

The final exam will cover material mainly from 04-06c and on, and projects 5-8.

Table of Contents

[Lecture 4: Processor Architecture - Y86-64 ISA, Logic Circuits](#)

[Lecture 4b: Processor Architecture - Y86-64 ISA, Logic Circuits](#)

[Assignment 5 - Y86-64](#)

[Lecture 5: Optimizing Program Performance - Limitation of Optimizing Compilers](#)

[Lecture 5b: Optimizing Program Performance - Modern Processors, Loop Unrolling, Branch Prediction](#)

[Assignment 6 - Optimization](#)

[Lecture 6: The Memory Hierarchy - Memory Technologies, Locality](#)

[Lecture 6b: The Memory Hierarchy - Caching in the Memory Hierarchy](#)

[Assignment 7 - Bitwise for Tag / Set / Offset](#)

[Lecture 6c: The Memory Hierarchy - Direct Mapped Caches, Set-Associative Caches](#)

[Assignment 8 - Cache Hits and Misses](#)

Instruction Set Architecture

- Assembly Language View
 - Processor state... Registers, memory...
 - Instructions
 - addq, pushq, ret....
 - How instructions are encoded as bytes
- Layers of Abstraction
 - Above: how to program machine
 - Processor executes instructions in a sequence
 - Below: what needs to be built
 - Use variety of tricks to make it run fast.
 - E.g. execute multiple instructions simultaneously

Y86-64 Processor State

- Program Registers
 - 15 registers (omit %r15). Each 64 bits
- Condition Codes
 - Single-bit flags set by arithmetic or logical instructions
 - ZF: Zero
 - SF: Negative
 - OF: Overflow
- Program Counter
 - Indicates address of next instructions
- Program Status
 - Indicates either normal operation or some error condition
- Memory
 - Byte-addressable storage array
 - Words stored in little-endian byte order

Y86-64 Instructions

- Format
 - 1 - 10 bytes of information read from memory
 - Can determine instruction length from first byte
 - Not as many instruction types, and simpler encoding than with x86-64
 - Each accesses and modifies some part(s) of the program state

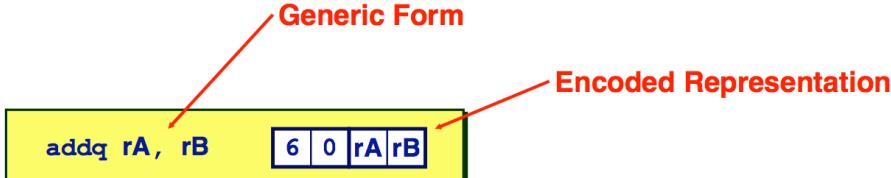
Encoding Registers

- Each register has 4-bit ID
- Register ID 15 (0xF) indicates "no register"

%rax	0	%r8	8
%rcx	1	%r9	9
%rdx	2	%r10	A
%rbx	3	%r11	B
%rsp	4	%r12	C
%rbp	5	%r13	D
%rsi	6	%r14	E
%rdi	7	No Register	F

Instruction Example

- Addition Instruction

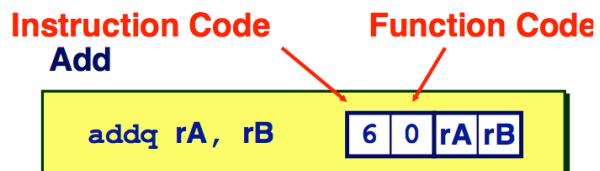


- Add value in register rA to that in register rB
 - Store result in register rB
 - Note that Y86-64 only allows addition to be applied to register data
- Set condition codes based on result
- E.g. `addq, %rax, %rsi` Encoding 60 06
- Two-byte encoding

- First indicates instruction type
- Second gives source and destination registers

Arithmetic and Logical Operations

- Add
- Subtract
- And
- Exclusive-Or
 - Refer to generically as "OPq"
 - Set condition codes as side effects



Subtract (rA from rB)



And



Exclusive-Or



Move Operations

- rrmovq
- irmovq
- removq
- mrmovq
 - Like the x86-64 move instruction
 - Simpler format for memory addresses
 - Give different names to keep them distinct

Register → Register



Immediate → Register



Register → Memory



Memory → Register



Move Instruction Examples

X86-64

movq \$0xabcd, %rdx

Y86-64

irmovq \$0xabcd, %rdx

Encoding: 30 82 cd ab 00 00 00 00 00 00

movq %rsp, %rbx

rrmovq %rsp, %rbx

Encoding: 20 43

movq -12(%rbp),%rcx

mrmovq -12(%rbp),%rcx

Encoding: 50 15 f4 ff ff ff ff ff ff ff ff

movq %rsi,0x41c(%rsp)

rmmovq %rsi,0x41c(%rsp)

Encoding: 40 64 1c 04 00 00 00 00 00 00 00

Jump Instructions

Jump Unconditionally

`jmp Dest`

7	0
---	---

 Dest

Jump When Less or Equal

`jle Dest`

7	1
---	---

 Dest

Jump When Less

`jl Dest`

7	2
---	---

 Dest

Jump When Equal

`je Dest`

7	3
---	---

 Dest

Jump When Not Equal

`jne Dest`

7	4
---	---

 Dest

Jump When Greater or Equal

`jge Dest`

7	5
---	---

 Dest

Jump When Greater

`jg Dest`

7	6
---	---

 Dest

Y86-64
Program

Stack

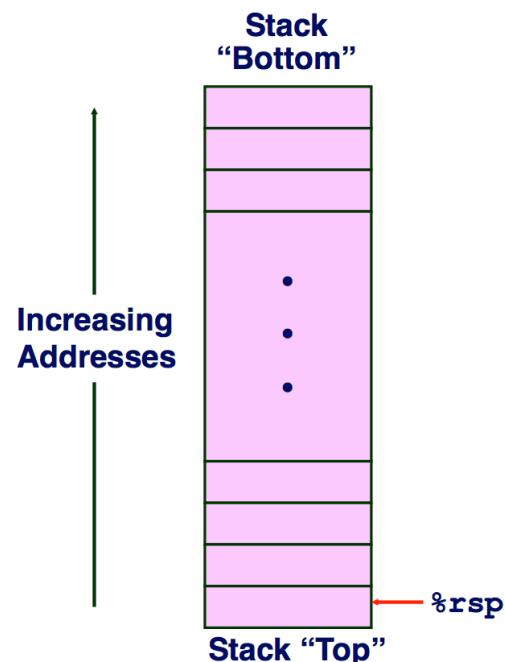
- Region of memory holding data
- Used in Y86-64 supporting procedure calls
- Stack top indicated by %rsp
 - Address of top stack element
- Stack grows toward lower addresses
 - Top element is at highest address in the stack
 - When pushing, must first decrement stack pointer
 - After popping, increment stack pointer

Stack Operations

- `pushq`
 - Decrement %rsp by 8
 - Store word from rA to memory at %rsp
 - Like x86-64

`pushq rA`

A	0	rA	F
---	---	----	---



- `popq`
 - Read word from memory at %rsp
 - Save in rA
 - Increment %rsp by 8
 - Like x86-64

`popq rA`

B	0	rA	F
---	---	----	---

Subroutine Call and Return

- call
 - Push address of next instruction onto stack
 - Start executing instructions at Dest
 - Like x86-64



- ret

- Pop value from stack
- Use as address for next instruction
- Like x86-64



- halt
 - Stop executing instructions
 - X86-64 has comparable instruction, but can't execute in user mode
 - We will use it to stop the simulator
 - Encoding ensures that program hitting memory initialized to zero will haunt



Y86-64 Sample Program Structure

- Program starts at address 0
- Must set up stack
 - Where located
 - Pointer values
 - Make sure don't overwrite code!
- Must initialize data

```
init:          # Initialization
    ...
    call Main
    halt

    .align 8      # Program data
array:
    ...

Main:         # Main function
    ...
    call len    ...

len:          # Length function
    ...
    .pos 0x100   # Placement of stack
Stack:
```

Overview of Logic Design

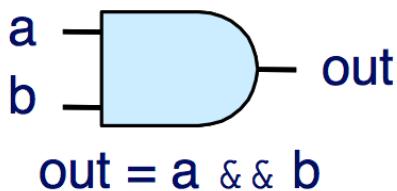
- Fundamental Hardware Requirements
 - Communication
 - How to get values from one place to another
 - Computation
 - Storage
- Bits are Our Friends
 - Everything expressed in terms of values 0 and 1
 - Communication
 - Low or high voltage on wire
 - Computation
 - Compute Boolean functions
 - Storage
 - Store bits of information

Computing with Logic Gates

- AND

* Outputs are Boolean functions of inputs

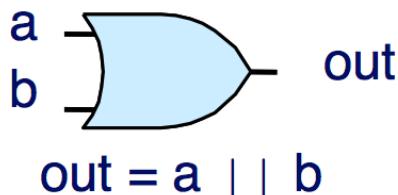
And



* Respond continuously to changes in inputs

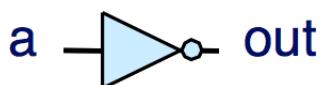
- OR

Or



- NOT

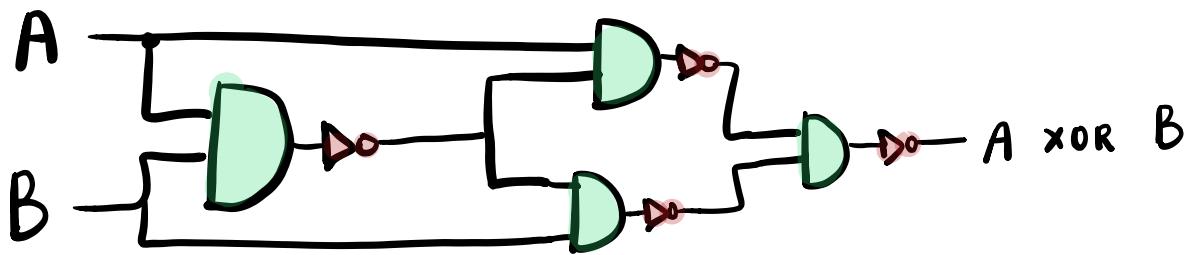
Not



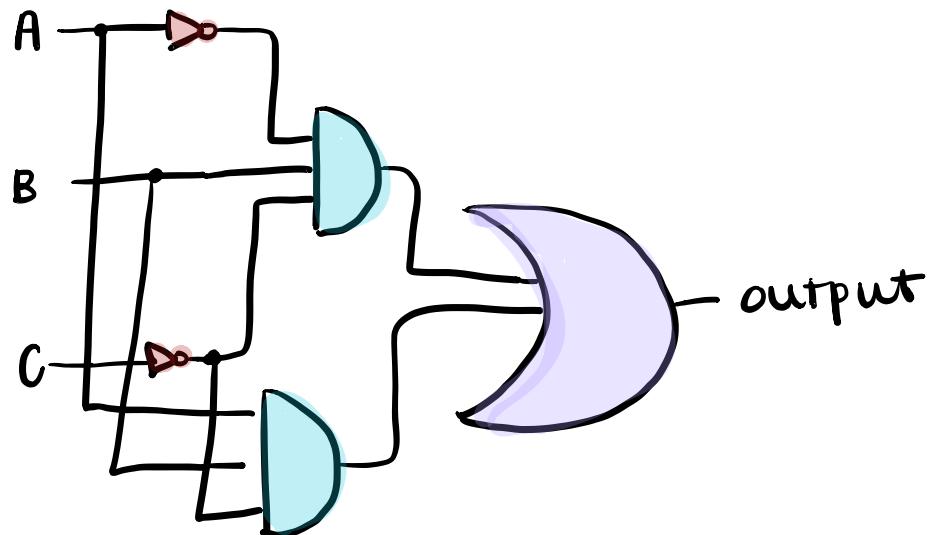
$$\text{out} = !a$$

Examples:

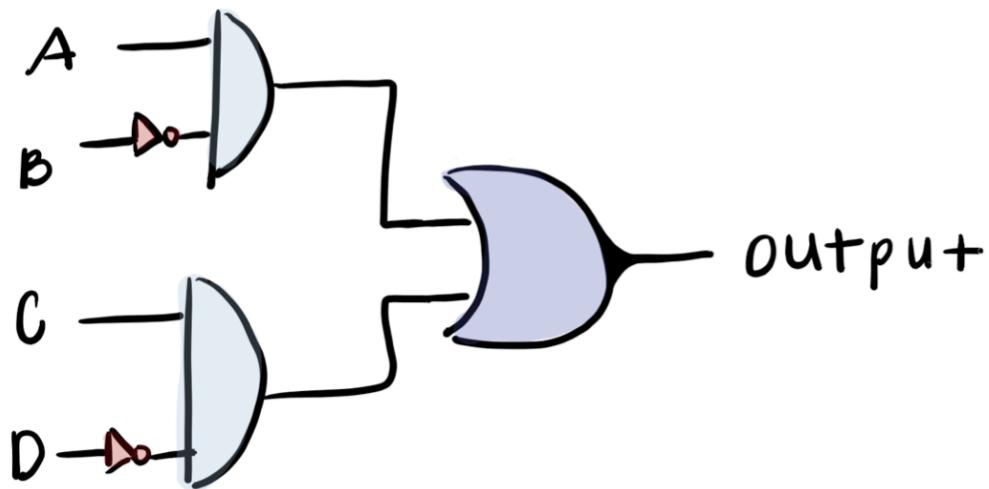
XOR



$$(\neg A \text{ && } B \text{ && } \neg C) \parallel (A \text{ && } B \text{ && } \neg C)$$



$$(A \text{ && } \neg B) \parallel (C \text{ && } \neg D)$$



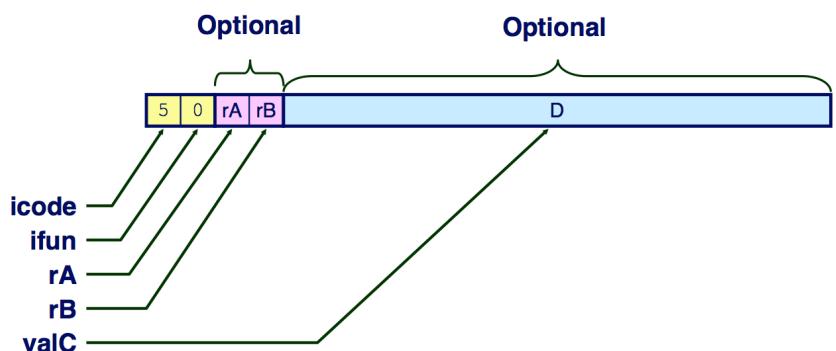
SEQ Stages

- Fetch
 - Read instruction from instruction memory
- Decode
 - Read program registers
- Execute
 - Compute value or address
- Memory
 - Read or write data
- Write back
 - Write program registers
- PC
 - Update program counter

Instruction Decoding

- Instruction Format
 - Instruction byte
 - Optional register byte
 - Optional constant word

icode : ifun
rA : rB
valC



Computed Values

Fetch

<code>icode</code>	Instruction code
<code>ifun</code>	Instruction function
<code>rA</code>	Instr. Register A
<code>rB</code>	Instr. Register B
<code>valC</code>	Instruction constant
<code>valP</code>	Incremented PC

Decode

<code>srcA</code>	Register ID A
<code>srcB</code>	Register ID B
<code>dstE</code>	Destination Register E
<code>dstM</code>	Destination Register M
<code>valA</code>	Register value A
<code>valB</code>	Register value B

Execute

- `valE` ALU result
- `Cnd` Branch/move flag

Memory

- `valM` Value from memory

Fetch Logic

- Predefined Blocks
 - PC: Register containing PC
 - Instruction memory: Read 10 bytes (PC to PC+9)
 - Signal invalid address
 - Split: Divide instruction byte into icode and ifun
 - Align: Get fields for rA, rB, and valC
- Control Logic
 - Need reigns: Does this instruction

Decode Logic

- Register File
 - Read ports A, B
 - Write ports E, M
 - Addresses are register IDs or 15 (0xF) (no access)
- Control Logic
 - srcA, srcB: read port addresses
 - dstE, dstM: write port addresses
- Signals
 - And: Indicate whether or not to perform conditional move
 - Computed in Execute stage

Execute Logic

- Units
 - ALU
 - Implements 4 required functions
 - Generates condition code values
 - CC
 - Register with 3 condition code bits
 - cond
 - Computes conditional
 - Jump/move flag
- Control Logic
 - Set CC: Should condition code register be loaded?
 - ALU A: Input A to ALU
 - ALU B: Input B to ALU
 - ALU fun: What function should ALU compute?

Memory Logic

- Memory
 - Reads and writes memory word
- Control Logic
 - Stat: What is instruction status?
 - Mem. Read: should word be read?
 - Mem. Write: should word be written?
 - Mem. Addr: select address
 - Mem. Data: select data

Writeback

- Register File
 - Write ports E, M
 - Addresses are register IDs or 15 (0xF) (no access)
- Control Logic
 - dstE, dstM: write port addresses
- Signals
 - Cnd: Indicate whether or not to perform conditional move
 - Computed in Execute stage

PC Update Logic

- New PC
 - Select next value of PC

Stage Computation: Arith/Log. Ops

	$OPq\ rA, rB$	
Fetch	$icode:ifun \leftarrow M_1[PC]$ $rA:rB \leftarrow M_1[PC+1]$ <p style="color:red; margin-left: 100px;"> ↙ means how many bytes M is. ↑ pointer to first byte of instruction </p> $valP \leftarrow PC+2$	Read instruction byte Read register byte Compute next PC
Decode	$valA \leftarrow R[rA]$ $valB \leftarrow R[rB]$	Read operand A Read operand B
Execute	$valE \leftarrow valB\ OP\ valA$ Set CC	Perform ALU operation Set condition code register
Memory		
Write back	$R[rB] \leftarrow valE$	Write back result
PC update	$PC \leftarrow valP$	Update PC

- Formulate instruction execution as sequence of simple steps
- Use same general form for all instructions

Stage Computation: rrmovq

	rrmovq rA, D(rB)	
Fetch	$\text{icode:ifun} \leftarrow M_1[\text{PC}]$ $\text{rA:rB} \leftarrow M_1[\text{PC+1}]$ $\text{valC} \leftarrow M_8[\text{PC+2}]$ $\text{valP} \leftarrow \text{PC+10}$	Read instruction byte Read register byte Read displacement D Compute next PC
Decode	$\text{valA} \leftarrow R[\text{rA}]$ $\text{valB} \leftarrow R[\text{rB}]$	Read operand A Read operand B
Execute	$\text{valE} \leftarrow \text{valB} + \text{valC}$	Compute effective address
Memory	$M_8[\text{valE}] \leftarrow \text{valA}$	Write value to memory
Write back		
PC update	$\text{PC} \leftarrow \text{valP}$	Update PC

■ Use ALU for address computation

Stage Computation: Jumps

	jXX Dest	
Fetch	$\text{icode:ifun} \leftarrow M_1[\text{PC}]$ $\text{valC} \leftarrow M_8[\text{PC+1}]$ $\text{valP} \leftarrow \text{PC+9}$	Read instruction byte Read destination address Fall through address
Decode		
Execute	$\text{Cnd} \leftarrow \text{Cond}(\text{CC}, \text{ifun})$	Take branch?
Memory		
Write back		
PC update	$\text{PC} \leftarrow \text{Cnd ? valC : valP}$	Update PC

- Compute both addresses
- Choose based on setting of condition codes and branch condition

Stage Computation: call

	call Dest	
Fetch	$\text{icode:ifun} \leftarrow M_1[\text{PC}]$ $\text{valC} \leftarrow M_8[\text{PC}+1]$ $\text{valP} \leftarrow \text{PC}+9$	Read instruction byte Read destination address Compute return point
Decode	$\text{valB} \leftarrow R[\%rsp]$	Read stack pointer
Execute	$\text{valE} \leftarrow \text{valB} + -8$	Decrement stack pointer
Memory	$M_8[\text{valE}] \leftarrow \text{valP}$	Write return value on stack
Write back	$R[\%rsp] \leftarrow \text{valE}$	Update stack pointer
PC update	$\text{PC} \leftarrow \text{valC}$	Set PC to destination

- Use ALU to decrement stack pointer
- Store incremented PC

Stage Computation: ret

	ret	
Fetch	$\text{icode:ifun} \leftarrow M_1[\text{PC}]$	Read instruction byte
Decode	$\text{valA} \leftarrow R[\%rsp]$ $\text{valB} \leftarrow R[\%rsp]$	Read operand stack pointer Read operand stack pointer
Execute	$\text{valE} \leftarrow \text{valB} + 8$	Increment stack pointer
Memory	$\text{valM} \leftarrow M_8[\text{valA}]$	Read return address
Write back	$R[\%rsp] \leftarrow \text{valE}$	Update stack pointer
PC update	$\text{PC} \leftarrow \text{valM}$	Set PC to return address

- Use ALU to increment stack pointer
- Read return address from memory

SEQ Summary

- Implementation
 - Express every instruction as series of simple steps
 - Follow same general flow for each instruction type
 - Assembly registers, memories, predesigned combinational blocks
 - Connect with control logic
- Limitations
 - Too slow to be practical
 - In one cycle, must propagate through instruction memory, register file, ALU, and data memory
 - Would need to run clock very slowly
 - Hardware units only active for fraction of clock cycle

Real-World Pipelines: Car Washes....

Sequential



Parallel



Pipelined



Idea

- **Divide process into independent stages**
- **Move objects through stages in sequence**
- **At any given times, multiple objects being processed**

Assignment 5 - Y86-64

```
1 # Author: Alyssa Kelley
2 # Assignment 5
3 # Source: I worked on this assignment with Miguel Hernan, Anne Glickenhouse, and Kristine Stecker.
4
5 # Execution begins at address 0
6     .pos 0
7     irmovq stack, %rsp      # This sets up the stack pointer
8     call main              # This executes the main program
9     halt                   # This terminates the program
10
11 # Array of 10 elements
12     .align 8
13 array:
14     .quad 0x0000000000000007
15     .quad 0x0000000000000001
16     .quad 0x0000000000000002
17     .quad 0x0000000000000003
18     .quad 0x0000000000000008
19     .quad 0x000000000000000A
20     .quad 0x0000000000000009
21     .quad 0x0000000000000006
22     .quad 0x0000000000000004
23     .quad 0x0000000000000005
24
25 main:
26     irmovq array, %rdi
27     irmovq $10, %rsi      # $10 is the length of the array stored in rsi
28     call selectionSort    # selectionSort(array)
29     ret
30
31 swap:
32     mrmmovq    (%rdi), %rax
33     mrmmovq    (%rsi), %rdx
34     rmmovq    %rdx, (%rdi)
35     rmmovq    %rax, (%rsi)
36     ret
37
38
39 selectionSort:
40
41     irmovq $1, %r12
42     rrmmovq %rsi, %r11
43     subq %r12, %r11      # This looks like the c code: len--
44     irmovq $0, %rcx      # This is setting i = 0 for the c code.
45     jmp L1
46
47
```

```

48 L2:
49     # Inner loop of the c code for selectionSort.
50
51     # C code: Get array[j]
52     rrmovq %rdx, %r14
53     addq %r14, %r14
54     addq %r14, %r14
55     addq %r14, %r14      # Adding 3 times = multiplied by 8 for the sizeof(long)
56     addq %rdi, %r14      # Offset from &array
57
58     # C code: Get array[min_idx]
59     rrmovq %rbx, %r8
60     addq %r8, %r8
61     addq %r8, %r8
62     addq %r8, %r8      # Adding 3 times = multiplied by 8
63     addq %rdi, %r8      # Offset from &array
64
65     # Dereference
66     mrmovq (%r14), %r14
67     mrmovq (%r8), %r8
68
69     pushq %r14
70     subq %r8, %r14      # C code: if(array[j] < array[min_idx])
71     popq %r14
72     jge L3
73     rrmovq %rdx, %rbx      # C code: min_idx = j
74
75 L3:
76     addq %r12, %rdx      # C code: j++
77     pushq %rdx
78     subq %rsi, %rdx
79     popq %rdx
80     jl L2
81
82     pushq %rcx
83     subq %rbx, %rcx      # C code: if(i==min_idx)
84     popq %rcx
85     je L4
86
87     # This is preparing for swap
88     pushq %rdi
89     pushq %rsi
90     rrmovq %rdi, %rsi
91
92     # This gets array[min_idx_idx]
93     rrmovq %rbx, %r13
94     addq %r13, %r13
95     addq %r13, %r13
96     addq %r13, %r13
97     addq %r13, %rdi
98
99     # This gets array[i]
100    rrmovq %rcx, %r13
101    addq %r13, %r13
102    addq %r13, %r13
103    addq %r13, %r13
104    addq %r13, %rsi
105
106    # C code: swap(array[min_idx],array[i])
107    call swap
108    popq %rsi
109    popq %rdi
110

```

```

111     L4:
112         addq %r12, %rcx      # C code: i++
113         pushq %rcx
114         subq %r11, %rcx      # C code: if(i < len)
115         popq %rcx
116         jl L1
117
118     L1:
119         # This is the outer loop for selectionSort
120         rrmovq %rcx, %rbx # min_idx = i
121         rrmovq %rcx, %rdx # j = i
122         addq %r12, %rdx    # j++
123         jmp L2
124     ret
125
126 # Stack starts here and grows to lower addresses
127 .pos 0x220
128 stack:

```

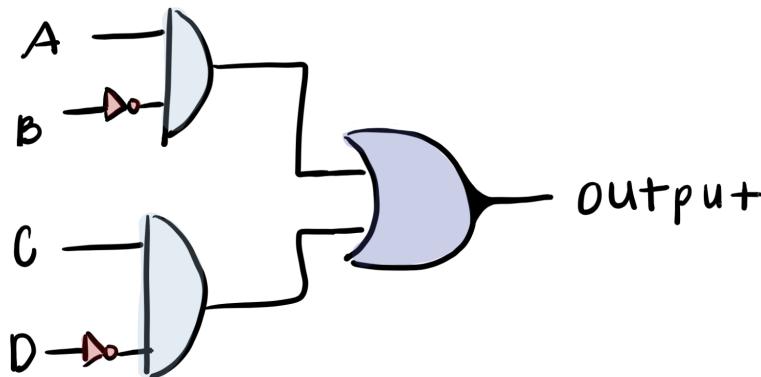
2. Draw a circuit with inputs A, B, C and D

$$(A \oplus !B) \text{II} (C \oplus !D)$$

NOT

AND

OR



3. Draw a diagram expressing how the 64-bit architecture could implement a new `iaddq` instruction to accomplish this functionality:

	<code>iaddq V, rB</code>	
Fetch	<code>iCode:ifun</code>	$M_1[PC]$
	<code>rA:rB</code>	$M_1[PC+1]$
	<code>val C</code>	$M_8[PC+2]$
	<code>Val P</code>	$PC + 10$
Decode	<code>val B</code>	$R[rB]$
Execute	$Val E \leftarrow Val C + Val B$	$Set cc$
Memory		
Write back	$R[rB] \leftarrow Val E$	
PC Update	$PC \leftarrow Val P$	

Lecture 5: Optimizing Program Performance - Limitation of Optimizing Compilers

Performance Realities

- There is more to performance than asymptotic complexity
- Constant factors matter too
 - Easily see 10 : 1 performance range depending on how code is written
 - Must optimize at multiple levels:
 - Algorithm, data representations, procedures and loops
- Must understand system to optimize performance
 - How programs are compiled and executed
 - How modern processors + memory systems operate
 - How to measure program performance and identify bottlenecks
 - How to improve performance without destroying case modularity and generality

Optimizing Compilers

- Provide efficient mapping of program to machine
 - Register allocation
 - Code selection and ordering (scheduling)
 - Dead code elimination
 - Eliminating minor inefficiencies
- Don't (usually) improve 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
- Have difficulty overcoming "optimization blockers"
 - Potential memory aliasing
 - Potential procedure side-effects

Limitations of Optimizing Compilers

- Operate under fundamental constraint
 - Must not cause any change in program behavior
 - Except, possibly when program making use of nonstandard language features
 - Often prevents it from making optimizations that would only affect behavior under pathological conditions
- Behavior that may be obvious to the programmer can be obfuscated by languages and coding styles
 - 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
 - Newer version of GCC do inter procedural analysis within individual files
 - But, not between code in different files
- Most analysis is based only on *static* information
 - Compiler has difficulty anticipating run-time inputs
- When in doubt, the compiler must be conservative

Generally Useful Optimizations

- Optimizations that you or the compiler should do regardless of processor / compiler
- Code Motion
 - Reduce frequency with which computation performed
 - If it will always produce same result
 - Especially moving code out of loop

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

Reduction in Strength

- Replace costly operation with simpler one
- Shift, add instead of multiply or divide
 - $16 * x \rightarrow x << 4$
 - Utility machine dependent
 - Depends on cost of multiply or divide instruction
- Recognize sequence of products

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

Share Common Subexpressions

- Reuse portions of expressions
- GCC will do this with -O1

Optimization Blocker #1: Procedure Calls

- Procedure to Convert String to Lower Case

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

Lower Case Conversion Performance

- Time quadruples when double string length
- Quadratic performance

Calling Strlen

- Strlen performance
 - Only way to determine length of string is to scan its entire length, looking for null character
- Overall performance, string of length N
 - N calls to strlen
 - Require times N, N-1, N-2, ... 1
 - Overall O(N^2) performance

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

Improving Performance

- Move call to strlen outside of loop
- Since result does not change from one iteration to another
- Form of code motion

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

Optimization Blocker: Procedure Calls

- Why couldn't compiler move strlen out of inner loop?
 - Procedure may have side effects
 - Alters global state each time called
 - Function may not return same value for given arguments
 - Depends on other parts of global state
 - Procedure lower could interact with strlen
- Warning:
 - Compiler treats procedure call as a black box
 - Weak optimizations near them
- Remedies
 - Use of inline functions
 - GCC does this with -O1
 - With single file
 - Do you own code motion

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

Exploiting Instruction-Level Parallelism

- Need general understanding of modern processor design
 - Hardware can execute multiple instructions in parallel
- Performance limited by data dependencies
- Simple transformations can yield dramatic performance improvement
 - Compilers often cannot make these transformations
 - Lack of associativity and distributivity in floating-point arithmetic

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
- Eliminated sources of overhead in loop

Loop Unrolling (2 x 1)

```
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

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 throughout limitations of execution units
 - overhead for short lengths
 - Finish off iterations sequentially

Assignment 6 - Optimization

```
1 // Author: Alyssa Kelley
2 // Assignment 6 - Part 1
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <time.h>
6
7 #define N 800
8
9 typedef int array_t[N][N];
10
11 int dim()
12 {
13     return N;
14 }
15 void old_f(array_t a, int x, int y)
16 {
17     for (int i = 0; i < dim(); ++i)
18     {
19         for (int j = 0; j < dim(); ++j)
20         {
21             a[i][j] = i * x * y + j; // Change this part!
22         }
23     }
24 }
25 /*
26 Task: Rewrite the above procedure f to minimize unnecessary function calls and multiplications. Note that you are
not required to replace the a[i][j] lookups with pointer arithmetic; focus on
27 the i * x * y + j equation with respect to removing unnecessary multiplications. Also write a main() function to
test your procedure. Name your source file 6-1.c.
28 */
29
30 void _f(array_t a, int x, int y)
31 {
32     int *beginning = &a[0][0]; // starting the array at the beginning location
33     int len = dim(); /* Creating the length as a variable so it does not have to waste that time
continuing to reference another function (dim()) */
34
35     int mul_xy = x*y; // Pulling this multiplication out of loop composition to further optimize
36     int mul_xyi;
37     int i = 0; // initializing here so it is less costly in the for loops
38     int j = 0;
39
40     for (i = 0; i < len; i++)
41     {
42         mul_xyi = mul_xy * i; // Pulling the multiplication outside of the inner for loop.
43         for (j = 0; j < len; j++)
44         {
45             //pointer = mul_xyi + j; // a[i][j] = i*x*y + j
46             a[i][j] = mul_xyi + j; // Only adding here because it is less costly.
47         }
48     }
49 }
```

```

51 void print_array(array_t a) // This is a function just to print the array in an ideal format.
52 {
53     printf("{ ");
54     for (int i = 0; i < N; ++i)
55     {
56         printf("{{");
57
58         for (int j = 0; j < N; ++j)
59         {
60             if (j < N-1) // Only printing a comma after the number if it is not the last one in the sub arrays.
61             {
62                 printf("%i, ", a[i][j]);
63             }
64             else
65             {
66                 printf("%i ", a[i][j]);
67             }
68         }
69
70         if (i < N-1)
71         {
72             printf("}, "); // Only printing a comma after each inner array if it is not the last array within the
73             // big array.
74         }
75         else
76         {
77             printf("}\n");
78         }
79     }
80     printf("}\n");
81 }
82 . . .
83 int main()
84 {
85     clock_t start, end;
86     double total_time;
87     array_t array_1;
88     array_t array_2;
89
90     int x = 2;
91     int y = 2;
92
93     start = clock();
94     for(int i=0; i<1000; i++)
95     {
96         old_f(array_1, x, y);
97     }
98 //old_f(a, x, y);
99     end = clock();
100    total_time = ((double)(end - start)) / CLOCKS_PER_SEC;
101    printf("This is the unoptimized run time for array a: %f\n", total_time); // you can divide here by 1000 to get
102    // an average
103    //print_array(a);
104
105    // Optimized time for array a
106    start = clock();
107    //_f(a, x, y);
108    for(int i=0; i<1000; i++)
109    {
110        _f(array_1, x, y);
111    }
112    end = clock();
113    //print_array(a);
114    total_time = ((double)(end - start)) / CLOCKS_PER_SEC;
115    printf("This is the optimized run time for array a: %f\n", total_time);
116    //print_array(a);

```

```
117     start = clock();
118     // old_f(b, x, y);
119     for(int i=0; i<1000; i++)
120     {
121         old_f(array_2, x, y);
122     }
123     end = clock();
124     total_time = ((double)(end - start)) / CLOCKS_PER_SEC;
125     //print_array(b);
126     printf("This is the unoptimized run time for array b: %f\n", total_time);
127
128     start = clock();
129     //_f(b, x, y);
130     for(int i=0; i<1000; i++)
131     {
132         _f(array_2, x, y);
133     }
134     end = clock();
135     total_time = ((double)(end - start)) / CLOCKS_PER_SEC;
136     //print_array(b);
137     printf("This is the optimaized run time: %f\n", total_time);
138
139     return 0;
140 }
```

```

1 // Author: Alyssa Kelley
2 // Assignment 6 - Part 2c
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <time.h>
6
7 #define N 20000 // Length of array
8
9 /*
10 # u in %rbx, v in %rax, length in %rcx, i in %rdx, sum in %xmm1
11 .L87:
12 movss (%rbx, %rdx, 4), %xmm0 # Get u[i]
13 mulss (%rax, %rdx, 4), %xmm0
14 adds %xmm0, %xmm1
15 addq $1, %rdx
16 cmpq %rcx, %rdx
17 j1 .L87
18 # Multiply by v[i] #Addtosum
19 # Increment i
20 # Compare i to length # If <, keep looping)
21 */
22
23 void inner(float *u, float *v, int length, float *dest)
24 { // This is the unoptimized function.
25     int i;
26     float sum = 0.0f;
27     for (i = 0; i < length; ++i)
28     {
29         sum += u[i] * v[i];
30     }
31     *dest = sum;
32 }
33
34 void inner2(float *u, float *v, int length, float *dest)
35 {
36     // four-way loop unrolling with four parallel accumulators (which is the contents of the for loop)
37     int i; // Initialize as much as possible before the for loop to get the most optimization.
38     float sum1 = 0.0f; // Setting them all to start at 0.
39     float sum2 = 0.0f;
40     float sum3 = 0.0f;
41     float sum4 = 0.0f;
42     int limit = length - 3; // Offsetting this to subtract 3 to reflect the 4 we are looking at at one time.
43
44     for(i = 0; i < limit; i += 4) // pull out 4 at a time as just mentioned above.
45     {
46         // implementing four-way loop unrolling. See Lecture 5b Slide # 14/18
47         sum1 = sum1 + u[i] * v[i]; // Creating a sum for each 4 we pull out.
48         sum2 = sum2 + u[i+1] * v[i+1];
49         sum3 = sum3 + u[i+2] * v[i+2];
50         sum4 = sum4 + u[i+3] * v[i+3];
51     }
52     // This is cleaning up the array, Source: Lecture 5b, Slide 18.
53     for(; i < length; i++)
54     {
55         sum1 = sum1 + (u[i] * v[i]);
56     }
57     *dest = sum1 + sum2 + sum3 + sum4; // destination is used to create a new array.
58 }
59

```

```

60 int main(int argc, char *argv[])
61 {
62     clock_t start1; /* Initializing the start and end times for both the inner
63     function (start1 / end1) and the optimized inner2 funciton (start2 / end2) */
64     clock_t start2;
65     clock_t end1;
66     clock_t end2;
67     float *u = malloc(sizeof(float) * N); // Mallocing
68     float *v = malloc(sizeof(float) * N);
69     float dest = 0;
70     int num_func_calls = 1000; // This is the amount of how many times we run the functions.
71
72     for (int i=0; i<N; i++)
73     {
74         u[i]=i; // This is filling the arrays as large as N which is a
75         v[i]=i; // global variable defined at the top.
76     }
77
78     start1 = clock();
79     for(int i=0; i<num_func_calls; i++) /* This is calling the function 1000 times (num_func_calls) so we can have a large, and
80                                         more accurate amount of data. You can then divide the total_time by 1000 to get
81                                         the average time it takes to run it once. */
82
83     {
84         inner(u, v, N, &dest);
85     }
86     end1 = clock();
87
88     double total_time1 = ((double)(end1 - start1)) / CLOCKS_PER_SEC; // This is calculating the time it took to execute the function.
89     printf("This is the unoptimized run time for array a: %f\n", total_time1/num_func_calls); // Average time for one funciton call
90     printf("This is the sum for inner: %f\n", dest);
91
92     start2 = clock();
93     for(int i=0; i<num_func_calls; i++)
94     {
95         inner2(u, v, N, &dest); // Optimized function.
96     }
97     end2 = clock();
98     double total_time2 = ((double)(end2 - start2)) / CLOCKS_PER_SEC;
99     printf("This is the unoptimized run time for array a: %f\n", total_time2/num_func_calls);
100    printf("This is the sum for inner2: %f\n", dest);
101
102    printf("This means this was %f times faster!\n", (total_time1-total_time2)*1000);
103
104    return 0;
105 }
106

```

The x86-64 assembly code for the inner loop is as follows:

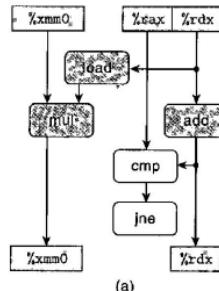
```
# u in %rbx, v in %rax, length in %rcx, i in %rdx, sum in %xmm1
.L87:
    movss (%rbx, %rdx, 4), %xmm0 # Get u[i]
    mulss (%rax, %rdx, 4), %xmm0 # Multiply by v[i]
    adds %xmm0, %xmm1 # Add to sum
    addq $1, %rdx # Increment i
    cmpq %rcx, %rdx # Compare i to length
    jl .L87 # If <, keep looping
```

Registers to include:

- rax, rbx
- rdx, rcx
- xmm0, xmm1

From the textbook:

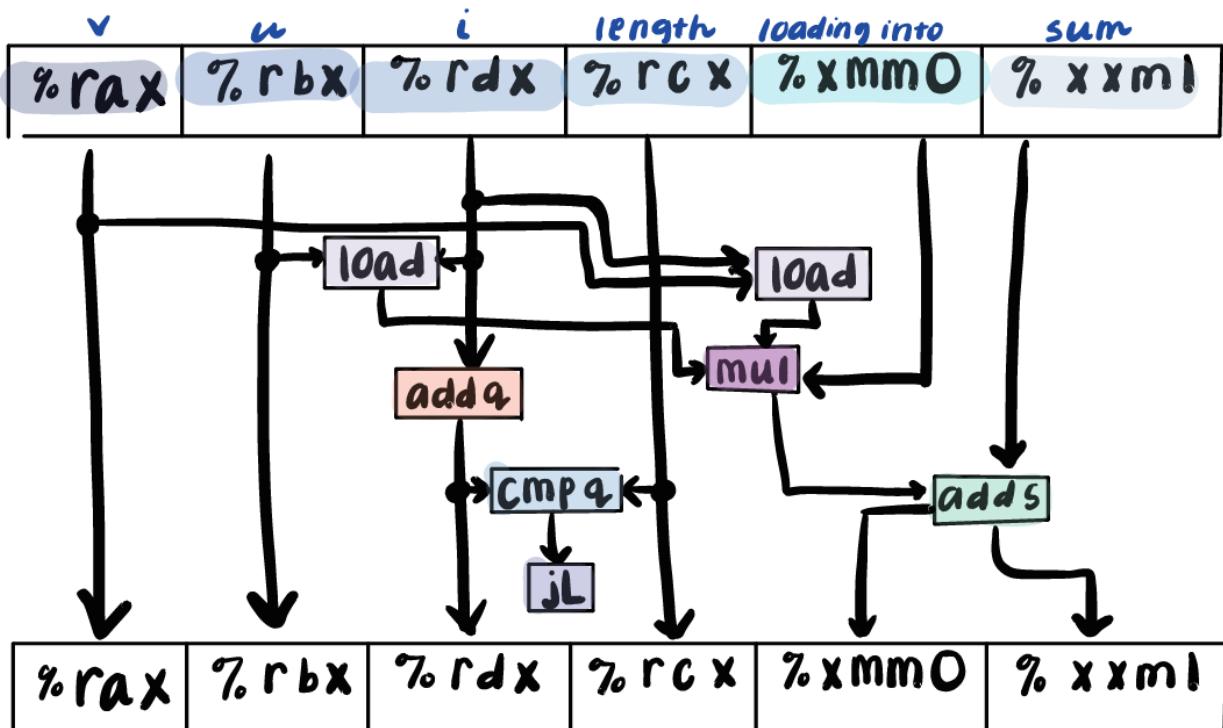
Figure 5.14
Abstracting combine4 operations as a data-flow graph. We rearrange the operators of Figure 5.13 to more clearly show the data dependencies (a), and then further show only those operations that use values from one iteration to produce new values for the next (b).



following this example.

Instructions:

load
addq, adds
cmpq



1)

2)

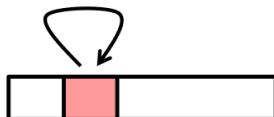
b. (20) Which operation(s) in the loop can NOT be pipelined? Why? What are the latencies of these operations? Based on this, what is the lower latency bound (in terms of CPE) of the procedure? Assume that *float* addition has a latency of 3, *float* multiplication has a latency of 5, and all integer operations have a latency of 1. Hint: think about which operation(s) depend on a result from the previous loop iteration. Write your answers in your solutions document.

- 1) The operations in the loop that cannot be pipelined are those that rely on the previous / returned values. An example of this in the scenario above would be the incrementer $i \rightarrow i + 1$, and the two add instructions use i to iterate through the loop and compare require i so these all can not be pipelined. To reiterate, the add instructions are what cannot be pipelined.
- 2) The two latencies for these operations are 1 for the integer latency (adds) and 3 for the float addition latency (addq).
- 3) The lower bound latency is 3.

Lecture 6: The Memory Hierarchy - Memory Technologies, Locality

Random-Access Memory (RAM)

- Key feature

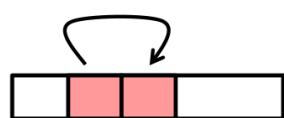


- RAM is traditionally packages as a chip
- Basic storage unit is normally a cell (one bit per cell)
- Multiple RAM chips from a memory
- RAM comes in two varieties
 - SRAM (Static RAM)

- DRAM (Dynamic RAM)

Locality

- Principle of Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently



- **Temporal locality:**

- Recently referenced items are likely to be referenced again in the near future

- **Spatial locality:**

- Items with nearby addresses tend to be referenced close together in time

Locality Example

- Data reference
 - Reference array elements in succession
 - Reference variable sum each iteration
- Instruction references
 - Reference instructions in sequence
 - Cycle through loop repeatedly

Spatial locality
Temporal locality

Spatial locality
Temporal locality

Qualitative Estimates of Locality

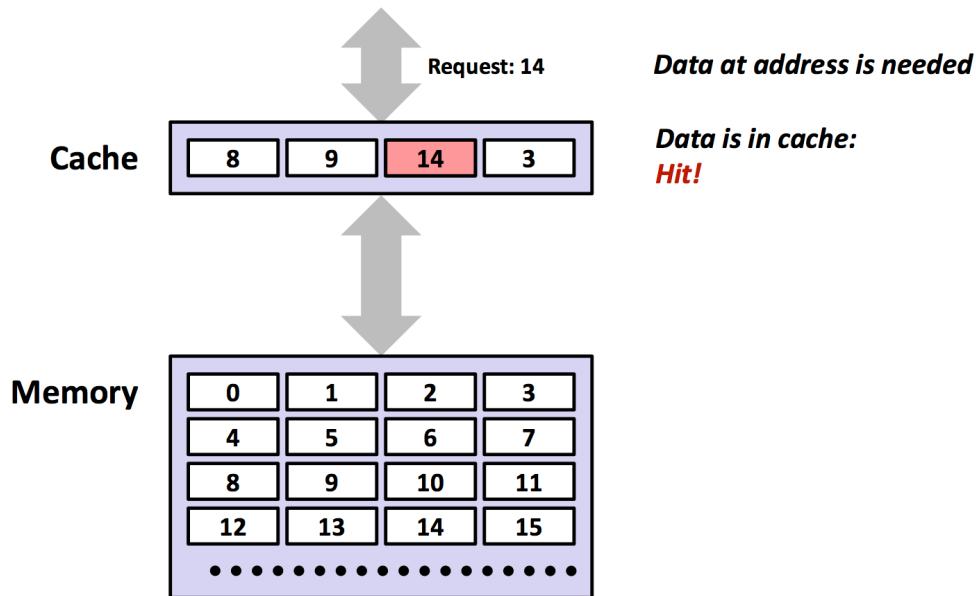
- Claim: Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer

Lecture 6b: The Memory Hierarchy - Caching in the Memory Hierarchy

Memory Hierarchies

- Some fundamental and enduring properties of hardware and software
 - The gap between CPU and main memory speed is widening
 - Fast storage technologies cost more per byte, have less capacity, and require more power (heat!)
 - 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 knowns as a **memory hierarchy**

General Cache Concepts: Hit



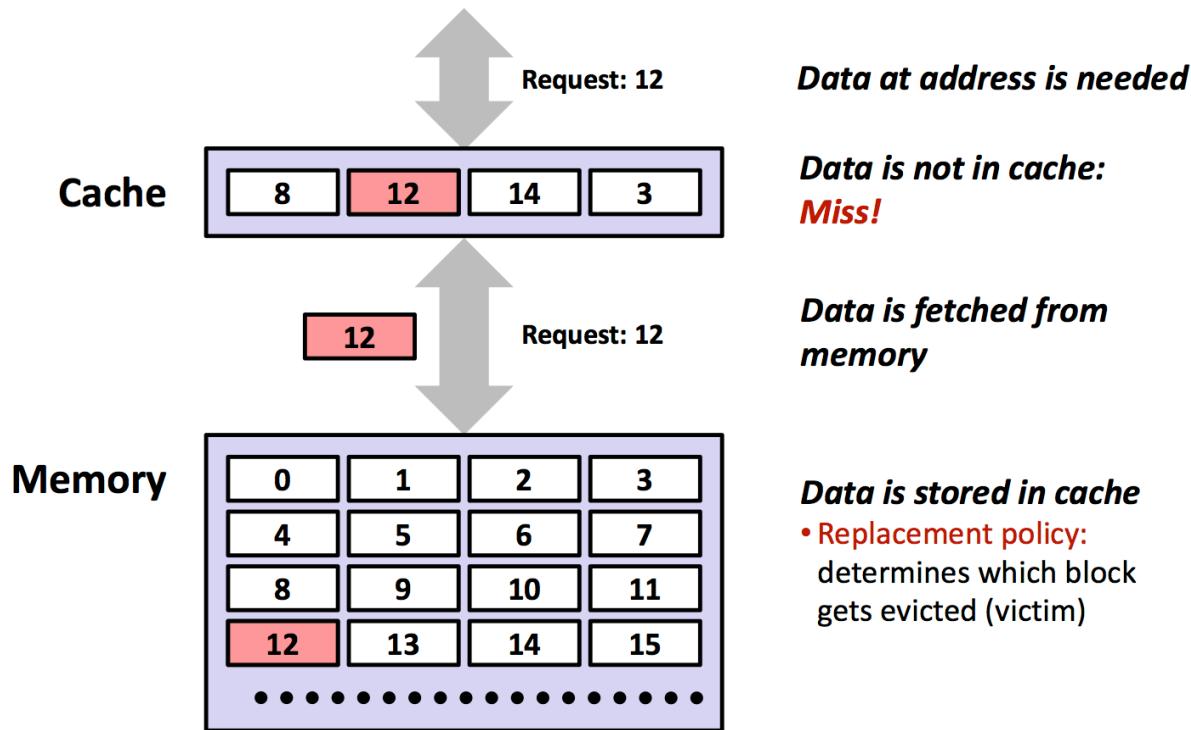
Cache

- Cache: A smaller, faster storage device that acts as a staging area for a subset of the data in a larger, slower device
- Fundamental idea of a memory hierarchy:
 - For each k , the faster, smaller device at level k serves as a cache for the larger, slower device at level $k + 1$
- Why do memory hierarchies work?
 - Because of locality, programs tend to access the data at level k more often than they access the data at level $k + 1$
 - Thus, the storage at level $k + 1$ can be slower, and thus larger and cheaper per bit

- Big idea: The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top

Hit!

General Cache Concepts: Miss



Miss!

	Tag	Set	Offset BLOCK	
1-way	3	2	3	direct map
2-way	4	1	3	2 way associative cache

	valid	tag
set 0	01	? 0000 *
	01	? 0010 *

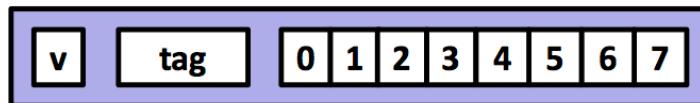
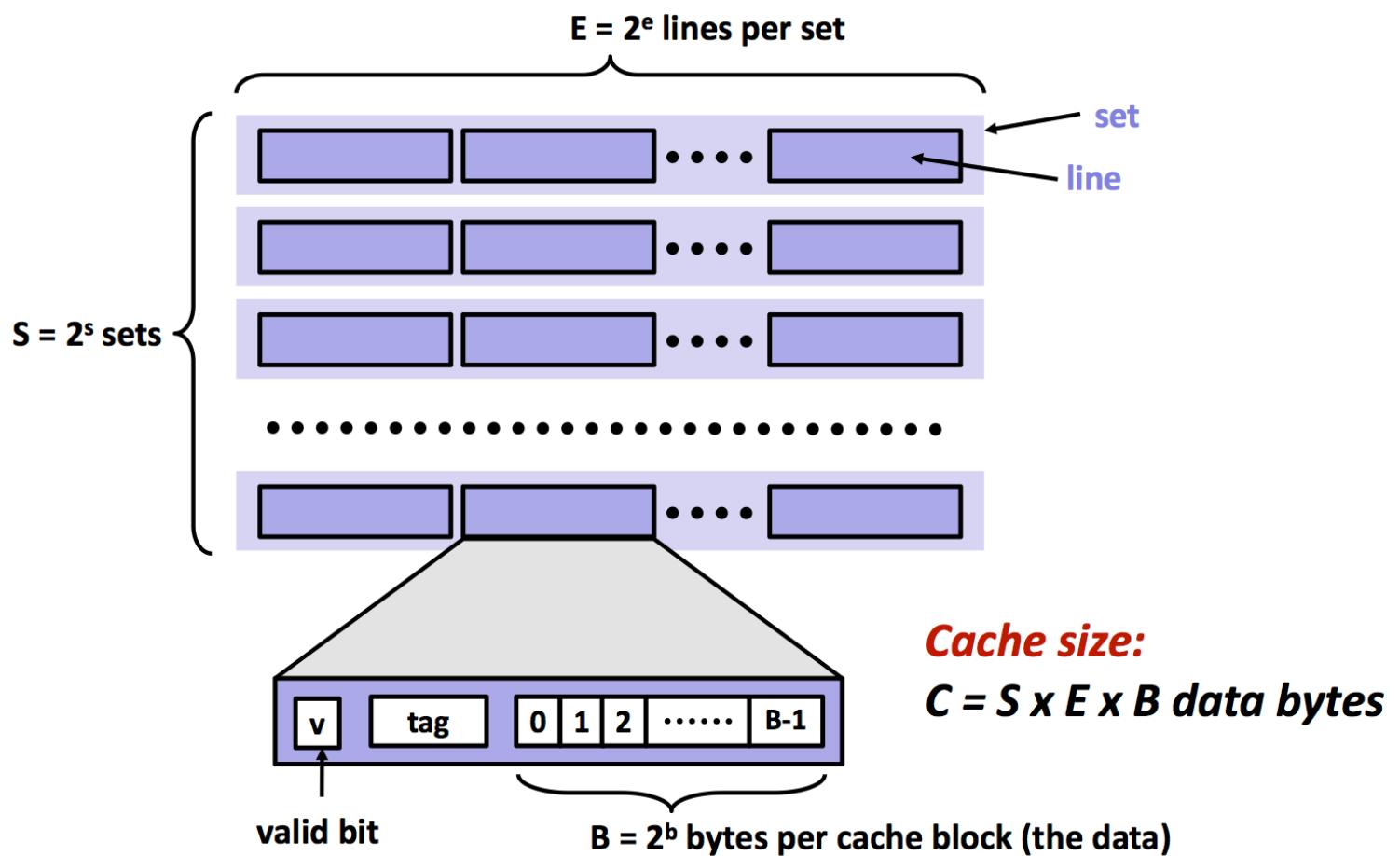
	valid	tag
set 1	01	? 0000
	0	?

addresses

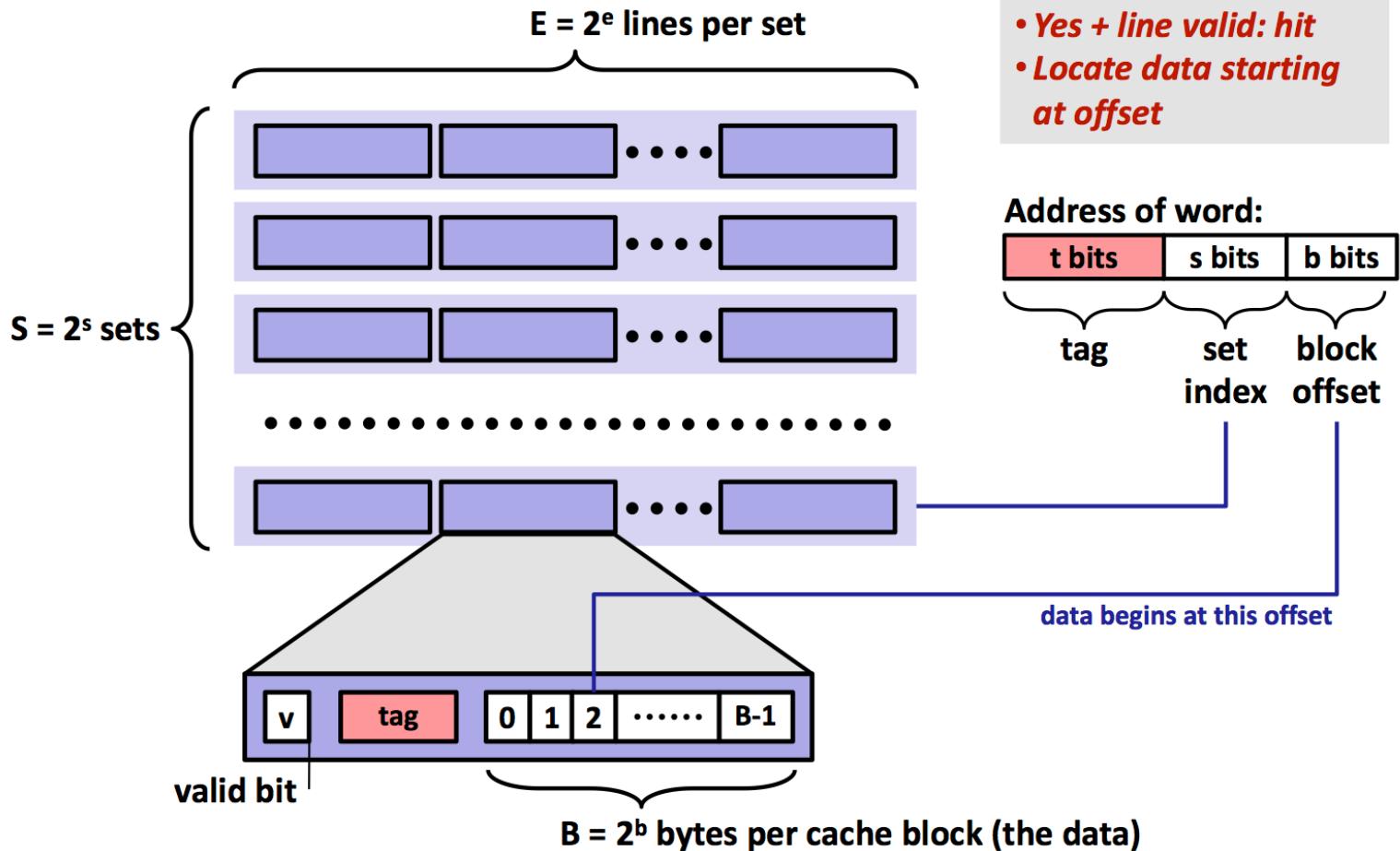
0x03	0000	0	010	miss (allways miss b/c nothing there)
0x04	0000	0	100	hit (don't load b/c it is already in there)
0x08	0000	1	000	miss
Example:	0010	0	110	miss
0x04	0000	0	100	hit
0x24	0010	0	100	hit

tag set offset

General Cache Organization (S, E, B)



Cache Read



- Every byte in memory belongs to a block:
 - Lowest order o "offset" bits specify the byte's position in its block
- Each block belongs to a set in the cache:
 - Next lowest order s "set" bits specify the cache set to which the block belongs
- Many many blocks in memory will belong to the same cache set, so we need a way to identify cached blocks:
 - Remaining t "tag" bits are used to determine if a cached block is the block containing the desired byte

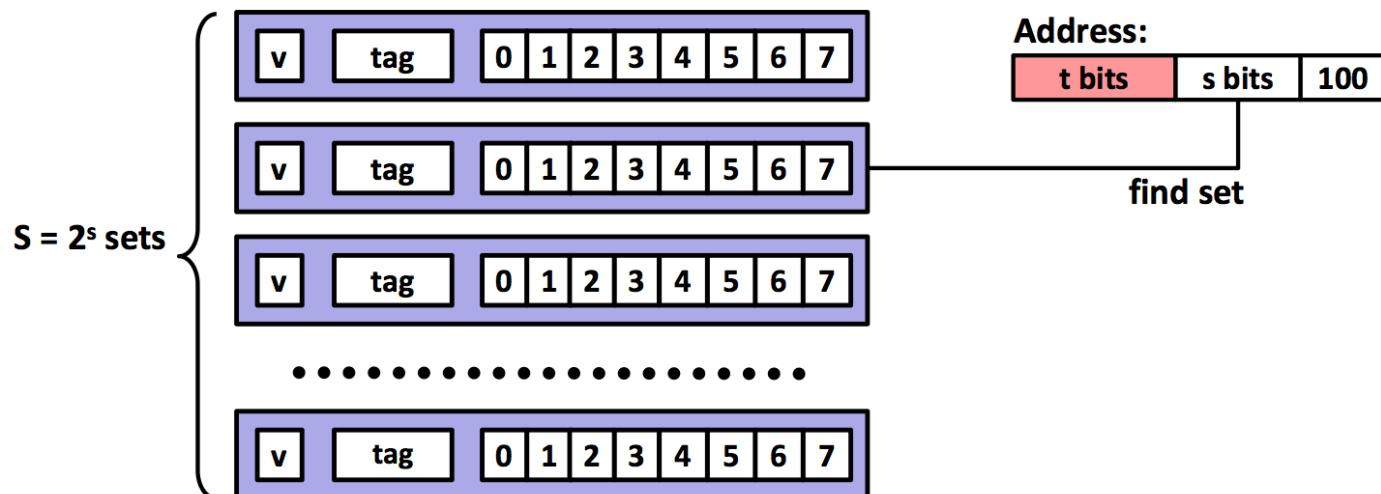
Cached Blocks

- Each cached block is represented in the cache by:
 - The bytes in the block
 - The block's identifying tag (the high order bits of all addresses in the block)
 - A valid bit (can't determine if the tag and block bytes were written intentionally by looking at the tag and block bytes alone)

Example: Direct Mapped Cache ($E = 1$)

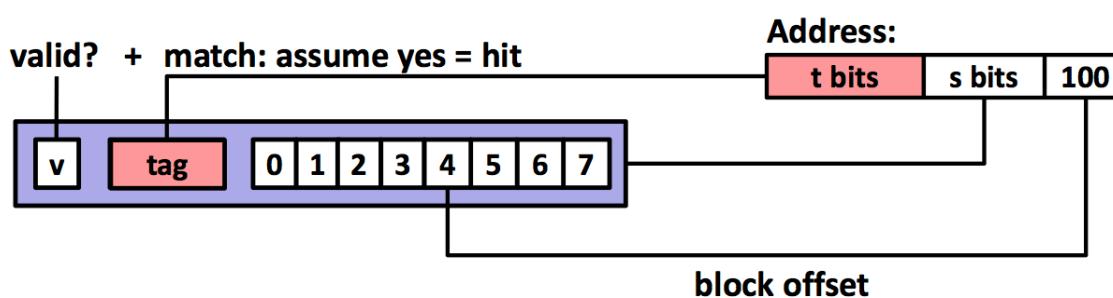
Direct mapped: One block per set

Assume: cache block size 8 bytes



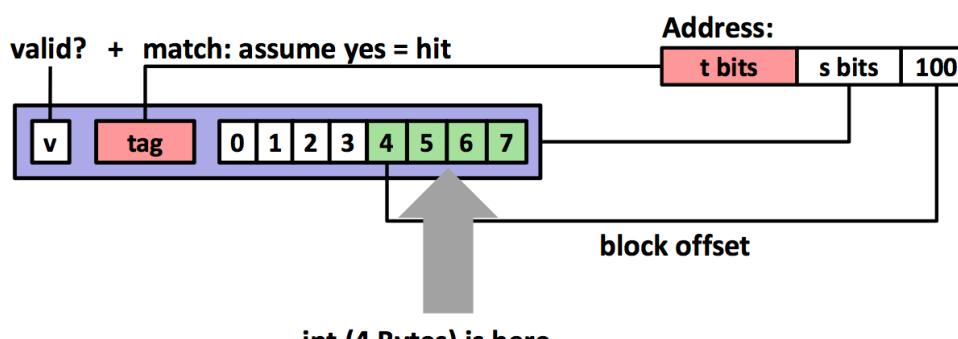
Direct mapped: One block per set

Assume: cache block size 8 bytes



Direct mapped: One block per set

Assume: cache block size 8 bytes



same tag, then it is a hit! If not, then it is a miss, and that missed information is stored in the cache set if there is room. If there is no room, it kicks out the oldest saved set and tag and replaces it with this recently missed one.

If no match: old line is evicted and replaced

--->
Basic
ally,
you
have
your
sets,
and in
those
sets,
you
will be
able
to
store
the
set
and
the
corres
pondi
ng
tag.

Then
you will
continu
e to
look at
the
next
addre
ses and
if that
saved
set has
the

Direct-Mapped Cache Simulation

t=1 s=2 b=1

x	xx	x
---	----	---

M=16 bytes, 4-bit addresses,
8-byte cache, B=2 bytes/block,
E=1 blocks/set, S=4 sets

This
will be
on the
final!!

Address trace (reads, one byte per read):

0	[0000 ₂],	miss
1	[0001 ₂],	hit
7	[0111 ₂],	miss
8	[1000 ₂],	miss
0	[0000 ₂]	miss

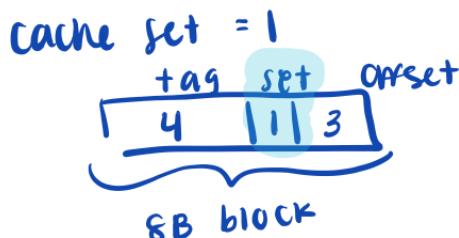
	v	Tag	Block
Set 0	1	0	M[0-1]
Set 1	0	?	?
Set 2	0	?	?
Set 3	1	0	M[6-7]

2. [25] Consider a 16B direct-mapped cache with 8B blocks and 2 sets for an 8-bit architecture (i.e., 256 bytes of memory):

c. (15) Consider the following sequence of memory accesses. For each address, show the tag, set, offset, and whether it resulted in hit or miss:

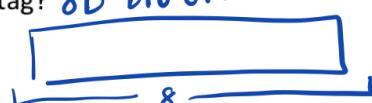
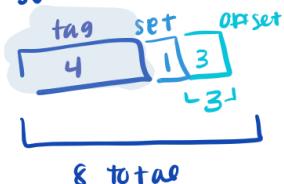
		tag	set	offset	hit / miss	
0000 0000	0x00	0000	0	000	miss	
0000 0100	0x04	0000	0	100	hit	
0000 1000	0x08	0000	1	000	miss	(nothing in set 1)
0000 0100	0x04	0000	0	100	hit	
0000 0000	0x00	0000	0	000	hit	

→ explanation for why hit / miss.



b. (5) How many address bits are used to represent the cache tag? 8B block

$$2^3 = 8 \\ \text{So this would be 3 bits for the offset}$$



Cache tag = 4

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <string.h>
4
5 // Author: Alyssa Kelley
6 // Assignment 7 - Part 2
7
8 /* Notes for 7-2:
9
10 0 x 1 2 3 4 5 6 7 8
11   |-----| |_-|_|_
12     tag      set offset
13 */
14
15 unsigned int getOffset(unsigned int address)
16 {
17     /* This is going to return the right most bit.
18      Example: 0x12345678 -> Offset = 8
19 */
20     unsigned int bitmask = 0xF; // Using the mask 1111 to get the last 4 bits to indicate the last number.
21     unsigned int offset = address & bitmask; // & these together will have 0 everything else except for the last 4 bits.
22     return offset; // Returning the right most 4 bits.
23 }
24
25 unsigned int getTag(unsigned int address)
26 {
27     /* This is going to return the left most 4*6 bits.
28      Example: 0x12345678 -> Tag = 123456
29 */
30     unsigned int bitmask = 0xFFFFFFF0; /* Using the bit mask 1111 1111 1111 1111 1111 1111 0000 0000
31     and the reason for this is to get only the first 6 integers and then the last two should be dropped. */
32     unsigned int set = (address & bitmask);
33     return set >> 8; // Shifting off 8 so get rid of the last 2 0'd out integers.
34 }
35
36 unsigned int getSet(unsigned int address)
37 {
38     unsigned int bitmask = 0x000000F0; /* Using the mask 0000 0000 0000 0000 0000 0000 1111 0000 and
39     this will get the second to last integer in the hex number. */
40     unsigned int tag = address & bitmask;
41     return tag >> 4; // This shifts off the last integer before the second to last
42 }
43
44 int main()
45 {
46     unsigned int address1 = 0x12345678;
47     unsigned int address2 = 0x87654321;
48
49     printf("0X%.8X: offest: %d - tag: %X - set: %X\n", address1, getOffset(address1), getTag(address1), getSet(address1));
50     printf("0X%.8X: offest: %d - tag: %X - set: %X\n", address2, getOffset(address2), getTag(address2), getSet(address2));
51
52     return 0;
53 }
54

```

Another
cache

example:

A. (15) What is the total number of memory reads? Why? It may help to think in terms of movq instructions.

Loop iteration - (for loop #1 = 16) * (for loop #2 = 16)

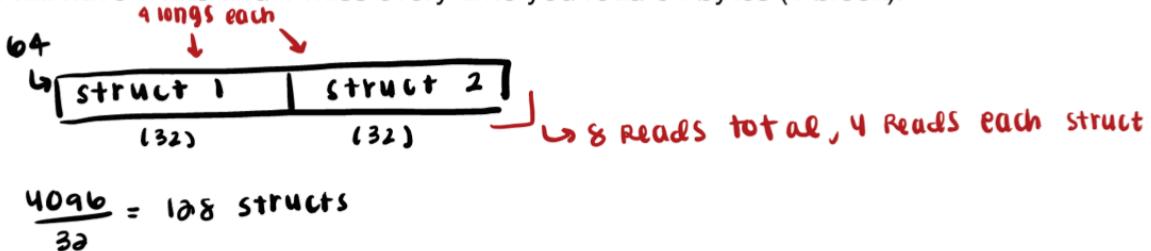
4 memory reads so $(4) * (16 * 16) = 1024$ total reads

Answer: 1024 total reads.

each instruction = 8 bytes and how many reads is how many times an instruction

B. (15) What is the total number of memory reads that miss in the cache? Why?

You will have 7 hits and 1 miss every time you read 64 bytes (1 block).



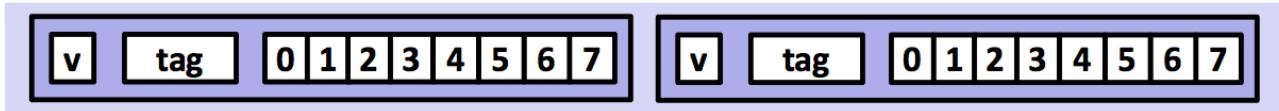
`sizeof(colorPoint) = 32 bytes`

Each block can hold 2 structs (because it is 64)

Explanation: Since we are reading 1024 total reads and every time we read 64 bytes then we have 1 miss so that would be a total of 16 misses and 112 hits per 1024 reads.

C. (10) What is the miss rate?

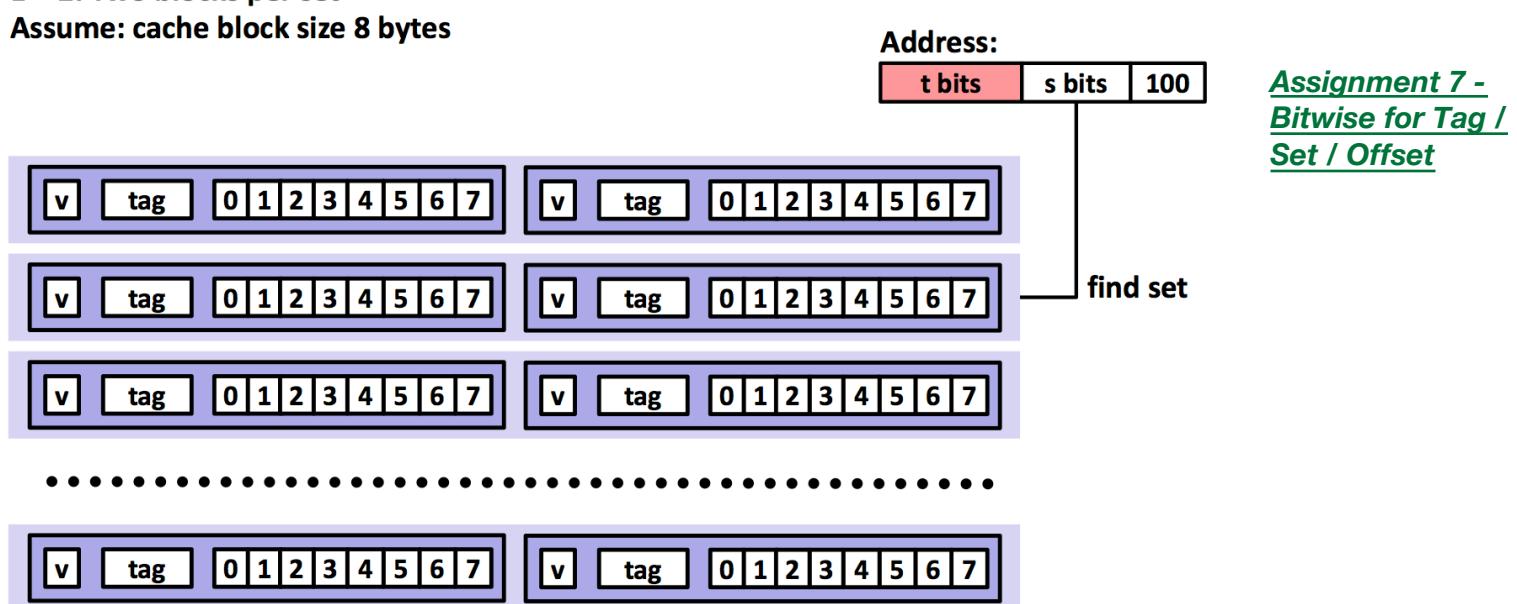
-> $16/128 = 1/8 = 12.5\%$ miss rate for 1024 total memory reads.



E-way Set Associative Cache (Here: E = 2)

E = 2: Two blocks per set

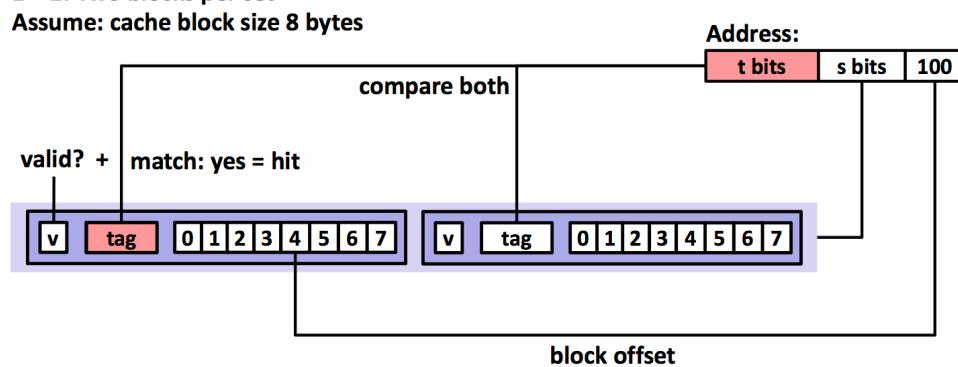
Assume: cache block size 8 bytes



E-way Set Associative Cache (Here: E = 2)

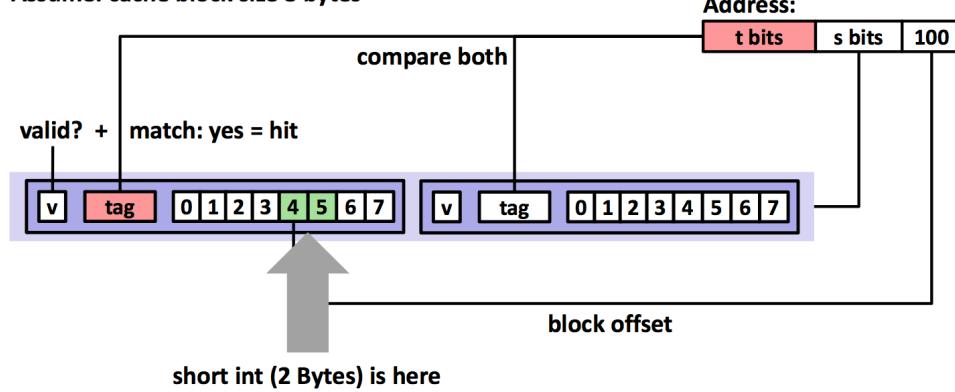
E = 2: Two blocks per set

Assume: cache block size 8 bytes



E-way Set Associative Cache (Here: E = 2)

E = 2: Two blocks per set
Assume: cache block size 8 bytes



[Lecture 6c: The Memory Hierarchy - Direct Mapped Caches, Set-Associative Caches](#)

If no match:

- One block in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

Improving Temporal Locality

2-Way Set Associative Cache Simulation

t=2 s=1 b=1

xx	x	x
----	---	---

M=16 bytes, 4-bit addresses,
8-byte cache, B=2 bytes/block,
E=2 blocks/set, S=2 sets

Address trace (reads, one byte per read):

0	[0000 ₂]	miss
1	[0001 ₂]	hit
7	[0111 ₂]	miss
8	[1000 ₂]	miss
0	[0000 ₂]	hit

	v	Tag	Block
Set 0	1	00	M[0-1]
	1	10	M[8-9]
Set 1	1	01	M[6-7]
	0	?	?

- Allow multiple blocks to be cached in each set:
 - Number of blocks in each set is called "set associativity"
 - Improving temporal locality by allowing a choice regarding which blocks get evicted
- But how do we efficiently search two tags and two valid bits in each set?
- And how do we decide which blocks get evicted?

```

1 // Author: Alyssa Kelley
2 // Assignment 8
3
4 // Note: Source code from Eric Wills, and I collaborated with
5 // Anne Glickenhaus and Kristine Stecker in CIS 314 as we worked on
6 // this project together.
7
8 #include <stdio.h>
9 #include <stdlib.h>
10
11 struct Line {
12     unsigned char data[4];
13     unsigned int tag;
14     unsigned char valid;
15 };
16
17 struct Cache {
18     struct Line *lines;
19     int numLines;
20 };
21
22 unsigned int getOffset(unsigned int address) {
23     // 4B blocks, so offset is bits 0-1
24     return address & 0x3;
25 }
26
27 unsigned int getSet(unsigned int address) {
28     // 16 sets, so offset is bits 2-6
29     return (address >> 2) & 0xF;
30 }
31
32 unsigned int getTag(unsigned int address) {
33     // Offset and set are 6 bits total, so tag is high-order 26 bits
34     return address >> 6;
35 }
36
37 struct Cache* mallocCache(int numLines) {
38     // TODO - malloc a pointer to a struct Cache, malloc a pointer to an array
39     // of struct Line instances (array length is numLines). Also initialize
40     // valid to 0 for each struct Line. Return the struct Cache pointer.
41     struct Cache* cache = (struct Cache *) malloc(sizeof(struct Cache)); // casting -> pointer
42     cache->numLines = numLines;
43     struct Line* lines = (struct Line *) malloc(sizeof(struct Line));
44     cache->lines = lines;
45     // to initialize cache so have everything set to 0 to begin with.
46     for(int i = 0; i< numLines; i++)
47     {
48         cache->lines[i].valid = 0;
49         cache->lines[i].tag = 0;
50         for(int j=0; j<4; j++)
51         {
52             cache->lines[i].data[j] = 0;
53         }
54     }
55     return cache;
56 }
57
58 void freeCache(struct Cache *cache) {
59     free(cache->lines);
60     free(cache);
61 }
62
63 void printCache(struct Cache *cache) {
64     // TODO - print all valid lines in the cache.
65     // This runs when you hit p.
66     for(int i = 0; i<cache->numLines; i++) // going through the lines in the cache
67     {
68         if(cache->lines[i].valid == 1) // making sure it is valid
69         {
70             printf("set: %d - tag: %d - valid: %d - data:", i, cache->lines[i].tag, cache->lines[i].valid);
71             for(int j=0; j<4; j++)
72             {
73                 printf(" %.2x", cache->lines[i].data[j]); // printing the cache data as two lower case hex values if it is valid.
74             }
75             printf("\n");
76         }
77     }
78 }
79

```

```

80 void readValue(struct Cache *cache, unsigned int address) {
81 // TODO - check the cache for a cached byte at the specified address.
82 // If found, indicate a hit and print the byte. If not found, indicate
83 // a miss due to either an invalid line (cold miss) or a tag mismatch
84 // (conflict miss).
85
86 /* Pass in cache and address and pull off set, tag, and offset from that
87 and the point is to go look for this value in the cache
88 Direct mapped so we are allowed to look in to get the line
89 directly with set on line 98.
90 */
91
92 // parse the address
93 unsigned int set = getSet(address);
94 unsigned int tag = getTag(address); // tag is like a key so you need to make sure the tags (aka keys) match up for it to be a match -- think of this
95     list hash tables
96 unsigned int offset = getOffset(address);
97
98 // Block we need
99 printf("looking for set: %d - tag: %d\n", set, tag);
100 struct Line* index = &cache->lines[set]; // Direct map so we go look in here and get the lines directly and then it points to it. No need for malloc
101     here.
102
103 // Making sure it is valid
104 if(index->valid == 0)
105 {
106     printf("no valid set found - miss! \n");
107 }
108 else if(index->tag == tag) // It is valid, and the tags match... it is a hit!
109 {
110     printf("found set: %d - tag: %d - valid: %u - data: ", set, index->tag, offset, index->valid); // this is the print block
111     printf(" %.2x", index->data[offset]); // printing the offset with two lower case hex values.
112     printf("\n");
113     printf("hit!\n");
114 }
115 else // it is valid, but it does not match - it is a miss.
116 {
117     // print block
118     printf("found set: %d - tag: %d offset: %u - valid: %d - data: ", set, index->tag, offset, index->valid);
119     // print data
120     printf(" %.2x", index->data[offset]);
121     printf("\n");
122     printf("tags don't match - miss!\n");
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 int main() {
154     struct Cache *cache = mallocCache(16);
155
156     // Loop until user enters 'q'
157     char c;
158     do {
159         printf("Enter 'r' for read, 'w' for write, 'p' to print, 'q' to quit: ");
160         scanf(" %c", &c);
161
162         if (c == 'r') {
163             printf("Enter 32-bit unsigned hex address: ");
164
165             unsigned int a;
166             scanf(" %x", &a);
167
168             readValue(cache, a);
169         } else if (c == 'w') {
170             printf("Enter 32-bit unsigned hex address: ");
171
172             unsigned int a;
173             scanf(" %x", &a);
174
175             printf("Enter 32-bit unsigned hex value: ");
176
177             unsigned int v;
178             scanf(" %x", &v);
179
180             // Get byte pointer to v
181             unsigned char *data = (unsigned char *)&v;
182
183             writeValue(cache, a, data);
184         } else if (c == 'p') {
185             printCache(cache);
186         }
187     } while (c != 'q');
188
189     freeCache(cache);
190 }
191
```

Concluding Observations

- Programmer can optimize for cache performance
 - How data structures are organized
 - How data are accessed
 - Nested loop structure
- All systems favor "cache friendly code"
 - Getting absolute optimum performance is very platform specific
 - Cache sizes, block size, associativities, etc.
 - Can get most of the advantage with generic code
 - Keep working set reasonably small (temporal locality)
 - Use small strides (spatial locality)

Assignment 8 - Cache Hits and Misses

