := $\bigcap_{n' \in pred[n]} out[n']$
- "Every node dominates itself."
 - $out[n] := in[n] \cup \{n\}$
- Formally: $\mathcal{L} = \text{set of nodes ordered by } \subseteq$
 - $T = \{all nodes\}$
 - $\ F_n(x) = x \ U \ \{n\}$
 - \square is \cap
- Easy to show monotonicity and that F_n distributes over meet.
 - So algorithm terminates and is MOP

Improving the Algorithm

- Dom[b] contains just those nodes along the path in the dominator tree from the root to b:
 - e.g. $Dom[8] = \{1, 2, 4, 8\}, Dom[7] = \{1, 2, 4, 5, 7\}$
 - There is a lot of sharing among the nodes
- More efficient way to represent Dom sets is to store the dominator *tree*.
 - doms[b] = immediate dominator of b
 - doms[8] = 4, doms[7] = 5
- To compute Dom[b] walk through doms[b]
- Need to efficiently compute intersections of Dom[a] and Dom[b]
 - Traverse up tree, looking for least common ancestor:
 - Dom[8] \cap Dom[7] = Dom[4]



• See: "A Simple, Fast Dominance Algorithm" Cooper, Harvey, and Kennedy

Completing Control-flow Analysis

- Dominator analysis identifies *back edges*:
 - Edge n \rightarrow h where h dominates n
- Each back edge has a *natural loop*:
 - h is the header
 - All nodes reachable from h that also reach n without going through h
- For each back edge $n \rightarrow h$, find the natural loop:
 - $\{n' \mid n \text{ is reachable from } n' \text{ in } G \{h\}\} \cup \{h\}$
- Two loops may share the same header: merge them
- Nesting structure of loops is determined by set inclusion
 - Can be used to build the control tree





Example Natural Loops

Control Tree:

The control tree depicts the nesting structure of the program.

Natural Loops

Uses of Control-flow Information

- Loop nesting depth plays an important role in optimization heuristics.
 - Deeply nested loops pay off the most for optimization.
- Need to know loop headers / back edges for doing
 - loop invariant code motion
 - loop unrolling
- Dominance information also plays a role in converting to SSA form
 - Used internally by LLVM to do register allocation.

Phi nodes Alloc "promotion" Register allocation

REVISITING SSA

Single Static Assignment (SSA)

- LLVM IR names (via **%uids**) *all* intermediate values computed by the program.
- It makes the order of evaluation explicit.
- Each **%uid** is assigned to only once
 - Contrast with the mutable quadruple form
 - Note that dataflow analyses had these kill[n] sets because of updates to variables...
- Naïve implementation of backend: map **%uids** to stack slots
- Better implementation: map **%uids** to registers (as much as possible)
- Question: How do we convert a source program to make maximal use of **%uids**, rather than alloca-created storage?
 - two problems: control flow & location in memory
- Then: How do we convert SSA code to x86, mapping **%uids** to registers?
 - Register allocation.

Alloca vs. %UID

• Current compilation strategy:

• Directly map source variables into **%uids**?

• Does this always work?

What about If-then-else?

• How do we translate this into SSA?

entry: %y1 = ... %x1 = ... %z1 = ... %p = icmp ... br i1 %p, label %then, label %else then: $x^2 = add i^3 y_1, 1$ br label %merge else: %x3 = mult i32 %y1, 2 merge: %z2 = %add i32 ???, 3

• What do we put for ???

Phi Functions

- Solution: φ functions
 - Fictitious operator, used only for analysis
 - implemented by Mov at x86 level
 - Chooses among different versions of a variable based on the path by which control enters the phi node.

 $\texttt{suid} \texttt{=} \texttt{phi} < \texttt{ty} > \texttt{v}_1, < \texttt{label}_1 >, \dots, \texttt{v}_n, < \texttt{label}_n >$

Phi Nodes and Loops

- Importantly, the **%uids** on the right-hand side of a phi node can be defined "later" in the control-flow graph.
 - Means that **%uids** can hold values "around a loop"
 - Scope of **%uids** is defined by dominance (discussed soon)

```
entry:
  %y1 = ...
  %x1 = ...
  br label %body
body:
  %x2 = phi i32 %x1, %entry, %x3, %body
  %x3 = add i32 %x2, %y1
  %p = icmp slt i32, %x3, 10
  br i1 %p, label %body, label %after
after:
  ...
```