



**HARVARD**

School of Engineering  
and Applied Sciences

# Program Optimization

*CS61, Lecture 9*

Prof. Stephen Chong

September 29, 2011

# Announcements

- Homework 3 (the Buffer Bomb) due in one week
- There will be no section on Columbus Day (Monday Oct 10)
  - Section TFs will reschedule, and info posted on website
  - Remember: you can attend any section

# CISC vs RISC

- CISC (“sisk”): Complex Instruction Set Computer
- RISC (“risk”): Reduced Instruction Set Computer
- Different philosophies regarding the design (and implementation) of ISAs

# CISC vs RISC

- CISC (“sisk”): Complex Instruction Set Computer
  - Historically first
  - Large instruction sets (x86 has several hundred)
  - Specialized instructions for high-level tasks
    - Instructions that are closer to what applications are wanting to do
    - Can provide hardware support for application-specific instructions
      - ▶ E.g., x86 contains instructions such as *LOOPZ label*, which decrements %cx (without modifying flags) and jumps to label if %cx is non-zero
  - Presents a clean interface to programmer
    - Hides implementation details such as pipelining

# CISC vs RISC

- RISC (“risk”): Reduced Instruction Set Computer
  - Philosophy developed in early 1980s
  - Small, simple, instruction sets (typically <100)
    - E.g., may have only base+displacement memory addressing
    - E.g., memory access only via load and store; ALU operations need register operands
    - E.g., no condition codes, only explicit test instructions
    - Often fixed length encodings for instructions
    - Leads to simple, efficient implementation
  - Reveals implementation details to programmer (e.g., pipe-lining)
    - E.g., certain instruction sequences may be prohibited
    - E.g., jump instruction may not take effect until after following instruction
      - ▶ Compiler must be aware of these restrictions, and can use them to optimize performance
  - ARM (originally “Acorn RISC Machine”) widely used in embedded devices

# Modern computers

- In most settings, neither CISC nor RISC clearly better
- RISC machines
  - Exposing implementation details made it difficult to use them, and difficult to evolve the ISA
  - Added more instructions
- CISC machines
  - Take advantage of RISC-like pipelines
    - Essentially translate CISC instructions into simpler RISC-like instructions
    - E.g., `addl %eax, 8(%esp)` broken up into a load from memory, followed by an addition, followed by a store to memory

# Today

- Program optimization
  - Overview
  - Code motion
  - Strength reduction
  - Common subexpressions
  - Optimization blockers
    - Procedure calls
    - Aliasing
  - Understanding modern processors
  - Loop unrolling
  - Tail recursion
  - Summary

# Getting the best performance

- There's more to performance than asymptotic complexity!
- Constant factors matter too
  - Easily see  $10\times$ – $100\times$  difference depending on how code is written
  - Must optimize at multiple levels:
    - algorithm structure (locality, instruction level parallelism, ...)
    - data representations (e.g., structs vs arrays)
    - coding style (e.g., unnecessary procedure calls, unrolling, reordering, ...)
- Must understand underlying system to optimize performance
  - How programs are compiled and executed
  - How to measure program performance and identify bottlenecks
  - How to improve performance while maintaining code modularity and generality



# Optimizing compilers (e.g., gcc)

- Compilers do a **lot** of optimization when generating machine code
- Use optimization flags when compiling
  - Default is no optimization (`-O0`)
  - Good choices for gcc: `-O2`, `-O3`, `-march=xxx`, `-m64`
  - Try different flags and maybe different compilers

# Optimizing compilers (e.g., gcc)

- Compilers are **good** at: mapping programs to machines
  - register allocation
  - instruction selection and ordering (scheduling)
  - dead code elimination
  - eliminating minor inefficiencies
- Compilers are **not good** at: improving asymptotic efficiency
  - up to programmer to select best overall algorithm
  - big-O savings are (often) more important than constant factors
  - but constant factors also matter
- Compilers are **not good** at: overcoming “optimization blockers”
  - potential memory aliasing
  - potential procedure side-effects

# Limitations of Optimizing Compilers

- When in doubt, the **compiler must be conservative**
- Must not change program behavior under any possible condition
  - Often prevents it from making optimizations when would only affect behavior under pathological conditions.
- Behavior that may be obvious to the programmer can be obfuscated by languages and coding styles
  - e.g., data ranges may be more limited than variable types suggest
- Most analysis is performed only within procedures
  - Whole-program analysis is too expensive in most cases
  - Not amenable to modular compilation
- Code analysis generally based only on **static** information
  - That is, whatever it can determine at compile time
  - Difficult (in general, undecidable) to determine run-time, or **dynamic**, behavior

# Machine-independent optimizations

- Some simple optimizations, regardless of specific machine or compiler
  - Code motion
  - Strength reduction
  - Common subexpressions
- For some instances of these optimizations, almost all compilers will perform them
- For other instances, very difficult for a compiler to perform them
  - You need to understand why

# Code motion

- **Key idea:** Move code to reduce the number of times it executes
- Most common case: move code out of loop

• E.g.

```
void set_row(long *a, long *b,
             long i, long n)
{
    long j;

    for (j = 0; j < n; j++) {
        a[n*i+j] = b[j];
    }
}
```

```
long j;
int ni = n*i;
for (j = 0; j < n; j++) {
    a[ni+j] = b[j];
}
```

- Moving code means  $n-1$  fewer multiplications!

# Compiler generated code motion

```
set_row:
    pushl   %ebp                # Setup
    movl   %esp, %ebp
    pushl   %esi
    pushl   %ebx
    movl   12(%ebp), %esi       # esi = b
    movl   20(%ebp), %ebx       # ebx = n
    testl  %ebx, %ebx          # is n <= 0?
    jle    .L26                # return
    movl   %ebx, %edx          # edx = n
    imull  16(%ebp), %edx      # edx = n*i
    movl   8(%ebp), %eax        # eax = a
    leal   (%eax,%edx,4), %edx  # edx = &(a[n*i])
    movl   $0, %ecx            # ecx = 0

.L25:
    movl   (%esi,%ecx,4), %eax  # eax = &(b[j])
    movl   %eax, (%edx)        # a[n*i+j] = b[j]
    addl   $1, %ecx            # j++
    addl   $4, %edx            # edx = next element of a
    cmpl  %ecx, %ebx          # j == n?
    jne    .L25                # if not, continue loop

.L26:
    popl   %ebx                # Finish
    popl   %esi
    popl   %ebp
    ret
```

```
long j;
int ni = n*i;
for (j = 0; j < n; j++) {
    a[ni+j] = b[j];
}
```



# Strength reduction



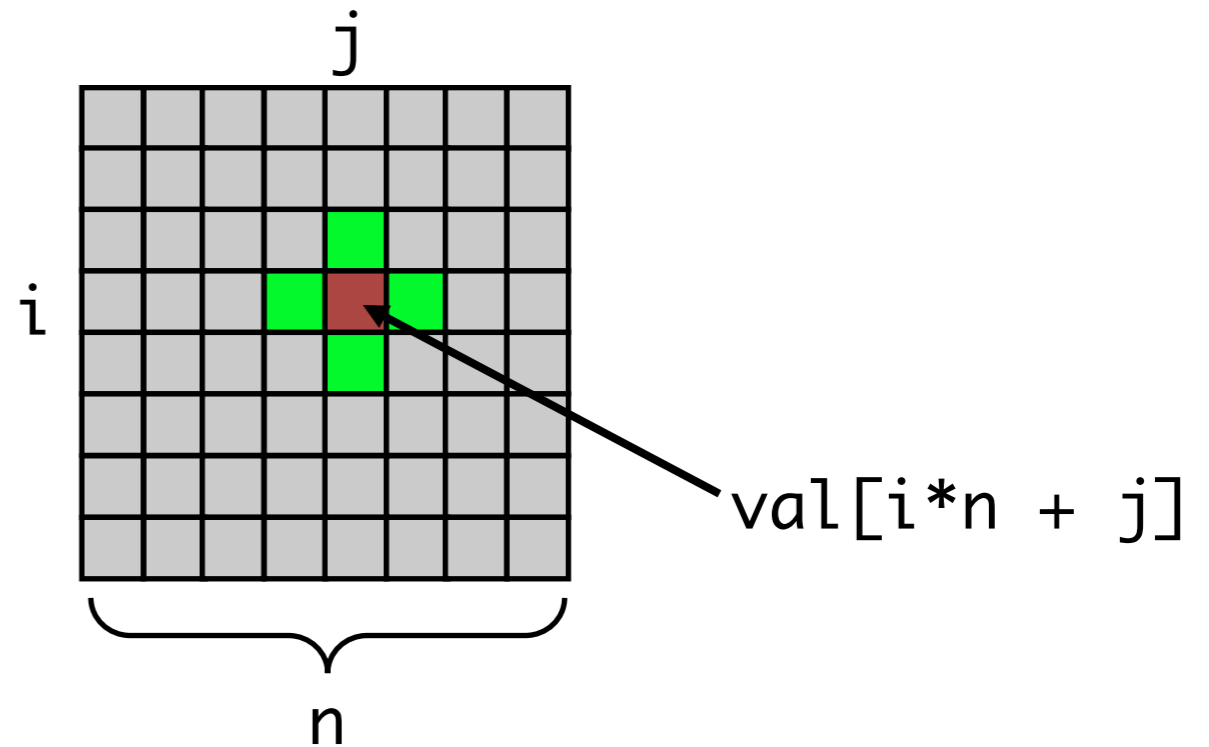
- **Key idea:** replace expensive operations with cheaper ones
- E.g., reduce a multiplication inside a loop to an addition
  - Addition of integers much faster than multiplication

```
/* sum column i of n x n array a */
int sum_col(int *a, int n, int i) {
    int s = 0;
    for (j = 0; j < n; j++) {
        s += a[n*j+i];
    }
    return s;
}
```

```
/* sum column i of n x n array a */
int sum_col(int *a, int n, int i) {
    int s = 0;
    int r = 0;
    for (j = 0; j < n; j++) {
        s += a[r+i];
        r += n;
    }
    return s;
}
```

# Share Common Subexpressions

- **Key idea:** reuse common portions of expressions
- Compilers often not very sophisticated in exploiting arithmetic properties



```
/* Sum neighbors of i,j */
up = val[(i-1)*n + j];
down = val[(i+1)*n + j];
left = val[i*n + j-1];
right = val[i*n + j+1];
sum = up + down + left + right;
```

```
int inj = i*n + j;
up = val[inj - n];
down = val[inj + n];
left = val[inj - 1];
right = val[inj + 1];
sum = up + down + left + right;
```



# Share Common Subexpressions

```
imull %eax, %esi
leal  (%ebx,%esi), %esi
leal  -1(%ecx), %edx
imull %eax, %edx
addl  %ebx, %edx
addl  $1, %ecx
imull %ecx, %eax
addl  %eax, %ebx
movl  (%edi,%ebx,4), %eax
addl  (%edi,%edx,4), %eax
movl  -4(%edi,%esi,4), %edx
addl  4(%edi,%esi,4), %edx
addl  %edx, %eax
```

3 multiplications

```
/* Sum neighbors of i,j */
up = val[(i-1)*n + j];
down = val[(i+1)*n + j];
left = val[i*n + j-1];
right = val[i*n + j+1];
sum = up + down + left + right;
```

```
imull %ecx, %edx
addl  16(%ebp), %edx
leal  0(,%edx,4), %edi
movl  %edx, %esi
subl  %ecx, %esi
movl  -4(%ebx,%edi), %eax
addl  (%ebx,%esi,4), %eax
addl  %edx, %ecx
movl  4(%ebx,%edi), %edx
addl  (%ebx,%ecx,4), %edx
addl  %edx, %eax
```

1 multiplication

```
int inj = i*n + j;
up = val[inj - n];
down = val[inj + n];
left = val[inj - 1];
right = val[inj + 1];
sum = up + down + left + right;
```

# Today

- Program optimization
  - Overview
  - Code motion
  - Strength reduction
  - Common subexpressions
  - Optimization blockers
    - Procedure calls
    - Aliasing
  - Understanding modern processors
  - Loop unrolling
  - Tail recursion
  - Summary

# Optimization Blocker: Procedure Calls

- Converting a string to lower case:

```
void lower(char *s) {
    int i;
    for (i = 0; i < mystrlen(s); i++)
        if (s[i] >= 'A' && s[i] <= 'Z')
            s[i] -= ('A' - 'a');
}
```

```
/* Return length of string s */
size_t mystrlen(const char *s) {
    size_t length = 0;
    while (*s != '\0') {
        s++;
        length++;
    }
    return length;
}
```

- What's wrong (performance-wise) with this code?

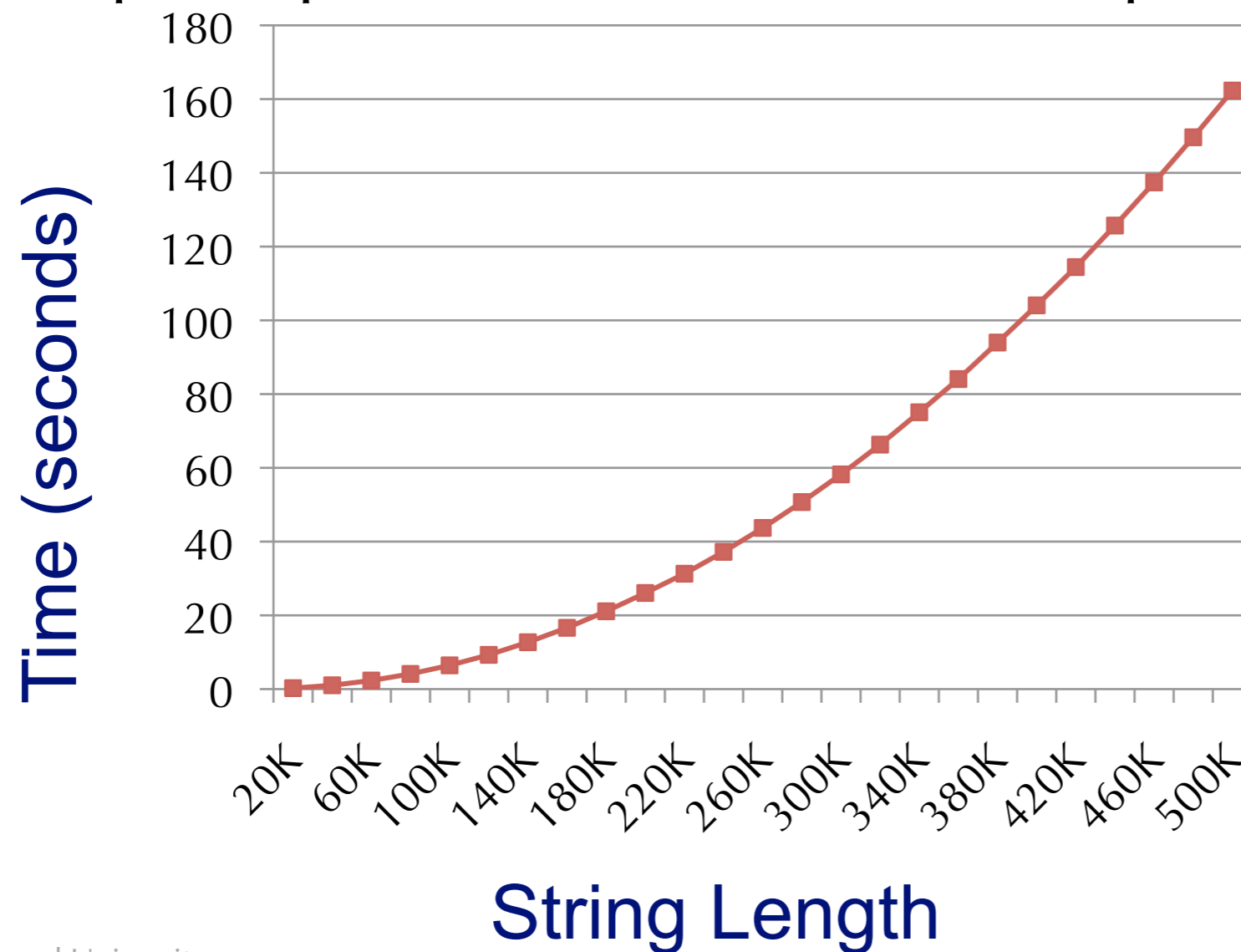
# Convert Loop To Goto Form

- `mystrlen` executed every iteration!
- `mystrlen()` performance
  - Must scan string looking for null character.
- Overall performance, string of length  $n$ 
  - $n$  calls to `mystrlen`
  - Each call requires  $n$  accesses (i.e., go through entire string)
  - Overall  $O(n^2)$  performance

```
void lower(char *s)
{
    int i = 0;
    if (i >= mystrlen(s))
        goto done;
loop:
    if (s[i] >= 'A' && s[i] <= 'Z')
        s[i] -= ('A' - 'a');
    i++;
    if (i < mystrlen(s))
        goto loop;
done:
}
```

# Lower Case Conversion Performance

- $O(n^2)$ 
  - Quadratic performance
  - Time quadruples when we double the input string length



# How to improve performance?

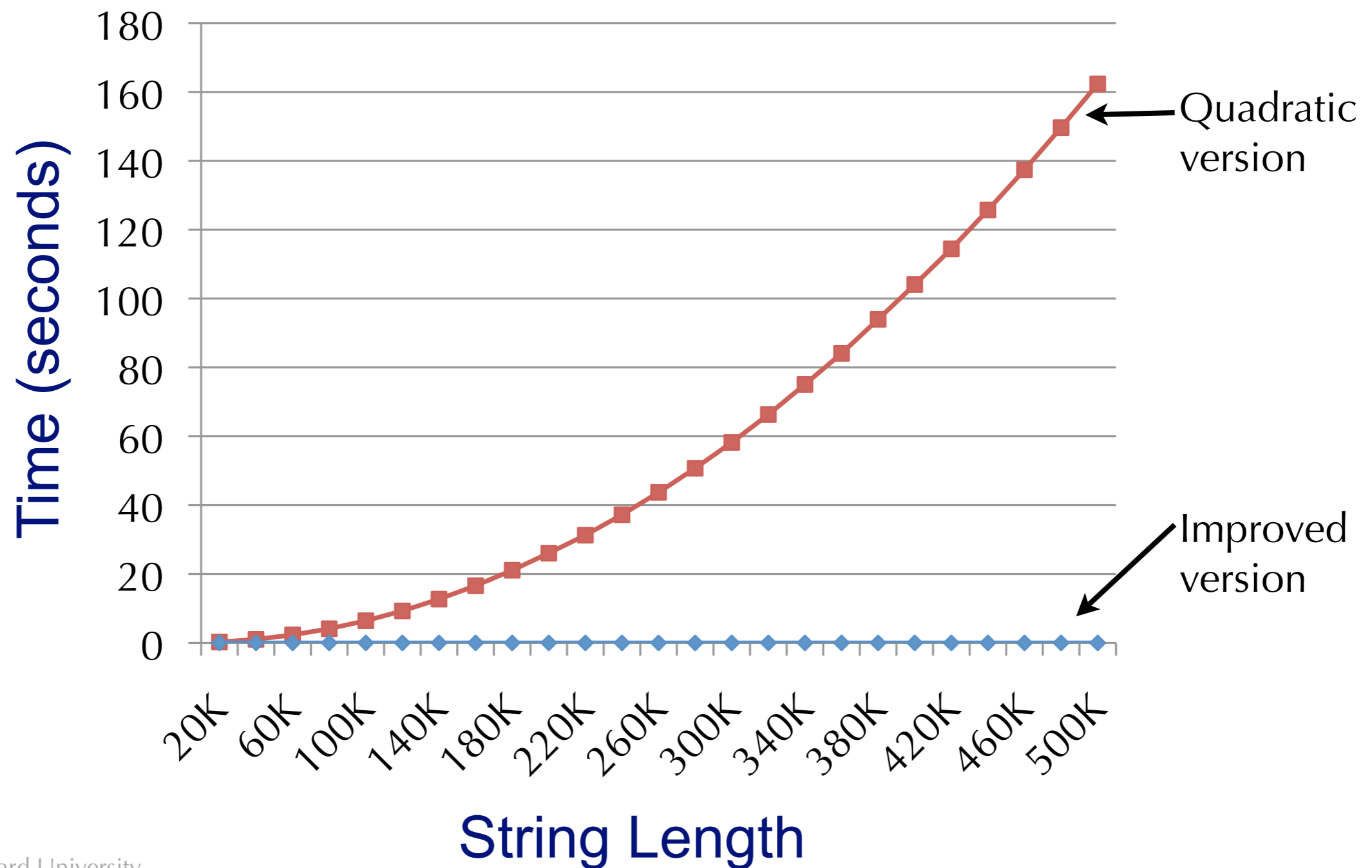
```
void lower(char *s) {
    int i;
    for (i = 0; i < mystrlen(s); i++)
        if (s[i] >= 'A' && s[i] <= 'Z')
            s[i] -= ('A' - 'a');
}
```

- Code motion!

- Move call to `mystrlen()` outside of loop
- OK because result does not change from one iteration to another

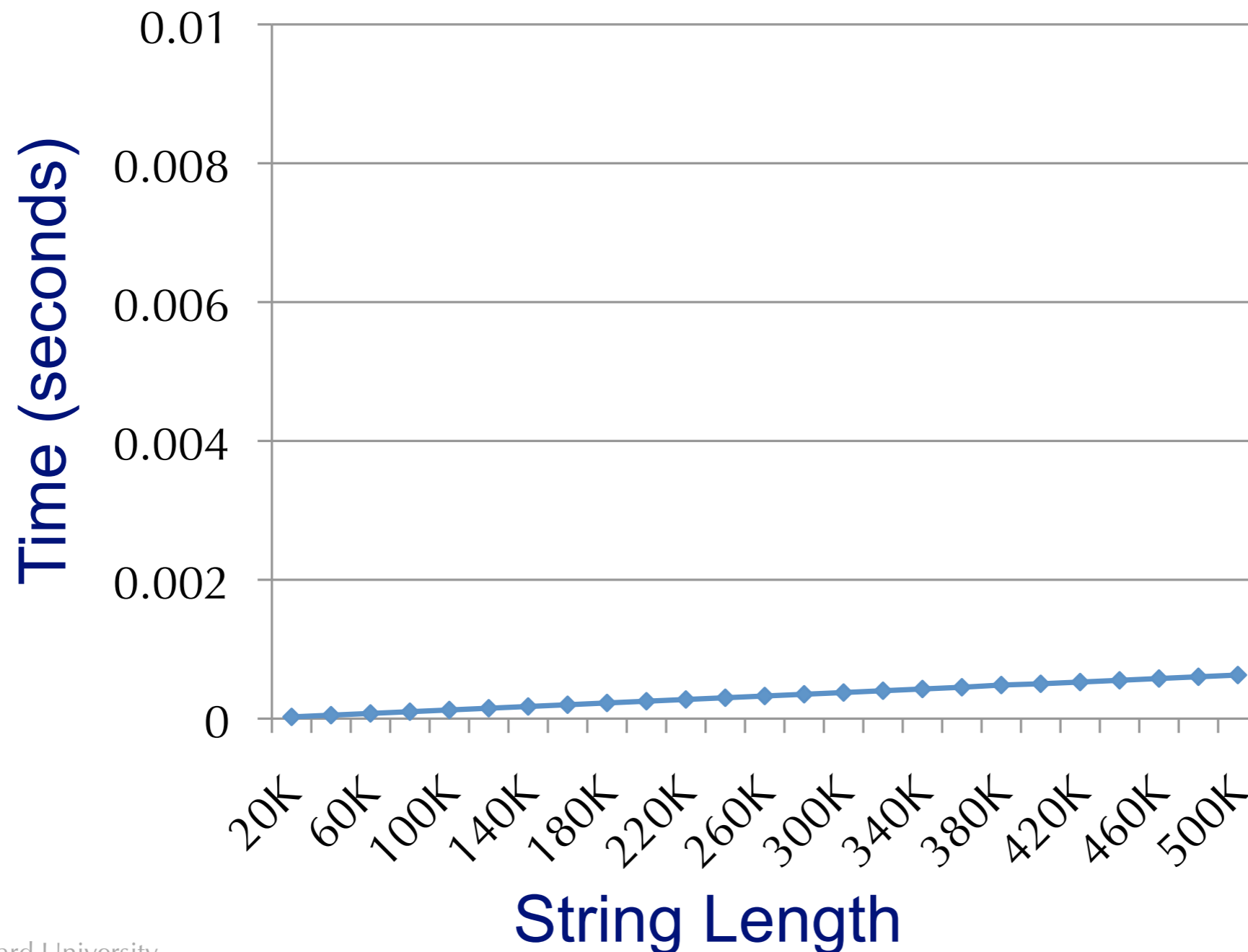
```
void lower(char *s)
{
    int i;
    int len = mystrlen(s);
    for (i = 0; i < len; i++)
        if (s[i] >= 'A' && s[i] <= 'Z')
            s[i] -= ('A' - 'a');
}
```

# Improved performance!



# Improved performance!

- Linear performance





# Optimization Blocker: Procedure Calls

- Why couldn't compiler move `mystrlen()` out of inner loop?
- The compiler treats procedure calls as a "black box"
  - Must be conservative!
- Procedure may be **nondeterministic**
  - Does not return same value each time it is called with same inputs
  - Output could depend on global state (not just its input parameters)
- Procedure may have **side effects**
  - Alters global state each time called

# Example: mystrlen with side effects

```
int lencnt = 0;
size_t mystrlen(const char *s)
{
    size_t length = 0;
    while (*s != '\0') {
        s++; length++;
    }
    lencnt += length;
    return length;
}
```

- Calling `mystrlen` once versus calling it  $n$  times has different behavior!
- (gcc does know about some “built in” functions, including `strlen` and other functions from the standard library. Can optimize knowing about the behavior of these functions)

# Potential remedies

- Do your own code motion
  - Rewrite code to move procedure call outside of the inner loop
- Use the `inline` keyword
  - Tells compiler that the function code can be inserted into the calling function
  - Allows compiler to optimize across caller and callee
  - Also done by default (for “simple” functions) when using `gcc -O3` (or use `-finline-functions`)

```
static inline size_t mystrlen(const char *s) {
    size_t length = 0;
    while (*s != '\0') {
        s++; length++;
    }
    return length;
}
```

# Optimization blocker: aliasing

```
void twiddle1(int *xp, int *yp) {  
    *xp += *yp;  
    *xp += *yp;  
}
```

```
void twiddle2(int *xp, int *yp) {  
    *xp += 2* *yp;  
}
```

- Are the two functions above equivalent?
  - If so, `twiddle2` looks more efficient. Compiler should optimize `twiddle1` so it looks like `twiddle2`, right?

# Optimization blocker: aliasing

```
void twiddle1(int *xp, int *yp) {  
    *xp += *yp;  
    *xp += *yp;  
}
```

```
void twiddle2(int *xp, int *yp) {  
    *xp += 2* *yp;  
}
```

- But what if `xp` and `yp` are equal?
  - e.g., `int foo = 42; twiddle1(&foo, &foo);`
  - `twiddle1` computes:
    - `foo += foo; // doubles foo`
    - `foo += foo; // doubles foo again`
  - `twiddle2` computes:
    - `foo += 2* foo; // triples foo`
- Not equivalent!!!

# Memory aliasing

- If two pointers point to the same memory location, they *alias* each other.
- Compiler must assume that pointers may alias each other
  - Must be conservative!
  - Severely limits optimizations
- Lesson: Reduce unnecessary memory accesses

# Reduce unnecessary memory accesses

- The following programs are not equivalent
  - Why?
- **prod\_array1** must access memory repeatedly
  - Compiler cannot remove these accesses
- **prod\_array2** can be compiled using a register for **res**
  - Much more efficient

```
void prod_array1(int *a, int n,
                int *dest) {
    int i;
    *dest = 1;
    for (i = 0; i < n; i++) {
        *dest = *dest * a[i];
    }
}
```

```
void prod_array2(int *a, int n,
                int *dest) {
    int i, res = 1;
    for (i = 0; i < n; i++) {
        res = res * a[i];
    }
    *dest = res;
}
```

# Today

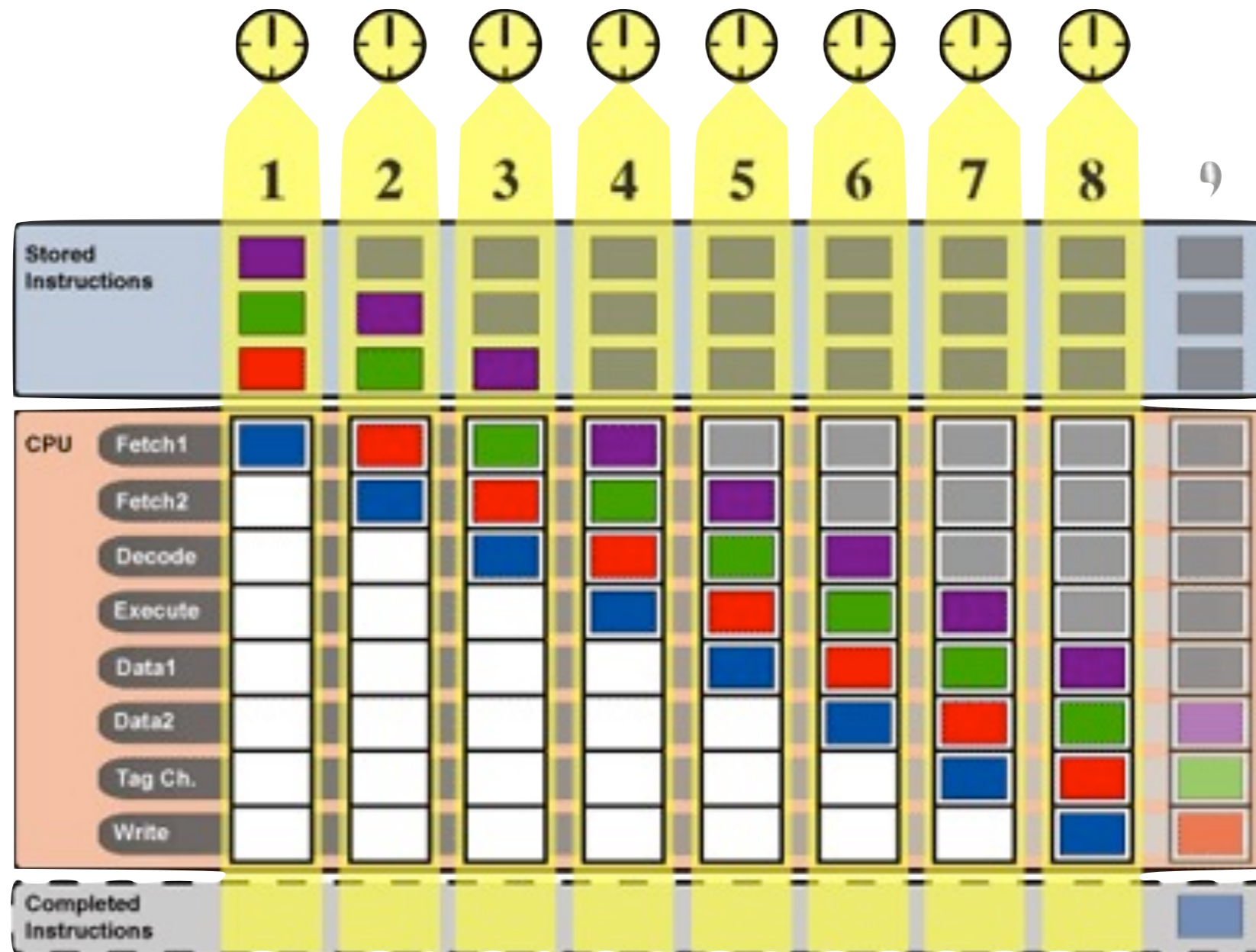
- Program optimization
  - Overview
  - Code motion
  - Strength reduction
  - Common subexpressions
  - Optimization blockers
    - Procedure calls
    - Aliasing
  - Understanding modern processors
  - Loop unrolling
  - Tail recursion
  - Summary



# Three kinds of parallelism

- Three kinds of parallelism supported by modern CPUs:
  - Pipelining
  - Superscalar
  - Multicore

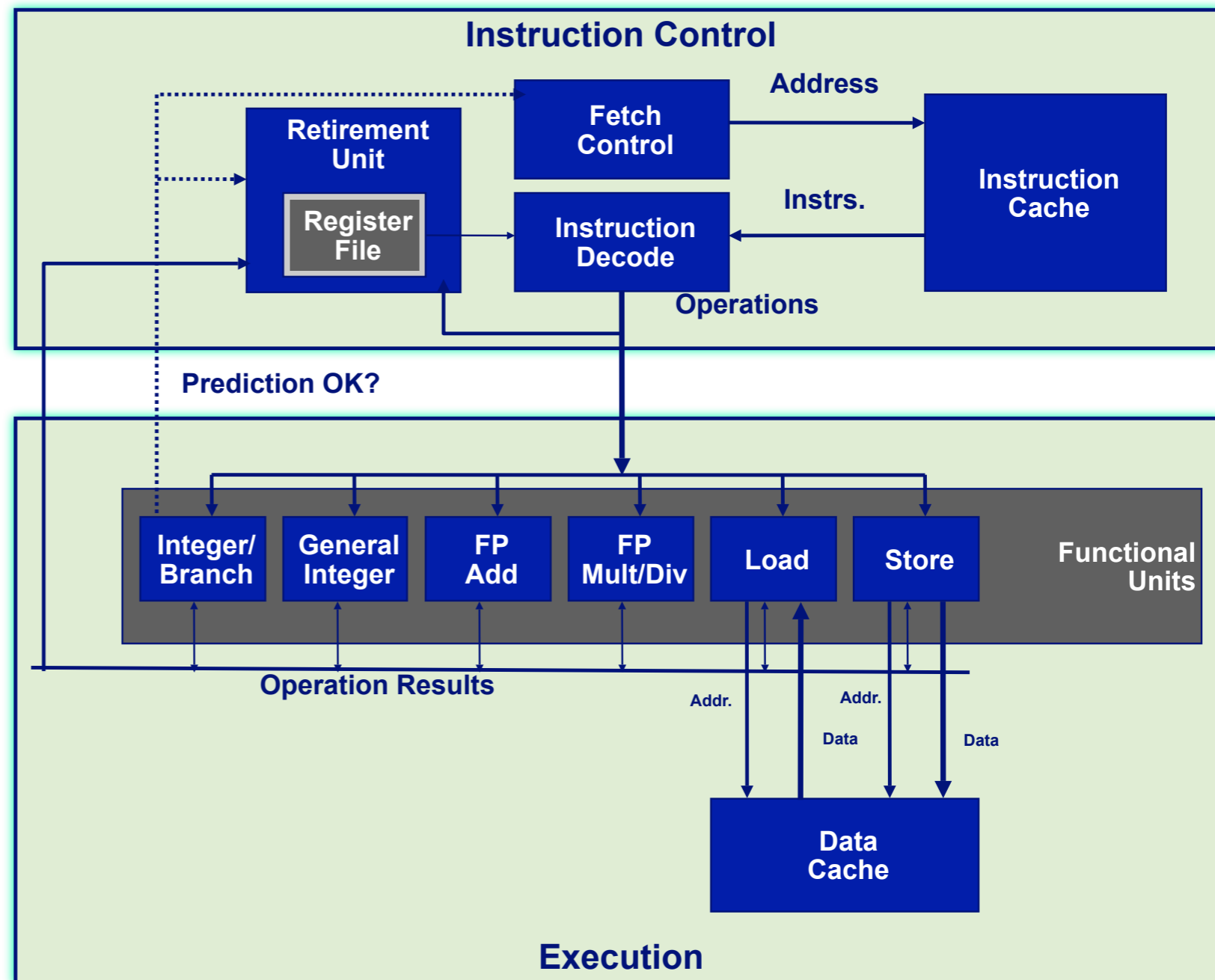
# Pipelining



<http://arstechnica.com/old/content/2004/09/pipelining-2.ars/4>

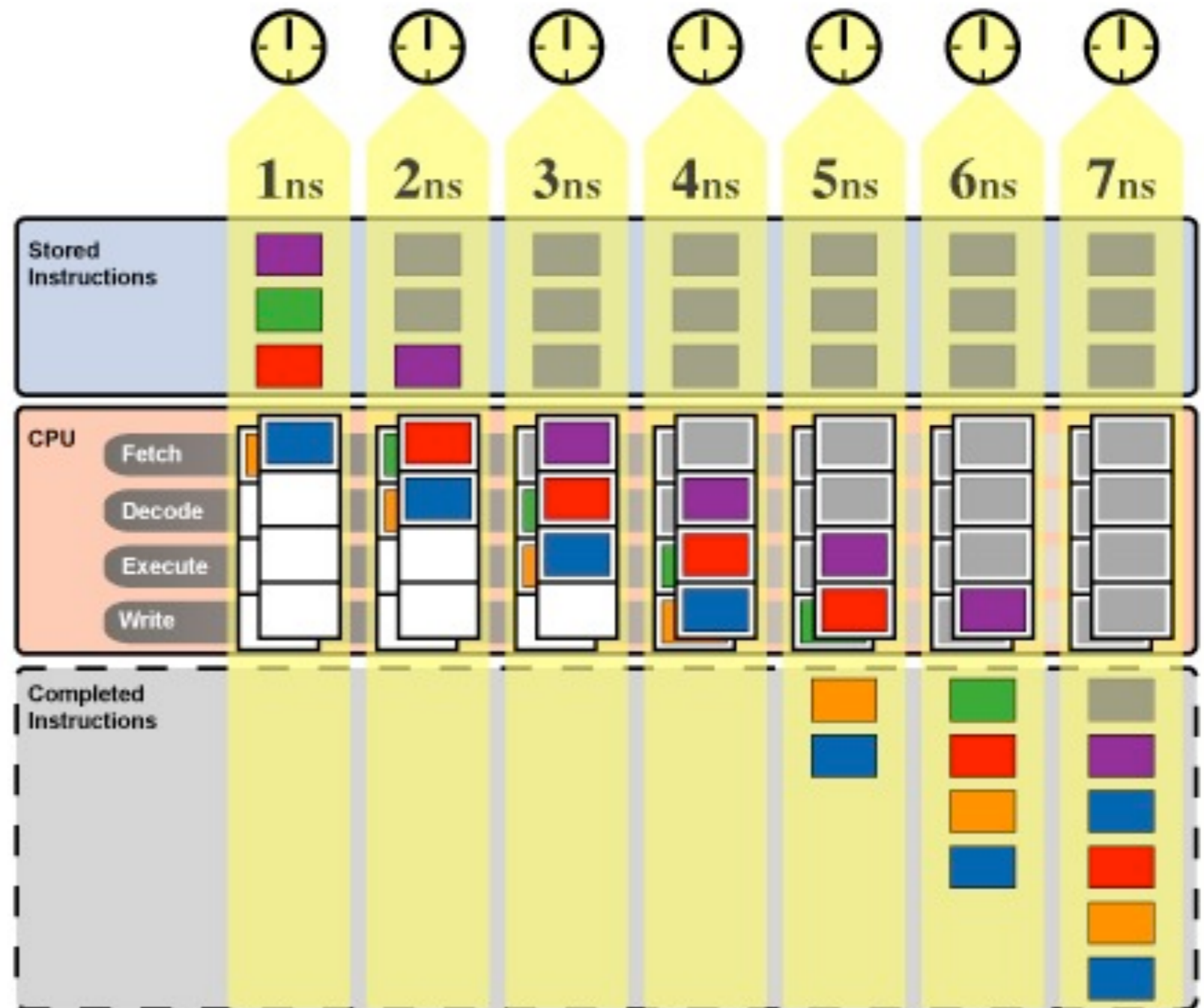
# Superscalar processors

- CPU has multiple functional units
  - Each can deal with different kinds of operations
  - Some overlap, e.g., most functional units can do integer arithmetic
- Each functional unit has its own pipeline
  - ⇒ Multiple pipelines executing in parallel



# Superscalar processors

- In one cycle can issue different instructions to different functional units
  - Hence in one cycle can complete more than one operation
  - Thus, “superscalar”
- Not the same as multicore!



# Multicore processors

- Each chip contains multiple separate processor cores
- Each core can run completely different code
- To take advantage (in a single program) of multiple cores must write *concurrent code*. More on this later in course...



# Today

- Program optimization
  - Overview
  - Code motion
  - Strength reduction
  - Common subexpressions
  - Optimization blockers
    - Procedure calls
    - Aliasing
  - Understanding modern processors
  - Loop unrolling
  - Tail recursion
  - Summary

# Loop unrolling

- Reduce number of iterations of loop by doing more work each iteration
  - Reduces number of loop index/comparison operations
  - Further transformations can enable additional speedup.

```
int prod_array(int *a, int n) {
    int i, result=1;
    for (i = 0; i < n; i++) {
        result *= a[i];
    }
    return result;
}
```

```
/* Note: assuming n is even! */
int prod_array2(int *a, int n) {
    int i, tmp1=1, tmp2=1;
    for (i = 0; i < n; i+=2) {
        result *= a[i];
        result *= a[i+1];
    }
    return result;
}
```

# Enhancing parallelism

- In unrolled version, multiplications must occur in sequence
  - Why?
- What if we used two accumulators?
  - No dependency between two multiplications
  - Can be run in parallel

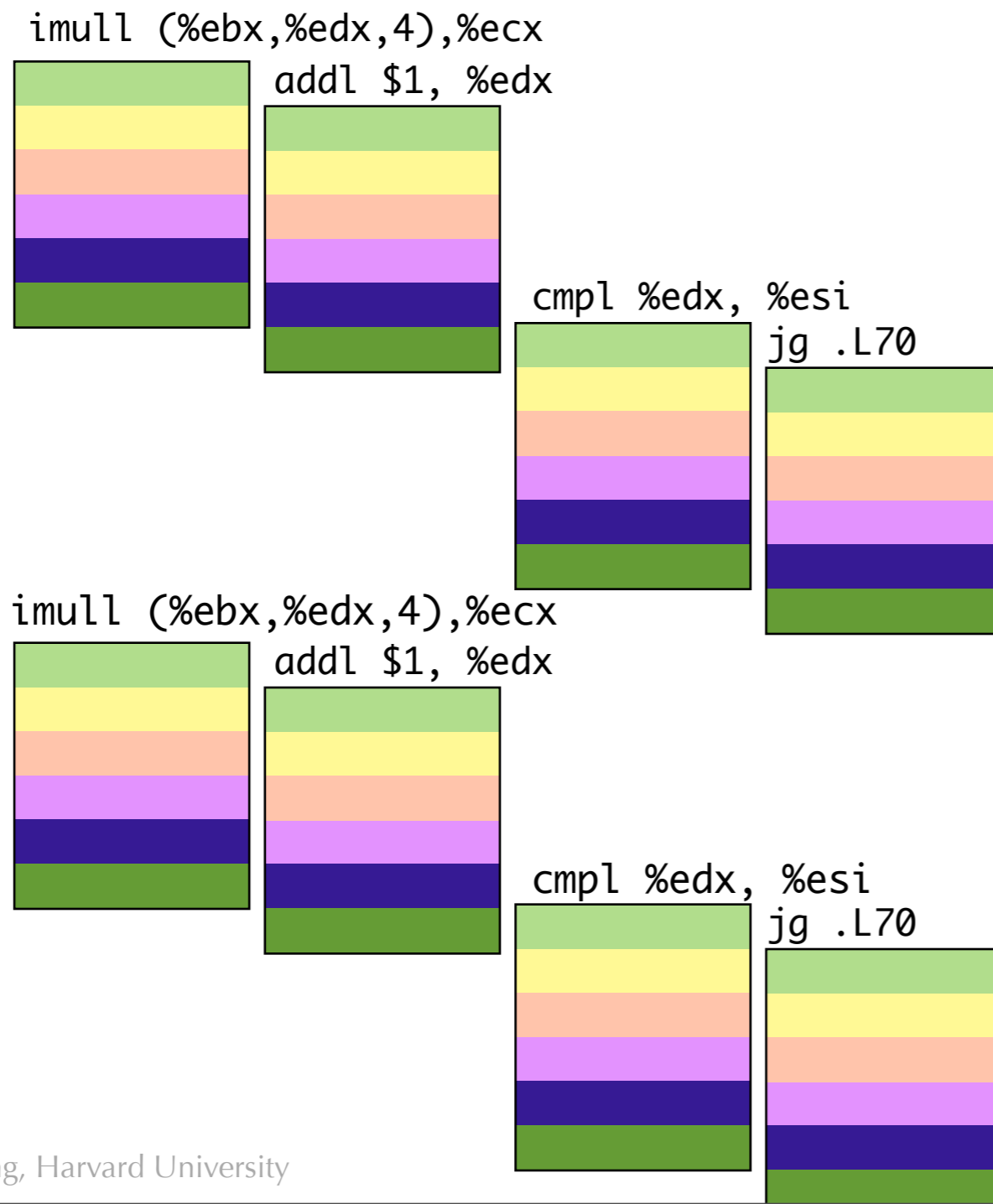
```
/* Note: assuming n is even! */
int prod_array2(int *a, int n) {
    int i, result=1;
    for (i = 0; i < n; i+=2) {
        result *= a[i];
        result *= a[i+1];
    }
    return result;
}
```

```
/* Note: assuming n is even! */
int prod_array2(int *a, int n) {
    int i, tmp1=1, tmp2=1;
    for (i = 0; i < n; i+=2) {
        tmp1 *= a[i];
        tmp2 *= a[i+1];
    }
    return tmp1 * tmp2;
}
```



# Visualizing unrolling

Time



# Original loop

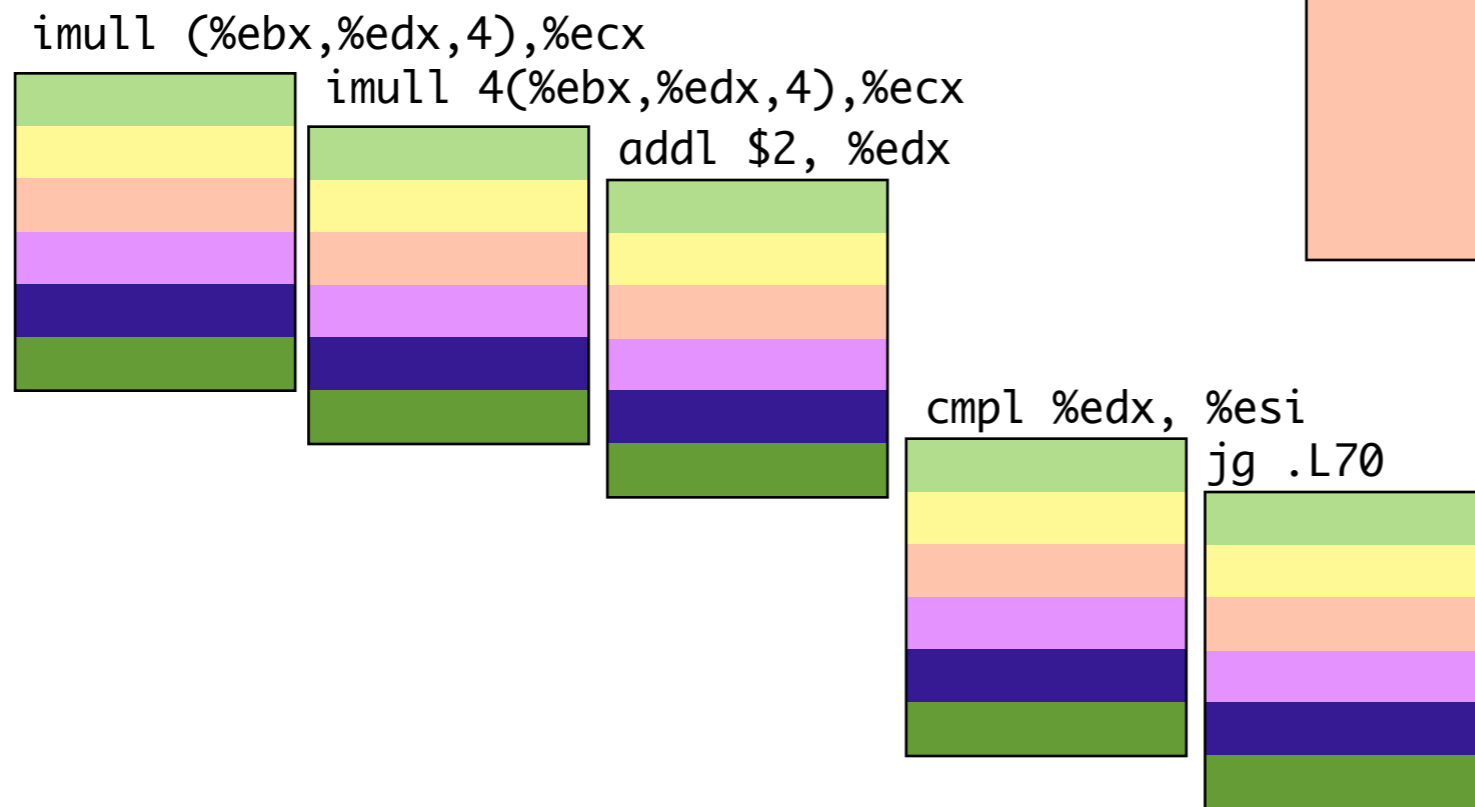
.L77:

```
imull (%ebx,%edx,4), %ecx  
addl $1, %edx  
cmpl %edx, %esi  
jg .L77
```

- In original loop, need to stall pipeline each iteration due to dependencies

# Visualizing unrolling

Time



# Unrolled loop

```
.L77:  
    imull    (%ebx,%edx,4), %ecx  
    imull    4(%ebx, %edx, 4), %eax  
    addl    $2, %edx  
    cml     %edx, %esi  
    jg     .L77
```

- Still need to stall pipeline each iteration, but each iteration does 2 multiplications!
- More parallelism for each iteration

# Today

- Program optimization
  - Overview
  - Code motion
  - Strength reduction
  - Common subexpressions
  - Optimization blockers
    - Procedure calls
    - Aliasing
  - Understanding modern processors
  - Loop unrolling
  - Tail recursion
  - Summary

# What does the stack look like?

```
int factorial(int fact_so_far, int n) {  
    if (n == 1) return fact_so_far;  
    return factorial(fact_so_far*n, n-1);  
}
```

- Suppose we called `factorial(1, 10)`
  - What does the stack look like?



# What does the stack look like?

```
int factorial(int fact_so_far, int n) {
    if (n == 1) return fact_so_far;
    return factorial(fact_so_far*n, n-1);
}
```

- Result of function returned in **%eax**
- In base case, we moved **fact\_so\_far** arg off stack and into **%eax**
- In the recursive case, after the **call** instruction, result is already in **%eax**
  - No further computation to do!
  - Just need to clean up stack

```
factorial:
    pushl   %ebp
    movl   %esp, %ebp
    subl   $12, %esp
    movl   8(%ebp), %eax
    movl   12(%ebp), %edx
    cmpl   $1, %edx
    je     .L4
    imull  %edx, %eax
    subl   $1, %edx
    movl   %edx, 4(%esp)
    movl   %eax, (%esp)
    call  factorial
.L4:
    leave
    ret
```

# Tail recursion

```
int factorial(int fact_so_far, int n) {  
    if (n == 1) return fact_so_far;  
    return factorial(fact_so_far*n, n-1);  
}
```

- Function factorial is **tail recursive**
  - A special case of recursion
  - The last thing it does is call itself, and return the result
  - No more computation to be done once the recursive call finishes
- No need for stack frame once recursive call finishes
- So why bother allocating a new stack frame for the recursive call?

```
factorial:  
    pushl    %ebp  
    movl    %esp, %ebp  
    subl    $12, %esp  
    movl    8(%ebp), %eax  
    movl    12(%ebp), %edx  
    cmpl    $1, %edx  
    je     .L4  
    imull   %edx, %eax  
    subl    $1, %edx  
    movl    %edx, 4(%esp)  
    movl    %eax, (%esp)  
    call   factorial  
  
.L4:  
    leave  
    ret
```

# Tail-call elimination

- **Tail-call elimination** removes the tail call
  - Converts it to a jmp
  - Fewer stack frames, fewer function call overhead

```
factorial:
    pushl   %ebp
    movl   %esp, %ebp
    subl   $12, %esp
    movl   8(%ebp), %eax
    movl   12(%ebp), %edx

.L5:
    cmpl   $1, %edx
    je     .L4
    imull  %edx, %eax
    subl   $1, %edx
    jmp    .L5

.L4:
    leave
    ret
```

```
factorial:
    pushl   %ebp
    movl   %esp, %ebp
    subl   $12, %esp
    movl   8(%ebp), %eax
    movl   12(%ebp), %edx
    cmpl   $1, %edx
    je     .L4
    imull  %edx, %eax
    subl   $1, %edx
    movl   %edx, 4(%esp)
    movl   %eax, (%esp)
    call   factorial

.L4:
    leave
    ret
```



# Finding tail recursion

- gcc pretty good at finding tail recursion
- Natural implementation of factorial is not tail recursive

```
int fact(int n) {  
    if (n == 1) return 1;  
    return n*fact(n-1);  
}
```

- After the recursive call returns, it performs computation on the result
- However, gcc compiles it to machine code with a `jmp` instead of a `call`



# Optimization recap

- Write code to help the compiler, and CPU, do their jobs well.
  - Remember: The compiler has to be conservative, but you might know better.
- High-level design
  - Choose appropriate algorithms and data structures
- Basic coding principles
  - Avoid optimization blockers
  - Eliminate unnecessary function calls and memory references
- Low-level optimization
  - Unroll loops to reduce overhead and enable further optimizations
  - Find ways to increase instruction-level parallelism
  - Code motion:
    - Move constant expressions outside of loops
    - Especially in the presence of function calls
  - Strength reduction
    - Use less expensive operations/functions when possible (Though, most compilers will do this for you!)

# Caveats

- Does this mean you should write crazy, convoluted, repetitive, but high performing code?
- Probably not.
- Need to balance maintainability/readability with performance
- Always clearly comment when you are doing something funky
  - State your assumptions: someone may change your code later and break it in subtle ways

# Caveats

*There is no doubt that the grail of efficiency leads to abuse. Programmers waste enormous amounts of time thinking about, or worrying about, the speed of noncritical parts of their programs, and these attempts at efficiency actually have a strong negative impact when debugging and maintenance are considered. We should forget about small efficiencies, say about 97% of the time: **premature optimization is the root of all evil.***

*Yet we should not pass up our opportunities in that critical 3%.*



Donald Knuth

*Structured Programming with goto Statements*  
Computing Surveys, Vol 6, No 4, December 1974

# How to find the 3%...

- Identifying and eliminating performance bottlenecks
  - Use a program profiler to find out where your program is spending its time
    - e.g., gprof
  - Speed up of program depends on how much you improved performance of component, and how significant component is

# Performance lab

- In a few weeks, we will release a performance lab
  - Optional lab, to explore performance improvement
  - Not part of course assessment
  - More details coming...