Lecture 20
CIS 341: COMPILERS

#### Announcements

- Project 5 Compiling objects in full Oat
  - Available from the course web pages
  - Updated oat.pdf fixes a few typos (mentioned on Piazza)
  - Due April 8<sup>th</sup>

- Final Exam:
  - Tuesday, April 30<sup>th</sup> noon-2:00 pm
  - Moore 216

A high-level tour of a variety of optimizations.

# **OPTIMIZATIONS**

Zdancewic CIS 341: Compilers

## 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
- High-level languages make avoiding redundant computation inconvenient or impossible
  - 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:
  - Time: improve execution speed
  - Space: reduce amount of memory needed
  - Power: lower power consumption (e.g. to extend battery life)

#### **Some caveats**

- Optimization are code transformations:
  - They can be applied at any stage of the compiler
  - They must be *safe* they can'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 optimization

## 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
- Example: Loop unrolling

```
- Idea: rewrite a loop like:
    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 codes space slows down whole program a tiny bit but speeds up the loop
  - 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.
  - These have a much bigger impact on performance that compiler optimizations.
  - Reduce # of operations
  - Reduce memory accesses
  - Minimize indirection it breaks working-set coherence
- *Then* turn on compiler optimizations
- Profile to determine program hot spots
- Evaluate whether the algorithm/data structure design works
- ...if so: "tweak" the source code until the optimizer does "the right thing" to the machine code

# Safety

- Whether an optimization is *safe* depends on the programming language semantics.
  - Languages that provide weaker guarantees to the programmer permit more optimizations, but have more ambiguity in their behavior.
  - e.g. In Java tail-call optimization (that turns recursive function calls into loops) is not valid.
  - e.g. In C, loading from initialized memory is undefined, so the compiler can do anything.
- Example: loop-invariant code motion
  - Idea: hoist invariant code out of a loop



- Is this more efficient?
- Is this safe?

#### **Constant Folding**

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

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

#### **Constant Folding Conditionals**



## **Algebraic Simplification**

- More general form of constant folding
  - Take advantage of mathematically sound simplification rules
- Identities:

—	a	*	1 -	a a		a	*	0	→	0	
_	a	+	0 -	🕨 a		a	_	0	→	a	
_	b	T	fal	se 🕇	b	b	æ	+ r	-11e	→	b

• Reassociation & commutativity:

- 
$$(a + 1) + 2 \rightarrow a + (1 + 2) \rightarrow a + 3$$

$$-(2 + a) + 4 \rightarrow (a + 2) + 4 \rightarrow a + (2 + 4) \rightarrow a + 6$$

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

—	а	*	4	→	a << 2
_	a	*	7	→	(a << 3) – a
—	a	/	32767	→	(a >> 15) + (a >> 30

- Note 1: must be careful with floating point (due to rounding)
- Note 2: iteration of these optimizations is useful... how much?

#### **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
  - This is a substitution operation

Example: int x = 5; int y = x \* 2; → int y = 5 \* 2; → int y = 10; → int z = a[y]; int z = a[y]; int z = a[10];

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

## **Copy Propagation**

- If one variable is assigned to another, replace uses of the assigned variable with the copied variable.
- Need to know where copies of the variable propagate.
- Interacts with the scoping rules of the language.



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

#### **Dead Code Elimination**

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

- A variable is *dead* if it is never used after it is defined.
  - Computing such *definition* and *use* information is an important component of compiler
- Dead variables can be created by other optimizations...

### **Unreachable/Dead Code**

- Basic blocks not reachable by any trace leading from the starting basic block are *unreachable* and can be deleted.
  - Performed at the canonical IR or assembly level
  - Improves cache, TLB performance
- 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
  - 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 optimizations (and code transformations in general) easier!

# Inlining

- Replace a call to a function with the body of the function itself with arguments rewritten to be local variables:
- Example in OAT code:

```
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; }</pre>
```

#### →

```
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* 
  - Example of what can go wrong?
- Best done at the AST or relatively high-level IR.
- When is it profitable?
  - Eliminates the stack manipulation, jump, etc.
  - Can increase code size.
  - Enables further optimizations

## **Code Specialization**

- Idea: create specialized versions of a function that is called from different places 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(); } // don't know which m
A a = new A(); f(a); // know it's A.m
B b = new B(); f(b); // know it's B.m
```

- **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.

## **LOOP OPTIMIZATIONS**

Zdancewic CIS 341: Compilers

## **Common Subexpression Elimination**

- In some sense it's the opposite of inlining: fold redundant computations together
- 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;
    }
```

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

```
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;
}
```

# **Loop Optimizations**

- Program hot spots often occur in loops.
  - Especially inner loops
  - Not always: consider operating systems code or compilers vs. a computer game or word processor
- 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
  - Invariant code not visible at the source level



## **Strength Reduction (revisited)**

- Strength reduction can work for loops too
- Idea: replace expensive operations (multiplies, divides) by cheap ones (adds and subtracts)
- For loops, create a dependent induction variable:

```
• Example:
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
}
```

#### Loop Unrolling (revisited)

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



- With k unrollings, eliminates (k-1)/k conditional branches
  - So for the above program, it eliminates <sup>3</sup>/<sub>4</sub> of the branches
- Space-time tradeoff:
  - Not a good idea for large S or small n
- Interacts with instruction caching, branch prediction