

# CS 61C: Great Ideas in Computer Architecture *SIMD*

Instructor:

Randy H. Katz

<http://inst.eecs.Berkeley.edu/~cs61c/fa13>

10/10/13

Fall 2013 -- Lecture #14

1

## Review

- Three C's of cache misses
  - Compulsory
  - Capacity
  - Conflict
- Amdahl's Law:  $\text{Speedup} = 1 / ((1-F) + F/S)$
- Flynn's Taxonomy: SISD/SIMD/MISD/MIMD

10/10/13

Fall 2013 -- Lecture #14

2

## SIMD Architectures

- *Data parallelism*: executing one operation on multiple data streams
- Example to provide context:
  - Multiplying a coefficient vector by a data vector (e.g., in filtering)
$$y[i] := c[i] \times x[i], 0 \leq i < n$$
- Sources of performance improvement:
  - One instruction is fetched & decoded for entire operation
  - Multiplications are known to be independent
  - Pipelining/concurrency in memory access as well

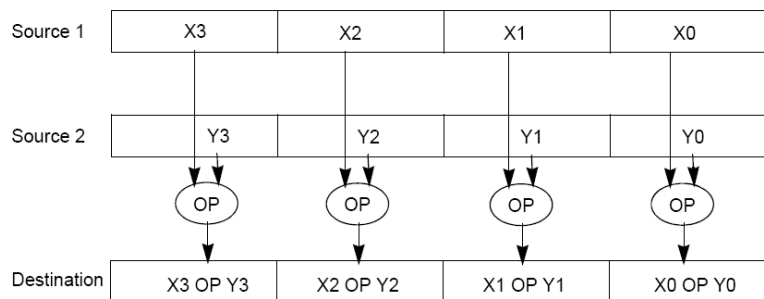
10/10/13

Fall 2013 -- Lecture #14

Slide 3

## “Advanced Digital Media Boost”

- To improve performance, Intel’s SIMD instructions
  - Fetch one instruction, do the work of multiple instructions
  - MMX (MultiMedia eXtension, Pentium II processor family)
  - SSE (*Streaming SIMD Extension, Pentium III and beyond*)



10/10/13

Fall 2013 -- Lecture #14

4

## Example: SIMD Array Processing

```
for each f in array
  f = sqrt(f)
```

```
for each f in array
```

```
{
  load f to the floating-point register
  calculate the square root
  write the result from the register to memory
}
```

```
{ for each 4 members in array
```

```
{
  load 4 members to the SSE register
  calculate 4 square roots in one operation
  store the 4 results from the register to memory
}
```

SIMD style

10/10/13

Fall 2013 -- Lecture #14

5

## Data-Level Parallelism and SIMD

- SIMD wants adjacent values in memory that can be operated in parallel
- Usually specified in programs as loops
 

```
for(i=1000; i>0; i=i-1)
  x[i] = x[i] + s;
```
- How can reveal more data-level parallelism than available in a single iteration of a loop?
- *Unroll loop* and adjust iteration rate

10/10/13

Spring 2013 -- Lecture #14

6

## Looping in MIPS

Assumptions:

- \$t1 is initially the address of the element in the array with the highest address
- \$f0 contains the scalar value s
- 8(\$t2) is the address of the last element to operate on

CODE:

```

Loop: 1. l.d      $f2,0($t1) ; $f2=array element
      2. add.d   $f10,$f2,$f0 ; add s to $f2
      3. s.d     $f10,0($t1) ; store result
      4. addui   $t1,$t1,#-8 ; decrement pointer 8 byte
      5. bne    $t1,$t2,Loop ; repeat loop if $t1 != $t2
  
```

10/10/13

Spring 2013 -- Lecture #14

7

## Loop Unrolled

```

Loop: l.d      $f2,0($t1)
      add.d   $f10,$f2,$f0
      s.d     $f10,0($t1)
      l.d     $f4,-8($t1)
      add.d   $f12,$f4,$f0
      s.d     $f12,-8($t1)
      l.d     $f6,-16($t1)
      add.d   $f14,$f6,$f0
      s.d     $f14,-16($t1)
      l.d     $f8,-24($t1)
      add.d   $f16,$f8,$f0
      s.d     $f16,-24($t1)
      addui   $t1,$t1,#-32
      bne    $t1,$t2,Loop
  
```

NOTE:

1. Only 1 Loop Overhead every 4 iterations
2. This unrolling works if  
loop\_limit(mod 4) = 0
3. (Different Registers eliminate stalls in pipeline; we'll see later in course)

10/10/13

Spring 2013 -- Lecture #14

8

## Loop Unrolled Scheduled

```

Loop:l.d  $f2,0($t1)
      l.d  $f4,-8($t1)
      l.d  $f6,-16($t1)
      l.d  $f8,-24($t1)
      add.d $f10,$f2,$f0
      add.d $f12,$f4,$f0
      add.d $f14,$f6,$f0
      add.d $f16,$f8,$f0
      s.d  $f10,0($t1)
      s.d  $f12,-8($t1)
      s.d  $f14,-16($t1)
      s.d  $f16,-24($t1)
      addui $t1,$t1,#-32
      bne  $t1,$t2,Loop
  
```

4 Loads side-by-side: Could replace with 4-wide SIMD Load

4 Adds side-by-side: Could replace with 4-wide SIMD Add

4 Stores side-by-side: Could replace with 4-wide SIMD Store

10/10/13

Spring 2013 -- Lecture #14

9

## Loop Unrolling in C

- Instead of compiler doing loop unrolling, could do it yourself in C

```
for(i=1000; i>0; i=i-1)
```

```
    x[i] = x[i] + s;
```

- Could be rewritten

What is downside of doing it in C?

```
for(i=1000; i>0; i=i-4) {
```

```
    x[i]  = x[i]  + s;
```

```
    x[i-1] = x[i-1] + s;
```

```
    x[i-2] = x[i-2] + s;
```

```
    x[i-3] = x[i-3] + s;
```

```
}
```

10/10/13

Spring 2013 -- Lecture #14

10

## Generalizing Loop Unrolling

- A loop of **n iterations**
- **k copies** of the body of the loop
- **Assuming  $(n \bmod k) \neq 0$**

Then we will run the loop with 1 copy of the body  **$(n \bmod k)$**  times and

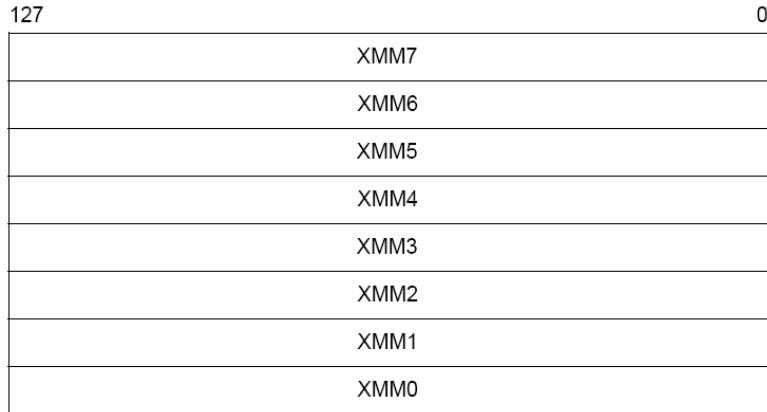
with k copies of the body  **$\text{floor}(n/k)$**  times

- (Will revisit loop unrolling again when get to pipelining later in semester)

## Intel SIMD Extensions

- MMX 64-bit registers, reusing floating-point registers [1992]
- SSE2/3/4, new 128-bit registers [1999]
- AVX, new 256-bit registers [2011]
  - Space for expansion to 1024-bit registers

## XMM Registers



- Architecture extended with eight 128-bit data registers: XMM registers
  - x86 64-bit address architecture adds 8 additional registers (XMM8 – XMM15)

10/10/13

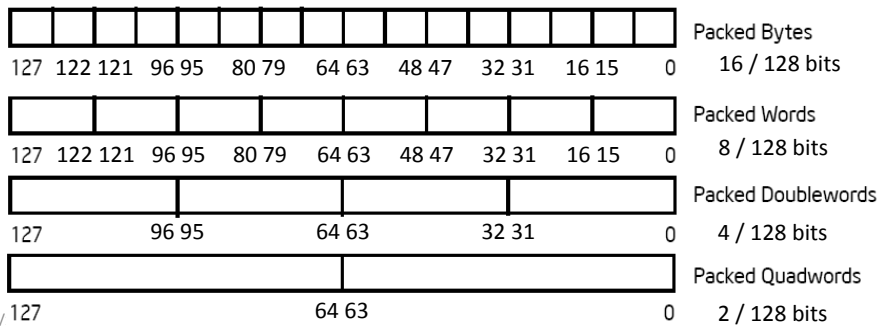
Fall 2013 -- Lecture #14

13

## Intel Architecture SSE2+ 128-Bit SIMD Data Types

- Note: in Intel Architecture (unlike MIPS) a word is 16 bits
  - Single-precision FP: Double word (32 bits)
  - Double-precision FP: Quad word (64 bits)

Fundamental 128-Bit Packed SIMD Data Types



10/127

## First SIMD Extensions: MIT Lincoln Labs TX-2, 1957

**ONE 36 BIT AE (36)**

|   |   |   |   |   |
|---|---|---|---|---|
|   | 4 | 3 | 2 | 1 |
| D |   |   |   |   |
| C |   |   |   |   |
| A |   |   |   |   |
| B |   |   |   |   |

**OPERAND WORD STRUCTURE**

|   |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|
| S |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|

**TWO 18 BIT AE'S (18,18)**

|   |   |   |   |   |  |   |   |   |   |
|---|---|---|---|---|--|---|---|---|---|
|   | 4 | 3 | 2 | 1 |  | 4 | 3 | 2 | 1 |
| D |   |   |   |   |  |   |   |   |   |
| C |   |   |   |   |  |   |   |   |   |
| A |   |   |   |   |  |   |   |   |   |
| B |   |   |   |   |  |   |   |   |   |

**OPERAND WORD STRUCTURE**

|   |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |   |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|
| S |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  | S |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|

**ONE 27 BIT & ONE 9 BIT AE (27,9)**

|   |   |   |   |   |  |   |   |   |   |
|---|---|---|---|---|--|---|---|---|---|
|   | 4 | 3 | 2 | 1 |  | 4 | 3 | 2 | 1 |
| D |   |   |   |   |  |   |   |   |   |
| C |   |   |   |   |  |   |   |   |   |
| A |   |   |   |   |  |   |   |   |   |
| B |   |   |   |   |  |   |   |   |   |

**OPERAND WORD STRUCTURE**

|   |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |   |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|
| S |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  | S |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|

**FOUR 9 BIT AE'S (9,9,9,9)**

|   |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |
|---|---|---|---|---|--|---|---|---|---|--|---|---|---|---|--|---|---|---|---|--|
|   | 4 | 3 | 2 | 1 |  | 4 | 3 | 2 | 1 |  | 4 | 3 | 2 | 1 |  | 4 | 3 | 2 | 1 |  |
| D |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |
| C |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |
| A |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |
| B |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |   |   |   |   |  |

**OPERAND WORD STRUCTURE**

|   |  |  |  |  |  |  |  |  |   |  |  |  |  |  |  |  |  |   |  |  |  |  |  |  |  |  |   |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|
| S |  |  |  |  |  |  |  |  | S |  |  |  |  |  |  |  |  | S |  |  |  |  |  |  |  |  | S |  |  |  |  |  |  |  |  |
|---|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|--|

ecture #14 15

## SSE/SSE2 Floating Point Instructions

|  | Data transfer                      | Arithmetic                    | Compare          |
|--|------------------------------------|-------------------------------|------------------|
| Move<br>does<br>both<br>load<br>and<br>store | MOV{A/U}{SS/PS/SD/PD} xmm, mem/xmm | ADD{SS/PS/SD/PD} xmm, mem/xmm | CMP{SS/PS/SD/PD} |
|  |                                    | SUB{SS/PS/SD/PD} xmm, mem/xmm |                  |
|  | MOV {H/L} {PS/PD} xmm, mem/xmm     | MUL{SS/PS/SD/PD} xmm, mem/xmm |                  |
|  |                                    | DIV{SS/PS/SD/PD} xmm, mem/xmm |                  |
|  |                                    | SQRT{SS/PS/SD/PD} mem/xmm     |                  |
|  |                                    | MAX {SS/PS/SD/PD} mem/xmm     |                  |
|  |                                    | MIN{SS/PS/SD/PD} mem/xmm      |                  |

xmm: one operand is a 128-bit SSE2 register  
 mem/xmm: other operand is in memory or an SSE2 register  
 {SS} Scalar Single precision FP: one 32-bit operand in a 128-bit register  
 {PS} Packed Single precision FP: four 32-bit operands in a 128-bit register  
 {SD} Scalar Double precision FP: one 64-bit operand in a 128-bit register  
 {PD} Packed Double precision FP, or two 64-bit operands in a 128-bit register  
 {A} 128-bit operand is aligned in memory  
 {U} means the 128-bit operand is unaligned in memory  
 {H} means move the high half of the 128-bit operand  
 {L} means move the low half of the 128-bit operand



## Example: Add Two Single-Precision Floating-Point Vectors

Computation to be performed:

```
vec_res.x = v1.x + v2.x;
vec_res.y = v1.y + v2.y;
vec_res.z = v1.z + v2.z;
vec_res.w = v1.w + v2.w;
```

mov a ps : **move** from mem to XMM register, memory aligned, **packed** single precision

add ps : **add** from mem to XMM register, **packed** single precision

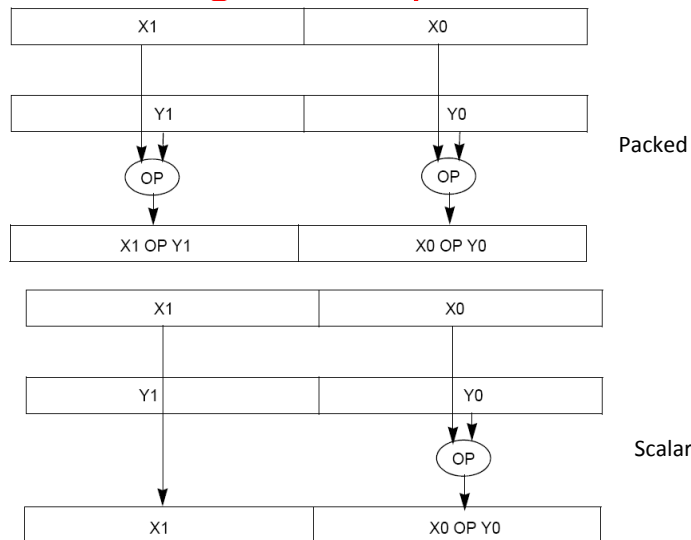
mov a ps : **move** from XMM register to mem, memory aligned, **packed** single precision

SSE Instruction Sequence:

(Note: Destination on the right in x86 assembly)

```
movaps address-of-v1, %xmm0
    // v1.w | v1.z | v1.y | v1.x -> xmm0
addps address-of-v2, %xmm0
    // v1.w+v2.w | v1.z+v2.z | v1.y+v2.y | v1.x+v2.x -> xmm0
movaps %xmm0, address-of-vec_res
```

## Packed and Scalar Double-Precision Floating-Point Operations



## Intel SSE Intrinsics

- Intrinsics are C functions and procedures for inserting assembly language into C code, including SSE instructions
  - With intrinsics, can program using these instructions indirectly
  - One-to-one correspondence between SSE instructions and intrinsics

10/10/13

Fall 2013 -- Lecture #14

19

## Example SSE Intrinsics

| Intrinsics:                                | Corresponding SSE instructions: |
|--|---------------------------------|
| • Vector data type:<br><code>_m128d</code> |                                 |
| • Load and store operations:               |                                 |
| <code>_mm_load_pd</code>                   | MOVAPD/aligned, packed double   |
| <code>_mm_store_pd</code>                  | MOVAPD/aligned, packed double   |
| <code>_mm_loadu_pd</code>                  | MOVUPD/unaligned, packed double |
| <code>_mm_storeu_pd</code>                 | MOVUPD/unaligned, packed double |
| • Load and broadcast across vector         |                                 |
| <code>_mm_load1_pd</code>                  | MOVSD + shuffling/duplicating   |
| • Arithmetic:                              |                                 |
| <code>_mm_add_pd</code>                    | ADDPD/add, packed double        |
| <code>_mm_mul_pd</code>                    | MULPD/multiple, packed double   |

10/09/2010

CS267 Lecture 7

20 20

## Example: 2 x 2 Matrix Multiply

Definition of Matrix Multiply:

$$C_{i,j} = (A \times B)_{i,j} = \sum_{k=1}^2 A_{i,k} \times B_{k,j}$$

$$\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix} \times \begin{bmatrix} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{bmatrix} = \begin{bmatrix} C_{1,1} = A_{1,1}B_{1,1} + A_{1,2}B_{2,1} & C_{1,2} = A_{1,1}B_{1,2} + A_{1,2}B_{2,2} \\ C_{2,1} = A_{2,1}B_{1,1} + A_{2,2}B_{2,1} & C_{2,2} = A_{2,1}B_{1,2} + A_{2,2}B_{2,2} \end{bmatrix}$$

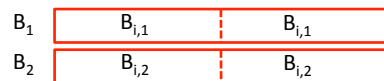
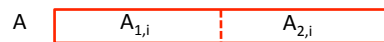
10/10/13

Fall 2013 -- Lecture #14

21

## Example: 2 x 2 Matrix Multiply

- Using the XMM registers
  - Two 64-bit doubles per XMM reg



10/10/13

Fall 2013 -- Lecture #14

22

## Example: 2 x 2 Matrix Multiply

- Initialization

|       |   |   |
|-------|---|---|
| $C_1$ | 0 | 0 |
| $C_2$ | 0 | 0 |

10/10/13

Fall 2013 -- Lecture #14

23

## Example: 2 x 2 Matrix Multiply

- Initialization

|       |   |   |
|-------|---|---|
| $C_1$ | 0 | 0 |
| $C_2$ | 0 | 0 |

- $i = 1$

|   |           |           |
|---|-----------|-----------|
| A | $A_{1,1}$ | $A_{2,1}$ |
|---|-----------|-----------|

`_mm_load_pd`: Load 2 doubles into XMM reg, Stored in memory in Column-major order

|       |           |           |
|-------|-----------|-----------|
| $B_1$ | $B_{1,1}$ | $B_{1,1}$ |
| $B_2$ | $B_{1,2}$ | $B_{1,2}$ |

`_mm_load1_pd`: SSE instruction that loads a double word and stores it in the high and low double words of the XMM register (duplicates value in both halves of XMM)

10/10/13

Fall 2013 -- Lecture #14

24

## Example: 2 x 2 Matrix Multiply

- First iteration intermediate result

|       |                    |                    |  |
|-------|--------------------|--------------------|--|
| $C_1$ | $0+A_{1,1}B_{1,1}$ | $0+A_{2,1}B_{1,1}$ | $c1 = \_mm\_add\_pd(c1, \_mm\_mul\_pd(a, b1));$<br>$c2 = \_mm\_add\_pd(c2, \_mm\_mul\_pd(a, b2));$<br>SSE instructions first do parallel multiplies<br>and then parallel adds in XMM registers |
| $C_2$ | $0+A_{1,1}B_{1,2}$ | $0+A_{2,1}B_{1,2}$ |  |

- $i = 1$

|     |           |           |  |
|-----|-----------|-----------|--|
| $A$ | $A_{1,1}$ | $A_{2,1}$ | $\_mm\_load\_pd$ : Stored in memory in<br>Column order |
|-----|-----------|-----------|--|

|       |           |           |   |
|-------|-----------|-----------|---|
| $B_1$ | $B_{1,1}$ | $B_{1,1}$ | $\_mm\_load1\_pd$ : SSE instruction that loads<br>a double word and stores it in the high and<br>low double words of the XMM register<br>(duplicates value in both halves of XMM) |
| $B_2$ | $B_{1,2}$ | $B_{1,2}$ |   |

10/10/13

Fall 2013 -- Lecture #14

25

## Example: 2 x 2 Matrix Multiply

- First iteration intermediate result

|       |                    |                    |  |
|-------|--------------------|--------------------|--|
| $C_1$ | $0+A_{1,1}B_{1,1}$ | $0+A_{2,1}B_{1,1}$ | $c1 = \_mm\_add\_pd(c1, \_mm\_mul\_pd(a, b1));$<br>$c2 = \_mm\_add\_pd(c2, \_mm\_mul\_pd(a, b2));$<br>SSE instructions first do parallel multiplies<br>and then parallel adds in XMM registers |
| $C_2$ | $0+A_{1,1}B_{1,2}$ | $0+A_{2,1}B_{1,2}$ |  |

- $i = 2$

|     |           |           |  |
|-----|-----------|-----------|--|
| $A$ | $A_{1,2}$ | $A_{2,2}$ | $\_mm\_load\_pd$ : Stored in memory in<br>Column order |
|-----|-----------|-----------|--|

|       |           |           |   |
|-------|-----------|-----------|---|
| $B_1$ | $B_{2,1}$ | $B_{2,1}$ | $\_mm\_load1\_pd$ : SSE instruction that loads<br>a double word and stores it in the high and<br>low double words of the XMM register<br>(duplicates value in both halves of XMM) |
| $B_2$ | $B_{2,2}$ | $B_{2,2}$ |   |

10/10/13

Fall 2013 -- Lecture #14

26

## Example: 2 x 2 Matrix Multiply

- Second iteration intermediate result

|       |                                   |                                   |   |
|-------|-----------------------------------|-----------------------------------|---|
| $C_1$ | $A_{1,1}B_{1,1} + A_{1,2}B_{2,1}$ | $A_{2,1}B_{1,1} + A_{2,2}B_{2,1}$ | $c1 = \_mm\_add\_pd(c1, \_mm\_mul\_pd(a, b1));$ |
| $C_2$ | $A_{1,1}B_{1,2} + A_{1,2}B_{2,2}$ | $A_{2,1}B_{1,2} + A_{2,2}B_{2,2}$ | $c2 = \_mm\_add\_pd(c2, \_mm\_mul\_pd(a, b2));$ |

SSE instructions first do parallel multiplies and then parallel adds in XMM registers

- $i = 2$

|     |           |           |   |
|-----|-----------|-----------|---|
| $A$ | $A_{1,2}$ | $A_{2,2}$ | $\_mm\_load\_pd$ : Stored in memory in Column order |
|-----|-----------|-----------|---|

|       |           |           |  |
|-------|-----------|-----------|--|
| $B_1$ | $B_{2,1}$ | $B_{2,1}$ | $\_mm\_load1\_pd$ : SSE instruction that loads a double word and stores it in the high and low double words of the XMM register (duplicates value in both halves of XMM) |
| $B_2$ | $B_{2,2}$ | $B_{2,2}$ |  |

## Live Example: 2 x 2 Matrix Multiply

Definition of Matrix Multiply:

$$C_{i,j} = (A \times B)_{i,j} = \sum_{k=1}^2 A_{i,k} \times B_{k,j}$$

|  |          |  |     |  |
|--|----------|--|-----|--|
| $\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$ | $\times$ | $\begin{bmatrix} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{bmatrix}$ | $=$ | $\begin{bmatrix} C_{1,1} = A_{1,1}B_{1,1} + A_{1,2}B_{2,1} & C_{1,2} = A_{1,1}B_{1,2} + A_{1,2}B_{2,2} \\ C_{2,1} = A_{2,1}B_{1,1} + A_{2,2}B_{2,1} & C_{2,2} = A_{2,1}B_{1,2} + A_{2,2}B_{2,2} \end{bmatrix}$ |
| $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$                         | $\times$ | $\begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}$                         | $=$ | $\begin{bmatrix} C_{1,1} = 1*1 + 0*2 = 1 & C_{1,2} = 1*3 + 0*4 = 3 \\ C_{2,1} = 0*1 + 1*2 = 2 & C_{2,2} = 0*3 + 1*4 = 4 \end{bmatrix}$   |

## Example: 2 x 2 Matrix Multiply (Part 1 of 2)

```
#include <stdio.h>
// header file for SSE compiler intrinsics
#include <emmintrin.h>

// NOTE: vector registers will be represented in
// comments as v1 = [ a | b ]
// where v1 is a variable of type __m128d and
// a, b are doubles

int main(void) {
    // allocate A,B,C aligned on 16-byte boundaries
    double A[4] __attribute__((aligned(16)));
    double B[4] __attribute__((aligned(16)));
    double C[4] __attribute__((aligned(16)));
    int lda = 2;
    int i = 0;
    // declare several 128-bit vector variables
    __m128d c1,c2,a,b1,b2;

    // Initialize A, B, C for example
    /* A = (note column order!)
       1 0
       0 1
    */
    A[0] = 1.0; A[1] = 0.0; A[2] = 0.0; A[3] = 1.0;

    /* B = (note column order!)
       1 3
       2 4
    */
    B[0] = 1.0; B[1] = 2.0; B[2] = 3.0; B[3] = 4.0;

    /* C = (note column order!)
       0 0
       0 0
    */
    C[0] = 0.0; C[1] = 0.0; C[2] = 0.0; C[3] = 0.0;
}
```

10/10/13

Fall 2013 -- Lecture #14

29

## Example: 2 x 2 Matrix Multiply (Part 2 of 2)

```
// used aligned loads to set
// c1 = [c_11 | c_21]
c1 = _mm_load_pd(C+0*lda);
// c2 = [c_12 | c_22]
c2 = _mm_load_pd(C+1*lda);

for (i = 0; i < 2; i++) {
    /* a =
       i = 0: [a_11 | a_21]
       i = 1: [a_12 | a_22]
    */
    a = _mm_load_pd(A+i*lda);
    /* b1 =
       i = 0: [b_11 | b_11]
       i = 1: [b_21 | b_21]
    */
    b1 = _mm_load1_pd(B+i*0*lda);
    /* b2 =
       i = 0: [b_12 | b_12]
       i = 1: [b_22 | b_22]
    */
    b2 = _mm_load1_pd(B+i+1*lda);

    /* c1 =
       i = 0: [c_11 + a_11*b_11 | c_21 + a_21*b_11]
       i = 1: [c_11 + a_21*b_21 | c_21 + a_22*b_21]
    */
    c1 = _mm_add_pd(c1,_mm_mul_pd(a,b1));
    /* c2 =
       i = 0: [c_12 + a_11*b_12 | c_22 + a_21*b_12]
       i = 1: [c_12 + a_21*b_22 | c_22 + a_22*b_22]
    */
    c2 = _mm_add_pd(c2,_mm_mul_pd(a,b2));
}

// store c1,c2 back into C for completion
_mm_store_pd(C+0*lda,c1);
_mm_store_pd(C+1*lda,c2);

// print C
printf("%g,%g\n%g,%g\n",C[0],C[2],C[1],C[3]);
return 0;
}
```

10/10/13

Fall 2013 -- Lecture #14

30

## Inner loop from gcc -O -S

```

L2: movapd  (%rax,%rsi), %xmm1 //Load aligned A[i,i+1]->m1
    movddup (%rdx), %xmm0     //Load B[j], duplicate->m0
    mulpd   %xmm1, %xmm0     //Multiply m0*m1->m0
    addpd   %xmm0, %xmm3     //Add m0+m3->m3
    movddup 16(%rdx), %xmm0   //Load B[j+1], duplicate->m0
    mulpd   %xmm0, %xmm1     //Multiply m0*m1->m1
    addpd   %xmm1, %xmm2     //Add m1+m2->m2
    addq    $16, %rax        // rax+16 -> rax (i+=2)
    addq    $8, %rdx        // rdx+8 -> rdx (j+=1)
    cmpq    $32, %rax       // rax == 32?
    jne     L2              // jump to L2 if not equal
    movapd  %xmm3, (%rcx)    //store aligned m3 into C[k,k+1]
    movapd  %xmm2, (%rdi)    //store aligned m2 into C[l,l+1]

```

10/10/13

Fall 2013 -- Lecture #14

31

## Performance-Driven ISA Extensions

- Subword parallelism, used primarily for multimedia applications
  - Intel MMX: multimedia extension
    - 64-bit registers can hold multiple integer operands
  - Intel SSE: Streaming SIMD extension
    - 128-bit registers can hold several floating-point operands
- Adding instructions that do more work per cycle
  - Shift-add: replace two instructions with one (e.g., multiply by 5)
  - Multiply-add: replace two instructions with one ( $x := c + a \times b$ )
  - Multiply-accumulate: reduce round-off error ( $s := s + a \times b$ )
  - Conditional copy: to avoid some branches (e.g., in if-then-else)

10/10/13

Fall 2013 -- Lecture #14

Slide 32



# Administrivia

10/10/13

Fall 2013 -- Lecture #147

33

# Midterm Review

Topics we've covered

10/10/13

Fall 2013 -- Lecture #14

34


## New-School Machine Structures (It's a bit more complicated!)

*Software*


- **Parallel Requests**  
Assigned to computer  
e.g., Search "Katz"
- **Parallel Threads**  
Assigned to core  
e.g., Lookup, Ads
- **Parallel Instructions**  
>1 instruction @ one time  
e.g., 5 pipelined instructions
- **Parallel Data**  
>1 data item @ one time  
e.g., Add of 4 pairs of words
- **Hardware descriptions**  
All gates functioning in parallel at same time

*Hardware*

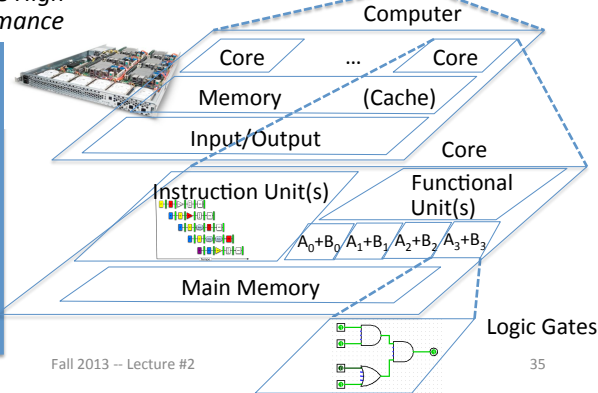
Warehouse Scale Computer



Smart Phone



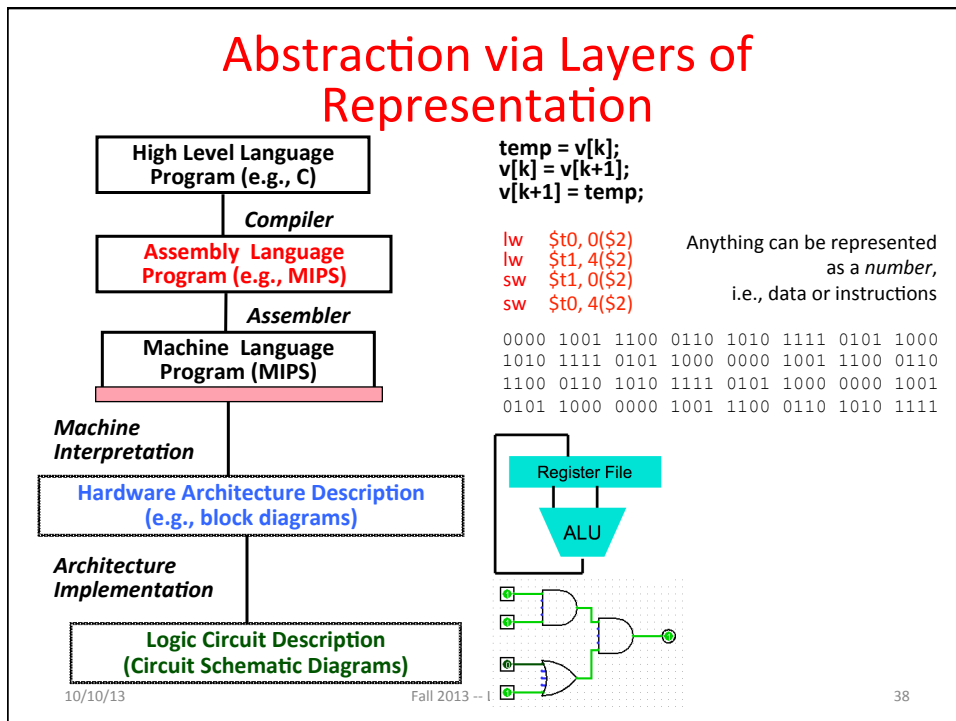
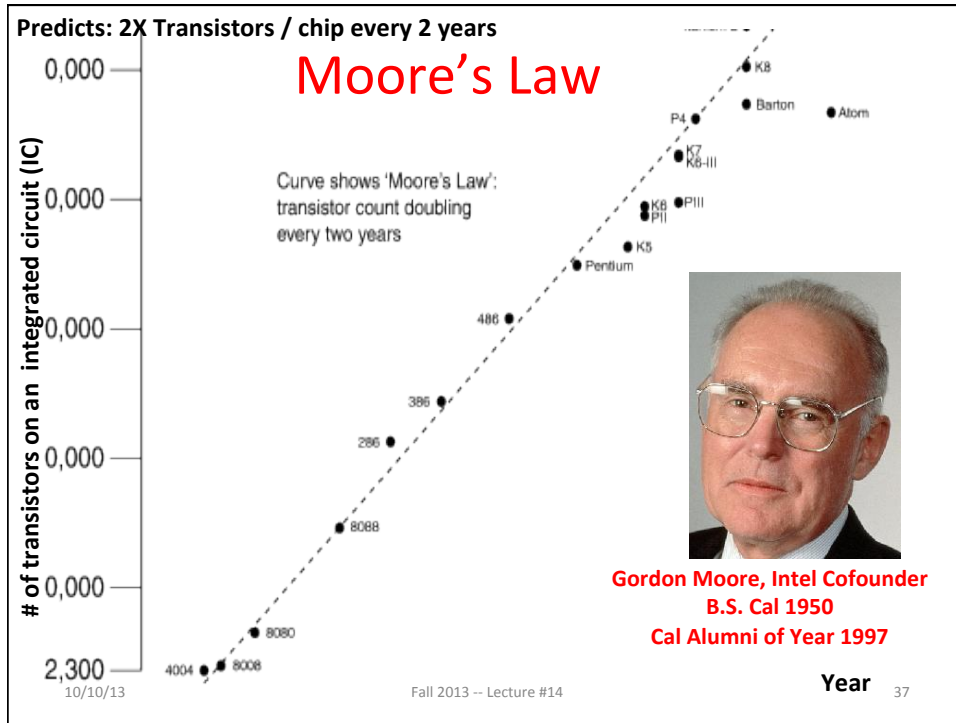
*Harness Parallelism & Achieve High Performance*



Fall 2013 -- Lecture #2

## Great Ideas in Computer Architecture

1. Design for Moore's Law
2. Abstraction to Simplify Design
3. Make the Common Case Fast
4. Dependability via Redundancy
5. Memory Hierarchy
6. Performance via Parallelism/Pipelining/Prediction



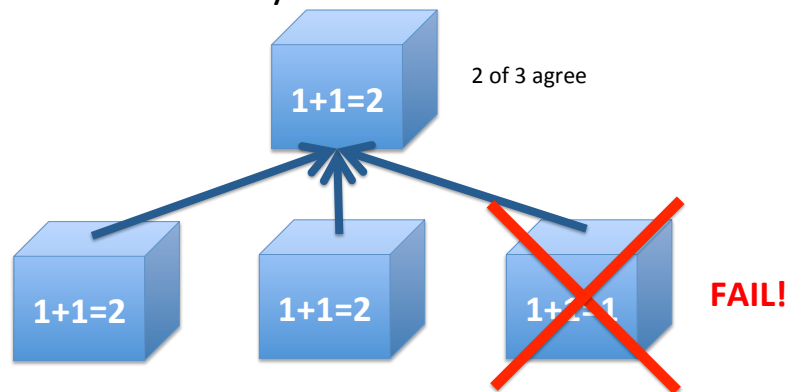
## Make the Common Case Fast

- In making a design tradeoff, favor the common over the infrequent case
- Don't spend time optimizing code that is run infrequently
- Choose your performance metric and use measurement to determine the common case



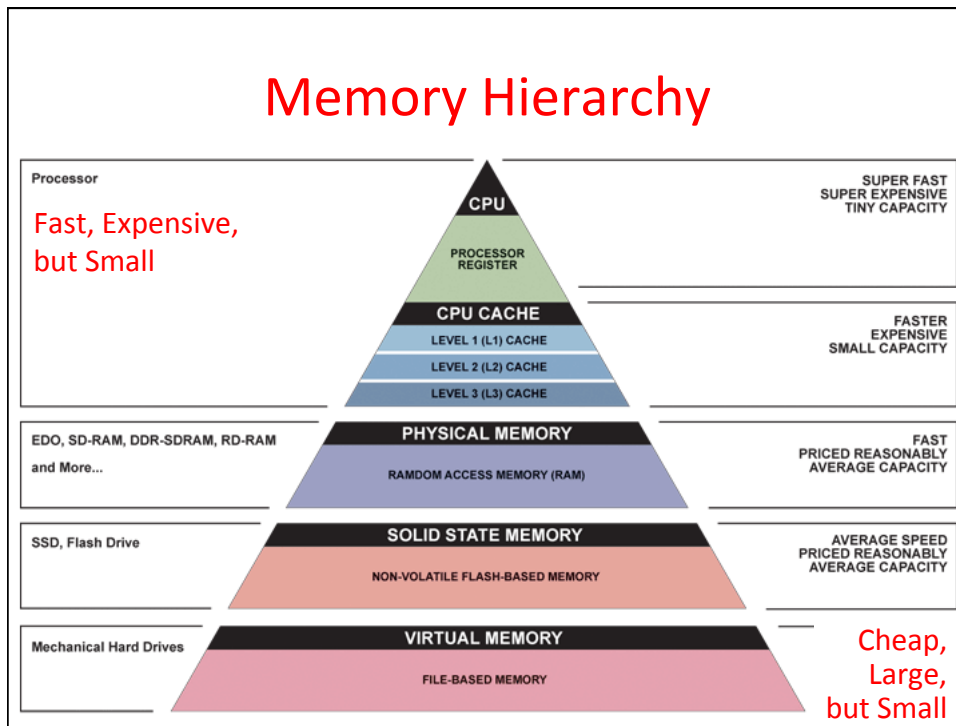
## Dependability via Redundancy

- Redundancy so that a failing piece doesn't make the whole system fail

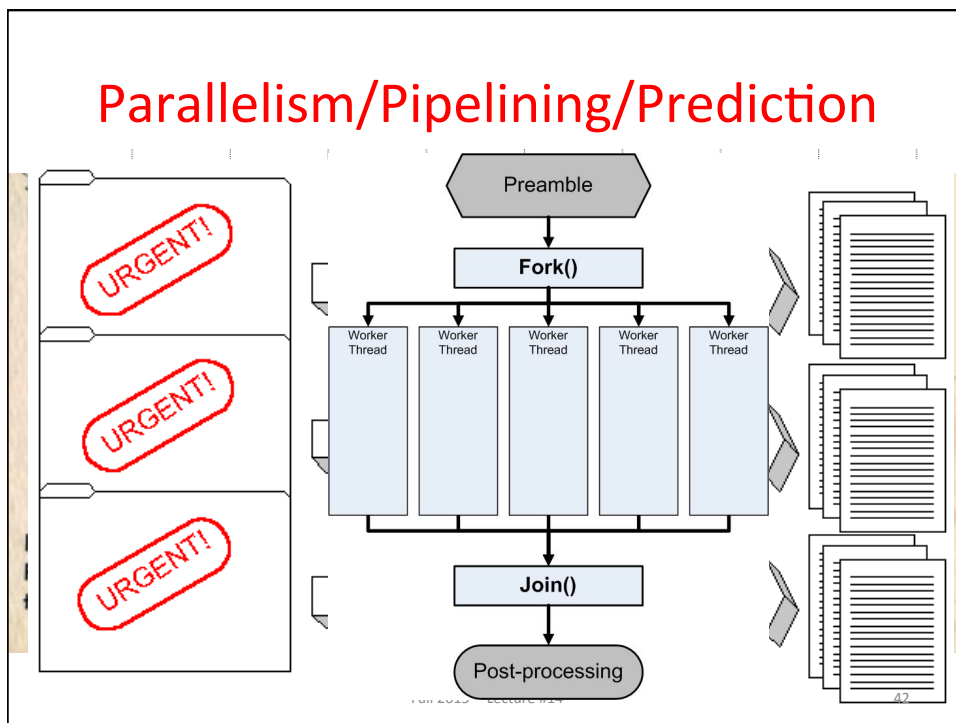


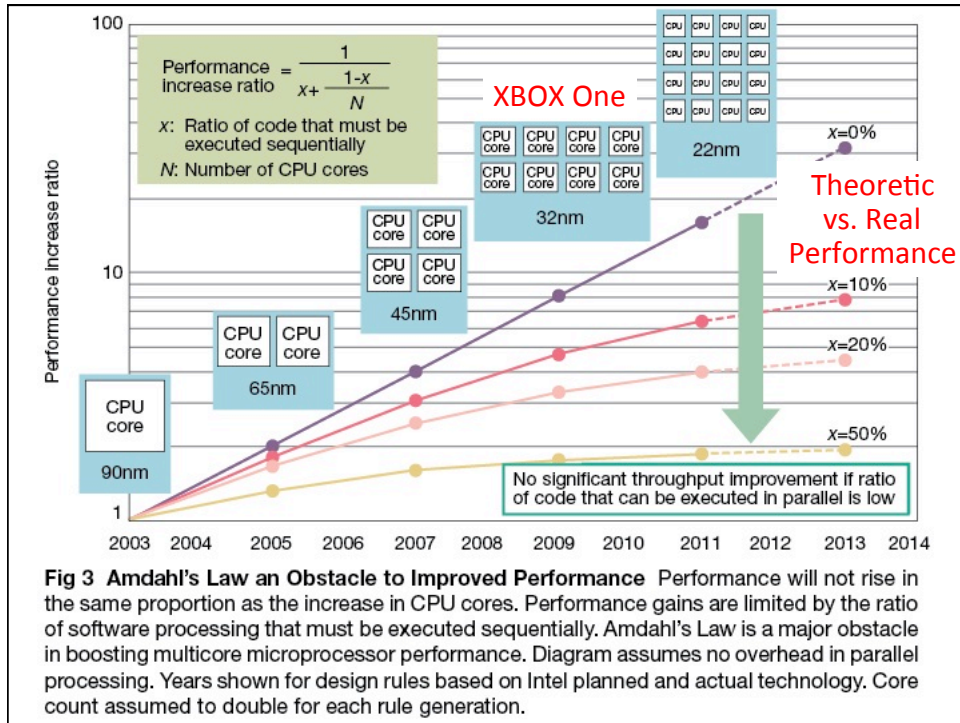
Increasing transistor density reduces the cost of redundancy

# Memory Hierarchy



# Parallelism/Pipelining/Prediction





## Warehouse-Scale Computers

- Power Usage Effectiveness
- Request-Level Parallelism
- MapReduce
- Handling failures
- Costs of WSC

# C Language and Compilation

- C Types, including Structs, Consts, Enums
- Arrays and strings
- C Pointers
- C functions and parameter passing

# MIPS Instruction Set

- ALU operations
- Loads/Stores
- Branches/Jumps
- Registers
- Memory
- Function calling conventions
- Stack

**OPCODES, BASE CONVERSION, ASCII SYMBOLS**

| OPCODE | Function | Format         | Op1 | Op2 | Op3 | Op4 | Op5 | Op6 | Op7 | Op8 | Op9 | Op10 | Op11 | Op12 | Op13 | Op14 | Op15 | Op16 | Op17 | Op18 | Op19 | Op20 | Op21 | Op22 | Op23 | Op24 | Op25 | Op26 | Op27 | Op28 | Op29 | Op30 | Op31 | Op32 | Op33 | Op34 | Op35 | Op36 | Op37 | Op38 | Op39 | Op40 | Op41 | Op42 | Op43 | Op44 | Op45 | Op46 | Op47 | Op48 | Op49 | Op50 | Op51 | Op52 | Op53 | Op54 | Op55 | Op56 | Op57 | Op58 | Op59 | Op60 | Op61 | Op62 | Op63 | Op64 | Op65 | Op66 | Op67 | Op68 | Op69 | Op70 | Op71 | Op72 | Op73 | Op74 | Op75 | Op76 | Op77 | Op78 | Op79 | Op80 | Op81 | Op82 | Op83 | Op84 | Op85 | Op86 | Op87 | Op88 | Op89 | Op90 | Op91 | Op92 | Op93 | Op94 | Op95 | Op96 | Op97 | Op98 | Op99 | Op100 | Op101 | Op102 | Op103 | Op104 | Op105 | Op106 | Op107 | Op108 | Op109 | Op110 | Op111 | Op112 | Op113 | Op114 | Op115 | Op116 | Op117 | Op118 | Op119 | Op120 | Op121 | Op122 | Op123 | Op124 | Op125 | Op126 | Op127 | Op128 | Op129 | Op130 | Op131 | Op132 | Op133 | Op134 | Op135 | Op136 | Op137 | Op138 | Op139 | Op140 | Op141 | Op142 | Op143 | Op144 | Op145 | Op146 | Op147 | Op148 | Op149 | Op150 | Op151 | Op152 | Op153 | Op154 | Op155 | Op156 | Op157 | Op158 | Op159 | Op160 | Op161 | Op162 | Op163 | Op164 | Op165 | Op166 | Op167 | Op168 | Op169 | Op170 | Op171 | Op172 | Op173 | Op174 | Op175 | Op176 | Op177 | Op178 | Op179 | Op180 | Op181 | Op182 | Op183 | Op184 | Op185 | Op186 | Op187 | Op188 | Op189 | Op190 | Op191 | Op192 | Op193 | Op194 | Op195 | Op196 | Op197 | Op198 | Op199 | Op200 | Op201 | Op202 | Op203 | Op204 | Op205 | Op206 | Op207 | Op208 | Op209 | Op210 | Op211 | Op212 | Op213 | Op214 | Op215 | Op216 | Op217 | Op218 | Op219 | Op220 | Op221 | Op222 | Op223 | Op224 | Op225 | Op226 | Op227 | Op228 | Op229 | Op230 | Op231 | Op232 | Op233 | Op234 | Op235 | Op236 | Op237 | Op238 | Op239 | Op240 | Op241 | Op242 | Op243 | Op244 | Op245 | Op246 | Op247 | Op248 | Op249 | Op250 | Op251 | Op252 | Op253 | Op254 | Op255 |     |
|--------|----------|----------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-----|
| 000000 | and      | r, r2, r1, imm | 0   | 1   | 2   | 3   | 4   | 5   | 6   | 7   | 8   | 9    | 10   | 11   | 12   | 13   | 14   | 15   | 16   | 17   | 18   | 19   | 20   | 21   | 22   | 23   | 24   | 25   | 26   | 27   | 28   | 29   | 30   | 31   | 32   | 33   | 34   | 35   | 36   | 37   | 38   | 39   | 40   | 41   | 42   | 43   | 44   | 45   | 46   | 47   | 48   | 49   | 50   | 51   | 52   | 53   | 54   | 55   | 56   | 57   | 58   | 59   | 60   | 61   | 62   | 63   | 64   | 65   | 66   | 67   | 68   | 69   | 70   | 71   | 72   | 73   | 74   | 75   | 76   | 77   | 78   | 79   | 80   | 81   | 82   | 83   | 84   | 85   | 86   | 87   | 88   | 89   | 90   | 91   | 92   | 93   | 94   | 95   | 96   | 97   | 98   | 99    | 100   | 101   | 102   | 103   | 104   | 105   | 106   | 107   | 108   | 109   | 110   | 111   | 112   | 113   | 114   | 115   | 116   | 117   | 118   | 119   | 120   | 121   | 122   | 123   | 124   | 125   | 126   | 127   | 128   | 129   | 130   | 131   | 132   | 133   | 134   | 135   | 136   | 137   | 138   | 139   | 140   | 141   | 142   | 143   | 144   | 145   | 146   | 147   | 148   | 149   | 150   | 151   | 152   | 153   | 154   | 155   | 156   | 157   | 158   | 159   | 160   | 161   | 162   | 163   | 164   | 165   | 166   | 167   | 168   | 169   | 170   | 171   | 172   | 173   | 174   | 175   | 176   | 177   | 178   | 179   | 180   | 181   | 182   | 183   | 184   | 185   | 186   | 187   | 188   | 189   | 190   | 191   | 192   | 193   | 194   | 195   | 196   | 197   | 198   | 199   | 200   | 201   | 202   | 203   | 204   | 205   | 206   | 207   | 208   | 209   | 210   | 211   | 212   | 213   | 214   | 215   | 216   | 217   | 218   | 219   | 220   | 221   | 222   | 223   | 224   | 225   | 226   | 227   | 228   | 229   | 230   | 231   | 232   | 233   | 234   | 235   | 236   | 237   | 238   | 239   | 240   | 241   | 242   | 243   | 244   | 245   | 246   | 247   | 248   | 249   | 250   | 251   | 252   | 253   | 254   | 255 |

**STANDARD**

IEEE 754 Single Precision and Double Precision Formats

**MEMORY ALLOCATION**

**DATA ALIGNMENT**

**EXCEPTION CONTROL, REGISTERS, CAUSE AND STATUS**

**EXCEPTION CODES**

**SIZE PREFIXES (10<sup>6</sup> for Disk, Communications; 2<sup>6</sup> for Memory)**

## Everything is a Number

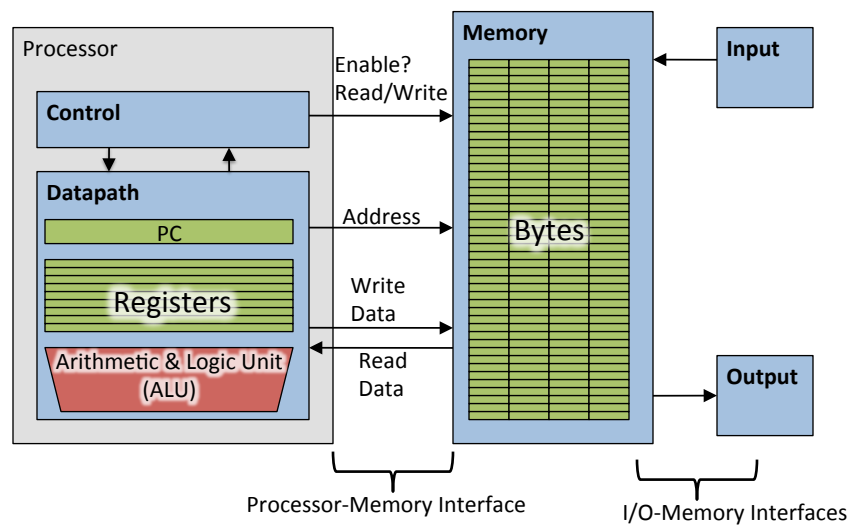
- Binary
- Signed versus Unsigned
- One's Complement/Two's Complement
- Floating-Point numbers
- Character Strings
- Instruction Encoding

10/10/13

Fall 2013 -- Lecture #14

47

## Components of a Computer



10/10/13

Fall 2013 -- Lecture #140

48



## Caches

- Spatial/Temporal Locality
- Instruction versus Data
- Block size, capacity
- Direct-Mapped cache
- 3 C's
- Cache-aware performance programming

10/10/13

Fall 2013 -- Lecture #14

49

## Parallelism

- SIMD/MIMD/MISD/SISD
- Amdahl's Law
- Strong vs. weak scaling
- Data-Parallel execution

10/10/13

Fall 2013 -- Lecture #14

50

## And in Conclusion, ...

- Intel SSE SIMD Instructions
  - Exploit data-level parallelism in loops
  - One instruction fetch that operates on multiple operands simultaneously
  - 128-bit XMM registers
- SSE Instructions in C
  - Embed the SSE machine instructions directly into C programs through use of intrinsics
  - Achieve efficiency beyond that of optimizing compiler