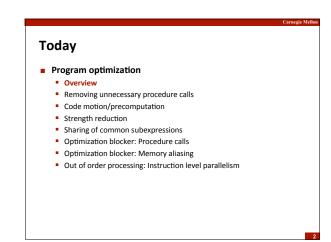
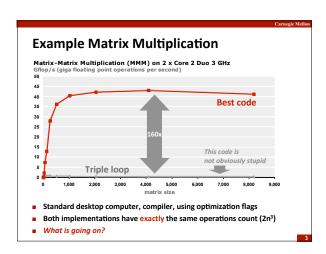
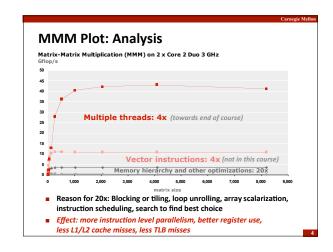
Introduction to Computer Systems
15-213/18-243, spring 2009
10<sup>th</sup> Lecture, Oct. 1<sup>st</sup>

Instructors:
Roger B. Dannenberg and Greg Ganger







Harsh Reality

In There's more to runtime performance than asymptotic complexity

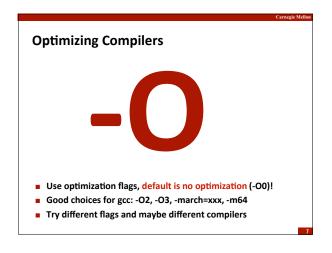
In One can easily loose 10x, 100x in runtime or even more

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

Compilers are good at: mapping program to machine
register allocation
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"

Compilers are not good 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

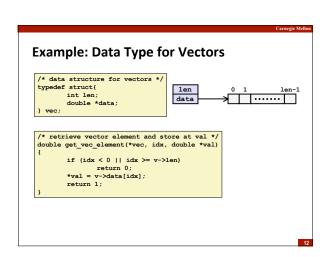
Whole-program analysis is too expensive in most cases

Most analysis is based only on static information

Compiler has difficulty anticipating run-time inputs

Today

Program optimization
Overview
Removing unnecessary procedure calls
Code motion/precomputation
Strength reduction
Strength reduction
Sharing of common subexpressions
Optimization blocker: Procedure calls
Optimization blocker: Memory aliasing
Out of order processing: Instruction level parallelism



```
Example: Summing Vector Elements
/* retrieve vector element and store at val */
double get_vec_element(*vec, idx, double *val)
                                                                   Bound check
   if (idx < 0 || idx >= v->len)
                                                                   unnecessary
   return 0;
*val = v->data[idx];
return 1;
                                                                  in sum_elements
                                                                   Why?
/* sum elements of vector */
double sum_elements(vec *v, double *res)
                                                              Overhead for every fp +:
                                                              One fct call
   int i;
n = vec_length(v);
*res = 0.0;
double val;
                                                              •One <
                                                              •One II
                                                              •One memory variable
   access
                                                              Slowdown:
                                                              probably 10x or more
```

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 for in superfast numerical library functions

Watch your innermost loop!

Get a feel for overhead versus actual computation being performed

Today

Program optimization
Overview
Removing unnecessary procedure calls
Code motion/precomputation
Strength reduction
Sharing of common subexpressions
Optimization blocker: Procedure calls
Optimization blocker: Memory aliasing
Out of order processing: Instruction level parallelism

```
Reduce frequency with which computation is performed

• If it will always produce same result

• Especially moving code out of loop

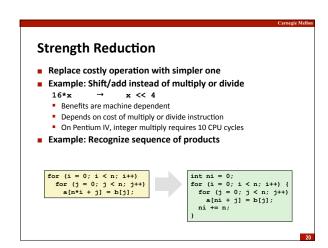
Sometimes also called precomputation

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];</pre>
```

Today

Program optimization
Overview
Removing unnecessary procedure calls
Code motion/precomputation
Strength reduction
Sharing of common subexpressions
Optimization blocker: Procedure calls
Optimization blocker: Memory aliasing
Out of order processing: Instruction level parallelism



Today

Program optimization

Overview

Removing unnecessary procedure calls

Code motion/precomputation

Strength reduction

Sharing of common subexpressions

Optimization blocker: Procedure calls

Optimization blocker: Memory aliasing

Out of order processing: Instruction level parallelism

Today

Program optimization

Overview

Removing unnecessary procedure calls

Code motion/precomputation

Strength reduction

Sharing of common subexpressions

Optimization blocker: Procedure calls

Optimization blocker: Memory aliasing

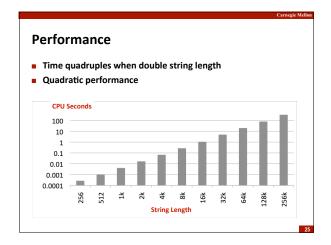
Out of order processing: Instruction level parallelism

Optimization Blocker #1: Procedure Calls

Procedure to convert string to lower case

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 213 lab submissions, Fall 1998



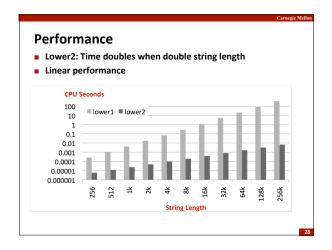
```
Why is That?

void lower(char *s) {
   int i;
   for (i = 0; i < strlen(s); i++)
        if (s[i] >= 'A' ' & s[i] <= '2')
        s[i] -= ('A' - 'a');
}

■ String length is called in every iteration!

■ And strlen is O(n), so lower is O(n²)

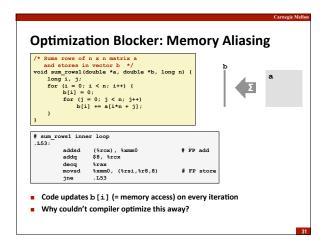
/* My version of strlen */
   size_t strlen(const char *s)
   {
        size_t length = 0;
        while ('s != '\0') {
            s++;
            length++;
        }
        return length;
}</pre>
```

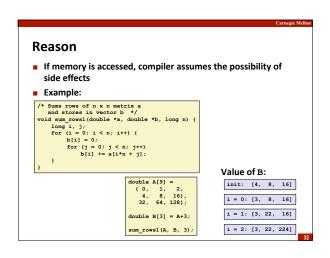


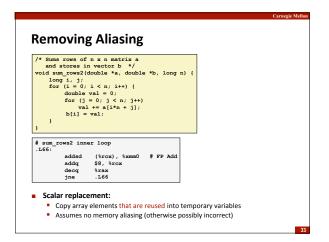
**Optimization Blocker: Procedure Calls** Why couldn't compiler move strlen out of inner loop? Procedure may have side effects Function may not return same value for given arguments Could depend on other parts of global state Procedure lower could interact with strlen ■ Compiler usually treats procedure call as a black box that cannot be analyzed Consequence: conservative in optimizations int lencnt = 0; Remedies: size\_t strlen(const char \*s) Inline the function if possible size\_t length = 0; while (\*s != '\0') { Do your own code motion s++; length++; return length;

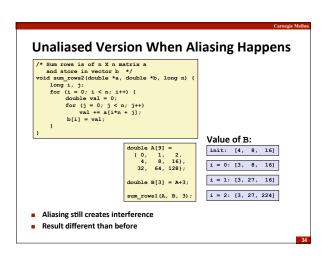
Today

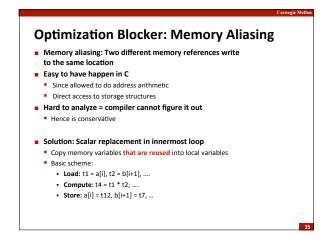
Program optimization
Overview
Removing unnecessary procedure calls
Code motion/precomputation
Strength reduction
Sharing of common subexpressions
Optimization blocker: Procedure calls
Optimization blocker: Memory aliasing
Out of order processing: Instruction level parallelism

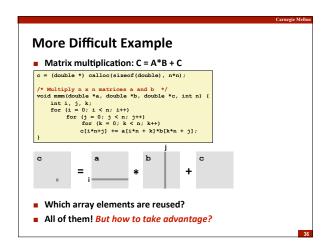


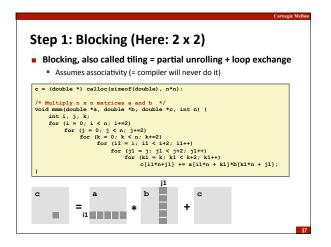


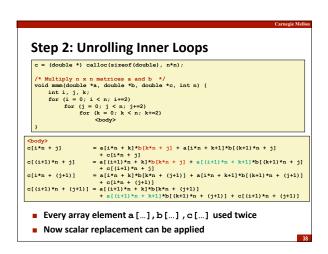


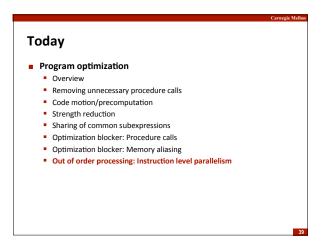


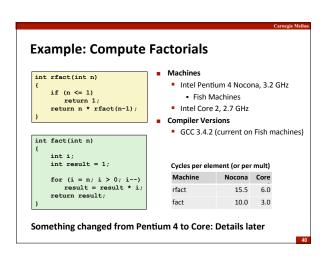


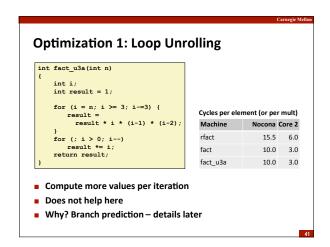


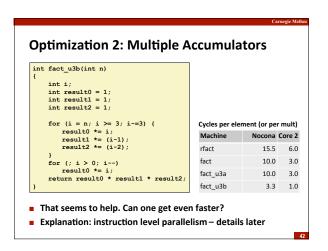


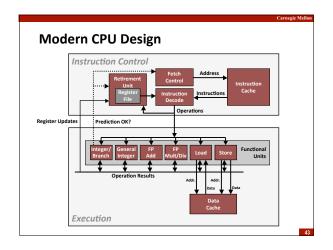






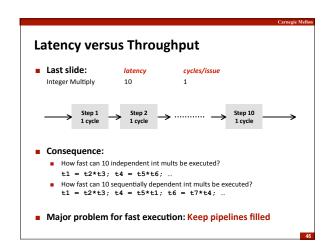






## Superscalar Processor Definition: A superscalar processor can issue and execute multiple instructions in one cycle. The instructions are retrieved from a sequential instruction stream and are usually scheduled dynamically. Benefit: without programming effort, superscalar processor can take advantage of the instruction level parallelism that most programs have Most CPUs since about 1998 are superscalar. Intel: since Pentium Pro

**Pentium 4 Nocona CPU** ■ Multiple instructions can execute in parallel 1 load, with address computation 1 store, with address computation 2 simple integer (one may be branch) 1 complex integer (multiply/divide) 1 FP/SSF3 unit 1 FP move (does all conversions) Some instructions take > 1 cycle, but can be pipelined Latency Cycles/Issue Load / Store Integer Multiply 10 Integer/Long Divide 36/106 36/106 Single/Double FP Multiply Single/Double FP Add Single/Double FP Divide 32/46



**Hard Bounds**  Latency and throughput of instructions Cycles/Issue Load / Store 1 Integer Multiply 10 Integer/Long Divide 36/106 36/106 Single/Double FP Multiply Single/Double FP Add Single/Double FP Divide 32/46 32/46 How many cycles at least if Function requires n int mults? Function requires n float adds? • Function requires n float ops (adds and mults)?

Performance in Numerical Computing

Numerical computing = computing dominated by floating point operations

Example: Matrix multiplication

Performance measure: Floating point operations per second (flop/s)

Counting point operations per second (flop/s)

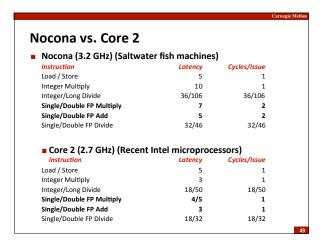
Counting only floating point adds and mults

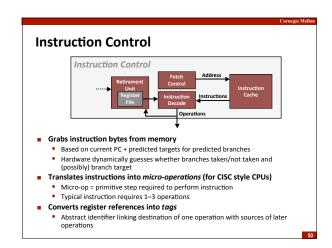
Higher is better

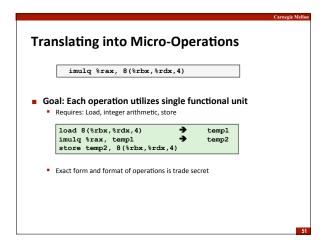
Like inverse runtime

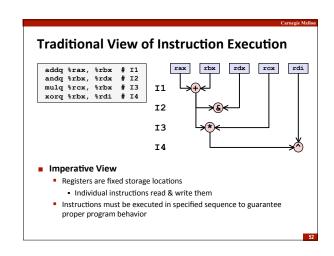
Theoretical scalar (no vector SSE) peak performance on fish machines?

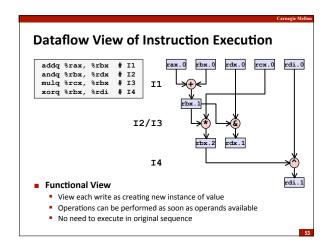
3.2 Gflop/s = 3200 Mflop/s. Why?

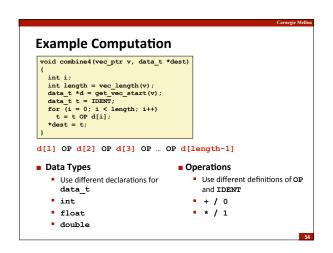


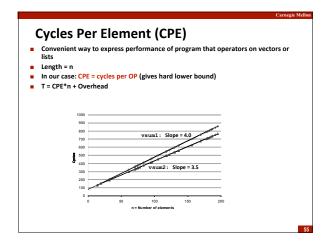


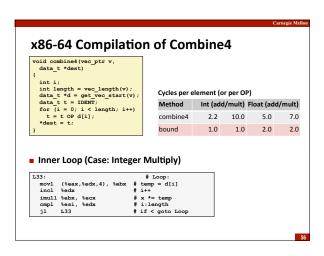


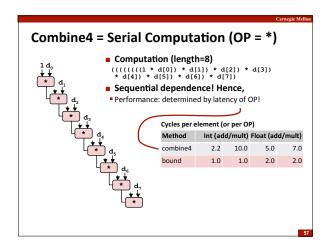


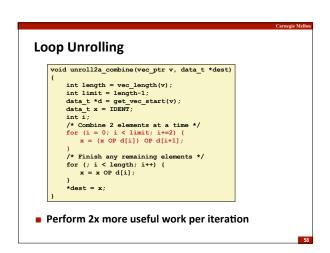










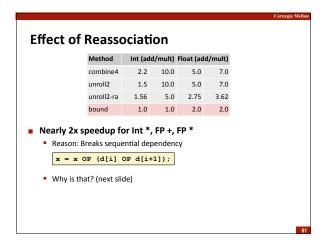


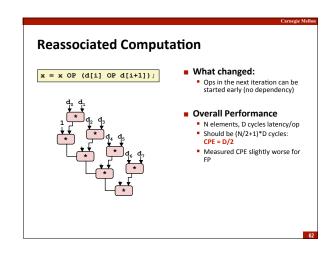
```
Loop Unrolling with Reassociation

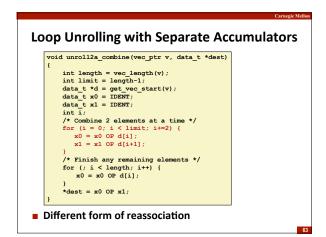
void unrol12aa_combine(vec_ptr v, data_t *dest)
{
    int length = vec_length(v);
    int limit = length-1;
    data_t *d = get_vec_start(v);
    data_t *d = get_vec_start(v);
    int i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x = x OP (d[i] OP d[i+1]);
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x = x OP d[i];
    }
    *dest = x;
}

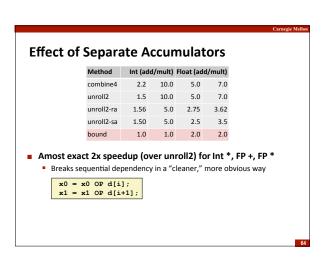
■ Can this change the result of the computation?

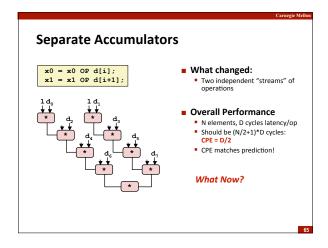
■ Yes, for FP. Why?</pre>
```

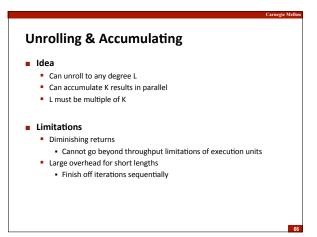


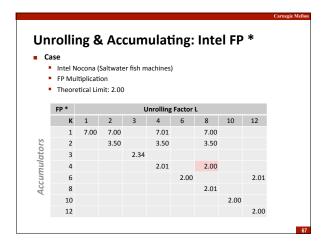


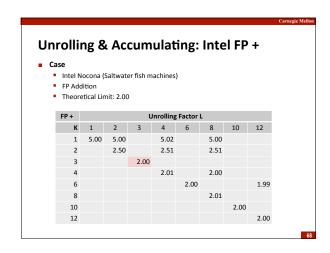


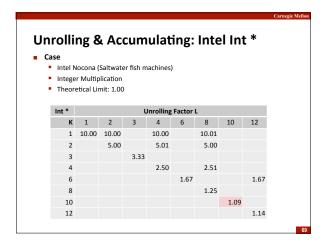


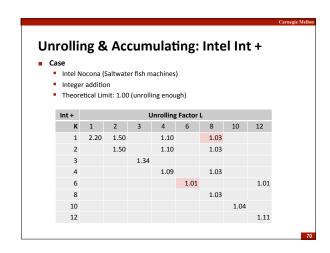


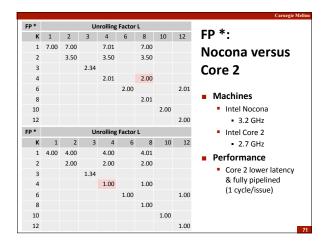










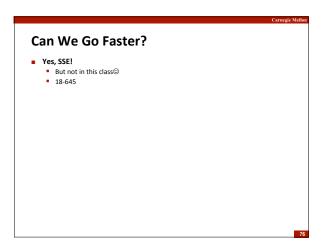


Int *			U	nrolling					
K	1	2	3	4	6	8	10	12	Nocona vs.
1	10.00	10.00		10.00		10.01			Core 2 Int *
2		5.00		5.01		5.00			Core 2 int
3			3.33						
4				2.50		2.51			■ Performance
6					1.67			1.67	Newer version of
8						1.25			GCC does
10							1.09		reassociation
12								1.14	
Int *	Unrolling Factor L								Why for int's and not for float's?
K	1	2	3	4	6	8	10	12	ioi iloat s:
1	3.00	1.50		1.00		1.00			
2		1.50		1.00		1.00			
3			1.00						
4				1.00		1.00			
6					1.00			1.00	
8						1.00			
10							1.00		
12								1.33	l

FP*			U	nrolling	Factor	Intel vs. AMD FP			
K	1	2	3	4	6	8	10	12	*
1	7.00	7.00		7.01		7.00			•
2		3.50		3.50		3.50			
3			2.34						Machines
4				2.01		2.00			<ul> <li>Intel Nocona</li> </ul>
6					2.00			2.01	■ 3.2 GHz
8						2.01			AMD Opteron
10							2.00		■ 2.0 GHz
12								2.00	
FP*			U	nrolling	Factor	L			<ul><li>Performance</li></ul>
K	1	2	3	4	6	8	10	12	<ul> <li>AMD lower latency</li> </ul>
1	4.00	4.00		4.00		4.01			& better pipelining
2		2.00		2.00		2.00			<ul> <li>But slower clock rate</li> </ul>
3			1.34						
4				1.00		1.00			
6					1.00			1.00	
8						1.00			
10							1.00		
12								1.00	

Int *			U	nrolling	Intel vs. AND				
К	1	2	3	4	6	8	10	12	
1	10.00	10.00		10.00		10.01			Int *
2		5.00		5.01		5.00			
3			3.33						■ Performance
4				2.50		2.51			<ul> <li>AMD multiplier</li> </ul>
6					1.67			1.67	much lower latency
В						1.25			<ul> <li>Can get high</li> </ul>
10							1.09		performance with less work
L2								1.14	
Int *			U	nrolling	Factor	L			<ul> <li>Doesn't achieve as good an optimum</li> </ul>
K	1	2	3	4	6	8	10	12	good an optimum
1	3.00	3.00		3.00		3.00			
2		2.33		2.0		1.35			
3			2.00						
4				1.75		1.38			
6					1.50			1.50	
8						1.75			
10							1.30		
12								1.33	1

Int+			Uı	nrolling	Factor	r L			Intel vs. AMD in
K	1	2	3	4	6	8	10	12	
1	2.20	1.50		1.10		1.03			+
2		1.50		1.10		1.03			
3			1.34						<ul><li>Performance</li></ul>
4				1.09		1.03			<ul> <li>AMD gets below 1.0</li> </ul>
6					1.01			1.01	<ul><li>Even just with</li></ul>
8						1.03			unrolling
10							1.04		<ul><li>Explanation</li></ul>
12								1.11	<ul><li>Both Intel &amp; AMD</li></ul>
Int+			Uı	nrolling	Factor	r L			can "double pump"
K	1	2	3	4	6	8	10	12	integer units
1	2.32	1.50		0.75		0.63			<ul> <li>Only AMD can load</li> </ul>
2		1.50		0.83		0.63			two elements / cycle
3			1.00						
4				1.00		0.63			
6					0.83			0.67	
8						0.63			
10							0.60		
								0.85	



## **Summary**

- Optimization comes from many directions:
  - Algorithm design: huge potential
  - Optimizing compilers: effective but conservative
  - Manual tuning: many techniques
  - Parallel computation: we'll talk about this later
- Understanding processors, memory, and compilers