

EE671: VLSI DESIGN

SPRING 2024/25

LAXMEESHA SOMAPPA
DEPARTMENT OF ELECTRICAL ENGINEERING
IIT BOMBAY
laxmeesha@ee.iitb.ac.in

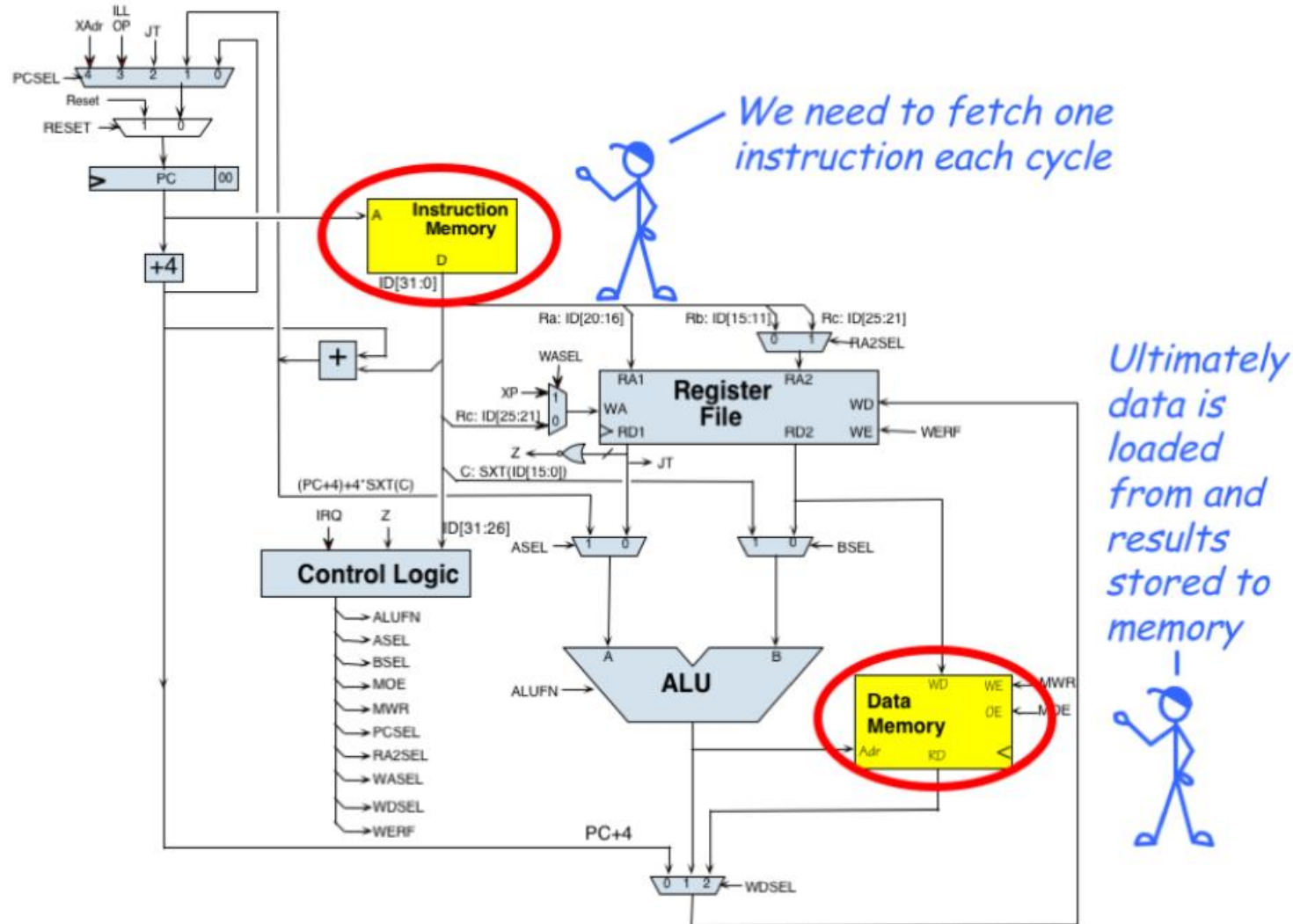


LECTURE – 31

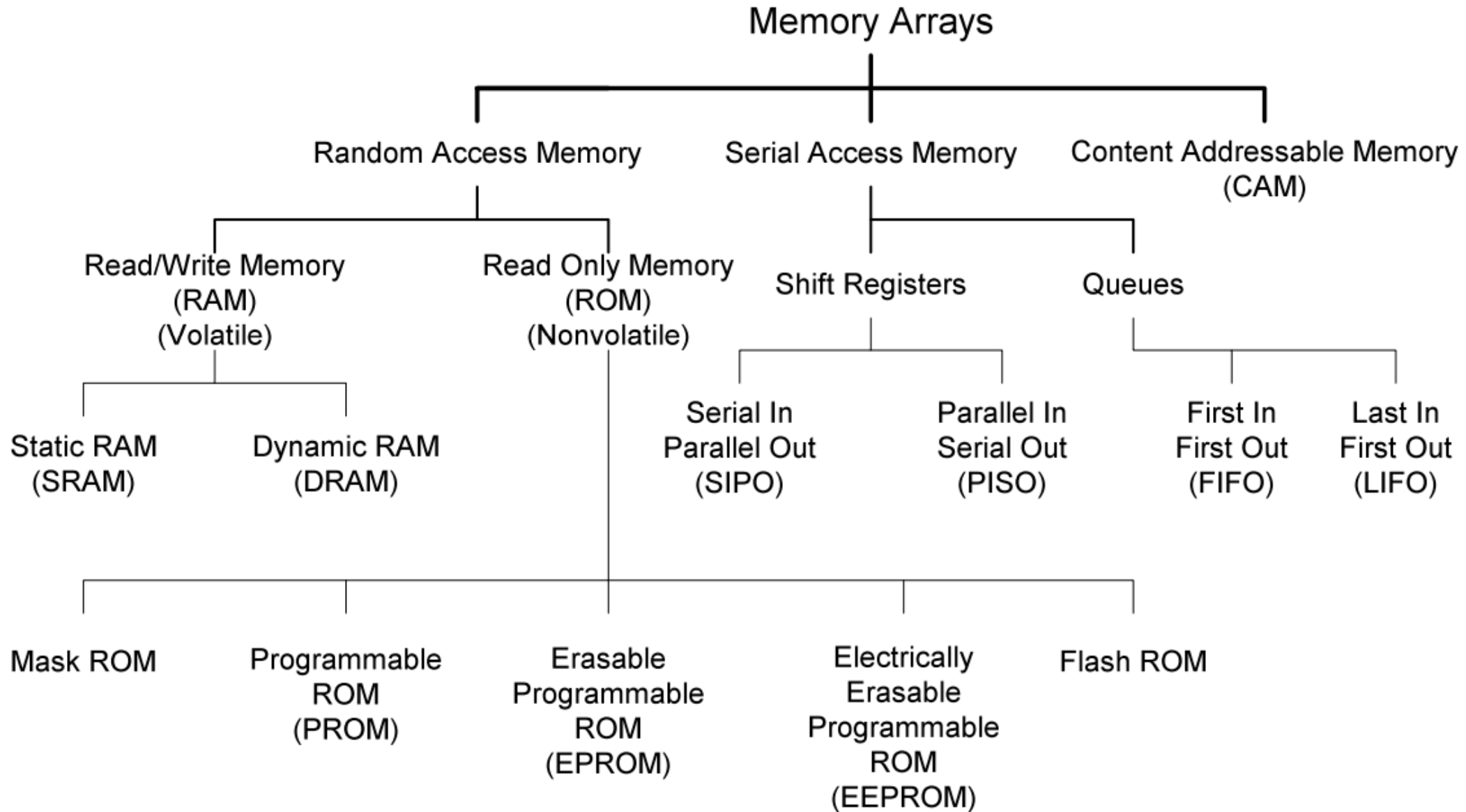
MEMORY – HIERARCHY & IMPLEMENTATION

A GENERAL PROCESSOR OVERVIEW

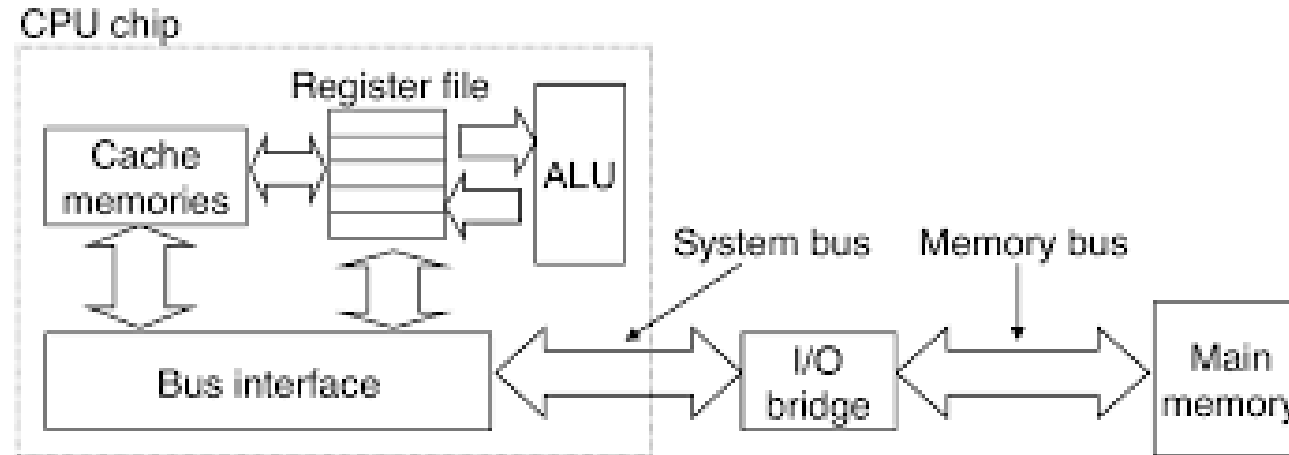
Our “~~Memory~~ Computing Machine”



MEMORY ARRAY OPTIONS



MEMORY HIERARCHY IN PROCESSORS



	Capacity	Latency	Cost/GB
Register	1000s of bits	20 ps	\$\$\$\$
SRAM	~10 KB-10 MB	1-10 ns	~\$1000
DRAM	~10 GB	80 ns	~\$10
Flash*	~100 GB	100 us	~\$1
Hard disk*	~1 TB	10 ms	~\$0.10

**Processor
Datapath**

**Memory
Hierarchy**

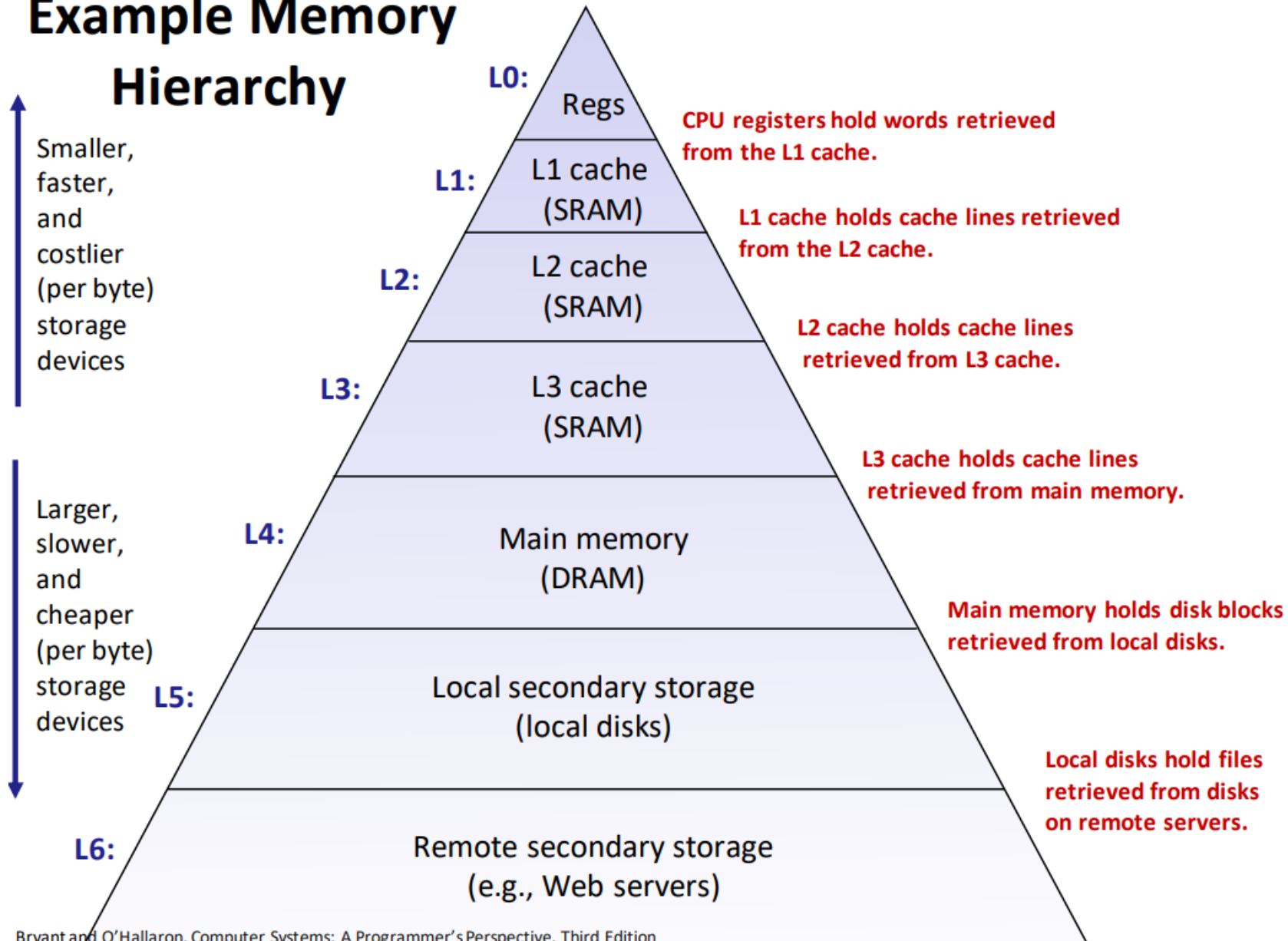
**I/O
subsystem**

* non-volatile (retains contents when powered off)



MEMORY HIERARCHY FULL PICTURE

Example Memory Hierarchy



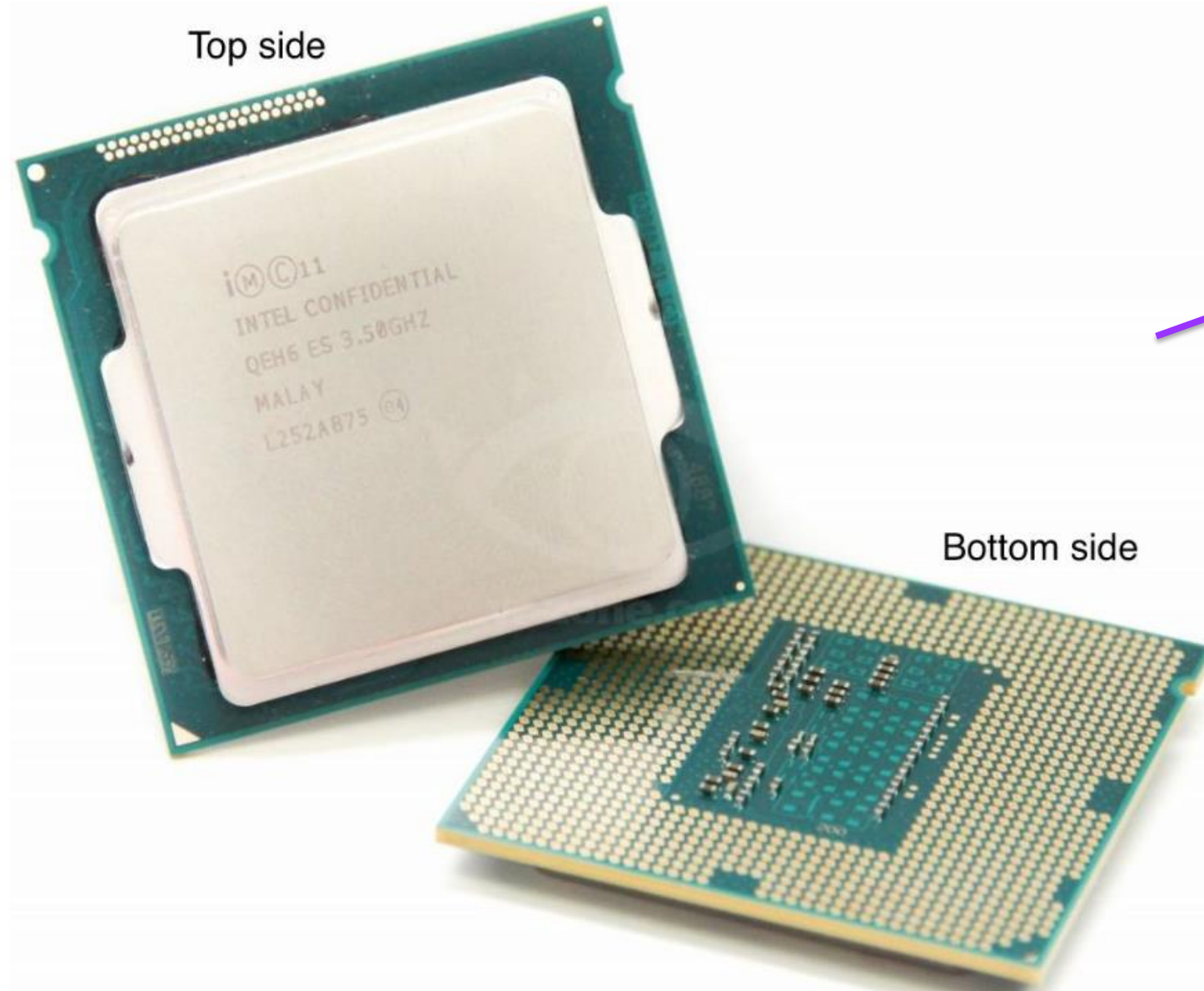
Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

laxmeesha@ee.iitb.ac.in

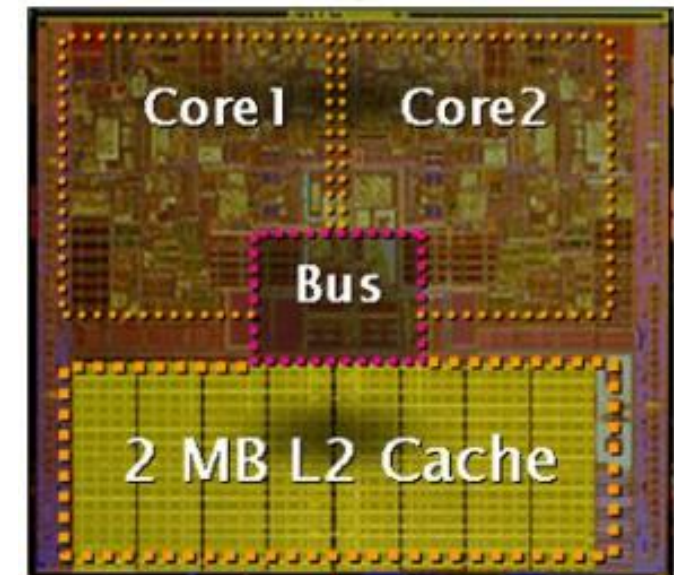
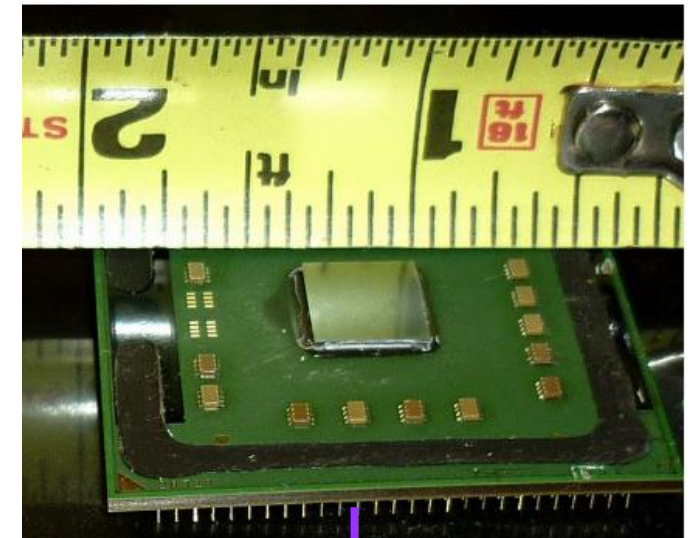


PROCESSOR PACKAGED

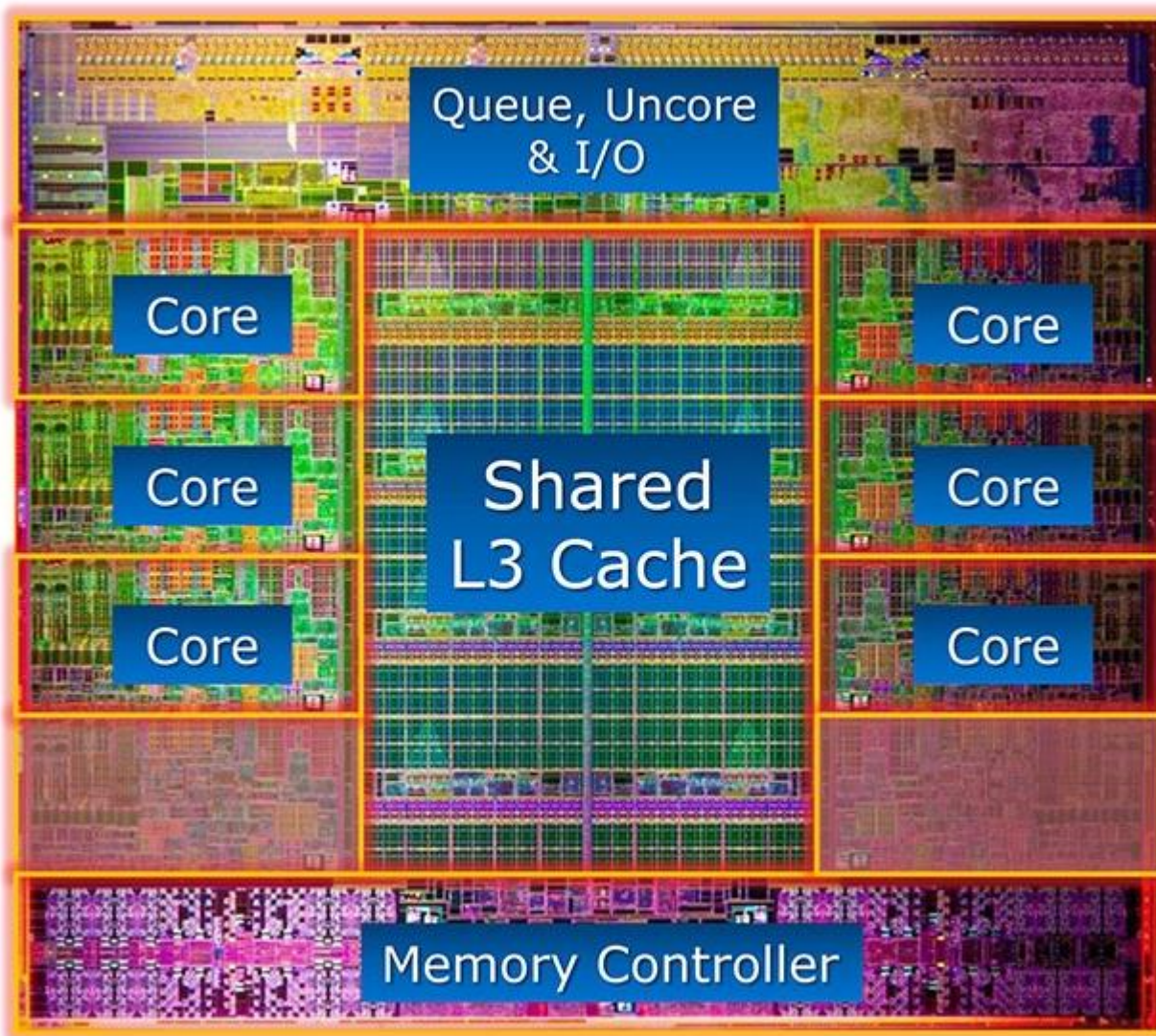
Top side



Bottom side



PROCESSOR CACHE AREA

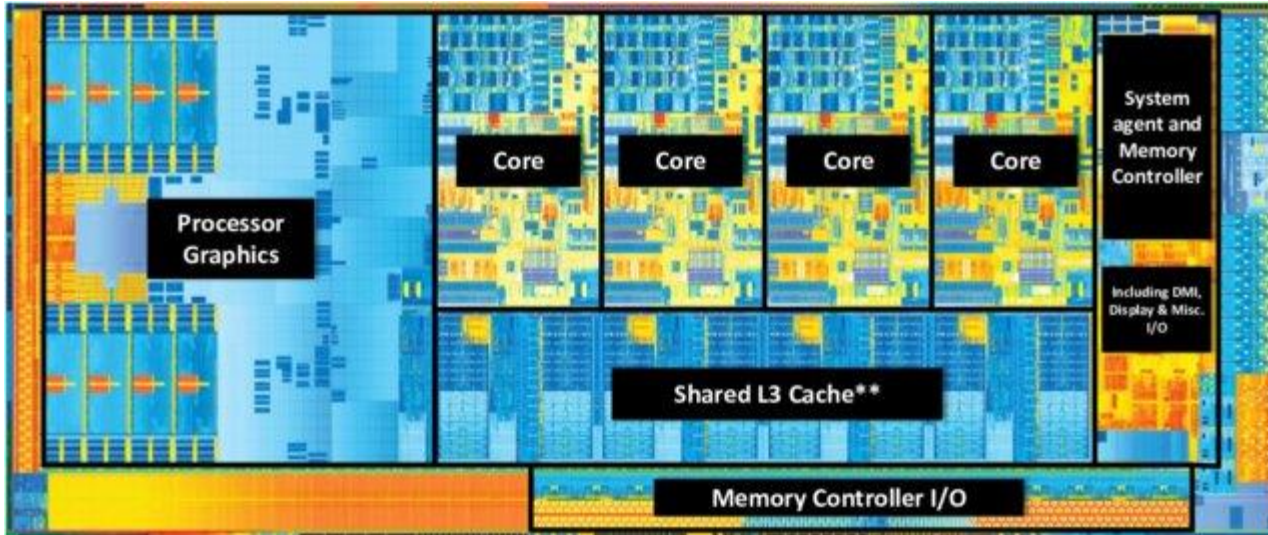


Intel® Core™ i7-3960X Processor Die Detail

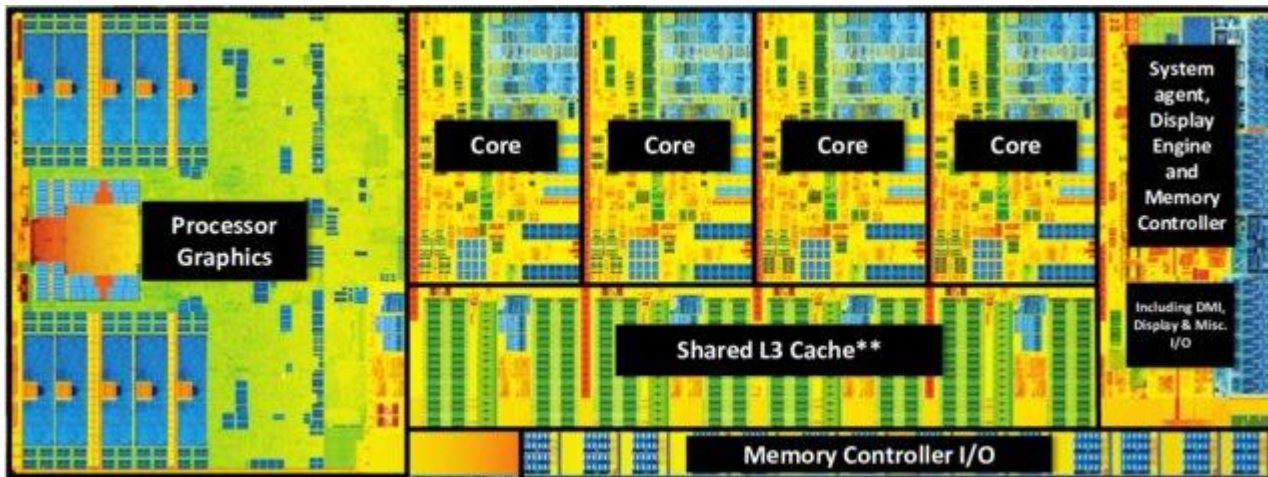
Takeaway:

- Cache (L3, L2 and L1) contributes > 60~70% of the chip area !!!
- Cache memory: has to be very dense to pack more
- Extremely compact layout !!!

PROCESSOR CACHE AREA



Quad-core Ivy Bridge



Quad-core Haswell

PROCESSOR CACHE AREA



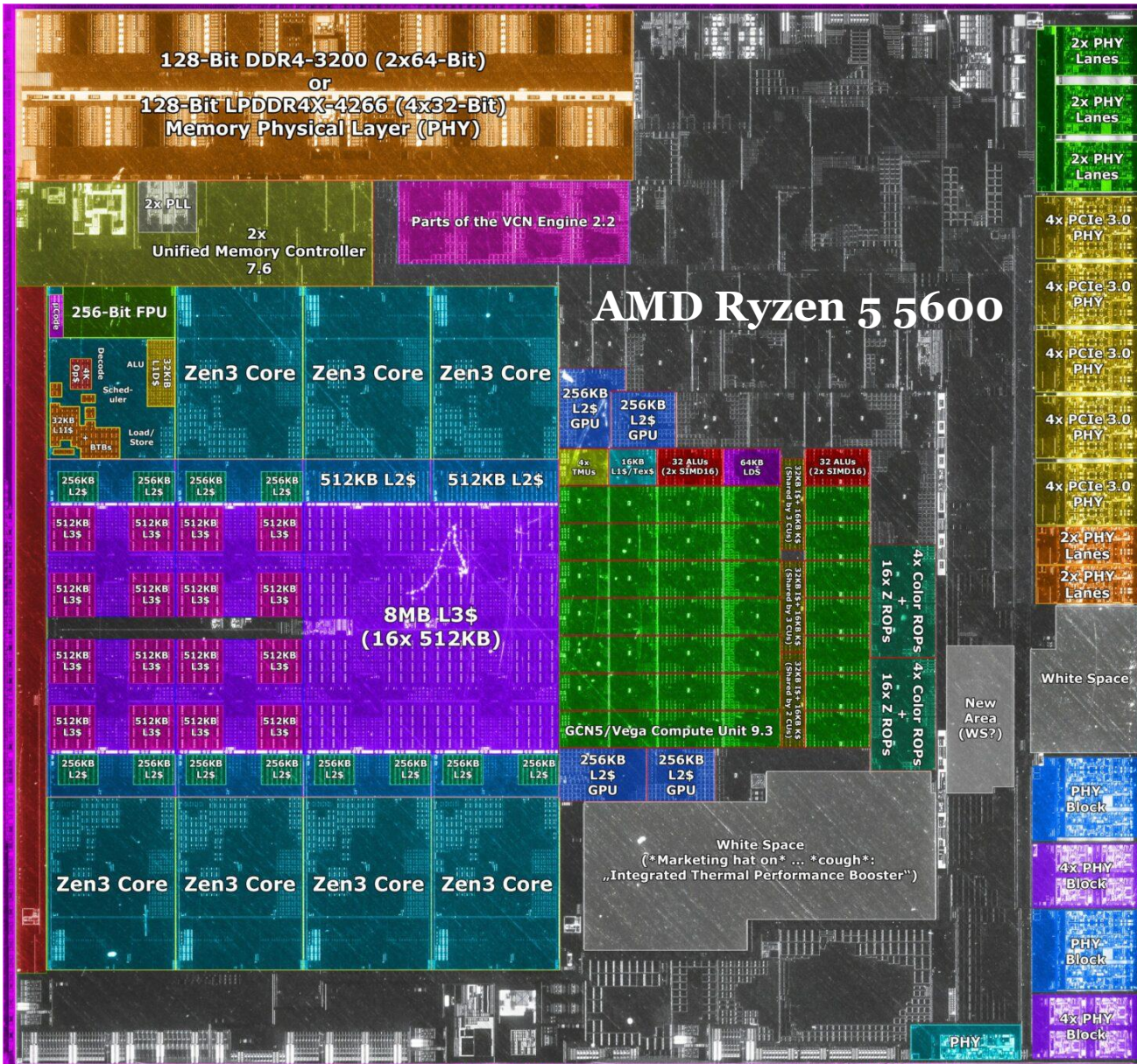
Intel Raptor Lake (i9)

PROCESSOR CACHE AREA



Intel Alder Lake

PROCESSOR CACHE AREA



Cezanne die shot from Fritzchens Fritz: <https://www.flickr.com/photos/130561288@N04/51375154375/>
Annotated by Locuza, August 2021

AMD 3D CHIPLET TECHNOLOGY

Structural silicon

64MB L3 cache die

Direct copper-to-copper bond

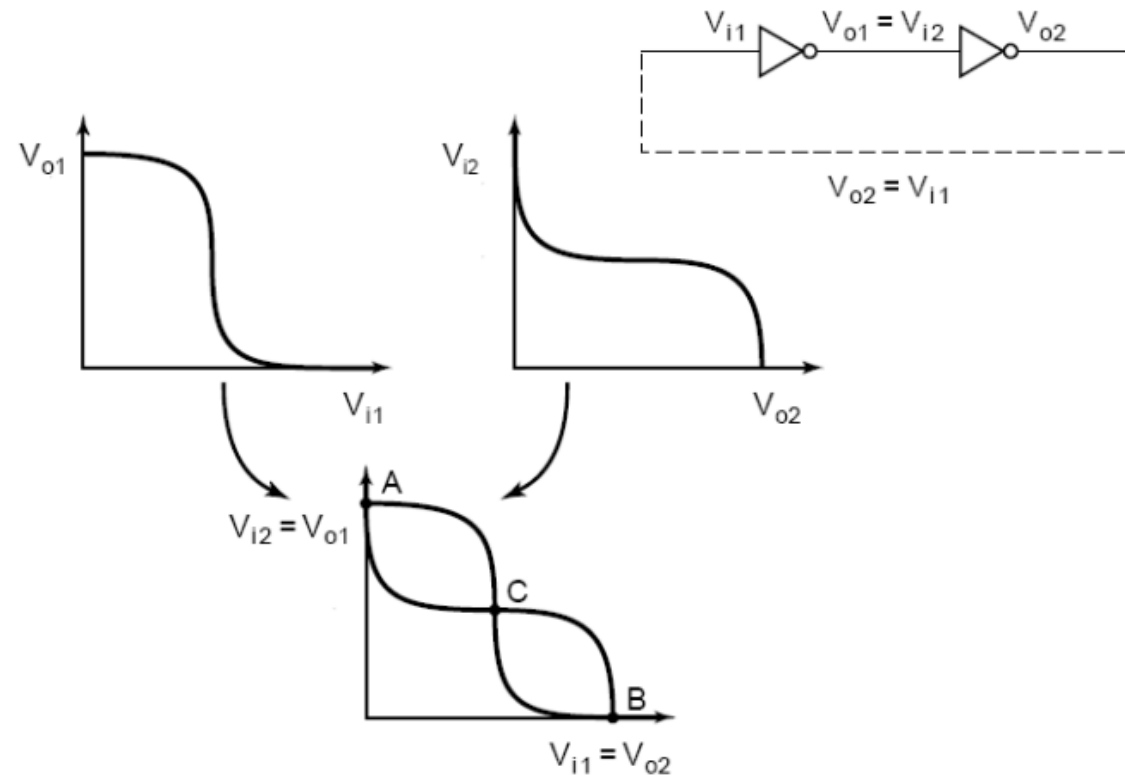
Through Silicon Vias (TSVs) for silicon-to-silicon communication

Up to 8-core "Zen 3" CCD

AMD Ryzen 7 Series 3D L3 Cache
– Releasing some time soon!

1- BIT MEMORY

- Recall our discussion on 1-bit memory to implement DFF (master and a slave)
- Back to back inverter – two states \rightarrow 1-bit memory

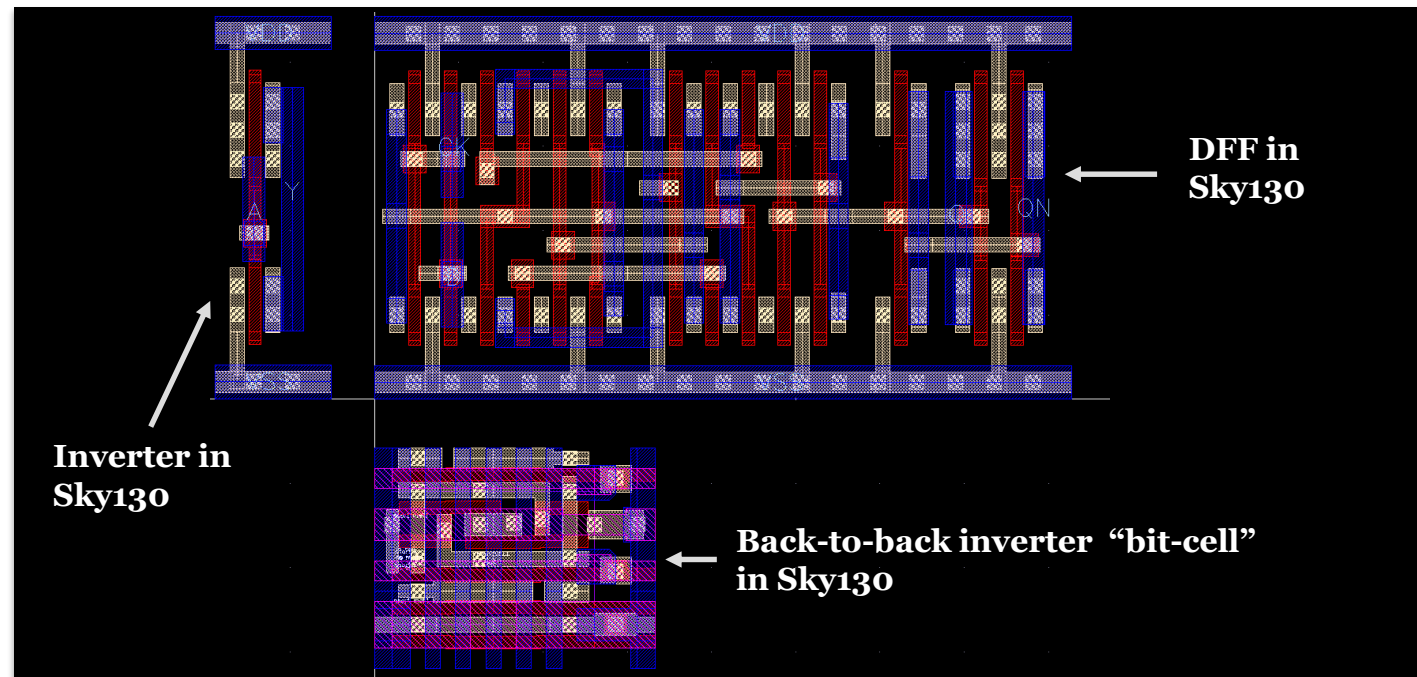


- Can we use a DFF to implement cache? – pros/cons?



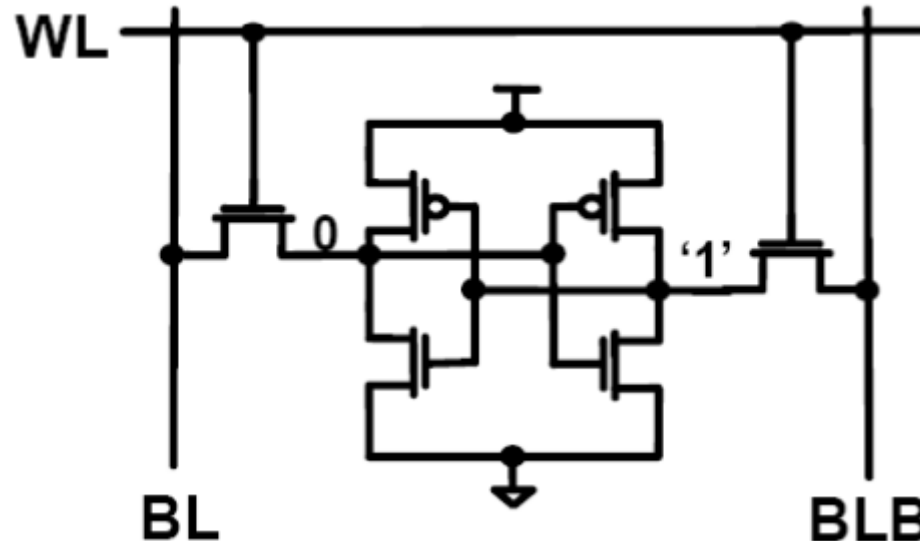
1- BIT MEMORY

- ❑ DFF: main usage was to implement 1-bit memory
- ❑ Cache: we will never access a single bit → access will be in bytes, half-word or word depending the micro-architecture of the processor
 - ❑ To implement a 1 byte memory → need 8 DFFs → is it worth it?
 - ❑ We have 2 back-to-back inverters per DFF and multiple inverters and TGs
 - ❑ DFF by default has a larger area to store a single bit memory !!!!!



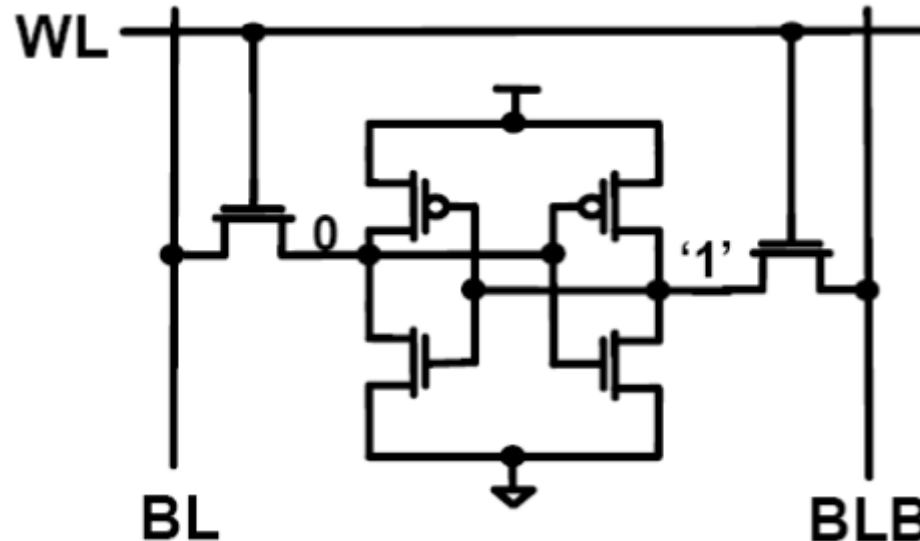
SRAM BIT-CELL

- ❑ Bit-cell: how to read/write/preserve a bit?
- ❑ Recall, in DFF we used TGs to perform read/write/preserve → can we reduce area further?
- ❑ Use NMOS only access transistors
- ❑ WL: word-line controls if read/write is to be performed. With $WL = 1$
 - ❑ If BL and BLB are driven to a value → we perform write
 - ❑ If BL and BLB are not driven → we perform read



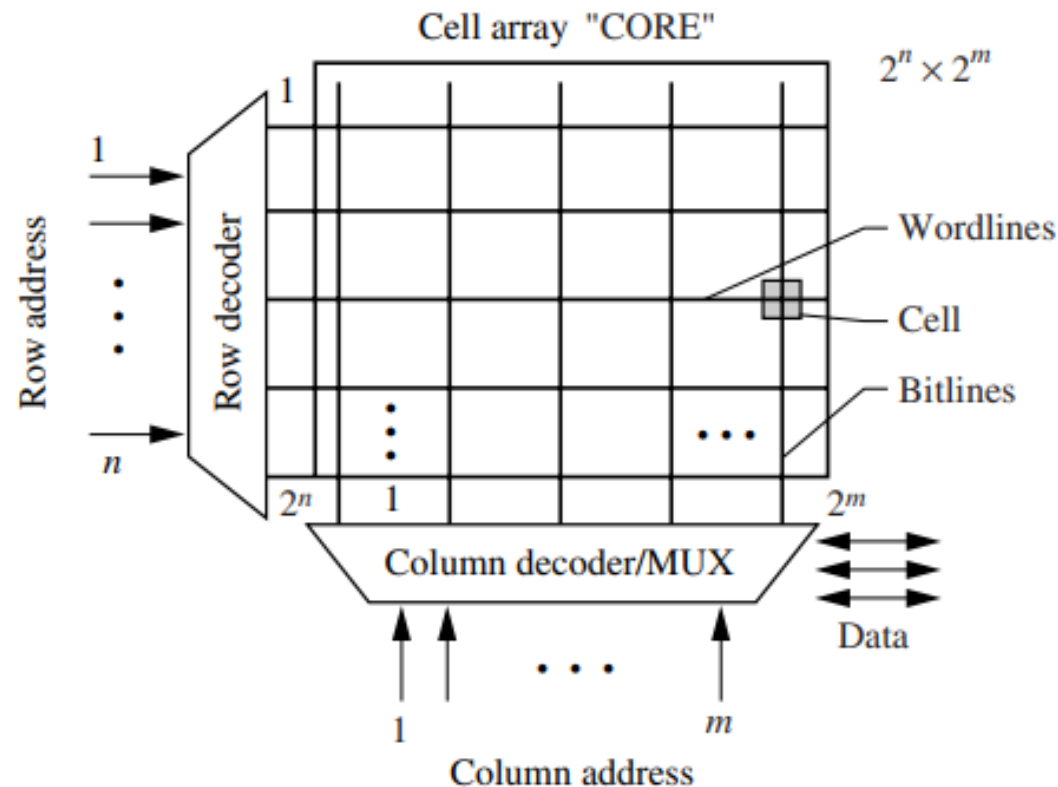
SRAM BIT-CELL

- ❑ SRAM bit-cell: Uses only 6 transistors to store 1-bit data → highly dense
- ❑ SRAM bit-cell: also called as 6T-cell
- ❑ Word-line controls if read/write is to performed
 - ❑ Word-line will have to be generated based on the memory address
- ❑ Bit-line: corresponds to data to be read or written (single-bit)



SRAM FULL PICTURE

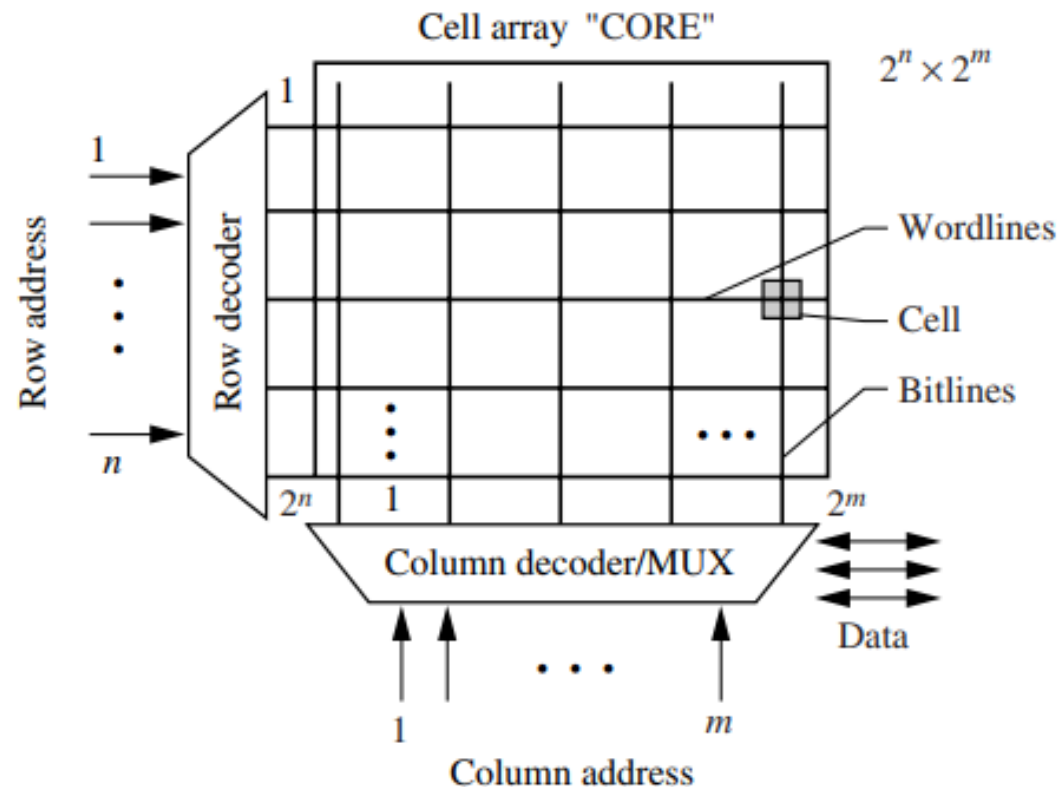
- ❑ Address decoding: split into row and column decoding
 - ❑ Split between row and column \rightarrow controls aspect ratio of memory
 - ❑ Say “n-bit” row address and “m-bit” column address
 - ❑ If $n > m \rightarrow$ more rows \rightarrow tall memory
 - ❑ If $m > n \rightarrow$ more columns \rightarrow wide memory



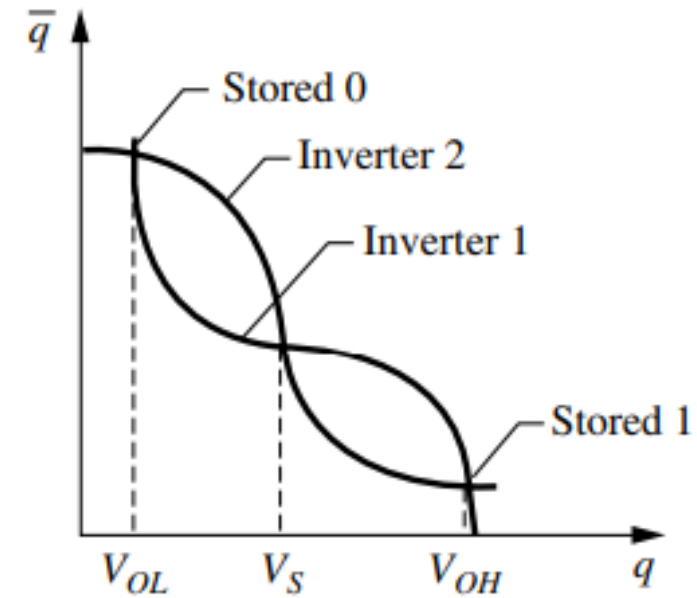
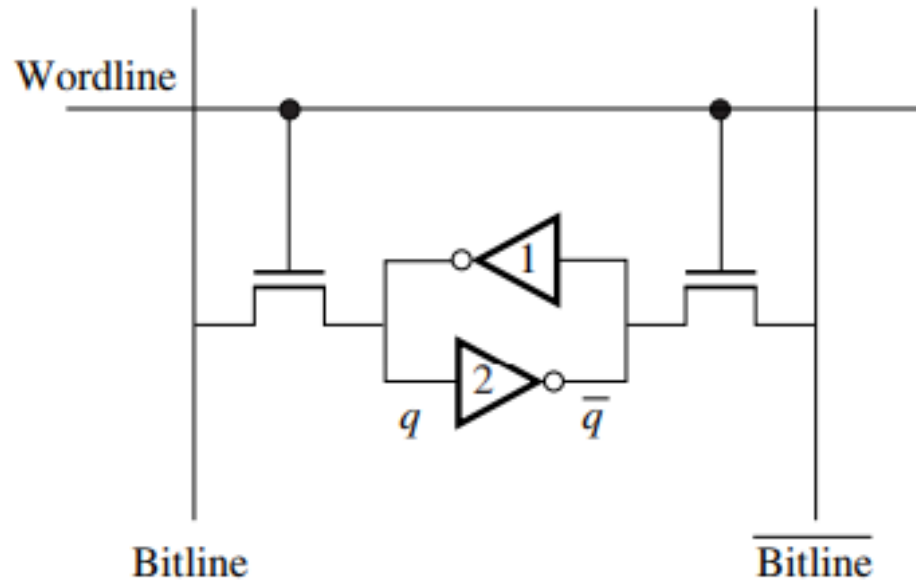
SRAM FULL PICTURE

Implications:

- If n is very large (tall memory) \rightarrow the shared bitline per column is longer \rightarrow this is a long metal in the layout \rightarrow large bit-line capacitance \rightarrow memory will be slower
- If m is large (wide memory) \rightarrow the wordline per row is longer \rightarrow large wordline capacitance (because of routing and also the 2 access transistors gate cap per column) \rightarrow Address decoder will drive larger load \rightarrow delay limited by logical effort !!!!



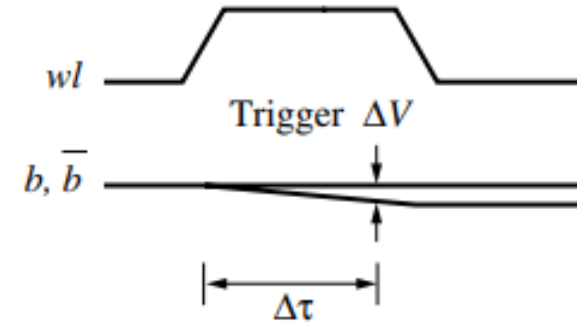
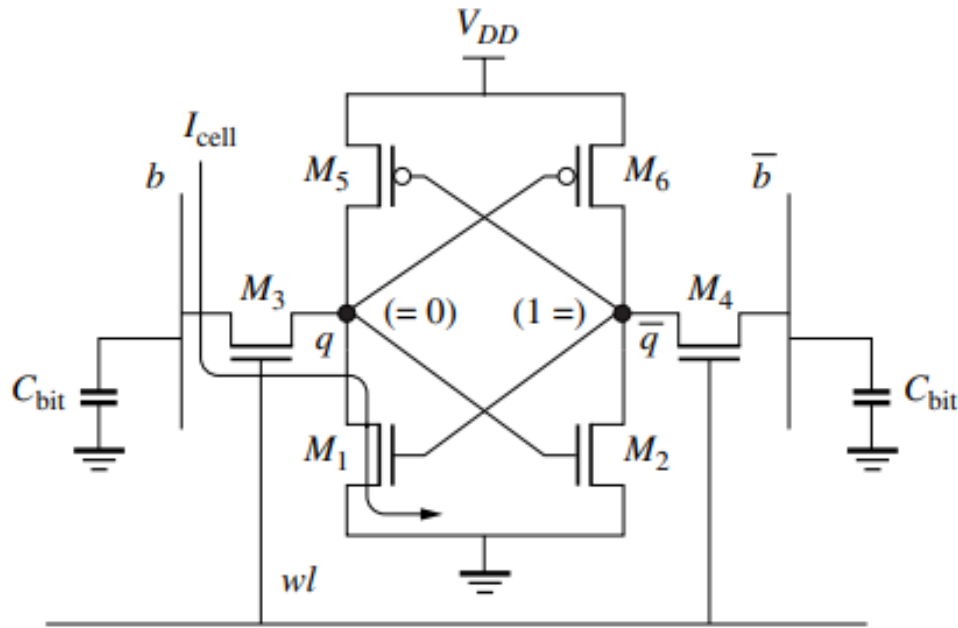
BACK TO THE BITCELL



□ Lets look at the read and write operations

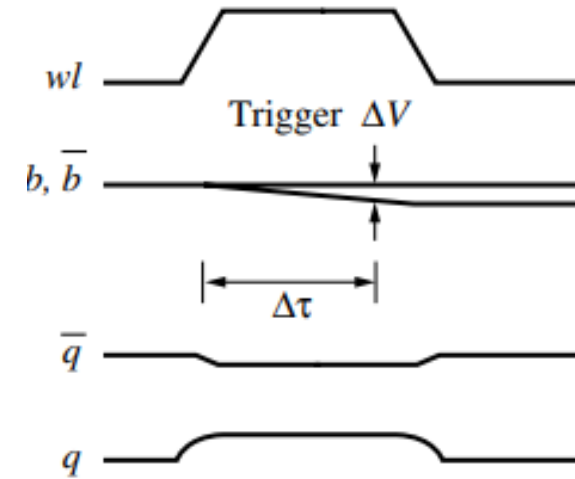
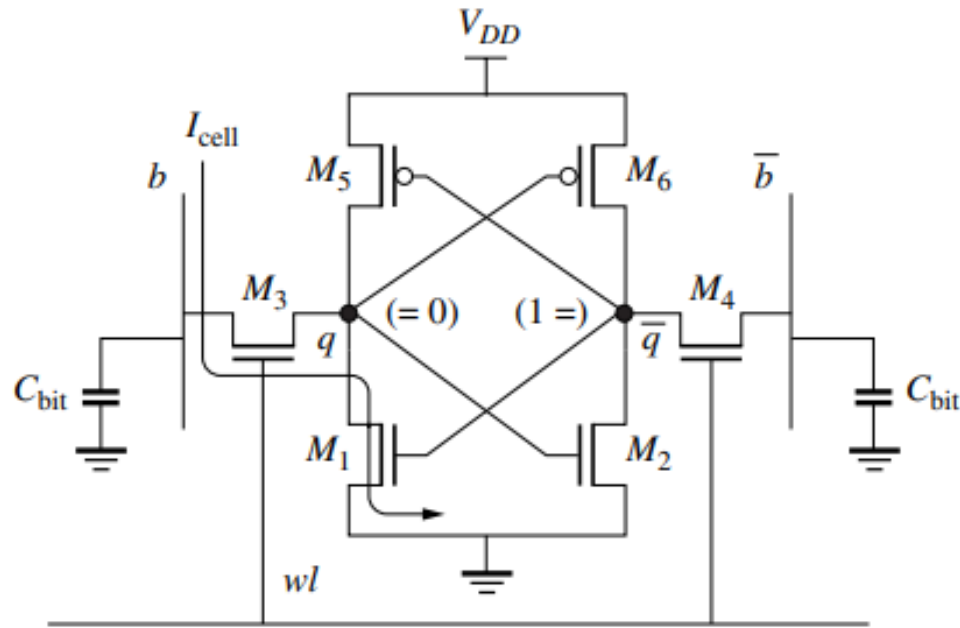


READ OPERATION



- ❑ Say bit “0” is already stored and we want to read
- ❑ Before performing any read → we pre-charge both BL and BLb to VDD → remove all prior read/write history on the lines → we cannot perform a read without pre-charging
- ❑ After pre-charge → WL is made high → capacitor discharge happens on one side!
- ❑ This small change in voltage will have to be amplified → will look at this later

LOOK AT READ OPERATION AGAIN



- ❑ M3 and M1 are actually fighting for node “q” → if M3 is strong it can turn on M2 !
- ❑ During read, we do not want to disturb the state of node “q”
 - ❑ i.e, change in node voltage “q” should not flip the previously stored data !
- ❑ To ensure the this, M3 and M1 will have to be sized accordingly
 - ❑ i.e., M1 should be stronger than M3 !! (M1 will have lower resistor than M3)

