CS4212: Compiler Design

Week 12: Code Optimizations and Dataflow Analysis

Ilya Sergey

ilya@nus.edu.sg

ilyasergey.net/CS4212/

Optimizations

- - Lots of redundant moves.
 - Lots of unnecessary arithmetic instructions.
- Consider this OAT program: lacksquare

```
int foo(int w) {
 var x = 3 + 5;
 var y = x * w;
                  frontend.ml
 var z = y - 0;
 return z * 4;
```

The code generated by our Oat compiler so far is pretty inefficient.



Optimized vs Non-Optimized Output

```
define i64 @foo(i64 %_w1) {
 %_w2 = alloca i64
 %_x5 = alloca i64
 %_y10 = alloca i64
 %_z14 = alloca i64
  store i64 %_w1, i64* %_w2
 %_bop4 = add i64 3, 5
  store i64 %_bop4, i64* %_x5
 %_x7 = load i64, i64* %_x5
 %_w8 = load i64, i64* %_w2
 %_bop9 = mul i64 %_x7, %_w8
  store i64 %_bop9, i64* %_y10
 %_y12 = load i64, i64* %_y10
 %_bop13 = sub i64 %_y12, 0
  store i64 %_bop13, i64* %_z14
 %_z16 = load i64, i64* %_z14
 %_bop17 = mul i64 %_z16, 4
  ret i64 %_bop17
```



%rbp %rsp, %rbp \$136, %rsp %rdi, %rax %rax, -8(%rbp) \$0 %rsp, -16(%rbp) \$0 %rsp, -24(%rbp) \$0 %rsp, -32(%rbp) %rsp, -40(%rbp) -8(%rbp), %rcx -16(%rbp), %rax %rcx, (%rax) \$3, %rax \$5, %rcx %rcx, %rax %rax, -56(%rbp) -56(%rbp), %rcx -24(%rbp), %rax %rcx, (%rax) -24(%rbp), %rax (%rax), %rcx %rcx, -72(%rbp) -16(%rbp), %rax (%rax), %rcx %rcx, -80(%rbp) -72(%rbp), %rax -80(%rbp), %rcx %rcx. %rax %rax, -88(%rbp) -88(%rbp), %rcx -32(%rbp), %rax %rcx, (%rax) -32(%rbp), %rax (%rax), %rcx %rcx, -104(%rbp) -104(%rbp), %rax \$0, %rcx %rcx, %rax %rax, -112(%rbp) -112(%rbp), %rcx -40(%rbp), %rax %rcx, (%rax) -40(%rbp), %rax (%rax), %rcx %rcx, -128(%rbp) -128(%rbp), %rax \$4, %rcx %rcx, %rax %rax, -136(%rbp) -136(%rbp), %rax %rbp, %rsp %rbp

.text .globl

pushq

mova

subq

movq

movq

pusho

movq

pushq

movq

pushq

movq

pushq

movq

movq

movq

movq

movq

movq

addq

movq

movq

movq movq

movq

movq

movq

movq

movq

movq

movq

movq

imulq

movq

movq

movq

movq

movq movq

movq

movq

movq subq

movq

movq

movq

movo

movq mova

movq

movq

movq

imula

movq

movq

movq

popq retq

foo:

_foo



- Code above generated by clang –O3
- Function foo may be inlined by the compiler, so it can be implemented by just one instruction!



Why do we need optimizations?

- To help programmers...
 - They write modular, clean, high-level programs
 - Compiler generates efficient, high-performance assembly _____
- Programmers don't write optimal code
- lacksquare- e.g. A[i][j] = A[i][j] + 1
- Architectural independence •
 - Optimal code depends on features not expressed to the programmer
 - Modern architectures *assume* optimization
- Different kinds of optimizations: \bullet
 - Time: improve execution speed
 - Space: reduce amount of memory needed
 - Power: lower power consumption (e.g. to extend battery life)

High-level languages make avoiding redundant computation inconvenient or impossible

Some Caveats

- Optimizations are code transformations:
 - They can be applied at any stage of the compiler
 - They must be *safe* (?)
 - they shouldn't change the meaning of the program. _____
- In general, optimizations require some program analysis:
 - To determine if the transformation really is safe
 - To determine whether the transformation is cost effective
- This course: most common and valuable performance optimizations
 - See Muchnick (optional text) for ~ 10 chapters about optimizations _____

Advanced COMPILER DESIGN IMPLEMENTATION

Steven S. Muchnick







When to apply optimization

- Inlining
- Function specialization
- Constant folding
- Constant propagation
- Value numbering
- Dead code elimination
- Loop-invariant code motion
- Common sub-expression elimination
- Strength Reduction
- Constant folding & propagation
- Branch prediction / optimization
- Register allocation
- Loop unrolling
- Cache optimization

Where to Optimize?

- Usual goal: improve time performance
- Problem: many optimizations trade space for time \bullet
- Example: Loop unrolling \bullet
 - Idea: rewrite a loop like (why?): for(int i=0; i<100; i=i+1) { s = s + a[i];– Into a loop like: for(int i=0; i<99; i=i+2){ s = s + a[i];s = s + a[i+1];
- **Tradeoffs:**
 - Increasing code space slows down whole program a tiny bit (extra instructions to manage) but speeds up the loop a lot
 - For frequently executed code with long loops: generally a win
 - Interacts with instruction cache and branch prediction hardware
- Complex optimizations may never pay off!

Writing Fast Programs In Practice

- Pick the right algorithms and data structures. ullet
 - These have a much bigger impact on performance that compiler optimizations.
 - Reduce # of operations
 - Reduce memory accesses _____
 - Minimize indirection
- Then turn on compiler optimizations
- Profile to determine program hot spots \bullet
- Evaluate whether the algorithm/data structure design works \bullet

... if so: "tweak" the source code until the optimizer does "the right thing" to the machine code

- - have more ambiguity in their behaviour.
- Example: loop-invariant code motion
 - Idea: hoist invariant code out of a loop

- Is this more efficient?
- Is this safe?



Whether an optimization is *safe* depends on the programming language semantics. Languages that provide weaker guarantees to the programmer permit more optimizations but

– e.g. In C, loading from initialized memory is undefined, so the compiler can do anything.





Constant Folding

•

int $x = (2 + 3) * y \rightarrow int x = 5 * y$ b & false \rightarrow false

- Performed at every stage of optimization... Why?
- Constant expressions can be created by translation or earlier optimizations ullet

Example: A[2] might be compiled to: $MEM[MEM[A] + 2 * 4] \rightarrow MEM[MEM[A] + 8]$

Idea: If operands are known at compile time, perform the operation statically.

Constant Folding Conditionals

if (true) S \rightarrow Sif (false) S \rightarrow ;if (true) S else S' \rightarrow Sif (false) S else S' \rightarrow Swhile (false) S \rightarrow ;if (2 > 3) S \rightarrow ;



Algebraic Simplification

- More general form of constant folding - Take advantage of mathematically sound simplification rules
- Identities: ullet
 - $a^* 1 \rightarrow a \qquad a^* 0 \rightarrow 0$
 - $-a+0 \rightarrow a \qquad a-0 \rightarrow a$
 - $b | false \rightarrow b$ $b \& true \rightarrow b$
- Reassociation & commutativity: - $(a + 1) + 2 \rightarrow a + (1 + 2) \rightarrow a + 3$ $-(2+a) + 4 \rightarrow (a+2) + 4 \rightarrow a + (2+4) \rightarrow a + 6$
- a*4 → a<<2 -a*7 \rightarrow (a << 3) - a $- a/32767 \rightarrow (a >> 15) + (a >> 30)$

• Strength reduction: (replace expensive op with cheaper op)

Note: must be careful with floating point (due to rounding) and integer arithmetic (due to overflow/underflow)

- •
- •
- This is a substitution operation •
- Example: • int x = 5; int y = x * 2; \rightarrow int y = 5 * 2; \rightarrow int y = 10; int $z = a[y]; \rightarrow int z = a[y]; \rightarrow int z = a[y]; \rightarrow int z = a[10];$
- •

Constant Propagation

If the value is known to be a constant, replace the use of the variable by the constant

Value of the variable must be propagated forward from the point of assignment

To be most effective, constant propagation should be interleaved with constant folding

Copy Propagation

- ulletthe copied variable.
- Need to know where copies of the variable propagate. ullet
- Interacts with the scoping rules of the language. \bullet
- Example: ullet
- $\mathbf{x} = \mathbf{y};$ $\mathbf{X} = \mathbf{y};$ if (x > 1) { \rightarrow if (y > 1) { x = y * f(y - 1);x = x * f(x - 1);} }
- \bullet

If one variable is assigned to another, replace uses of the assigned variable with

Can make the first assignment to x *dead* code (that can be eliminated).

Dead Code Elimination

- ullet
- x = y * y// x is dead! // x never used \rightarrow X = Z * ZX = Z * Z
- A variable is *dead* if it is never used after it is defined. \bullet - Computing such *definition* and *use* information is an important component of compiler
- Dead variables can be created by other optimizations... ullet

If a side-effect free statement can never be observed, it is safe to eliminate the statement.

Unreachable/Dead Code

- unreachable and can be deleted.
 - Performed at the IR or assembly level
- Dead code: similar to unreachable blocks. – A value might be computed but never subsequently used.
- Code for computing the value can be dropped •
- But only if it's pure, i.e. it has no externally visible side effects •

 - optimizations (and code transformations in general) easier!

Basic blocks not reachable by any trace leading from the starting basic block are

- Externally visible effects: raising an exception, modifying a global variable, going into an infinite loop, printing to standard output, sending a network packet, launching a rocket Note: Pure functional languages (e.g. Haskell) make reasoning about the safety of

- Example in Oat code: \bullet

```
int g(int x) { return x + pow(x); }
int pow(int a) { int b = 1; int n = 0;
                 while (n < a) \{b = 2 * b\};
                 return b; }
\rightarrow
int g(int x) {
 int a = x; int b = 1; int n = 0;
 while (n < a) \{b = 2 * b\}; tmp = b;
 return x + tmp;
```

- May need to rename variable names to avoid *name capture* \bullet Example of what can go wrong? ____
- Best done at the AST or relatively high-level IR. \bullet
- When is it profitable? \bullet
 - Eliminates the stack manipulation, jump, etc. _____
 - Can increase code size. _____
 - Enables further optimizations _____



Replace a call with the body of the function itself with arguments rewritten to be local variables:

int g(int x) (1 + f(x))Int f(int a) (a + x) \rightarrow const int x = 3; int g(int x) (1 + (int a = x; a + x))

Loop Optimizations

- Most program execution time occurs in loops.
 The 90/10 rule of thumb holds here too. (90% of the execution time is spent in 10% of the code)
- Loop optimizations are very important, effective, and numerous
 Also, concentrating effort to improve loop body code is usually a win

Loop Invariant Code Motion (revisited)

- Another form of redundancy elimination.
- If the result of a statement or expression does not change during the loop and it's pure, it can be hoisted outside the loop body.
- Often useful for array element-addressing code • so-called invariant code

for (i = 0; i < a.length; i++) { /* a not modified in the body */ } t = a.length;for (i =0; i < t; i++) { /* same body as above */

Hoisted loopinvariant expression

Strength Reduction (revisited)

- Strength reduction can work for loops too \bullet
- \bullet
- For loops, create a dependent induction variable: \bullet
- Example: \bullet for (int i = 0; i<n; i++) { a[i*3] = 1; } // stride by 3



```
int j = 0;
for (int i = 0; i<n; i++) {
  a[j] = 1;
  j = j + 3; // replace multiply by add
```

Idea: replace expensive operations (multiplies, divides) by cheap ones (adds and subtracts)

Loop Unrolling (revisited)

Branches can be expensive, unroll loops to avoid them. lacksquarefor (int i=0; i < n; i++) { S }

for (int i=0; i < n-3; i+=4) $\{S;S;S;S\};$ for (; i<n; i++) { S } // left over iterations

- With k unrollings, eliminates (k-1)/k conditional branches \bullet
 - So for the above program, it eliminates ³/₄ of the branches
- Space-time tradeoff: ullet
 - Not a good idea for large S or small n

Code Specialization

- Idea: create specialized versions of a with different arguments.
- Example: specialize function f in:

class A implements I { int m() {...} } class B implements I { int m() {...} } int f(I x) { x.m(); } // A a = new A(); f(a); // B b = new B(); f(b); //

- f_A would have code specialized to dispatch to A.m
- f_B would have code specialized to dispatch to B.m
- You can also *inline* methods when the run-time type is known statically
 Often just one class implements a method.

Idea: create specialized versions of a function that is called from different places

- } }
 // don't know which m
- // know it's A.m
- // know it's B.m

Common Subexpression Elimination (CSE)

- In some sense it's the opposite of inlining: fold redundant computations together \bullet
- Example:

a[i] = a[i] + 1 compiles to: $[a + i^{*}4] = [a + i^{*}4] + 1$

Common subexpression elimination removes the redundant add and multiply:

 $t = a + i^{*}4; [t] = [t] + 1$

For safety, you must be sure that the shared expression *always* has the same value in both places! •



Unsafe Common Subexpression Elimination

Example: consider this OAT function: •

```
unit f(int[] a, int[] b, int[] c) {
   int j = ...; int i = ...; int k = ...;
   b[j] = a[i] + 1;
   c[k] = a[i];
   return;
```

•

```
unit f(int[] a, int[] b, int[] c) {
   int j = ...; int i = ...; int k = ...;
   t = a[i];
   b[j] = t + 1;
   c[k] = t;
   return;
```

The optimization that shares the expression a[i] is unsafe... why?

Time for a short break?





Motivating Code Analyses

- There are lots of things that might influence the safety/applicability of an optimization
 What algorithms and data structures can help?
- How do you know what is a loop?
- How do you know an expression is invariant (constant)?
- How do you know if an expression has no side effects?
- How do you keep track of where a variable is defined?
- How do you know where a variable is used?
- How do you know if two reference values may be aliases of one another?

Moving Towards Register Allocation

- – These are the %uids you should be very familiar with by now.
- Current compilation strategy: •
 - Each %uid maps to a stack location.
 - This yields programs with many loads/stores to memory.
 - Very inefficient.
- Ideally, we'd like to map as many %uid's as possible into registers. •
 - Eliminate the use of the alloca instruction?
 - Only 16 max registers available on 64-bit X86

 - This means that a register must hold more than one slot
- When is this safe? lacksquare

The Oat compiler currently generates as *many* temporary variables as it needs

- %rsp and %rbp are reserved and some have special semantics, so really only 10 or 12 available

Liveness

- Observation: %uid1 and %uid2 can be assigned to the same register if their \bullet values will not be needed at the same time.
 - What does it mean for an %uid to be "needed"?
 - Ans: its contents will be used as a source operand in a later instruction.
- Such a variable is called "*live*" ullet

• Two variables can share the same register if they are *not* live at the same time.



- We can already get some coarse liveness information from variable scoping. ullet
- Consider the following OAT program: \bullet

```
int f(int x) {
  var a = 0;
  if (x > 0) {
     var b = x * x;
     a = b + b;
  var c = a * x;
  return c;
```

- - c's scope is disjoint from b's scope
- ullet

Scope vs. Liveness

Note that due to Oat's scoping rules, variables b and c can never be live at the same time.

So, we could assign b and c to the same alloca'ed slot and potentially to the same register.

Consider this program: ullet



- But, a, b, c are never live at the same time. - So they can share the same stack slot / register

But Scope is too Coarse

• The scopes of a, b, c, x all overlap – they're all in scope at the end of the block.

Live Variable Analysis

- A variable v is *live* at a program point if v is defined before the program point and used after it.
- Liveness is defined in terms of where variables are *defined* and where variables are *used*
- Liveness analysis: Compute the live variables between each statement. - May be conservative (i.e. it may claim a variable is live when it isn't) so because that's a safe approximation – To be useful, it should be more *precise* than simple scoping rules.
- Liveness analysis is one example of *dataflow analysis* \bullet Other examples: Available Expressions, Reaching Definitions, Constant-Propagation Analysis, ... ____



Control-flow Graphs Revisited

- lacksquare
- Recall that a basic block is a sequence of instructions such that:
 - There is a distinguished, labeled entry point (no jumps into the middle of a basic block)
 - There is a (possibly empty) sequence of non-control-flow instructions —
 - The block ends with a single control-flow instruction (jump, conditional branch, return, etc.)
- A control flow graph \bullet
 - Nodes are blocks
 - ____
 - There are no "dangling" edges there is a block for every jump target.

For the purposes of dataflow analysis, we use the *control-flow graph* (CFG) intermediate form.

There is an edge from B1 to B2 if the control-flow instruction of B1 might jump to the entry label of B2

Dataflow over CFGs

- ullet
 - Different implementation tradeoffs in practice...



For precision, it is helpful to think of the "fall through" between sequential instructions as an edge of the control-flow graph too.

"Exploded" CFG

Liveness is Associated with Edges



- This is useful so that the same register can be used for different temporaries in the same statement.
- Example: a = b + 1
- Compiles to:


Uses and Definitions

- Every instruction/statement *uses* some set of variables ullet– i.e. reads from them
- Every instruction/statement *defines* some set of variables – i.e. writes to them
- For a node/statement s define:
 - use[s] : set of variables used by s
 - def[s] : set of variables defined by s
- Examples: •

 - a = b + c $use[s] = \{b,c\}$ $def[s] = \{a\}$

- a = a + 1 $use[s] = {a} def[s] = {a}$

Liveness, Formally

- A variable v is *live* on edge e if: There is
 - a node n in the CFG such that use[n] contains v, and
 - a directed path from e to n such that for every statement s' on the path, def[s'] does not contain v
- The first clause says that v will be used on *some* path starting from edge e. The second clause says that v won't be redefined on that path before the use.
- Questions:
 - How to compute this efficiently?
 - How to use this information (e.g., for register allocation)?
 - How does the choice of IR affect this? (e.g. LLVM IR uses SSA, so it doesn't allow redefinition \Rightarrow simplify liveness analysis)

Simple, inefficient algorithm

- a directed path from e to n passing through no def of v."
- Backtracking Algorithm: •
 - For each variable v...
 - until either v is defined or a previously visited node has been reached.
 - Mark the variable v live across each edge traversed.
- Why inefficient?
- Because it explores the same paths many times • (for different uses and different variables)

• "A variable v is live on an edge e if there is a node n in the CFG using it and

- Try all paths from *each use* of v, tracing *backwards* through the control-flow graph

Dataflow Analysis

- *Idea*: compute liveness information for all variables simultaneously. – Keep track of sets of information about each node
- - Equations based on "obvious" constraints.
- Solve the equations by iteratively converging on a solution. •
 - Start with a "rough" approximation to the answer
 - Refine the answer at each iteration
 - Keep going until no more refinement is possible: a *fixpoint* has been reached
- dataflow analysis

Approach: define equations that must be satisfied by any liveness determination.

• This is an instance of a general framework for computing program properties:

Dataflow Value Sets for Liveness

- Nodes are program statements, so:
 - use[n] : set of variables used by n
 - def[n] : set of variables defined by n
 - in[n] : set of variables live on entry to n
 - out[n] : set of variables live on exit from n
- Associate in[n] and out[n] with the "collected" information about incoming/outgoing edges
- For Liveness: what constraints are there among these sets?
- Clearly:

 $in[n] \supseteq use[n]$

What other constraints?



Other Dataflow Constraints

- We have: $in[n] \supseteq use[n]$
 - "A variable must be live on entry to n if it is used by n"
- Also: $in[n] \supseteq out[n] def[n]$ lacksquare
 - "If a variable is live on exit from n, and n doesn't define it, it is live on entry to n"
 - Note: here '-' means "set difference"
- And: $out[n] \supseteq in[n']$ if $n' \in succ[n]$ •
 - "If a variable is live on entry to a successor node of n, it must be live on exit from n."



Iterative Dataflow Analysis

- Find a solution to those constraints by starting from a rough guess. - Start with: $in[n] = \emptyset$ and $out[n] = \emptyset$
- The guesses don't satisfy the constraints:
 - in[n] \supseteq use[n]
 - in[n] \supseteq out[n] def[n]
 - $out[n] \supseteq in[n']$ if $n' \in succ[n]$
- Idea: iteratively re-compute in [n] and out [n] where forced to by the constraints. - Each iteration will add variables to the sets in[n] and out[n] (i.e. the live variable sets will increase monotonically)
- We stop when in[n] and out[n] satisfy these equations: (which are derived from the constraints above) What are they?
 - in[n] := use[n] \cup (out[n] def[n])
 - out[n] := $\bigcup_{n' \in \text{succ}[n]} \text{in}[n']$

Complete Liveness Analysis Algorithm

for all n, $in[n] := \emptyset$, $out[n] := \emptyset$ for all n end end

- Finds a *fixpoint* of the in and out equations. • - The algorithm is guaranteed to terminate... Why?
- Why do we start with \emptyset ?

repeat until no change in 'in' and 'out'

```
out[n] := \bigcup_{n' \in succ[n]} in[n']
in[n] := use[n] \cup (out[n] - def[n])
```

Example Liveness Analysis

• Example flow graph:



Example Liveness Analysis

Each iteration update: $out[n] := \bigcup_{n' \in succ[n]} in[n']$ $in[n] := use[n] \cup (out[n] - def[n])$

Iteration 1:
in[2] = x
in[3] = e
in[4] = x
in[5] = e,x
in[6] = x
in[7] = x
in[8] = z
in[9] = y

(showing only updates that make a change)

out:



> Iteration 2: out[1] = xin[1] = xout[2] = e,xin[2] = e,xout[3] = e,xin[3] = e,xout[5] = xout[6] = xout[7] = z,yin[7] = x, z, yout[8] = xin[8] = x, zout[9] = xin[9] = x,y



> Iteration 3: ulletout[1] = e,xout[6] = x,y,zin[6] = x,y,zout[7] = x,y,zout[8] = e,xout[9] = e,x

> > out: e,

> Iteration 4: lacksquareout[5] = x,y,zin[5] = e,x,z

> Iteration 5: out[3] = e,x,z

> > Done!

out: e,x

Improving the Iterative Algorithm

- Can we do better?
- Observe: the only way information propagates from one node to another is using: $out[n] := \bigcup_{n' \in succ[n]} in[n']$
 - This is the only rule that involves more than one node
- If a node's successors haven't changed, then the node itself won't change.
- Idea for an improved version of the algorithm:
 Keep track of which node's successors have changed

A Worklist Algorithm

Use a FIFO queue of nodes that might need to be updated. •

for all n, $in[n] := \emptyset$, $out[n] := \emptyset$ w = new queue with all nodesrepeat until w is empty let n = w.pop()old in = in[n] out[n] := $\bigcup_{n' \in succ[n]} in[n']$ $in[n] := use[n] \cup (out[n] - def[n])$ if (old in != in[n]), // if in[n] has changed for all m in pred[n], w.push(m) // add to worklist end

// pull a node off the queue // remember old in[n]

How about another break?

Other Dataflow Analyses

Generalising Dataflow Analyses

- The kind of iterative constraint solving used for liveness applies to other kinds of analyses.
 - Reaching Definitions analysis
 - Available Expressions analysis
 - Alias Analysis
 - **Constant Propagation** _____
 - These analyses follow the same approach as for liveness (accumulating values until fixpoint is reached). ____
- To see these as an instance of the same kind of algorithm, the next few examples to work over a canonical intermediate instruction representation called *quadruples* (op, arg1, arg2, and result)
 - Allows easy definition of def[n] and use[n]
 - A slightly "looser" variant of LLVM's IR that doesn't require the "static single assignment"
 - i.e. it has *mutable* local variables
 - We will use LLVM-IR-like syntax

Reaching Definitions

Reaching Definition Analysis

- Question: what uses in a program does a given variable definition reach?
- Unlike liveness, we are interested in *different* definitions of the same variable. \bullet
- This analysis is used for constant propagation & copy propagation lacksquare

 - available expressions analysis (next)
- Input: Quadruple CFG
- may reach the beginning (resp. end) of node n

– If only one definition reaches a particular use, can replace use by the definition (for constant propagation). Copy propagation additionally requires that the copied value still has its same value – computed using an

Output: in[n] (resp. out[n]) is the set of *nodes* defining some variable such that the definition

Example of Reaching Definitions

•

Results of computing reaching definitions on this simple CFG:

Reaching Definitions Step 1

- Define the sets of interest for the analysis lacksquare
 - Let defs[a] be the set of *nodes* (statements) that define the variable a •
- Define gen[n] and kill[n] as follows: \bullet

•	Quadruple forms n:	gen[n]	kill[
	a = b op c	{n}	defs
	a = load b	{n}	defs
	store b, a	Ø	Ø
	$a = f(b_1,, b_n)$	{n}	defs
	$f(b_1,,b_n)$	Ø	Ø
	br L	Ø	Ø
	braL1 L2	Ø	Ø
	return a	Ø	Ø

ullet

[n] $s[a] - \{n\}$

s[a] - {n}

 $s[a] - \{n\}$

gen[n] are node's definitions; kill[n] are the nodes, whose definitions are "shadowed" by n

Reaching Definitions Step 2

- Define the constraints that a reaching definitions solution must satisfy.
- $out[n] \supseteq gen[n]$
- $in[n] \supseteq out[n']$ if n' is in pred[n]
 - "The definitions that reach the beginning of a node include those that reach the exit of its any predecessor"
- out[n] \cup kill[n] \supseteq in[n]

 - Equivalently: $out[n] \supseteq in[n] kill[n]$

"The definitions that reach the end of a node at least include the definitions generated by the node"

"The definitions that come in to a node either reach the end of the node or are killed by it."

Reaching Definitions Step 3

- Convert constraints to iterated update equations: •
 - $in[n] := \bigcup_{n' \in pred[n]} out[n']$
 - $out[n] := gen[n] \cup (in[n] kill[n])$
- Algorithm: initialise in[n] and out[n] to \emptyset • – Iterate the update equations until a fixed point is reached - Why does it terminate?
- The algorithm terminates because in [n] and out [n] increase only monotonically lacksquare– At most to a maximum set that includes all variable definitions in the program
- •

The algorithm is *precise* because it finds the *smallest* sets that satisfy the constraints.

- Idea: want to perform common subexpression elimination: • -a = x + 1 a = x + 1
 - When is it safe?
- (i.e., x hasn't been assigned). - "x+1" is an available expression
- Dataflow values:
 - in[n] = set of *nodes* whose values are available on entry to n
 - out[n] = set of*nodes*whose values are available on exit of n

Available Expressions

This transformation is safe if x+1 means computes the same value at both places

Available Expressions Step 1

- Define the sets of values ullet
- Define gen[n] and kill[n] as follows: •
 - Quadruple forms n: gen[n] ullet $\{n\}$ - kill[n] a = b op ca = load b $\{n\}$ - kill[n] store b, a Ø
 - br L Ø br a L1 L2 Ø
 - $a = f(b_1, ..., b_n)$ Ø

$$f(b_1,...,b_n)$$
Øreturn aØ

- gen[n] node itself represents new available expression
- kill[n] nodes whose expressions no longer available after n

• Let uses [a] be the set of *nodes* that use the variable a in their expressions

```
kill[n]
uses a
uses [a]
                      Note the need for
                      "may alias" information...
uses[[x]]
(for all x that may equal a)
\emptyset
Ø
uses[a] \cup uses[[x]]
     (for all x)
                              Note that functions are
                 (for all x)
uses[[x]]
                              assumed to be impure...
Ø
```


Available Expressions Step 2

- Define the constraints that an available expressions solution must satisfy. •
- $out[n] \supseteq gen[n]$
 - "The expressions made available by n that reach the end of the node"
- $in[n] \subseteq out[n']$ if n' is in pred[n]
 - every predecessor"
- out[n] \cup kill[n] \supseteq in[n]

 - Equivalently: $out[n] \supseteq in[n] kill[n]$

• "The expressions available at the beginning of a node include those that reach the exit of

"The expressions available on entry either reach the end of the node or are killed by it."

Note similarities and differences with constraints for "reaching definitions".

Available Expressions Step 3

- Convert constraints to iterated update equations:
 - $in[n] := \bigcap_{n' \in pred[n]} out[n']$
 - $out[n] := gen[n] \cup (in[n] kill[n])$
- Unlike previous algorithms, this one is "shrinking" the set of desired facts
- Algorithm: initialise in[n] and out[n] to {set of all nodes}
 - Iterate the update equations until a fixed point is reached
 - Why does the algorithm terminate?
- The algorithm terminates because in [n] and out [n] decrease only monotonically
 - At most to a minimum of the empty set

The algorithm is precise because it finds the *largest* sets that satisfy the constraints.

General Dataflow Analysis Framework

Comparing Dataflow Analyses

- Look at the update equations in the inner loop of the analyses
- Liveness:
 - Let gen[n] = use[n] and kill[n] = def[n]
 - $out[n] := U_{n' \in succ[n]}in[n']$
 - $in[n] := gen[n] \cup (out[n] kill[n])$
- Reaching Definitions:
 - $in[n] := U_{n' \in pred[n]}out[n']$
 - $out[n] := gen[n] \cup (in[n] kill[n])$
- Available Expressions:
 - $in[n] := \bigcap_{n' \in pred[n]} out[n']$
 - $out[n] := gen[n] \cup (in[n] kill[n])$

(backward)

(forward)

(forward)

Common Features

- All of these analyses have a *domain* over which they solve constraints. •
 - Liveness, the domain is sets of variables – Reaching defns., Available exprs. the domain is sets of nodes
- Each analysis has a notion of gen[n] and kill[n] ullet- Used to explain how information *propagates* across a node: what is added, what is removed.
- Each analysis is propagates information either forward or backward ullet- Forward: in[n] defined in terms of predecessor nodes' out[]
- Backward: out[n] defined in terms of successor nodes' in[]
- Each analysis has a way of aggregating (combining) information from in/out flow • – Liveness & reaching definitions take union (\cup)
- - Available expressions uses intersection (\cap)
 - Union expresses a property that holds for some path (existential)
 - Intersection expresses a property that holds for *all* paths (universal)

(Forward) Dataflow Analysis Framework

A forward dataflow analysis can be characterized by:

- A domain of dataflow values \mathcal{L}
 - e.g. $\int f$ = the powerset of all variables
 - Think of $\ell \in \mathcal{L}$ as a property, then " $z \in \ell$ " means "z has the property"
- 2. For each node n, a flow function $F_n : \mathcal{L} \to \mathcal{L}$
 - So far we've seen $F_n(\ell) = gen[n] \cup (\ell kill[n])$
 - So: out[n] = $F_n(in[n])$
 - "If ℓ is a property that holds before the node n, then $F_n(\ell)$ holds after n"
- A combining operator ⊓ 3.
 - "If we know either ℓ_1 or ℓ_2 holds on entry to node n, we know at most $\ell_1 \sqcap \ell_2$ "
 - $in[n] := \prod_{n' \in pred[n]} out[n']$

Generic Iterative (Forward) Analysis

repeat until no change for all n

- $in[n] := \prod_{n' \in pred[n]} out[n']$ $out[n] := F_n(in[n])$

end

end

- Here, $\top \in \mathcal{L}$ ("top") represents having the "maximum" amount of information. – Having "more" information enables more optimizations "Maximum" amount could be inconsistent with the constraints, so we can't keep it. :-(– Iteration refines the answer, eliminating inconsistencies

- for all n, $in[n] := \top$, $out[n] := \top$

Structure of *L*

- The domain has structure that reflects the "amount" of information for each dataflow value.
- Some dataflow values are more informative than others:
 - Write $\ell_1 \subseteq \ell_2$ whenever ℓ_2 provides at least as much information as ℓ_1 .
 - The dataflow value ℓ_2 is "better" for enabling optimizations.
- Example 1: for available expressions analysis, larger sets of nodes are more informative. - Having a larger set of nodes (equivalently, expressions) available means that there is more opportunity for
- common subexpression elimination.
 - So: $\ell_1 \subseteq \ell_2$ if and only if $\ell_1 \subseteq \ell_2$
- Example 2: for liveness analysis, *smaller* sets of variables are more informative. – Having smaller sets of variables live across an edge means that there are fewer conflicts
 - for register allocation assignments.
 - So: $\ell_1 \sqsubseteq \ell_2$ if and only if $\ell_1 \supseteq \ell_2$

Las a Partial Order

- \mathcal{L} is a *partial order* defined by the ordering relation \sqsubseteq .
- A partial order is an ordered set.
- Some of the elements might be *incomparable*.
 That is, there might be ℓ₁, ℓ₂ ∈ L such that neither ℓ₁ ⊑ ℓ₂ nor ℓ₂ ⊑ ℓ₁
- Properties of a partial order:
 - Reflexivity: $\ell \sqsubseteq \ell$
 - *Transitivity*: $\boldsymbol{\ell}_1 \sqsubseteq \boldsymbol{\ell}_2$ and $\boldsymbol{\ell}_2 \sqsubseteq \boldsymbol{\ell}_3$ implies $\boldsymbol{\ell}_1 \sqsubseteq \boldsymbol{\ell}_2$
 - Anti-symmetry: $\ell_1 \sqsubseteq \ell_2$ and $\ell_2 \sqsubseteq \ell_1$ implies $\ell_1 = \ell_2$
- Examples:
 - Integers ordered by \leq
 - Types ordered by <:
 - Sets ordered by \subseteq or \supseteq

Subsets of $\{a,b,c\}$ ordered by \subseteq

 $\boldsymbol{\ell}_2$ $\sqsubseteq \boldsymbol{\ell}_2$ l

Partial orders are often presented as a Hasse diagram.



order ⊑ is ⊆ meet ⊓ is ∩

join L is U

- The *combining* operator \sqcap is called the "meet" operation.
- It constructs the greatest lower bound: •
 - $-\ell_1 \sqcap \ell_2 \sqsubseteq \ell_1$ and $\ell_1 \sqcap \ell_2 \sqsubseteq \ell_2$ "the meet is a lower bound"
 - If $\ell \subseteq \ell_1$ and $\ell \subseteq \ell_2$ then $\ell \subseteq \ell_1 \sqcap \ell_2$ "there is no greater lower bound"
- Dually, the \sqcup operator is called the "join" operation. •
- It constructs the *least upper bound*: ullet
 - $-\ell_1 \sqsubseteq \ell_1 \sqcup \ell_2$ and $\ell_2 \sqsubseteq \ell_1 \sqcup \ell_2$ "the join is an upper bound"
 - If $\ell_1 \sqsubseteq \ell$ and $\ell_2 \sqsubseteq \ell$ then $\ell_1 \sqcup \ell_2 \sqsubseteq \ell$ "there is no smaller upper bound"
- A partial order that has all meets and joins is called a *lattice*. ullet– If it has just meets, it's called a meet semi-lattice.

Meets and Joins

Another Way to Describe the (Forward) Algorithm

- Algorithm repeatedly computes (for each node n):
 - $out[n] := F_n(in[n])$
- Equivalently: $out[n] := F_n(\prod_{n' \in pred[n]} out[n'])$ By definition of in[n]
- Let $x_n = out[n]$

 - $\mathbf{F}(\mathbf{X}) = (F_1(\prod_{i \in pred[1]} out[j]), F_2(\prod_{i \in pred[2]} out[j]), ..., F_n(\prod_{i \in pred[n]} out[j]))$
- Any solution to the constraints is a *fixpoint* **X** of **F** - i.e. F(X) = X

We can write this as a simultaneous update of the vector of out[n] values:

- Let $X = (x_1, x_2, ..., x_n)$ it's a vector of points in \mathcal{L} corresponding to CFG nodes



Iteration Computes Fixpoints

- Let $\mathbf{X}_0 = (\top, \top, ..., \top)$
- Each loop through the algorithm apply **F** to the old vector: • $\mathbf{X}_1 = \mathbf{F}(\mathbf{X}_0)$ $\mathbf{X}_2 = \mathbf{F}(\mathbf{X}_1)$

 $\mathbf{F}^{k+1}(\mathbf{X}) = \mathbf{F}(\mathbf{F}^{k}(\mathbf{X}))$ ullet

• • •

- A fixpoint is reached when $F^k(X) = F^{k+1}(X)$ •
 - That's when the algorithm stops.
- Wanted: a maximal fixpoint \bullet

– Because that one is more informative/useful for performing optimizations

Monotonicity & Termination

- Each flow function F_n maps lattice elements to lattice elements; to be sensible is should be monotonic: ullet
- $F: \mathcal{I} \to \mathcal{I}$ is monotonic iff: $\ell_1 \sqsubseteq \ell_2$ implies that $F(\ell_1) \sqsubseteq F(\ell_2)$ - Intuitively: "If you have more information entering a node, then you have more information leaving the node."
- Monotonicity lifts point-wise to the function: $\mathbf{F}: \mathcal{L}^n \to \mathcal{L}^n$ \bullet - vector $(x_1, x_2, \dots, x_n) \sqsubseteq (y_1, y_2, \dots, y_n)$ iff $x_i \sqsubseteq y_i$ for each i
- Note that **F** is consistent: $F(X_0) \sqsubseteq X_0$ ullet

– So each iteration moves at least one step down the lattice (for some component of the vector) $- \ldots \sqsubseteq \mathsf{F}(\mathsf{F}(\mathsf{X}_0)) \sqsubseteq \mathsf{F}(\mathsf{X}_0) \sqsubseteq \mathsf{X}_0$

Therefore, # steps needed to reach a fixpoint is at most the height H of \mathcal{L} times the number of nodes: $O(H_n)$ — height of the lattice

Building Lattices?

Information about individual nodes or variables can be lifted pointwise: \bullet - If \mathcal{L} is a lattice, then so is $\{f: X \rightarrow \mathcal{L}\}$ where $f \sqsubseteq g$ if and only if $f(x) \sqsubseteq g(x)$ for all $x \in X$.

 \bullet - Could pick a lattice based on subtyping:

– Or other information:

Aliased Unaliased

Points in the lattice are sometimes called dataflow "facts"

Like *types*, the dataflow lattices are *static approximations* to the dynamic behavior:



More on Fixpoint Solutions

Remember constructing LL(1) parse tables ullet

Т	\longmapsto	S\$
S	\longmapsto	ES'
S'	\longmapsto	8
S'	\longmapsto	+ S
E	\longmapsto	number (S)

- First(T) = First(S)
- First(S) = First(E)
- $First(S') = \{ + \}$
- First(E) = { number, '(' }
- Follow(S') = Follow(S)

	number	+	(
Т	\mapsto S\$		⊢→S\$	
S	$\mapsto E S'$		$\mapsto E S'$	
S'		$\mapsto + S$		H
E	⊢ num.		\mapsto (S)	

• Follow(S) = { \$, ')' } \cup Follow(S')



Then: we want the *least* solution to this system of set equations... a fixpoint computation. More on these later in the course.

Now: This solution is obtained by starting from taking all First/Follow as \emptyset and then iterating the equations until *fixpoint* is reached.



Dataflow Analysis: Summary

- Many dataflow analyses fit into a common framework. •
- Key idea: *iterative solution* of a system of equations over a *lattice* of *facts* (constraints). • – Iteration terminates if flow functions are *monotonic*.

 - Solution is obtained as the greatest fixpoint is reached via the meet operation (\Box) .
- In the literature, sometimes the definition of the analysis lattice is reversed: • – The most useful/precise information is represented by the bottom element (\bot) - Solution is obtained as the *least* fixpoint via iterative application of *join* operator (\Box) – The two definitions are equivalent modulo the (semi-)lattice *direction*.

Next Lecture (Finally!)

- Register Allocation
- Modern research directions in PL and Compilers
- Wrap-Up