

## Recap of Last Class

- Generally useful optimizations
  - Code motion/precomputation
  - Strength reduction
  - Common subexpression elimination and sharing
  - Removing unnecessary procedure calls
  - Loop unrolling
- Optimization blockers
  - Procedure calls
  - Memory aliasing
- Modern processors and instruction-level parallelism

48

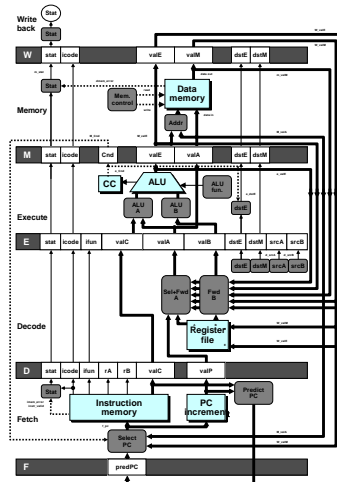
## Hazards

- Structural (e.g., instruction/data fetch)
- Data (e.g., add followed by sub reading dst register of add)
- Control (e.g., jeq)

10/28/2020

49

Figure 4.52:  
Pipelined  
Implementation  
as Discussed in  
Class



50

## Loop Unrolling

<pre> irmovl \$5, %edx irmovl \$80, %ebx Loop:   rmmovl array_base(%ebx), %eax   addl  %edx, %eax   rmmovl %eax, array_base(%ebx)   rrmovl %ebx, %esi   addl  \$-4, %esi   jne   Loop   ... </pre>	<pre> irmovl \$5, %edx irmovl \$80, %ebx Loop:   rmmovl array_base(%ebx), %eax   addl  %edx, %eax   rmmovl %eax, array_base(%ebx)   addl  \$-4, %ebx   jne   Loop   ... </pre>
--	--

51

## Instruction-level Parallelism

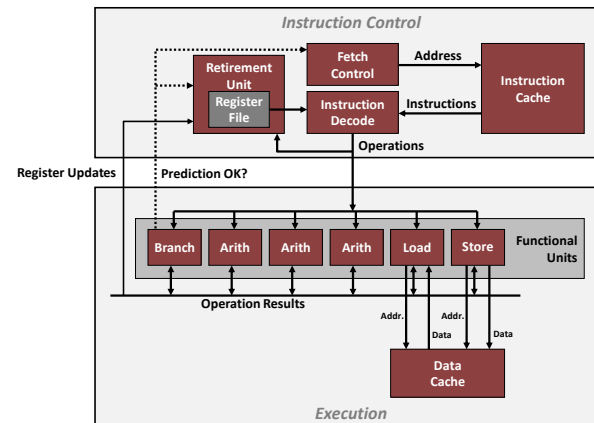
- Pipelining/super-pipelining
- Out-of-order execution
- Super-scalar
  - multiple instructions per pipeline stage
  - Dependences handled in hardware
- Very Large Instruction Word (VLIW)
  - Multiple instructions per pipeline stage
  - Dependences taken care of by compiler

10/28/2020

52

52

## Modern CPUs: Pipelined, Superscalar, Out-of-Order



53

## Types of Data Hazards (Dependencies)

- RAW – true dependence
- WAR – anti-dependence
- WAW – output dependence

10/28/2020

54

54

## Enabling Out-Of-Order Execution

- Tomasulo's generalized data forwarding algorithm in the IBM360
  - Register renaming
  - Reservation stations
  - Common data bus to broadcast data to reservation stations

10/28/2020

55

55

## Advanced Processor Design Techniques

- Trace caches
- Register renaming – eliminate WAW, WAR hazards
  - Register map table
  - Free list
  - Active list
- Speculative execution
- Value prediction
- Branch prediction
  - Branch history tables for prediction
  - Branch stack to save state prior to branch
  - Branch mask to determine instructions that must be squashed

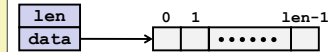
10/28/2020

56

56

## Benchmark Example: Data Type for Vectors

```
/* data structure for vectors */
typedef struct{
    size_t len;
    data_t *data;
} vec;
```



### •Data Types

- Use different declarations for data\_t
- int
- long
- float
- double

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

57

## Benchmark Computation

```
void combin1(vec_ptr v, data_t *dest)
{
    long int i;
    *dest = IDENT;
    for (i = 0; i < vec_length(v); i++) {
        data_t val;
        get_vec_element(v, i, &val);
        *dest = *dest OP val;
    }
}
```

Compute sum or product of vector elements

### •Data Types

- Use different declarations for data\_t
- int
- long
- float
- double

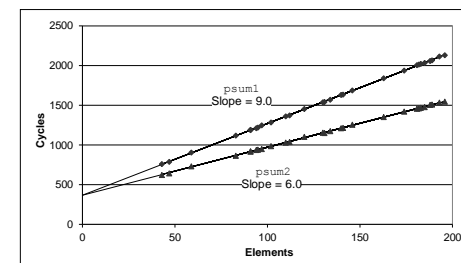
### •Operations

- Use different definitions of OP and IDENT
- + / 0
- \* / 1

58

## Cycles Per Element (CPE)

- Convenient way to express performance of program that operates on vectors or lists
- Length = n
- In our case: **CPE = cycles per OP**
- $T = CPE * n + \text{Overhead}$ 
  - CPE is slope of line



59

## Benchmark Performance

```
void combine1(vec_ptr v, data_t *dest)
{
    long int i;
    *dest = IDENT;
    for (i = 0; i < vec_length(v); i++) {
        data_t val;
        get_vec_element(v, i, &val);
        *dest = *dest OP val;
    }
}
```

Compute sum or product of vector elements

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Combine1 unoptimized	22.68	20.02	19.98	20.18
Combine1 -O1	10.12	10.12	10.17	11.14

60

## Basic Optimizations

```
void combine4(vec_ptr v, data_t *dest)
{
    long i;
    long length = vec_length(v);
    data_t *d = get_vec_start(v);
    data_t t = IDENT;
    for (i = 0; i < length; i++)
        t = t OP d[i];
    *dest = t;
}
```

- Move vec\_length out of loop
- Avoid bounds check on each cycle
- Accumulate in temporary

61

## Effect of Basic Optimizations

```
void combine4(vec_ptr v, data_t *dest)
{
    long i;
    long length = vec_length(v);
    data_t *d = get_vec_start(v);
    data_t t = IDENT;
    for (i = 0; i < length; i++)
        t = t OP d[i];
    *dest = t;
}
```

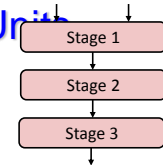
Method	Integer		Double FP	
	Add	Mult	Add	Mult
Combine1 -O1	10.12	10.12	10.17	11.14
Combine4	1.27	3.01	3.01	5.01

- Eliminates sources of overhead in loop

62

## Pipelined Functional Units

```
long mult_eg(long a, long b, long c) {
    long p1 = a*b;
    long p2 = a*c;
    long p3 = p1 * p2;
    return p3;
}
```



	Time						
	1	2	3	4	5	6	7
Stage 1	a*b	a*c			p1*p2		
Stage 2		a*b	a*c			p1*p2	
Stage 3			a*b	a*c			p1*p2

- Divide computation into stages
- Pass partial computations from stage to stage
- Stage i can start on new computation once values passed to i+1
- E.g., complete 3 multiplications in 7 cycles, even though each requires 3 cycles

63

## Haswell CPU

- 8 Total Functional Units
- Multiple instructions can execute in parallel
  - 2 load, with address computation
  - 1 store, with address computation
  - 4 integer
  - 2 FP multiply
  - 1 FP add
  - 1 FP divide
- Some instructions take > 1 cycle, but can be pipelined

Instruction	Latency	Cycles/Issue
Load / Store	4	1
Integer Multiply	3	1
<b>Integer/Long Divide</b>	<b>3-30</b>	<b>3-30</b>
Single/Double FP Multiply	5	1
Single/Double FP Add	3	1
<b>Single/Double FP Divide</b>	<b>3-15</b>	<b>3-15</b>

64

## x86-64 Compilation of Combine4

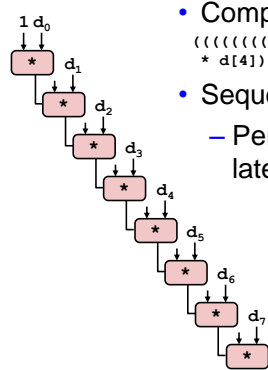
- Inner Loop (Case: Integer Multiply)

```
.L519:                # Loop:
imull  (%rax,%rdx,4), %ecx  # t = t * d[i]
addq  $1, %rdx             # i++
cmpq  %rdx, %rbp          # Compare length:i
jg    .L519                # If >, goto Loop
```

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Combine4	1.27	3.01	3.01	5.01
Latency Bound	1.00	3.00	3.00	5.00

65

## Combine4 = Serial Computation (OP = \*)



- Computation (length=8)
 
$$(((((((1 * d[0]) * d[1]) * d[2]) * d[3]) * d[4]) * d[5]) * d[6]) * d[7])$$
- Sequential dependence
  - Performance: determined by latency of OP

66

## Loop Unrolling (2x1)

```
void unroll2a_combine(vec_ptr v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x = IDENT;
    long 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;
}
```

- Perform 2x more useful work per iteration

67

## Effect of Loop Unrolling

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Combine4	1.27	3.01	3.01	5.01
Unroll 2x1	1.01	3.01	3.01	5.01
Latency Bound	1.00	3.00	3.00	5.00

- Helps integer add
  - Achieves latency bound
- Others don't improve. *Why?*
  - Still sequential dependency

```
x = (x OP d[i]) OP d[i+1];
```

68

## Loop Unrolling with Separate Accumulators (2x2)

```
void unroll2a_combine(vec_ptr v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x0 = IDENT;
    data_t x1 = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x0 = x0 OP d[i];
        x1 = x1 OP d[i+1];
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x0 = x0 OP d[i];
    }
    *dest = x0 OP x1;
}
```

- Different form of reassociation

72

## Effect of Separate Accumulators

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Combine4	1.27	3.01	3.01	5.01
Unroll 2x1	1.01	3.01	3.01	5.01
Unroll 2x1a	1.01	1.51	1.51	2.51
Unroll 2x2	0.81	1.51	1.51	2.51
Latency Bound	1.00	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50

- Int + makes use of two load units

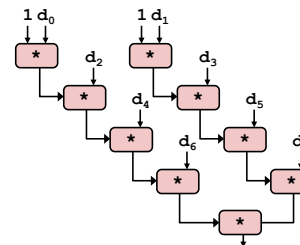
```
x0 = x0 OP d[i];
x1 = x1 OP d[i+1];
```

- 2x speedup (over unroll2) for Int \*, FP +, FP \*

73

## Separate Accumulators

```
x0 = x0 OP d[i];
x1 = x1 OP d[i+1];
```



### What changed:

- Two independent "streams" of operations

### Overall Performance

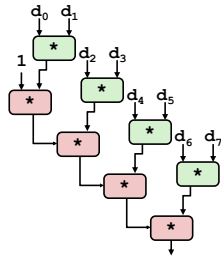
- N elements, D cycles latency/op
- Should be  $(N/2+1)*D$  cycles:  
**CPE = D/2**
- CPE matches prediction!

*What Now?*

74

## 2x1a: Reassociated Computation

```
x = x OP (d[i] OP d[i+1]);
```



- What changed:
  - Ops in the next iteration can be started early (no dependency)
- Overall Performance
  - N elements, D cycles latency/op
  - $(N/2+1)*D$  cycles:  
**CPE = D/2**

75

## Unrolling & Accumulating

- Idea
  - Can unroll to any degree L
  - Can accumulate K results in parallel
  - L must be multiple of K
- Limitations
  - Diminishing returns
    - Cannot go beyond throughput limitations of execution units
  - Large overhead for short lengths
    - Finish off iterations sequentially
  - Register spilling

76

## Achievable Performance

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Best	0.54	1.01	1.01	0.52
Latency Bound	1.00	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50

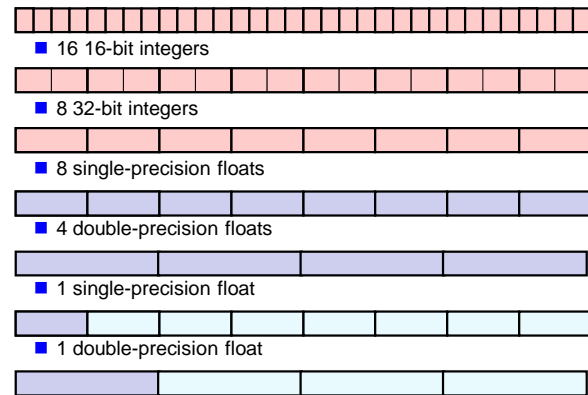
- Limited only by throughput of functional units
- Up to 42X improvement over original, unoptimized code

79

## Programming with AVX2

YMM Registers

- 16 total, each 32 bytes
- 32 single-byte integers
- 16 16-bit integers
- 8 32-bit integers
- 8 single-precision floats
- 4 double-precision floats
- 1 single-precision float
- 1 double-precision float



80

## SIMD Operations

- SIMD Operations: Single Precision  
`vaddsd %ymm0, %ymm1, %ymm1`

- SIMD Operations: Double Precision  
`vaddpd %ymm0, %ymm1, %ymm1`

81

## Using Vector Instructions

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Scalar Best	0.54	1.01	1.01	0.52
Vector Best	0.06	0.24	0.25	0.16
Latency Bound	0.50	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50
Vec Throughput Bound	0.06	0.12	0.25	0.12

- Make use of AVX Instructions
  - Parallel operations on multiple data elements
  - See Web Aside OPT:SIMD on CS:APP web page

82

## So far ...

- We have assumed that memory is one monolithic block that returns data within one cycle from when it is requested

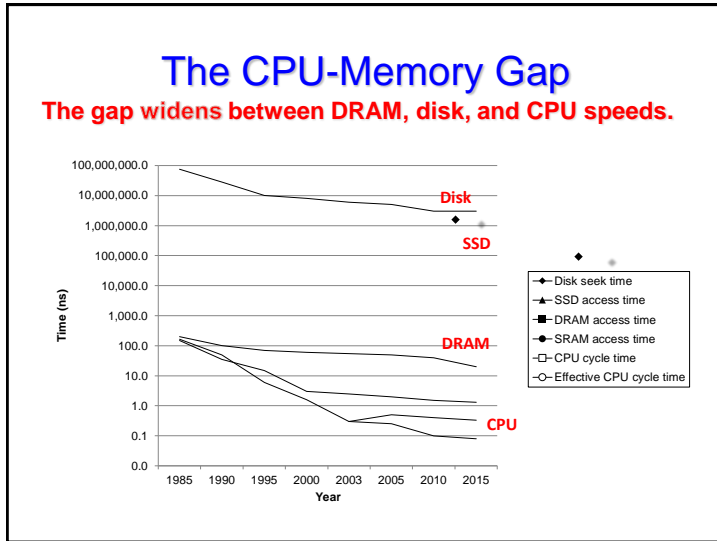
83

## Time Scales

- Absolute Time
  - Typically use nanoseconds
    - $10^{-9}$  seconds
  - Time scale of computer instructions
- Clock Cycles
  - Most computers controlled by high frequency clock signal
  - Examples
    - 100 MHz
      - $10^8$  cycles per second
      - Clock period = 10ns
    - 2 GHz
      - $2 \times 10^9$  cycles per second
      - Clock period = 0.5ns

84





85

## Locality

- Principle of Locality:
  - Programs tend to reuse data and instructions near those they have used recently, or that were recently referenced themselves.
  - Temporal locality:** Recently referenced items are likely to be referenced in the near future.
  - Spatial locality:** Items with nearby addresses tend to be referenced close together in time.

**Locality Example:**

```

sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;

```

- Data**
  - Reference array elements in succession (stride-1 reference pattern): **Spatial locality**
  - Reference `sum` each iteration: **Temporal locality**
- Instructions**
  - Reference instructions in sequence: **Spatial locality**
  - Cycle through loop repeatedly: **Temporal locality**

86

## Locality Example

- Claim:** Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer.
- Question:** Does this function have good locality?

```

int sumarrayrows(int a[M][N])
{
    int i, j, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum
}

```

87

## Locality Example

- Question:** Does this function have good locality?

```

int sumarraycols(int a[M][N])
{
    int i, j, sum = 0;

    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum
}

```

88

## Locality Example

- **Question:** Can you permute the loops so that the function scans the 3-d array `a[]` with a stride-1 reference pattern (and thus has good spatial locality)?

```
int sumarray3d(int a[M][N][N])
{
    int i, j, k, sum = 0;

    for (i = 0; i < N; i++)
        for (j = 0; j < N; j++)
            for (k = 0; k < M; k++)
                sum += a[k][i][j];

    return sum
}
```

89

## Memory Hierarchies

- Some fundamental and enduring properties of hardware and software:
  - Fast storage technologies cost more per byte and have less capacity.
  - The gap between CPU and main memory speed is widening.
  - Well-written programs tend to exhibit good locality.
- These fundamental properties complement each other beautifully.
- They suggest an approach for organizing memory and storage systems known as a **memory hierarchy**.

90