

# Computer Systems

CS107

Cynthia Lee

# Today's Topics

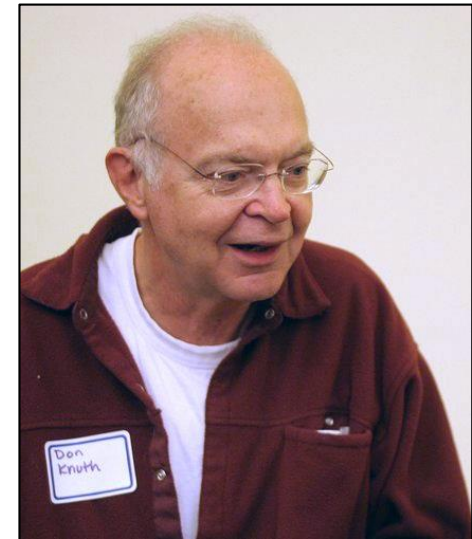
- Code optimization!

# Optimization reality check

Don't let it be your Waterloo.

## Optimization Reality Check

- Optimization is really exciting
- ...but it's easy to be overzealous or misguided about it.
- “We should forget about small inefficiencies, say about 97% of the time: premature optimization is the root of all evil.” – Donald Knuth
- “More computing sins are committed in the name of efficiency (without necessarily achieving it) than for any other single reason—including blind stupidity.” – W.A. Wulf
- “Bottlenecks occur in surprising places, so don't try to second guess and put in a speed hack until you have proven that's where the bottleneck is.” – Rob Pike



## Optimization Reality Check

- Most of what you need to do with optimization can be summarized in 3 easy steps:
  
- Step 1:
  - › If doing something seldom and only on small inputs, do whatever is simplest to code, understand, and debug
  
- Step 2:
  - › If doing things thing a lot, or on big inputs, make the algorithm's Big-O cost reasonable
  
- Step 3:
  - › **Let gcc do its magic from there**

## gcc optimization levels

- Today, we'll be comparing two levels of optimization in the gcc compiler:
  - › `gcc -O0` //mostly just literal translation of C
  - › `gcc -O2` //enable nearly all reasonable optimizations
  - › (we use `-Og`, like `-O0` but with less needless use of the stack)
- There are other custom and more aggressive levels of optimization, e.g.:
  - › `-O3` //more aggressive than `O2`, trade size for speed
  - › `-Os` //optimize for size
  - › `-Ofast` //disregard standards compliance (!!)
- Exhaustive list of gcc optimization-related flags:
  - › <https://gcc.gnu.org/onlinedocs/gcc/Optimize-Options.html>

## Example: gcc performance optimization

- Just a standard matrix multiply, triply-nested for loop:

```
static void mmm(double a[][DIM], double b[][DIM],
               double c[][DIM], int n) {
    for (int i = 0; i < n; i++)
        for (int j = 0; j < n; j++)
            for (int k = 0; k < n; k++)
                c[i][j] += a[i][k]*b[k][j];
}
```

## Measuring performance: Example code

```
> ./mult // -O0 (no optimization)
matrix multiply 25^2: cycles 0.44M
matrix multiply 50^2: cycles 3.13M
matrix multiply 100^2: cycles 24.80M
> ./mult_opt // -O2 (with optimization)
matrix multiply 25^2: cycles 0.11M (opt)
matrix multiply 50^2: cycles 0.47M (opt)
matrix multiply 100^2: cycles 3.67M (opt)
```

When I said, “Let gcc do its magic from there,” this is what I meant by magic!





# Kinds of optimization

Some main categories of optimization that are undertaken by the compiler

## Kinds of optimization

- Constant folding
- Common subexpression elimination
- Dead code
- Strength reduction
- Code motion
- Tail recursion
- Loop unrolling

## Kinds of optimization

- **Constant folding**

- › Precalculate constants at compile time where possible:

- `int volume = WIDTH * HEIGHT * DEPTH;`
- `double radius = sqrt(area) / 3.14;`

- Common subexpression elimination

- Dead code

- Strength reduction

- Code motion

- Tail recursion

- Loop unrolling

- **Discussion question:**

- › **What is a consequence of this for you as a programmer?**

What should you do (or do differently) now that you know about compilers doing constant folding for you?

# Constant folding

## Before:

```
00000000040098d <CF>:
40098d: 55
40098e: 48 89 e5
400991: 41 54
400993: 53
400994: 48 83 c4 80
400998: 89 bd 7c ff ff ff
40099e: c7 45 ec 07 01 00 00
4009a5: 8b 45 ec
4009a8: 6b c0 55
4009ab: 89 45 e8
4009ae: 48 b8 00 00 00 00 00
4009b5: 00 00 40
4009b8: 48 89 85 70 ff ff ff
4009bf: f2 0f 10 85 70 ff ff
4009c6: ff
4009c7: e8 b4 fe ff ff
4009cc: f2 0f 2c c0
4009d0: 89 45 e4
4009d3: 8b 45 ec
4009d6: 0f af 85 7c ff ff ff
4009dd: 89 c3
4009df: b8 15 00 00 00
4009e4: 99
4009e5: f7 7d e4
4009e8: 89 c2
4009ea: 8b 45 ec
4009ed: 01 d0
4009ef: 4c 63 e0
4009f2: bf b8 1b 40 00
4009f7: e8 d4 fd ff ff
```

## After:

```
000000000400d80 <CF>:
400d80: 69 c7 07 01 00 00    imul  $0x107,%edi,%eax
400d86: 05 61 6d 00 00      add   $0x6d61,%eax
400d8b: c3                  retq
```

```
push  %rbp
mov   %rsp,%rbp
push  %r12
push  %rbx
add   $0xfffffffffffffff80,%rsp
mov   %edi,-0x84(%rbp)
movl  $0x107,-0x14(%rbp)
mov   -0x14(%rbp),%eax
imul  $0x55,%eax,%eax
mov   %eax,-0x18(%rbp)
movabs $0x4000000000000000,%rax

mov   %rax,-0x90(%rbp)
movsd -0x90(%rbp),%xmm0

callq 400880 <sqrt@plt>
cvttsd2si %xmm0,%eax
mov   %eax,-0x1c(%rbp)
mov   -0x14(%rbp),%eax
imul  -0x84(%rbp),%eax
mov   %eax,%ebx
mov   $0x15,%eax
cld
idivl -0x1c(%rbp)
mov   %eax,%edx
mov   -0x14(%rbp),%eax
add   %edx,%eax
movslq %eax,%r12
mov   $0x401bb8,%edi
callq 4007d0 <strlen@plt>
```

\*\* Notice it also got rid of the stack frame handling overhead

## Kinds of optimization

- Constant folding
- **Common subexpression elimination**
  - › Prevent recalculation of the same thing many times by doing it once and saving result

```
int a = (param2 + 0x107);
int b = param1 * (param2 + 0x107) + a;
return a * (param2 + 0x107) + b * (param2 + 0x107);
```
- Dead code
- Strength reduction
- Code motion
- Tail recursion
- Loop unrolling

## CSE

### Before:

0000000000400a1c <CSE>:

```
400a1c:      55
400a1d:      48 89 e5
400a20:      89 7d ec
400a23:      89 75 e8
400a26:      8b 45 e8
400a29:      05 07 01 00 00
400a2e:      89 45 fc
400a31:      8b 45 e8
400a34:      05 07 01 00 00
400a39:      0f af 45 ec
400a3d:      89 c2
400a3f:      8b 45 fc
400a42:      01 d0
400a44:      89 45 f8
400a47:      8b 45 f8
400a4a:      8b 55 fc
400a4d:      01 c2
400a4f:      8b 45 e8
400a52:      05 07 01 00 00
400a57:      0f af c2
400a5a:      5d
400a5b:      c3
```

### After:

0000000000400d90 <CSE>:

```
400d90:      81 c6 07 01 00 00      add    $0x107,%esi
400d96:      0f af fe              imul   %esi,%edi
400d99:      8d 04 77              lea    (%rdi,%rsi,2),%eax
400d9c:      0f af c6              imul   %esi,%eax
400d9f:      c3                    retq
```

```
push    %rbp
mov     %rsp,%rbp
mov     %edi,-0x14(%rbp)
mov     %esi,-0x18(%rbp)
mov     -0x18(%rbp),%eax
add     $0x107,%eax
mov     %eax,-0x4(%rbp)
mov     -0x18(%rbp),%eax
add     $0x107,%eax
imul   -0x14(%rbp),%eax
mov     %eax,%edx
mov     -0x4(%rbp),%eax
add     %edx,%eax
mov     %eax,-0x8(%rbp)
mov     -0x8(%rbp),%eax
mov     -0x4(%rbp),%edx
add     %eax,%edx
mov     -0x18(%rbp),%eax
add     $0x107,%eax
imul   %edx,%eax
pop     %rbp
retq
```

## Kinds of optimization

- Constant folding
- Common subexpression elimination
- **Dead code**
  - › Remove code that doesn't serve a purpose:

```
while (false) {  
    j = func(j);  
    printf("This loop can't happen!\n");  
    i++;  
    if (i == 1000000) break;  
}
```
- Strength reduction
- Code motion
- Tail recursion
- Loop unrolling

# DC

## Before:

000000000400a5c <DC>:

```
400a5c: 55
400a5d: 48 89 e5
400a60: 48 83 ec 20
400a64: 89 7d ec
400a67: 89 75 e8
400a6a: 8b 45 ec
400a6d: 3b 45 e8
400a70: 7d 17
400a72: 8b 45 ec
400a75: 3b 45 e8
400a78: 7e 0f
400a7a: bf c0 1b 40 00
400a7f: b8 00 00 00 00
400a84: e8 57 fd ff ff
400a89: c7 45 fc 00 00 00 00
400a90: eb 04
400a92: 83 45 fc 01
400a96: 81 7d fc e7 03 00 00
400a9d: 7e f3
400a9f: 8b 45 ec
400aa2: 3b 45 e8
400aa5: 75 06
400aa7: 83 45 ec 01
400aab: eb 04
400aad: 83 45 ec 01
400ab1: 83 7d ec 00
400ab5: 75 07
400ab7: b8 00 00 00 00
400abc: eb 03
400abe: 8b 45 ec
400ac1: c9
```

```
push %rbp
mov %rsp,%rbp
sub $0x20,%rsp
mov %edi,-0x14(%rbp)
mov %esi,-0x18(%rbp)
mov -0x14(%rbp),%eax
cmp -0x18(%rbp),%eax
jge 400a89 <DC+0x2d>
mov -0x14(%rbp),%eax
cmp -0x18(%rbp),%eax
jle 400a89 <DC+0x2d>
mov $0x401bc0,%edi
mov $0x0,%eax
callq 4007e0 <printf@plt>
movl $0x0,-0x4(%rbp)
jmp 400a96 <DC+0x3a>
addl $0x1,-0x4(%rbp)
cmpl $0x3e7,-0x4(%rbp)
jle 400a92 <DC+0x36>
mov -0x14(%rbp),%eax
cmp -0x18(%rbp),%eax
jne 400aad <DC+0x51>
addl $0x1,-0x14(%rbp)
jmp 400ab1 <DC+0x55>
addl $0x1,-0x14(%rbp)
cmpl $0x0,-0x14(%rbp)
jne 400abe <DC+0x62>
mov $0x0,%eax
jmp 400ac1 <DC+0x65>
mov -0x14(%rbp),%eax
leaveq
```

## After:

000000000400da0 <DC>:

```
400da0: 8d 47 01 lea 0x1(%rdi),%eax
400da3: c3 retq
400da4: 66 66 66 2e 0f 1f 84 data32 data32
```



## Kinds of optimization

- Constant folding
- Common subexpression elimination
- Dead code
- **Strength reduction**
  - › Change divide to multiply, multiply to add or shift, and mod to and
  - › Avoids using instructions that cost many cycles (multiply and divide)

```
int doge_years = human_years * 7;
```

- Code motion
- Tail recursion
- Loop unrolling



## Kinds of optimization

- Constant folding
- Common subexpression elimination
- Dead code
- Strength reduction
- **Code motion**
  - › Move code out of a loop if possible

```
for (int i = 0; i < n; i++) {  
    sum += arr[i] + foo * (bar + 3);  
}
```
- Tail recursion
- Loop unrolling

## Kinds of optimization

- Constant folding
  - Common subexpression elimination
  - Dead code
  - Strength reduction
  - Code motion
  - **Tail recursion**
    - › Compiler notices some simple recursion patterns that could be more efficiently implemented using iteration (i.e. a loop) to avoid function call and return overhead
- ```
long factorial(int n) {  
    if (n<=1) return 1;  
    else return n * factorial(n-1);  
}
```
- Loop unrolling

## Kinds of optimization

- Constant folding
- Common subexpression elimination
- Dead code
- Strength reduction
- Code motion
- Tail recursion
- **Loop unrolling**
  - › Do  $n$  loop iterations' worth of work per actual loop iteration, so we save ourselves from doing the loop overhead (test and jump) every time, and instead incur overhead only every  $n$ -th time

```
for (int i=0; i<=n-4; i+=4) {  
    sum += arr[i];  
    sum += arr[i+1];  
    sum += arr[i+2];  
    sum += arr[i+3];
```

```
} // after the loop handle any leftovers
```

## For loop construction

# for-loop literal translation

Initialization

Test

Branch past loop if fails

Body # loop unrolling

Increment Initialization

jmp to Test jmp to Test

Body

Increment

Body

Increment

Body

Increment

Test

Branch to Body if succeeds

# for-loop gcc actually emits

Initialization

jmp to Test

Body

Increment

Test

Branch to Body if succeeds

# Measuring performance

Lab preview

## Two techniques for measuring performance—learn more in lab

- Wall clock time
  - › `gettimeofday()` // `#include <sys/time.h>`
  
- Cycle counting
  - › RTC: real-time clock counts elapsed cycles of the CPU
  - › Available on some hardware
  - › Will reveal some cases where an instruction takes more than one cycle (and a few interesting cases where more than one instruction was able to execute per clock cycle)
    - Recall that multiply instruction is more expensive (takes more cycles) than add
  - › See `/afs/ir/class/cs107/samples/lect15/fcyc.h`