

Number Representation

2's Complement hex to decimal

Ex. $0x52 = 0b1010010 \rightarrow 0b10101101 + 1 = 0b10101110$
 $\rightarrow 128 - (32 + 8 + 4 + 2) = 82$

Two's complement addition

$(-2) + (-4) = (-6) \rightarrow 1110 + 1100 = 11010$
 \rightarrow leading 1 means negative \rightarrow flip bits = 0101 \rightarrow add 1 = 0110

Counting

Unsigned binary \rightarrow Count = 2^n

- Range of unsigned representation is $[0, 2^n - 1]$
- Negative: $[-2^{n-1} - 1, -1]$ Positive: $[0, 2^{n-1} - 1]$
- Together: $[-2^{n-1} - 1, 2^{n-1} - 1]$

One's complement

- Zero 00...0 and 11...1 = ± 0
- Negative: $[-2^{n-1} - 1, 0]$ Positive: $[0, 2^{n-1} - 1]$
- Together: $[-2^{n-1} - 1, 2^{n-1} - 1]$
- Like one's complement, but shift negative numbers by 1
 \rightarrow get rid of 2 zeros
- Zero 00...0
- Negative: $[-2^{n-1}, 0]$, most negative: 100...0
- Positive: $[0, 2^{n-1}]$, most positive: 011...1
- Together: $[-2^{n-1}, 2^{n-1} - 1]$

Two's complement overflow will never occur when adding two numbers of opposite signs.

Bias N = $-(2^{n-1} - 1)$

- Non-biased encoding \rightarrow biased encoding (Subtract the bias)
- Biased encoding \rightarrow non-biased encoding (Add the bias)
 - Zero 011...1
 - Most negative: 00...0 \rightarrow $-(2^{n-1} - 1)$, Most positive: 11...1 \rightarrow 2^{n-1}
- Range of numbers can be represented $[B, 2^n - 1] \rightarrow$ it's like unsigned notation $[0, 2^n - 1]$ + bias
- IEC standard prefixes

• 2XY Kibi - Ki - 2^{10} Mebi - Mi - 2^{20} Gibi - Gi - 2^{30} Tebi - Ti - 2^{40} Pebi - Pi - 2^{50} Exbi - Ei - 2^{60} Zebi - Zi - 2^{70} Yobi - Yi - 2^{80}	• X tells us the prefix Y tells us the number that precedes the prefix (ex: 2 Ti) X = 0 => no prefix X = 1 => kibi X = 2 => mebi X = 3 => gibi X = 4 => tebi X = 5 => gibi X = 6 => exbi X = 7 => zebi X = 8 => yobi	Y = 0 => 1 Y = 1 => 2 Y = 2 => 4 Y = 3 => 8 Y = 4 => 16 Y = 5 => 32 Y = 6 => 64 Y = 7 => 128 Y = 8 => 256
---	--	---

- | | |
|---|--|
| • Write the following in IEC notation: 275
• X = 7 => zebi
• Y = 5 => 32
• Answer: 32 Zi | • Write the following as a power of 2: 64 Ti
• Ti => X = 4
• Y = log(64) = 6
• Answer: 2 ⁶ |
|---|--|

C Intro - Basics

C Pre-Processor (CPP) (begin with #) \rightarrow #define is string replacement

- C source files first pass through CPP before compiler sees code: Foo.c
 \rightarrow CPP \rightarrow foo.i \rightarrow compiler

what evaluates to FALSE in C? \rightarrow 0 and NULL

C Intro - Pointers, Arrays, Strings

- `sizeof(void) = word size = 4 bytes`
- * \rightarrow value of the pointer, & \rightarrow address of a value
- Pointers to functions \rightarrow int (*fn) (void *, void *) = &fn,
 - (*fn) (x, y) call the function
- Null pointers (0s) \rightarrow Cannot read or write a null pointer, cannot dereference (*p) a null pointer
- Pointers and Structures \rightarrow Passed by reference
- `sizeof()` operator \rightarrow Not a function is a static compile-time operation
- If we want a function to change a pointer:
- Segmentation Fault (`SIGSEGV`): invalid access to a valid memory
- Bus Error (`SIGBUS`): access to an invalid address

Copy a string

- a = malloc(sizeof(char) * (strlen(b) + 1));
`strcpy(a, b)` or `strncpy(a, b, strlen(b) + 1)`; \rightarrow strncpy safe, but doesn't copy null terminator

what would `strlen(a)` return?

- `char a = 'f'` \rightarrow N/A (not a string)
- `char a[] = {f, b, b, b, b, b}` \rightarrow N/A (no null terminator)
- `char* a = "foobar"` \rightarrow 3

Notes

- `(*ip)++ without ()` is `*(ip++)` (access ip before increment vs. increment before access)
- `ptr = &a[0]` is the same as `ptr = a`
- `*(&a + i) = a[i]` (next element in array)
- `&a[i] = a + i`
- only print null terminated string
- `int** mat = (int**)calloc(n, sizeof(int*));`
`for (int i = 0; i < n; i++) {mat[i] = (int*)calloc(m, sizeof(int))};`

char a[5] = {a, b, c, d, e}

- a = &a
- a+1 \rightarrow increment one char over
- &a+1 \rightarrow increment 5 chars over

```
int x = 5;
int *y = &x;
y = y + 2;
char* z = "bears!";
z = z + 2;
```

incrementing sizeof(pointer_type)!

Decimal	Binary(0b)	Hex(0x)	Two's Complement Number Line	One's Complement Number Line	Sign and Magnitude Number Line
0	0000	0	100 101 110 111 000 001 010 011	100 101 110 111 000 001 010 011	111 110 101 100 000 001 010 011
1	0001	1	-4 -3 -2 -1 0 1 2 3	-3 -2 -1 0 +0 1 2 3	-3 -2 -1 0 +0 1 2 3
2	0010	2			
3	0011	3			
4	0100	4			
5	0101	5			
6	0110	6			
7	0111	7			
8	1000	8			
9	1001	9			
10	1010	A			
11	1011	B			
12	1100	C			
13	1101	D			
14	1110	E			
15	1111	F			

Memory Management

- Stack Return address, arguments, space for local variables; Stack pointer indicates start of stack frame
- Heap space requested for dynamic data via malloc() or calloc(), resizes dynamically grows upward
- Static data global variables and constants declared outside functions;
 - Doesn't grow or shrink, loaded when program starts, can be modified.
- Code loaded when program starts, does not change
- 0xxxxx 00000 is reserved and unwriteable/unreadable \rightarrow crush on null pointer access

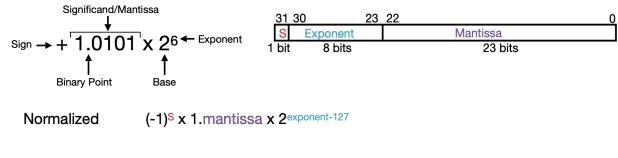
Managing the Heap

- void *malloc (size_t, n) - O(1)
- void *calloc (size_t nmem, size_t size) - O(n)
- void free (void *p)
- void *realloc (void *p, size_t size)
 - If p is NULL, then behaves like malloc
 - If size is 0, then behaves like free \rightarrow realloc (p, 0) otherwise it "increases the size"

Bitwise operations

- a >> b If a is signed sign extend, if unsigned zero extend (Not the same as dividing by 2^b due to rounding)
- x & 1 = x & 0 = 0, x & 1 = 1, x & 0 = x; x ^ 1 = ~x; x ^ 0 = x
- Turning bits ON \rightarrow 1; turning bits OFF \rightarrow 0; flipping bits \rightarrow ^

IEEE 754 Floating-Point Standard



Normalized $(-1)^S \times 1.\text{mantissa} \times 2^{\text{exponent}-127}$

Denormalized $(-1)^S \times 0.\text{mantissa} \times 2^{-126}$

Exponent = 0 and we need to shift the binary point over by 1 to get an implicit leading 0

Exponent (biased notation)

- Range of exponent [-126, 127]
- Bias formula $N = -(2^{n-1} - 1) \rightarrow$ 8 exponent bits
- bias = $-(2^{(8-1)-1}) = -127$

Range of floating-point values

- Smallest positive number we can represent

- “Caller” calls a function
 - Saves “caller” saved registers
 - “t” registers (if needed)
 - “a” registers (if needed)
 - ra (if needed)
 - Passes arguments through “a” registers
 - Can expect “s” registers to be saved
- “Callee” is called by another function
 - Handles epilogue and prologue
 - Decrements the stack
 - Saves callee-saved registers (“s” registers)
 - Restores register values before returning and increments the stack
 - Returns value (if any) through “a” registers

R-format

- Opcode (7) + funct7 (7) + funct3 (3), rs1 (5), rs2 (5), rd (5)
- Opcode 0110011

I-format

- imm[11o] holds values in range [-2048 (base 10), +2047]
- Always sign-extended to 32-bits before use

$$\begin{array}{l} a = b + c \Rightarrow add a, b, c \\ *x = y \Rightarrow sw y \text{ or } \\ \text{int } *x = ly \Rightarrow lw x, y \end{array}$$

Load Byte

- Load Instructions**
- Load Byte
 - Load Byte Unsigned (0-extended)
 - Load Halfword (lh)
 - Load Halfword Unsigned (lhu)
 - lw x10 (rd), 12 (offset) (x15) (rs)

S-format

- 2 sources reg no destination, No sign extension
- Store Halfword (sh)
 - sh x14, 36 (x5) → store the lower 16 bits of x14 into memory at address [x5]+36, no sign extension

B-Format {comparison? {reg1? {reg2? {label?}}

- Range of instructions: 2's complement range: $[-2^{n-1}, 2^{n-1} - 1]$
 - with 12 bits $\pm 2^{11}$ bytes away from the PC
 - Instructions are 4-bytes, so can jump $\pm 2^9$ instructions away from the current instruction

Branch Instruction Addressing

- Encoded immediate = 0boooo 0ooo 0011
- Actual immediate offset = 0booo 0ooo 0ooo 1100
- 3 instructions (12 bytes) away
- Can jump $\pm 2^{11}$ instructions ($\pm 2^{13}$ bytes) away from current PC

nx16-bit Instructions

- las bit is always 0b0 → discard this 1 bit instead of 2
- Range of bytes we can jump to $\pm 2^{12}$
- Range of 32-bit instructions we can jump to $\pm 2^{10}$

J-format jal rd, label

- 20 bits for immediate
- Discard last bit
- Can jump by $\pm 2^{19}$ bytes, $\pm 2^{18}$ 32-bit instructions

JALR is an I-Type Instruction

- jal rd, rs, imm → PC = [rs] + imm
- Usually used to return from a function (ret)
- imm[11o] holds range [-2048, +2047]
- Immediate always sign-extended to 32-bits before use
- Unlike JAL, must include the last 0 because we are using I-format

PC Relative Address

- PC → PC + offset
- Does the value in branch immediate field change if we change the location of the code in memory?
 - If moving all of code, then no (PC-relative offsets)

U-Format**Load Upper Immediate (LUI)**

- Destination register = immediate $\ll 12$
- lui x5, 0xAFBCDE → x5 = 0xAFBCDE000
- LUI + ADDI to create long immediates
 - lui x10, 0x87654 → # x10 = 0x87654000
 - addi x10, 0xE000 → x10 = 0x87654E000
- Load immediate (LI) → pseudo-instruction
 - li x10, 0x87654321
- Corner case How to set 0xABCDDEEE?

$$\begin{array}{rcl} \text{lui } x10, \text{ABCDE } \# x10 & = & 0xABCDDEEE \\ \text{addi } x10, \text{0EEE } \# x10 & = & 0xABCDDEEE \\ & & \downarrow \\ & & \begin{array}{c} \text{0xABCDDE000} \\ \text{0xFFFFFFFEE} \\ \text{0xABCDDEEE} \\ \text{+0xABCDDE} \\ \text{+0x0000} \\ \hline \text{0xABCDDEEE} \end{array} \\ & & \downarrow \\ & & \begin{array}{c} \text{0xABCDDE} \\ \text{0xFFFFFFF} \\ \text{0xABCDDEEE} \\ \text{+0xE000} \\ \hline \text{0xABCDDEEE} \end{array} \end{array}$$

Add Upper Immediate PC (AUIPC)

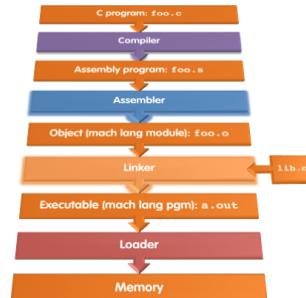
- rd = PC + (immediate $\ll 12$)
- auipc x5, 0xAFBCDE → x5 = PC + 0xAFBCDE000
- # Call function at any 32-bit absolute address

$$\begin{array}{l} \text{lui } x5, \langle \text{hi20bits} \rangle \\ \text{jalr ra, x5, \langle \text{lo12bits} \rangle} \end{array}$$

Jump PC-relative with 32-bit offset

$$\begin{array}{l} \text{auipc x5, \langle \text{hi20bits} \rangle} \\ \text{jalr ra, x5, \langle \text{lo12bits} \rangle} \end{array}$$
Compiler, Assembler, Linker, Loader (CALL)**CALL Chain**

- Compiler converts a single HLL file into a single assembly language file.
- Assembler removes pseudo-instructions, converts what it can to machine language, and creates a checklist for the linker (relocation table). A .s file becomes a .o file.
 - Does 2 passes to resolve addresses, handling internal forward references
- Linker combines several .o files and resolves absolute addresses.
 - Enables separate compilation, libraries that need not be compiled, and resolves remaining addresses
- Loader loads executable into memory and begins execution.

**Interpretation vs. Compilation**

- **Interpreter** directly executes a program in the source language
 - Easier to write (closer to high-level)
 - Slower
 - Provides instruction set independence
- **Translator** converts a program from the source language to an equivalent program in another language
 - Both INTERPRETER & COMPILER are translators
 - More efficient

Compiler

- output contains labels, pseudo-instructions, fully relative addressing
- 1) Lexer, 2) Parser, 3) Semantic Analysis and Optimization, 4) Code Generation

Assembler **RELATIVE ADDRESSING**

- A dumb compiler for assembly language
- Read and uses Directives
- .text: Subsequent items put in user text segment (machine code)
- .data: Subsequent items put in user data segment (binary rep of data in source file)
- .globl sym: Declares sym global and can be referenced from other files
- .string str: Store the string str in memory and null-terminate it
- .word w1..wn: Store the n 32-bit quantities in successive memory words

- Replace pseudo-instructions

Pseudo	Real
nop	addi x0, x0, 0
not rd, rs	xori rd, rs, -1
beq rs, rs, offset	beq rs, x0, offset
bgt rs, rt, offset	bit rt, rs, offset
j offset	jal x0, offset
ret	jalr x0, x1, offset
call offset (if too big for just a jal)	auipc x6, offset[31:12] jalr x1, x6, offset[11:0]
tail offset (if too far for a j)	auipc x6, offset[31:12] jalr x0, x6, offset[11:0]

Producing Machine Language

Simple Case

- Arithmetic, logical, shifts ... → convert into binary representations
- Branches → PC-Relative
 - "Forward Reference" problem → 2 passes
 - 1st pass: remember position of labels (as expanding pseudo-instructions)
 - 2nd pass uses label positions to generate code

Jumps (j and jal)

- Jumps within a file are PC relative (we can compute)
- Jumps to other files, we can't compute

Reference to Static Data

- la (load address) = lui (when absolute address) or auipc (when relocatable block) + addi
- Can't be determined immediately so create two tables

Symbol Table (non-static labels in present file)

- List of items in the file that may be used by other files,
 - Labels function calling
 - Data anything in the .data section; variables which may be accessed across files

Relocation Table (external, static reference)

- List of items that file needs to address of later
 - Any external label jumped to jal (external, including lib files)
 - Any piece of data in static section (such as la instruction)

Object File Format

- **object file header:** size and position of the other pieces of the object file
- **text segment:** the machine code
- **data segment:** binary representation of the static data in the source file
- **relocation information:** identifies lines of code that need to be fixed up later
- **symbol table:** list of this file's labels and static data that can be referenced
- **debugging information**

Linker ABSOLUTE ADDRESSING

- sum(*.o) → *.exe
- Enable separate compilation of files
- Process:
 1. Take text segments from each .o file and put them together
 2. Take data segments from each .o file, put them together, and concatenate them onto the end of text segments
 3. Resolve references (go through relocation table, handle each entry → fill in all absolute addresses like printf)

Three Types of Addresses

- **PC-Relative Addressing** (beq, bne, jal) → never relocate
- **External Function Reference** (usually jal) → always relocate
- **Static Data Reference** (often auipc, addi) → always relocate

Absolute Addresses in RISC-V

- which instructions need relocation editing?
 - Jump and link (jal) → ONLY for external
 - Loads and stores (lw, sw) to variables in static area, relative to the global pointer
 - Conditional branches → don't need to worry about since PC-relative addressing preserved even if code moves.

Resolving References

- Linker assumes first word of first text segment is at address 0x04000000
- Liner knows:
 - length of each text and data segments
 - ordering of text and data segments
- Liner calculates:
 - Absolute address of each label to be jumped to and each piece of data being referenced
- To resolve references:
 - Search for reference (data or label) in all "user" symbol tables → if not found, search library files → once absolute address is determined, fill in the machine code appropriately
- output: executable file containing **text** and **data** and **header**

Loader

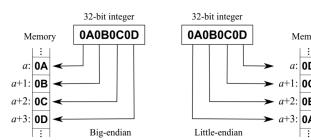
- Reads executable file's header to determine size of text and data segments
- Creates new address space for program large enough to hold text, data, and stack segments
- Copies instructions and data from executable file into new address space
- Copies arguments passed to the program onto the stack
- Initializes machine registers
 - Stack pointer assigned address of 1st free stack location
- Jumps to start-up routine that copies program's argument from stack to registers and sets the PC
- Executable files are stored on disk
- **Linker is actually 2 parts: linker in the compiler, linker in the loader**

Static vs. Dynamically Linked Libraries

Statically-linked Approach

- Library is part of the executable, so if library updates, we don't get the new update.
- Includes the entire library even if not all of it will be used.
- Executable is self-contained

Dynamically Linked Libraries (DLL)



SDS

Digital → discrete 1, 0

Analog → continuous

Logic Gate	Symbol	Description	Boolean
AND		Output is at logic 1 when and only when all its inputs are at logic 1, otherwise the output is at logic 0	X = A·B
OR		Output is at logic 1 when one or more are at logic 1. If all inputs are at logic 0, output is at logic 0	X = A+B
NAND		Output is at logic 0 when and only when all its inputs are at logic 1, otherwise the output is at logic 1	X = Ā·B̄
NOR		Output is at logic 0 when one or more of its inputs are at logic 1. If all the inputs are at logic 0, the output is at logic 1	X = Ā+B̄
XOR		Output is at logic 1 when one and Only one of its inputs is at logic 1. Otherwise it is logic 0	X = A⊕B
XNOR		Output is at logic 0 when one and only one of its inputs is at logic 1. Otherwise it is logic 1. Similar to XOR but inverted	X = Ā⊕B̄
NOT		Output is at logic 0 when its only input is at logic 1, and at logic 1 when its only input is at logic 0. That's why it is called an INVERTER	X = Ā

Logic Gates

$$\text{AND} \rightarrow \text{out} = A \cdot B$$

$$\text{OR} \rightarrow \text{out} = A + B; x \cdot 1 = 1, x + 0 = x$$

$$\text{XOR} \rightarrow \text{out} = A \wedge B = A\bar{B} + \bar{A}B$$

$$\text{NOT} \rightarrow \sim A$$

$$\text{NAND} \rightarrow \sim(A \wedge B) = \bar{A} \cdot \bar{B} + \bar{A}B + A\bar{B}$$

$$\text{NOR} \rightarrow \sim(A \vee B)$$

$$\text{XNOR} \rightarrow \sim(A \wedge B) = (\bar{A} \cdot \bar{B}) + A\bar{B}$$

Make XOR with AND and OR

A	B	out
0	0	0
0	1	1
1	0	1
1	1	0

$$\text{out} = \bar{A}\bar{B} + AB$$

A	B	out
0	0	1
0	1	0
1	0	0
1	1	1

Make XNOR

$$\text{out} = \bar{A}\bar{B} + AB$$

$$(A + B)(A + C) = A + BC$$

$$A + AB = A + B$$

$$A + B + \bar{A}C = AC$$

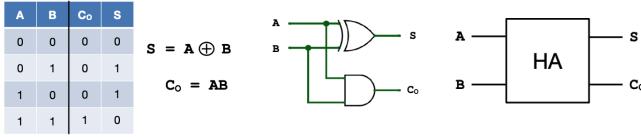
Laws of Boolean Algebra

Name	AND Form	OR form
Commutative	$AB = BA$	$A + B = B + A$
Associative	$AB(C) = A(BC)$	$A + (B + C) = (A + B) + C$
Identity	$1A = A$	$0 + A = A$
Null	$0A = 0$	$1 + A = 1$
Absorption	$A(A + B) = A$	$A + AB = A$
Distributive	$(A + B)(A + C) = A + BC$	$A(B + C) = AB + AC$
Idempotent	$A(A) = A$	$A + A = A$
Inverse	$A(\bar{A}) = 0$	$A + \bar{A} = 1$
De Morgan's	$\bar{AB} = \bar{A} + \bar{B}$	$\bar{A} + \bar{B} = \bar{A}(\bar{B})$

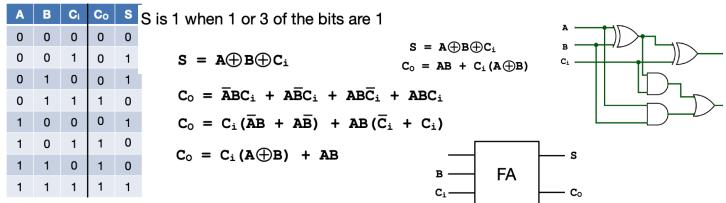
How to Build an Adder

Half Adder

- How to add two bits together using logic gates?

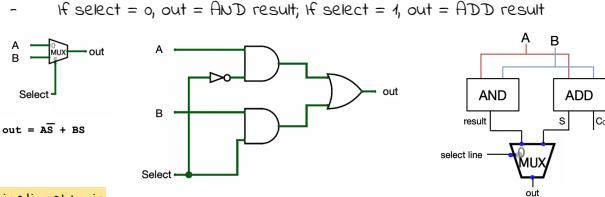


Full Adder (single bit)



Arithmetic Logic Unit (ALU)

2:1 Multiplexors



Combinational Logic

- As soon as the inputs are available, the output starts being computed.
- output depends only on the current input

D Flip Flops (state element)

- Rising edge triggered → stores D to Q the instant the clock goes from 0 to 1
- Falling edge triggered → stores D to Q the instant the clock goes from 1 to 0

Write enable

- when WE = 0, contents of register stays the same on the clock trigger
- when WE = 1, contents of register updates on the clock trigger

Combinational Logic Delay

- The amount of time that it takes for a value to propagate through the combinational logic

Maximum allowable hold time

- clk-to-q (reg) + shortest CL time
- The amount of time that it takes for the input to B to change after the trigger

Minimum acceptable clock cycle (critical path)

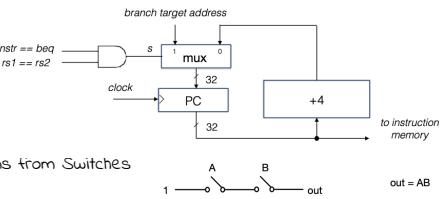
- clk-to-q (reg) + longest CL time + setup time (reg)
- Min clock cycle = The time it takes for the input of one state element to reach the input of the next state element.
- Minimum clock period (ns) = critical path
- Maximum clock frequency (Hz) = $1 / (\text{min. clk perio})$

*** Hold time < CLK-to-Q delay + min CL delay

*** CLK Period ≥ CLK-to-Q delay + critical path + setup time

*** clk-to-q ≥ hold

PC Register



Transistors

Boolean Expressions from Switches

$$(A + B) = \overline{\overline{AB}}$$

$$(A + B) = \overline{AB}$$

$$AB + CD = (\overline{AB})(\overline{CD})$$

DeMorgan's Law

$$(A + B) = \overline{AB}$$

$$AB + CD = (\overline{AB})(\overline{CD})$$

NMOS → Gate = 1, switch closed, gate = 0, switch open

- Not good at passing 1s, so hook source to 0

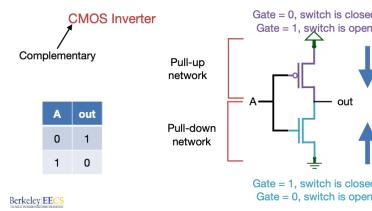
PMOS → Gate = 0, switch closed, gate = 1, switch open

- Not good at passing 0s, so hook source to 1

Vdd → supply voltage (Logic 1)

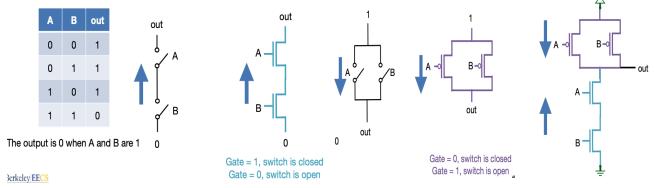
Ground → Logic 0

Building an Inverter

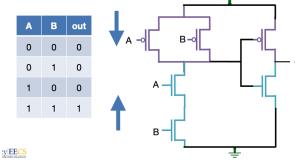


Building a NAND gate

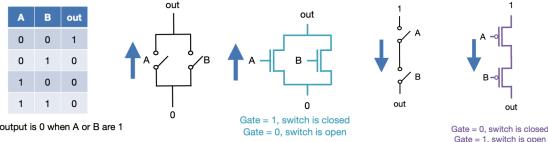
- output is 0 when A and B are 1, output is 1 when A and B are 0



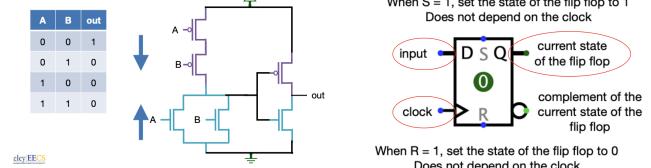
Building an AND gate → Invert a NAND gate



Building a NOR gate



Building an OR gate → Invert a NOR gate



Transistors

Inverter

- 2 transistors

NAND gate

- 4 transistors

NOR gate

- 4 transistors

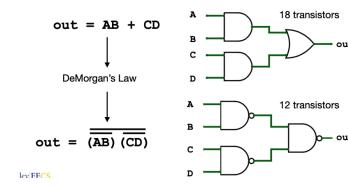
AND gate

- 6 transistors

OR gate

- 6 transistors

Converting ANDs and ORs to NANDs and NORs



Datapath Phases:

- IF: send address to IMEM, read IMEM at addr
 - ID: generate control signals from instruction, immediates, read registers from RegFile
 - EX: ALU operations, branch comparisons
 - MEM: read/write to DMEM
 - WB: write back either PC+4, result of ALU, or data from memory to RegFile
- $Lw = \text{clk-to-q}(PC) + t\text{MEMRead} + t\text{ImmGen} + t\text{Mux} + t\text{ALU} + t\text{Mux} + t\text{RegFileSetup}$
- $Sw = \text{clk-to-q} + \text{memread} + \max(\text{regfileread}, \text{immgen}) + \text{mux} + \text{alu} + \text{memwrite}$
- $\text{Branches} = \text{clk-to-q}(PC) + t\text{MEMRead} + \max(t\text{RegFileRead}, t\text{ImmGen}) + \max(t\text{BranchComp}, t\text{Mux} + t\text{ALU}) + t\text{PCSetup}$
- $JAL = \text{clk-to-q}(PC) + t\text{MEMRead} + t\text{ImmGen} + t\text{Mux} + t\text{ALU} + \max((t\text{Mux}(WBSel)) + t\text{RegFilSetup}), (t\text{Mux}(PCSel)) + t\text{PCSetup})$

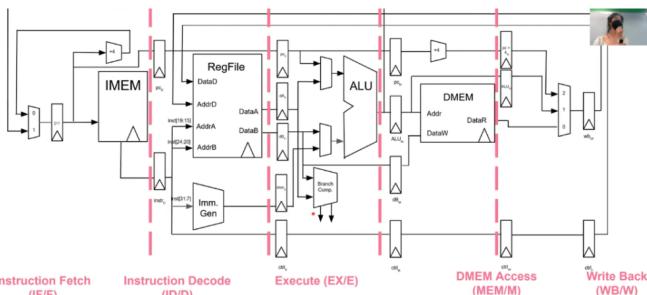
Control Logic "Truth Table" (incomplete)

Inst[31:0]	BrEq	BrLT	PCSel	ImmSel	BrUn	ASel	BSel	ALUSel	MemRW	RegWEn	WBSel	
add	*	*	+4	-	-	Reg	Reg	Add	Read	1	ALU	means "for all values"
sub	*	*	+4	-	-	Reg	Reg	Sub	Read	1	ALU	- means "don't care, use any value"
(R-R Op)	*	*	+4	-	-	Reg	Reg	(Op)	Read	1	ALU	
addi	*	*	+4	I	-	Reg	Imm	Add	Read	1	ALU	
lw	*	*	+4	I	-	Reg	Imm	Add	Read	1	Mem	
sw	*	*	+4	S	-	Reg	Imm	Add	Write	0	-	
beq	0	*	+4	B	-	PC	Imm	Add	Read	0	-	
bne	1	*	ALU	B	-	PC	Imm	Add	Read	0	-	
bne	0	*	ALU	B	-	PC	Imm	Add	Read	0	-	
bne	1	*	+4	B	-	PC	Imm	Add	Read	0	-	
blt	*	1	ALU	B	0	PC	Imm	Add	Read	0	-	
bltu	*	1	ALU	B	1	PC	Imm	Add	Read	0	-	
jair	*	*	ALU	I	-	Reg	Imm	Add	Read	1	PC+4	
jal	*	*	ALU	J	-	PC	Imm	Add	Read	1	PC+4	
auipc	*	*	+4	U	-	PC	Imm	Add	Read	1	ALU	

12

Instr	BrEq	PCSel	ImmSel	BrUn	ASel	BSel	ALUSel	MemRW	RegWEn	WBSel
add	X	0	X	X	0	0	Add	0	1	1
lw	X	0	I	X	0	1	Add	0	1	0
bge	0/1	0/1	SB	0	1	1	Add	0	0	X
sw	X	0	S	X	0	1	Add	1	0	X
auipc	X	0	U	X	1	1	Add	0	1	1
jal	X	1	UJ	X	1	1	Add	0	1	2

	IF	ID	EX	MEM	WB
add	X	X	X		X
ori	X	X	X		X
lw	X	X	X	X	
sw	X	X	X		
beq	X	X	X		
jal	X	X	X		X
bltu	X	X	X		



- Non-pipelined clock cycle: CP of datapath
- pipelined clock cycle: max(CP of each stage)
- pipelining increases latency (max delay * numcycle) and throughput (1/cycletime)

$$\frac{\text{time}}{\text{program}} = \frac{\text{instructions}}{\text{program}} \cdot \frac{\text{cycles}}{\text{instruction}} \cdot \frac{\text{time}}{\text{cycle}}$$

- increase performance by

	Instructions / program	Cycles/instruction	Time/cycle
A	decrease	decrease	same
B	same	increase	decrease
C	same	same	decrease
D	increase	decrease	increase

Single cycle fastest clock time: $t_{clk} \geq t_{PC} \text{ clk-to-q} + t_{IMEM} \text{ read} + t_{RF} \text{ read} + t_{MUX} + t_{ALU} + t_{DMEM} \text{ read} + t_{MUX} + t_{RF} \text{ setup}$

Hazards * reordering code also helps

Structural Hazard - conflict of resource

- solutions 1) stall, 2) more hardware; 3) separate memories

Register Access Hazard (R-Type) - register access

- solution separate ports

Register Access Policy: WB updates value, then ID reads new value

- ALU results -> 1) stalling, 2) forwarding (after EX)

```
if (rs1(n + 1) == rd(n) || rs2(n + 1) == rd(n) & RegWen(n) == 1) {
    forward ALU output of instruction n
}
```

Data Hazard (Load) wait until MEM

- solution load delay slots

Control Hazard (branch)

- branch prediction to reduce branch penalties -> keep a cache, if seen, backward

IF: PC clk-to-q + IMEM read + Reg setup

ID: Reg clk-to-q + max(RF read, ImmGen) + Reg setup

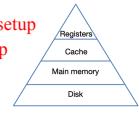
EX: Reg clk-to-q + max(ALU computation, Branch comparator) + mux

Branch comparator: PC clk-to-q + branch comp.

ALU computation: Reg clk-to-q + mux + ALU + Reg setup

MEM: Reg clk-to-q + max(DMEM read, DMEM write) + Reg setup

WB: Reg clk-to-q + mux + RF setup



Fully Associative Caches * only compulsory & capacity miss (no conflict miss)

- temporal locality vs. spatial locality
- Byte offset bits = log2(line size)
- Tag bits = address bits - offset bits

Write through policy

- write to cache and memory at the same time

Write back policy

- write data in cache and set dirty bit to 1, write to memory when this line gets evicted

Direct-mapped cache

- $\text{offset} = \log_2(\text{line size})$, $\text{index} = \log_2(\# \text{lines})$, tag = address - 1 - o

Write-allocate vs. no write-allocate

- write allocate on write miss, bring the line into cache and update the line
- No write-allocate on a write miss, only update memory, not bring the line
- write through, no write-allocate
 - write hits: update cache and main memory
 - write miss: not bring to cache, update memory
- write back, write-allocate:
 - Read and write miss: bring line into memory
 - write: only update cache, set dirty bit to 1

Approximate LRU - Clock Algorithm

- Each line has 1 referenced bit, set to 1 each time the line is accessed
- when evicting, start at the first entry in the cache and check reference bit
 - If bit is 1, set it to 0 and move to the next line
 - If bit is 0, evict this line
- The next time need to evict a line, start searching at the point where we left off

When is other eviction policies better than LRU

- 4 cache lines, iterating through an array to access ABCDE, each element maps to a different line in cache
- LRU has 0% hit rate, MRU has 75% hit rate, random has non-zero hit rate

Set Associative Cache

- $\text{offset} = \log_2(\text{line size})$, $\text{index} = \log_2(\text{cache size}/(\text{block size} * \text{associativity}))$, $\text{index} = \log_2(\text{num of sets})$, tag = address - 1 - o

Compulsory Miss → first access

Capacity Miss → blocks replaced and later retrieved

Conflict Miss → set-associative or direct mapped, won't occur in fully associative

Average Memory Access Time (AMAT)

$$\text{AMAT} = \text{hit time} + \text{miss rate} * \text{miss penalty}$$

$$\text{AMAT} = L1 \text{ hit time} + L1 \text{ miss rate} * (L2 \text{ hit time} + L2 \text{ miss rate} * L2 \text{ miss penalty})$$

Local Hit Rate: # hits at this level / # accesses to this level

Global Hit Rate: # hits at this level / total number of accesses

$$L1 \text{ Miss Rate} = \text{global MR} / \text{local MR}$$

Associativity & AMAT

- Associativity increase → hit time increase, miss rate decrease (less conflict misses)
- Associativity changes, miss penalty mostly unchanged

Number of Entries & AMAT

- #entries increase → hit time increase (read data from larger memory)
- miss rate decrease (inc capacity, dec conflict miss)
- Miss penalty unchanged
- Increase in hit time overcome improvement in hit rate at some point, decrease performance

Block Size & AMAT

- Block size increase → hit time mostly unchanged (slightly reduced as #tags reduced)
- miss rate decrease at first due to spatial locality, then increase due to increased conflict misses due to fewer blocks in cache
- Miss penalty increase
- Increase cache size doesn't always improve hit rate (depends on nature of program).
- Increase associativity decreases index bits, increase tag bits
- Decrease associativity: decrease hit time (check less spots), decrease hit rate (increase conflict misses)
- Increase cache size: increase hit rate, increase hit time (longer to search more data)
- "program run indefinitely" = miss rate is 1
- miss penalty → determined by block size

Multilevel Cache

- L2 bigger than L1, so higher hit rate, longer to access (farther fr processor)
- If line in L1 dirty when evicted, update copy in L2
- L2 reduces L1 miss penalty

Parallelism

- Matrix multiplication: $C_{ij} = \sum_k a_{ik} b_{kj}$
- Row major: $a_{ij} = a[i^{\text{th}} \text{ row}, j^{\text{th}} \text{ col}]$, Column major: $a_{ij} = a[i^{\text{th}} \text{ col}, j^{\text{th}} \text{ row}]$
- $C[i^{\text{th}} \text{ row}, j^{\text{th}} \text{ col}] = \sum_k a[i^{\text{th}} \text{ row}, k^{\text{th}} \text{ col}] \times b[k^{\text{th}} \text{ col}, j^{\text{th}} \text{ row}]$

		Data Streams	
		Single	Multiple
Instruction Streams	Single	SISD: Intel Pentium 4	SIMD: SSE instructions of x86
Instruction Streams	Multiple	MISD: No examples today	MIMD: Intel Xeon e5345 (Clovertown)

Amdahl's Law

- Speedup w/ E = Exec time w/o E / Exec time w/ E
- Fraction F, factor S
- Execution time w/ E = Execution Time w/o E * $[1-F+F/S]$
- Speedup w/E = $1 / [1-F+F/S]$

- Load balancing every processor do the same amount of work

Two ways to use a multiprocessor

- Deliver high throughput for independent jobs via job-level parallelism
- Improve run time of a single parallel-processing program

Threads

- A sequential flow of instructions that performs some task.
- Each thread has a PC+processor registers and accesses the shared memory of the process
- Each core provides one or more hardware threads
- Common intel chip supports 2 threads/core
- OS multiplexes multiple software threads onto available hardware threads

Hardware multithreading

- Two copies of PC and registers inside processor hardware
- Multithreading → better utilization, share integer adders, floatingpoint units, all caches, memory controller

- Multicore → duplicate processors, share outer caches, memory controller

* OMP thread or OS (software) threads

- Set num threads: `omp_set_num_threads(x)`
- Get num threads: `num_th = omp_get_num_threads()`
- Get thread id: `th_ID = omp_get_thread_num()`
- Load reserved `lr, rd, rs` (load word pointed to by rs into rd)
- Store conditional `sc rd, rs1, rs2` (store values in rs2 into memory pointed to by rs1 only if reservation is valid and set the status in rd. return 0 if success, return nonzero if location has changed)

Test-and-Set

```
li t2, 1
Try: lr t1, s1
      bne t1, x0, Try
      sc t0, s1, t2
      bne t0, x0, Try
```

Locked:

```
# critical sec
```

Unlock:

```
void sw x0, 0(s1) ←
```

```
const long num_steps = 10;
double step = 1.0/(double)num_steps;
double sum[NUM_THREADS];
for (int i=0; i<NUM_THREADS; i++) sum[i] = 0;
omp_set_num_threads(NUM_THREADS);
#pragma omp parallel
{
    int id = omp_get_thread_num();
    for (int i=id; i<num_steps; i+=NUM_THREADS) {
        double x = (i+0.5)*step;
        sum[id] += 4.0*step/(1.0+x*x);
        printf("i=%d, id=%d\n", i, id);
    }
    double pi = 0;
    for (int i=0; i<NUM_THREADS; i++) pi += sum[i];
    printf ("pi = %6.12f\n", pi);
}
```

for (int x = start; x < start + len/num_threads && x < len; x++)
if ((i / (n/omp_get_num_threads())) == omp_get_thread_num())

<code>_m128i_mm_set1_epi32(int i)</code>	sets the four signed 32-bit integers to i
<code>_m128i_mm_loadu_si128(_m128i *p)</code>	returns 128-bit vector stored at pointer p
<code>_m128i_mm_add_epi32(_m128i a, _m128i b)</code>	returns vector (a0+b0, a1+b1, a2+b2, a3+b3)
<code>_m128i_mm_mullo_epi32(_m128 a, _m128 b)</code>	returns vector (a0*b0, a1*b1, a2*b2, a3*b3)
<code>_m128i_mm_cmplt_epi32(_m128 a, _m128 b)</code>	returns: vector((a0>b0)?0xffffffff:0, (a1>b1)?0xffffffff:0, (a2>b2)?0xffffffff:0, (a3>b3)?0xffffffff:0)
<code>void_mm_storeu_si128(_m128i *p, _m128 a)</code>	stores 128-bit vector a at pointer p

Loop Unrolling

```
void multMat4( int n, float *A, float *B, float *C ) {
    int i,j,k;
    /* This is jki loop order. */
    for( j = 0; j < n; j++ )
        for( k = 0; k < n; k++ )
            for( i = 0; i < n; i++ )
                C[i+j*n] += A[i+k*n]*B[k+j*n];
}

void multMat6( int n, float *A, float *B, float *C ) {
    int i,j,k;
    /* This is kji loop order. */
    for( k = 0; k < n; k++ )
        for( j = 0; j < n; j++ )
            for( i = 0; i < n; i++ )
                C[i+j*n] += A[i+k*n]*B[k+j*n];
}

void transpose_blocking(int n, int blocksize, int *dst, int *src) {
    // YOUR CODE HERE
    for (int i = 0; i < n; i += blocksize){
        for (int j = 0; j < n; j += blocksize) {
            for (int x = i; x < blocksize + i; x++){
                for (int y = j; y < blocksize + j; y++) {
                    if (x < n && y < n){
                        dst[y+x*n] = src[x+y*n];
                    }
                }
            }
        }
    }
}
```

A. (0.5 pt) L1 Number of PTEs

64

L1 PTEs * # of PTEs in L1 PT = 1 * 2^6 = 64

B. (0.5 pt) L2 Number of PTEs

4096

L2 PTEs * size of L2 PT = 2^6 * 2^6 = 2^12 = 4096

```

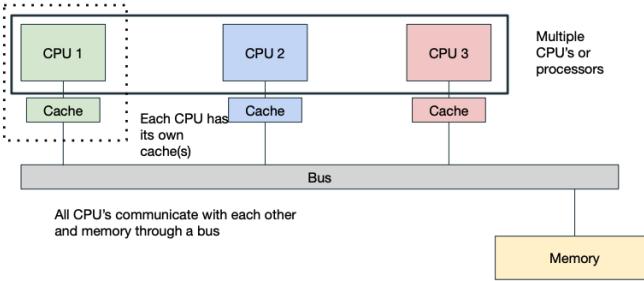
long long int sum_simd(int vals[NUM_ELEMS]) {
    clock_t start = clock();
    _m128i _127 = _mm_set1_epi32(127); // This is a vector with 127s in
    long long int result = 0; // This is where you should put your final
    /* DO NOT MODIFY ANYTHING ABOVE THIS LINE (in this function) */

    for(unsigned int w = 0; w < OUTER_ITERATIONS; w++) {
        /* YOUR CODE GOES HERE */
        _m128i sum = _mm_setzero_si128();
        for (unsigned int i = 0; i < NUM_ELEMS / 4 * 4; i += 4) {
            _m128i curr = _mm_loadu_si128((__m128 *) (vals + i));
            _m128i mask = _mm_cmplt_epi32(curr, _127);
            mask = _mm_and_si128(mask, curr);
            sum = _mm_add_epi32(sum, mask);
        }
        int temp[4];
        _mm_storeu_si128((__m128i *) temp, sum);
        result += temp[0] + temp[1] + temp[2] + temp[3];
        /* Hint: you'll need a tail case. */
        for (unsigned int i = NUM_ELEMS / 4 * 4; i < NUM_ELEMS; i++) {
            if (vals[i] >= 128) {
                result += vals[i];
            }
        }
    }
}

```

Shared Memory Multiprocessor (SMP)

- Two or more identical CPUs/Cores, single shared coherent memory



Let's go through an example to motivate why we need cache coherence!

One bank of memory is shared by all CPU's

M: Modified

- I have the **only** copy, and can write, and it's dirty
- If this gets evicted, I need to flush the entry

O: Owned

- I have the **official** copy, and can write, and it's dirty
- When writing, I have to tell everyone else I'm writing (and it now turns Modified)

E: Exclusive

- I have the **only** copy

S: Shared

- I have a copy and can read away...

I: Invalid (duh)

- Shared up-to-date data, other caches may have a copy (valid bit set, shared bit set)

- Modified up-to-date data, changed (dirty), no other cache has a copy, OK to write, memory out-of-date (valid and dirty bit is set)

- Exclusive up-to-date data, no other cache has a copy, OK to write, memory up-to-date

- If other cache reads this line, the state becomes shared

- If write to this line, state becomes modified, no need to broadcast

- Owner up-to-date data, other caches may have a copy (must be in shared state), memory not up-to-date

- This cache supplies data on read instead of going to memory

Invalid

Valid	Dirty	Shared	State
0	0	0	Invalid
0	0	1	Invalid
0	1	0	Invalid
0	1	1	Invalid
1	0	0	Exclusive
1	0	1	Shared
1	1	0	Modified
1	1	1	Owned

Operating System

- Provides **isolation** between running processes
- Provides **interaction** with the outside world (external devices)

Illusion - provide clean, easy-to-use abstractions of physical resources

Referee - manage protection, isolation, and sharing of resources

Context Switch

- Switching from executing one program to another
- Allow multiple processes to run on the same processor
- The OS determines when to context switch
 - The OS takes control of the CPU from the current process; saves the state of the current process; loads the state of the next process; hands over the CPU to the next process

Protection

- OS isolates processes from each other; OS isolates itself from other processes.

Dual Mode Operation

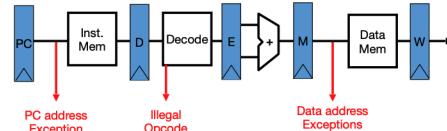
- 1) Kernel Mode ("supervisor"), 2) **User Mode** (OS mostly run in)
- Switching between user mode and kernel mode
 - System calls, interrupts, exceptions

System Calls (syscall)

- Allows the program to request a service from the operating system
 - Creating and deleting files; reading and writing files; accessing external devices like a scanner; ecalls in RISC-V
- Similar to function calls except it's executed by the kernel

Interrupts vs. Exceptions

- Interrupts**
 - Caused by an event *external* to the current running program
 - Ex: Key press
 - Asynchronous to the current program
 - Does not need to be handled immediately, but should be handled soon
- Exceptions**
 - Caused by an event *during* the execution of the current program
 - Ex: illegal instruction, divide by zero
 - Synchronous
 - Must be handled immediately



Traps

- All instructions before the faulting instruction must complete
- All instructions after the faulting instruction must be flushed
- The faulting instruction must be flushed
- Execution of the trap handler begins

Trap Handler

- code that services the interrupt or exception
- **Save the state of the current program (save all of the registers)**
- Determine what caused the exception or interrupt
- Handle exception of interrupt
 - Either continue execution of the program
 - Restore the state of the program
 - Return control to the program
 - or terminate the program
 - Terminate the program (free resources)
 - Schedule a new program

Which path to choose?

- Continue execution: interrupts, certain memory exceptions
- Terminate: illegal instruction, certain illegal memory accesses

Kernel

- Core of OS, manage resources (scheduling, memory, IO)
- Not in kernel UI, networking

Boot

1. The BIOS (Basic Input/Output System) runs
 - Power-on-self-test (POST)
 - The BIOS finds and executes the bootloader
2. The bootloader loads part of the operating system
3. The operating system initializes services, drivers, etc
4. Launch a process that waits for an input in a loop

Bootstrapping: A chain of stages, in which at each stage, a smaller, simpler program loads and then executes the larger, more complicated program of the next stage (Wikipedia)

Berkeley EECS

Executing a program

1. Loader loads program into memory
2. Loader sets argc and argv
3. OS jumps to main and transfers control to the process

Virtual Memory

Page size = 2^{offset}

#Page table entry = 2^{VPN}

PTE = PPN + metadata (account word alignment)

PTE size = $2^{\text{PTE (in bytes)}} * \#PTE = 2^{\text{PTE (in bytes)}} * 2^{\text{VPN}}$

#page = PT size / page size

PPN bits = $\log_2(\text{num pages}) = \log_2(\text{RAM size})$

Why VM: 1) RAM not big enough, 2) contiguous data, 3) protection b/w processes

- Q: How many bits would there be for the VPN, PPN, and page offset on a 32-bit machine with 8GB of RAM and 16KB pages?
 - Number of page offset bits = $\log_2(16\text{ KB}) = \log_2(2^4 * 2^{10}) = 14$
 - Number of VPN bits = $32 - 14 = 18$
 - Number of PPN bits = $\log_2(8\text{ GB}) - 14 = \log_2(2^3 * 2^{30}) - 14 = 33 - 14 = 19$
- How many virtual pages and physical pages would I have in a 32-bit system where I have a page size of 8KB and 8GB RAM?
 - Size of virtual address space: 2^{32}
 - Number of virtual pages: $2^{32} / 8\text{ K} = 2^{32} / (2^3 * 2^{10}) = 2^{19}$
 - Size of physical address space: 8G = $(2^3 * 2^{33}) = 2^{33}$
 - Number of physical pages: $2^{33} / 8\text{ K} = 2^{33} / (2^3 * 2^{10}) = 2^{20}$
- Number of entries in the page table depend on **Virtual Address Space**.
- The **number of entries decreases** as the size of the pages **increases**.
- The **size of each entry decreases** as the size of the pages **increases**.
- With VM, **2 accesses** to make each time to access a piece of data (read page table, use translated address to read data)

Page Fault Exception * page is not in RAM

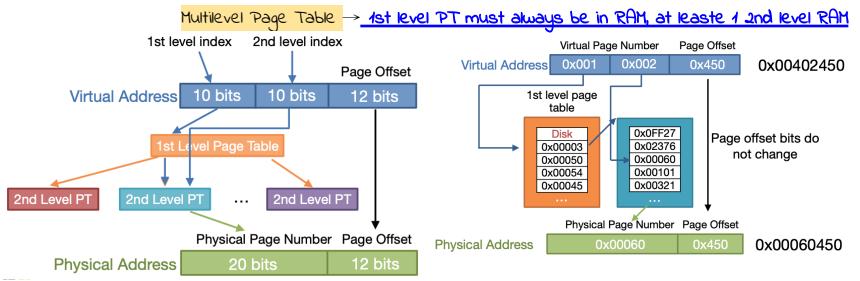
1. Hardware (CPU) generates a **page fault exception**
2. The hardware jumps to the OS page fault handler to fix it
3. The OS chooses a page to evict from **RAM** (if no more free space)
4. If the page is **dirty**, it needs to be written back to **disk** first
5. The OS updates the corresponding **PTE** to point to **disk**
6. The OS brings in the page we wanted disk and puts it in **RAM**
7. The OS updates the **PTE** of the new page
8. The OS jumps back to the instruction that caused the page fault (This time it won't cause a page fault since the page has been loaded.)

Page Table Metadata

- valid, dirty
- Page protection: read, write, execute
 - Memory protection fault or segmentation fault

Translation Lookaside Buffer (TLB) → cache for page table

- Must be small to be fast
- Usually fully associative cache, random or FIFO replacement policy
- On a **context switch**, TLB is **flushed**.



- Q: If I'm running 10 applications on my 32-bit computer with 4GB RAM, 4KB pages and 1-level page tables, how much of my RAM is consumed by page tables? (size of PTE is 4 bytes)

- VA size = 32 bits
- PA size = $\log_2(4\text{ GB}) = \log_2(2^2 * 2^{30}) = 32$ bits
- # bits in page offset = $\log_2(4\text{ KB}) = \log_2(2^2 * 2^{10}) = 12$
- # bits in VPN = $32 - 12 = 20$
- # entries in page table = 2^{20}
- size of each entry is ~4 bytes
- size of one page table = $2^{20} * 2^2 = 2^{22}$
- total RAM consumed by pages tables = $10 * 2^{22} \text{ bytes} = 40 \text{ MB}$

- Q: If I'm running 10 applications on my 32-bit computer with 4GB RAM, 4KB pages and a 2-level page table, how much of my RAM is consumed by 1st level page tables? (size of PTE is 4 bytes)

- # bits in VPN = $32 - 12 = 20 = 10$ bits for level 1 + 10 bits for level 2
- Size of 1st level page table = page size = 4KB
- total RAM consumed by 1st level pages tables = $10 * 2^{12} \text{ bytes} = 40 \text{ KB}$

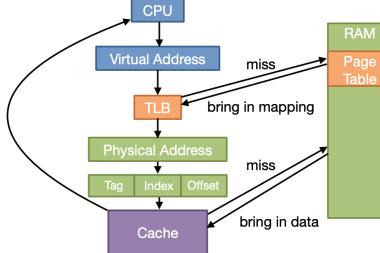
- **TLB contains mappings from VPN to the final PPN**

Synonyms → Different virtual addresses can map to the same physical address

Homonyms → The same virtual address can map to different physical addresses

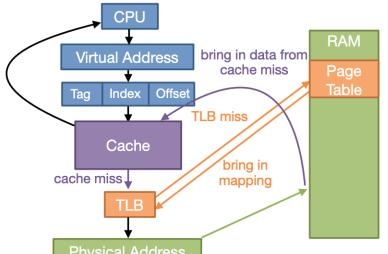
- For **PIPT**, no need to flush the cache since flushing the TLB will force cache replacements (since the new process will point to a different physical address, which is not in the cache).
- For **VIPT**, you need to flush the cache since a different process might use the same virtual addresses.
- For **VIPT**, no need to flush the cache since the tag comparison will fail anyway due to TLB flush (new set of physical addresses = new tags = forced cache replacements like in PIPT).
- In all cases, **TLB always gets flushed when context switching**.

Physically Indexed, Physically Tagged Cache



- Pros: Can use the same cache for multiple processes
- Cons: must translate address before accessing the cache
- Switch process: 1) update PT address register, 2) clear TLB valid bits
- AMAT: translation (TLB, PT) + data access (cache)

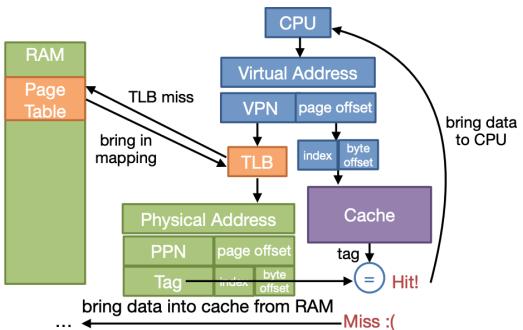
Virtually indexed, virtually Tagged Cache



- Cannot use the same virtual cache for different processes (synonyms and homonyms)
- Solve this issue by flushing the cache on context switch

Virtually indexed, Physically Tagged Caches (VIPT)

- Use physical address to access cache → avoid synonym and homonym
 - Size of cache is limited by number of page offset bits
 - How to make cache larger with the same page size? → increase associativity
- Q: With 4kB pages, how many bytes can a direct-mapped (1-way) VIPT cache store?
- 4kB
 - We can only use the page offset bits (12 bits for 4kB pages) to index into the cache. So the index can only address 12 bits of address, or 4kB of data



Memory Mapped I/O

- writes to device, hardware intercepts, no actual memory access.

Polling

- CPU periodically checks the status register to see if there is any data to receive
- occur while performing a context switch

- Let's say it takes **400 cycles** to complete a polling operation
- We have to poll a mouse 30 times per second. Our clock runs at **1 GHz** (1 billion cycle/second).
- How much of our processor time is spent polling?

$$\begin{aligned} \text{Number of cycles per second spent polling} \\ 30 \text{ polls/sec} * 400 \text{ cycles/poll} = 12K \text{ cycles/sec} \\ 12K \text{ cycles/sec} / 10^9 \text{ cycles/sec} = 0.0012\% \end{aligned}$$

- I want to read from a disk that can transfer data at a rate of **16 MB/sec**
- I can only accept **16B** per poll
- The clock runs at **1 GHz** (1 billion cycle/second).
- What percent of my processor time is spent polling?

$$\begin{aligned} \text{Number of polls needed per second to ensure that no data is missed} &\longrightarrow (16 \text{ MB/sec}) / (16 \text{ B/poll}) = 1M \text{ polls/sec} \\ \text{Number of cycles per second spent polling} &\longrightarrow 1M \text{ polls/sec} * 400 \text{ cycles/poll} = 400M \text{ cycles/sec} \quad \text{Unacceptable} \\ 400M \text{ cycles/sec} / 10^9 \text{ cycles/sec} &= 40\% \end{aligned}$$

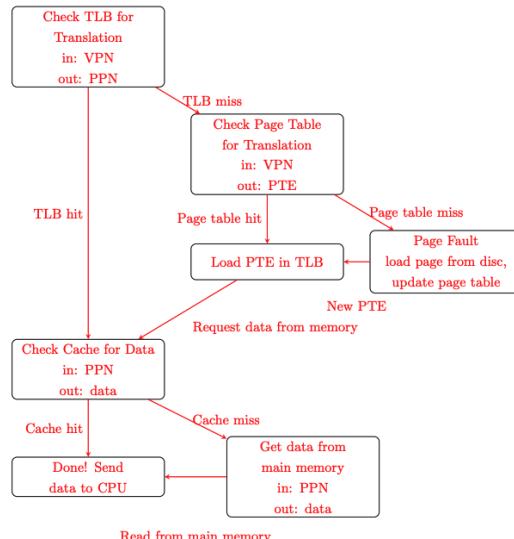
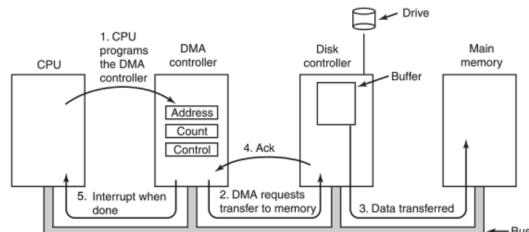
Interrupts

- Polling wastes resources when no data is available from I/O device
- Interrupts stop current program, transfer control to trap handler in OS
- Save program counter and user registers

	Interrupts	Polling
If no I/O activity	No wasted cycles	Lots of wasted cycles
If lots of I/O activity	Expensive – saving/restoring state	Less expensive - only poll when we are already context switching
Better suited for events that are	<ul style="list-style-type: none"> asynchronous (unsure when event will occur) urgent infrequent 	<ul style="list-style-type: none"> synchronous (occurs at fixed intervals) not urgent frequent (majority of polls are a hit)
Examples	Disk	Keyboard, mouse

Direct Memory Access (DMA)

- DMA engine (hardware), controls movement of data b/w memory and I/O devices, allow CPU to do other work
- Initiate transfer, CPU setup memory address to place data, # of bytes, I/O device #, direction of transfer
- Incoming Data
 - Receive interrupt from device
 - CPU takes interrupt, initiates transfer
 - DMA engine handle the transfer (CPU free)
 - Upon completion, DMA engine interrupt the CPU again
- Outgoing Data
 - CPU initiate transfer, confirms external device ready
 - CPU initiates transfer
 - DMA engine handle transfer
 - DMA engine interrupt CPU to signal completion



Network

Software send steps

- Application copies data to OS buffer
- OS calculates checksum
- OS sends DMA request to network interface hardware, start

Software receive steps

- Network interface copies data from network interface hardware to OS buffer, triggers interrupt
- OS calculates checksum, if OK send ACK, if not, delete message (sender resend)
- If OK, OS copies data to user address space, signal application to continue

Hierarchy of layers - network "stack":

- Application (chat client, game, etc.)
- Transport (TCP, UDP)
- Network (IP)
- Data Link Layer (ethernet)
- Physical Link (copper, wireless, etc.)



- Protocol packet structure and control commands to manage communication
- Protocol families (suites) a set of cooperating protocols that implement the network stack
- Encapsulation carry higher level info within lower level envelope

Layer 4 Protocol Transmission Control Protocol (TCP)

- Connection based, ports, reliable, in order, bytestream, congestion control

Universal Datagram Protocol (UDP)

- No connections, datagram, unreliable, out of order
- Use UDP than TCP when:
 - Can't tolerate latency in initiating a connection
 - old data is useless data (voice and video games), TCP pause data stream in case of error

Types of Faults in Digital Designs

- Design bugs, manufacturing defects, runtime failures
- *** Dependability via Redundancy redundancy → error detection

Spatial Redundancy

replicated data/info/hardware to handle hard and soft failures

Temporal Redundancy

retry to handle soft failures

Dependability Measures

- Reliability: Mean Time To Failure (MTTF)
- Service interruption: Mean Time to Repair (MTTR)
- Mean time between failures $MTBF = MTTF + MTTR$
- Availability = $MTTF / (MTTF + MTTR)$
- Improving availability
 - Increase MTTF (more reliable hw/sf + fault tolerance)
 - Decrease MTTR (improve tools for diagnosis and repair)
- Annualized Failure rate $AFR = C * \text{hw} * 8760 \text{ hr/yr} / MTTF = 1 / MTTF$

1 nine: 90% => 36 days of repair/year

- Eduroam Reliability?

2 nines: 99% => 3.6 days of repair/year

3 nines: 99.9% => 526 minutes of repair/year

4 nines: 99.99% => 53 minutes of repair/year

5 nines: 99.999% => 5 minutes of repair/year

Error Detection/Correction Codes (EDC/ECC)

Parity Bit → error detection coding

- Non-zero parity check indicates error occurred
- Even # of errors on different bits not detected
- odd # of errors are detected

X = 0101 0101

4 ones, even parity now

Write to memory:

0101 0101 0 **to keep parity even**

Y = 0110 0111

5 ones, odd parity now

Write to memory:

0110 0111 1 **to make parity even**

• Read X from memory

0101 0101 0

• 4 ones => even parity, so no error

• Read X from memory

1101 0101 0

• 5 ones => odd parity, so error

Hamming ECC

Single Error Detection (SEC) - min hamming dist = 3

- Use more parity bits to pinpoint bit(s) in error, so they can be corrected.

- Example: Single error correction (SEC) on 4-bit data
 - use 3 parity bits, with 4-data bits results in 7-bit code word
 - 3 parity bits sufficient to identify any one of 7 code word bits
 - overlap the assignment of parity bits so that a single error in the 7-bit word can be corrected

- Procedure: group parity bits so they correspond to subsets of the 7 bits:
 - p₁ protects bits 1,3,5,7
 - p₂ protects bits 2,3,6,7
 - p₃ protects bits 4,5,6,7

berkeleyECS

1	2	3	4	5	6	7
p ₁	p ₂	d ₁	p ₃	d ₂	d ₃	d ₄

Note:
number bits
from left to
right.

Bit position number

001 = 1₁₀

011 = 3₁₀

101 = 5₁₀

111 = 7₁₀

010 = 2₁₀

011 = 3₁₀

110 = 6₁₀

111 = 7₁₀

100 = 4₁₀

101 = 5₁₀

110 = 6₁₀

111 = 7₁₀

1	2	3	4	5	6	7
p ₁	p ₂	d ₁	p ₃	d ₂	d ₃	d ₄

- Note: parity bits occupy power-of-two bit positions in code-word.

- On writing to memory:

- parity bits are assigned to force even parity over their respective groups.

- On reading from memory:

- check bits (c₃, c₂, c₁) are generated by finding the parity of the group and its parity bit. If an error occurred in a group, the corresponding check bit will be 1, if no error the check bit will be 0.
- check bits (c₃, c₂, c₁) form the position of the bit in error.

- Example: c = c₃c₂c₁ = 101

- error in 4,5,6, or 7 (by c₃=1)
- error in 1,3,5, or 7 (by c₁=1)
- no error in 2, 3, 6, or 7 (by c₂=0)

- Therefore error must be in bit 5.

- Note the check bits point to 5

- By our clever positioning and assignment of parity bits, the check bits always address the position of the error!

- c=000 indicates no error

Double Error Detection (DEC)

- Adding an extra parity bit covering the entire word can provide double error detection

1	2	3	4	5	6	7	8
p ₁	p ₂	d ₁	p ₃	d ₂	d ₃	d ₄	p ₄

- On reading, the C bits are computed (as usual) plus the parity over the entire word, P:

C=0 P=0, no error

C!=0 P=1, correctable single error

C!=0 P=0, a double error occurred

C=0 P=1, an error occurred in p₄ bit

- Overhead involved in single error correction code:

- let p be the total number of parity bits and d the number of data bits in a p+d bit word.

- If p error correction bits are to point to the error bit (p+d cases) plus indicate that no error exists (1 case), we need:

2^p >= p + d + 1,

thus p >= log(p + d + 1)

for large d, p approaches log(d)

1.1 How many bits do we need to add to 0011₂ to allow single error correction?

m parity bits can cover bits 1 through 2^m - 1, of which 2^m - m - 1 are data bits. Thus, to cover 4 data bits, we need 3 parity bits.

Bit position	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Encoded data bits	p ₁	p ₂	d ₁	p ₃	d ₂	d ₃	p ₄	d ₄	d ₅	d ₆	d ₇	d ₈	d ₉	d ₁₀	d ₁₁
Parity bit coverage	p ₁	X	X	X	X	X	X	X	X	X	X	X	X	X	X
p ₂		X	X												
p ₄				X	X	X	X								
p ₈								X	X	X	X	X	X	X	X

RAID Redundant Arrays of Disks

Raid 0: Striping (actually no redundancy)

- Spread data across multiple disks

- Improves bandwidth linearly, doesn't help latency

Raid 1: Disk Mirroring/Shielding (online sparing)

- Each disk is fully duplicated onto its mirror (high availability)
- writes go to disk and mirror, read from original disk unless failure
- Most expensive solution, 2x capacity overhead

Raid 3: Single parity Disk

- Disk drives themselves code data and detect failures
- Reconstruction of data can be done with P disk if know which failed

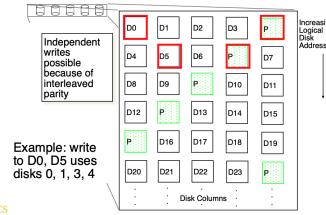
- writes change data disk and P disk
- P contains parity of other disks per stripe. If disk fails, use P and other disks to find missing info.

RAID 4 High Vo Rate Parity

- interleave data in data block rather than bit level → more parallelism
- Reconstruction of data can be done with single parity disk
- Reading (w/o) fault involve only data disk
- Writing involves data disk + parity disk (bottleneck for write)

RAID 5 High Vo Rate Interleaved Parity

- independent writes possible because of interleaved parity

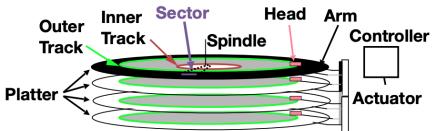


RAID 6 Add another parity block per stripe

- 2 blocks per stripe rather than 1
- sacrifice capacity for increased redundancy
- Array can tolerate 2 disk failures and continue operating

	Configuration	Pro/Good for	Con/Bad for
RAID 0	Split data across multiple disks	No overhead, fast read / write	Reliability
RAID 1	Mirrored Disks: Extra copy of data	Fast read / write, Fast recovery	High overhead → expensive
RAID 2	Hamming ECC: Bit-level striping, one disk per parity group	Smaller overhead	Redundant check disks
RAID 3	Byte-level striping with single parity disk.	Smallest overhead to check parity	Need to read all disks, even for small reads, to detect errors
RAID 4	Block-level striping with single parity disk.	Higher throughput for small reads	Still slow small writes (A single check disk is a bottleneck)
RAID 5	Block-level striping, parity distributed across disks.	Higher throughput of small writes	The time to repair a disk is so long that another disk might fail in the meantime.

Disk



Disk Access Time = Seek time + rotation time + transfer time + controller overhead

Transfer time = data size / transfer rate

Rotation time = 1/rpm * 6000 * 0.5

Transfer time = transfer rate convert to ms * 0.5

- Seek Time = time to position the head assembly at the proper cylinder
- Rotation Time = time for the disk to rotate to the point where the first sectors of the block to access reach the head
- Transfer Time = time taken by the sectors of the block and any gaps between them to rotate past the head

Flash Memory Solid State Drive (SSD)

- bandwidth similar to spinning disk, spinning has better storage density and cost density
- no seek time, no additional latency for random access vs. sequential access of a block
- HDD access time measured in milliseconds, SSD in microseconds

Power Usage Effectiveness (PUE)

$$\text{PUE} = \frac{\text{Total power}}{\text{IT power}}$$

Sources speculate Google has over 1 million servers. Assume each of the 1 million servers draw an average of 200W, the PUE is 1.5, and that Google pays an average of 6 cents per kilowatt-hour for datacenter electricity. Estimate Google's annual power bill for its datacenters.

$$1.5 \cdot 10^6 \text{ servers} \cdot 0.2 \text{ kW/server} \cdot \$0.06/\text{kW-hr} \cdot 8760 \text{ hrs/yr} \approx \$157.68 \text{ M/year}$$

Google reduced the PUE of a 50,000-machine datacenter from 1.5 to 1.25 without decreasing the power supplied to the servers. What's the cost savings per year?

$$\text{PUE} = \frac{\text{Total building power}}{\text{IT equipment power}} \implies \text{Savings} \propto (\text{PUE}_{\text{old}} - \text{PUE}_{\text{new}}) \cdot \text{IT equipment power}$$

$$(1.5 - 1.25) \cdot 50000 \text{ servers} \cdot 0.2 \text{ kW/server} \cdot \$0.06/\text{kW-hr} \cdot 8760 \text{ hrs/yr} \approx \$1.314 \text{ M/year}$$

Cloud

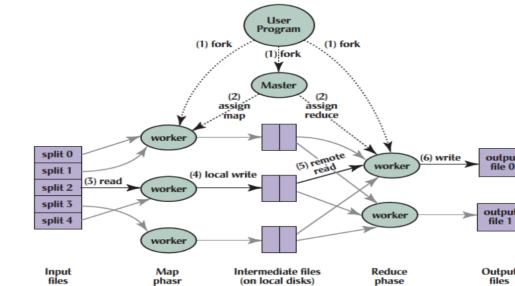
- shared platform with illusion of isolation
- attraction of low-cost cycles
- elastic service

Request-Level Parallelism (RLP)

MapReduce

- Apache Hadoop, Apache Spark

```
file.flatMap(lambda line: line.split())
           .map(lambda word: (word, 1))
           .reduceByKey(lambda a, b: a + b)
```



4.1 (Summer 2013 Q1) Solve for the maximum controller overhead to meet the following specifications: We need disk latency under 18 ms while reading 800 B of data. The hard drive spins at 6000 rev/min with a seek time of 2.5 ms and transfer rate of 80 KB/s (SI prefix). Include units!

controller overhead \leq disk latency - seek time - rotation time - transfer time
 rotation time = $0.5/\text{rpm} \cdot (60 \text{ sec}/\text{min}) = 0.005 \text{ s} = 5 \text{ ms}$
 transfer time = data size/transfer rate = $0.01 \text{ s} = 10 \text{ ms}$
 controller overhead $\leq 18 - 2.5 - 5 - 10 \text{ ms} = 0.5 \text{ ms}$
 max controller overhead $\rightarrow 0.5 \text{ ms}$

4.2 (Fall 2015 QF3) We have a hard drive with a controller overhead of 5 ms. The disk has 12000 cylinders, and it takes 2 ms to cross 1000 cylinders. The drive rotates at 2400 RPM, and we want to copy half a MB of data. Our hard drive has a transfer rate of 500 MB/s. What is the access time of a read from disk?

We want to add up the times that it takes to access this disk. We add them up in the following order: Seek + Rotation + Transfer + Controller. The number of tracks in the seek time is always number of tracks/3, and rotation time is averaged by taking half of the time to rotate around a disk.

$$12000/3 \cdot 2/1000 + 1/2400 \cdot 60 \cdot 1000 \cdot 1/2 + 1/500 \cdot 1000 + 5 \text{ ms} = 15.25 \text{ ms}$$

g) You have an SSD which can transfer data in 32-byte chunks at a rate of 64 MB/second. No transfer can be missed. If we have a 4GHz processor, which takes 200 cycles for a polling operation, what fraction of time does the processor spend polling the SSD drive for data? Leave your answer in the box provided as a percentage.

10%

SHOW YOUR WORK

$$64 \text{ MB/sec} / 32 \text{ B/poll} = 2 \text{ MB/poll/sec}$$

$$2 \text{ MB/sec} * 200 \text{ cycles/poll} = 400 \text{ cycles/sec}$$

$$4 \text{ GHz clock} \Rightarrow 4000 \text{ cycles/sec}; 400 \text{ c/s} / 4000 \text{ c/s} = 10\%$$

"redundancy through parity."

2, 3, 4, 5

"parity is bit-striped."

2

"a pro is having a small overhead."

0, 2, 3

"parity is byte-striped."

3

"fast small reads are a pro."

0, 1, 4

"parity is block-level striped."

4, 5

"fast small writes are a pro."

0, 1, 5

"disks play a part in increasing reliability."

1, 2, 3, 4, 5

"higher throughput is a pro."

4, 5