# Organisation und Architektur von Rechnern

Lecture 15

**Instructor:**

Reinhard v. Hanxleden

http://www.informatik.uni-kiel.de/rtsys/teaching/v-sysinf2

*These slides are used with kind permission from the Carnegie Mellon University*

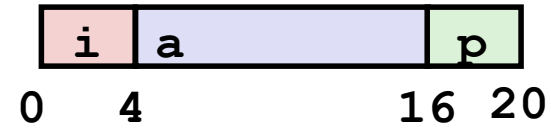# The 5 Minute Review Session

1. **How is the PC predicted (different cases)?**

2. **What are *pipelining hazards*?**

3. **What is *pipeline stalling*?**

4. **How does the pipeline handle mispredicted branches?**

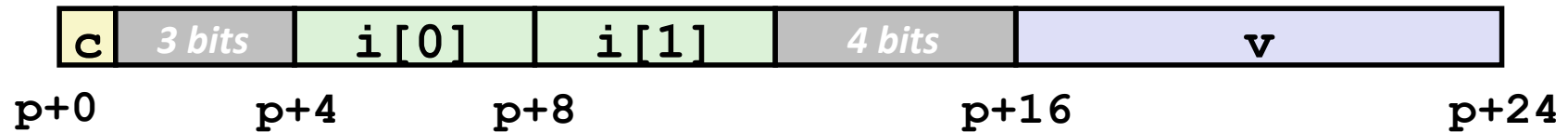5. **What is *data forwarding*, why is it used?**

# Last Time

- **Structures**

```
struct rec {
    int i;
    int a[3];
    int *p;
};
```
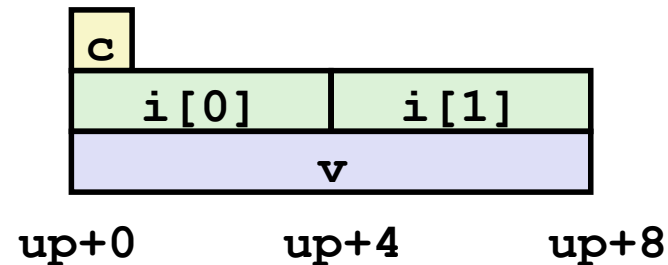
**Memory Layout**

| i | a | | p |
|---|---|---|---|

0   4          16 20

- **Alignment**

```
struct S1 {
    char c;
    int i[2];
    double v;
} *p;
```

| c | *3 bits* | i[0] | i[1] | *4 bits* | v |
|---|---|---|---|---|---|

p+0      p+4      p+8            p+16          p+24

- **Unions**

```
union U1 {
    char c;
    int i[2];
    double v;
} *up;
```

| c |
|---|

| i[0] | i[1] |
|---|---|

| v |
|---|

up+0      up+4      up+8

# Last Time

- **Floating point**
  - x87 (getting obsolete)

    %st(3)
    %st(2)
    %st(1)
    %st(0)

  **128 bit = 2 doubles = 4 singles**

  - x86-64 (SSE3 and later)

    %xmm0

    %xmm15

  - Vector mode and scalar mode

    addps $+$

    addss $+$

# Today

- **Memory layout**

- **Program optimization**
  - Overview
  - Removing unnecessary procedure calls
  - Code motion/precomputation
  - Strength reduction
  - Sharing of common subexpressions
  - Optimization blocker: Procedure calls

# IA32 Linux Memory Layout

- **Stack**
  - Runtime stack (8MB limit)

- **Heap**
  - Dynamically allocated storage
  - When call `malloc()`, `calloc()`, `new()`

- **Data**
  - Statically allocated data
  - E.g., arrays & strings declared in code

- **Text**
  - Executable machine instructions
  - Read-only

FF

Stack

8MB

Heap

Data

Text

Upper 2 hex digits
= 8 bits of address

08

00

6

# Memory Allocation Example

```
char big_array[1<<24];   /*   16 MB */
char huge_array[1<<28]; /* 256 MB */

int beyond;
char *p1, *p2, *p3, *p4;

int useless() {  return 0; }

int main()
{
 p1 = malloc(1 <<28);  /* 256 MB */
 p2 = malloc(1 << 8);  /* 256 B  */
 p3 = malloc(1 <<28);  /* 256 MB */
 p4 = malloc(1 << 8);  /* 256 B  */
 /* Some print statements ... */
}
```

*Where does everything go?*

FF

| Stack |
| Heap |
| Data |
| Text |

08
00

7

# IA32 Example Addresses

*address range ~$2^{32}$*

| | |
|---|---|
| `$esp` | `0xffffbcd0` |
| `p3` | `0x65586008` |
| `p1` | `0x55585008` |
| `p4` | `0x1904a110` |
| `p2` | `0x1904a008` |
| `&p2` | `0x18049760` |
| `beyond` | `0x08049744` |
| `big_array` | `0x18049780` |
| `huge_array` | `0x08049760` |
| `main()` | `0x080483c6` |
| `useless()` | `0x08049744` |
| `final malloc()` | `0x006be166` |

`malloc()` **is dynamically linked**
**address determined at runtime**

8

*not drawn to scale*

FF

Stack

80

Heap

Data

Text

08
00

# x86-64 Example Addresses

*address range ~$2^{47}$*

*not drawn to scale*

| | |
|---|---|
| `$rsp` | `0x7ffffff8d1f8` |
| `p3` | `0x2aaabaadd010` |
| `p1` | `0x2aaaaaadc010` |
| `p4` | `0x000011501120` |
| `p2` | `0x000011501010` |
| `&p2` | `0x000010500a60` |
| `beyond` | `0x000000500a44` |
| `big_array` | `0x000010500a80` |
| `huge_array` | `0x000000500a50` |
| `main()` | `0x000000400510` |
| `useless()` | `0x000000400500` |
| `final malloc()` | `0x00386ae6a170` |

`malloc()` **is dynamically linked**
**address determined at runtime**

**9**

Memory layout diagram:
- 00007F — Stack (arrow pointing down)
- (arrow pointing up)
- 000030 — Heap
- Data
- Text
- 000000

# C operators

| Operators | Associativity |
|---|---|
| `()` `[]` `->` `.` | left to right |
| `!` `~` `++` `--` `+` `-` `*` `&` `(type)` `sizeof` | right to left |
| `*` `/` `%` | left to right |
| `+` `-` | left to right |
| `<<` `>>` | left to right |
| `<` `<=` `>` `>=` | left to right |
| `==` `!=` | left to right |
| `&` | left to right |
| `^` | left to right |
| `|` | left to right |
| `&&` | left to right |
| `||` | left to right |
| `?:` | right to left |
| `=` `+=` `-=` `*=` `/=` `%=` `&=` `^=` `!=` `<<=` `>>=` | right to left |
| `,` | left to right |

- **`->` has very high precedence**
- **`()` has very high precedence**
- **monadic `*` just below**

# C Pointer Declarations: Test Yourself!

`int *p`                 p is a pointer to int

`int *p[13]`

`int *(p[13])`

`int **p`                p is a pointer to a pointer to an int

`int (*p)[13]`

`int *f()`               f is a function returning a pointer to int

`int (*f)()`             f is a pointer to a function returning int

`int (*(*f())[13])()`

`int (*(*x[3])())[5]`    x is an array[3] of pointers to functions
                         returning pointers to array[5] of ints

# C Pointer Declarations (Check out guide)

`int *p`                    p is a pointer to int

`int *p[13]`                p is an array[13] of pointer to int

`int *(p[13])`              p is an array[13] of pointer to int

`int **p`                   p is a pointer to a pointer to an int

`int (*p)[13]`              p is a pointer to an array[13] of int
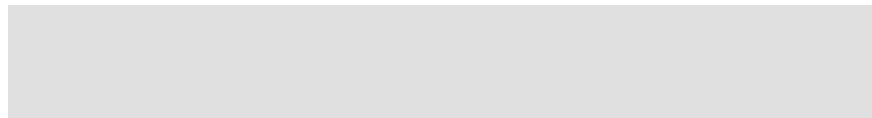
`int *f()`                  f is a function returning a pointer to int

`int (*f)()`                f is a pointer to a function returning int

`int (*(*f())[13])()`       f is a function returning ptr to an array[13] of pointers to functions returning int

`int (*(*x[3])())[5]`       x is an array[3] of pointers to functions returning pointers to array[5] of ints

# Avoiding Complex Declarations

- Use `typedef` to build up the declaration

- Instead of `int (*(*x[3])())[5]`:

  ```
  typedef int fiveints[5];

  typedef fiveints* p5i;

  typedef p5i (*f_of_p5is)();

  f_of_p5is x[3];
  ```

- `x` is an array of 3 elements, each of which is a pointer to a function returning an array of 5 ints

# Today

- **Memory layout**

- **Buffer overflow, worms, and viruses**

- **Program optimization**
  - Overview
  - Removing unnecessary procedure calls
  - Code motion/precomputation
  - Strength reduction
  - Sharing of common subexpressions
  - Optimization blocker: Procedure calls

# Internet Worm and IM War

- **November, 1988**
  - Internet Worm attacks thousands of Internet hosts.
  - How did it happen?

# Internet Worm and IM War

- **November, 1988**
  - Internet Worm attacks thousands of Internet hosts.
  - How did it happen?
- **July, 1999**
  - Microsoft launches MSN Messenger (instant messaging system).
  - Messenger clients can access popular AOL Instant Messaging Service (AIM) servers

# Internet Worm and IM War (cont.)

- **August 1999**
  - Mysteriously, Messenger clients can no longer access AIM servers.
  - Microsoft and AOL begin the IM war:
    - AOL changes server to disallow Messenger clients
    - Microsoft makes changes to clients to defeat AOL changes.
    - At least 13 such skirmishes.
  - How did it happen?

- **The Internet Worm and AOL/Microsoft War were both based on *stack buffer overflow* exploits!**
  - many Unix functions do not check argument sizes.
  - allows target buffers to overflow.

# String Library Code

- **Implementation of Unix function `gets()`**

```
/* Get string from stdin */
char *gets(char *dest)
{
    int c = getchar();
    char *p = dest;
    while (c != EOF && c != '\n') {
        *p++ = c;
        c = getchar();
    }
    *p = '\0';
    return dest;
}
```

- No way to specify limit on number of characters to read

- **Similar problems with other Unix functions**

- **`strcpy`**: Copies string of arbitrary length

- **`scanf, fscanf, sscanf,`** when given `%s` conversion specification

# Vulnerable Buffer Code

```
/* Echo Line */
void echo()
{
    char buf[4];   /* Way too small! */
    gets(buf);
    puts(buf);
}
```

```
int main()
{
  printf("Type a string:");
  echo();
  return 0;
}
```

```
unix>./bufdemo
Type a string:1234567
1234567
```

```
unix>./bufdemo
Type a string:12345678
Segmentation Fault
```

```
unix>./bufdemo
Type a string:123456789ABC
Segmentation Fault
```

# Buffer Overflow Disassembly

```
080484f0 <echo>:
 80484f0:   55                      push    %ebp
 80484f1:   89 e5                   mov     %esp,%ebp
 80484f3:   53                      push    %ebx
 80484f4:   8d 5d f8                lea     0xfffffff8(%ebp),%ebx
 80484f7:   83 ec 14                sub     $0x14,%esp
 80484fa:   89 1c 24                mov     %ebx,(%esp)
 80484fd:   e8 ae ff ff ff          call    80484b0 <gets>
 8048502:   89 1c 24                mov     %ebx,(%esp)
 8048505:   e8 8a fe ff ff          call    8048394 <puts@plt>
 804850a:   83 c4 14                add     $0x14,%esp
 804850d:   5b                      pop     %ebx
 804850e:   c9                      leave
 804850f:   c3                      ret

 80485f2:   e8 f9 fe ff ff          call    80484f0 <echo>
 80485f7:   8b 5d fc                mov 0xfffffffc(%ebp),%ebx
 80485fa:   c9                      leave
 80485fb:   31 c0                   xor     %eax,%eax
 80485fd:   c3                      ret
```

# Buffer Overflow Stack

*Before call to gets*

| |
|---|
| Stack Frame for **main** |
| |
| Return Address |
| Saved %**ebp** ← %**ebp** |
| |
| [3] [2] [1] [0] **buf** |
| Stack Frame for **echo** |

```
/* Echo Line */
void echo()
{
    char buf[4];   /* Way too small! */
    gets(buf);
    puts(buf);

}
```

```
echo:
    pushl %ebp              # Save %ebp on stack
    movl  %esp, %ebp
    pushl %ebx              # Save %ebx
    leal  -8(%ebp),%ebx     # Compute buf as %ebp-8
    subl  $20, %esp         # Allocate stack space
    movl  %ebx, (%esp)      # Push buf on stack
    call  gets              # Call gets

    . . .
```

# Buffer Overflow Stack Example

```
unix> gdb bufdemo
(gdb) break echo
Breakpoint 1 at 0x8048583
(gdb) run
Breakpoint 1, 0x8048583 in echo ()
(gdb) print /x $ebp
$1 = 0xffffc638
(gdb) print /x *(unsigned *)$ebp
$2 = 0xffffc658
(gdb) print /x *((unsigned *)$ebp + 1)
$3 = 0x80485f7
```

**Before call to gets**

| Stack Frame for **main** |
| --- |
| Return Address |
| Saved %**ebp** |
| |
| [3][2][1][0]  **buf** |
| Stack Frame for **echo** |

**Before call to gets**

0xffffc658

| Stack Frame for **main** | | | |
| --- | --- | --- | --- |
| 08 | 04 | 85 | f7 |
| ff | ff | c6 | 58 |

0xffffc638

| | | | |
| --- | --- | --- | --- |
| xx | xx | xx | xx |

**buf**

| Stack Frame for **echo** | | | |

```
80485f2: call 80484f0 <echo>
80485f7: mov  0xfffffffc(%ebp),%ebx # Return Point
```

22

# Buffer Overflow Example #1

*Before call to gets*

| | | | |
|---|---|---|---|
| Stack Frame for **main** | | | | `0xffffc658` |

| 08 | 04 | 85 | f7 |
|---|---|---|---|
| ff | ff | c6 | 58 |

`0xffffc638`

| | | | |
|---|---|---|---|
| xx | xx | xx | xx |

**buf**

| | | | |
|---|---|---|---|
| Stack Frame for **echo** | | | |

*Input 1234567*

| | | | |
|---|---|---|---|
| Stack Frame for **main** | | | | `0xffffc658` |

| 08 | 04 | 85 | f7 |
|---|---|---|---|
| ff | ff | c6 | 58 |

`0xffffc638`

| 00 | 37 | 36 | 35 |
|---|---|---|---|
| 34 | 33 | 32 | 31 |

**buf**

| | | | |
|---|---|---|---|
| Stack Frame for **echo** | | | |

**Overflow buf, but no problem**

# Buffer Overflow Example #2

**Before call to gets**

| | | | |
|---|---|---|---|
| Stack Frame for **main** | | | | `0xffffc658`

| 08 | 04 | 85 | f7 |
|---|---|---|---|
| ff | ff | c6 | 58 |

`0xffffc638`

| | | | |
|---|---|---|---|
| xx | xx | xx | xx |

`buf`

| | | | |
|---|---|---|---|
| Stack Frame for **echo** | | | |

**Input 12345678**

| | | | |
|---|---|---|---|
| Stack Frame for **main** | | | | `0xffffc658`

| 08 | 04 | 85 | f7 |
|---|---|---|---|
| ff | ff | c6 | 00 |
| 38 | 37 | 36 | 35 |
| 34 | 33 | 32 | 31 |

`0xffffc638`

`buf`

| | | | |
|---|---|---|---|
| Stack Frame for **echo** | | | |

**Base pointer corrupted**

```
 . . .
804850a:  83 c4 14    add      $0x14,%esp   # deallocate space
804850d:  5b          pop      %ebx         # restore %ebx
804850e:  c9          leave                 # movl %ebp, %esp; popl %ebp
804850f:  c3          ret                   # Return
```
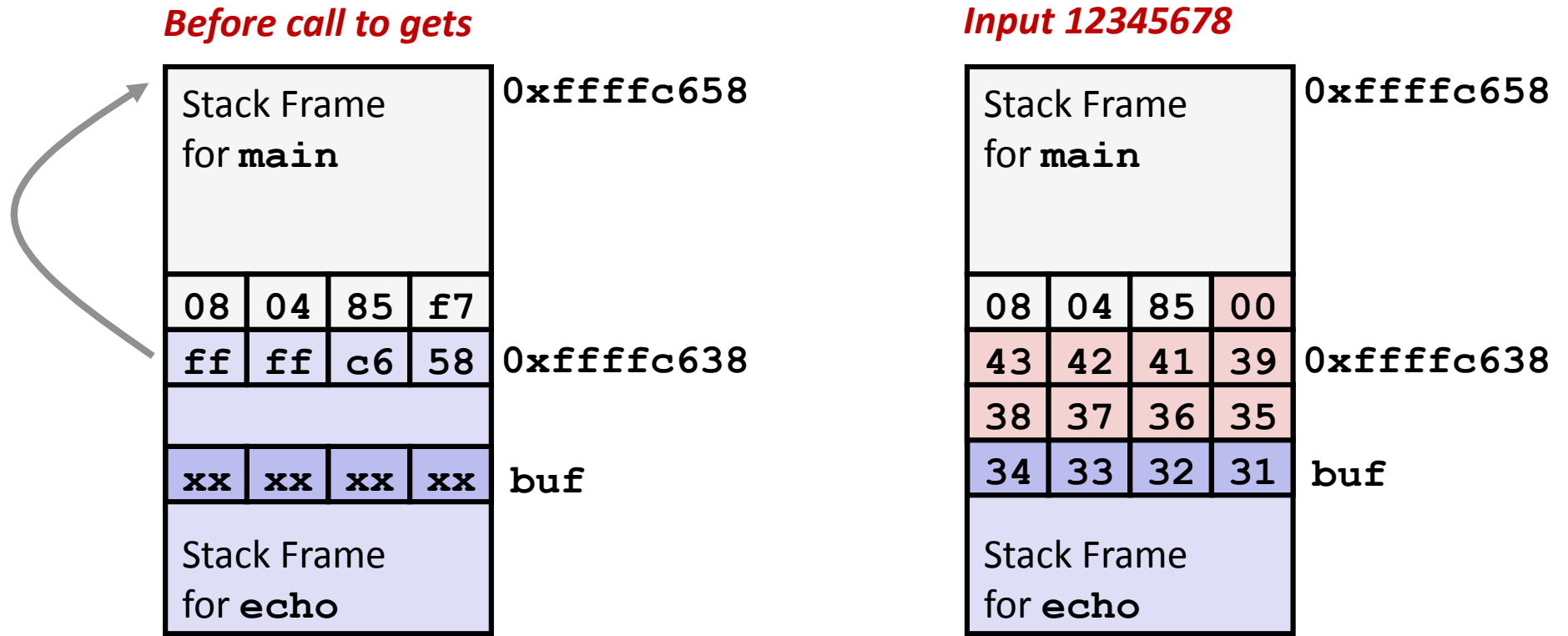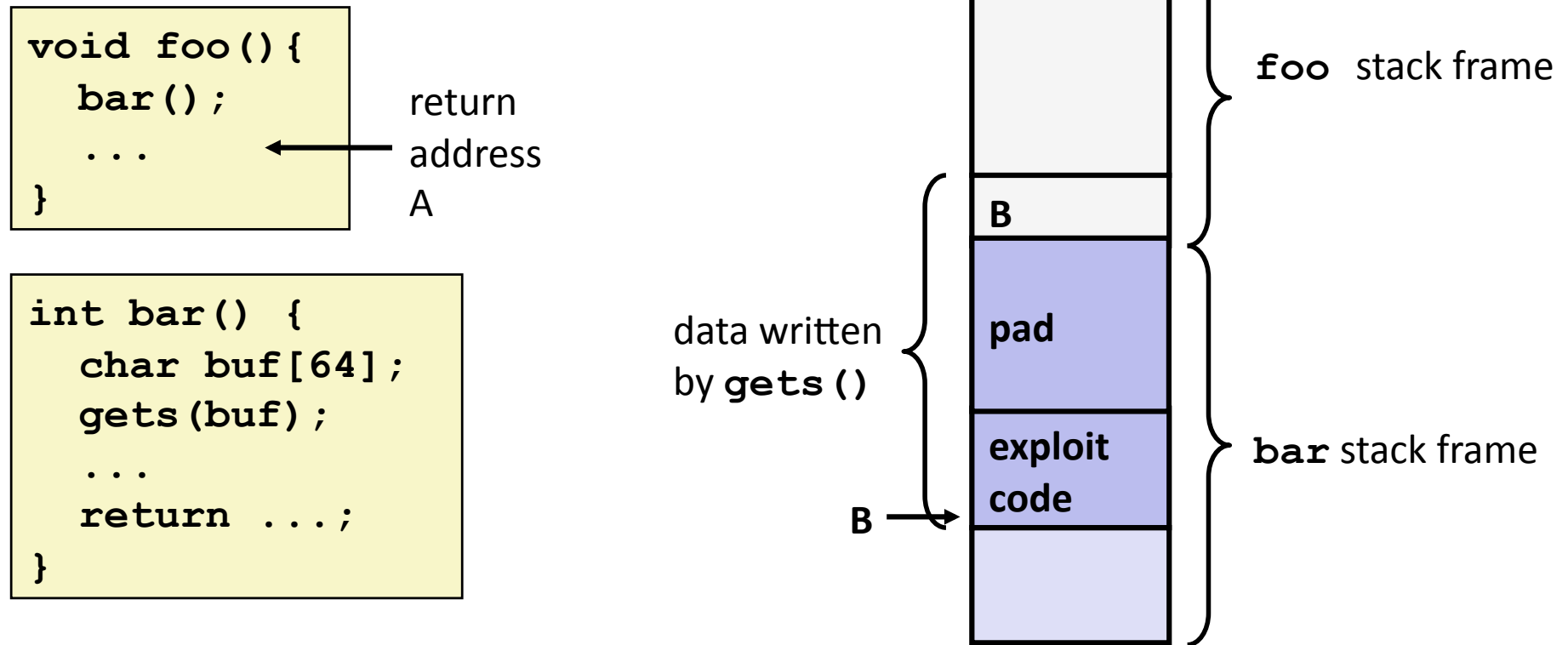
# Buffer Overflow Example #3

**Before call to gets**

| | | | |
|---|---|---|---|
| Stack Frame for **main** | | | |

`0xffffc658`

| 08 | 04 | 85 | f7 |
|----|----|----|----|
| ff | ff | c6 | 58 |

`0xffffc638`

| | | | |
|---|---|---|---|

| xx | xx | xx | xx |
|----|----|----|----|

`buf`

| | | | |
|---|---|---|---|
| Stack Frame for **echo** | | | |

**Input 12345678**

| | | | |
|---|---|---|---|
| Stack Frame for **main** | | | |

`0xffffc658`

| 08 | 04 | 85 | 00 |
|----|----|----|----|
| 43 | 42 | 41 | 39 |

`0xffffc638`

| 38 | 37 | 36 | 35 |
|----|----|----|----|
| 34 | 33 | 32 | 31 |

`buf`

| | | | |
|---|---|---|---|
| Stack Frame for **echo** | | | |

**Return address corrupted**

```
80485f2:  call  80484f0 <echo>
80485f7:  mov   0xfffffffc(%ebp),%ebx # Return Point
```

# Malicious Use of Buffer Overflow

Stack after call to `gets()`

```
void foo(){
   bar();
   ...
}
```

return
address
A

```
int bar() {
   char buf[64];
   gets(buf);
   ...
   return ...;
}
```

data written
by `gets()`

B

foo stack frame

B

pad

exploit
code

bar stack frame

- Input string contains byte representation of executable code
- Overwrite return address with address of buffer
- When `bar()` executes `ret`, will jump to exploit code

# Exploits Based on Buffer Overflows

- ***Buffer overflow bugs allow remote machines to execute arbitrary code on victim machines***

- **Internet worm**
  - Early versions of the finger server (fingerd) used `gets()` to read the argument sent by the client:
    - `finger droh@cs.cmu.edu`
  - Worm attacked fingerd server by sending phony argument:
    - `finger "exploit-code  padding  new-return-address"`
    - exploit code: executed a root shell on the victim machine with a direct TCP connection to the attacker.

# Exploits Based on Buffer Overflows

- *Buffer overflow bugs allow remote machines to execute arbitrary code on victim machines*

- **IM War**
  - AOL exploited existing buffer overflow bug in AIM clients
  - exploit code: returned 4-byte signature (the bytes at some location in the AIM client) to server.
  - When Microsoft changed code to match signature, AOL changed signature location.

```
Date: Wed, 11 Aug 1999 11:30:57 -0700 (PDT)
From: Phil Bucking <philbucking@yahoo.com>
Subject: AOL exploiting buffer overrun bug in their own software!
To: rms@pharlap.com

Mr. Smith,

I am writing you because I have discovered something that I think you
might find interesting because you are an Internet security expert with
experience in this area. I have also tried to contact AOL but received
no response.

I am a developer who has been working on a revolutionary new instant
messaging client that should be released later this year.
...
It appears that the AIM client has a buffer overrun bug. By itself
this might not be the end of the world, as MS surely has had its share.
But AOL is now *exploiting their own buffer overrun bug* to help in
its efforts to block MS Instant Messenger.
....
Since you have significant credibility with the press I hope that you
can use this information to help inform people that behind AOL's
friendly exterior they are nefariously compromising peoples' security.

Sincerely,
Phil Bucking
Founder, Bucking Consulting
philbucking@yahoo.com
```
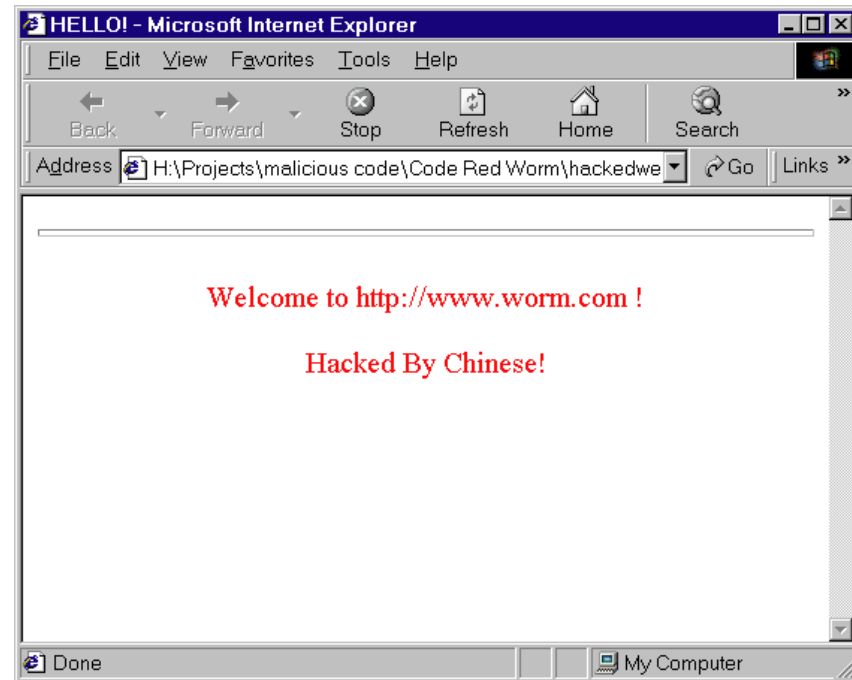
*It was later determined that this email originated from within Microsoft!*

# Code Red Worm

- **History**
  - June 18, 2001. Microsoft announces buffer overflow vulnerability in IIS Internet server
  - July 19, 2001. over 250,000 machines infected by new virus in 9 hours
  - White house must change its IP address. Pentagon shut down public WWW servers for day

- **When We Set Up CS:APP Web Site**
  - Received strings of form

```
GET /default.ida?
   NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN....NNNNNNN
   NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
   %u9090%u6858%ucbd3%u7801%u9090%u6858%ucbd3%u7801%u
   9090%u6858%ucbd3%u7801%u9090%u9090%u8190%u00c3%u00
   03%u8b00%u531b%u53ff%u0078%u0000%u00=a
HTTP/1.0" 400 325 "-" "-"
```

# Code Red Exploit Code

- **Starts 100 threads running**
- **Spread self**
  - Generate random IP addresses & send attack string
  - Between 1st & 19th of month
- **Attack www.whitehouse.gov**
  - Send 98,304 packets; sleep for 4-1/2 hours; repeat
    - Denial of service attack
  - Between 21st & 27th of month
- **Deface server's home page**
  - After waiting 2 hours



Welcome to http://www.worm.com !

Hacked By Chinese!

# Code Red Effects

- **Later Version Even More Malicious**
  - Code Red II
  - As of April, 2002, over 18,000 machines infected
  - Still spreading

- **Paved Way for NIMDA**
  - Variety of propagation methods
  - One was to exploit vulnerabilities left behind by Code Red II

- **ASIDE (security flaws start at home)**
  - .rhosts used by Internet Worm
  - Attachments used by MyDoom  (1 in 6 emails Monday morning!)

# Avoiding Overflow Vulnerability

```
/* Echo Line */
void echo()
{
    char buf[4];   /* Way too small!
*/
    fgets(buf, 4, stdin);
    puts(buf);
}
```

- **Use library routines that limit string lengths**
  - **fgets** instead of **gets**
  - strncpy instead of **strcpy**
  - Don't use **scanf** with **%s** conversion specification
    - Use **fgets** to read the string
    - Or use **%ns** where **n** is a suitable integer

# System-Level Protections

- **Randomized stack offsets**
  - At start of program, allocate random amount of space on stack
  - Makes it difficult for hacker to predict beginning of inserted code

- **Nonexecutable code segments**
  - In traditional x86, can mark region of memory as either "read-only" or "writeable"
    - Can execute anything readable
  - Add explicit "execute" permission

```
unix> gdb bufdemo
(gdb) break echo

(gdb) run
(gdb) print /x $ebp
$1 = 0xffffc638

(gdb) run
(gdb) print /x $ebp
$2 = 0xffffbb08

(gdb) run
(gdb) print /x $ebp
$3 = 0xffffc6a8
```
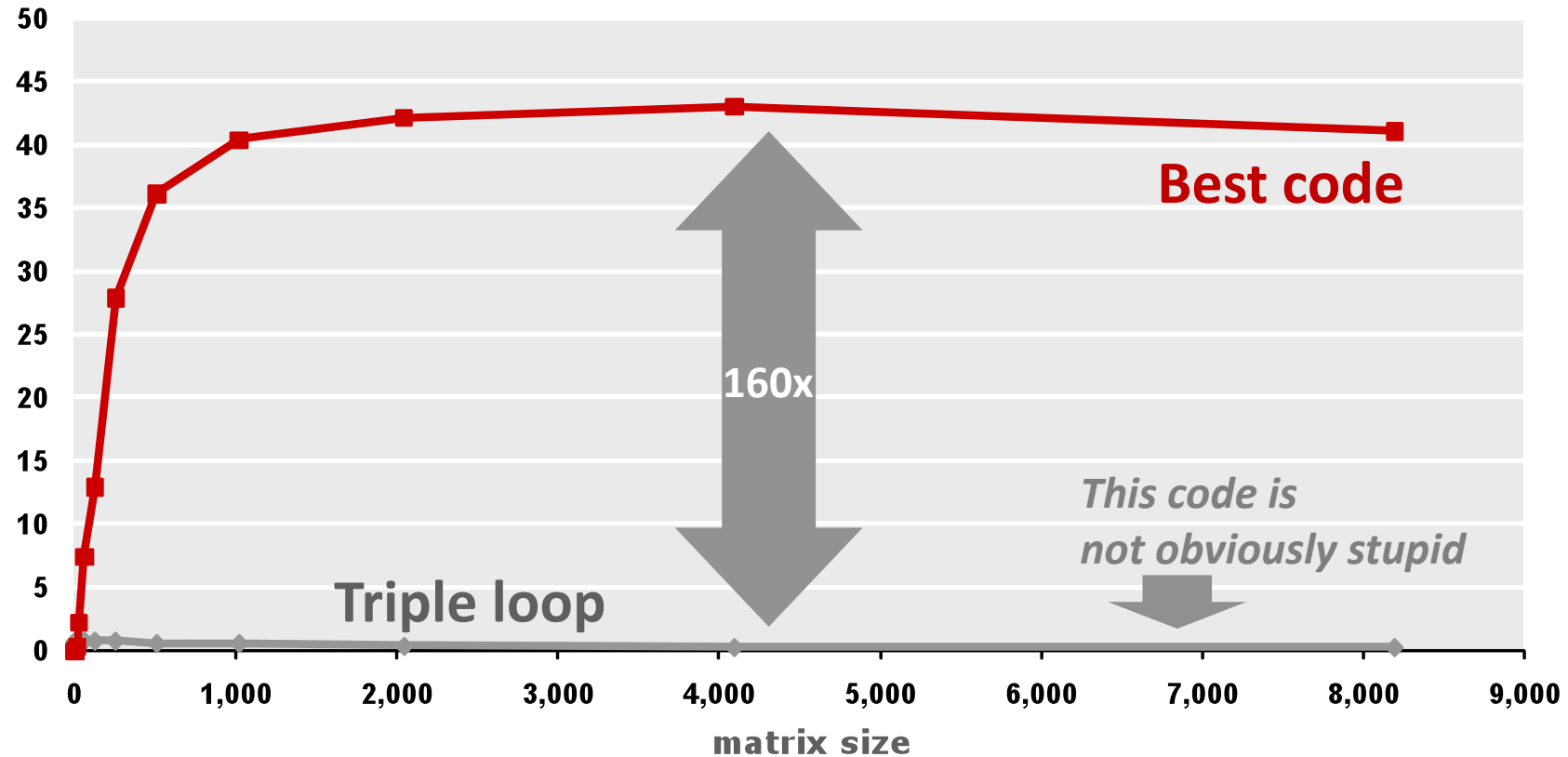
# Worms and Viruses

- **Worm: A program that**
  - Can run by itself
  - Can propagate a fully working version of itself to other computers

- **Virus: Code that**
  - Add itself to other programs
  - Cannot run independently

- **Both are (usually) designed to spread among computers and to wreak havoc**

# Today

- **Memory layout**

- **Program optimization**
  - **Overview**
  - Removing unnecessary procedure calls
  - Code motion/precomputation
  - Strength reduction
  - Sharing of common subexpressions
  - Optimization blocker: Procedure calls

# Example Matrix Multiplication

**Matrix-Matrix Multiplication (MMM) on 2 x Core 2 Duo 3 GHz**
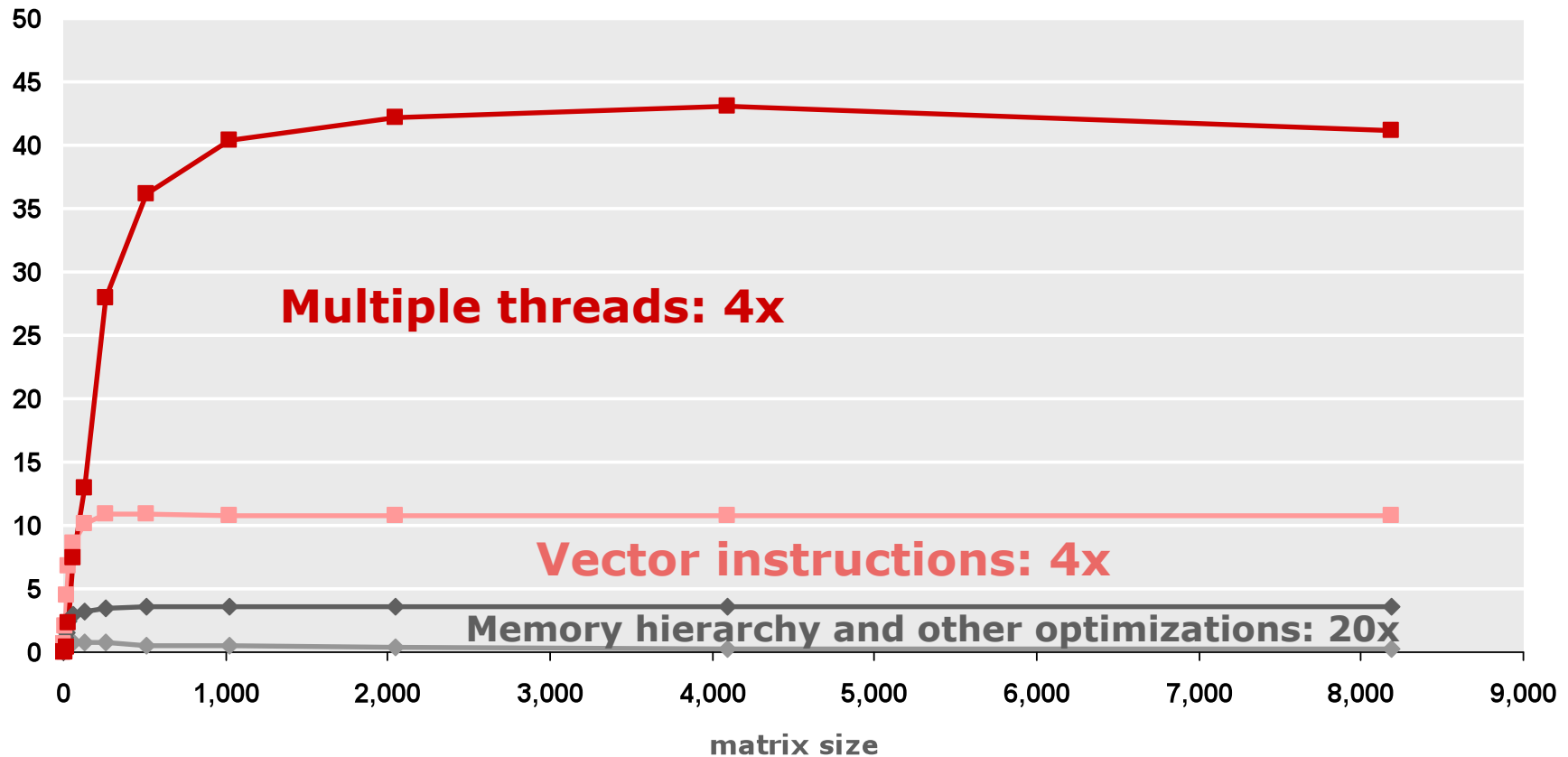
Gflop/s (giga floating point operations per second)



**Best code**

**160x**

*This code is not obviously stupid*

**Triple loop**

matrix size

- Standard desktop computer, compiler, using optimization flags
- Both implementations have **exactly** the same operations count ($2n^3$)
- *What is going on?*

# MMM Plot: Analysis

**Matrix-Matrix Multiplication (MMM) on 2 x Core 2 Duo 3 GHz**

**Gflop/s**



**Multiple threads: 4x**

**Vector instructions: 4x**

**Memory hierarchy and other optimizations: 20x**

matrix size

- **Reason for 20x: Blocking or tiling, loop unrolling, array scalarization, instruction scheduling, search to find best choice**

- *Effect: more instruction level parallelism, better register use, less L1/L2 cache misses, less TLB misses*

# Harsh Reality

- *There's more to runtime performance than asymptotic complexity*


- *One can easily loose 10x, 100x in runtime or even more*


- **What matters:**
  - Constants (100n and 5n is both O(n), but ....)
  - Coding style (unnecessary procedure calls, unrolling, reordering, …)
  - Algorithm structure (locality, instruction level parallelism, …)
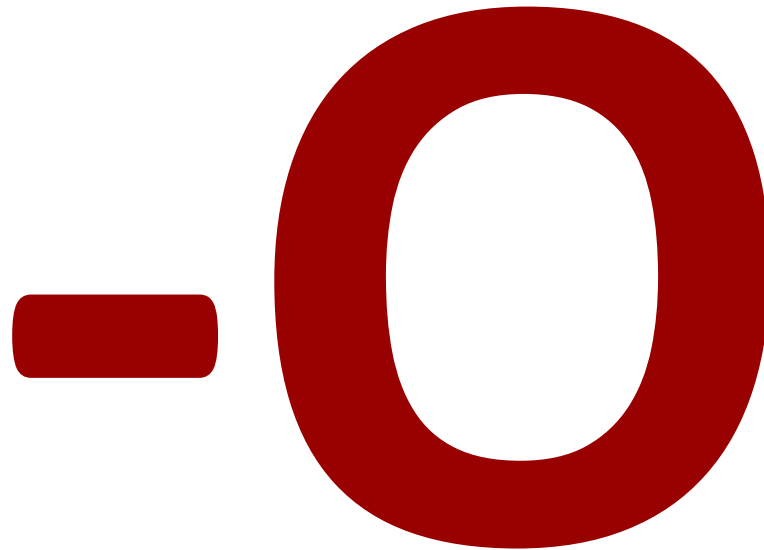  - Data representation (complicated structs or simple arrays)

# Harsh Reality

- **Must optimize at multiple levels:**
  - Algorithm
  - Data representations
  - Procedures
  - Loops

- **Must understand system to optimize performance**
  - How programs are compiled and executed
    - Execution units, memory hierarchy
  - How to measure program performance and identify bottlenecks
  - How to improve performance without destroying code modularity and generality

# Optimizing Compilers

-O

- **Use optimization flags, default is no optimization (-O0)!**
- **Good choices for gcc: -O2, -O3, -march=xxx, -m64**
- **Try different flags and maybe different compilers**

# Example

```
double a[4][4];
double b[4][4];
double c[4][4]; # set to zero

/* Multiply 4 x 4 matrices a and b  */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < 4; i++)
        for (j = 0; j < 4; j++)
            for (k = 0; k < 4; k++)
                c[i*4+j] += a[i*4 + k]*b[k*4 + j];
}
```

- **Compiled without flags:**
  **~1300 cycles**

- **Compiled with –O3 –m64 -march=… –fno-tree-vectorize**
  **~150 cycles**

- **Core 2 Duo, 2.66 GHz**

# Optimizing Compilers

- **Compilers are <span style="color:red">good</span> at: mapping program to machine**
  - register allocation
  - code selection and ordering (scheduling)
  - dead code elimination
  - eliminating minor inefficiencies

- **Compilers are <span style="color:red">not good</span> 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 <span style="color:red">not good</span> at: overcoming "optimization blockers"**
  - potential memory aliasing
  - potential procedure side-effects

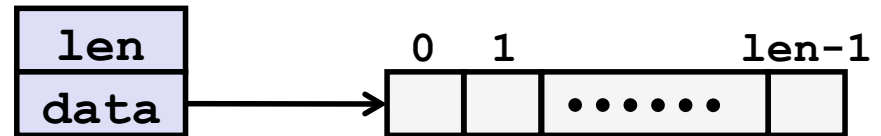# Limitations of Optimizing Compilers

- ***If in doubt, the compiler is conservative***

- **Operate under fundamental constraints**
  - 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

- **Most analysis is based only on *static* information**
  - Compiler has difficulty anticipating run-time inputs

# Today

- **Memory layout**

- **Program optimization**
  - Overview
  - **Removing unnecessary procedure calls**
  - Code motion/precomputation
  - Strength reduction
  - Sharing of common subexpressions
  - Optimization blocker: Procedure calls
  - Optimization blocker: Memory aliasing

# Example: Data Type for Vectors

```
/* data structure for vectors */
typedef struct{
        int len;
        double *data;
} vec;
```



```
/* retrieve vector element and store at val */
double get_vec_element(*vec, idx, double *val)
{
        if (idx < 0 || idx >= v->len)
                return 0;
        *val = v->data[idx];
        return 1;
}
```

# Example: Summing Vector Elements

```c
/* retrieve vector element and store at val */
double get_vec_element(*vec, idx, double *val)
{
    if (idx < 0 || idx >= v->len)
            return 0;
    *val = v->data[idx];
    return 1;
}
```

**Bound check unnecessary in sum_elements** *Why?*

```c
/* sum elements of vector */
double sum_elements(vec *v, double *res)
{
    int i;
    n = vec_length(v);
    *res = 0.0;
    double val;

    for (i = 0; i < n; i++) {
        get_vec_element(v, i, &val);
            *res += val;
    }
    return res;
}
```

**Overhead for every fp +:**
- One fct call
- One <
- One >=
- One ||
- One memory variable access

**Slowdown:**
**probably 10x or more**

# Removing Procedure Call

```
/* sum elements of vector */
double sum_elements(vec *v, double *res)
{
  int i;
  n = vec_length(v);
  *res = 0.0;
  double val;

  for (i = 0; i < n; i++) {
    get_vec_element(v, i, &val);
        *res += val;
  }
  return res;
}
```

```
/* sum elements of vector */
double sum_elements(vec *v, double *res)
{
  int i;
  n = vec_length(v);
  *res = 0.0;
  double *data = get_vec_start(v);

  for (i = 0; i < n; i++)
        *res += data[i];
  return res;
}
```

# Removing Procedure Calls

- **Procedure calls can be very expensive**

- **Bound checking can be very expensive**

- **Abstract data types can easily lead to inefficiencies**
  - Usually avoided in superfast numerical library functions

- **Watch your innermost loop!**

- **Get a feel for overhead versus actual computation being performed**

# Today

- **Memory layout**

- **Program optimization**
  - Overview
  - Removing unnecessary procedure calls
  - **Code motion/precomputation**
  - Strength reduction
  - Sharing of common subexpressions
  - Optimization blocker: Procedure calls
  - Optimization blocker: Memory aliasing

# Code Motion

- **Reduce frequency with which computation is performed**
  - If it will always produce same result
  - Especially moving code out of loop
- **Sometimes also called precomputation**

```
void set_row(double *a, double *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];
```

# Compiler-Generated Code Motion

```c
void set_row(double *a, double *b,
    long i, long n)
{
    long j;
    for (j = 0; j < n; j++)
        a[n*i+j] = b[j];
}
```

```c
long j;
long ni = n*i;
double *rowp = a+ni;
for (j = 0; j < n; j++)
    *rowp++ = b[j];
```

```
set_row:
        xorl     %r8d, %r8d           #   j = 0
        cmpq     %rcx, %r8            #   j:n
        jge      .L7                  #   if >= goto done
        movq     %rcx, %rax           #   n
        imulq    %rdx, %rax           #   n*i outside of inner loop
        leaq     (%rdi,%rax,8), %rdx  #   rowp = A + n*i*8
.L5:                                  # loop:
        movq     (%rsi,%r8,8), %rax   #   t = b[j]
        incq     %r8                  #   j++
        movq     %rax, (%rdx)         #   *rowp = t
        addq     $8, %rdx             #   rowp++
        cmpq     %rcx, %r8            #   j:n
        jl       .L5                  #   if < goto loop
.L7:                                  # done:
        rep ; ret                     #   return
```

# Today

- **Memory layout**

- **Program optimization**
  - Overview
  - Removing unnecessary procedure calls
  - Code motion/precomputation
  - **Strength reduction**
  - Sharing of common subexpressions
  - Optimization blocker: Procedure calls
  - Optimization blocker: Memory aliasing

# Strength Reduction

- **Replace costly operation with simpler one**

- **Example: Shift/add instead of multiply or divide**

  $$16*x \quad \rightarrow \quad x << 4$$

  - Utility machine dependent
  - Depends on cost of multiply or divide instruction
  - On Pentium IV, integer multiply requires 10 CPU cycles

- **Example: Recognize sequence of products**

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

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

# Today

- **Memory layout**

- **Program optimization**
  - Overview
  - Removing unnecessary procedure calls
  - Code motion/precomputation
  - Strength reduction
  - **Sharing of common subexpressions**
  - Optimization blocker: Procedure calls
  - Optimization blocker: Memory aliasing

# Share Common Subexpressions

- **Reuse portions of expressions**

- **Compilers often not very sophisticated in exploiting arithmetic properties**

*3 mults: i\*n, (i–1)\*n, (i+1)\*n*

```
/* 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;
```

*1 mult: i\*n*

```
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;
```

```
leaq   1(%rsi), %rax  # i+1
leaq   -1(%rsi), %r8  # i-1
imulq  %rcx, %rsi     # i*n
imulq  %rcx, %rax     # (i+1)*n
imulq  %rcx, %r8      # (i-1)*n
addq   %rdx, %rsi     # i*n+j
addq   %rdx, %rax     # (i+1)*n+j
addq   %rdx, %r8      # (i-1)*n+j
```

```
imulq   %rcx, %rsi  # i*n
addq    %rdx, %rsi  # i*n+j
movq    %rsi, %rax  # i*n+j
subq    %rcx, %rax  # i*n+j-n
leaq    (%rsi,%rcx), %rcx # i*n+j+n
```

# Today

- **Memory layout**

- **Program optimization**
  - Overview
  - Removing unnecessary procedure calls
  - Code motion/precomputation
  - Strength reduction
  - Sharing of common subexpressions
  - **Optimization blocker: Procedure calls**
  - Optimization blocker: Memory aliasing
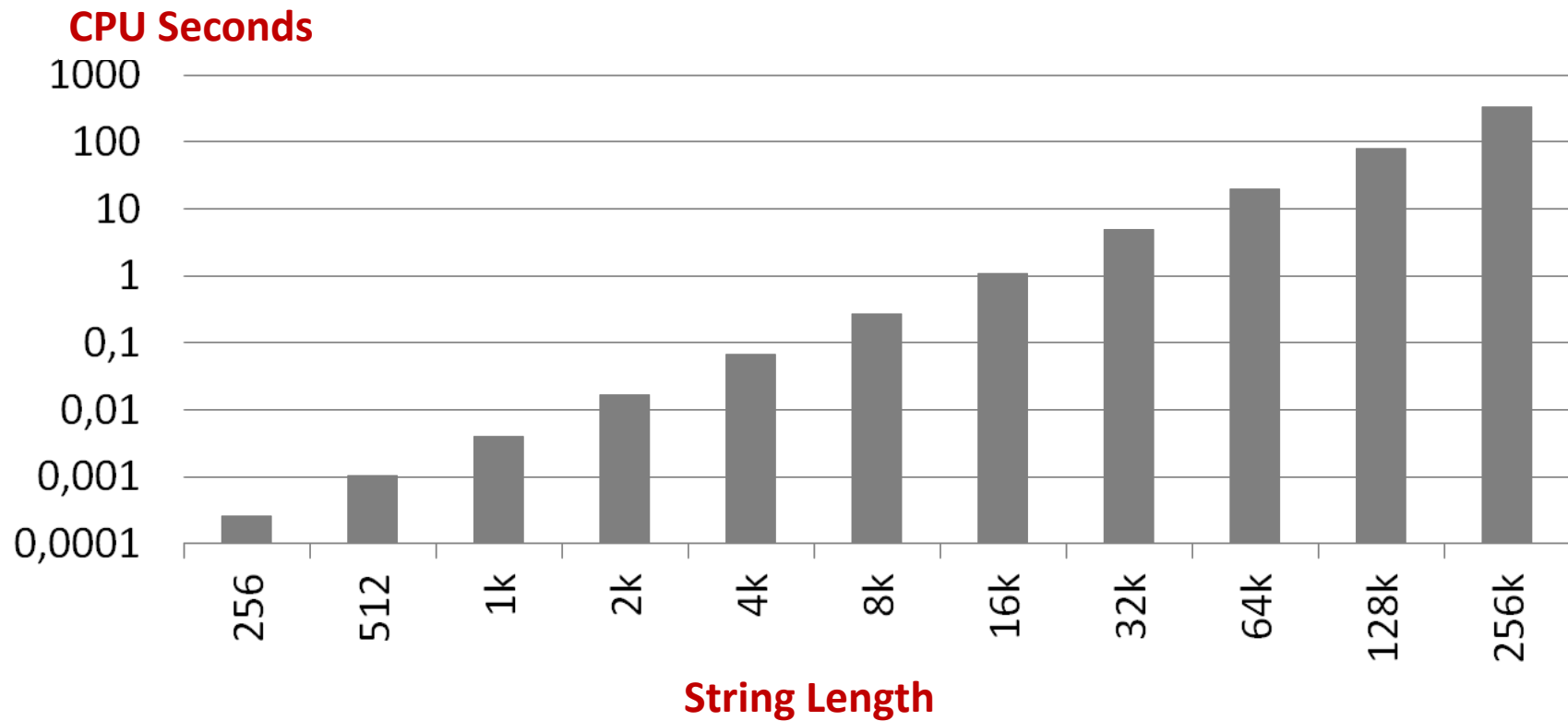
# Optimization Blocker #1: Procedure Calls

- **Procedure to convert string to lower case**

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

*Extracted from actual lab submissions*

# Performance

- **Time quadruples when double string length**
- **Quadratic performance**

# Why is That?

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

- **String length is called in every iteration!**
  - And `strlen` is O(n), so `lower` is O($n^2$)

```
/* My version of strlen */
size_t strlen(const char *s)
{
    size_t length = 0;
    while (*s != '\0') {
        s++;
        length++;
    }
    return length;
}
```
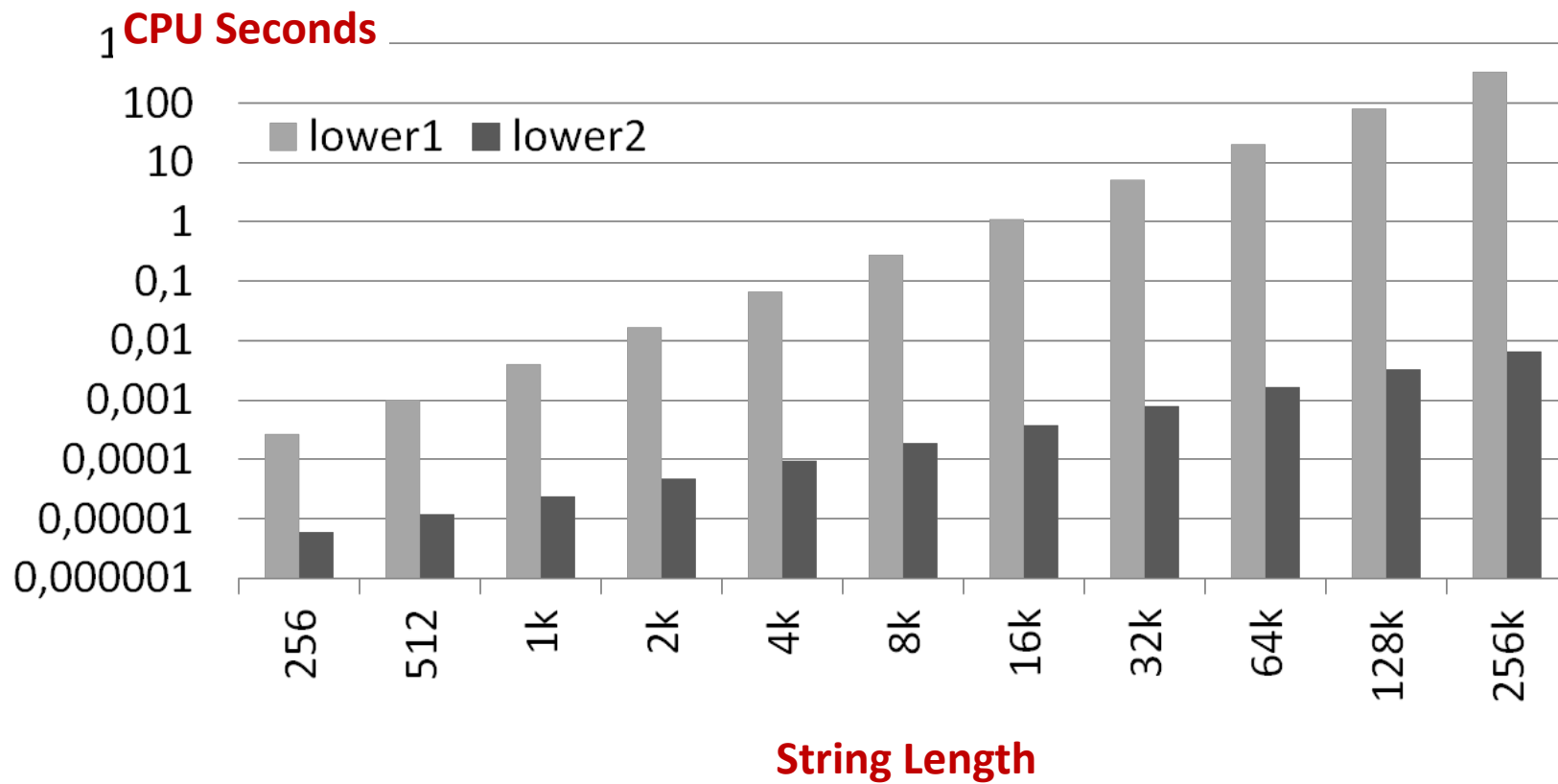
# Improving Performance

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

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

- **Move call to `strlen` outside of loop**

- **Since result does not change from one iteration to another**

- **Form of code motion/precomputation**

# Performance

- **Lower2: Time doubles when double string length**
- **Linear performance**

# Optimization Blocker: Procedure Calls

- **Why couldn't compiler move `strlen` out of inner loop?**
  - Procedure may have side effects
  - Function may not return same value for given arguments
    - Could depend on other parts of global state
    - Procedure `lower` could interact with `strlen`

- **Compiler usually treats procedure call as a black box that cannot be analyzed**
  - Consequence: conservative in optimizations

- **Remedies:**
  - Inline the function if possible
  - Do your own code motion

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