

CS 110

Computer Architecture

Lecture 24:

Review for Midterm

Instructors:

Sören Schwertfeger & Chundong Wang

<https://robotics.shanghaitech.edu.cn/courses/ca/20s/>

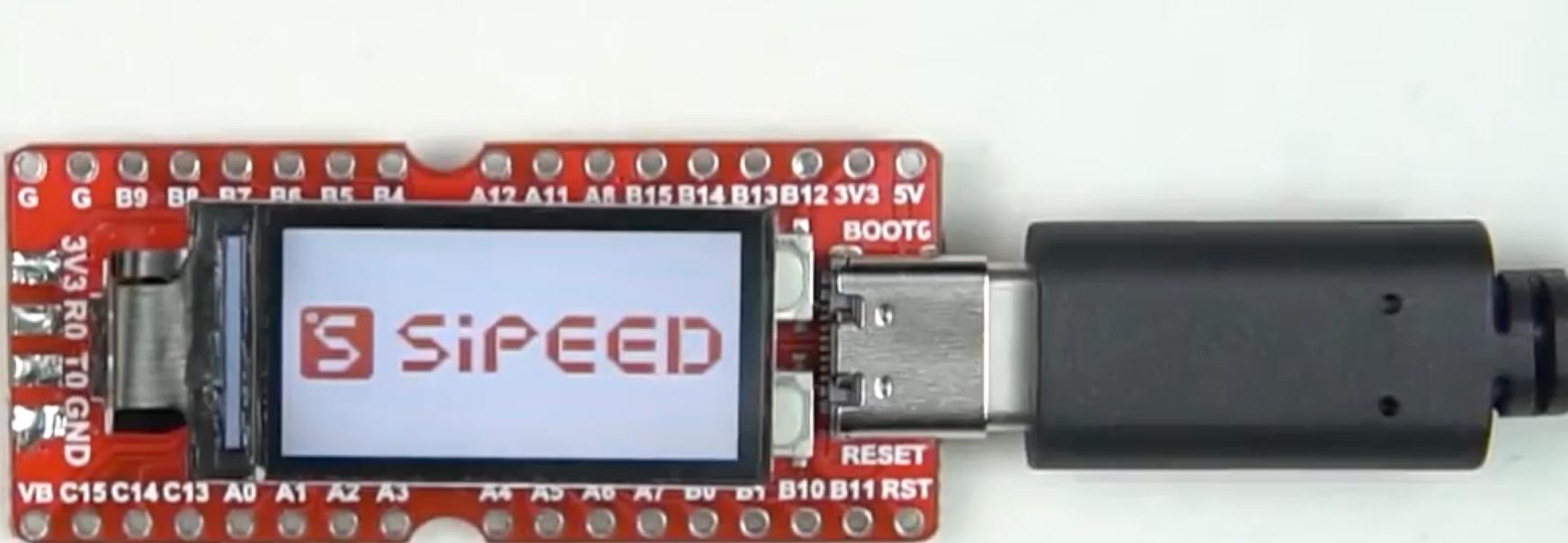
School of Information Science and Technology SIST

ShanghaiTech University

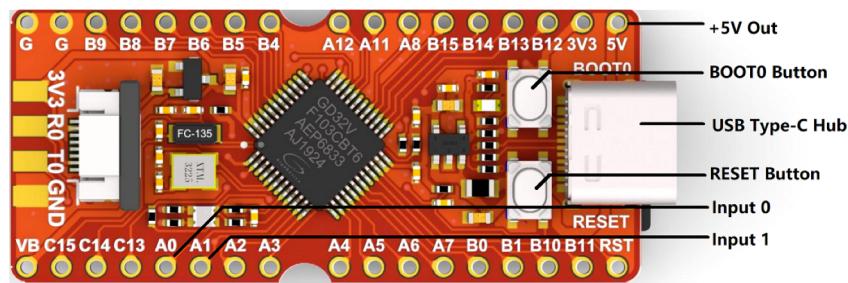
Slides based on UC Berkley's CS61C

Project 4

- Program on a RISC-V CPU!
- Sipeed Longan Nano Development Board with
 - RISC-V 32bit CPU!
 - Comes with a Screen!

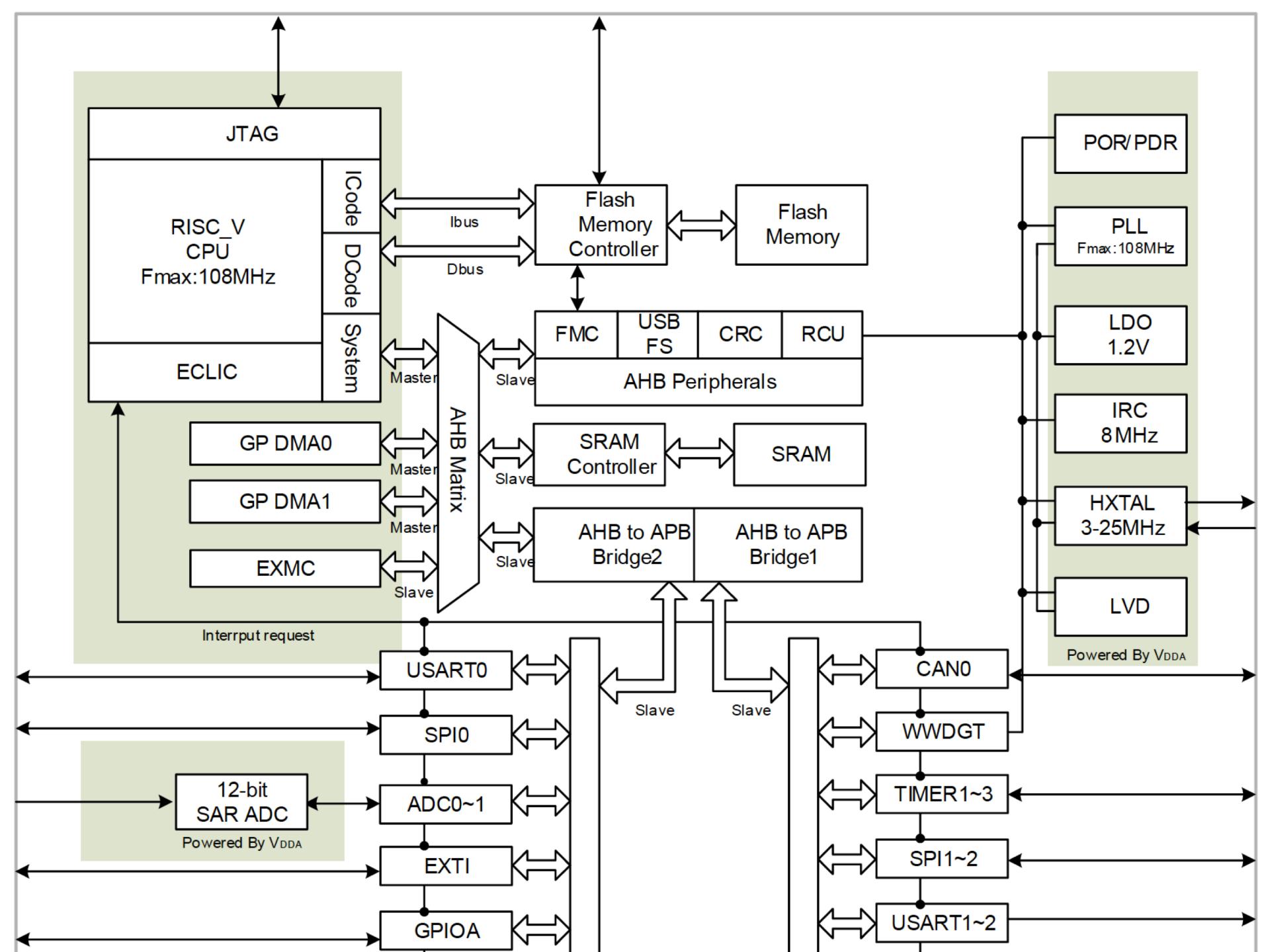


32 bit RISC-V CPU



- Gigadevince GD32VF103CBT6 **GD32VF103 “Bumblebee Core”**
- ISA: RV32IMAC
 - Integer; Multiplication & Division; Atomic Operations; Compressed (16bit) instructions
 - Single-cycle hardware multiplier and Multi-cycles hardware divider
- 108MHz
- 128kb Flash
- 32kb SRAM
- Cache?
 - No Cache needed: 0 wait states for Flash and SRAM!
(CPU speed is low while using modern memory)
- Supports misaligned memory access operations (Load/Store)
- Virtual Memory?
 - No Virtual Memory – everything (Flash, SRAM, devices) is mapped to the 32bit address space.

- DMA?
 - Yes!
 - peripheral to memory, memory to peripheral, memory to memory DMA modes
- Interrupts?
 - Yes!
 - Supports the RISC-V architecturally defined software, timer and external interrupts.
 - Dozens of external interrupt sources
 - Programmable 16 interrupt levels and priorities
- Power Saving Modes?
 - Yes!
 - Sleep (core clock off – interrupts run)
 - Deep-Sleep (most interrupts off)
 - Standby (SRAM and registers are lost – save state to flash! Few interrupts available)
- Pipelining?
 - 2-stage pipeline!
- Simple dynamic branch predictor
- Instruction fetch unit (IFU) can prefetch the following two instructions to mask the instruction memory access latency
- Support Machine Mode and User Mode



Memory Map

Pre-defined Regions	Bus	Address	Peripherals		
External device		0xA000 0000 - 0xA000 0FFF	EXMC - SWREG		
External RAM	AHB	0x9000 0000 - 0x9FFF FFFF	Reserved		
		0x7000 0000 - 0x8FFF FFFF	Reserved		
			EXMC - NOR/PSRAM/SRAM		
		0x6000 0000 - 0x6FFF FFFF	M		
Peripheral	AHB	0x5000 0000 - 0x5003 FFFF	USBFS		
		0x4008 0000 - 0x4FFF FFFF	Reserved		
		0x4002 3C00 - 0x4002 3FFF	Reserved		
		0x4002 3800 - 0x4002 3BFF	Reserved		
		0x4002 3400 - 0x4002 37FF	Reserved		
		0x4002 3000 - 0x4002 33FF	CRC		
	APB2	UX4UU1 3CUU - UX4UU1 3FFFF	Reserved		
		0x4001 3800 - 0x4001 3BFF	USART0		
		0x4001 3400 - 0x4001 37FF	Reserved		
		0x4001 3000 - 0x4001 33FF	SPI0		
		0x4001 2C00 - 0x4001 2FFF	TIMER0		
		0x4001 2800 - 0x4001 2BFF	ADC1		
		0x4001 2400 - 0x4001 27FF	ADC0		
		0x4001 2000 - 0x4001 23FF	Reserved		
		0x4001 1C00 - 0x4001 1FFF	Reserved		
		0x4001 1800 - 0x4001 1BFF	GPIOE		
		0x4001 1400 - 0x4001 17FF	GPIOD		
		0x4001 1000 - 0x4001 13FF	GPIOC		
		0x4001 0C00 - 0x4001 0FFF	GPIOB		
		0x4001 0800 - 0x4001 0BFF	GPIOA		
	Code	0x4001 0400 - 0x4001 07FF	EXTI		
		0x4001 0000 - 0x4001 03FF	AFIO		
		0x4000 7C00 - 0x4000 7FFF	Reserved		
		0x4000 7800 - 0x4000 7BFF	Reserved		
		0x4000 7400 - 0x4000 77FF	DAC		
		0x4000 7000 - 0x4000 73FF	PMU		
		0x4000 6C00 - 0x4000 6FFF	BKP		
		0x4000 6800 - 0x4000 6BFF	CAN1		
		0x4000 6400 - 0x4000 67FF	CAN0		
		0x4000 6000 - 0x4000 63FF	Shared USB/CAN SRAM 512 bytes		
	AHB	0x4000 5C00 - 0x4000 5FFF	USB device FS registers		
		0x4000 5800 - 0x4000 5BFF	I2C1		
		0x4000 5400 - 0x4000 57FF	I2C0		
		0x4000 5000 - 0x4000 53FF	UART4		
		0x4000 4C00 - 0x4000 4FFF	UART3		
	SRAM	0x4000 4800 - 0x4000 4BFF	USART2		
		0x4000 4400 - 0x4000 47FF	USART1		
		0x4000 4000 - 0x4000 43FF	Reserved		
		0x4000 3C00 - 0x4000 3FFF	SPI2/I2S2		
		0x4000 3800 - 0x4000 3BFF	SPI1/I2S1		
		0x4000 3400 - 0x4000 37FF	Reserved		
		0x4000 3000 - 0x4000 33FF	FWDGT		
		0x4000 2C00 - 0x4000 2FFF	WWDT		
		0x4000 2800 - 0x4000 2BFF	RTC		
		0x4000 2400 - 0x4000 27FF	Reserved		
	AHB	0x4000 2000 - 0x4000 23FF	Reserved		
		0x4000 1C00 - 0x4000 1FFF	Reserved		
		0x4000 1800 - 0x4000 1BFF	Reserved		
		0x4000 1400 - 0x4000 17FF	TIMER6		
		0x4000 1000 - 0x4000 13FF	TIMER5		
		0x4000 0C00 - 0x4000 0FFF	TIMER4		
		0x4000 0800 - 0x4000 0BFF	TIMER3		
		0x4000 0400 - 0x4000 07FF	TIMER2		
		0x4000 0000 - 0x4000 03FF	TIMER1		
		0x2007 0000 - 0x3FFF FFFF	Reserved		
	Code	0x2006 0000 - 0x2006 FFFF	Reserved		
		0x2003 0000 - 0x2005 FFFF	Reserved		
		0x2002 0000 - 0x2002 FFFF	Reserved		
		0x2001 C000 - 0x2001 FFFF	Reserved		
		0x2001 8000 - 0x2001 BFFF	Reserved		
		0x2000 5000 - 0x2001 7FFF	SRAM		
		0x2000 0000 - 0x2000 4FFF			
		0x1FFF F810 - 0x1FFF FFFF	Reserved		
		0x1FFF F800 - 0x1FFF F80F	Option Bytes		
		0x1FFF B000 - 0x1FFF F7FF	Boot loader		
	AHB	0x1FFF 7A10 - 0x1FFF AFFF	Reserved		
		0x1FFF 7800 - 0x1FFF 7A0F	Reserved		
		0x1FFF 0000 - 0x1FFF 77FF	Reserved		
		0x1FFE C010 - 0x1FFE FFFF	Reserved		
		0x1FFE C000 - 0x1FFE C00F	Reserved		
		0x1001 0000 - 0x1FFE BFFF	Reserved		
		0x1000 0000 - 0x1000 FFFF	Reserved		
		0x083C 0000 - 0x0FFF FFFF	Reserved		
		0x0802 0000 - 0x083B FFFF	Reserved		
		0x0800 0000 - 0x0801 FFFF	Main Flash		
	Code	0x0030 0000 - 0x07FF FFFF	Reserved		
		0x0010 0000 - 0x002F FFFF			
		0x0002 0000 - 0x000F FFFF			
		0x0000 0000 - 0x0001 FFFF	Aliased to Main Flash or Boot loader		

2.3.1.

Flash memory architecture

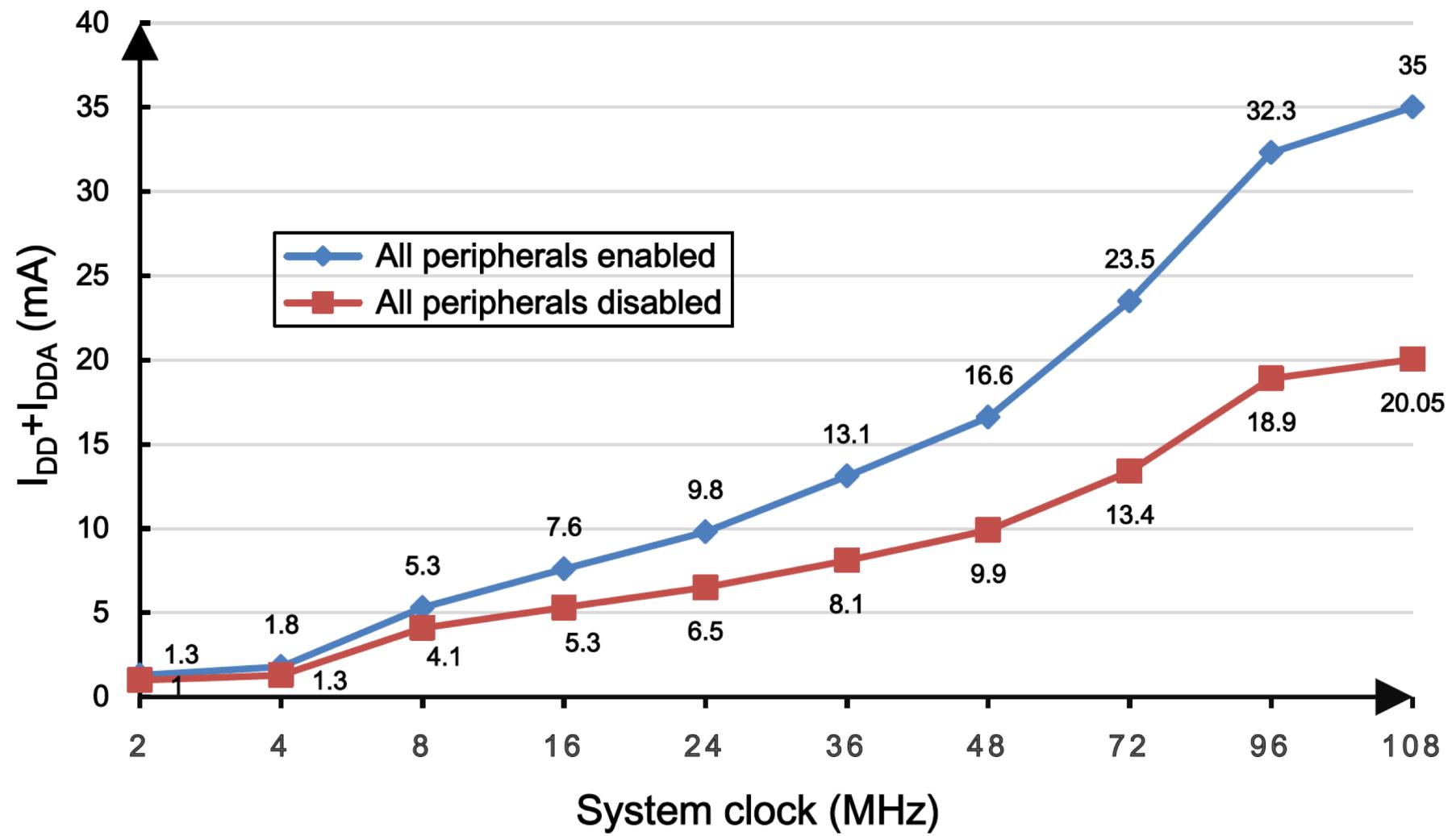
The flash memory consists of up to 128 KB main flash organized into 128 pages with 1 KB capacity per page and a 18 KB Information Block for the Boot Loader. The main flash memory contains a total of up to 128 pages which can be erased individually. The [Table 2-1. Base address and size for flash memory](#) shows the details of flash organization.

Table 2-1. Base address and size for flash memory

Block	Name	Address Range	size (bytes)
Main Flash Block	Page 0	0x0800 0000 - 0x0800 03FF	1KB
	Page 1	0x0800 0400 - 0x0800 07FF	1KB
	Page 2	0x0800 0800 - 0x0800 0BFF	1KB

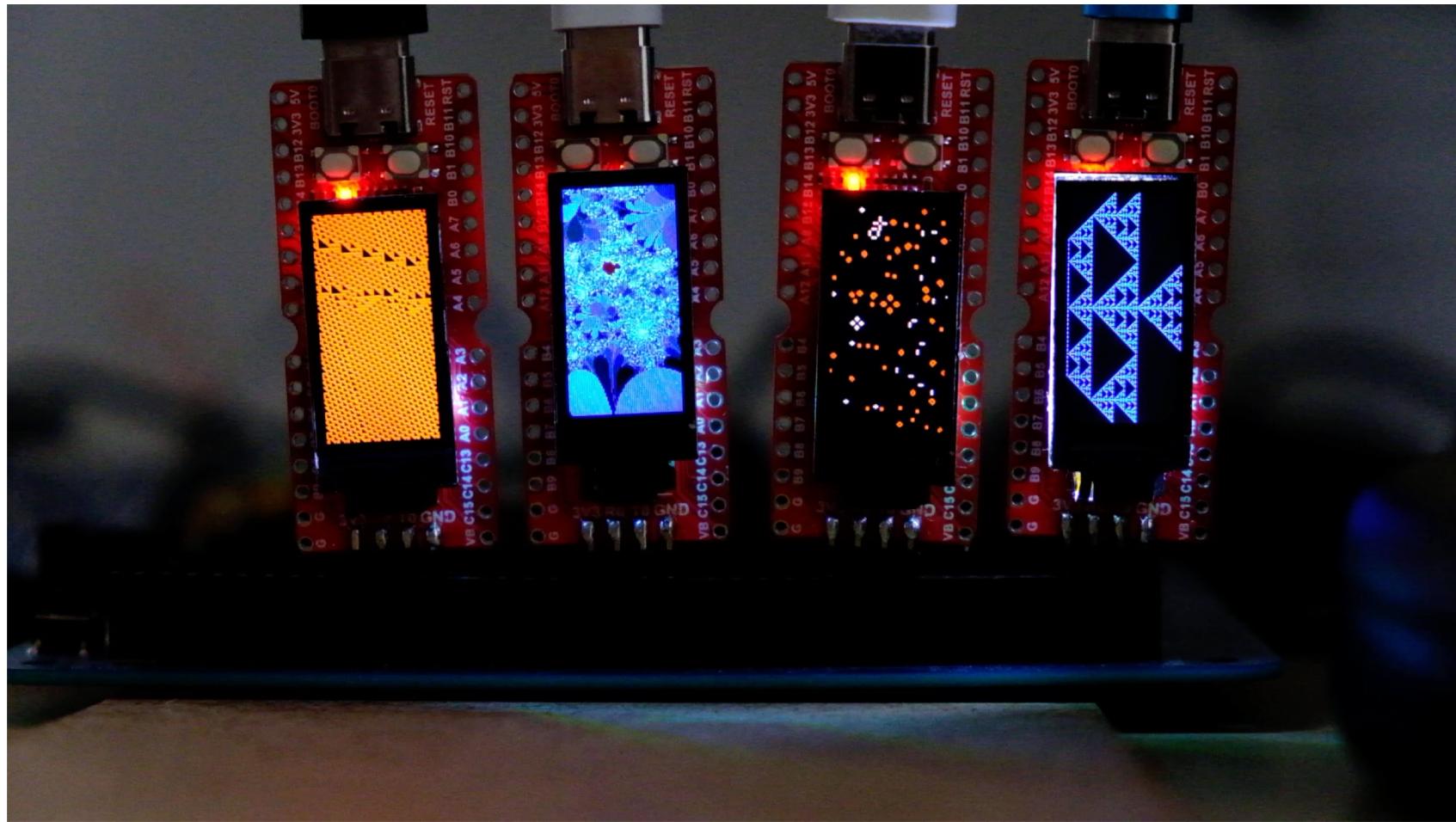
	Page 127	0x0801 FC00 - 0x0801 FFFF	1KB
Information Block	Boot loader area	0x1FFF B000- 0x1FFF F7FF	18KB
Option bytes Block	Option bytes	0x1FFF F800 - 0x1FFF F80F	16B

Figure 4-2. Typical supply current consumption in Run mode



- $V_{dd} = 3.3V \Rightarrow 35\text{mA} @ 108\text{MHz} \Rightarrow 0.12\text{W}$

Cool Projects



Project 4

- Program in C and RISC-V on the Logan Nano
- Implement Pong Game
 - OR some better game if you want
- Call C functions from RISC-V!
 - Good – we do NOT need to program everything by hand on hardware – use provided C library!
- Lab 11: get familiar with the Logan Nano
- Lab 14 in week 18 (“2nd final week”): Demo and checkup of Project 4

We provide each project group with:

- 1 x Sipeed Logan Nano with screen and housing RMB 34.8
<https://item.taobao.com/item.htm?id=601743142093>
- 2 x push buttons <https://detail.tmall.com/item.htm?id=554574318222>



Sipeed Longan Nano RISC-V GD32VF103CBT6 单片机 开发板

价格 ¥34.80 388 65

优惠 淘金币可抵1.04元 累计评论 交易成功

配送 广东深圳至 上海浦东新区 快递 ¥12.00

颜色分类 C8 64KB FLASH 20KB RAM

CB 128KB FLASH 32KB RAM 0.96lcd+外壳

数量 1 件(库存52件)

立即购买 加入购物车

承诺 7天无理由

支付 集分宝

- You need to provide your own USB-C cable!
- You are free to buy your own hardware (e.g. your own Logan Nano, buttons, potentiometers, speaker!?) for Project 4...
- Code still goes to gitlab!

Links:

- <https://longan.sipeed.com/en/>
- <https://www.gigadevice.com/microcontroller/gd32vf103cbt6/>
- [https://github.com/nucleisys/Bumblebee Core Doc](https://github.com/nucleisys/Bumblebee_Core_Doc)
- <https://docs.platformio.org/en/latest/platforms/gd32v.html>
- Lab 11: <https://robotics.shanghaitech.edu.cn/courses/ca/20s/labs/11/>

Midterm

- Date: Tuesday, May. 26
- Time: 10:15- 12:15 (normal lecture slot++)
 - Be there latest 10:00 – we start 10:15 sharp!
- Venue: 4 rooms – check on egate which room you are!:
 - SPST1-503
 - SPST1-201
 - SPST1-501
 - SIST1A-106
- Closed book:
 - You can bring two A4 pages with notes (both sides; in English): Write your Chinese and **Pinyin** name on the top! Handwritten by you!
 - Final: you can bring **three** A4 pages
 - You will be provided with the RISC-V "green sheet"
 - No other material allowed!

Midterm I

- Wear your Corona mask! =>
- Switch cell phones **off!**
(not silent mode – off!)
 - Put them in your bags.
- Bags under the table. Nothing except paper, pen, 1 drink, 1 snack, your student ID card on the table!
- No other electronic devices are allowed!
 - No ear plugs, music, smartwatch...
- Anybody touching any electronic device will **FAIL** the course!
- Anybody found cheating (copy your neighbors answers, additional material, ...) will **FAIL** the course!



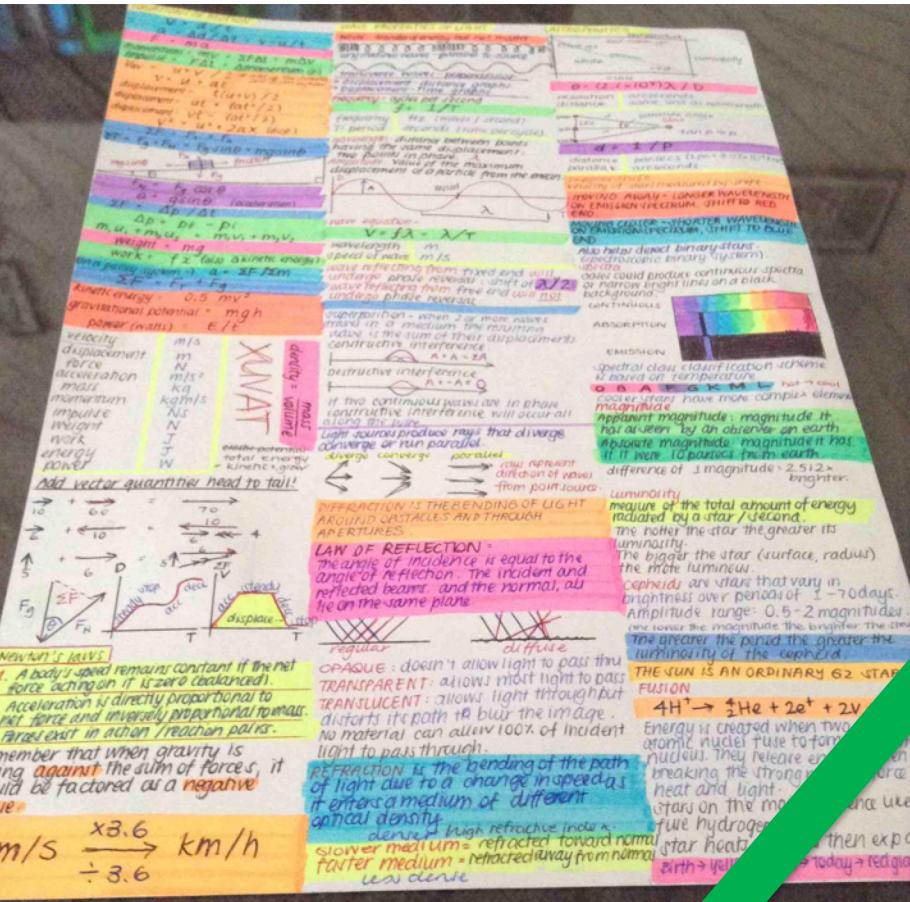


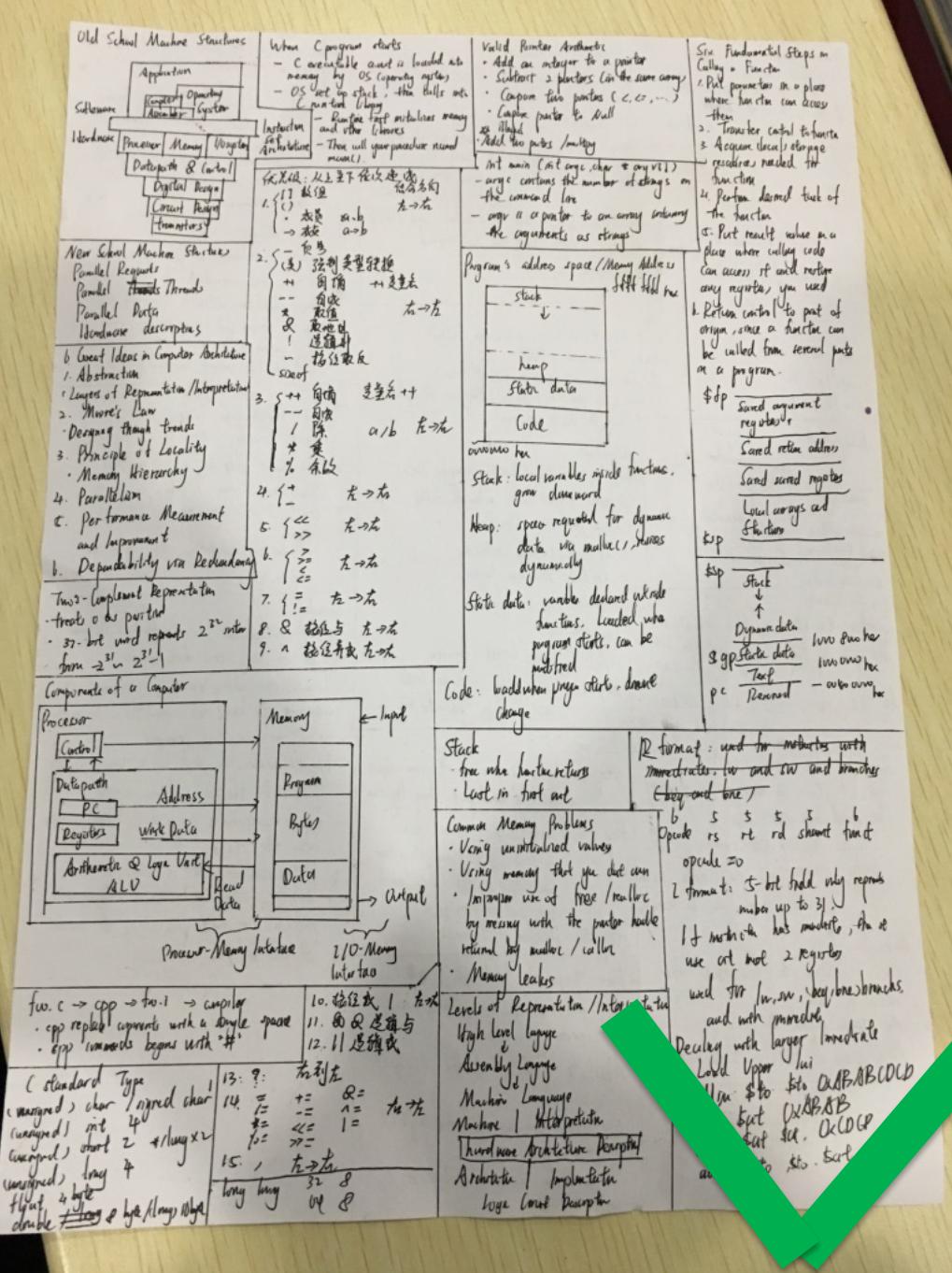






FUNCTIONS OF SEVERAL VARIABLES	$Z = f(x, y)$, $w = f(x, y, z)$	DOMAIN: ALLOWED (x, y) , (x, y, z) RANGES: 25.9° 'S
LEVEL CURVES (2-D)	$f(x_1, y_1) = k$ (CONST.)	$\text{f-DEFINITION OF CONTINUITY}$
CONTOURE MAPS (2-D)	$w = f(x_1, x_2, \dots, x_n)$	LET L BE A FUNCTION OF 2 VARIABLES DEFINED ON A DISK IN CENTER (a, b) , EXCEPT POSSIBLY AT (a, b) . THEN $\lim_{(x,y) \rightarrow (a,b)} f(x, y) = L$
WEIGHTS $f(x, y, z) = \text{CONST.}$	$\text{CONTINUOUS TO } f(x_1, x_2, \dots, x_n)$	IF FOR EVERY $\epsilon > 0$, THERE IS A CORRESPONDING $\delta > 0$ ST. $ f(x, y) - L < \epsilon$ whenever $\sqrt{(x-a)^2 + (y-b)^2} < \delta$
SURFACE LAYERS (3-D)	1. w AS A FUNCTION OF n REAL VARIABLES x_1, x_2, \dots, x_n	IF THE LIMIT AS A FUNCTION APPROACHED A POINT (a, b) ALONG TWO DIFFERENT PATHS IS NOT THE SAME, THE LIMIT DOES NOT EXIST. \square
PARTIAL DERIVATIVES	2. w AS A FUNCTION OF A SINGLE VARIABLE x .	$f(x, y)$ IS CONTINUOUS AT (a, b) IF THE LIMIT OF $f(x, y)$ AS $(x, y) \rightarrow (a, b)$ EXISTS.
$B = f(x, y)$	Derivatives w respect to two variables while holding the other variables constant	COMPOSITE FUNCTIONS OF CONTINUOUS FUNCTIONS ARE CONTINUOUS, AS ARE SUMS AND PRODUCTS
$f_x(x, y) = f_x = \frac{\partial f}{\partial x} = \frac{\partial f(x, y)}{\partial x}$	SAME HELDS FOR FUNCTIONS OF MORE THAN TWO VARIABLES	AND COMPOSITIONS OF TANGENT PLANES TO SURFACES
$f_y(x, y) = f_y = \frac{\partial f}{\partial y} = \frac{\partial f(x, y)}{\partial y}$		$z = f(x, y)$ AT (x_0, y_0, z_0) EVALUATED AT A POINT
SECOND PARTIAL DERIVATIVES	CLAIRAUT'S THEOREM	$z = f(x_0, y_0, z_0)$
$f_{xx} = \frac{\partial}{\partial x} \left(\frac{\partial f}{\partial x} \right) = \frac{\partial^2 f}{\partial x^2} = \frac{\partial^2 f}{\partial x^2}$	IF f_{xy} AND f_{yx} ARE BOTH CONTINUOUS	$z = f(x_0, y_0, z_0) + f_y(x_0, y_0) \Delta y + f_{yy}(x_0, y_0) \Delta y^2$
$f_{xy} = \frac{\partial}{\partial y} \left(\frac{\partial f}{\partial x} \right) = \frac{\partial^2 f}{\partial y \partial x} = \frac{\partial^2 f}{\partial y \partial x}$	$f_{xy} = f_{yx}$	TOTAL DIFFERENTIAL EQUATION: $dz = f'(x, y) dx + g'(x, y) dy$
$f_{yx} = \frac{\partial}{\partial x} \left(\frac{\partial f}{\partial y} \right) = \frac{\partial^2 f}{\partial x \partial y} = \frac{\partial^2 f}{\partial x \partial y}$	LAPLACE'S EQUATION: $\frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} = 0$ ETC...	$dz = f_x(x, y) dx + f_y(x, y) dy = \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy$
$f_{yy} = \frac{\partial}{\partial y} \left(\frac{\partial f}{\partial y} \right) = \frac{\partial^2 f}{\partial y^2} = \frac{\partial^2 f}{\partial y^2}$	EQUATION: $\frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} = 0$ ETC...	INCREMENT: $\Delta x, \Delta y, \Delta z$ DIFFERENTIALS, dx, dy, dz
THE CHAIN RULE	SINGLE VARIABLE: $y = g(t)$, $x = g(t)$, $10 \leq g(t) \leq 15$	FOR SMALL BOX, $\Delta x \approx dx$, $\Delta y \approx dy$
	$y'(t) = f'(g(t)) \cdot g'(t)$	IF f AND g ARE CONTINUOUS, $\Delta z \approx dz$
CASE 1: $z = f(x, y)$, $x = g(t)$, $y = h(t)$	$dz = \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy$	(Δz CHANGE IN HEIGHT OF SURFACE (z, x)) \approx CHANG IN HEIGHT OF THE TANGENT PLANE (z, x)
$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$ OR $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial t} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$	IF SOME PARTIAL	$\Delta z = f(g(t_0) + \Delta t) - f(g, t_0)$ THEOREM
CASE 2: $z = f(x, y)$, $x = g(t, s)$, $y = h(t, s)$, $t = f(g(t, s), h(s))$	$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial t} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial s} \frac{\partial x}{\partial s}$	$\Delta z = f_x(g(t_0, s_0) + \Delta t, h(s_0) + \Delta s) + f_y(g(t_0, s_0) + \Delta t, h(s_0) + \Delta s) \Delta t$ WHERE t_0 AND s_0 ARE FUNCTIONS OF dx AND dy FROM APPROACH O AS $(dx, dy) \rightarrow (0, 0)$
$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$	THREE APPROXIMATE CAN USE SLOPES OF SURFACE, BUT DON'T DO THAT	IF t IS DIFFERENTIABLE AT (x, y) FOR $z = f(x, y)$ DEF.
CHAIN RULE: GENERAL VERSION: $z = f(x_1, \dots, x_n)$	$x = g_1(t_1, \dots, t_m)$	DEPENDENCY DIAGRAMS CASE 1 & 2
$\frac{\partial f}{\partial x_i} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t_1} + \dots + \frac{\partial f}{\partial x_n} \frac{\partial x_n}{\partial t_m}$	$y = g_2(t_1, \dots, t_m)$	YOU CAN FIND DERIVATIVES FOR ALL THE FUNDAMENTAL FUNCTIONS AND INDEPENDENT VARIABLES
IMPLICIT DIFFERENTIATION	You can always solve for y and diff. w.r.t. x	$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_1} + \dots + \frac{\partial z}{\partial x_m} \frac{\partial x_m}{\partial t_m}$ MEANING LINE AND dy/dx
$\frac{\partial f}{\partial x} = -\frac{F_x}{F_y}$	$\frac{\partial z}{\partial y} = -\frac{F_x}{F_y}$	THE GRADIENT VECTOR: $z = f(x, y)$ AT P
$F(x, y) = 0$ FOR (x, y) , $F_x, F_y \neq 0$	$\frac{\partial z}{\partial x} = -\frac{F_y}{F_x}$	$\nabla f(x, y) = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right) = \left(f_x, f_y \right) = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$
TANGENT PLANE TO A LEVEL SURFACE	$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_1} + \dots + \frac{\partial z}{\partial x_m} \frac{\partial x_m}{\partial t_m}$	DIRECTIONAL DERIVATIVES D_f , $\frac{\partial f}{\partial u}$
$F_x(x-y) + F_y(y-x) = F_0$	$\frac{\partial z}{\partial y} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial t_1} + \dots + \frac{\partial z}{\partial y_m} \frac{\partial y_m}{\partial t_m}$	$D_f f(x, y) = f_x(x, y) + f_y(x, y)$ SAME AS $\frac{\partial f}{\partial x}$
NORMAL LINE TO A LEVEL SURFACE	$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_1} + \dots + \frac{\partial z}{\partial x_m} \frac{\partial x_m}{\partial t_m}$	$D_f^u f(x, y) = \nabla f(x, y) \cdot u$ 3 VARIABLES
$F_x(x-y) + F_y(y-x) = F_0$	$\frac{\partial z}{\partial y} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial t_1} + \dots + \frac{\partial z}{\partial y_m} \frac{\partial y_m}{\partial t_m}$	$D_f^u f(x, y) = \nabla f(x, y) \cdot u$ WHEN f IS IN THE SAME DIR. AS u
SPECIAL CASE: $2 = f(x_1, y_1)$	$F(x_1, y_1, z) = F_0$	$D_f^u f(x, y) = \nabla f(x, y) \cdot u$ $u = f(x, y)$ (1, 0, 0)
$F(x_1, y_1, z) = F_0$	$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_1} + \dots + \frac{\partial z}{\partial x_m} \frac{\partial x_m}{\partial t_m}$	@ THE GRADIENT VECTOR POINTS IN THE DIRECTION OF STEEPEST ASCENT OR DESCENT (W/SURFACE)
THEN $F_x = -1$, $F_y = 1$, $F_z = 0$	$\frac{\partial z}{\partial y} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial t_1} + \dots + \frac{\partial z}{\partial y_m} \frac{\partial y_m}{\partial t_m}$	@ THE GRADIENT VECTOR IS ORTHOGONAL TO THE LEVEL CURVES OF A SURFACE
MATRIX MAXIMUM AND MINIMUM VALUES: $Z = f(x, y)$	$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_1} + \dots + \frac{\partial z}{\partial x_m} \frac{\partial x_m}{\partial t_m}$	@ f HAS AS MANY COMPONENTS, AS f HAS INDEPENDENT VARIABLES. N^m (f_1, f_2, \dots, f_N)
$f(x, y) = 0$ $f_y(x, y) = 0$ $\nabla f(x, y) = 0$ IS NECESSARY BUT NOT SUFFICIENT TO GUARANTEE A MAX. OR MIN.	$\frac{\partial z}{\partial y} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial t_1} + \dots + \frac{\partial z}{\partial y_m} \frac{\partial y_m}{\partial t_m}$	@ TO FIND THE MAXIMAL (AND LATER) TANGENT PLANE TO A SURFACE, LET THAT SURFACE BE THE LEVEL SET OF SOME HIGHER DIMENSIONAL FUNCTION, THEN THE GRADIENT OF THE HIGHER f FUNCTION IS \perp TO YOUR SURFACE
$\text{SET } f_x = 0, f_y = 0, \nabla f = 0$ TO FIND CRITICAL PTS. (ALWAYS DOUBLE CHECK THIS)	$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_1} + \dots + \frac{\partial z}{\partial x_m} \frac{\partial x_m}{\partial t_m}$	$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_1} + \dots + \frac{\partial z}{\partial x_m} \frac{\partial x_m}{\partial t_m}$ LET $W^2 = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$ THE LEVEL SET $W^2 = 0$
THEN APPLY THE 2ND DERIVATIVE TEST: $\text{POSITIVE } \frac{\partial^2 f}{\partial x^2} \text{ AND } \frac{\partial^2 f}{\partial y^2} \text{ AND } \frac{\partial^2 f}{\partial x \partial y} < 0$	$\frac{\partial z}{\partial y} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial t_1} + \dots + \frac{\partial z}{\partial y_m} \frac{\partial y_m}{\partial t_m}$	$\text{SO } \nabla f \cdot \nabla f = (f_x, f_y) \cdot (f_x, f_y) = 1$ IS A NORMAL VECTOR TO THE $\frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} + 2 \frac{\partial^2 f}{\partial x \partial y} = 2 = 1$
FINDING ABSOLUTE MAX. AND. MIN. FOR f ON A CLOSED, BOUNDED		
1. Find values of f at the critical points of f in D		
2. Find the extreme values of f at the boundary of D		
3. The largest value from 1, 2, 3 is the ABS. MAX., the smallest is ABS. MIN.		
MATRIX MAXIMIZING AND MINIMIZING: Set f as a function of two variables of the form $z = f(x, y)$ and then Do the usual routine		





THE PURPOSE OF THIS REFERENCE GUIDE IS TO WALK THROUGH THE PROCESS OF BOOTING THE SIFT WORKSTATION, CREATING A TIMELINE ("SUPER" OR "MICRO") AND REVIEWING IT.



Content

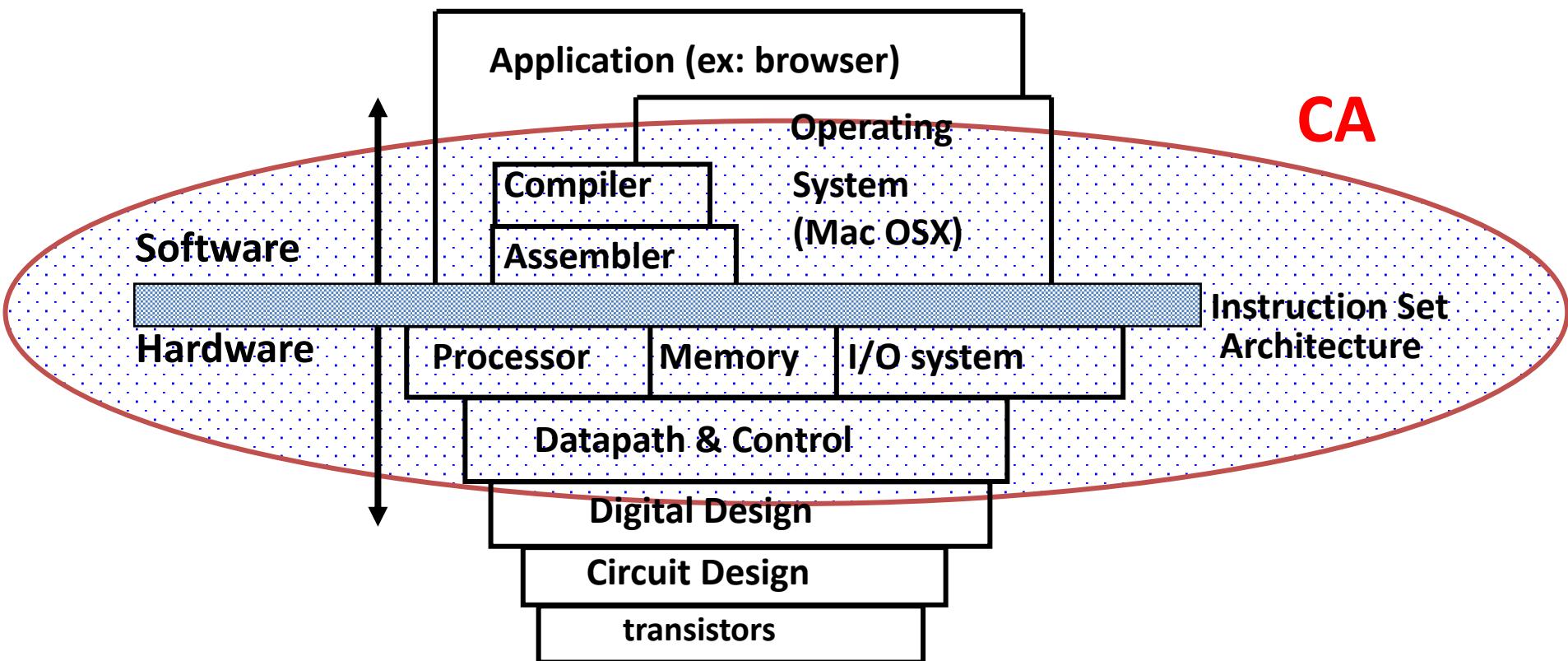
- Main topics: Everything till (including) Lecture 16
 - Number representation (int & float (Lecture 17!))
 - C
 - CALL
 - RISC-V
 - SDS; Datapath & Control
 - Pipelining & Superscalar
 - Caches
- Plus general "Computer Architecture" knowledge
- Disclaimer: In this review, important topics for CA are covered. It does not indicate that other topics from lectures 1-16 will not be covered in the exams, nor does it mean that everything written here will be covered.

New School Computer Architecture (1/3)

Personal
Mobile
Devices

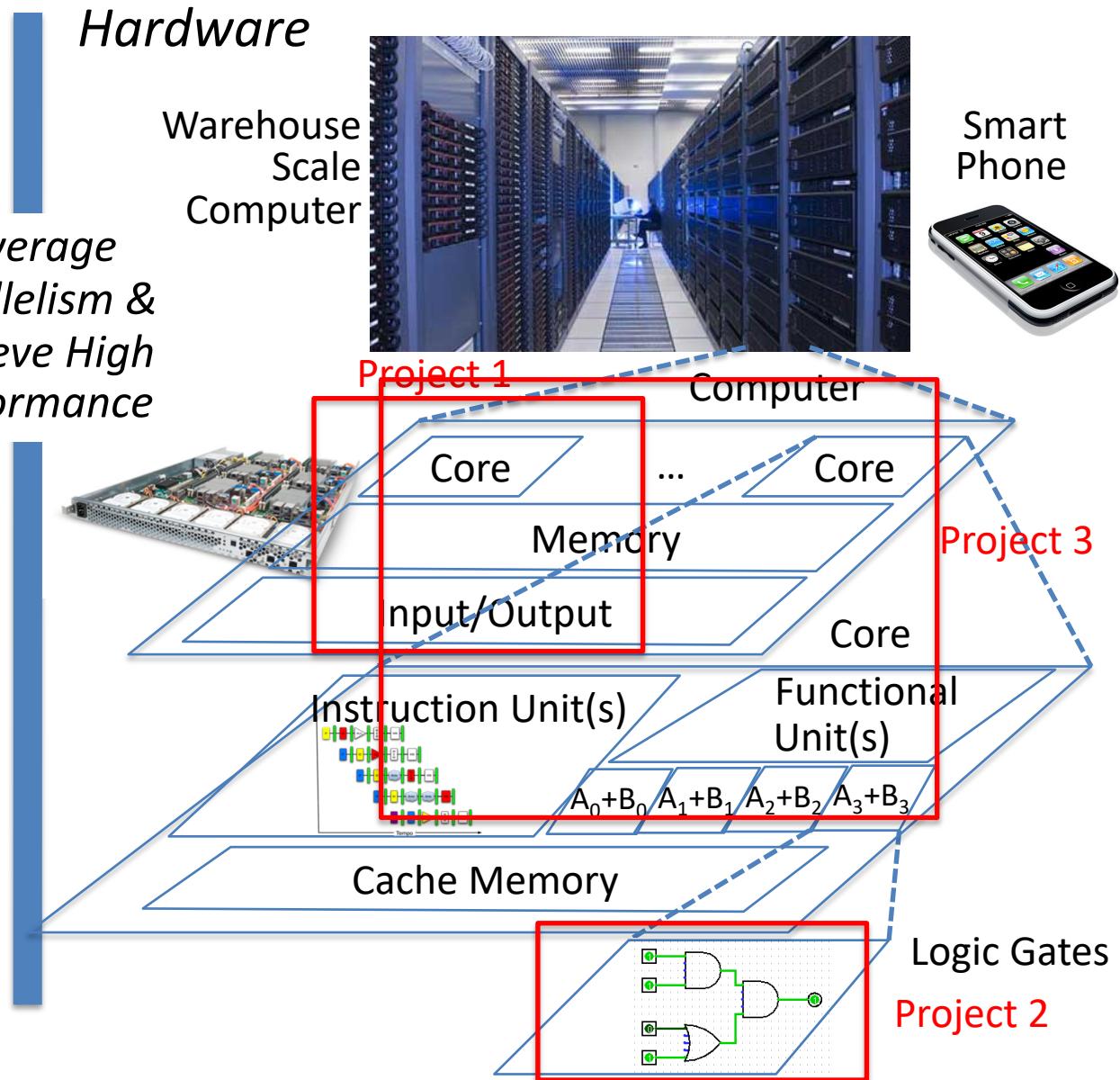


Old Machine Structures



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

- | Software | Hardware |
|--|---|
| • Parallel Requests
Assigned to computer
e.g., Search "Katz" | Warehouse Scale Computer |
| • Parallel Threads
Assigned to core
e.g., Lookup, Ads | Leverage Parallelism & Achieve High Performance |
| • Parallel Instructions
>1 instruction @ one time
e.g., 5 pipelined instructions | Computer |
| • Parallel Data
>1 data item @ one time
e.g., Add of 4 pairs of words | Core |
| • Hardware descriptions
All gates functioning in parallel at same time | Memory |
| • Programming Languages | Input/Output |



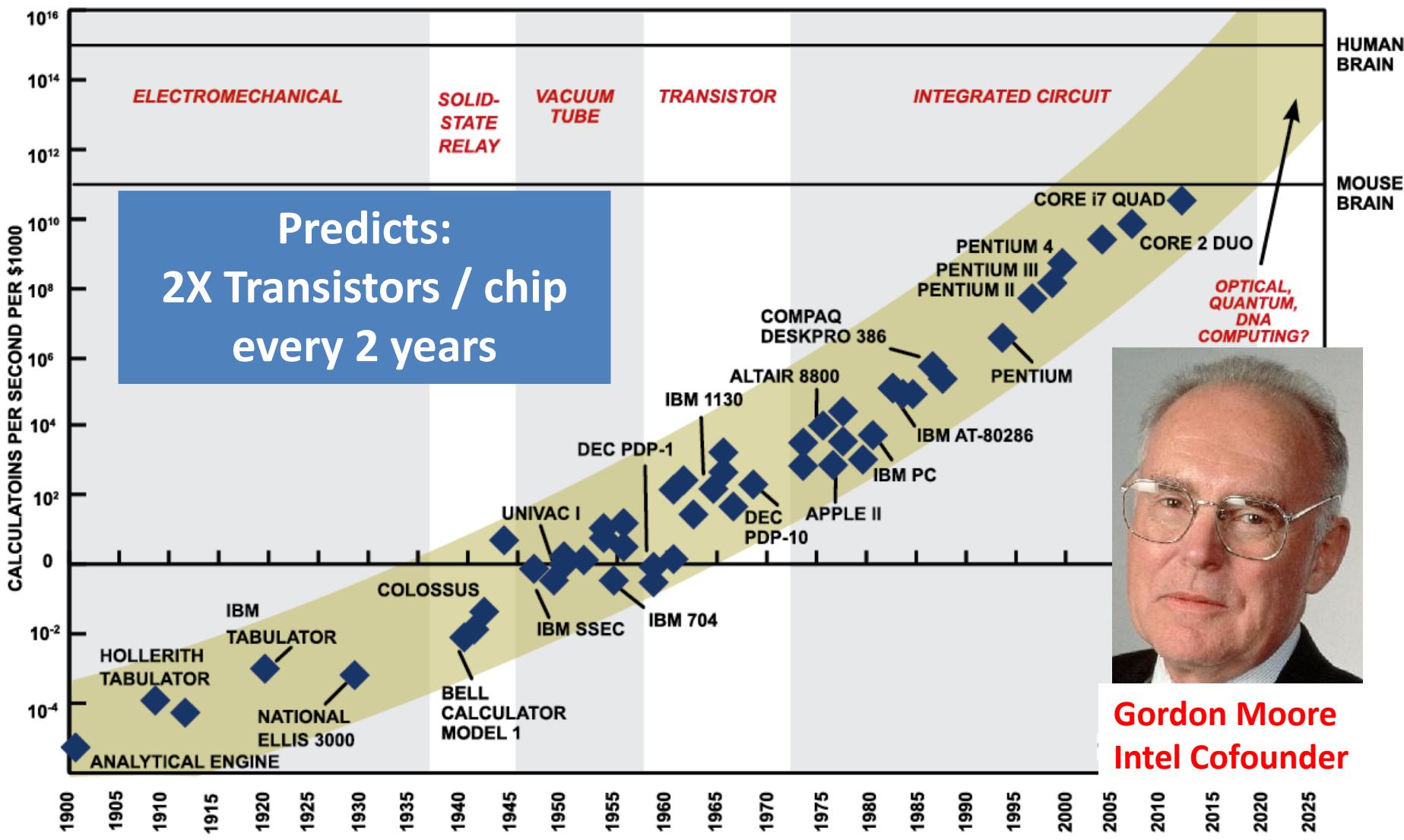
CA is NOT about C Programming

- It's about the hardware-software interface
 - What does the programmer need to know to achieve the highest possible performance
- Languages like C are closer to the underlying hardware, unlike languages like Python!
 - Allows us to talk about key hardware features in higher level terms
 - Allows programmer to explicitly harness underlying hardware parallelism for high performance: “programming for performance”

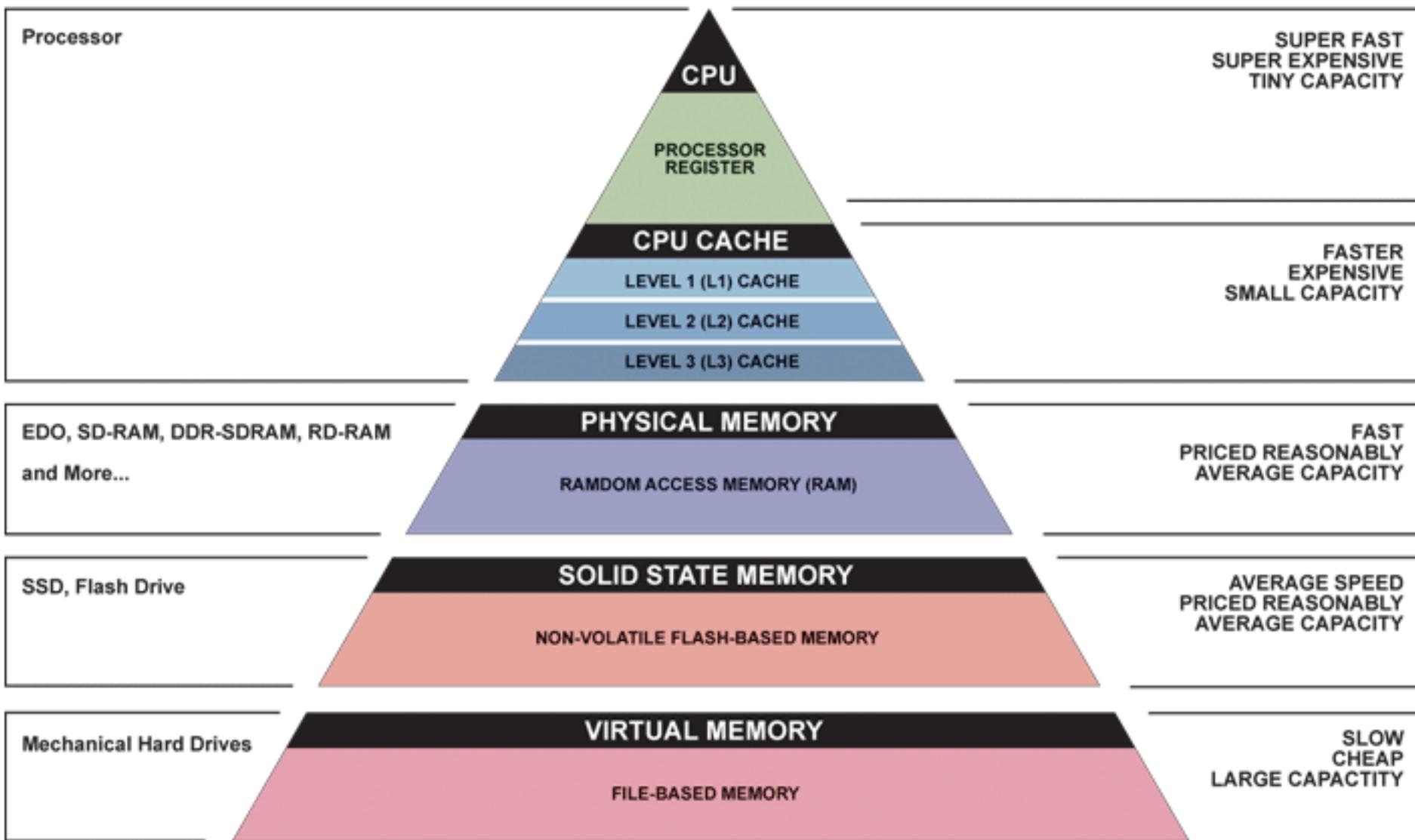
Great Ideas in Computer Architecture

1. *Design for Moore's Law*
 - *Higher capacities caches and DRAM*
2. Abstraction to Simplify Design
3. Make the Common Case Fast
4. *Dependability via Redundancy*
 - *Parity, SEC/DEC*
5. *Memory Hierarchy*
 - *Caches, TLBs*
6. *Performance via Parallelism/Pipelining/Prediction*
 - *Data-level Parallelism*

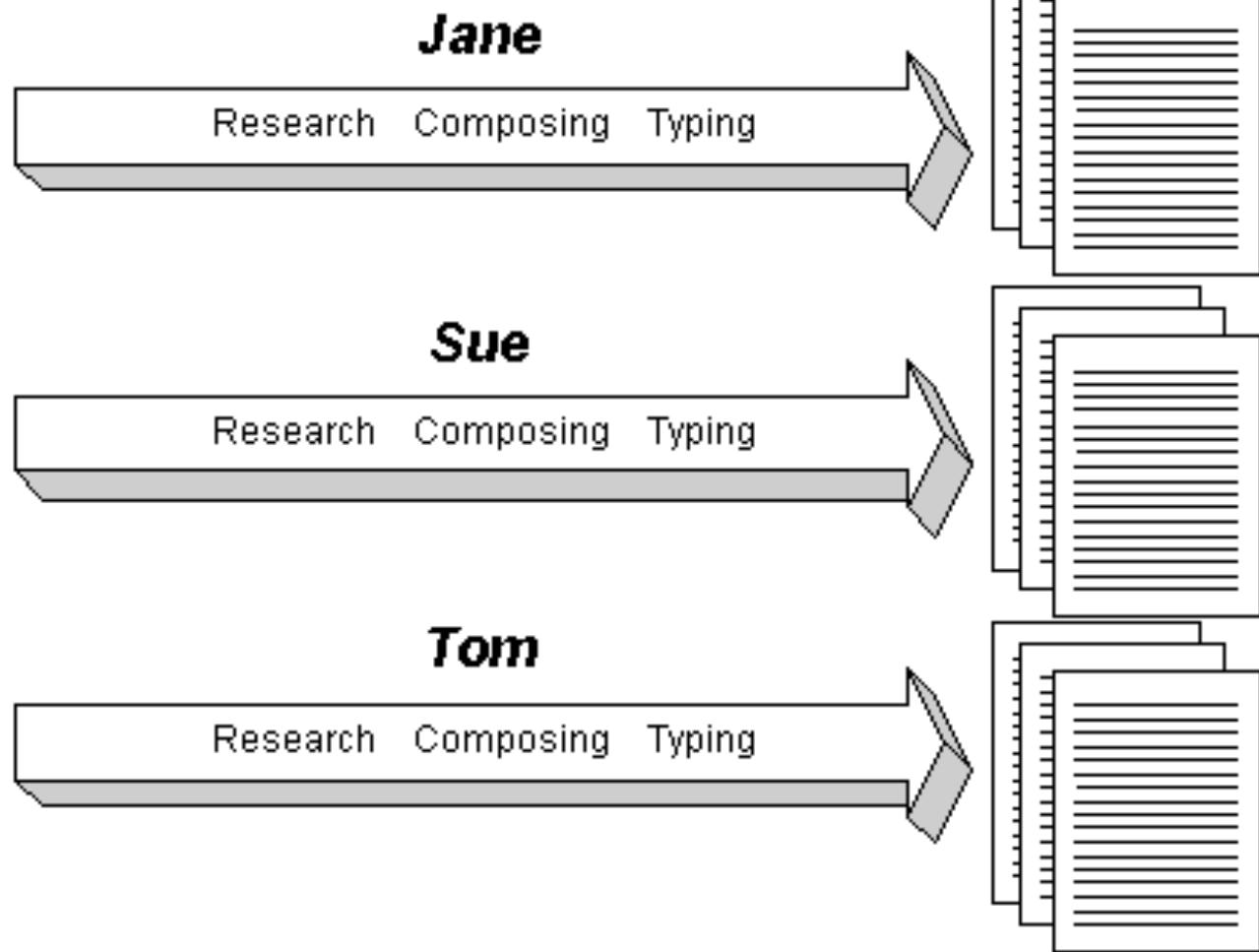
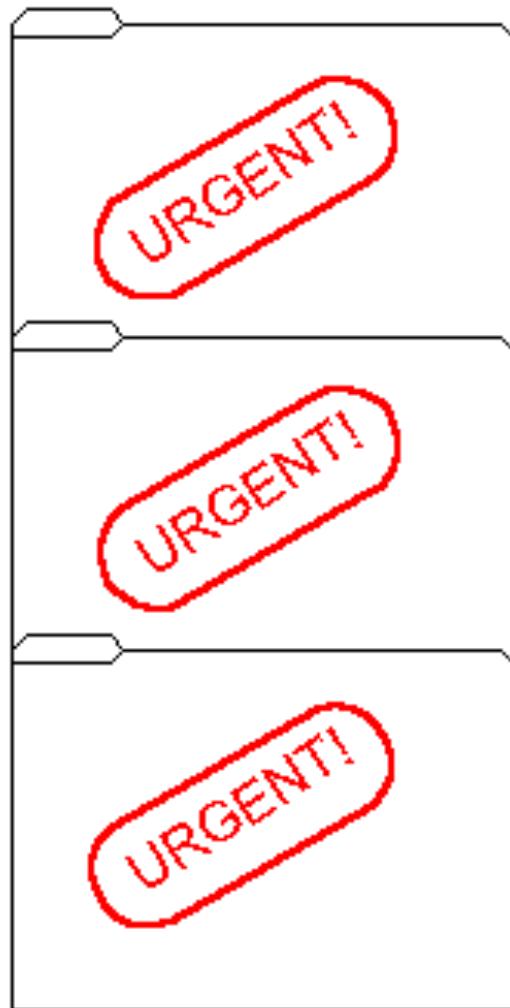
#2: Moore's Law



Great Idea #3: Principle of Locality/ Memory Hierarchy



Great Idea #4: Parallelism

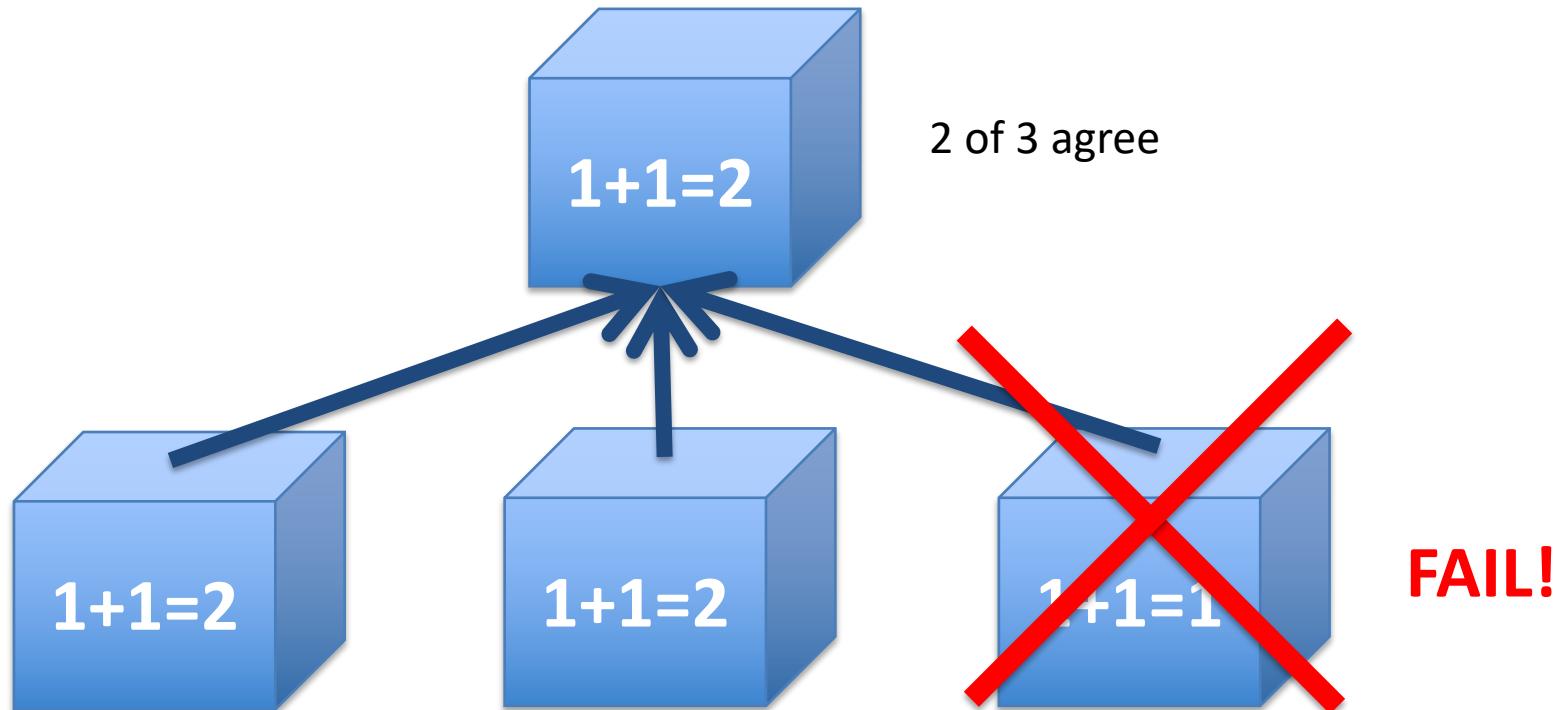


Great Idea #5: Performance Measurement and Improvement

- Tuning application to underlying hardware to exploit:
 - Locality
 - Parallelism
 - Special hardware features, like specialized instructions (e.g., matrix manipulation)
- Latency
 - How long to set the problem up
 - How much faster does it execute once it gets going
 - It is all about *time to finish*

Great Idea #6: Dependability via Redundancy

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



Increasing transistor density reduces the cost of redundancy

Key Concepts

- Inside computers, everything is a number
- But numbers usually stored with a fixed size
 - 8-bit bytes, 16-bit half words, 32-bit words, 64-bit double words, ...
- Integer and floating-point operations can lead to results too big/small to store within their representations: *overflow/underflow*

Number Representation

Signed Integers and Two's-Complement Representation

- Signed integers in C; want $\frac{1}{2}$ numbers < 0 , want $\frac{1}{2}$ numbers > 0 , and want one 0
- *Two's complement* treats 0 as positive, so 32-bit word represents 2^{32} integers from $-2^{31} (-2,147,483,648)$ to $2^{31}-1 (2,147,483,647)$
 - Note: one negative number with no positive version
 - Book lists some other options, all of which are worse
 - Every computer uses two's complement today
- *Most-significant bit* (leftmost) is the *sign bit*, since 0 means positive (including 0), 1 means negative
 - Bit 31 is most significant, bit 0 is least significant

Two's-Complement Integers

Sign Bit

$$0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2 = 0_{10}$$

$$0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0001_2 = 1_{10}$$

$$0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0010_2 = 2_{10}$$

...

...

$$0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1101_2 = 2,147,483,645_{10}$$

$$0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1110_2 = 2,147,483,646_{10}$$

$$0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111_2 = 2,147,483,647_{10}$$

$$1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2 = -2,147,483,648_{10}$$

$$1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0001_2 = -2,147,483,647_{10}$$

$$1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0010_2 = -2,147,483,646_{10}$$

...

...

$$1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1101_2 = -3_{10}$$

$$1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1110_2 = -2_{10}$$

$$1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111_2 = -1_{10}$$

Ways to Make Two's Complement

- For N-bit word, complement to 2_{ten}^N
 - For 4 bit number $3_{\text{ten}} = 0011_{\text{two}}$, two's complement (i.e. -3_{ten}) would be

$$16_{\text{ten}} - 3_{\text{ten}} = 13_{\text{ten}} \text{ or } 10000_{\text{two}} - 0011_{\text{two}} = 1101_{\text{two}}$$

- Here is an easier way:

- Invert all bits and add 1

3_{ten} 0011_{two}

Bitwise complement 1100_{two}

$$\begin{array}{r} + \\ \hline -3_{\text{ten}} & 1101_{\text{two}} \end{array}$$

- Computers actually do it like this, too

Two's-Complement Examples

- Assume for simplicity 4 bit width, -8 to +7 represented

$$\begin{array}{r} 3 \ 0011 \\ +2 \ 0010 \\ \hline 5 \ 0101 \end{array}$$

$$\begin{array}{r} 3 \ 0011 \\ +(-2) \ 1110 \\ \hline 1 \ 1\ 0001 \end{array}$$

$$\begin{array}{r} -3 \ 1101 \\ +(-2) \ 1110 \\ \hline -5 \ 1\ 1011 \end{array}$$

*Overflow when
magnitude of result
too big small to fit
into result
representation*

$$\begin{array}{r} 7 \ 0111 \\ +1 \ 0001 \\ \hline -8 \ 1000 \end{array}$$

$$\begin{array}{r} 7 \ 0111 \\ +(-1) \ 1111 \\ \hline +7 \ 1\ 0111 \end{array}$$

Carry into MSB =
Carry Out MSB

Overflow!

Overflow!

Carry into MSB ≠
Carry Out MSB

Carry in = carry from less significant bits

Carry out = carry to more significant bits

Suppose we had a 5-bit word. What integers can be represented in two's complement?

- 32 to +31
- 0 to +31
- 16 to +15
- 15 to +16

Suppose we had a 5-bit word. What integers can be represented in two's complement?

- 32 to +31
- 0 to +31
- 16 to +15
- 15 to +16

Conclusion

- Floating Point lets us:
 - Represent numbers containing both integer and fractional parts; makes efficient use of available bits.
 - Store **approximate** values for very large and very small #s.
 - **IEEE 754 Floating-Point Standard** is most widely accepted attempt to standardize interpretation of such numbers (Every desktop or server computer sold since ~1997 follows these conventions)

Can store NaN,
 $\pm \infty$

- ## • Summary (single precision):

31 30 23 22 0

S	Exponent	Significand
---	----------	-------------

- $(-1)^S \times (1 + \text{Significand}) \times 2^{(\text{Exponent}-127)}$

- Double precision identical, except with exponent bias of 1023 (half, quad similar)

Float: Special Numbers Summary

- Reserve exponents, significands:

Exponent	Significand	Object
0	0	0
0	<u>nonzero</u>	<u>Denorm</u>
1-254	anything	+/- fl. pt. #
255	0	+/- ∞
255	<u>nonzero</u>	<u>NaN</u>

C Programming

Quiz: Pointers

```
void foo(int *x, int *y)
{  int t;
   if ( *x > *y ) { t = *y; *y = *x; *x = t; }
}
int a=3, b=2, c=1;
foo(&a, &b);
foo(&b, &c);
foo(&a, &b);
printf("a=%d b=%d c=%d\n", a, b, c);
```

- Result is:
- A: a=3 b=2 c=1
 - B: a=1 b=2 c=3
 - C: a=1 b=3 c=2
 - D: a=3 b=3 c=3
 - E: a=1 b=1 c=1

Arrays and Pointers

```
int
foo(int array[],
     unsigned int size)
{
    ...
    printf("%d\n", sizeof(array));
}

int
main(void)
{
    int a[10], b[5];
    int c[] = {1, 3, 2, 5, 6};
    ... foo(a, 10) ... foo(c, 5) ...
    printf("%d\n", sizeof(c));
}
```

What does this print (64bit) ? 8

... because **array** is really
a pointer (and a pointer is
architecture dependent, but
8 on 64bit machines!)

What does this print? 40

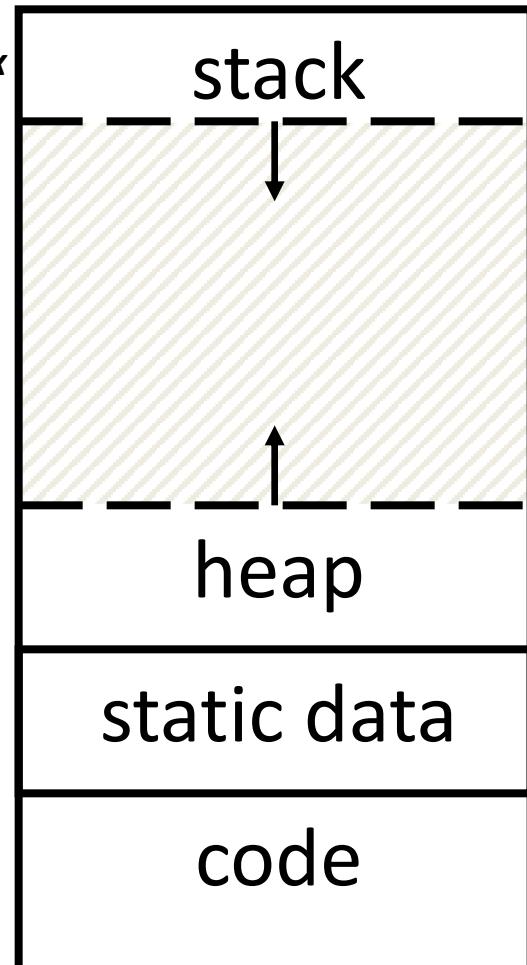
C Memory Management

- Program's *address space* contains 4 regions:
 - **stack**: local variables inside functions, grows downward
 - **heap**: space requested for dynamic data via `malloc()`; resizes dynamically, grows upward
 - **static data**: variables declared outside functions, does not grow or shrink. Loaded when program starts, can be modified.
 - **code**: loaded when program starts, does not change

Memory Address
(32 bits assumed here)

$\sim FFFF\ FFFF_{hex}$

$\sim 0000\ 0000_{hex}$



The Stack

- Every time a function is called, a new frame is allocated on the stack
- Stack frame includes:
 - Return address (who called me?)
 - Arguments
 - Space for local variables
- Stack frames contiguous blocks of memory; stack pointer indicates start of stack frame
- When function ends, stack frame is tossed off the stack; frees memory for future stack frames
- We'll cover details later for RISC-V processor

```
fooA() { fooB(); }
fooB() { fooC(); }
fooC() { fooD(); }
```

fooA frame

fooB frame

fooC frame

fooD frame

Stack Pointer →

Faulty Heap Management

- What is wrong with this code?
- Memory leak!

```
int foo() {  
    int *value = malloc(sizeof(int));  
    *value = 42;  
    return *value;  
}
```

And In Conclusion, ...

- Pointers are an abstraction of machine memory addresses
- Pointer variables are held in memory, and pointer values are just numbers that can be manipulated by software
- In C, close relationship between array names and pointers
- Pointers know the type of the object they point to (except void *)
- Pointers are powerful but potentially dangerous

RISC-V

Addition and Subtraction of Integers (3/4)

- How to do the following C statement?

`a = b + c + d - e; // a: x10; b: x1; c: x2, e: x3; f: x4`

- Break into multiple instructions

`add x10, x1, x2 # a_temp = b + c`

`add x10, x10, x3 # a_temp = a_temp + d`

`sub x10, x10, x4 # a = a_temp - e`

- Notice: A single line of C may break up into several lines of RISC-V.

- Notice: Everything after the hash mark on each line is ignored (comments).

Question:

We want to translate $*x = *y + 1$ into RISC-V
(x, y int pointers stored in: $s0\ s1$)

A: addi $s0, s1, 1$

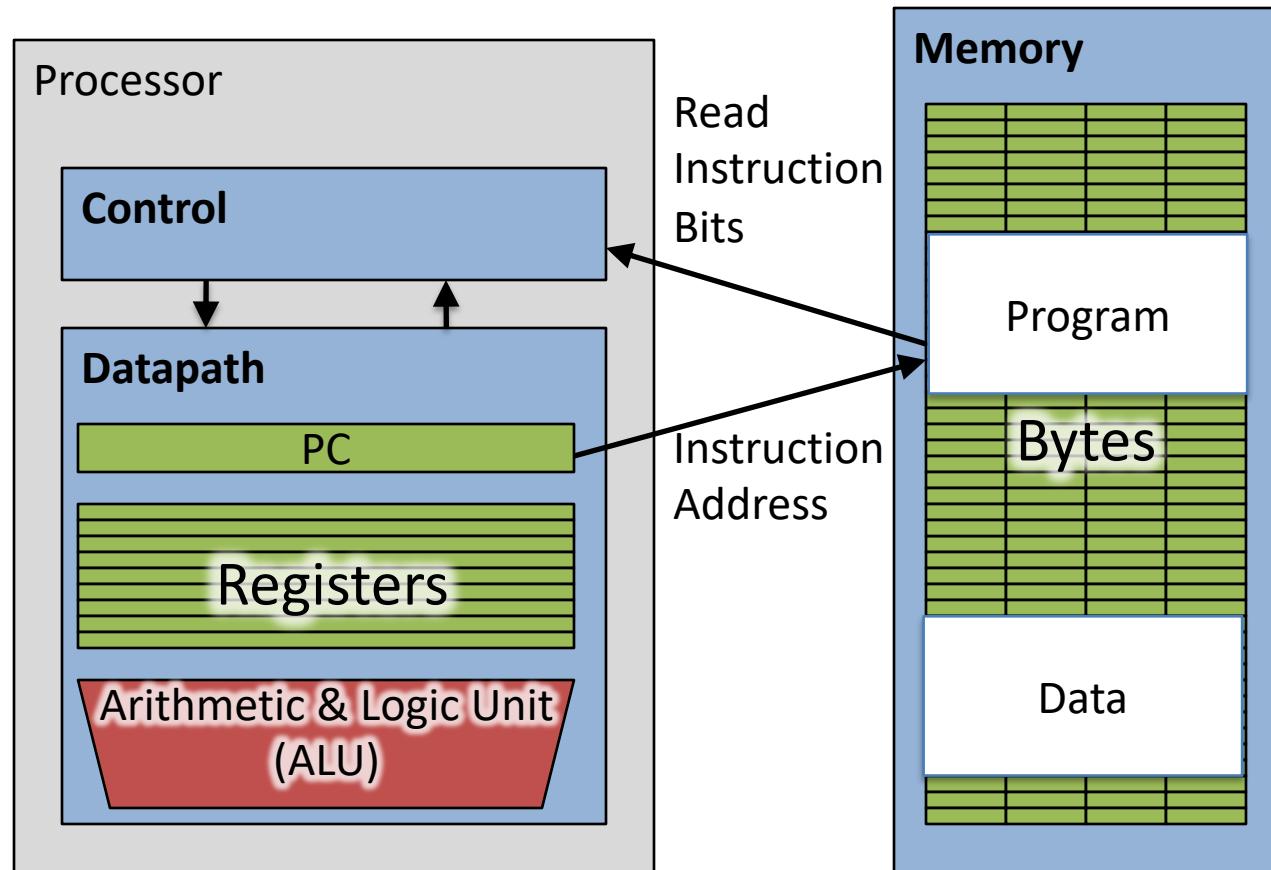
B: lw $s0$; sw $s1$; $1\{s1\}$ $0\{s0\}$

C: lw $t0, 0(s1)$
addi $t0, t0, 1$
sw $t0, 0(s0)$

D: sw $t0, 0(s1)$
addi $t0, t0, 1$
lw $t0, 0(s0)$

E: lw $s0, 1\{t0\}$
sw $s1, 0\{t0\}$

Executing a Program



- The PC (program counter) is internal register inside processor holding byte address of next instruction to be executed.
- Instruction is fetched from memory, then control unit executes instruction using datapath and memory system, and updates program counter (default is add +4 bytes to PC, to move to next sequential instruction)

Question!

```
          addi s0,zero,0  
Start:   slt  t0,s0,s1  
          beq  t0,zero,Exit  
          sll  t1,s0,2  
          add  t1,t1,s5  
          lw   t1,0(t1)  
          add  s4,s4,t1  
          addi s0,s0,1  
          j   Start
```

Exit:

- What is the code above?
- A: while loop
 - B: do ... while loop
 - C: for loop
 - D: A or C
 - E: Not a loop

Question

- What value does x12 have at the end?

- Answer:

$$x12 = 16$$

```
1 #include <stdio.h>
2
3 int main (){
4     int x10 = 7;
5     int x12 = 0;
6     do{
7         int x14 = x10 & 1;
8         if(x14)
9             x12 += x10;
10
11         x10--;
12     } while (x10 != 0);
13     printf("%d", x12);
14 }
```

```
addi  x10, x0 , 0x07
add   x12, x0 , x0
label_a:
andi  x14, x10, 1
beq   x14, x0 , label_b
add   x12, x10, x12
label_b:
addi  x10, x10, -1
bne   x10, x0 , label_a
```

RISC-V Function Call Conventions

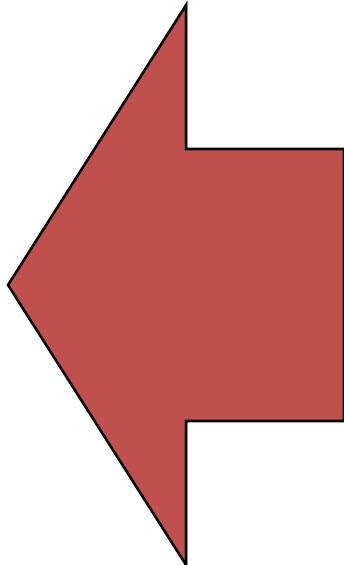
- Registers faster than memory, so use them
- Give names to registers, conventions on how to use them
- a₀-a₇ (x₁₀-x₁₇) : eight *argument* registers to pass parameters and return values (a₀-a₁)
- ra: one *return address* register to return to the point of origin (x₁)
- Also s₀-s₁ (x₈-x₉) and s₂-s₁₁ (x₁₈-x₂₇) : saved registers (more about those later)

Instruction Support for Functions (1/4)

```
... sum(a,b); ... /* a, b: s0, s1 */  
}  
int sum(int x, int y) {  
    return x+y;  
}
```

address (shown in decimal)

1000
1004
1008
1012
1016
...
2000
2004



In RV32, instructions are 4 bytes, and stored in memory just like data. So here we show the addresses of where the programs are stored.

Instruction Support for Functions (2/4)

C

```
... sum(a,b); ... /* a, b: s0, s1 */
}
int sum(int x, int y) {
    return x+y;
}
```

address (shown in decimal)

1000	add	a0, s0, x0	# x = a
1004	mv	a1, s1	# y = b
1008	addi	ra, zero, 1016	# ra=1016
1012	j	sum	# jump to sum
1016	...		# next instruction
...			
2000	sum:	add a0, a0, a1	
2004	jr	ra	# new instr. "jump register"

Instruction Support for Functions (3/4)

```
... sum(a,b); ... /* a,b:$s0,$s1 */  
}  
int sum(int x, int y) {  
    return x+y;  
}
```

- Question: Why use **jr** here? Why not use **j**?
- Answer: **sum** might be called by many places, so we can't return to a fixed place. The calling proc to **sum** must be able to say "return here" somehow.

2000 sum: add a0, a0, a1

2004 jr ra # new instr. "jump register"

Instruction Support for Functions (4/4)

- Single instruction to jump and save return address:
jump and link (**jal**)

- Before:

```
1008 addi ra, zero, 1016 # $ra=1016
1012 j sum                # goto sum
```

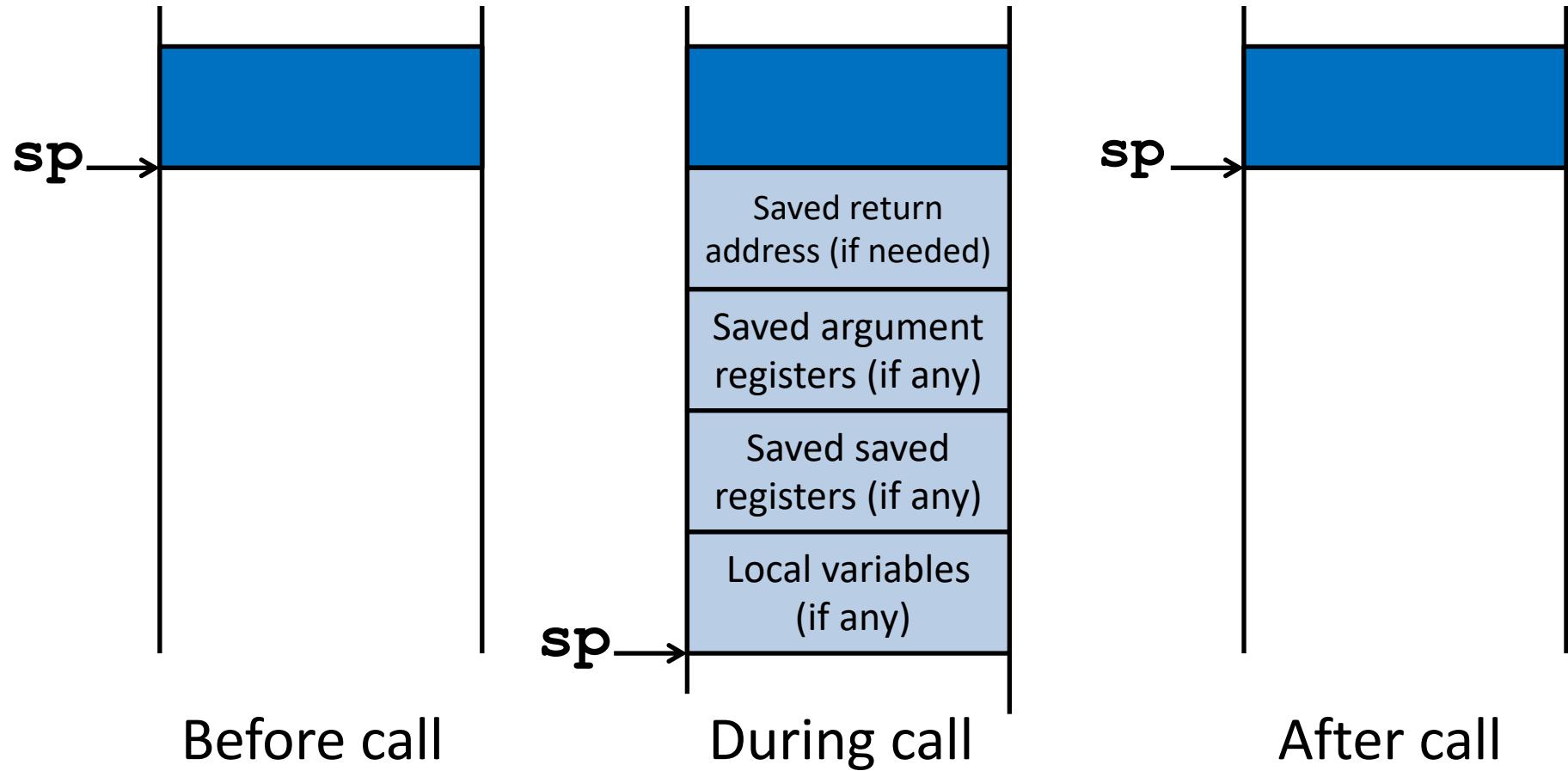
- After:

```
1008 jal sum    # ra=1012, goto sum
```

- Why have a **jal**?

- Make the common case fast: function calls very common.
- Reduce program size
- Don't have to know where code is in memory with **jal**!

Stack Before, During, After Function



Using the Stack (1/2)

- We have a register **sp** which always points to the last used space in the stack.
- To use stack, we decrement this pointer by the amount of space we need and then fill it with info.
- So, how do we compile this?

```
int sumSquare(int x, int y) {  
    return mult(x,x)+ y;  
}
```

Using the Stack (2/2)

```
int sumSquare(int x, int y) {  
    return mult(x,x)+ y; }
```

sumSquare:

“push” addi sp, sp, -8 # space on stack
 sw ra, 4(sp) # save ret addr
 sw a1, 0(sp) # save y
 mv a1, a0 # mult(x,x)
 jal mult # call mult
 lw a1, 0(sp) # restore y

“pop” add a0, a0, a1 # mult() + y
 lw ra, 4(sp) # get ret addr
 addi sp, sp, 8 # restore stack
 jr ra

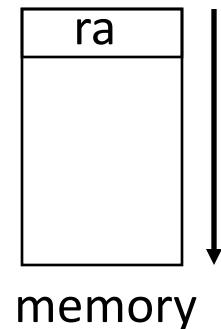
mult: ...

Basic Structure of a Function

Prologue

```
entry_label:  
addi sp,sp, -framesize  
sw ra, framesize-4(sp) # save ra  
save other regs if need be
```

Body ... (call other functions...)



Epilogue

```
restore other regs if need be  
lw ra, framesize-4(sp) # restore $ra  
addi sp, sp, framesize  
jr ra
```

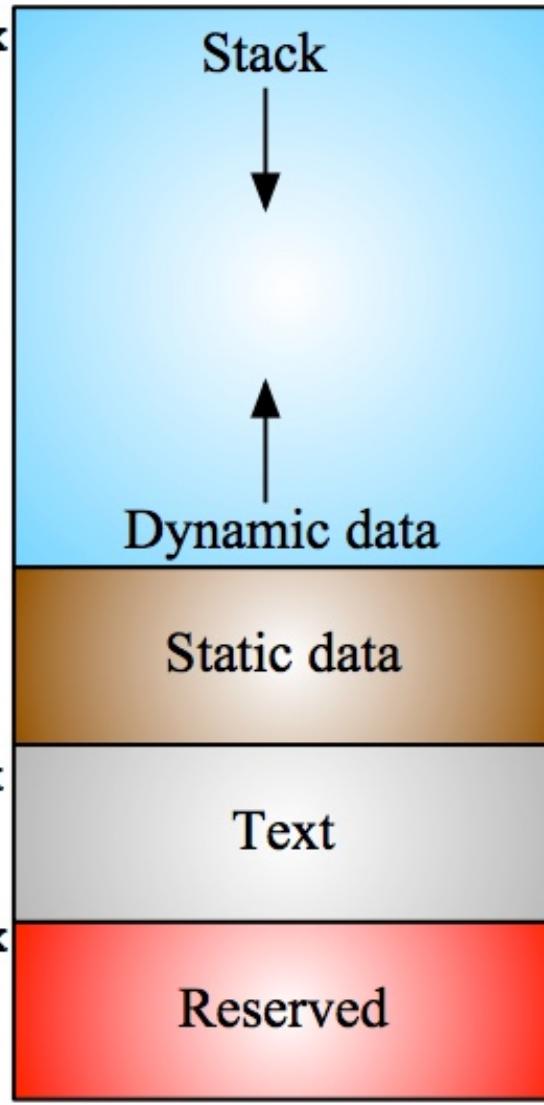
RV32 Memory Allocation

$sp = bfff\ fff0_{hex}$

$1000\ 0000_{hex}$

$pc = 0001\ 0000_{hex}$

0



RISC-V ISA so far...

- Registers we know so far (All of them!)
 - a0-a7 for function arguments, a0-a1 for return values
 - sp, stack pointer, ra return address
 - s0-s11 saved registers
 - t0-t6 temporaries
 - zero
- Instructions we know:
 - Arithmetic: add, addi, sub
 - Logical: sll, srl, slli, srli, slai, and, or, xor, andi, ori, xori
 - Decision: beq, bne, blt, bge
 - Unconditional branches (jumps): j, jr
 - Functions called with **jal**, return with **jr ra**.
- The stack is your friend: Use it to save anything you need. Just leave it the way you found it!

12 Shift Instructions...

- Two versions of all shift instructions. Shift amount via:
 - Register
 - Immediate
- (On RV64: additional “word” version of instruction: only works on first 32bit of 64bit register)
- Shift Left
- Shift Right Arithmetic: Fill upper bits with **msb**
- Shift Right Logic: Fill upper bits with 0’s

sll, sllw	R	Shift Left (Word)	$R[rd] = R[rs1] \ll R[rs2]$	1)
slli, slliw	I	Shift Left Immediate (Word)	$R[rd] = R[rs1] \ll \text{imm}$	1)
sra, sraw	R	Shift Right Arithmetic (Word)	$R[rd] = R[rs1] \gg R[rs2]$	1,5)
srai, sraiw	I	Shift Right Arith Imm (Word)	$R[rd] = R[rs1] \gg \text{imm}$	1,5)
srl, srlw	R	Shift Right (Word)	$R[rd] = R[rs1] \gg R[rs2]$	1)
srli, srliw	I	Shift Right Immediate (Word)	$R[rd] = R[rs1] \gg \text{imm}$	1)

Notes: 1) The Word version only operates on the rightmost 32 bits of a 64-bit registers
5) Replicates the sign bit to fill in the leftmost bits of the result during right shift

Summary of RISC-V Instruction Formats

31	30	25	24	21	20	19	15	14	12	11	8	7	6	0
funct7		rs2		rs1		funct3			rd		opcode			R-type
		imm[11:0]		rs1		funct3			rd		opcode			I-type
imm[11:5]		rs2		rs1		funct3		imm[4:0]			opcode			S-type
imm[12 10:5]		rs2		rs1		funct3	imm[4:1 11]				opcode			B-type
		imm[31:12]							rd		opcode			U-type
		imm[20 10:1 11]		imm[19:12]					rd		opcode			J-type

CORE INSTRUCTION FORMATS

	31	27	26	25	24	20	19	15	14	12	11	7	6	0
R		funct7		rs2		rs1		funct3		rd		Opcode		
I			imm[11:0]			rs1		funct3		rd		Opcode		
S		imm[11:5]		rs2		rs1		funct3		imm[4:0]		opcode		
SB		imm[12 10:5]		rs2		rs1		funct3		imm[4:1 11]		opcode		
U			imm[31:12]							rd		opcode		
UJ			imm[20 10:1 11 19:12]							rd		opcode		

Question

- What is correct encoding of add x4, x3, x2 ?
A: 4021 8233_{hex}
B: 0021 82b3_{hex}
C: 4021 82b3_{hex}
D: 0021 8233_{hex}
E: 0021 8234_{hex}

31	25 24	20 19	15 14	12 11	7 6	0	
0000000	rs2	rs1	000	rd	0110011		add
0100000	rs2	rs1	000	rd	0110011		sub
0000000	rs2	rs1	100	rd	0110011		xor
0000000	rs2	rs1	110	rd	0110011		or
0000000	rs2	rs1	111	rd	0110011		and

Branch Example, Determine Offset

- RISC-V Code:

```
Loop: beq    x19,x10,End  
      add    x18,x18,x10  
      addi   x19,x19,-1  
      j      Loop  
End: # target instruction
```

A vertical list of numbers 0, 1, 2, 3, 4. A curly brace on the left groups the first four numbers. An arrow points from the number 4 to the word "Count" on the right, which is followed by the text "instructions from branch".

- Branch offset = **4×32-bit instructions = 16 bytes**
- (Branch with offset of 0, branches to itself)

Branch Example, Determine Offset

- RISC-V Code:

Loop: **beq** **x19,x10,End**
add x18,x18,x10
addi x19,x19,-1
j Loop
End: # target instruction

0 Count
1 instructions
2 from branch
3
4



???????	01010	10011	000	?????	1100011
imm	rs2=10	rs1=19	BEQ	imm	BRANCH

Branch Example, Encode Offset

- RISC-V Code:

Loop: **beq** **x19,x10,End**
add x18,x18,x10
addi x19,x19,-1
j Loop

End: # target instruction



offset = 16 bytes = 8x2 bytes

???????	01010	10011	000	?????	1100011
imm	rs2=10	rs1=19	BEQ	imm	BRANCH

RISC-V Immediate Encoding

Instruction encodings, $\text{inst}[31:0]$

31	30	25 24	20 19	15 14	12 11	8 7 6	0	
		funct7	rs2	rs1	funct3	rd	opcode	R-type
		imm[11:0]		rs1	funct3	rd	opcode	I-type
		imm[11:5]	rs2	rs1	funct3	imm[4:0]	opcode	S-type
		imm[12 10:5]	rs2	rs1	funct3	imm[4:1 11]	opcode	B-type

32-bit immediates produced, $\text{imm}[31:0]$

31	25 24	12	11	10	5	4	1	0	
-inst[31]-				inst[30:25]	inst[24:21]	inst[20]			I-imm.
-inst[31]-				inst[30:25]	inst[11:8] 	inst[7]			S-imm.
-inst[31]- 	inst[7]	inst[30:25]	inst[11:8]				0		B-imm.

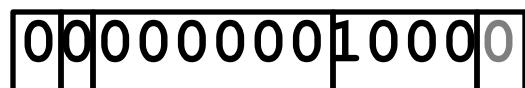
Upper bits sign-extended from $\text{inst}[31]$ always

Only bit 7 of instruction changes role in immediate between S and B

Branch Example, complete encoding

beq x19,x10, offset = 16 bytes

13-bit immediate, imm[12:0], with value 16



imm[0] discarded,
always zero

imm[12]



imm[10:5] rs2=10

rs1=19

BEQ

imm[4:1]

BRANCH

LUI to Create Long Immediates

- LUI writes the upper 20 bits of the destination with the immediate value, and clears the lower 12 bits.
- Together with an ADDI to set low 12 bits, can create any 32-bit value in a register using two instructions (LUI/ADDI).

```
LUI x10, 0x87654      # x10 = 0x87654000  
ADDI x10, x10, 0x321# x10 = 0x87654321
```

One Corner Case

How to set 0xDEADBEEF?

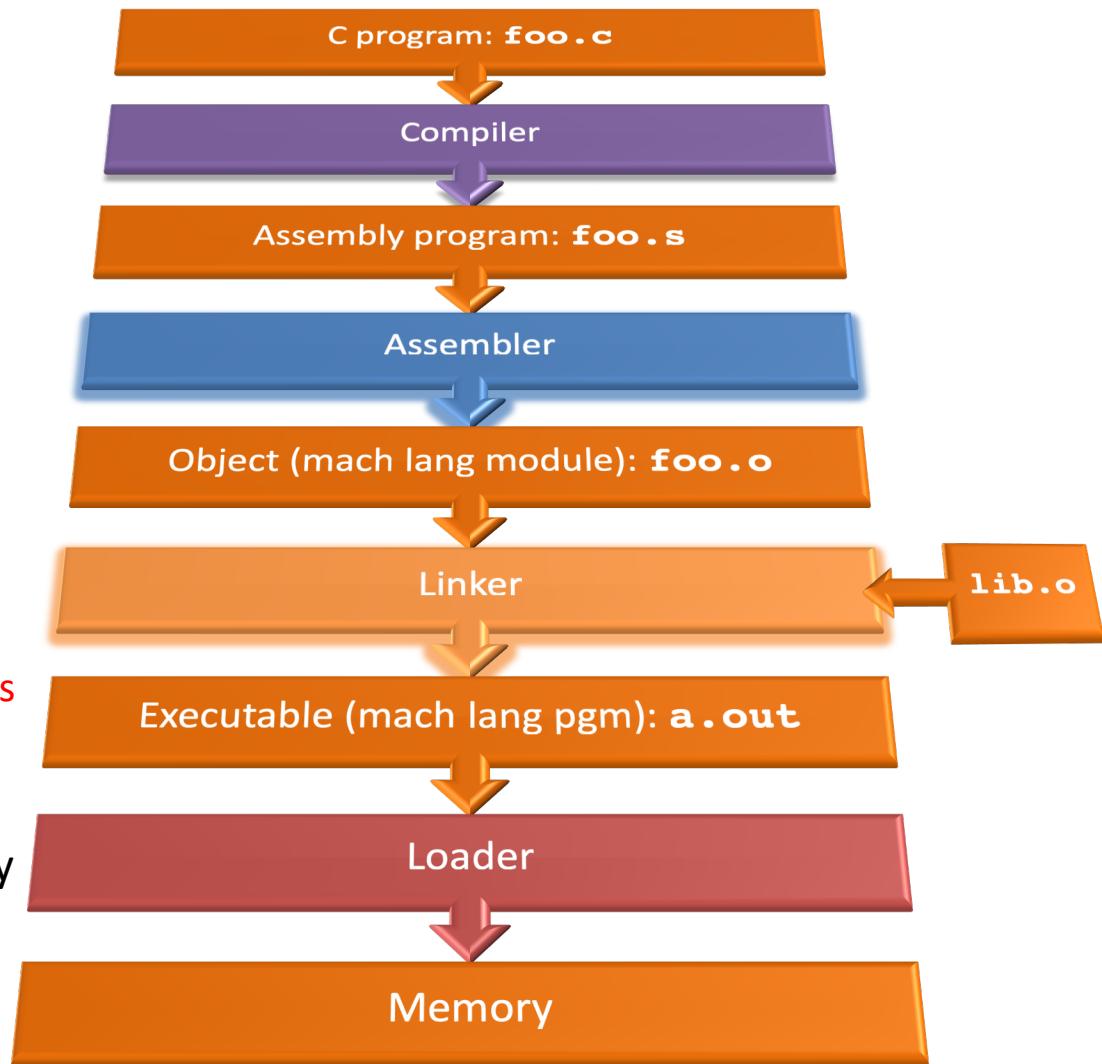
```
LUI x10, 0xDEADB          # x10 = 0xDEADB000
```

```
ADDI x10, x10, 0xEEF    # x10 = 0xDEADAEEF
```

ADDI 12-bit immediate is always sign-extended, if top bit is set,
will subtract 1 from upper 20 bits

Steps in compiling a C program

- 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.



Question

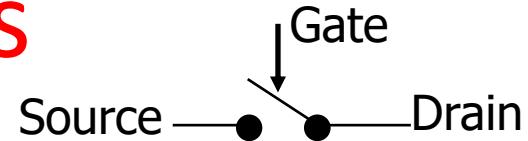
At what point in process are all the machine code bits generated for the following assembly instructions:

- 1) add 6, 7, 8
- 2) jal fprintf

- A: 1) & 2) After compilation
- B: 1) After compilation, 2) After assembly
- C: 1) After assembly, 2) After linking
- D: 1) After assembly, 2) After loading
- E: 1) After compilation, 2) After linking

SDS/ RISC-V Pipeline

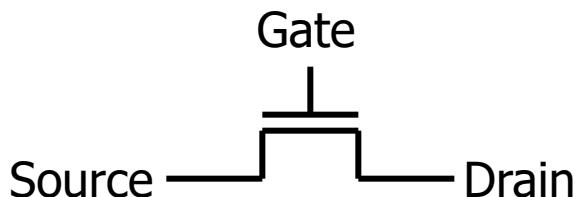
CMOS Transistors



- Three terminals: source, gate, and drain

- Switch action:

if voltage on gate terminal is (some amount) higher/lower than source terminal then conducting path established between drain and source terminals (switch is closed)



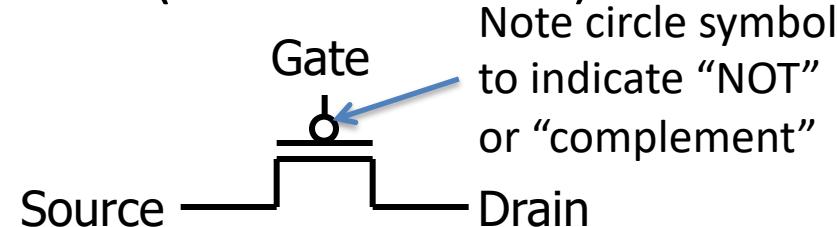
n-channel transistor

off when voltage at Gate is low

on when:

voltage (Gate) > voltage (Threshold)

(**High** resistance when gate voltage **Low**, **Low** resistance when gate voltage **High**)



p-channel transistor

on when voltage at Gate is low

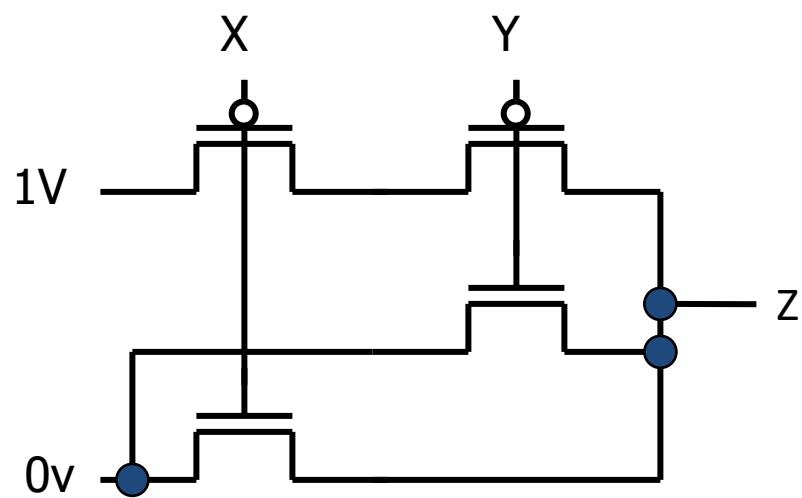
off when:

voltage (Gate) > voltage (Threshold)

(**Low** resistance when gate voltage **Low**, **High** resistance when gate voltage **High**)

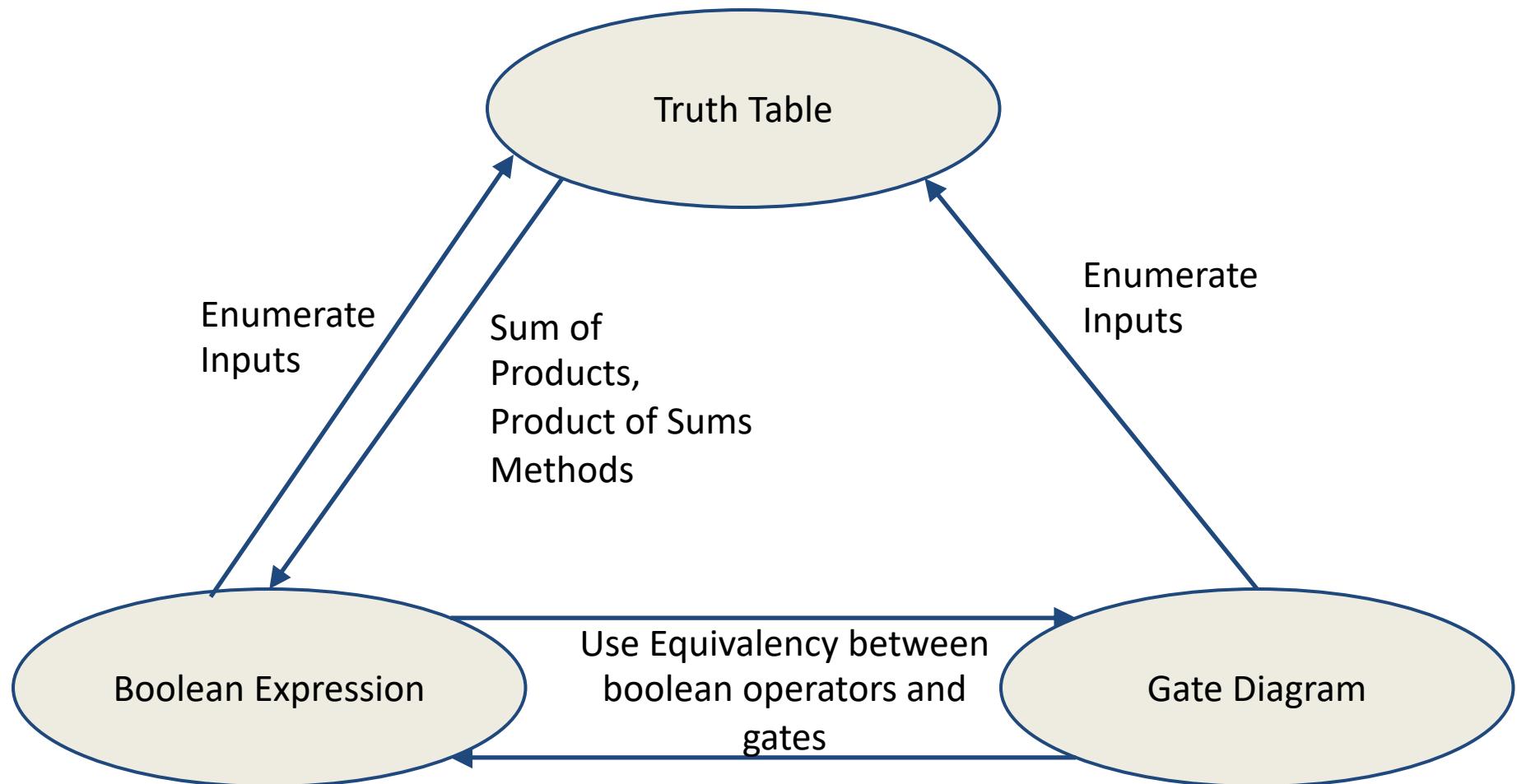
Field-Effect Transistor (FET) => CMOS circuits use a combination of p-type and n-type metal-oxide-semiconductor field-effect transistors =>
MOSFET

Question



X	Y	Z		
A	B	C	D	Volts
0 Volt	0 Volt	0 Volt	1	Volts
0 Volt	1 Volt	0 Volt	1	Volts
1 Volt	0 Volt	0 Volt	1	Volts
1 Volt	1 Volt	0 Volt	0	Volts

Representations of Combinational Logic (groups of logic gates)



Laws of Boolean Algebra

$$X \bar{X} = 0$$

$$X 0 = 0$$

$$X 1 = X$$

$$X X = X$$

$$X Y = Y X$$

$$(X Y) Z = X (Y Z)$$

$$X (Y + Z) = X Y + X Z$$

$$X Y + X = X$$

$$\bar{X} Y + X = X + Y$$

$$\overline{XY} = \bar{X} + \bar{Y}$$

$$X + \bar{X} = 1$$

$$X + 1 = 1$$

$$X + 0 = X$$

$$X + X = X$$

$$X + Y = Y + X$$

$$(X + Y) + Z = X + (Y + Z)$$

$$X + Y Z = (X + Y) (X + Z)$$

$$(X + Y) X = X$$

$$(\bar{X} + Y) X = X Y$$

$$\overline{X + Y} = \bar{X} \bar{Y}$$

Complementarity

Laws of 0's and 1's

Identities

Idempotent Laws

Commutativity

Associativity

Distribution

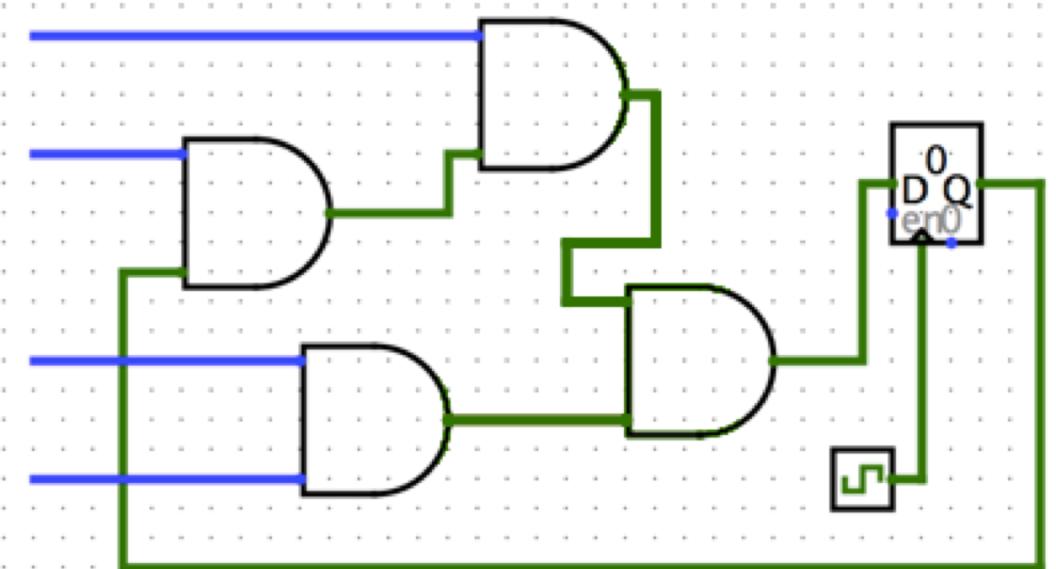
Uniting Theorem

Uniting Theorem v. 2

DeMorgan's Law

Question

Piazza: "Lecture 9 Freq Poll"



Clock->Q 1ns
Setup 1ns
Hold 1ns
AND delay 1ns

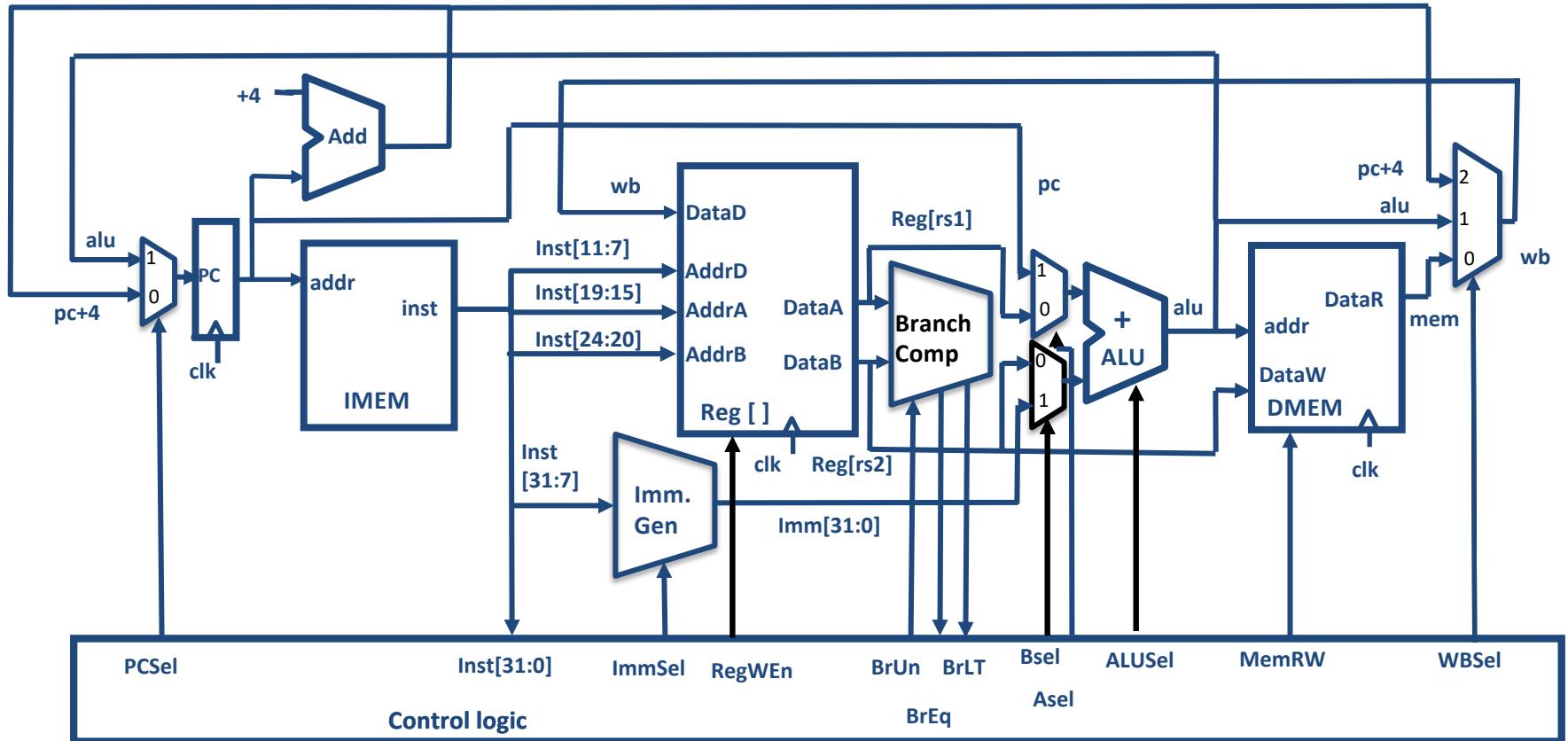
What is maximum clock frequency?

- A: 5 GHz
- B: 500 MHz
- C: 200 MHz
- D: 250 MHz
- E: 1/6 GHz

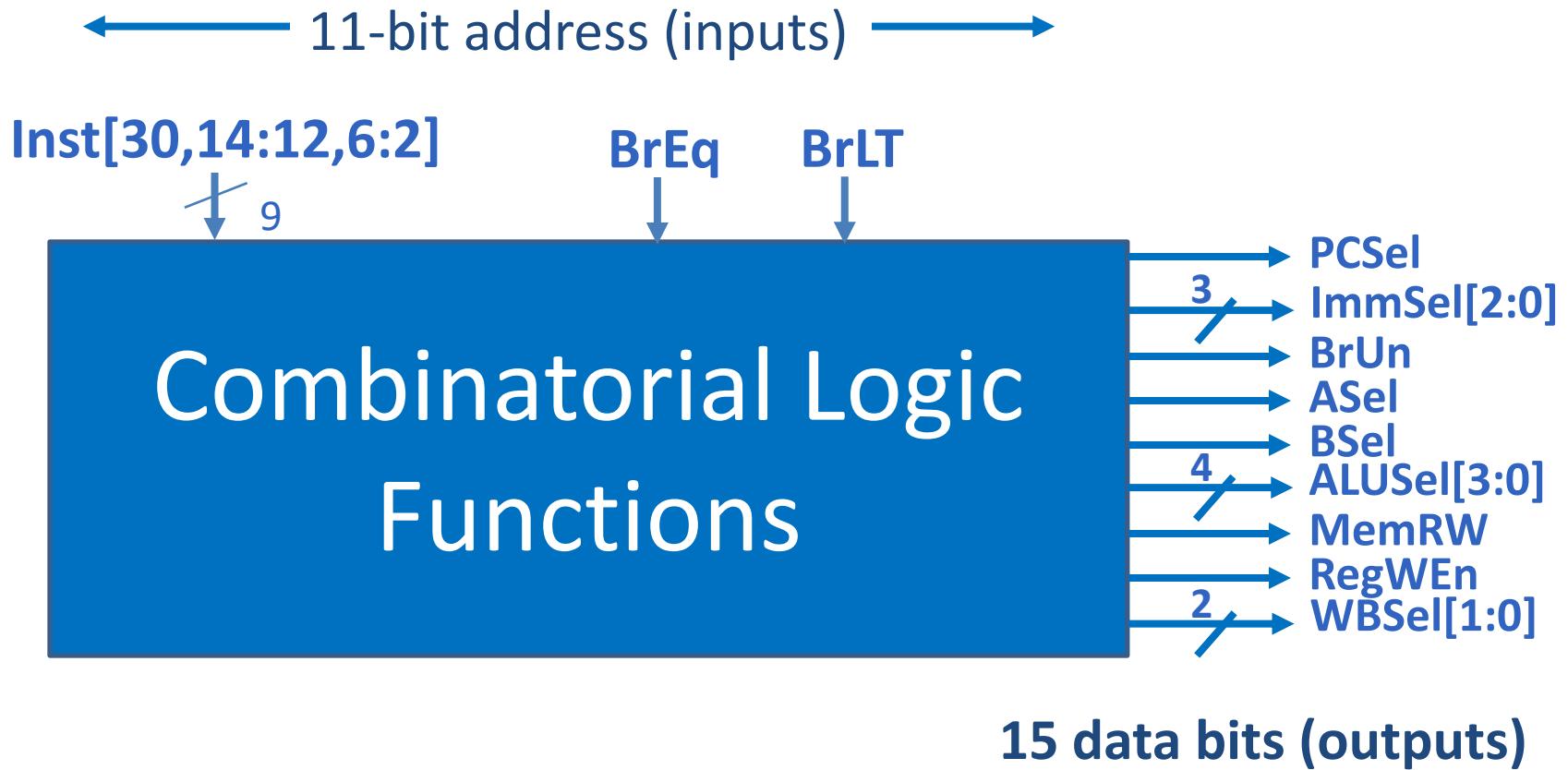
Keywords

- Mux
- Register
- FlipFlop
- Adder
- ALU
- AND; OR; XOR; NOT
- NAND; NOR

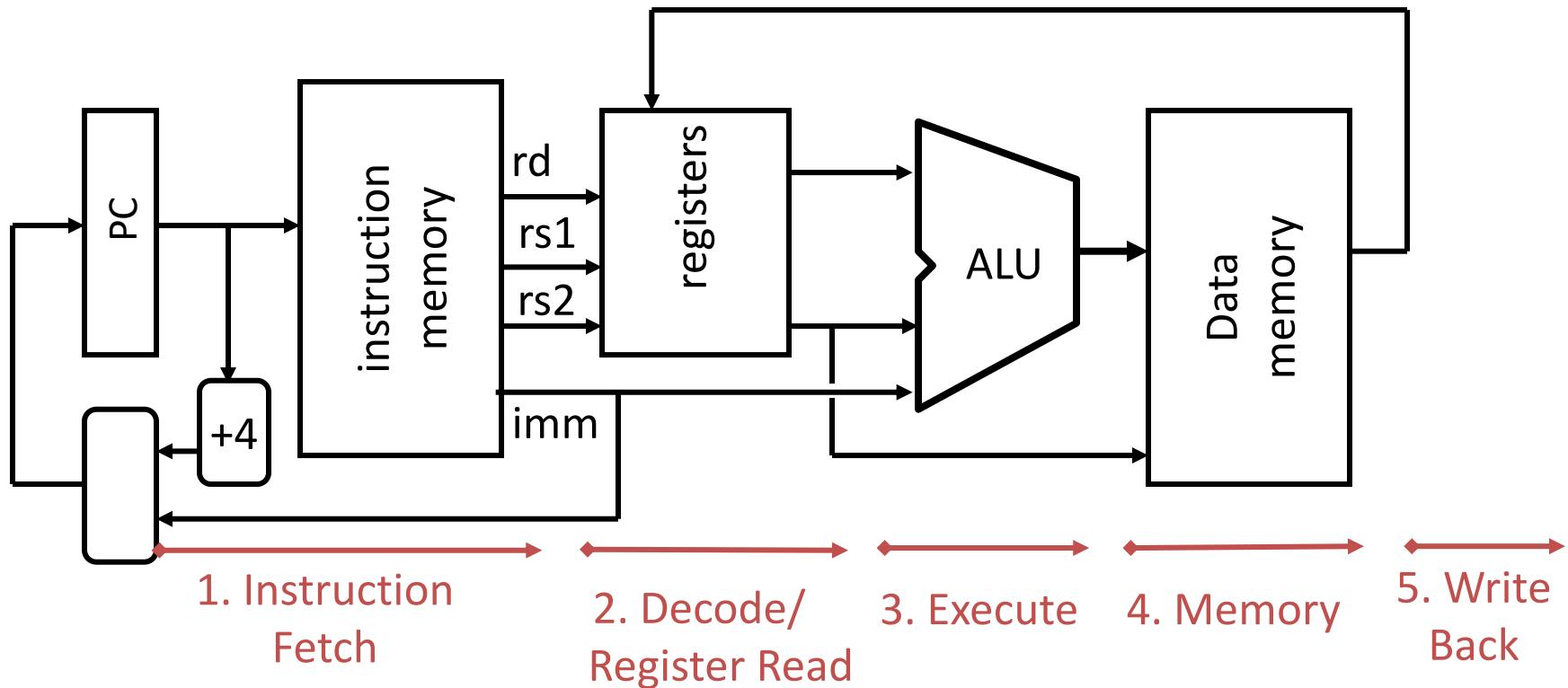
Complete RV32I Datapath!



Control Block Design



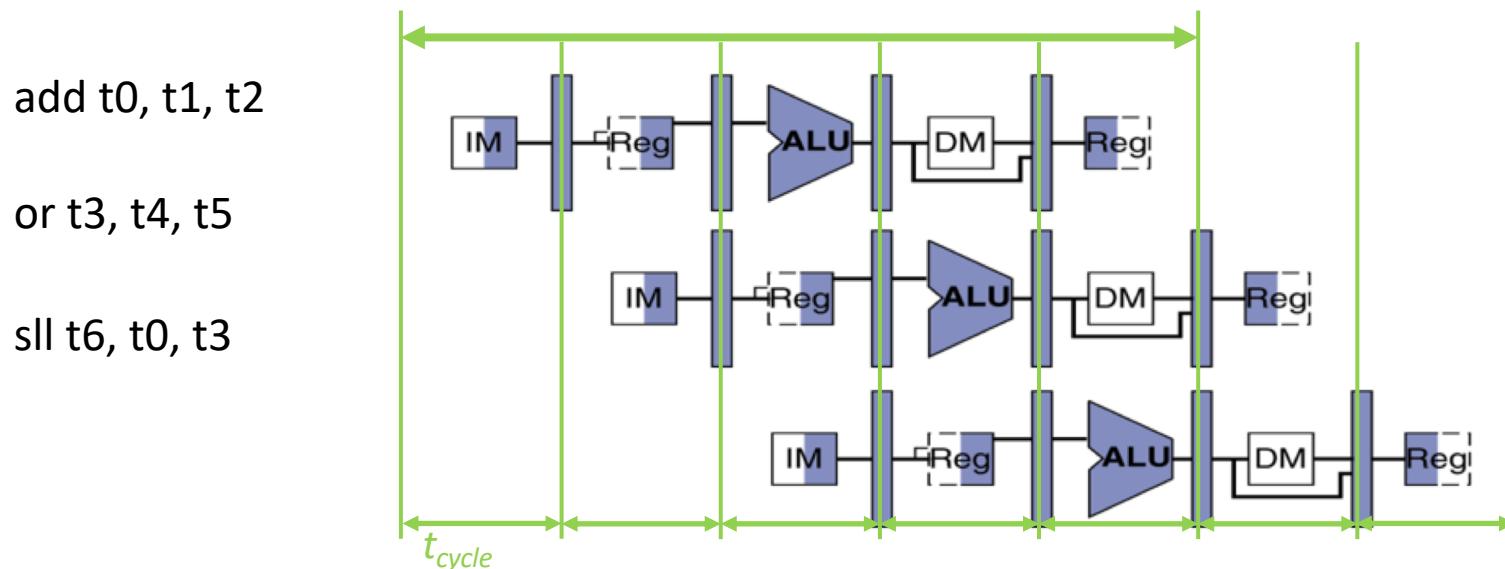
Single Cycle Datapath



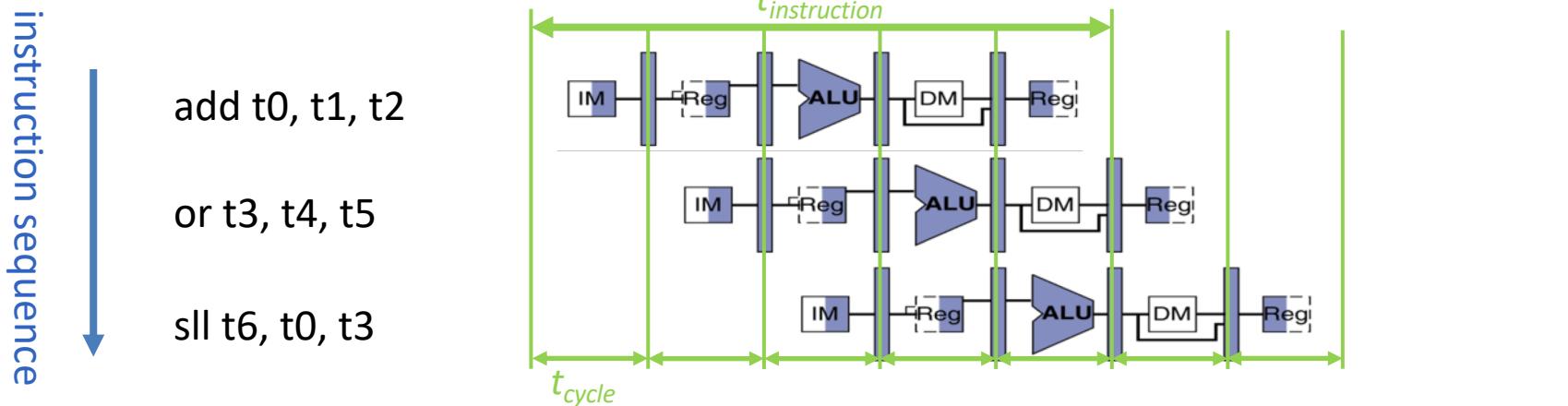
Pipelining with RISC-V

Phase	Pictogram	t_{step} Serial	t_{cycle} Pipelined
Instruction Fetch	IM	200 ps	200 ps
Reg Read	Reg	100 ps	200 ps
ALU	ALU	200 ps	200 ps
Memory	DM	200 ps	200 ps
Register Write	Reg	100 ps	200 ps
$t_{instruction}$	IM → Reg → ALU → DM → Reg	800 ps	1000 ps

instruction sequence ↓



Pipelining with RISC-V



	Single Cycle	Pipelining
Timing	$t_{step} = 100 \dots 200 \text{ ps}$	$t_{cycle} = 200 \text{ ps}$
	Register access only 100 ps	All cycles same length
Instruction time, $t_{instruction}$	$= t_{cycle} = 800 \text{ ps}$	1000 ps
CPI (Cycles Per Instruction)	~1 (ideal)	~1 (ideal), >1 (actual)
Clock rate, f_s	$1/800 \text{ ps} = 1.25 \text{ GHz}$	$1/200 \text{ ps} = 5 \text{ GHz}$
Relative speed	1 x	4 x

Sequential vs Simultaneous

What happens sequentially, what happens simultaneously?

add t0, t1, t2

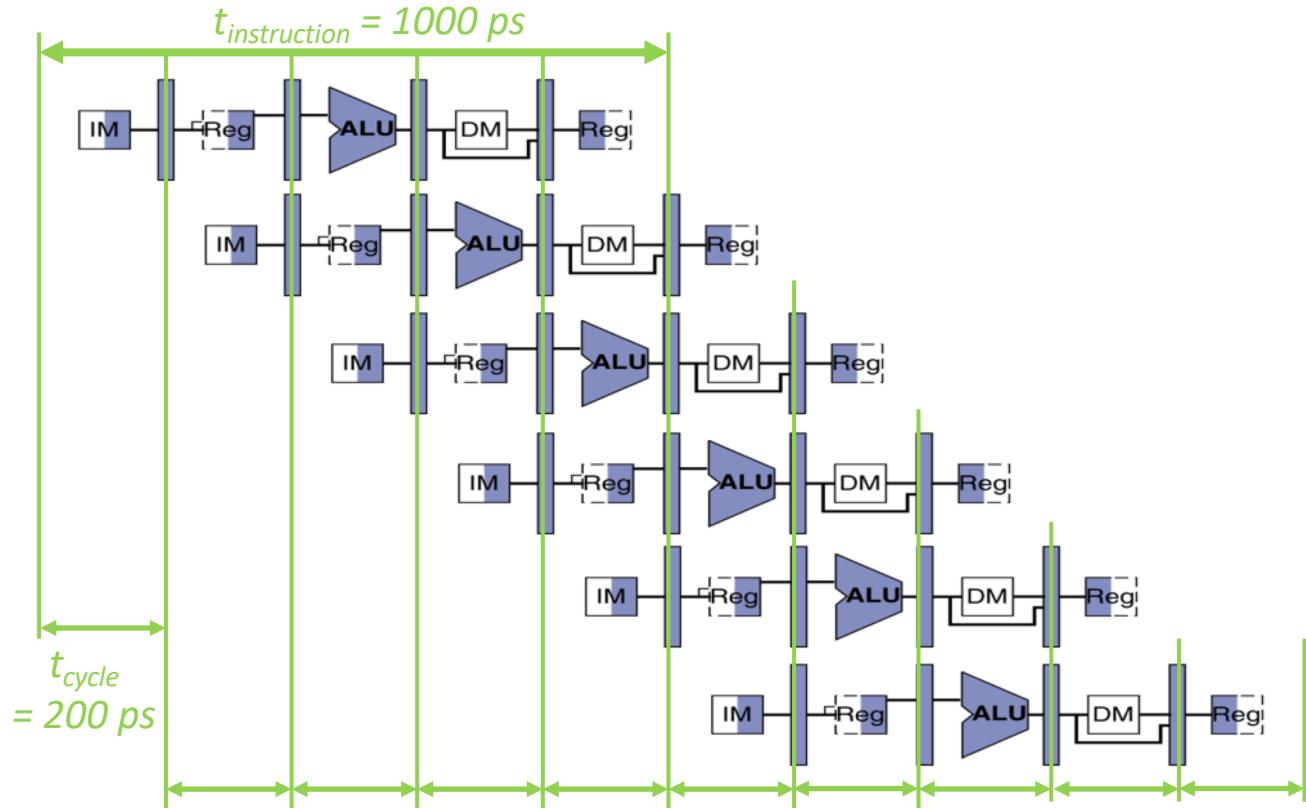
or t3, t4, t5

sll t6, t0, t3

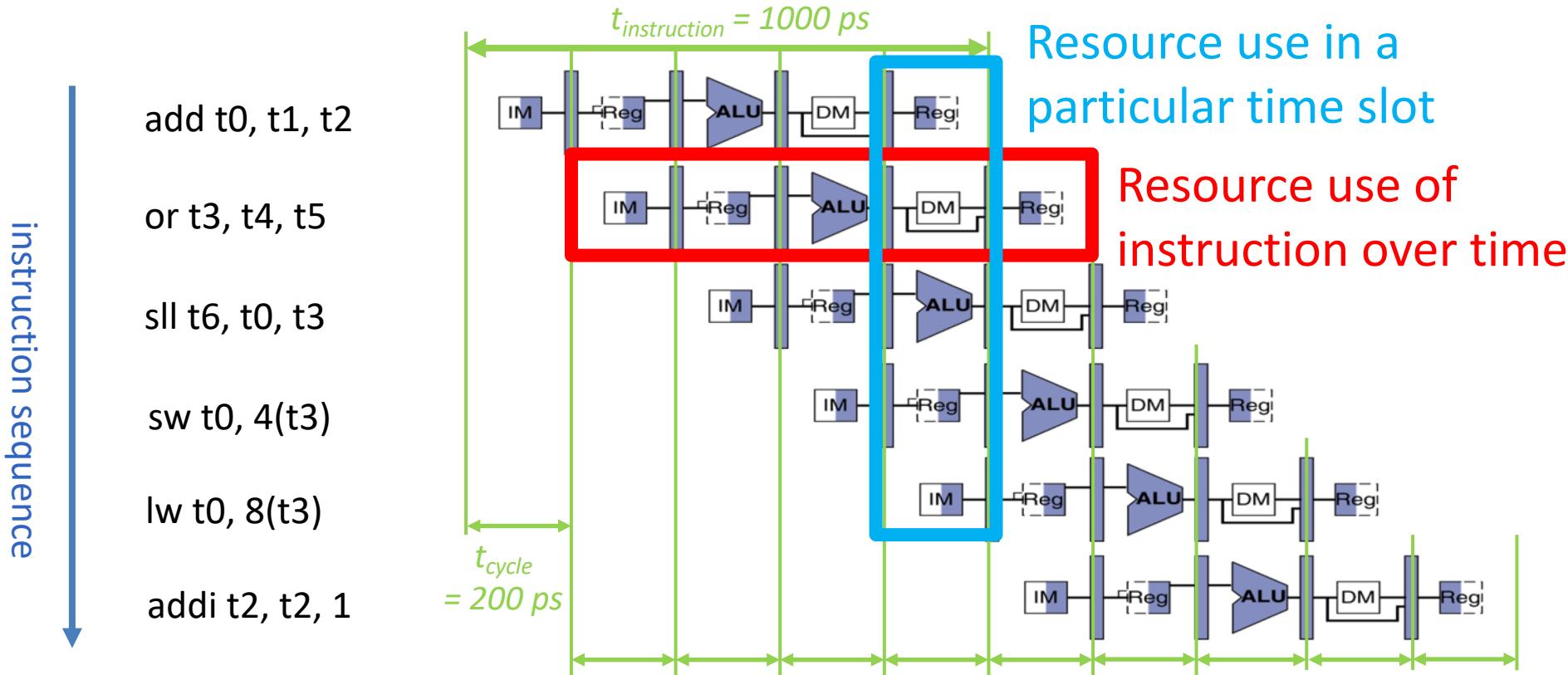
sw t0, 4(t3)

lw t0, 8(t3)

addi t2, t2, 1



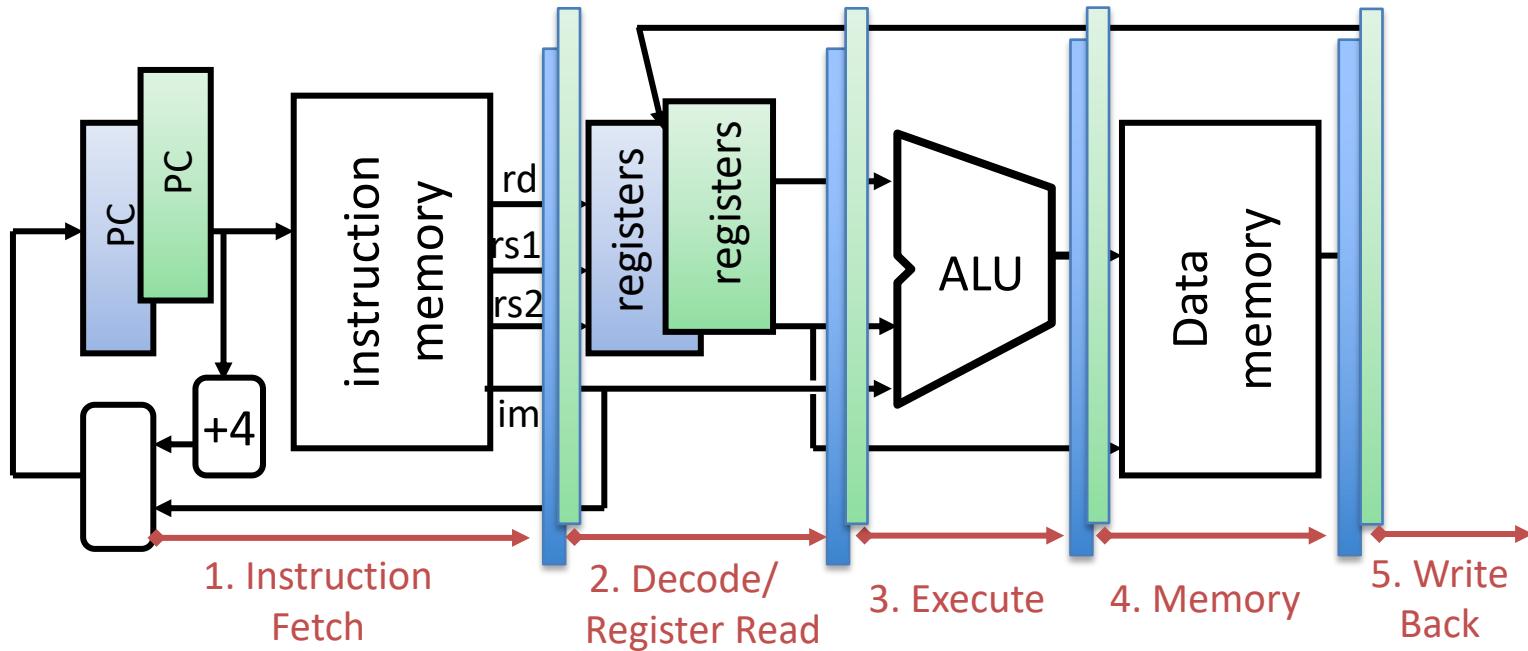
RISC-V Pipeline



Hazards!

- Keywords:
 - Structural Hazards
 - Data Hazards
 - Control Hazards
- Forwarding/ Bypassing
- Stall/ nop/ load delay slot
- Kill instruction
- Branch prediciton

Hyper-threading (simplified)



- Duplicate all elements that hold the state (registers)
- Use the same CL blocks
- Use muxes to select which state to use every clock cycle
- => run 2 independent processes
 - No Hazards: registers different; different control flow; memory different;
Threads: memory hazard should be solved by software (locking, mutex, ...)
- Speedup?
 - No obvious speedup; Complex pipeline: make use of CL blocks in case of unavailable resources (e.g. wait for memory)

Superscalar

- “Iron Law” of Processor Performance to estimate speed
- Complex Pipelines
 - Multiple Functional Units => Parallel execution
 - Static Multiple Issues (VLIW)
 - E.g. 2 instructions per cycle
 - Dynamic Multiple Issues (Superscalar)
 - Re-order instructions
 - Issue Buffer; Re-order Buffer; Commit Unit
 - Re-naming of registers

“Iron Law” of Processor Performance

$$\frac{\text{Time}}{\text{Program}} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Time}}{\text{Cycle}}$$

CPI = Cycles Per Instruction

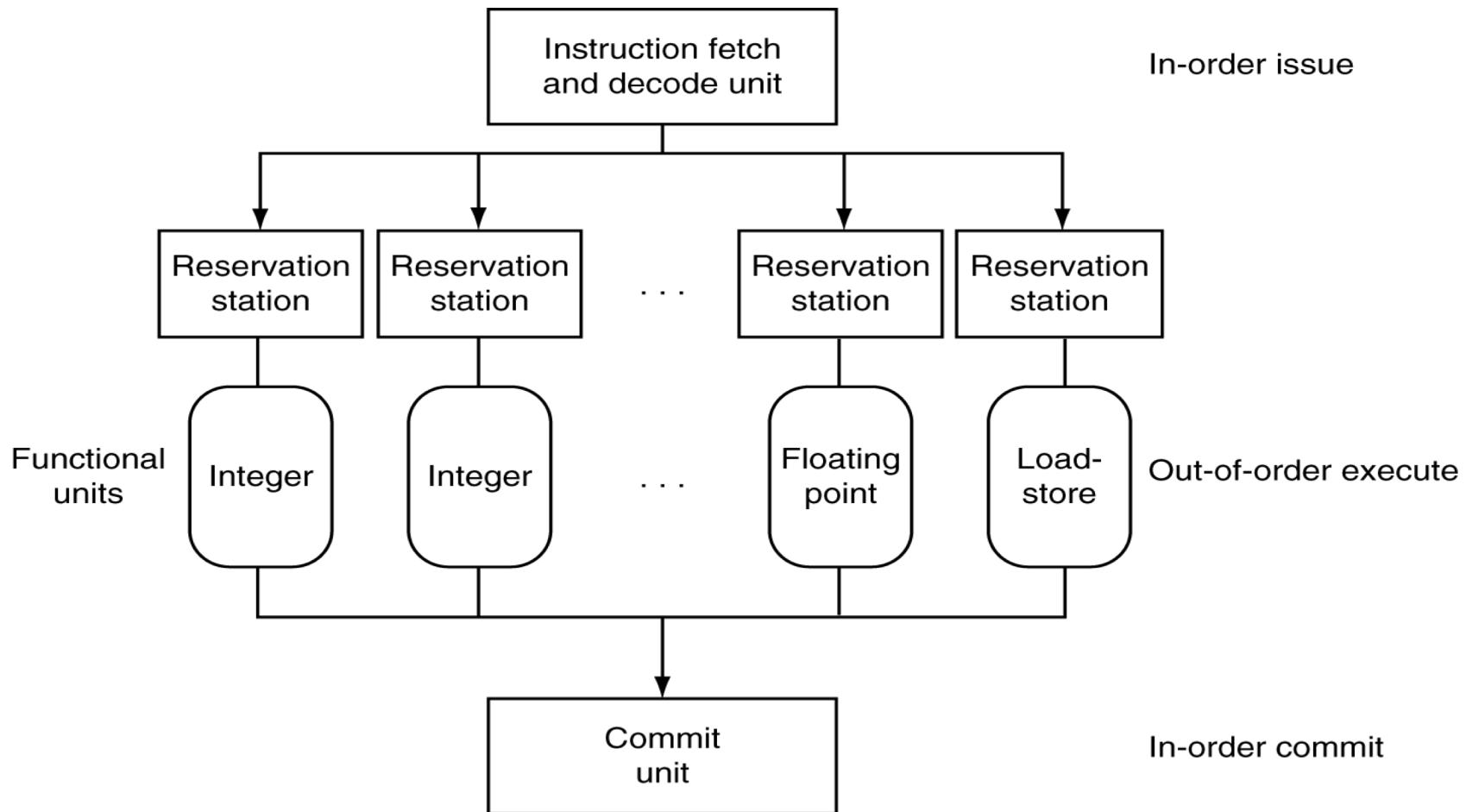
Can time Can count Can look up

$$\text{CPI} = \frac{\text{Cycles}}{\text{Instruction}} = \frac{\text{Time}}{\text{Program}} \div \left(\frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Time}}{\text{Cycle}} \right)$$

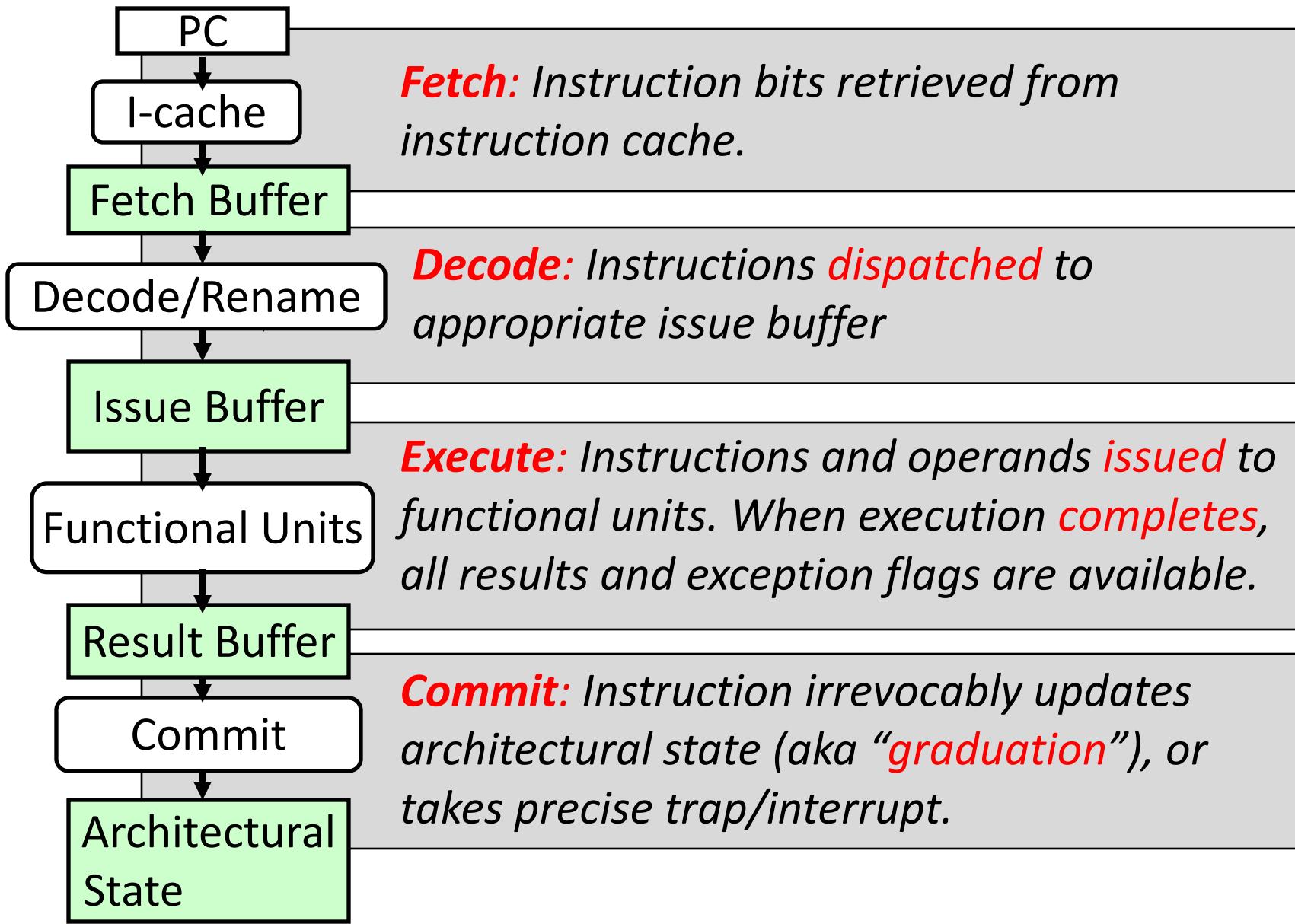
Example (RISC processor)

Op	Freq _i	CPI _i	Prod	(% Time)
ALU	50%	1	.5	(23%)
Load	20%	5	1.0	(45%)
Store	10%	3	.3	(14%)
Branch	20%	2	.4	(18%)
<u>Instruction Mix</u>		2.2		(Where time spent)

Superscalar Processor

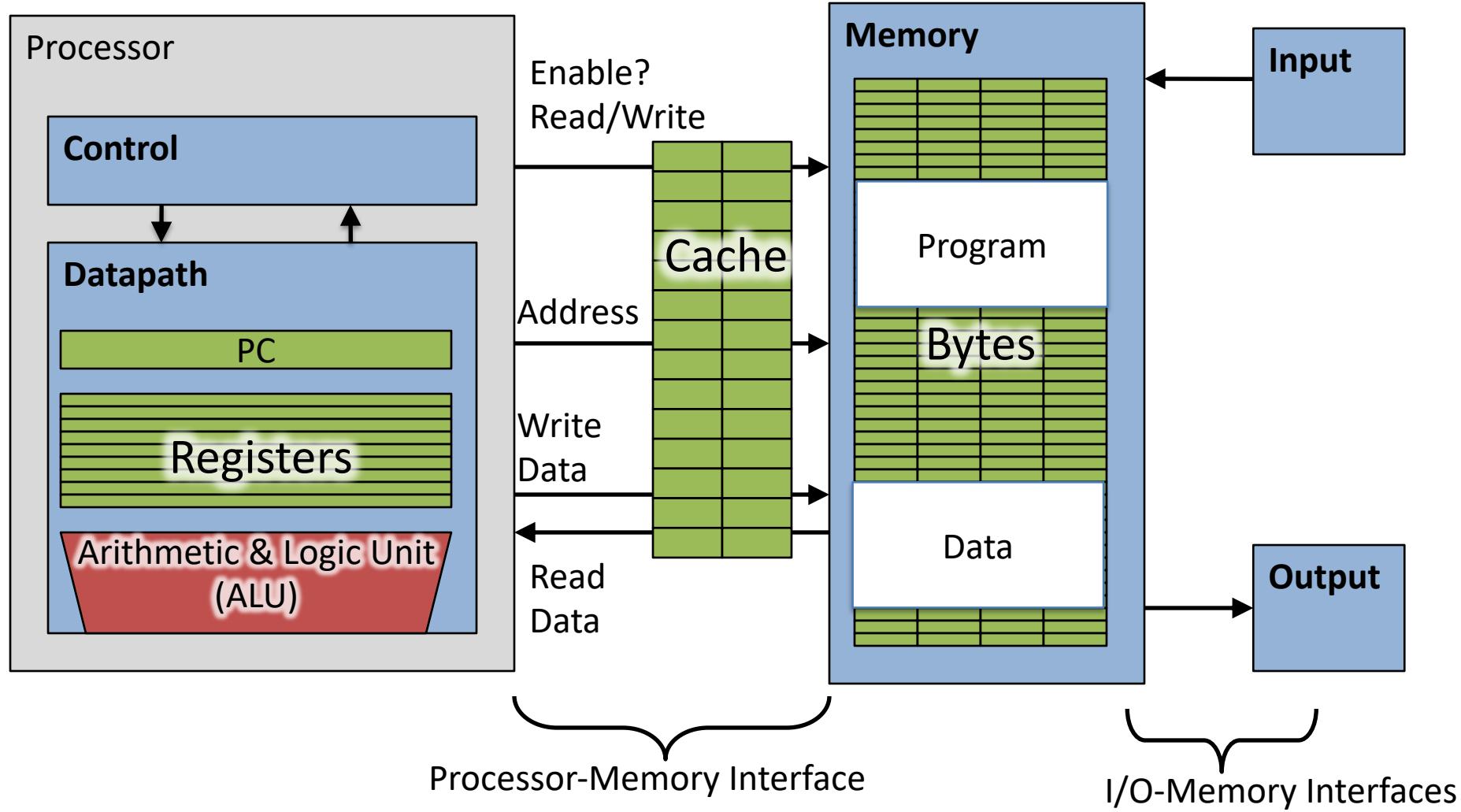


Phases of Instruction Execution

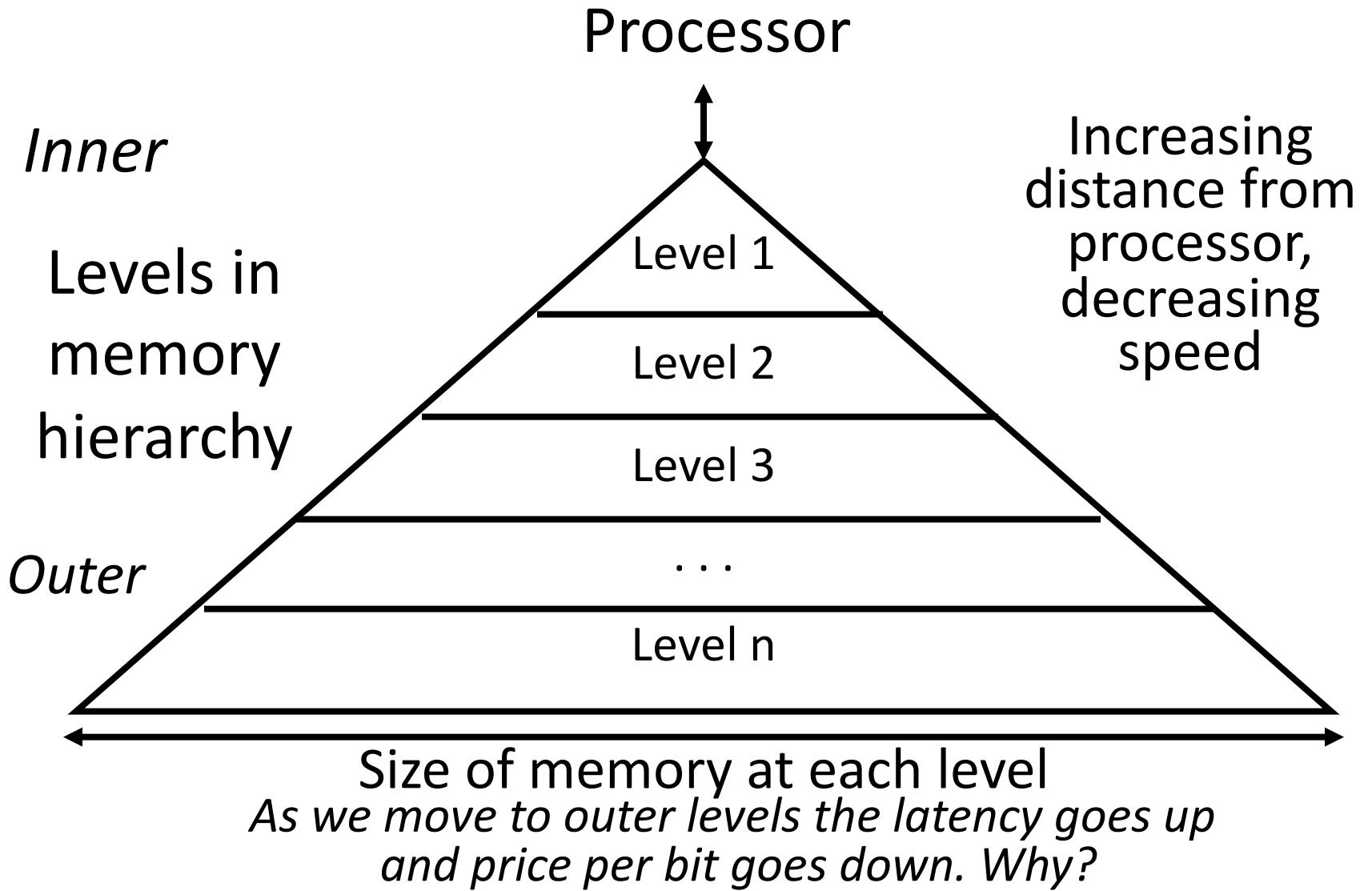


Intro to Caches

Adding Cache to Computer



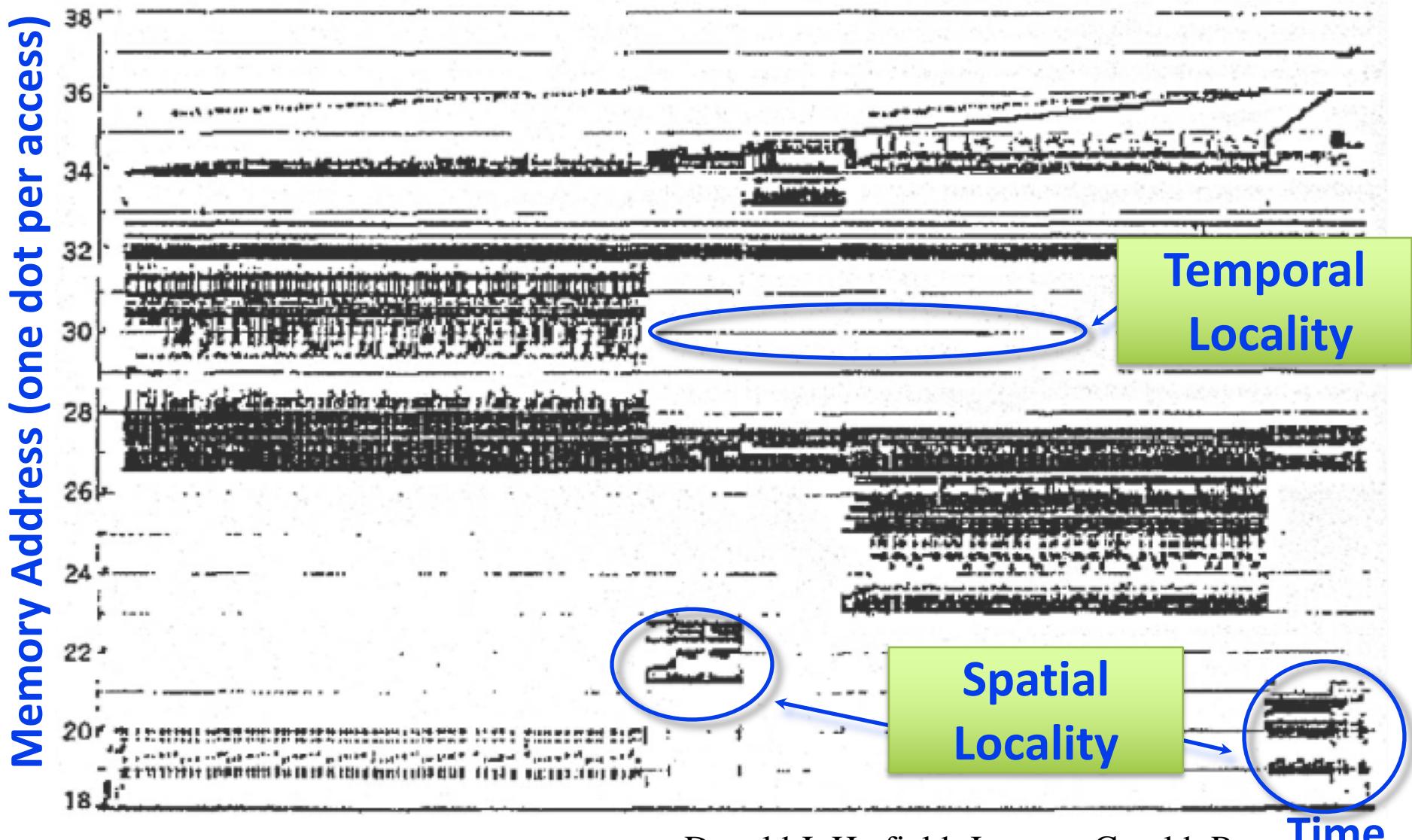
Big Idea: Memory Hierarchy



Big Idea: Locality

- *Temporal Locality* (locality in time)
 - Go back to same book on desktop multiple times
 - If a memory location is referenced, then it will tend to be referenced again soon
- *Spatial Locality* (locality in space)
 - When go to book shelf, pick up multiple books on J.D. Salinger since library stores related books together
 - If a memory location is referenced, the locations with nearby addresses will tend to be referenced soon

Memory Reference Patterns



Donald J. Hatfield, Jeanette Gerald: Program Restructuring for Virtual Memory. IBM Systems Journal 10(3): 168-192 (1971)

Principle of Locality

- *Principle of Locality*: Programs access small portion of address space at any instant of time (spatial locality) and repeatedly access that portion (temporal locality)
- What program structures lead to temporal and spatial locality in instruction accesses?
- In data accesses?

Cache Philosophy

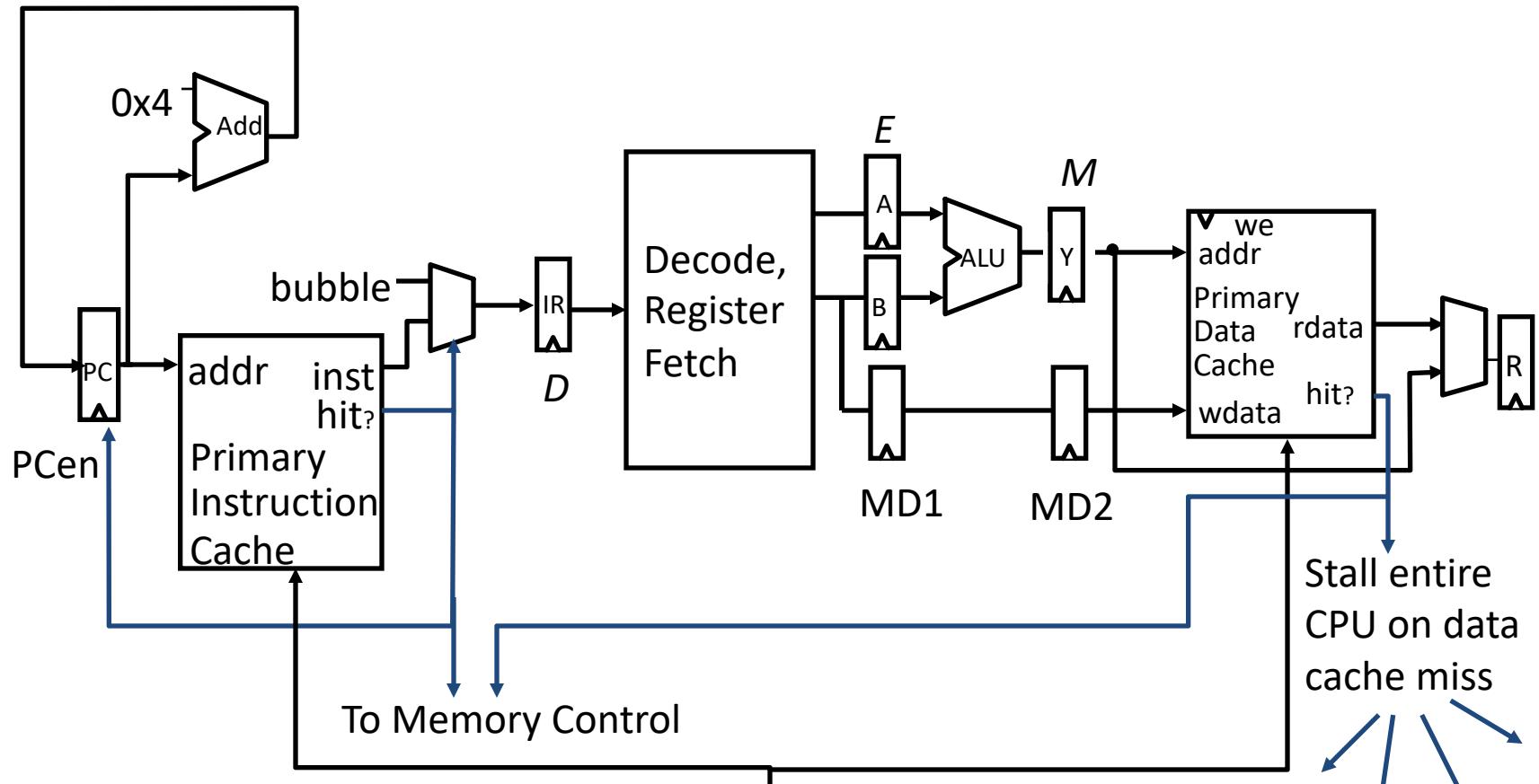
- Programmer-invisible hardware mechanism to give illusion of speed of fastest memory with size of largest memory
 - Works fine even if programmer has no idea what a cache is
 - However, performance-oriented programmers today sometimes “reverse engineer” cache design to design data structures to match cache

Cache Terms I

- **Cache:**
 - A small and fast memory used to increase the performance of accessing a big and slow memory
 - Uses ***temporal locality***: The tendency to reuse data in the same space over time
 - Uses ***spacial locality***: The tendency to use data at addresses near
- Cache ***hit***: The address being fetched is in the cache
- Cache ***miss***: The address being fetched is not in the cache
- ***Valid bit***: Is a particular entry valid
- Cache ***flush***: Invalidate all entries

CPU-Cache Interaction

(5-stage pipeline)



Cache Refill Data from Lower Levels of
Memory Hierarchy

Cache Terms II

- Cache ***level***:
 - The order in the memory hierarchy: L1\$ is closest to the processor
 - L1 caches may only hold data (Data-cache, D\$) or instructions (Instruction Cache, I\$)
 - Most L2+ caches are "unified", can hold both instructions and data
- Cache ***capacity***:
 - The total # of bytes in the cache
- Cache ***line*** or cache ***block***:
 - A single entry in the cache
- Cache ***block size***:
 - The number of bytes in each cache line

Cache Terms III

Associativity

- **Number of cache lines:**
 - Cache capacity / block size:
- Cache **associativity**:
 - The number of possible cache lines a given address may exist in.
 - Also the number of comparison operations needed to check for an element in the cache
 - **Direct mapped**: A data element can only be in one possible location ($N=1$)
 - **N-way set associative**: A data element can be in one of N possible positions
 - **Fully associative**: A data element can be at any location in the cache.
 - Associativity == # of lines
- Total # of cache lines == capacity of cache/line size
- Total # of lines in a set == # ways == N == associativity
- Total # of sets == # of cache lines / associativity

Victim Cache

- Conflict misses are a pain, but...
 - Perhaps a little associativity can help without having to be a fully associative cache
- In addition to the main cache...
 - Optionally have a very small (16-64 entry) ***fully associative*** "victim" cache
- Whenever we evict a cache entry
 - Don't just get rid of it, put it in the victim cache
- Now on cache misses...
 - Check the victim cache first, if it is in the victim cache you can just reload it from there

Cache Terms IV

Parts of the Address

- Address is divided into | **TAG** | **INDEX** | **OFFSET** |
- **Offset:**
 - The lowest bits of the memory address which say where data exists within the cache line.
 - It is $\log_2(\text{line/block size})$
 - So for a cache with 64B blocks it is 6 bits
- **Index:**
 - The portion of the address which says where in the cache an address may be stored
 - Takes $\log_2(\# \text{ of cache lines / associativity})$ bits
 - So for a 4 way associative cache with 512 lines it is 7 bits
- **Tag:** The portion of the address which must be stored in the cache to check if a location matches
 - # of bits of address - (# of bits for index + # of bits for offset)
 - So with 64b addresses it is 51b...

Cache Terms V

Writing

- ***Eviction:***
 - The process of removing an entry from the cache
- ***Write Back:***
 - A cache which only writes data up the hierarchy when a cache line is evicted
 - Instead set a ***dirty bit*** on cache entries
 - All i7 caches are ***write back***
- ***Write Through:***
 - A cache which always writes to memory
- ***Write Allocate:***
 - If writing to memory ***not in the cache*** fetch it first
 - i7 L2 is Write Allocate
- ***No Write Allocate:***
 - Just write to memory without a fetch
 - i7 L1 is no write allocate

Replacement Policy

In an associative cache, which line from a set should be evicted when the set becomes full?

- Random
- Least-Recently Used (LRU)
 - LRU cache state must be updated on every access
 - True implementation only feasible for small sets (2-way)
 - Pseudo-LRU binary tree often used for 4-8 way
- First-In, First-Out (FIFO) a.k.a. Round-Robin
 - Used in highly associative caches
- Not-Most-Recently Used (NMRU)
 - FIFO with exception for most-recently used line or lines

This is a second-order effect. Why?

Replacement only happens on misses

Cache Terms VI

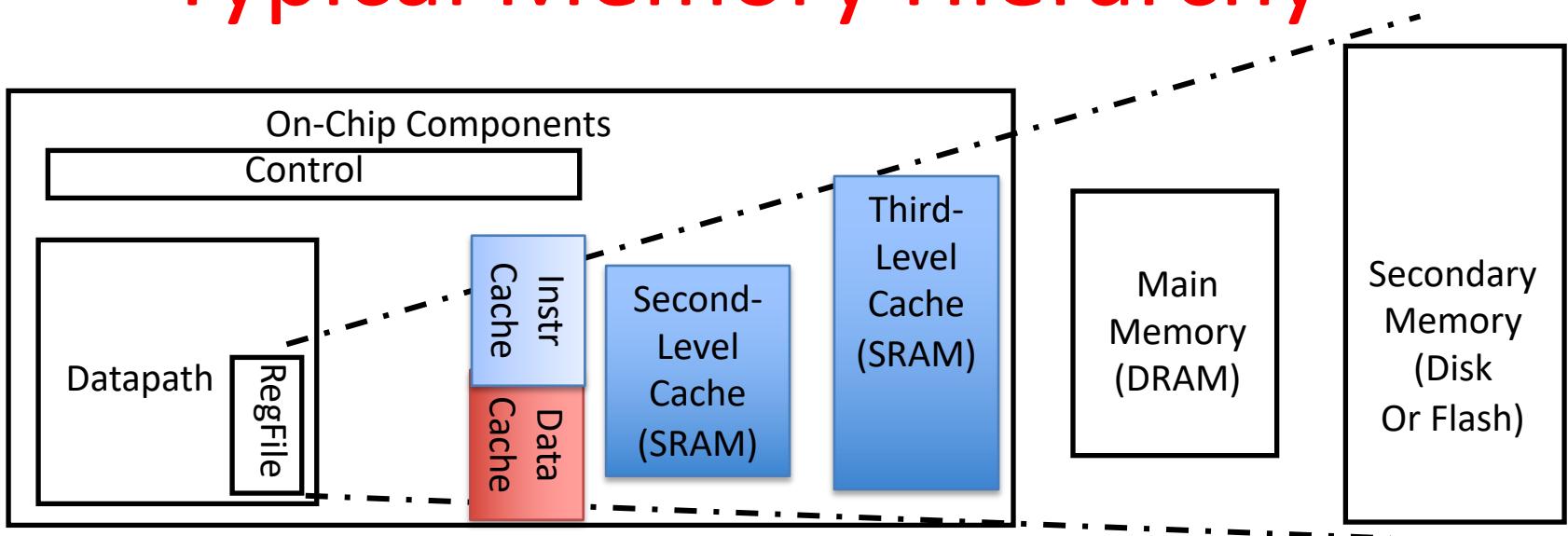
Cache Performance

- ***Hit Time:***
 - Amount of time to return data in a given cache: depends on the cache
 - i7 L1 hit time: 4 clock cycles
- ***Miss Penalty:***
 - Amount of ***additional*** time to return an element if its not in the cache: depends on the cache
- ***Miss Rate:***
 - Fraction of a ***particular program's*** memory requests which miss in the cache
- Average Memory Access Time (***AMAT***):
 - Hit time + Miss Rate * Miss Penalty

Understanding Cache Misses: The 3Cs

- **Compulsory** (cold start or process migration, 1st reference):
 - First access to block impossible to avoid; small effect for long running programs
 - Solution: increase block size (increases miss penalty; very large blocks could increase miss rate)
- **Capacity:**
 - Cache cannot contain all blocks accessed by the program
 - Solution: increase cache size (may increase access time)
- **Conflict (*collision*):**
 - *Multiple memory locations mapped to the same cache location*
 - *Solution 1: increase cache size*
 - *Solution 2: increase associativity (may increase access time)*

Typical Memory Hierarchy



Speed (cycles):	½'s	1's	10's	100's	1,000,000's
Size (bytes):	100's	10K's	M's	G's	T's
Cost/bit:	highest				lowest

- Principle of locality + memory hierarchy presents programmer with ≈ as much memory as is available in the *cheapest* technology at the ≈ speed offered by the *fastest* technology

Local vs. Global Miss Rates

- *Local miss rate* – the fraction of references to one level of a cache that miss
- Local Miss rate $L2\$ = \frac{\$L2\ Misses}{L1\$ Misses}$
- *Global miss rate* – the fraction of references that miss in all levels of a multilevel cache
 - $L2\$$ local miss rate \gg than the global miss rate
- Global Miss rate $= L2\$ Misses / \text{Total Accesses}$
 $= (\frac{L2\$ Misses}{L1\$ Misses}) \times (\frac{L1\$ Misses}{\text{Total Accesses}})$
 $= \text{Local Miss rate } L2\$ \times \text{Local Miss rate } L1\$$
- $\text{AMAT} = \text{Time for a hit} + \text{Miss rate} \times \text{Miss penalty}$
- $\text{AMAT} = \text{Time for a } L1\$ \text{ hit} + (\text{local}) \text{ Miss rate } L1\$ \times (\text{Time for a } L2\$ \text{ hit} + (\text{local}) \text{ Miss rate } L2\$ \times L2\$ \text{ Miss penalty})$

In Conclusion, Cache Design Space

- Several interacting dimensions
 - Cache size
 - Block size
 - Associativity
 - Replacement policy
 - Write-through vs. write-back
 - Write-allocation
- Optimal choice is a compromise
 - Depends on access characteristics
 - Workload
 - Use (I-cache, D-cache)
 - Depends on technology / cost
- Simplicity often wins

