

Lecture 14:

Assignment Project Exam Help

<https://tutorcs.com>

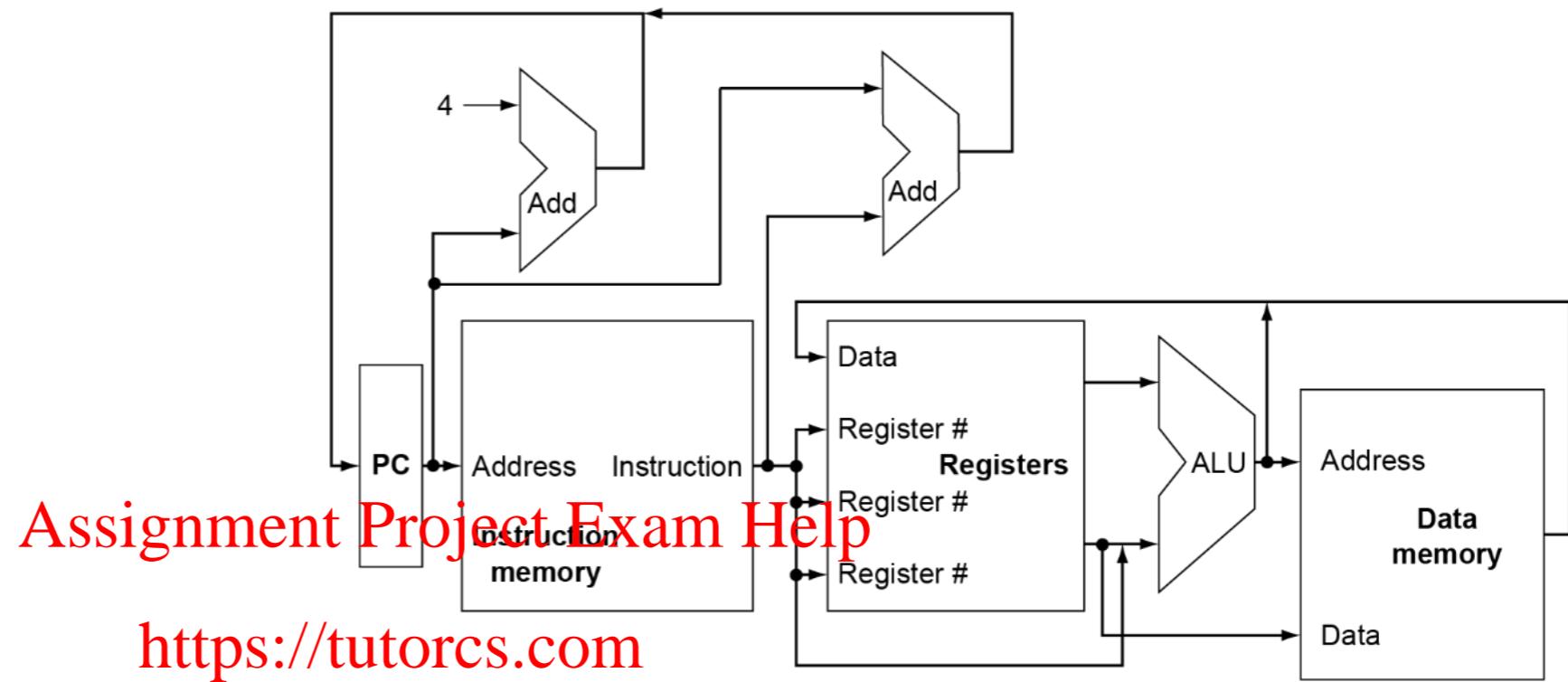
WeChat: cstutorcs

The Memory System 1/3

Introduction to Computer Architecture

UC Davis EEC 170, Fall 2019

The fundamental problem



- **We want a big fast memory** WeChat: cstutorcs
- **We have:**
 - **Big slow memory**
 - **Small fast memory**
- **How do we solve this?**

Principle of Locality

- Programs access a small proportion of their address space at any time
- Temporal locality
 - Items accessed recently are likely to be accessed again soon
Assignment Project Exam Help
 - e.g., instructions in a loop, induction variables
<https://tutorcs.com>
- Spatial locality
 - Items near those accessed recently are likely to be accessed soon
 - E.g., sequential instruction access, array data

Taking Advantage of Locality

- Memory hierarchy
- Store everything on disk
- Copy recently accessed (and nearby) items from disk to smaller DRAM memory
 - Main memory
- Copy more recently accessed (and nearby) items from DRAM to smaller SRAM memory
 - Cache memory attached to CPU

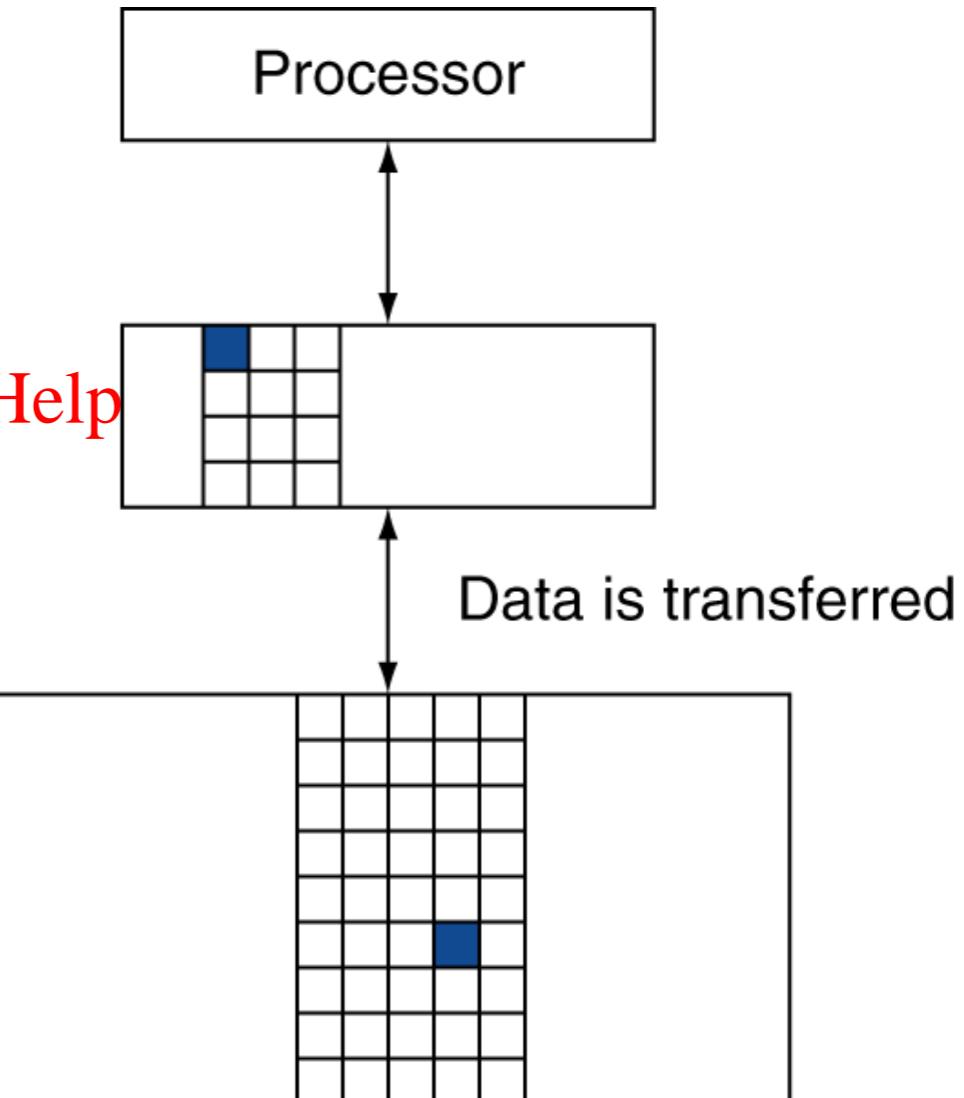
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

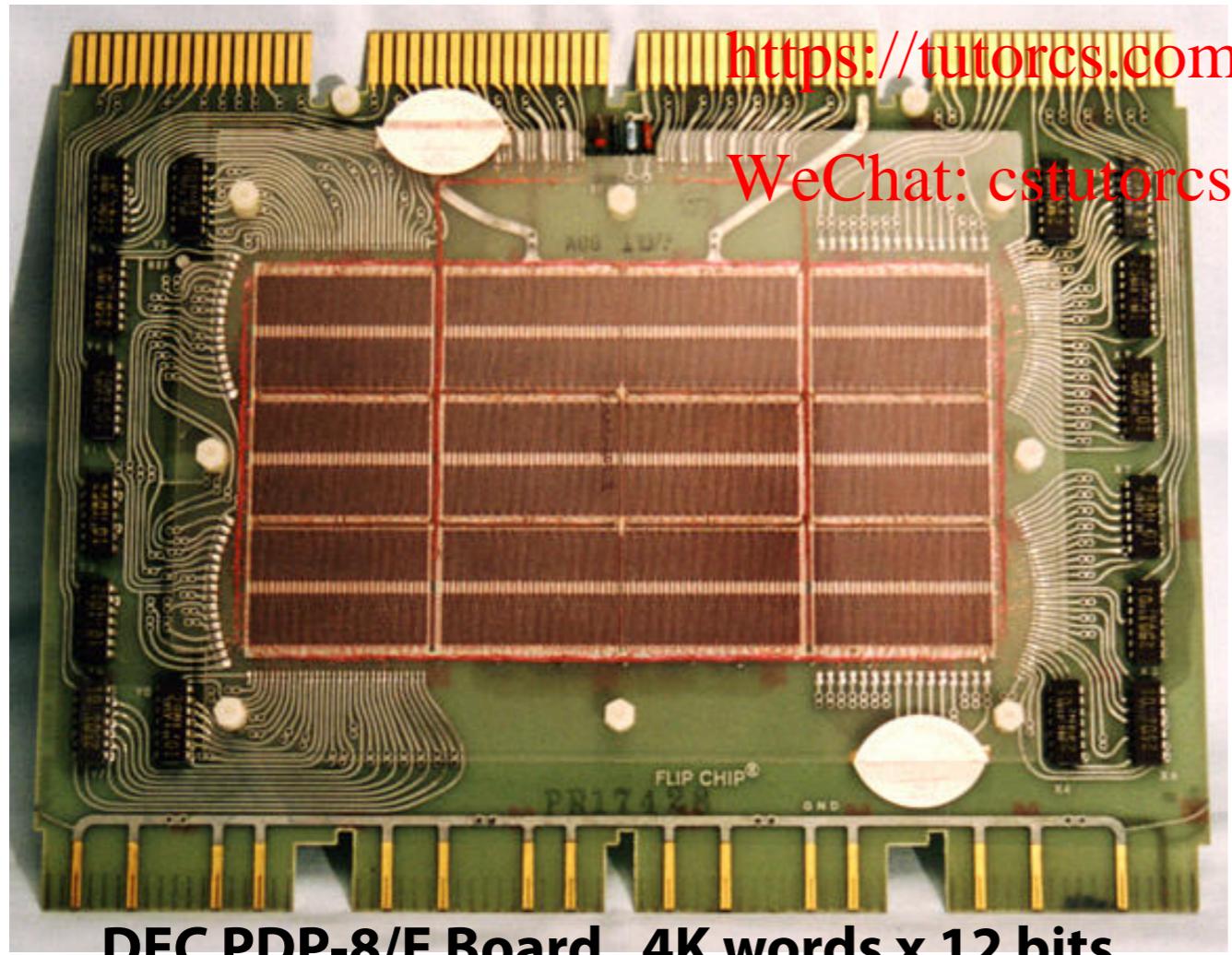
Memory Hierarchy Levels

- **Block (aka line): unit of copying**
 - May be multiple words
- **If accessed data is present in upper level**
 - **Hit: access satisfied by upper level**
Assignment Project Exam Help
<https://tutorcs.com>
 - **Hit ratio: hits/accesses**
- **If accessed data is absent**
WeChat: cstutorcs
 - **Miss: block copied from lower level**
 - Time taken: miss penalty
 - **Miss ratio: misses/accesses**
 $= 1 - \text{hit ratio}$
 - Then accessed data supplied from upper level



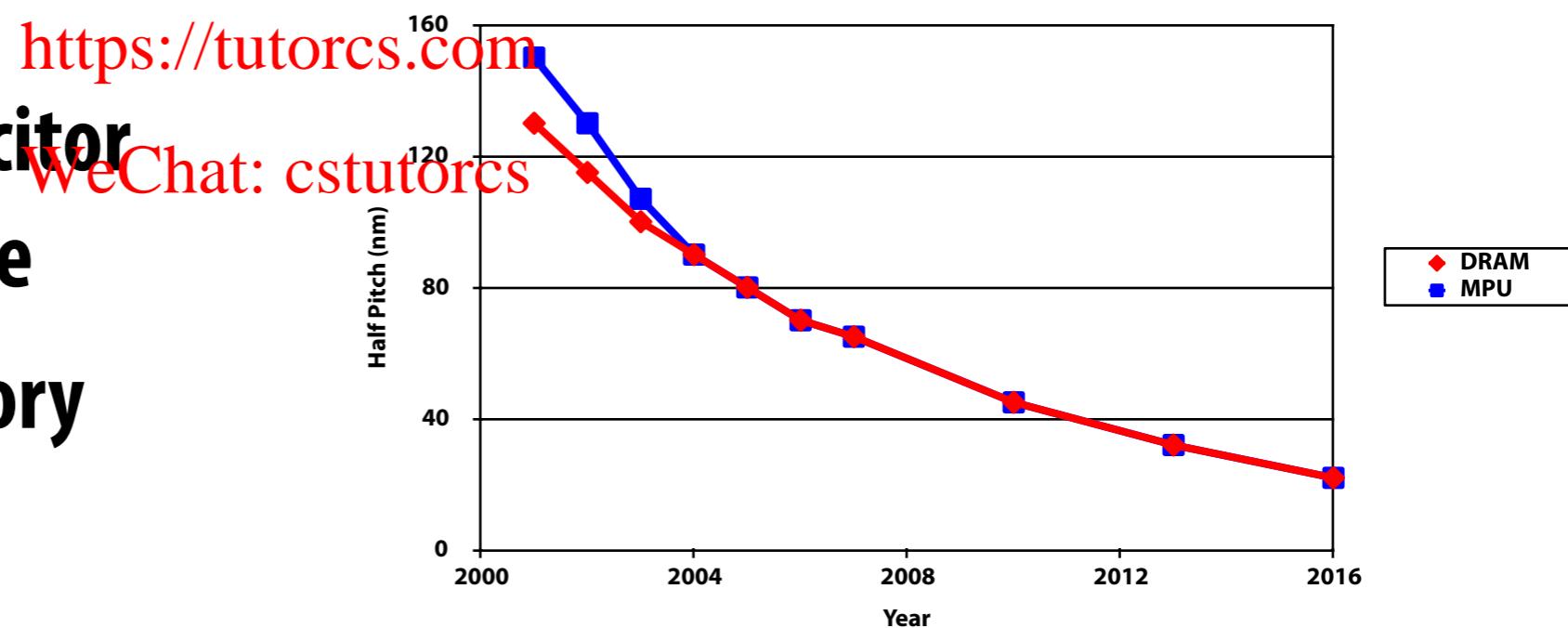
Core Memory

- Core memory was first large scale reliable main memory
 - invented by Forrester in late 40s at MIT for Whirlwind project
- Bits stored as magnetization polarity on small ferrite cores threaded onto 2 dimensional grid of wires
- Coincident current pulses on X and Y wires would write cell and also sense original state (destructive reads)
 - Assignment Project Exam Help
 - Robust, non-volatile storage
 - Used on space shuttle computers until recently
 - Cores threaded onto wires by hand (25 billion a year at peak production)
 - Core access time ~ 1ms



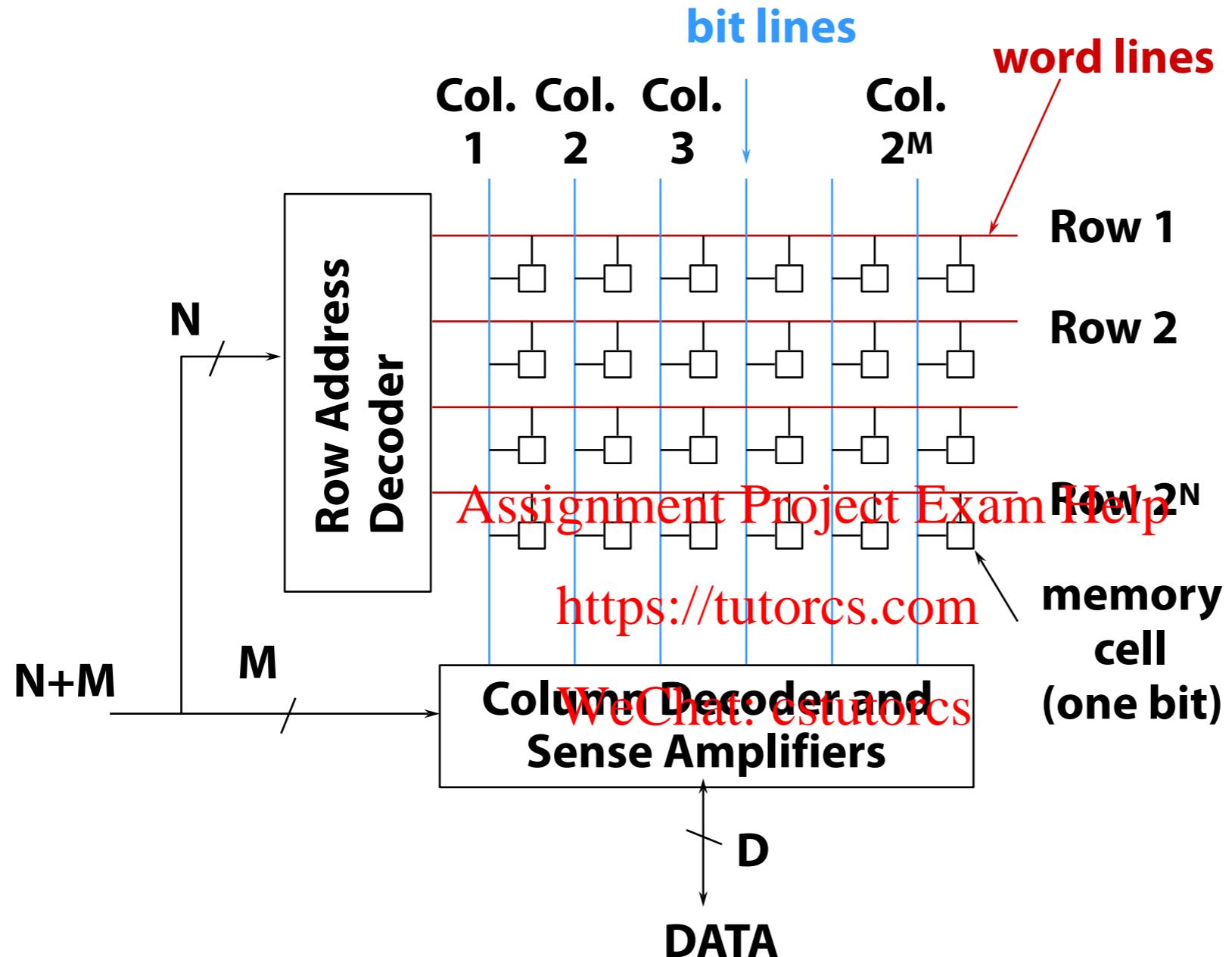
Semiconductor Memory, DRAM

- Semiconductor memory began to be competitive in early 1970s
 - Intel formed to exploit market for semiconductor memory
- First commercial DRAM was Intel 1103
 - 1Kbit of storage on single chip
 - charge on a capacitor used to hold value
- Semiconductor memory quickly replaced core in 1970s



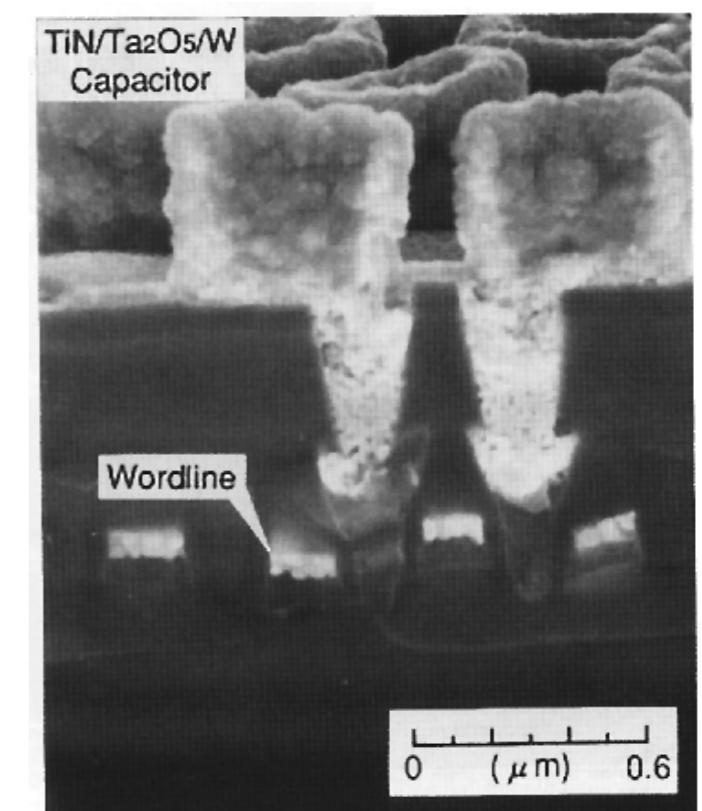
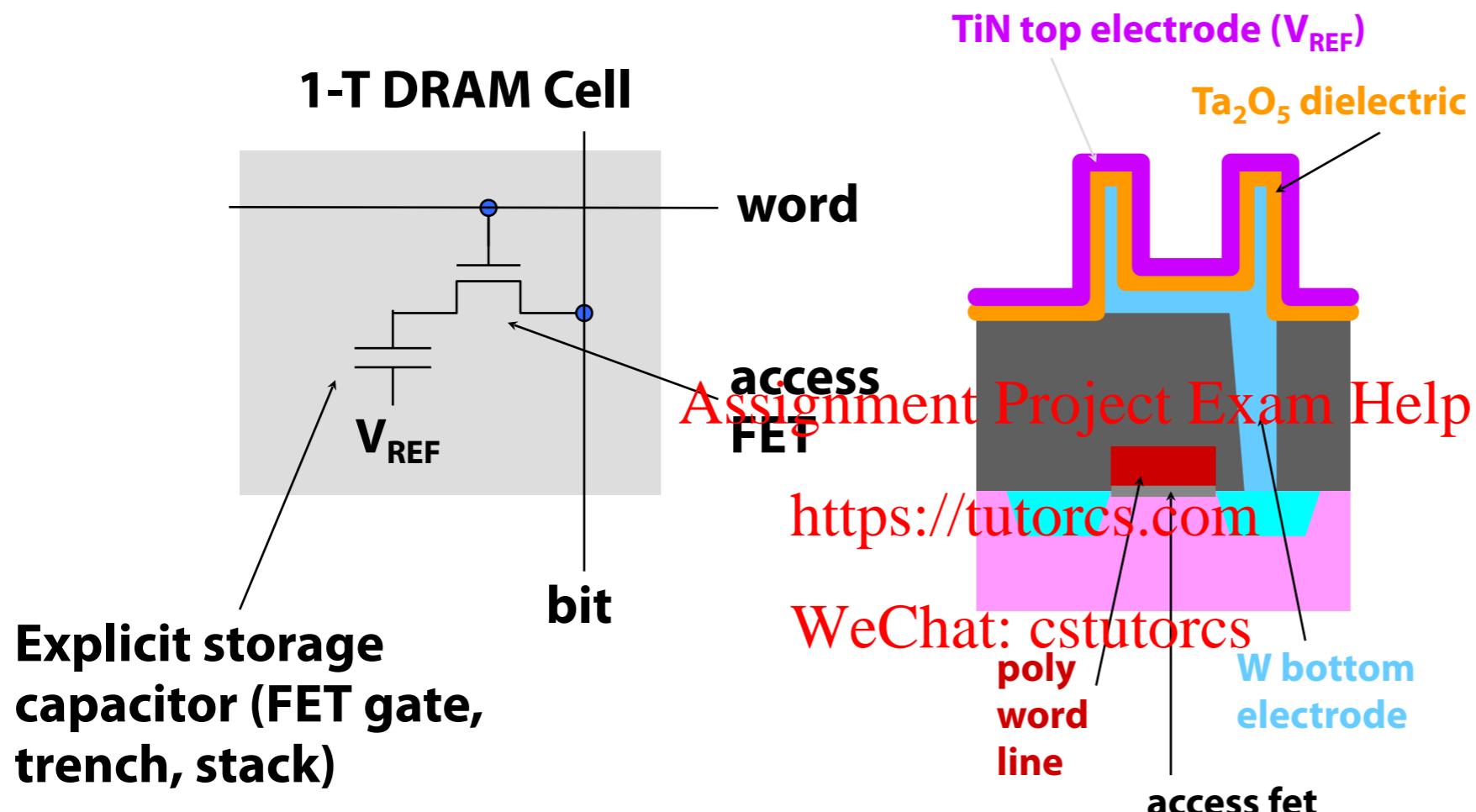
Int'l Tech Roadmap 2001

DRAM Architecture

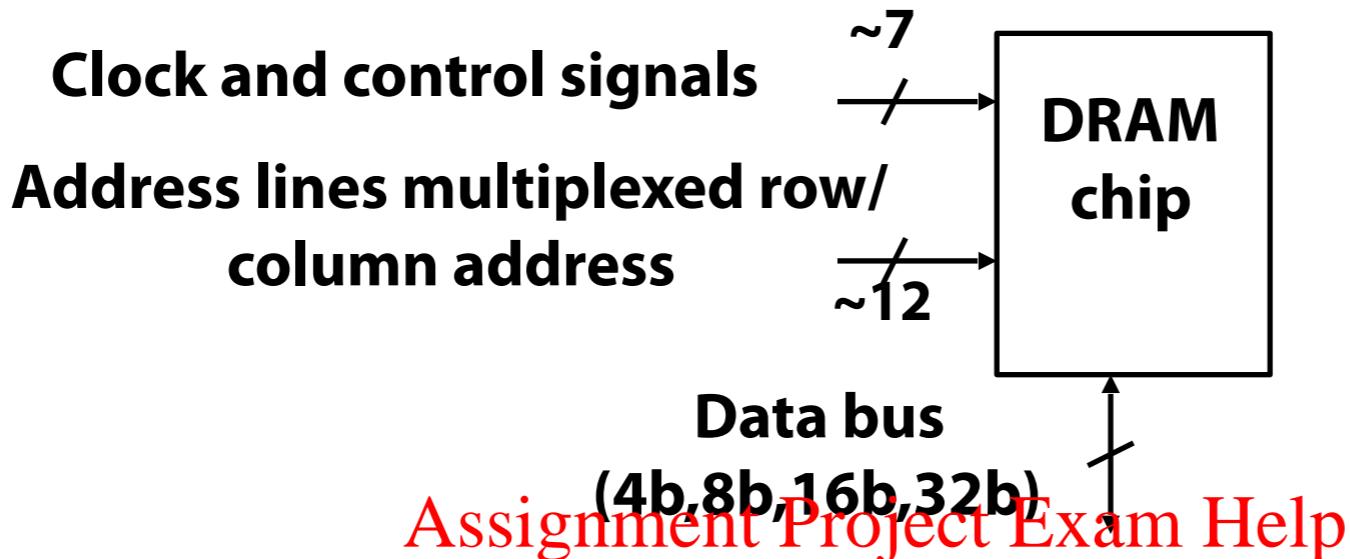


- Bits stored in 2-dimensional arrays on chip
- Modern chips have around 4 logical banks on each chip
 - each logical bank physically implemented as many smaller arrays

One Transistor Dynamic RAM



DRAM Packaging



- **DIMM (Dual Inline Memory Module) contains multiple chips with clock/control/address signals connected in parallel (sometimes need buffers to drive signals to all chips)**
- **Data pins work together to return wide word (e.g., 64-bit data bus using 16x4-bit parts)**



72-pin SO DIMM

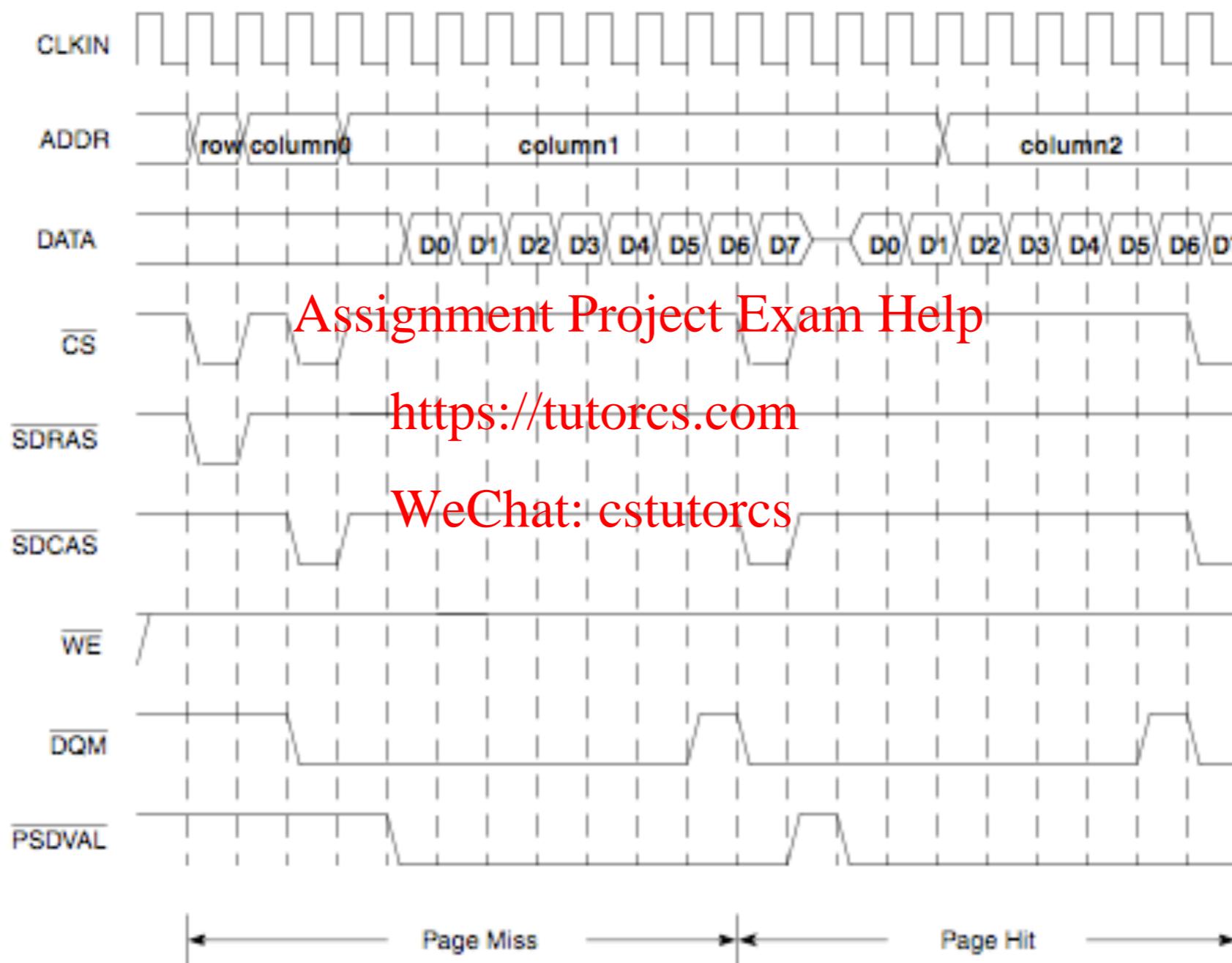


168-pin DIMM

DRAM Operation

- Three steps in read/write access to a given bank
- Precharge
 - charges bit lines to known value
- Row access (RAS)
 - decode row address, enable addressed row (often multiple Kb)
 - bitlines share charge with storage cell
 - small change in voltage detected by sense amplifiers which latch whole row
 - sense amplifiers drive bitlines <https://cstutorcs.com> to charge storage cell
- Column access (CAS)
 - decode column address to select small number of sense amplifier latches (4, 8, 16, or 32 bits depending on DRAM package)
 - on read, send latched bits out to chip pins
 - on write, change sense amplifier latches which then charge storage cells to required value
 - can perform multiple column accesses on same row without another row access (burst mode)
- Each step has a latency of around 15–20ns in modern DRAMs
- Various DRAM standards (DDR, RDRAM) have different ways of encoding the signals for transmission to the DRAM, but all share the same core architecture

Burst Read (Freescale MPC8260 SDRAM)



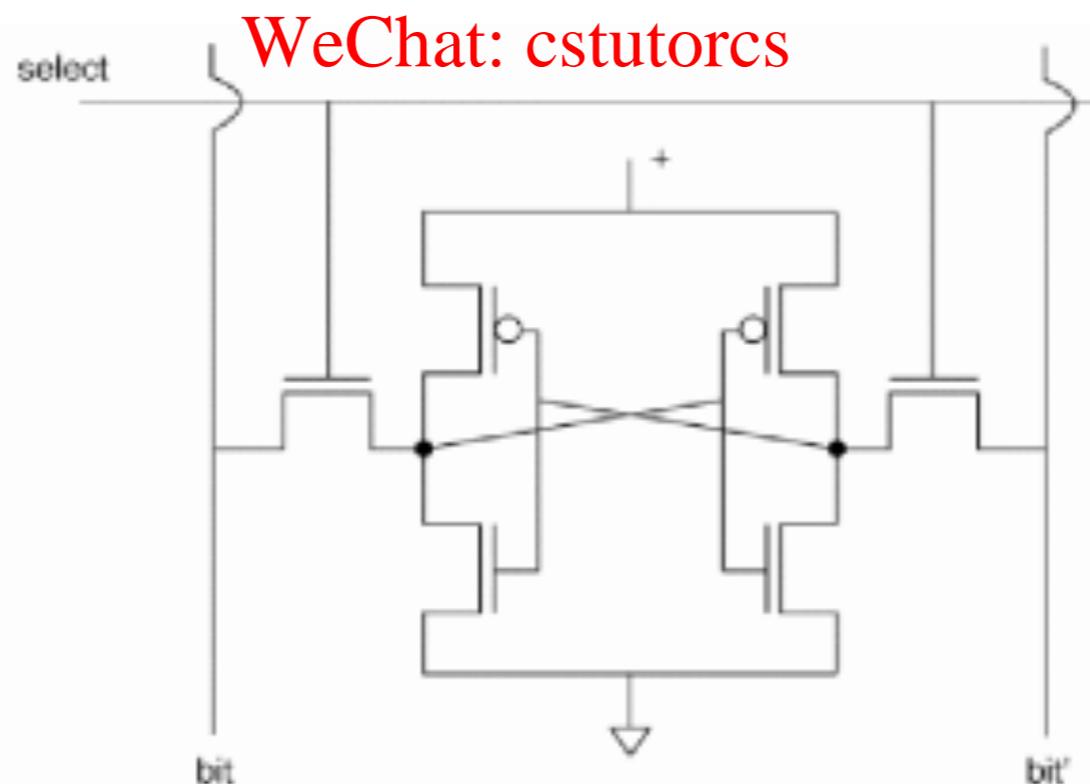
PSDMR[ACTTORW] = 010, PSDMR[CL] = 2, PSDMR[BL] = 1 (Burst Length = 8)

SRAM Cell

- SRAM cell takes 6 transistors
- Cross-coupled inverters
- “select” is “word line”, selects bit
- Bit lines either read result or write result

Assignment Project Exam Help

<https://tutorcs.com>

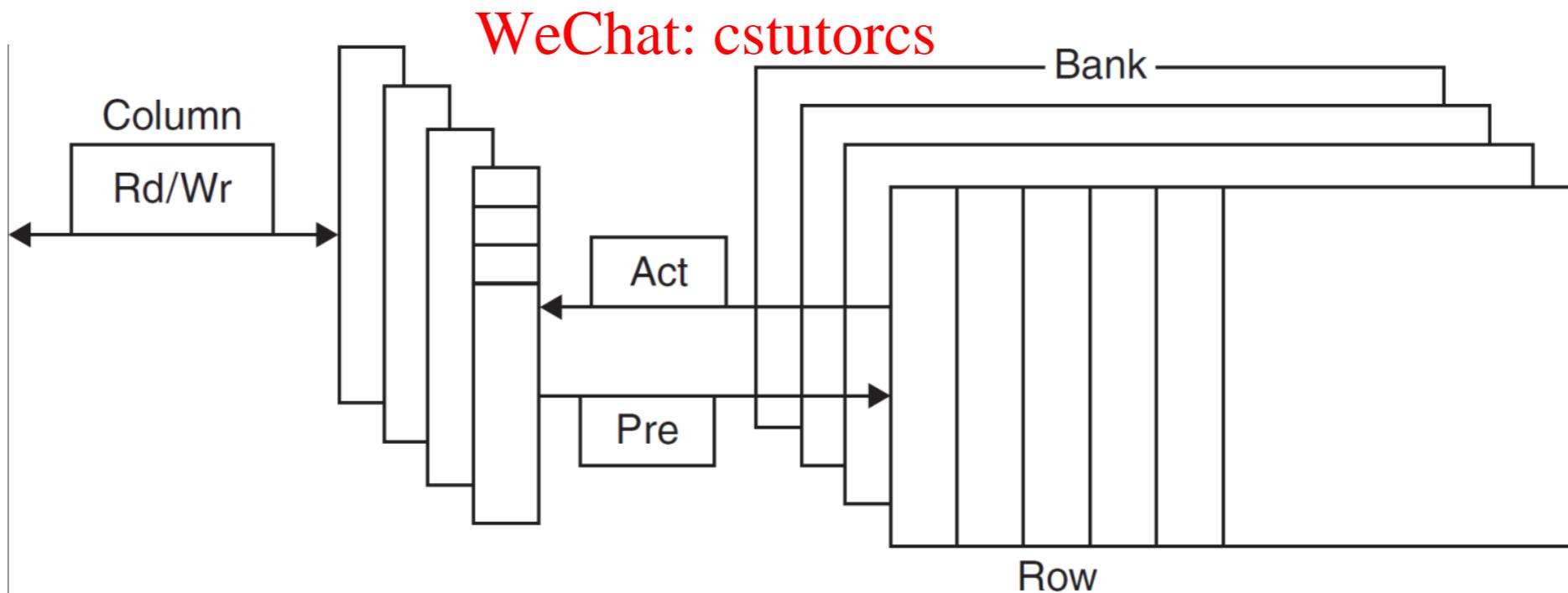


Memory Technology

- **Static RAM (SRAM)**
 - 0.5ns – 2.5ns, \$2000 – \$5000 per GB
- **Dynamic RAM (DRAM)**
 - 50ns – 70ns, \$20 – \$75 per GB
Assignment Project Exam Help
- **Magnetic disk** <https://tutorcs.com>
 - 5ms – 20ms, \$0.20 – \$2 per GB
WeChat: cstutorcs
- **Ideal memory**
 - Access time of SRAM
 - Capacity and cost/GB of disk

DRAM Technology

- Data stored as a charge in a capacitor
 - Single transistor used to access the charge
 - Must periodically be refreshed
 - Read contents and write back
Assignment Project Exam Help
 - Performed on a DRAM “row”
<https://tutorcs.com>



Advanced DRAM Organization

- Bits in a DRAM are organized as a rectangular array
 - DRAM accesses an entire row
 - Burst mode: supply successive words from a row with reduced latency
- Double data rate (DDR) DRAM
 - Transfer on rising and falling clock edges
- Quad data rate (QDR) DRAM
 - Separate DDR inputs and outputs

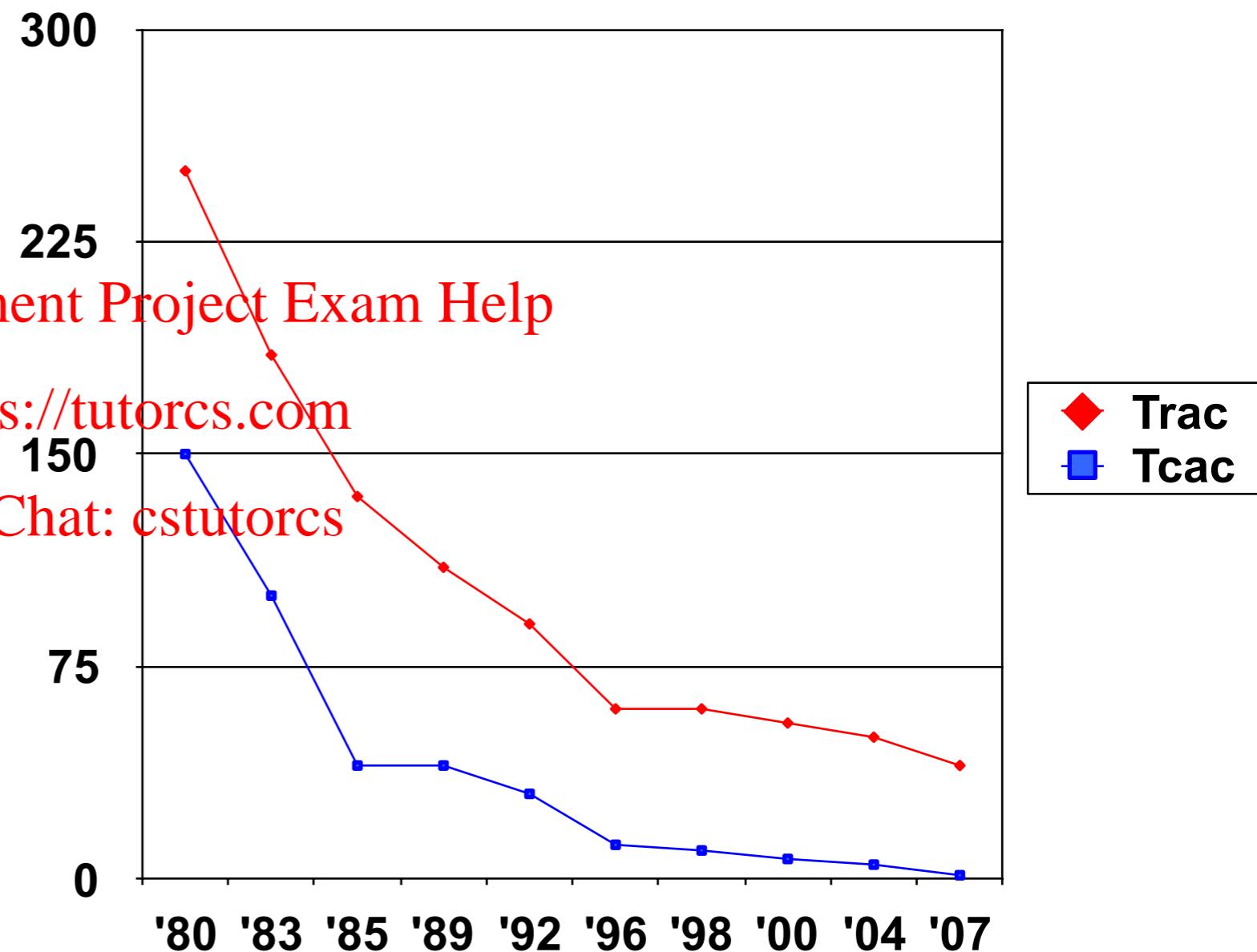
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: estutorcs

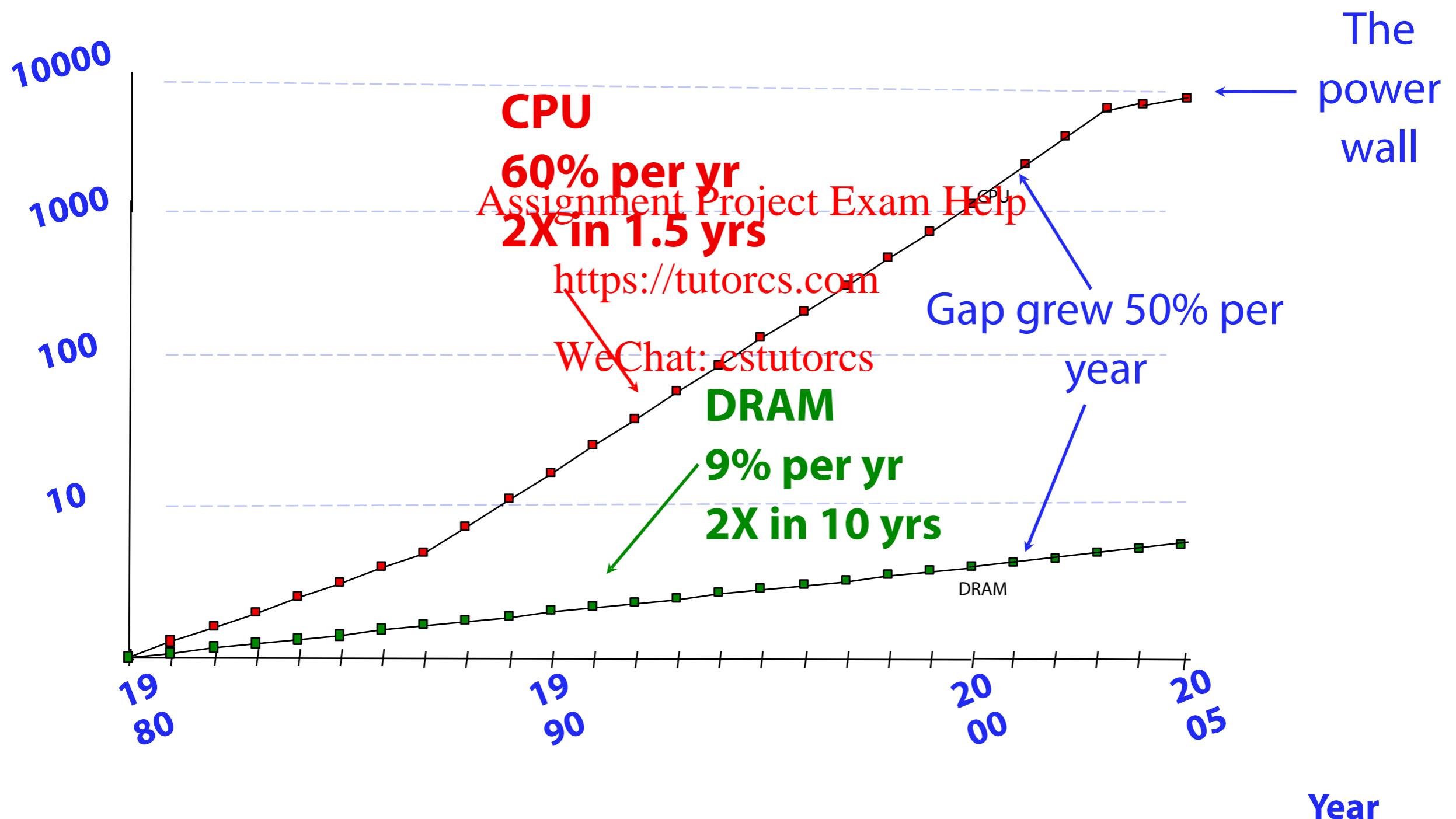
DRAM Generations

Year	Capacity	\$/GB
1980	64Kbit	\$1500000
1983	256Kbit	\$500000
1985	1Mbit	\$200000
1989	4Mbit	\$50000
1992	16Mbit	\$15000
1996	64Mbit	\$10000
1998	128Mbit	\$4000
2000	256Mbit	\$1000
2004	512Mbit	\$250
2007	1Gbit	\$50



1980–2003, CPU speed outpaced DRAM ...

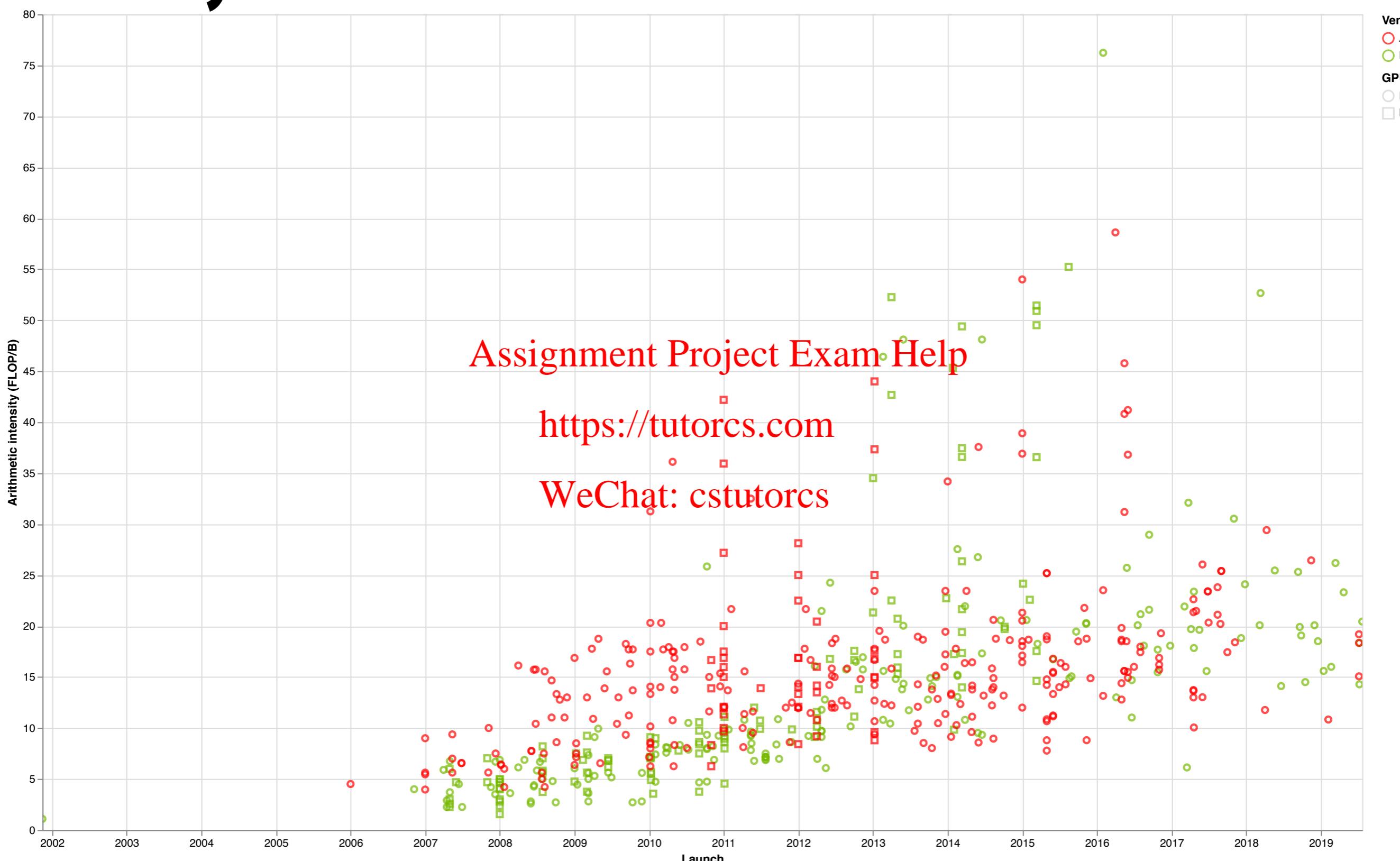
Performance
(1/latency)



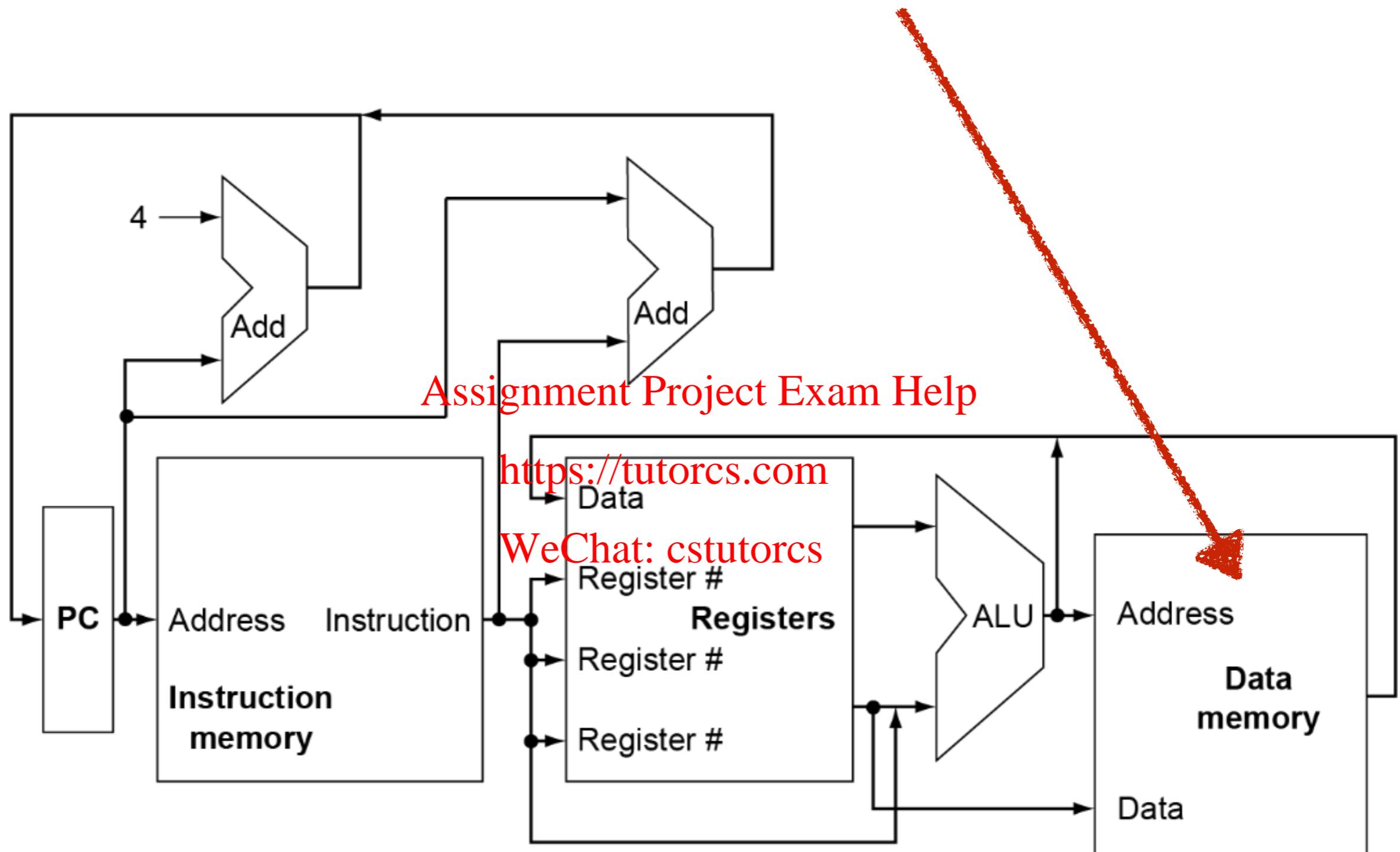
Today's Situation: Microprocessor

- Rely on caches to bridge gap
- Microprocessor-DRAM performance gap
 - time of a full cache miss in instructions executed
 - 1st Alpha (7000): $340 \text{ ns}/5.0 \text{ ns} = 68 \text{ clks} \times 2 \text{ or } 136$
Assignment Project Exam Help
instructions
 - 2nd Alpha (8400): $266 \text{ ns}/3.3 \text{ ns} = 80 \text{ clks} \times 4 \text{ or } 320$
WeChat: cstutorcs
instructions
 - 3rd Alpha (t.b.d.): $180 \text{ ns}/1.7 \text{ ns} = 108 \text{ clks} \times 6 \text{ or } 648$
instructions
 - 1/2X latency x 3X clock rate x 3X Instr/clock -> ~5X

Today's situation: GPU



What if one load cost 648 instructions?



DRAM Performance Factors

■ Row buffer

- Allows several words to be read and refreshed in parallel

■ Synchronous DRAM

- Allows for consecutive accesses in bursts without needing to send each address

Assignment Project Exam Help

- Improves bandwidth

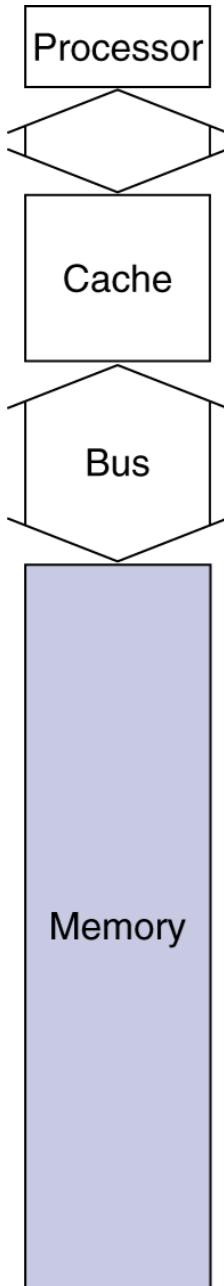
<https://tutorcs.com>

WeChat: cstutorcs

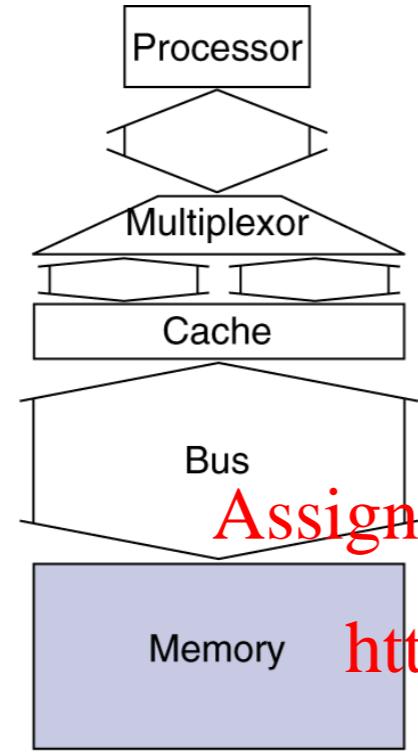
■ DRAM banking

- Allows simultaneous access to multiple DRAMs
- Improves bandwidth

Increasing Memory Bandwidth



a. One-word-wide
memory organization

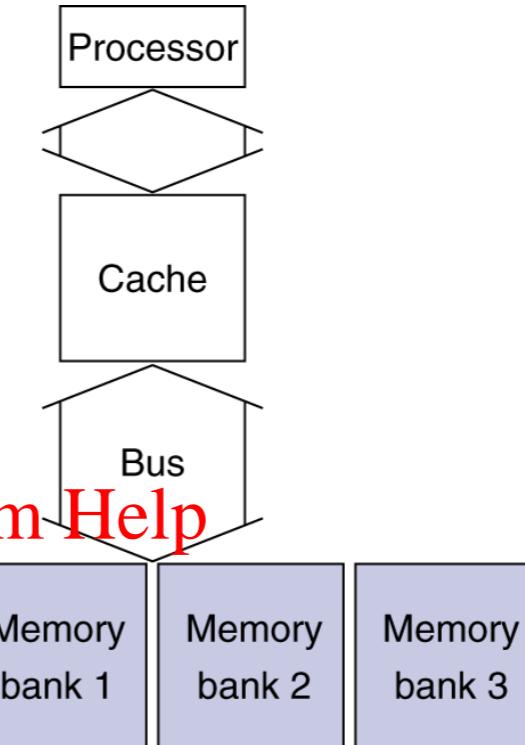


b. Wider memory organization

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutores



c. Interleaved memory organization

■ 4-word wide memory

- **Miss penalty = $1 + 15 + 1 = 17$ bus cycles**
- **Bandwidth = $16 \text{ bytes} / 17 \text{ cycles} = 0.94 \text{ B/cycle}$**

■ 4-bank interleaved memory

- **Miss penalty = $1 + 15 + 4 \times 1 = 20$ bus cycles**
- **Bandwidth = $16 \text{ bytes} / 20 \text{ cycles} = 0.8 \text{ B/cycle}$**

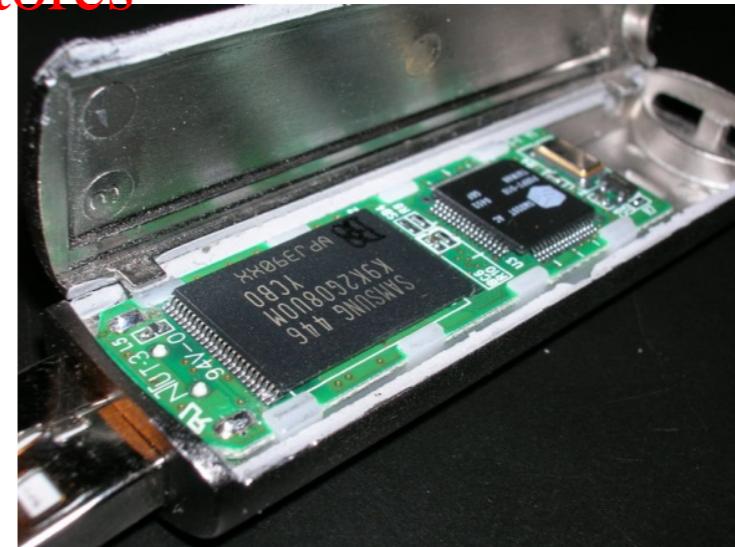
Flash Storage

■ Nonvolatile semiconductor storage

- 100× – 1000× faster than disk
 - Especially with random reads
- Smaller, lower power, more robust
Assignment Project Exam Help
- But more \$/GB (between disk and DRAM)
<https://tutotorcs.com>



WeChat: cstutorcs



Flash Types

- NOR flash: bit cell like a NOR gate
 - Random read/write access
 - Used for instruction memory in embedded systems
- NAND flash: bit cell like a NAND gate
 - Denser (bits/area), but block-at-a-time access
 - Cheaper per GB
 - Used for USB keys, media storage, ...
- Flash bits wear out after 1000's of accesses
 - Not suitable for direct RAM or disk replacement
 - Wear leveling: remap data to less used blocks

Disk Storage

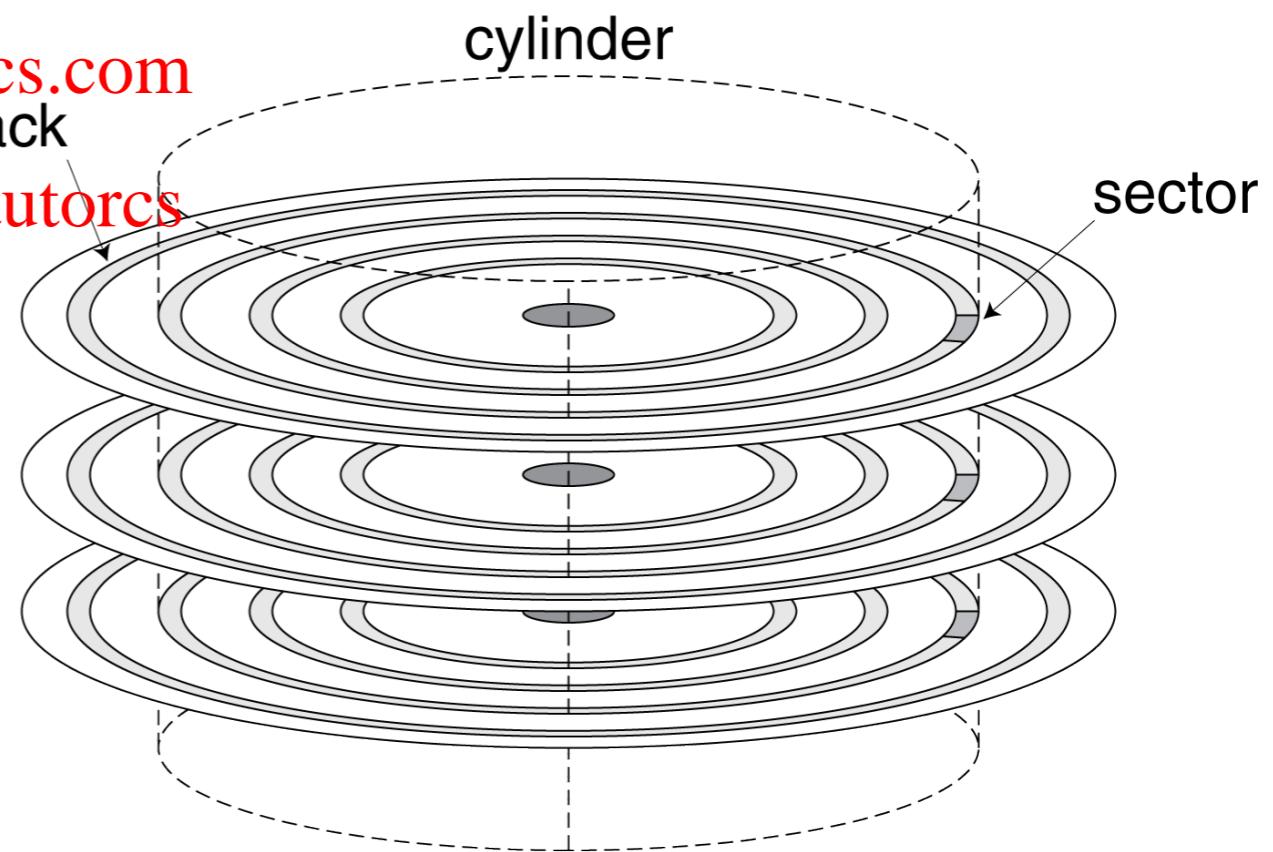
- Nonvolatile, rotating magnetic storage



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



Disk Sectors and Access

- **Each sector records**
 - Sector ID
 - Data (512 bytes, 4096 bytes proposed)
 - Error correcting code (ECC)
 - Used to hide defects and recording errors
 - Synchronization fields
- **Access to a sector involves** Assignment Project Exam Help WeChat: cstutorcs
 - Queuing delay if other accesses are pending
 - Seek: move the heads
 - Rotational latency
 - Data transfer
 - Controller overhead

Disk Access Example

- Given
 - 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk
- Average read time
 - 4ms seek time
 - + $\frac{1}{2} / (15,000/60) = 2\text{ms rotational latency}$
 - + $512 / 100\text{MB/s} = 0.005\text{ms transfer time}$
 - + 0.2ms controller delay
 - = 6.2ms
- If actual average seek time is 1ms
 - Average read time = 3.2ms

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Disk Performance Issues

- Manufacturers quote average seek time
 - Based on all possible seeks
 - Locality and OS scheduling lead to smaller actual average seek times
- Smart disk controller allocate physical sectors on disk
 - Present logical sector interface to host
 - SCSI, ATA, SATA
- Disk drives include caches
 - Prefetch sectors in anticipation of access
 - Avoid seek and rotational delay

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutors

The Goal: illusion of large, fast, cheap memory

■ Fact:

- Large memories are slow
- Fast memories are small

■ Goal:

Assignment Project Exam Help

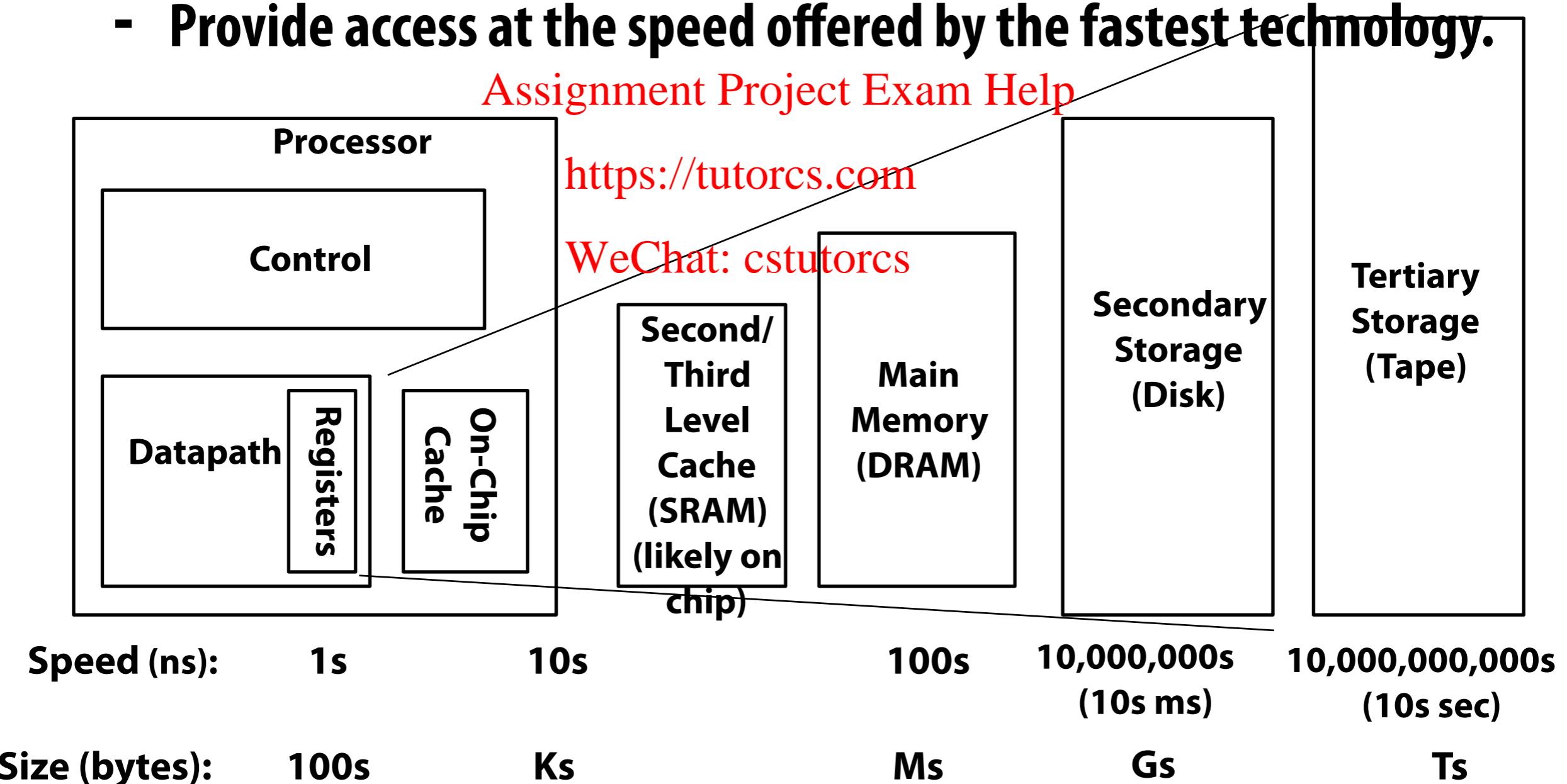
- Large, fast memory <https://tutorcs.com>

■ How do we create a memory that is large, cheap and fast (most of the time)?

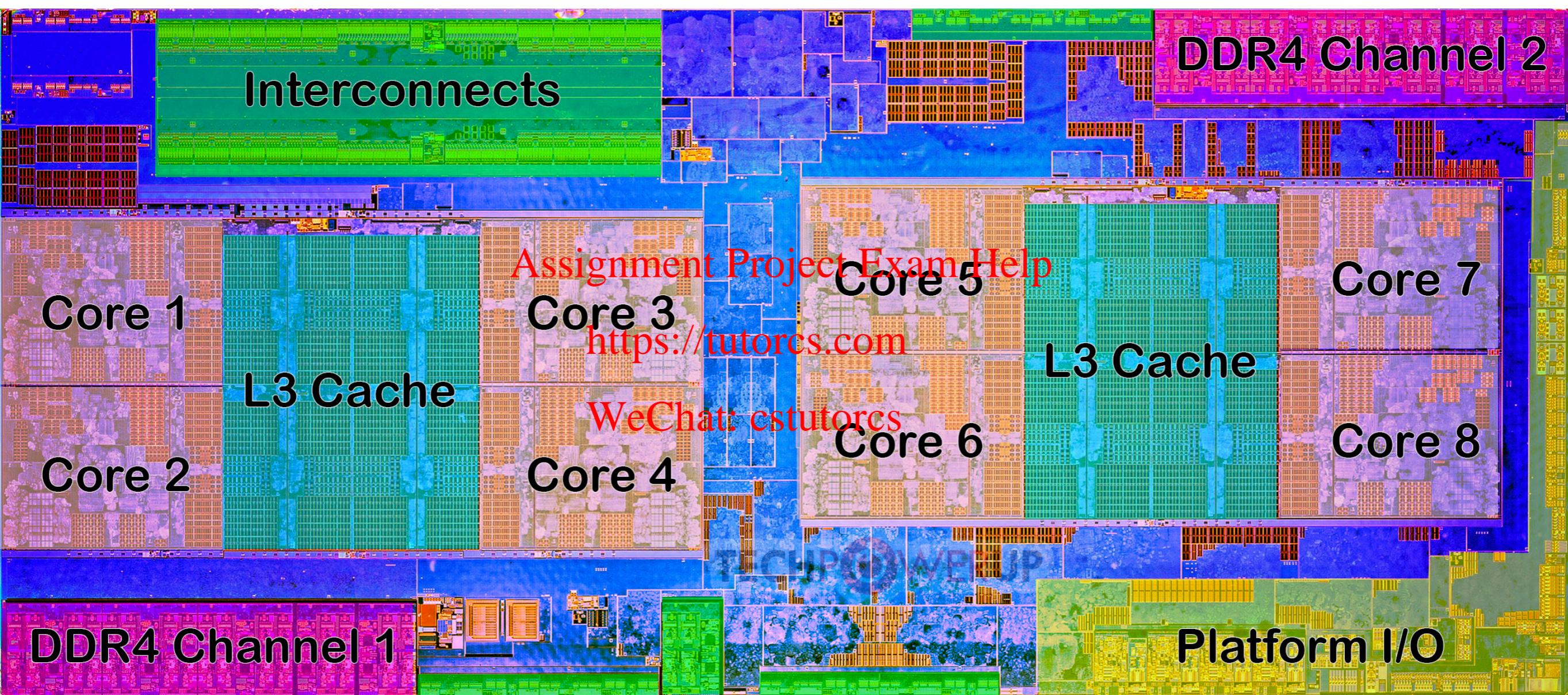
- Hierarchy
 - Have both large/slow and small/fast memories!
- Parallelism

Memory Hierarchy of a Modern Computer System

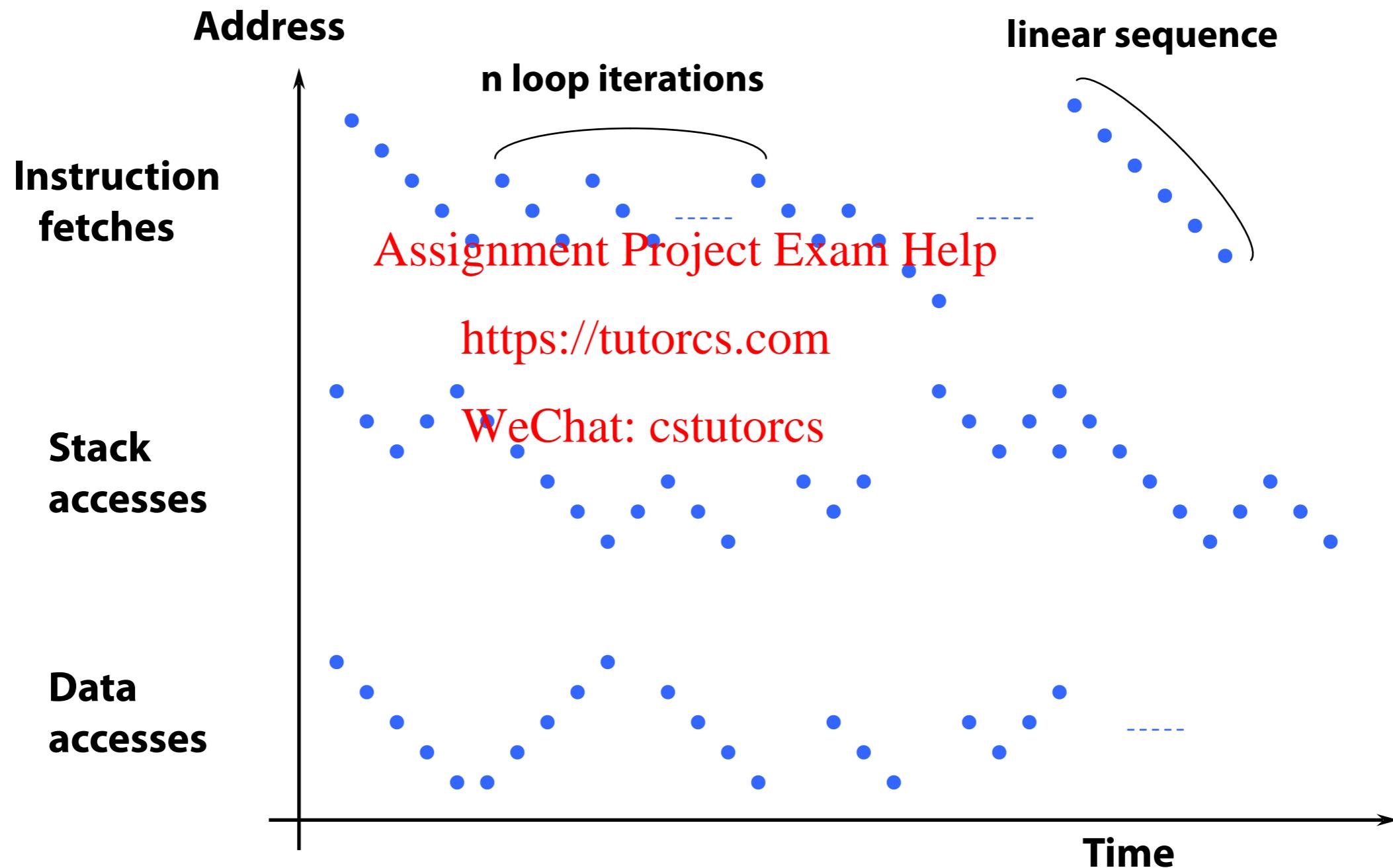
- By taking advantage of the principle of *locality*:
 - Present the user with as much memory as is available in the cheapest technology.
 - Provide access at the speed offered by the fastest technology.



AMD Ryzen



Typical Memory Reference Patterns



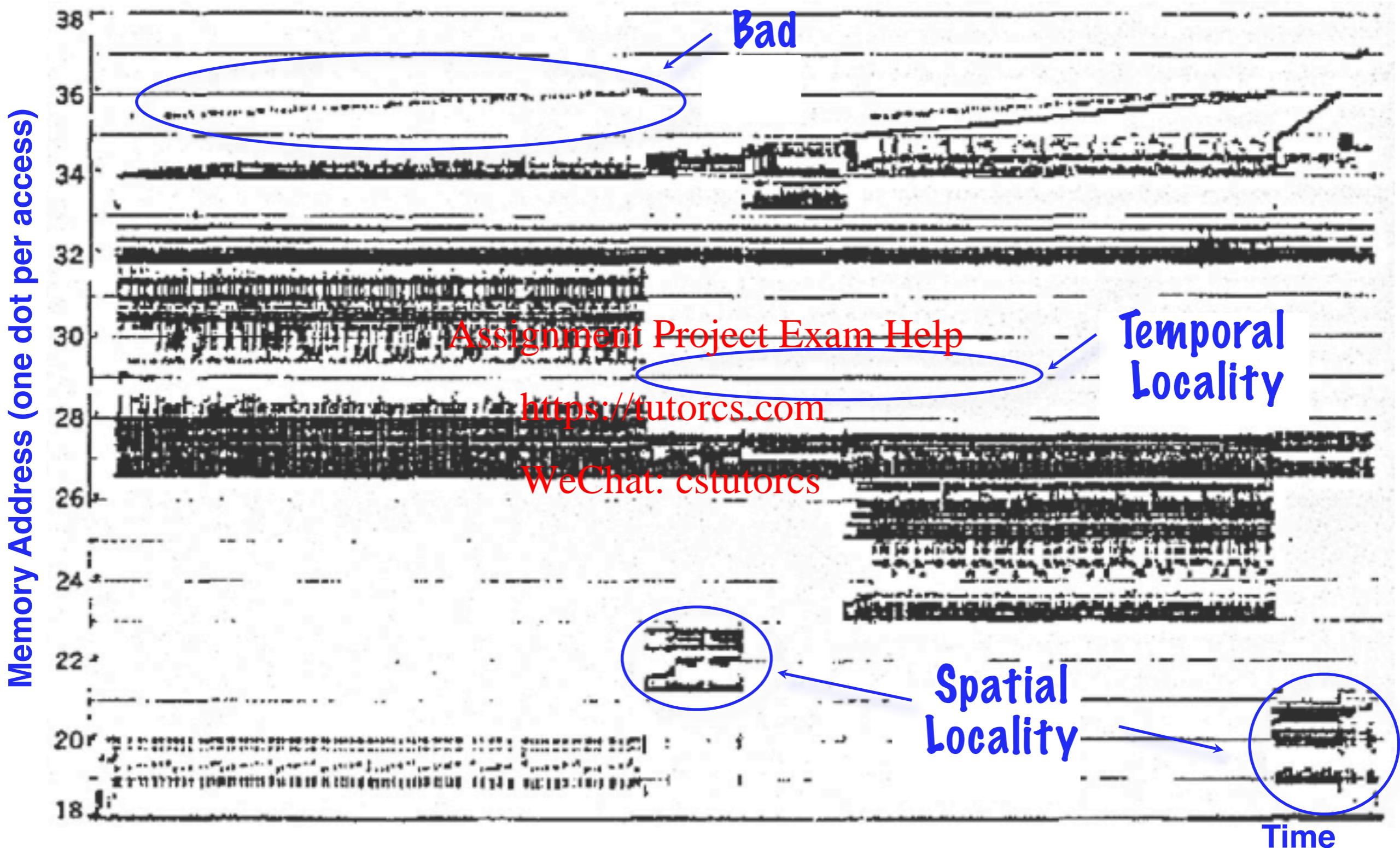
Locality Properties of Patterns

- Two locality properties of memory references:
 - **Temporal Locality:** If a location is referenced it is likely to be referenced again in the near future.
Assignment Project Exam Help
<https://tutorcs.com>
 - **Spatial Locality:** If a location is referenced it is likely that locations near it will be referenced in the near future.

Caches

- Caches exploit both properties of patterns.
 - Exploit *temporal* locality by remembering the contents of recently accessed locations.
Assignment Project Exam Help
<https://tutorcs.com>
 - Exploit *spatial* locality by remembering blocks of contents of recently accessed locations.
WeChat: cstutorcs

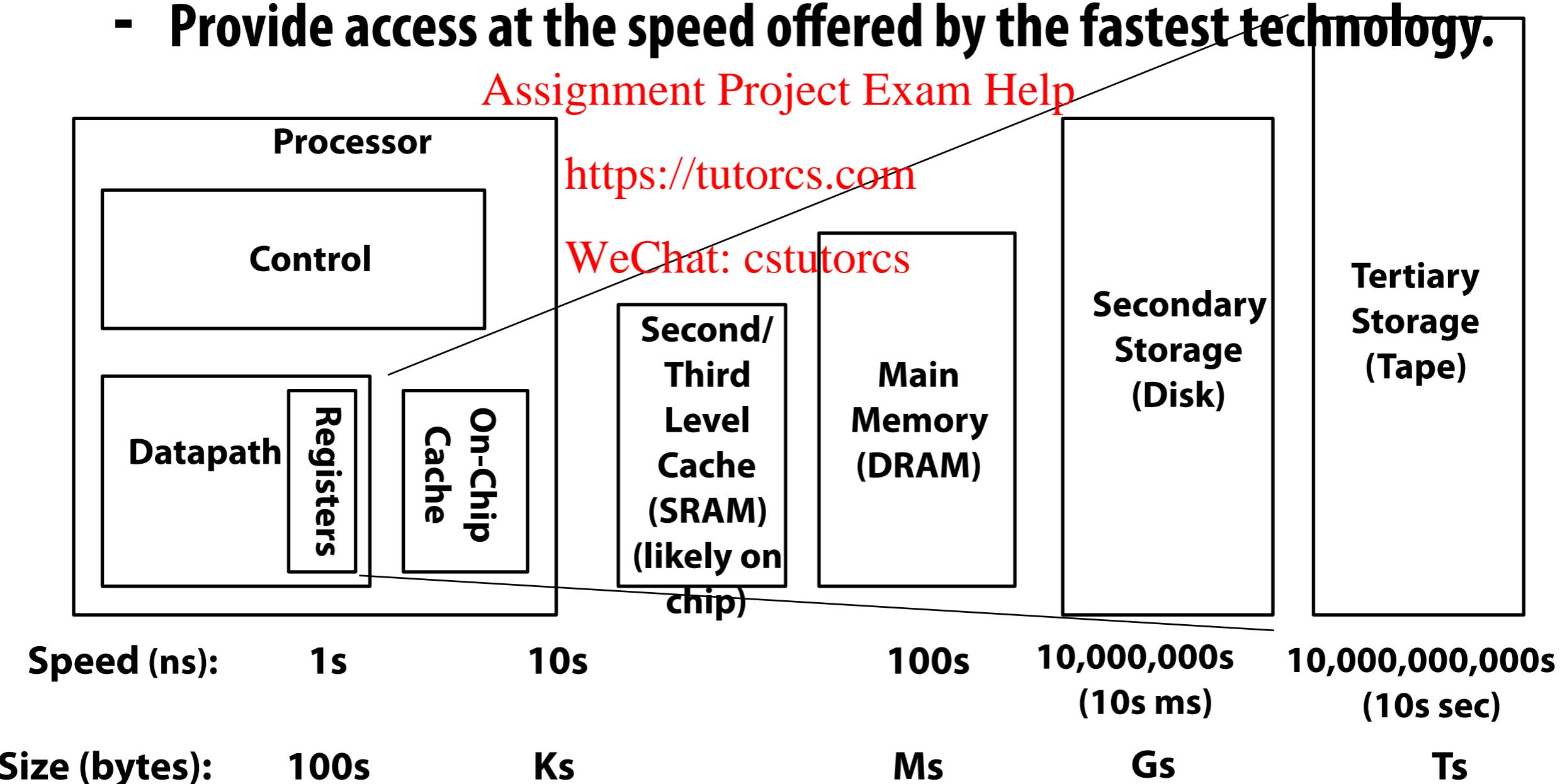
Programs with locality cache well ...



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

Memory Hierarchy of a Modern Computer System

- By taking advantage of the principle of *locality*:
 - Present the user with as much memory as is available in the cheapest technology.
 - Provide access at the speed offered by the fastest technology.

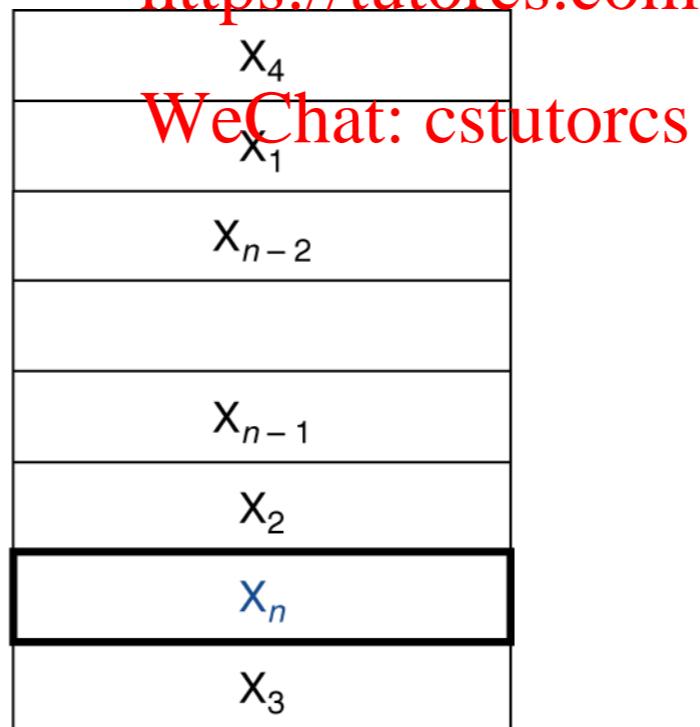
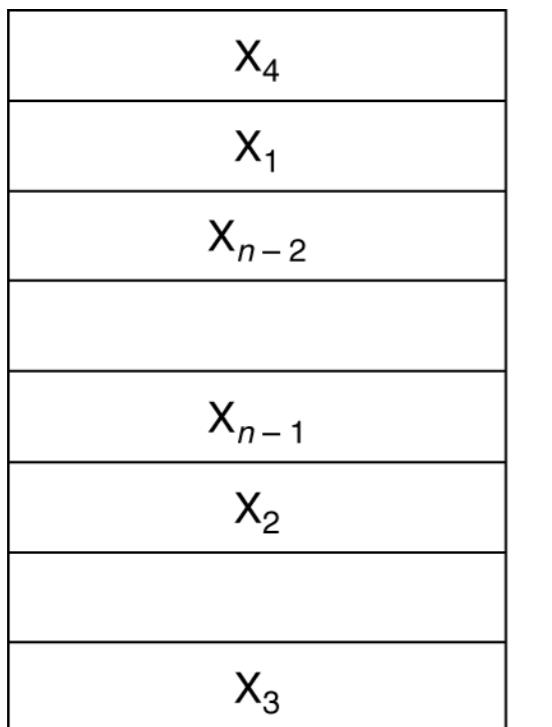


Cache Memory

- Cache memory
 - The level of the memory hierarchy closest to the CPU
- Given accesses X_1, \dots, X_{n-1}, X_n

Assignment Project Exam Help

<https://tutorcs.com>



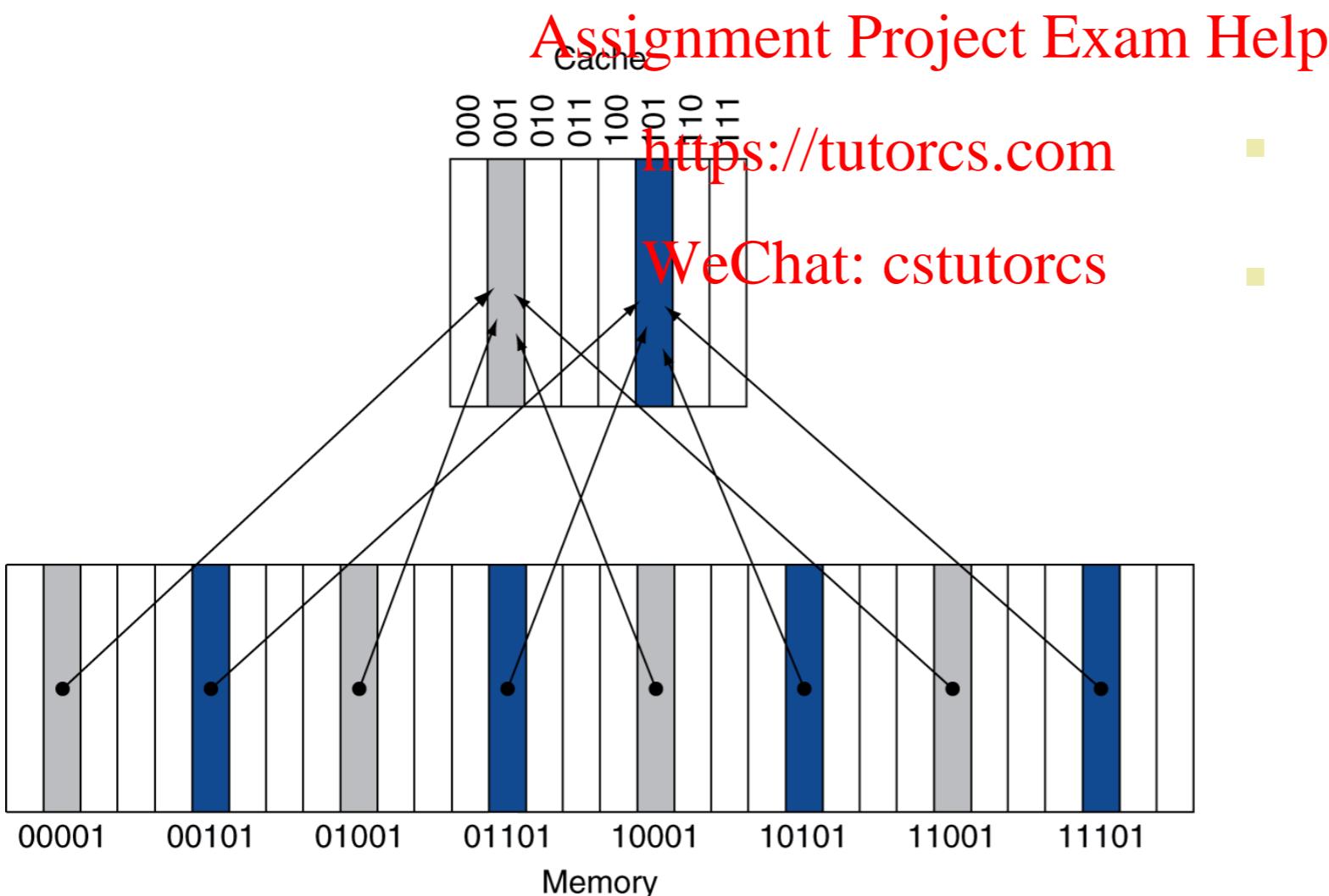
a. Before the reference to X_n b. After the reference to X_n

How do we know if the data is present?

Where do we look?

Direct Mapped Cache

- Location determined by address
- Direct mapped: only one choice
 - (Block address) modulo (#Blocks in cache)



- *#Blocks is a power of 2*
- *Use low-order address bits*

Tags and Valid Bits

- How do we know which particular block is stored in a cache location?

- Store block address as well as the data
 - Actually, only need the high-order bits
Assignment Project Exam Help
 - Called the tag
<https://tutorcs.com>

- What if there is no data in a location?
WeChat: estutorcs

- Valid bit: 1 = present, 0 = not present
 - Initially 0

Cache Example

- **8-blocks, 1 word/block, direct mapped**
- **Initial state**

Index	V	Tag	Data
000	N	Assignment Project Exam Help https://tutorcs.com	
001	N	WeChat: cstutorcs	
010	N		
011	N		
100	N		
101	N		
110	N		
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Miss	110

Assignment Project Exam Help
<https://tutorcs.com>

Index	V	Tag	Data
000	N		
001	N	WeChat: cstutorcs	
010	N		
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
26	11 010	Miss	010

Assignment Project Exam Help
<https://tutorcs.com>

Index	V	Tag	Data
000	N		
001	N	WeChat: cstutorcs	
010	Y	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Hit	110
26	11 010	Hit	010

Assignment Project Exam Help
<https://tutorcs.com>

Index	V	Tag	Data
000	N		
001	N	WeChat: cstutorcs	
010	Y	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
16	10 000	Miss	000
3	00 011	Miss	011
16	10 000	Hit	000

Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N	WeChat: cstutorcs	
010	Y	11	Mem[11010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

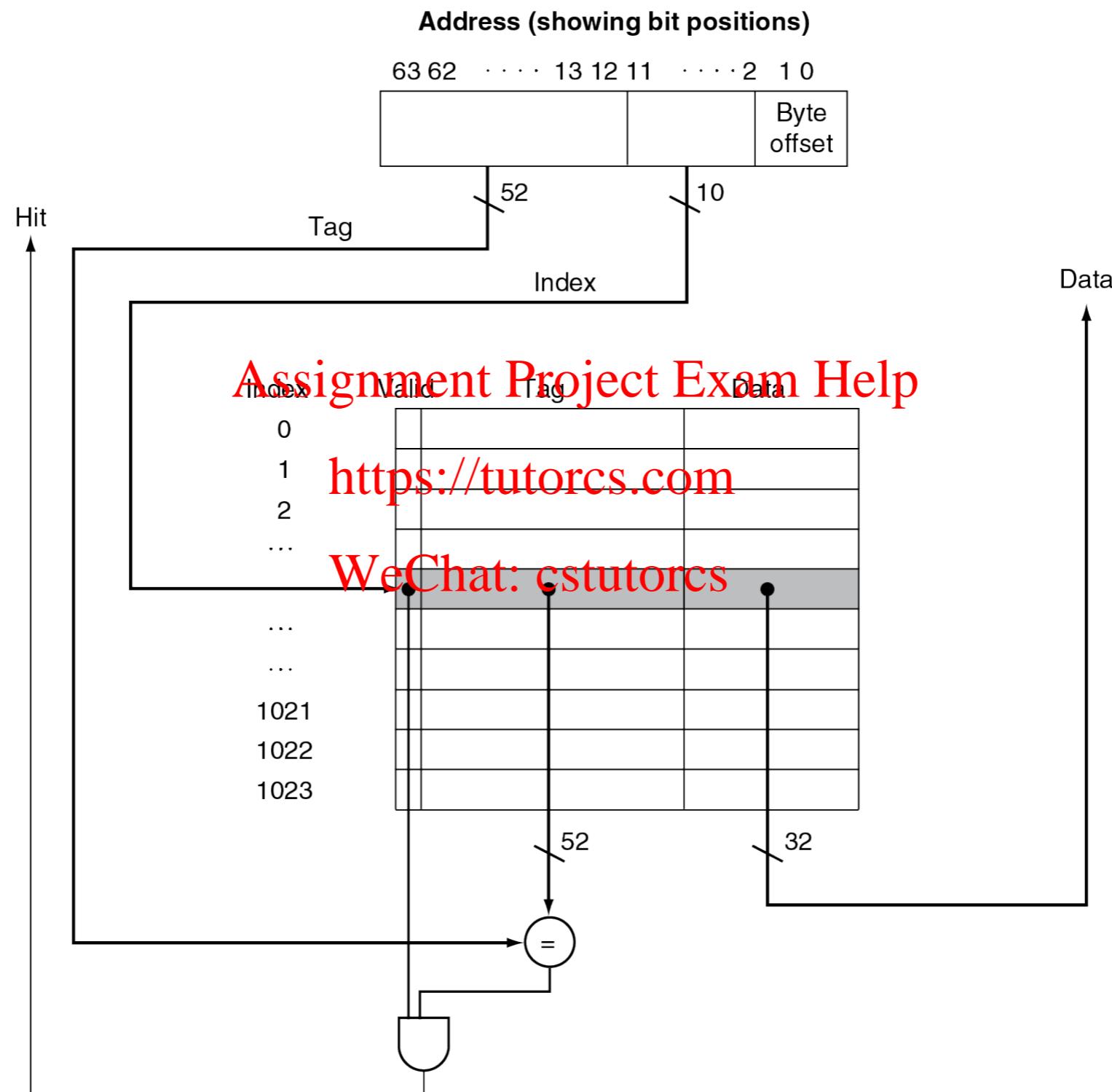
Cache Example

Word addr	Binary addr	Hit/miss	Cache block
18	10 010	Miss	010

Assignment Project Exam Help
<https://tutorcs.com>

Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N	WeChat: cstutorcs	
010	Y	10	Mem[10010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

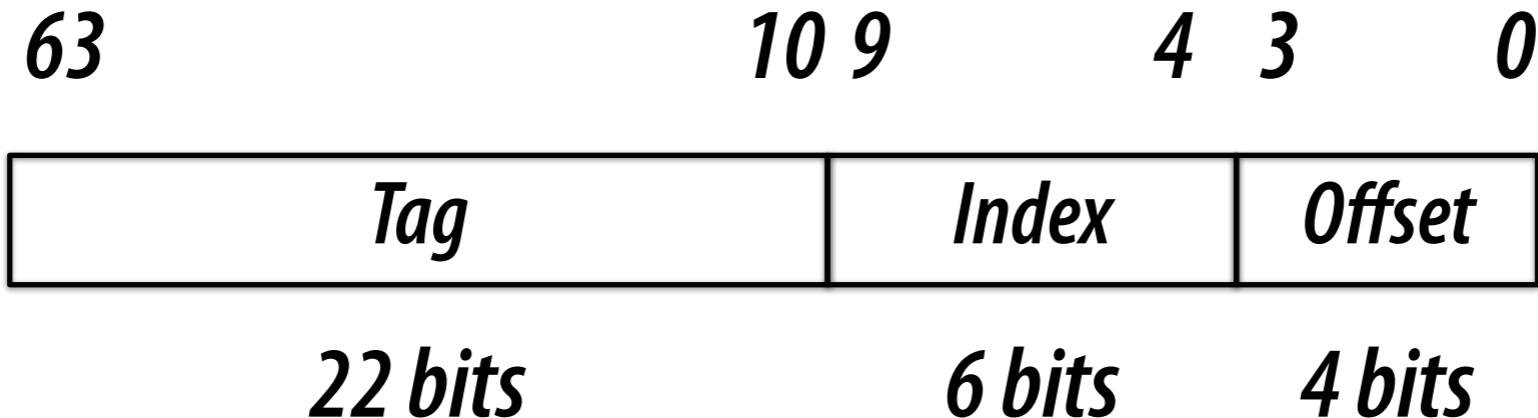
Address Subdivision



Example: Larger Block Size

- 64 blocks, 16 bytes/block
 - To what block number does address 1200 map?
- Block address = $\lfloor 1200/16 \rfloor = 75$
 - I am the 75th 16-byte block in the address space
- Block number = 75 modulo 64 = 11

WeChat: cstutorcs



Block Size Considerations

- Larger blocks should reduce miss rate
 - Due to spatial locality
- But in a fixed-sized cache
 - Larger blocks \Rightarrow fewer of them
Assignment Project Exam Help
 - More competition <https://tutorcs.com> \Rightarrow increased miss rate
WeChat: cstutorcs
 - Larger blocks \Rightarrow pollution
- Larger miss penalty
 - Can override benefit of reduced miss rate
 - Early restart and critical-word-first can help

Cache Misses

- On cache hit, CPU proceeds normally
 - This is (hopefully) the common case
 - This is one or a few cycles at most
- On cache miss
 - Stall the CPU pipeline <https://tutorcs.com>
 - Fetch block from ~~next level of hierarchy~~ <https://chat.csail.mit.edu>
 - Instruction cache miss
 - Restart instruction fetch
 - Data cache miss
 - Complete data access

Assignment Project Exam Help

Write-Through

- On data-write hit, could just update the block in cache
 - But then cache and memory would be inconsistent
- Write through: also update memory
- But makes writes take longer
 - e.g., if base CPI = 1, 10% of instructions are stores, write to memory takes 100 cycles
 - Effective CPI = $1 + 0.1 \times 100 = 11$
- Solution: write buffer
 - Holds data waiting to be written to memory
 - CPU continues immediately
 - Only stalls on write if write buffer is already full

Write-Back

- Alternative: On data-write hit, just update the block in cache
 - Keep track of whether each block is dirty
- When a dirty block is replaced
 - Write it back to memory
Assignment Project Exam Help
 - Can use a write buffer to allow replacing block to be read first
<https://tutorcs.com>
WeChat: cstutorcs

Write Allocation

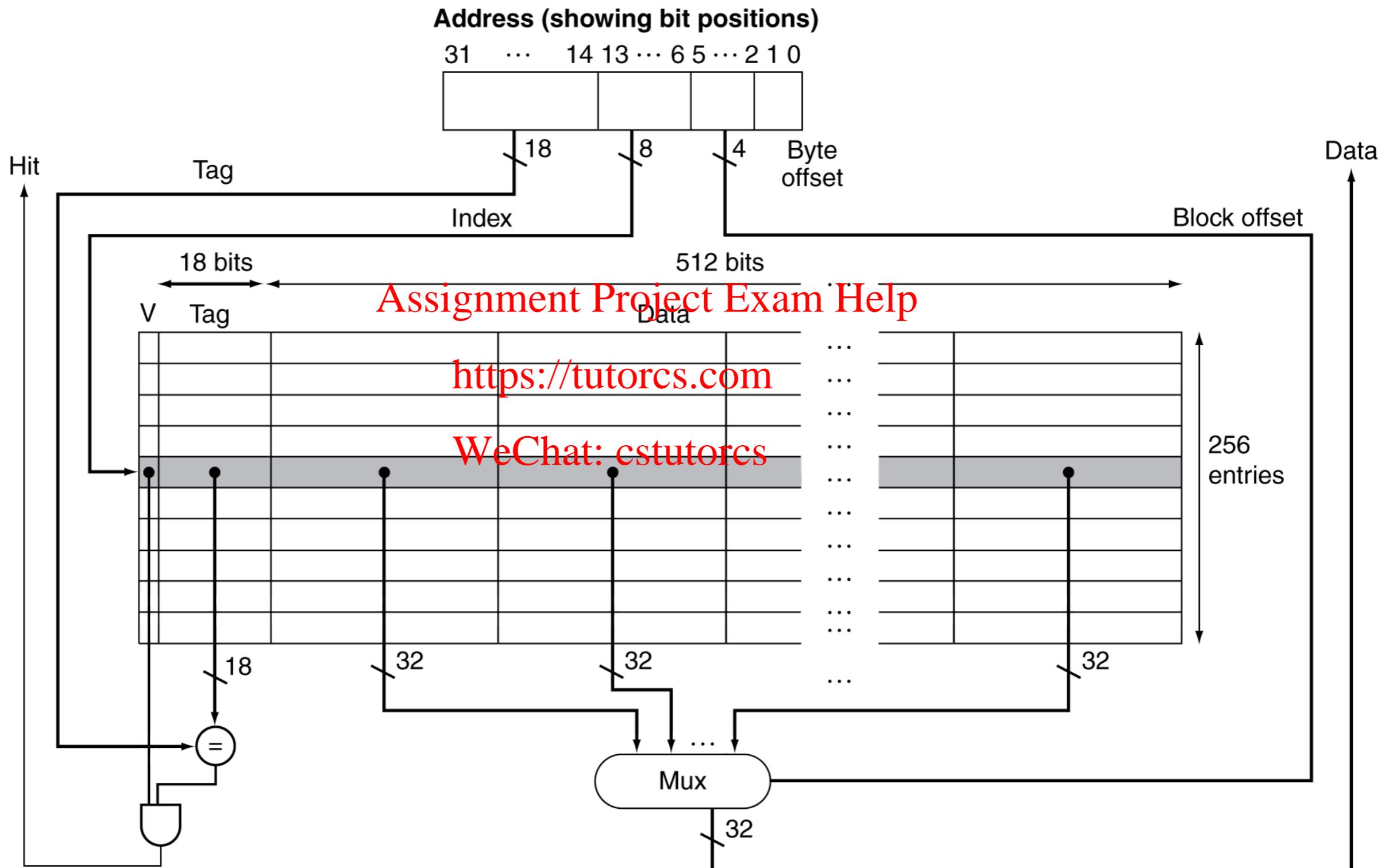
- What should happen on a write miss?
- Alternatives for write-through
 - Allocate on miss: fetch the block
 - Write around: don't fetch the block
 - Assignment Project Exam Help
 - Since programs often write a whole block before reading it (e.g., initialization)
https://tutorcs.com
WeChat: cstutorcs
- For write-back
 - Usually fetch the block

Example: Intrinsity FastMATH

- Embedded MIPS processor
 - 12-stage pipeline
 - Instruction and data access on each cycle
- Split cache: separate I-cache and D-cache
 - Each 16KB: 256 blocks \times 16 words/block
 - D-cache: write-through or write-back
- SPEC2000 miss rates
 - I-cache: 0.4%
 - D-cache: 11.4%
 - Weighted average: 3.2%

Assignment Project Exam Help
<https://tutorcs.com>

Example: Intrinsity FastMATH



Main Memory Supporting Caches

■ Use DRAMs for main memory

- Fixed width (e.g., 1 word)
- Connected by fixed-width clocked bus
 - Bus clock is typically slower than CPU clock

■ Example cache block read

- 1 bus cycle for address transfer
- 15 bus cycles per DRAM access
- 1 bus cycle per data transfer

■ For 4-word block, 1-word-wide DRAM

- Miss penalty = $1 + 4 \times 15 + 4 \times 1 = 65$ bus cycles
- Bandwidth = 16 bytes / 65 cycles = 0.25 B/cycle

Summary

- **Two Different Types of Locality:**
 - **Temporal Locality (Locality in Time): If an item is referenced, it will tend to be referenced again soon.**
 - **Spatial Locality (Locality in Space): If an item is referenced, items whose addresses are close by tend to be referenced soon.**
- **By taking advantage of the principle of locality:**
 - **Present the user with as much memory as is available in the cheapest technology.**
 - **Provide access at the speed offered by the fastest technology.**
- **DRAM is slow but cheap and dense:**
 - **Good choice for presenting the user with a BIG memory system**
- **SRAM is fast but expensive and not very dense:**
 - **Good choice for providing the user FAST access time.**

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs