



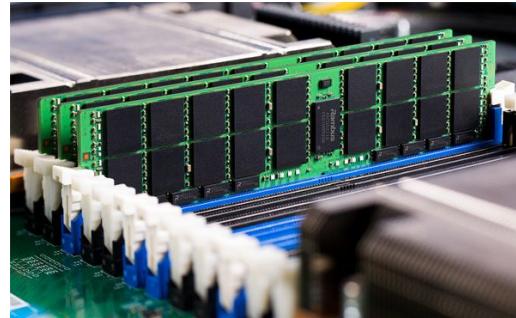
**EECS 151/251A**  
**Spring 2024**  
**Digital Design and Integrated Circuits**

Instructor:  
John Wawrzynek

Lecture 17:  
Memory Circuits  
and Blocks, Part 2

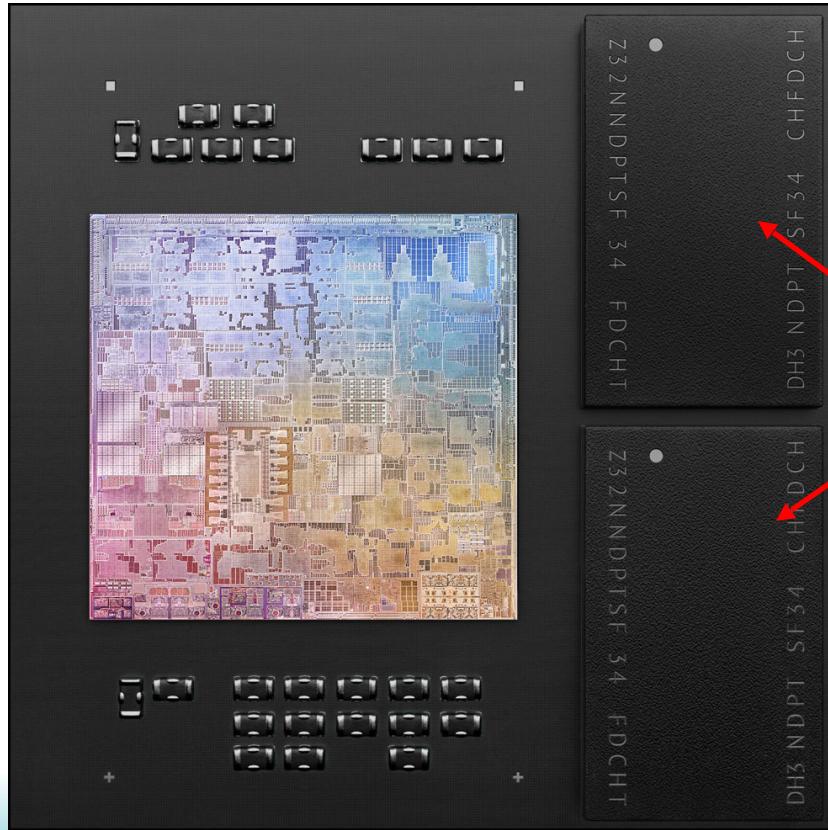


## DRAM



- ❑ Gets used for off-chip large inexpensive memories.
- ❑ Most commonly not compatible with logic processes. Requires special IC processing.

# DRAM Packaging, Apple M1

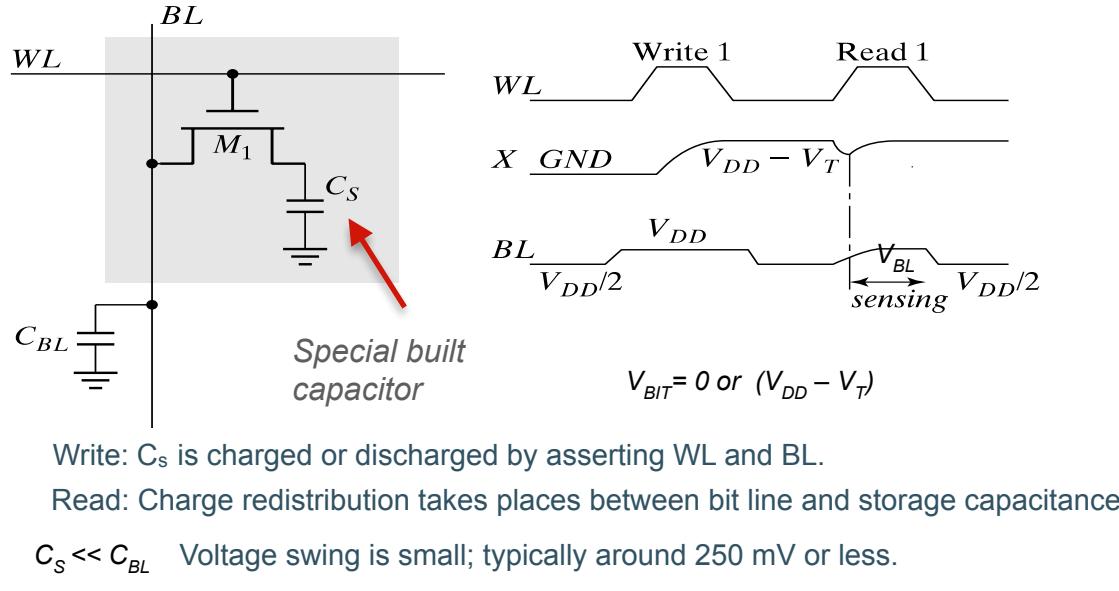


Two DRAM chips  
on same package  
as system SoC

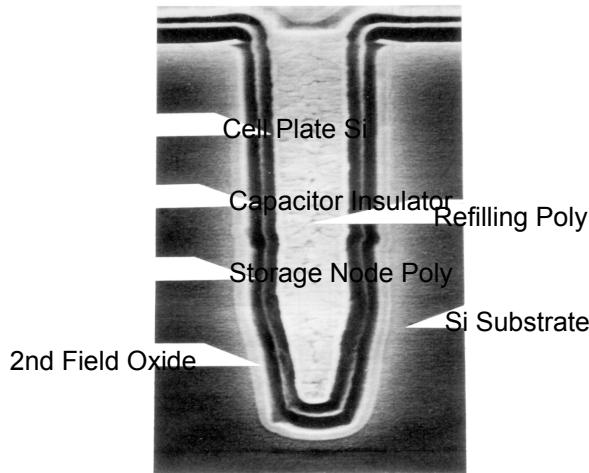
- 128b databus,  
running at 4.2Gb/s
- 68GB/s bandwidth

# 1-Transistor DRAM Cell

- For sufficient  $C_s$  in small area, **special IC process is used**
- Cell reading is destructive, therefore read operation always is followed by a write-back
- Cell loses charge (leaks away in ms - highly temperature dependent), therefore cells occasionally need to be “refreshed” using read/write cycle

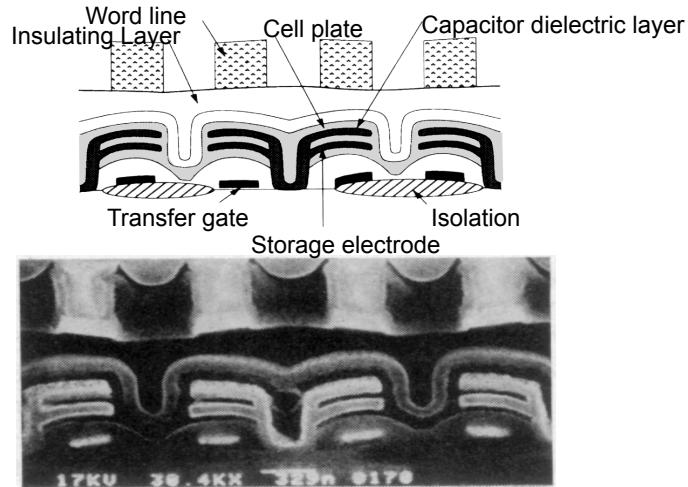


# Advanced 1T DRAM Cells



Trench Cell

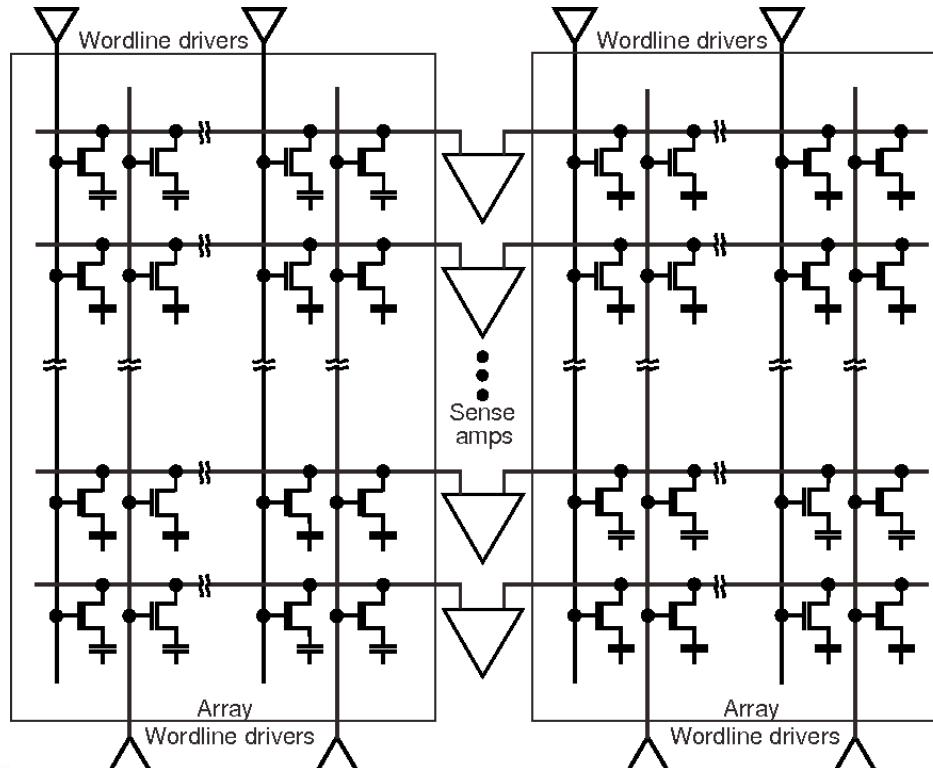
*No longer common*



Stacked-capacitor Cell

*Common*

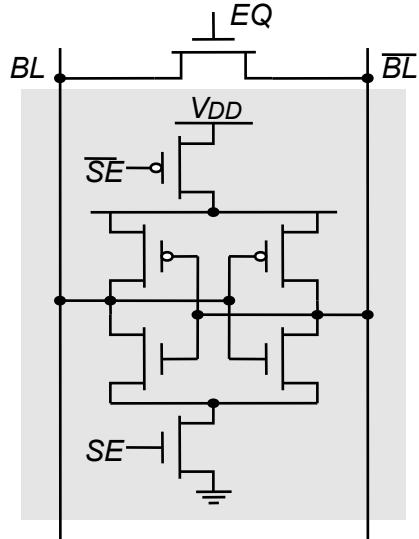
# DRAM Sub-array Organization



- ❑ Unlike SRAM, DRAM is “single-sided” read/write
- ❑ But sense-amps like dual-rail
- ❑ “split bit lines” is standard practice

Figure 2.6 from The DRAM Array | Semantic Scholar

# Latch-Based Sense Amplifier (DRAM)



- Bit lines equalized, with EQ, and precharged to  $V_{dd}/2$
- Sense amp initialized to its meta-stable point with EQ
- Once adequate voltage gap created, sense amp enabled with SE
- Positive feedback quickly forces output to a stable operating point.
- With row select kept on, cell gets “refreshed”
- Bit line from inactive array below/above used as reference for differential sensing.



# Memory Blocks

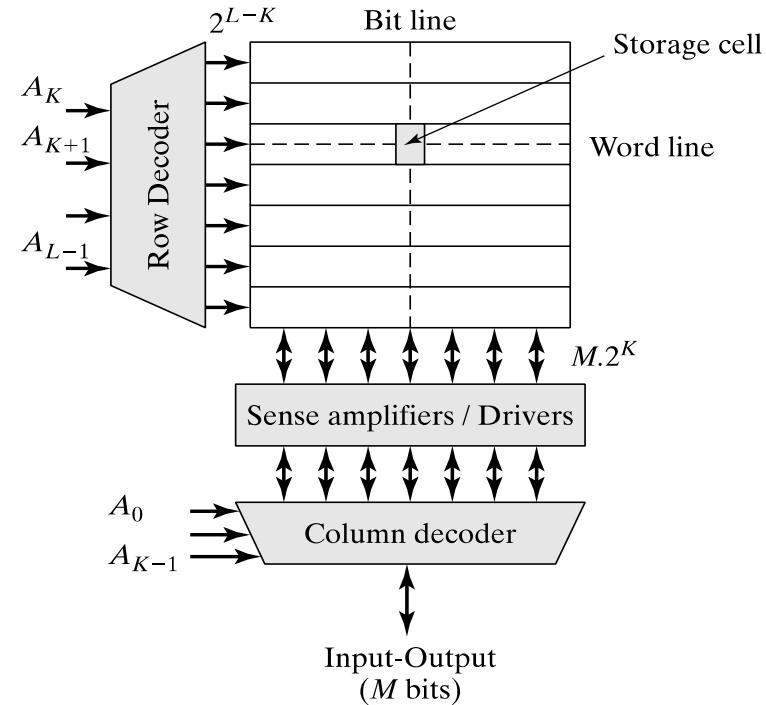
- Multi-ported RAM
- Combining Memory blocks
- FIFOs
- FPGA memory blocks
- Caches
- Memory Blocks in the Project



## *Multi-ported memory*

# Memory Architecture Review

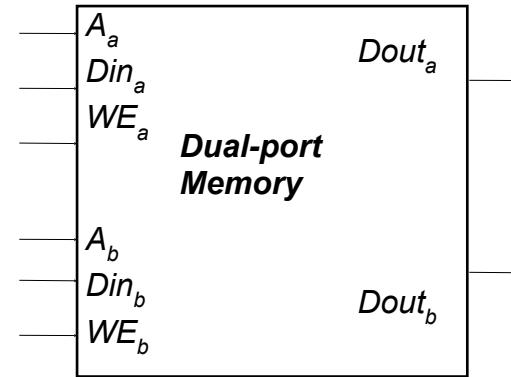
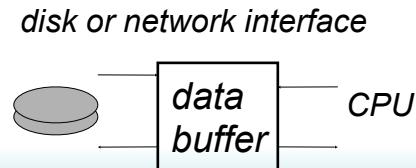
- ❑ **Word lines** used to select a row for reading or writing
- ❑ **Bit lines** carry data to/from periphery
- ❑ **Core aspect ratio** keep close to 1 to help balance delay on word line versus bit line
- ❑ **Address bits** are divided between the two decoders
- ❑ **Row decoder** used to select word line
- ❑ **Column decoder** used to select one or more columns for input/output of data



# Multi-ported Memory

## □ Motivation:

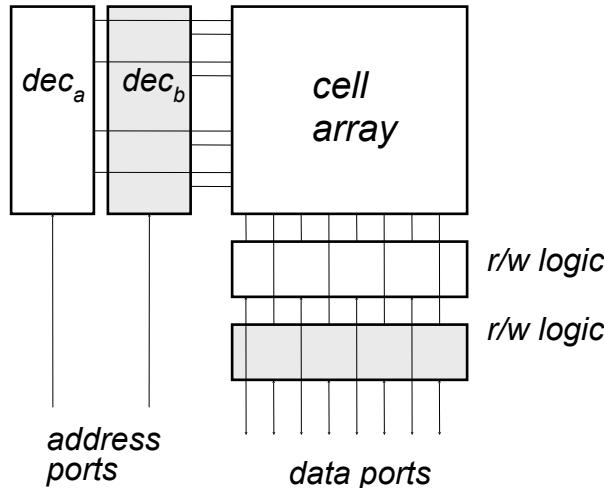
- Consider CPU core register file:
  - 1 read or write per cycle limits processor performance.
  - Complicates pipelining. Difficult for different instructions to simultaneously read or write regfile.
  - Common arrangement in pipelined CPUs is 2 read ports and 1 write port.
  - Another example: I/O data buffering:



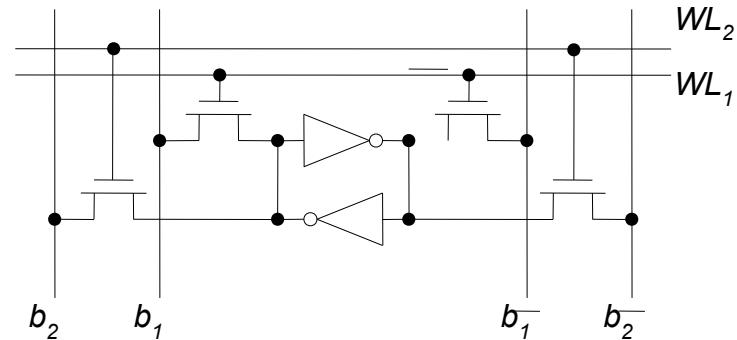
- dual-porting allows both sides to simultaneously access memory at full bandwidth.

# Dual-ported Memory Internals

- Add decoder, another set of read/write logic, bits lines, word lines:



- *Example cell: SRAM*

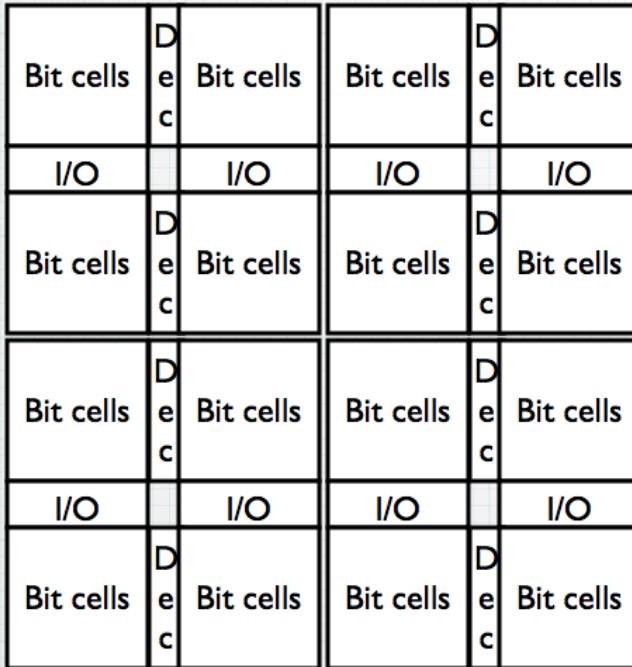


- Repeat everything but cross-coupled inverters.
- This scheme extends up to a couple more ports, then need to add additional transistors.

## *Combining Memory Blocks*



# Building Larger Memories

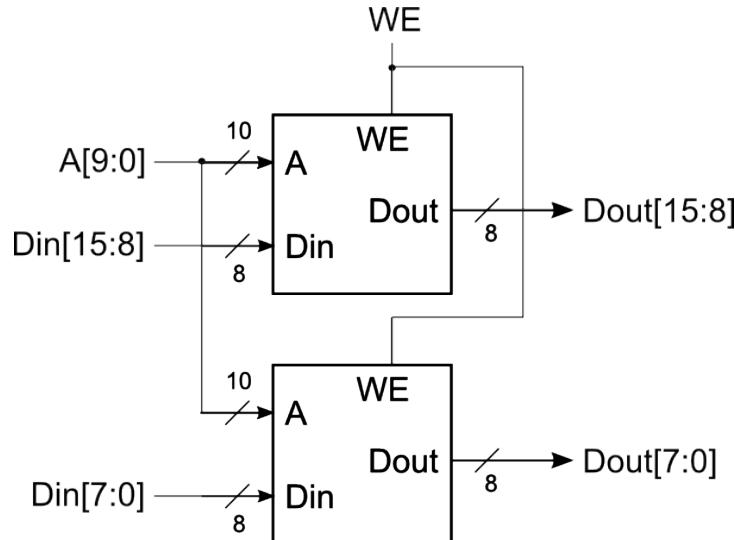


- Large arrays constructed by tiling multiple leaf arrays, sharing decoders and I/O circuitry
  - e.g., sense amp attached to arrays above and below
- Leaf array limited in size to 128-256 bits in row/column due to RC delay of wordlines and bitlines
- Also to reduce power by only activating selected sub-bank
- In larger memories, delay and energy dominated by I/O wiring

# Cascading Memory-Blocks

How to make larger memory blocks out of smaller ones.

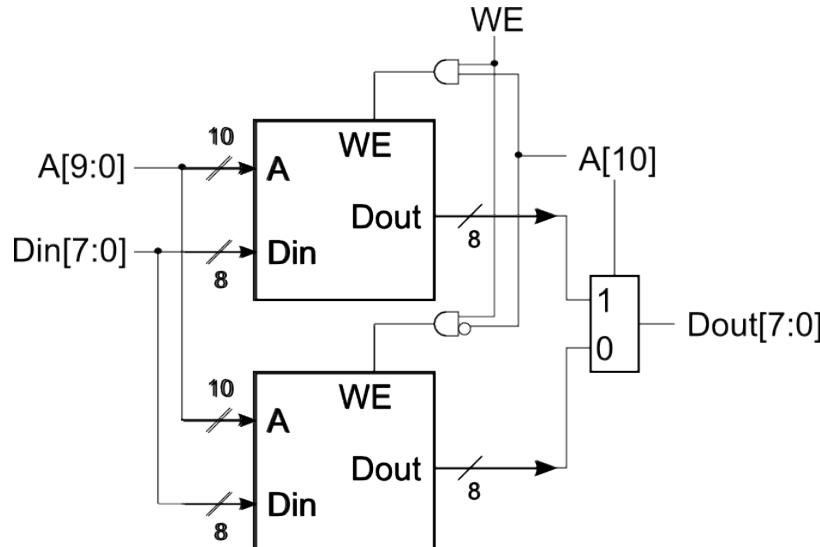
Increasing the width. Example: given 1Kx8, want 1Kx16



# Cascading Memory-Blocks

How to make larger memory blocks out of smaller ones.

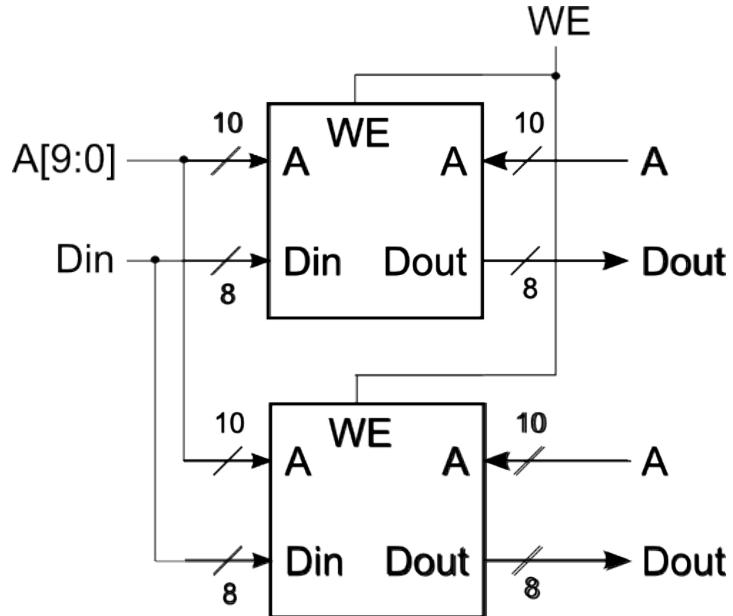
Increasing the depth. Example: given  $1K \times 8$ , want  $2K \times 8$



# Adding Ports to Primitive Memory Blocks

Adding a read port to a simple dual port (SDP) memory.

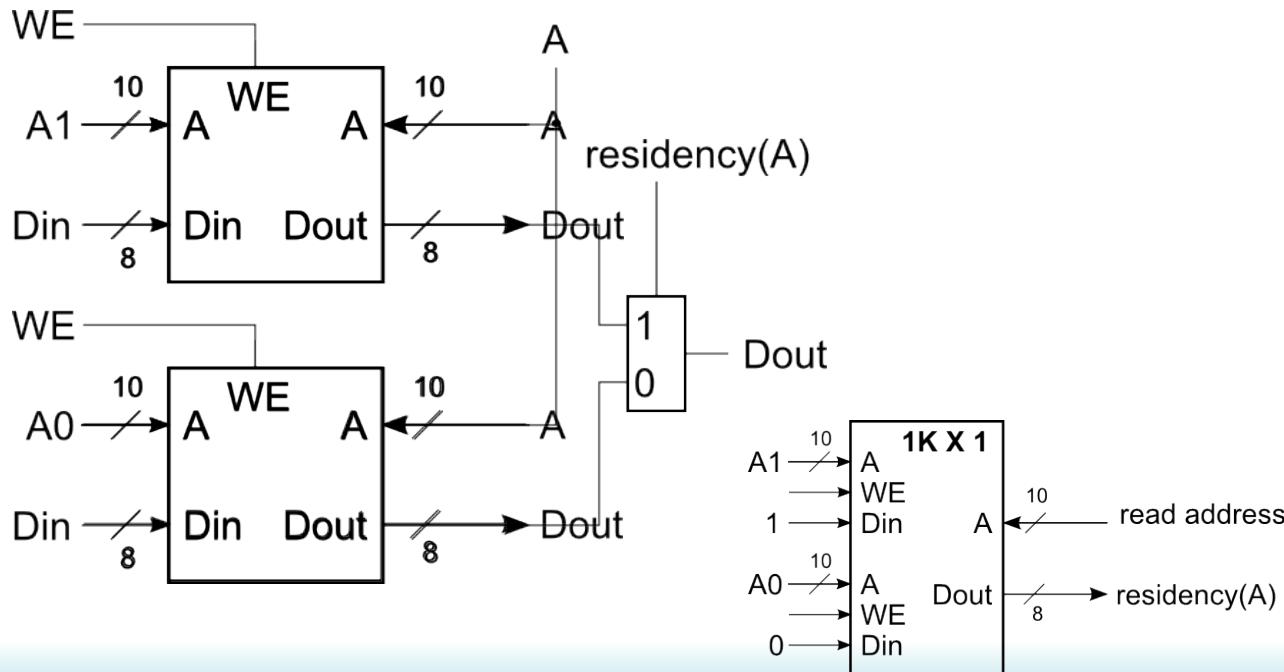
Example: given 1Kx8 SDP, want 1 write & 2 read ports.



# Adding Ports to Primitive Memory Blocks

How to add a write port to a simple dual port memory.

Example: given 1Kx8 SDP, want 1 read & 2 write ports.

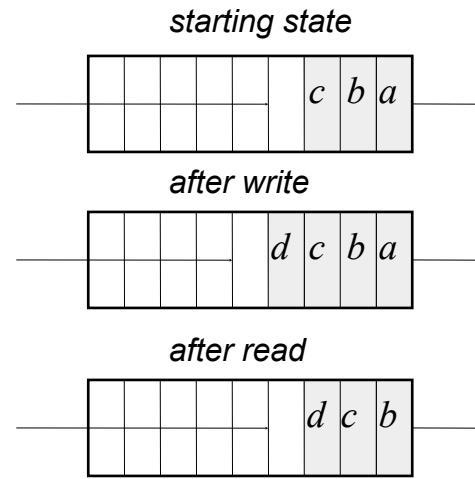




*FIFOs*

# First-in-first-out (FIFO) Memory

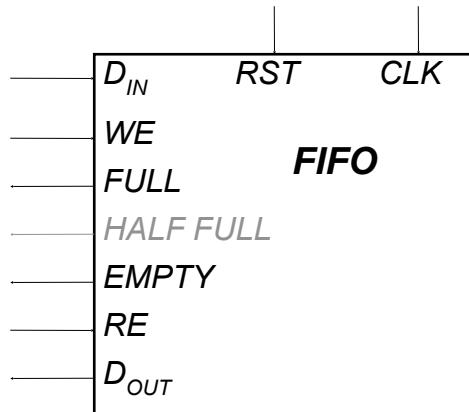
- ❑ Used to implement queues.
- ❑ These find common use in computers and communication circuits.
- ❑ Generally, used to “decouple” actions of producer and consumer:



- Producer can perform many writes without consumer performing any reads (or vis versa). However, because of finite buffer size, on average, need equal number of reads and writes.
- Typical uses:
  - interfacing I/O devices.  
Example network interface.  
Data bursts from network, then processor bursts to memory buffer (or reads one word at a time from interface). Operations not synchronized.
  - Example: Audio output.  
Processor produces output samples in bursts (during process swap-in time). Audio DAC clocks it out at constant sample rate.

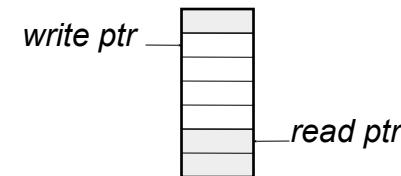
# FIFO Interfaces

- ❑ After write or read operation, FULL and EMPTY indicate status of buffer.
- ❑ Used by external logic to control its own reading from or writing to the buffer.
- ❑ FIFO resets to EMPTY state.
- ❑ HALF FULL (or other indicator of partial fullness) is optional.

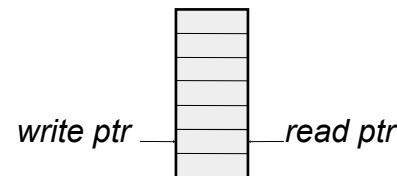


## "Circular buffer" implementation:

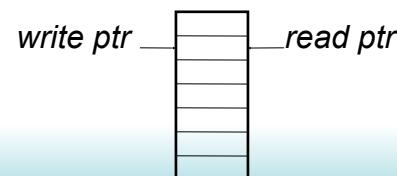
- Address pointers are used internally to keep next write position and next read position into a dual-port memory.



- If pointers equal after write  $\Rightarrow$  FULL:



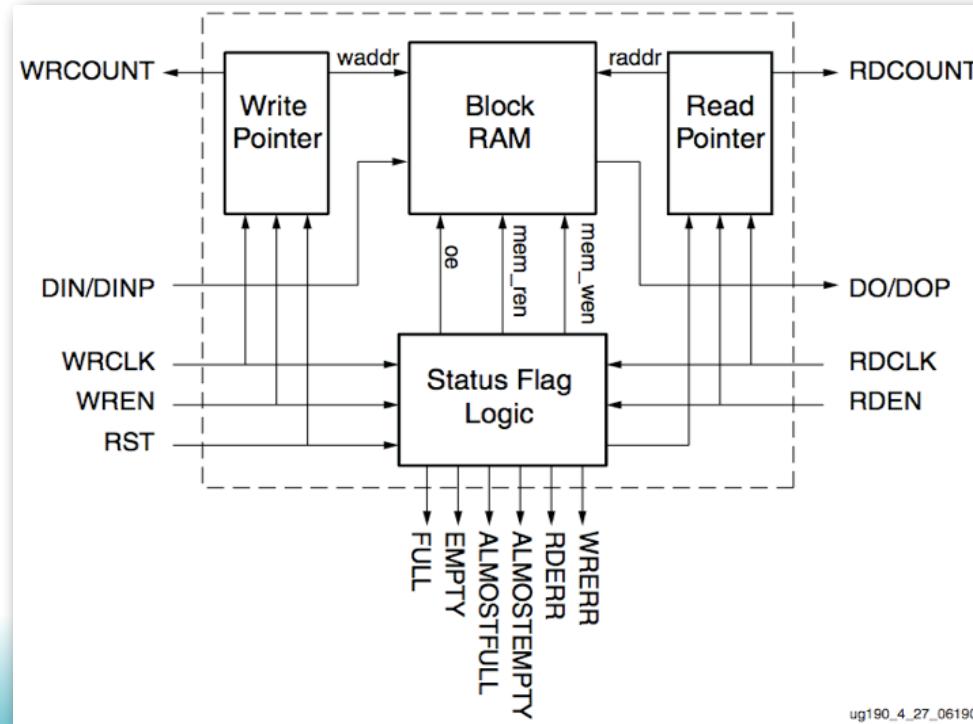
- If pointers equal after read  $\Rightarrow$  EMPTY:



Note: pointer incrementing is done "mod size-of-buffer"

# Xilinx Virtex5 FIFOs

- ❑ Virtex5 BlockRAMS include dedicated circuits for FIFOs.
- ❑ Details in User Guide (ug190).
- ❑ Takes advantage of separate dual ports and independent ports clocks.



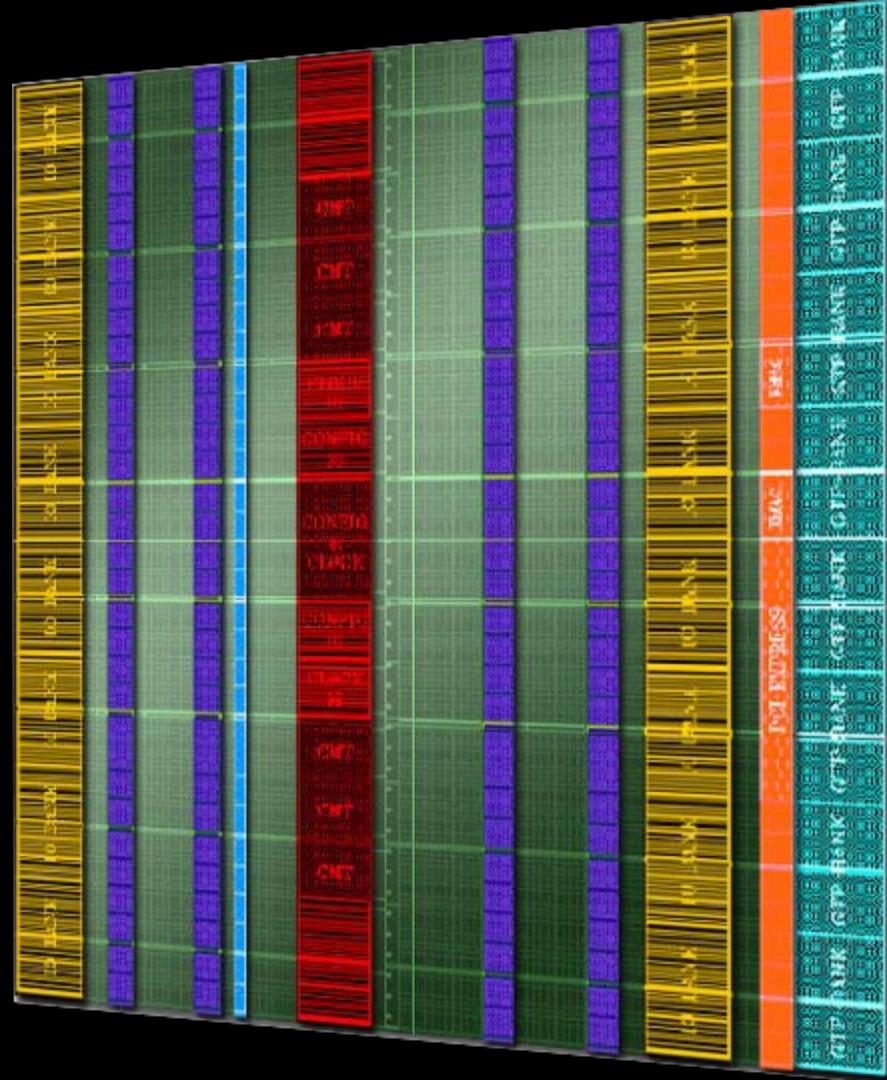


# *Memory on FPGAs*

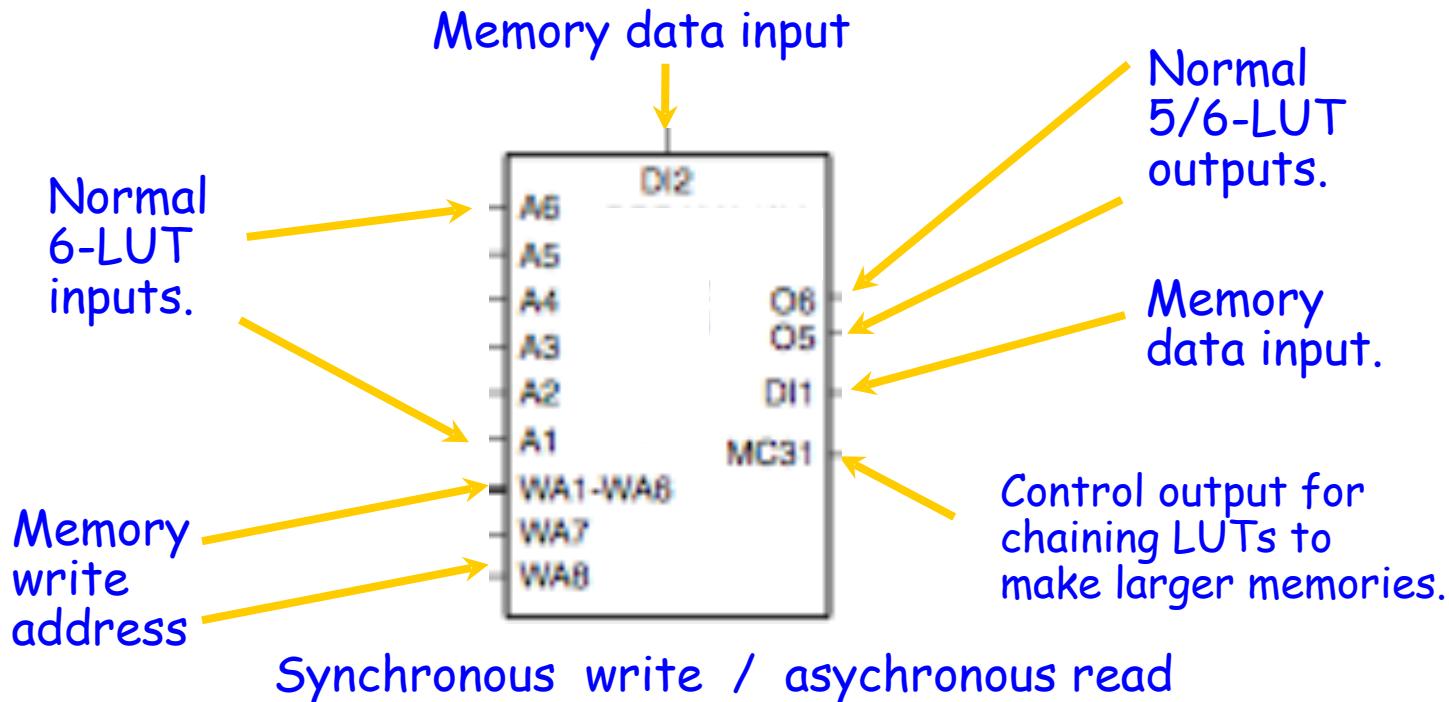
Virtex-5 LX110T  
memory blocks.

Distributed RAM  
using LUTs among  
the CLBs.

Block RAMs in  
four columns.

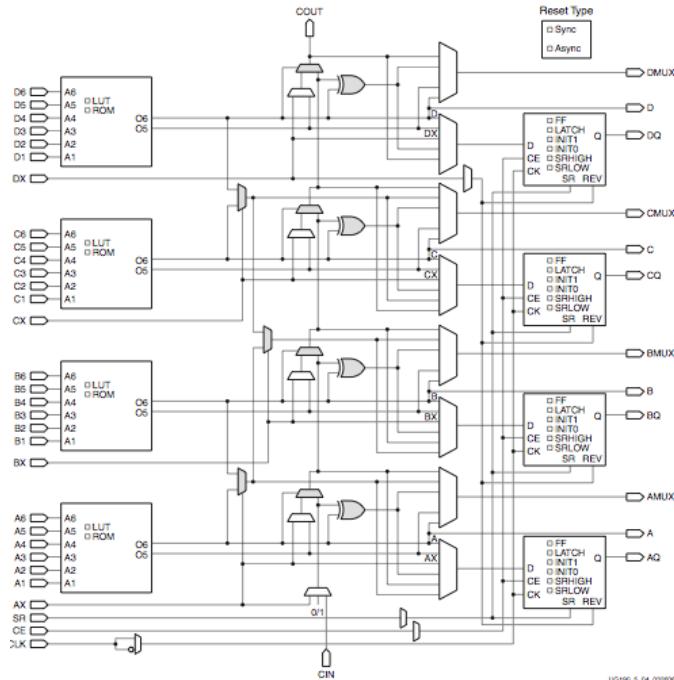


# A SLICEM 6-LUT ... (“distributed RAM” aka “LUT RAM”)



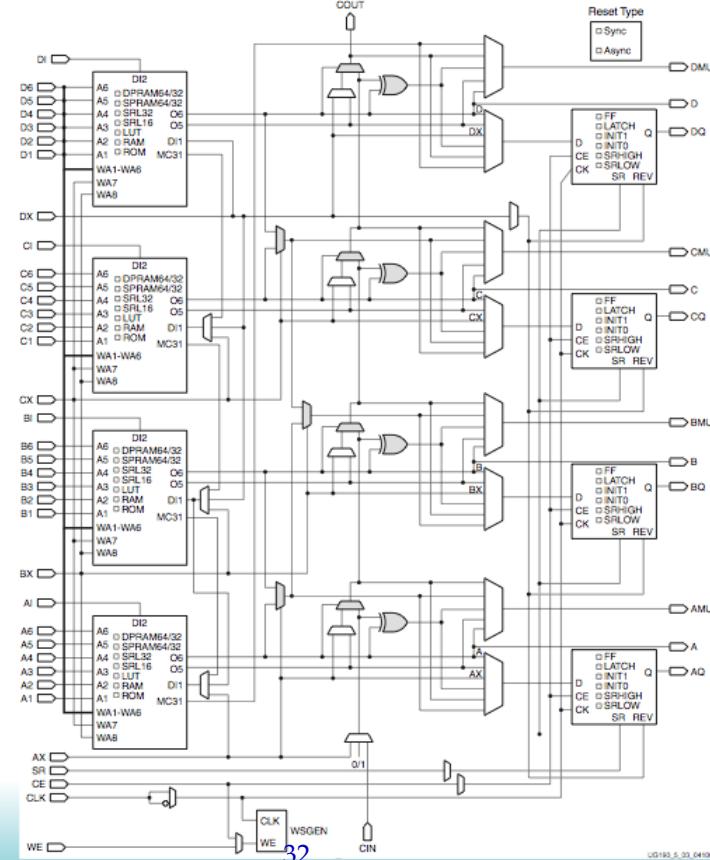
# SLICEL vs SLICEM ...

SLICEL

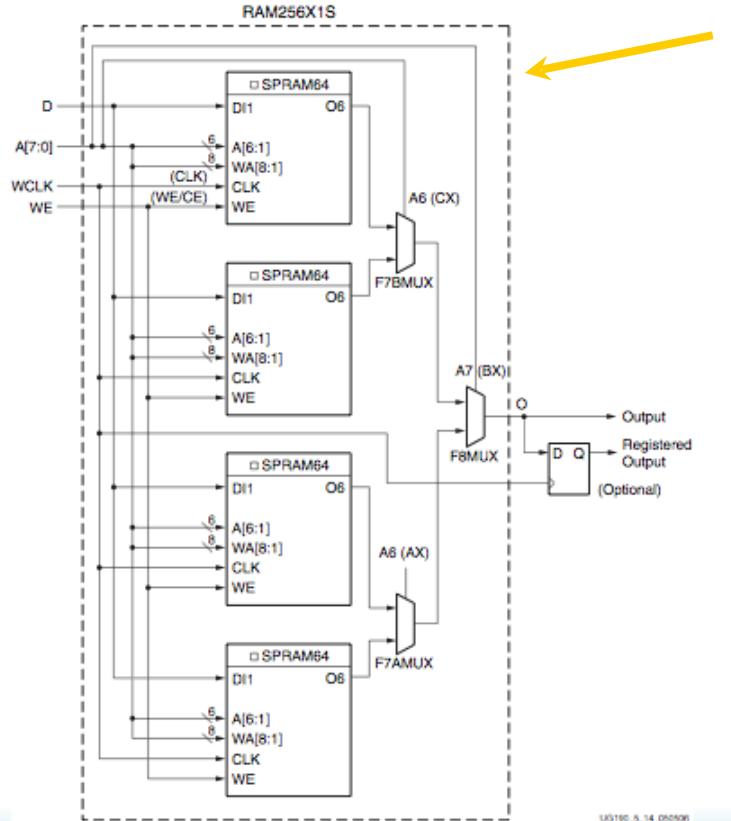


SLICEM adds memory features to LUTs, + muxes.

SLICEM



# Example Distributed RAM (LUT RAM)



Example configuration:  
Single-port 256b x 1,  
registered output.

Figure 5-14: Distributed RAM (RAM256X1S)

# Distributed RAM Timing

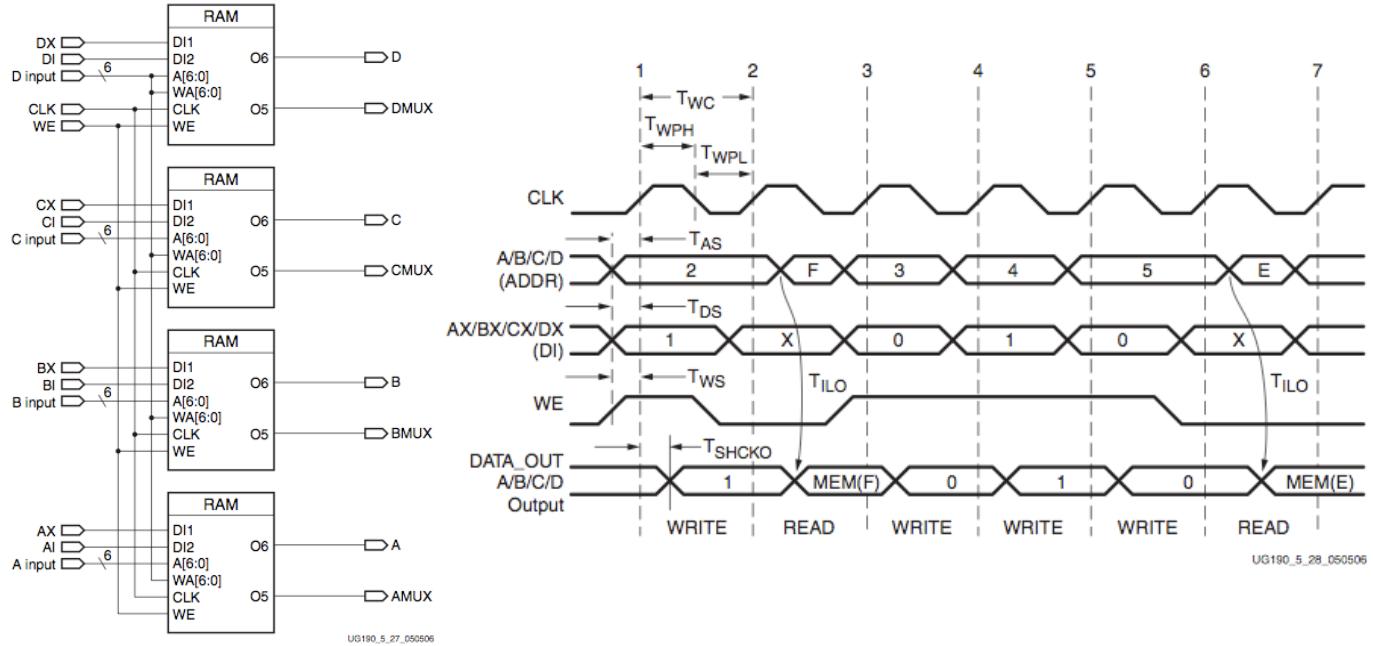
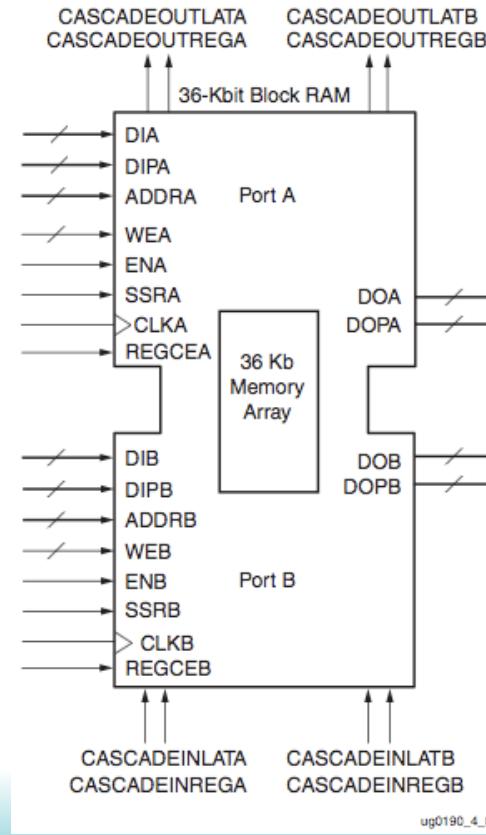


Figure 5-27: Simplified Virtex-5 FPGA SLICEM Distributed RAM

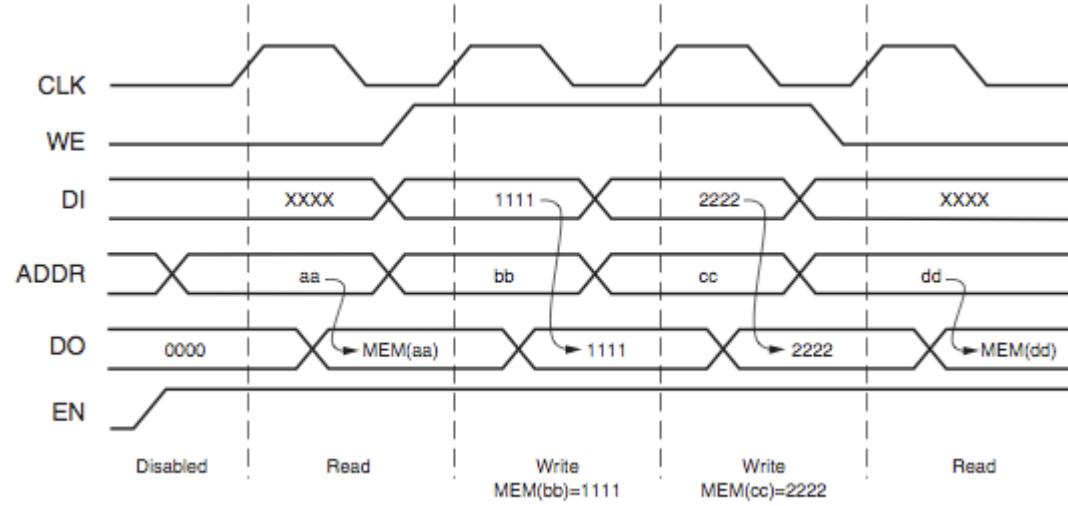
# Block RAM Overview



- ❑ 36K bits of data total, can be configured as:
  - 2 independent 18Kb RAMs, or one 36Kb RAM.
- ❑ Each 36Kb block RAM can be configured as:
  - 64Kx1 (when cascaded with an adjacent 36Kb block RAM), 32Kx1, 16Kx2, 8Kx4, 4Kx9, 2Kx18, or 1Kx36 memory.
- ❑ Each 18Kb block RAM can be configured as:
  - 16Kx1, 8Kx2, 4Kx4, 2Kx9, or 1Kx18 memory.
- ❑ Write and Read are synchronous operations.
- ❑ The two ports are symmetrical and totally independent (can have different clocks), sharing only the stored data.
- ❑ Each port can be configured in one of the available widths, independent of the other port. The read port width can be different from the write port width for each port.
- ❑ The memory content can be initialized or cleared by the configuration bitstream.

# Block RAM Timing

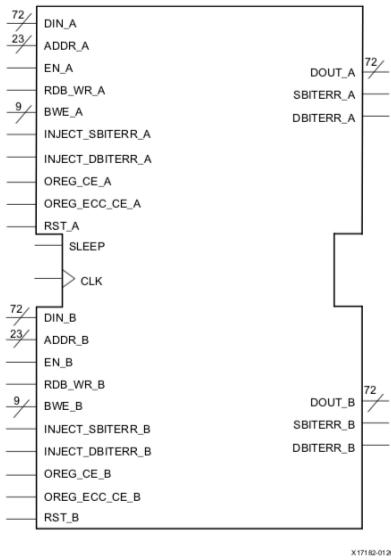
- ❑ Optional output register, would delay appearance of output data by one cycle.
- ❑ Maximum clock rate, roughly 400MHz.



ug190\_4\_03\_032206

# Ultra-RAM Blocks

Table 2-1: Block RAM and UltraRAM Comparison



Feature	Block RAM	UltraRAM
Clocking	Two clocks	Single clock
Built-in FIFO	Yes	No
Data width	Configurable (1, 2, 4, 9, 18, 36, 72)	Fixed (72-bits)
Modes	SDP and TDP	Two ports, each can independently read or write (a superset of SDP)
ECC	64-bit SECDED Supported in 64-bit SDP only (one ECC decoder for port A and one ECC encoder for port B)	64-bit SECDED One set of complete ECC logic for each port to enable independent ECC operations (ECC encoder and decoder for both ports)
Cascade	<ul style="list-style-type: none"> <li>Cascade output only (input cascade implemented via logic resources)</li> <li>Cascade within a single clock region</li> </ul>	<ul style="list-style-type: none"> <li>Cascade both input and output (with global address decoding)</li> <li>Cascade across clock regions in a column</li> <li>Cascade across several columns with minimal logic resources</li> </ul>
Power savings	One mode via manual signal assertion	One mode via manual signal assertion

Figure 2-1: UltraRAM URAM288\_BASE Primitive

# State-of-the-Art - Xilinx FPGAs



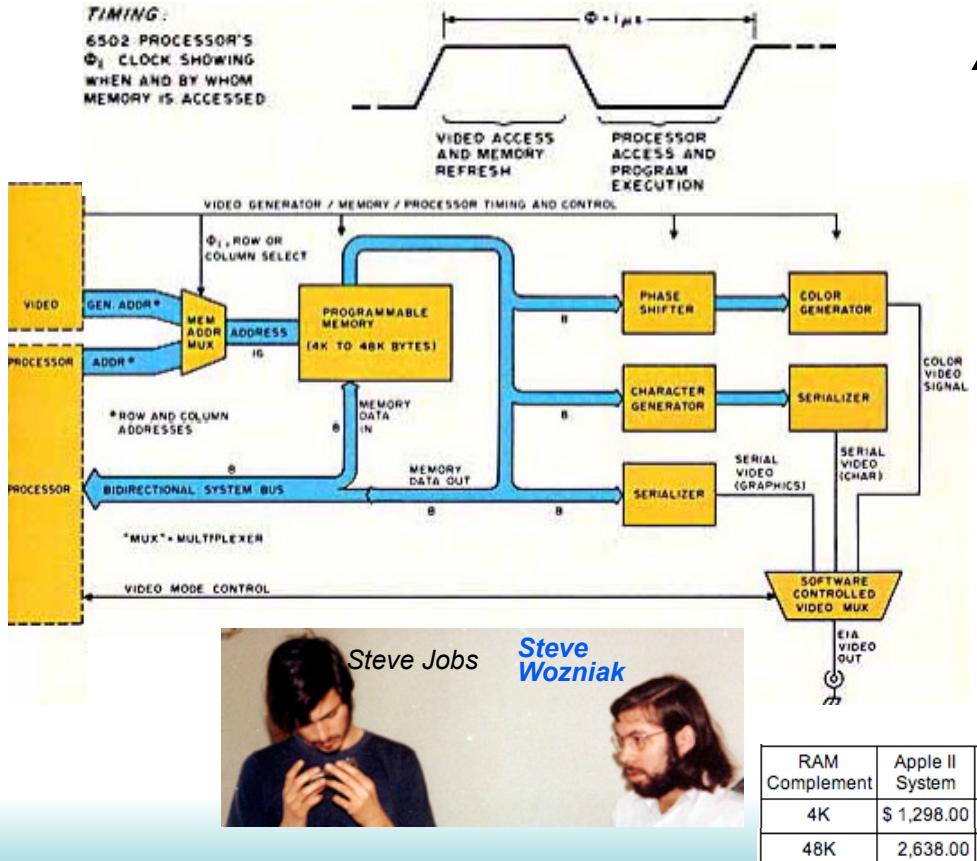
## Virtex Ultra-scale

Device Name	VU3P	VU5P	VU7P	VU9P	VU11P	VU13P	VU27P	VU29P	VU31P	VU33P	VU35P	VU37P
System Logic Cells (K)	862	1,314	1,724	2,586	2,835	3,780	2,835	3,780	962	962	1,907	2,852
CLB Flip-Flops (K)	788	1,201	1,576	2,364	2,592	3,456	2,592	3,456	879	879	1,743	2,607
CLB LUTs (K)	394	601	788	1,182	1,296	1,728	1,296	1,728	440	440	872	1,304
Max. Dist. RAM (Mb)	12.0	18.3	24.1	36.1	36.2	48.3	36.2	48.3	12.5	12.5	24.6	36.7
Total Block RAM (Mb)	25.3	36.0	50.6	75.9	70.9	94.5	70.9	94.5	23.6	23.6	47.3	70.9
UltraRAM (Mb)	90.0	132.2	180.0	270.0	270.0	360.0	270.0	360.0	90.0	90.0	180.0	270.0
HBM DRAM (GB)	—	—	—	—	—	—	—	—	4	8	8	8
HBM AXI Interfaces	—	—	—	—	—	—	—	—	32	32	32	32
Clock Mgmt Tiles (CMTs)	10	20	20	30	12	16	16	16	4	4	8	12
DSP Slices	2,280	3,474	4,560	6,840	9,216	12,288	9,216	12,288	2,880	2,880	5,952	9,024
Peak INT8 DSP (TOP/s)	7.1	10.8	14.2	21.3	28.7	38.3	28.7	38.3	8.9	8.9	18.6	28.1
PCIe® Gen3 x16	2	4	4	6	3	4	1	1	0	0	1	2
PCIe Gen3 x16/Gen4 x8 / CCIX <sup>[1]</sup>	—	—	—	—	—	—	—	—	4	4	4	4
150G Interlaken	3	4	6	9	6	8	6	8	0	0	2	4
100G Ethernet w/ KR4 RS-FEC	3	4	6	9	9	12	11	15	2	2	5	8
Max. Single-Ended HP I/Os	520	832	832	832	624	832	520	676	208	208	416	624
GTY 32.75Gb/s Transceivers	40	80	80	120	96	128	32	32	32	32	64	96
GTM 58Gb/s PAM4 Transceivers							32	48				
100G / 50G KP4 FEC							16 / 32	24 / 48				
Extended <sup>[2]</sup>	-1	-2	-2L	-3	-1	-2	-2L	-3	-1	-2	-2L	-3
Industrial	-1	-2	-1	-2	-1	-2	-1	-2	-	-	-	-



# Caches

# 1977: DRAM faster than microprocessors



Apple II (1977)  
CPU: 1000 ns  
DRAM: 400 ns

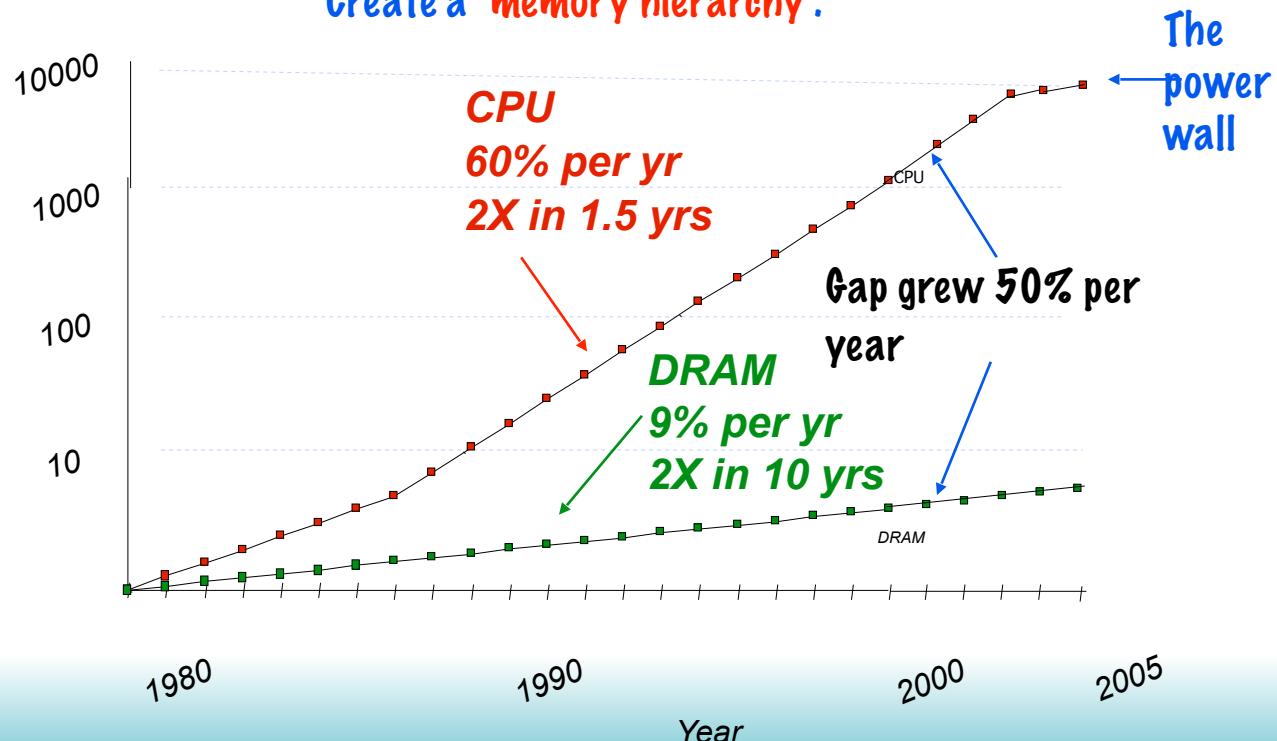


# 1980-2003, CPU speed outpaced DRAM ...

Q. How did architects address this gap?

Performance  
(1/latency)

A. Put smaller, faster "cache" memories between CPU and DRAM.  
Create a "memory hierarchy".

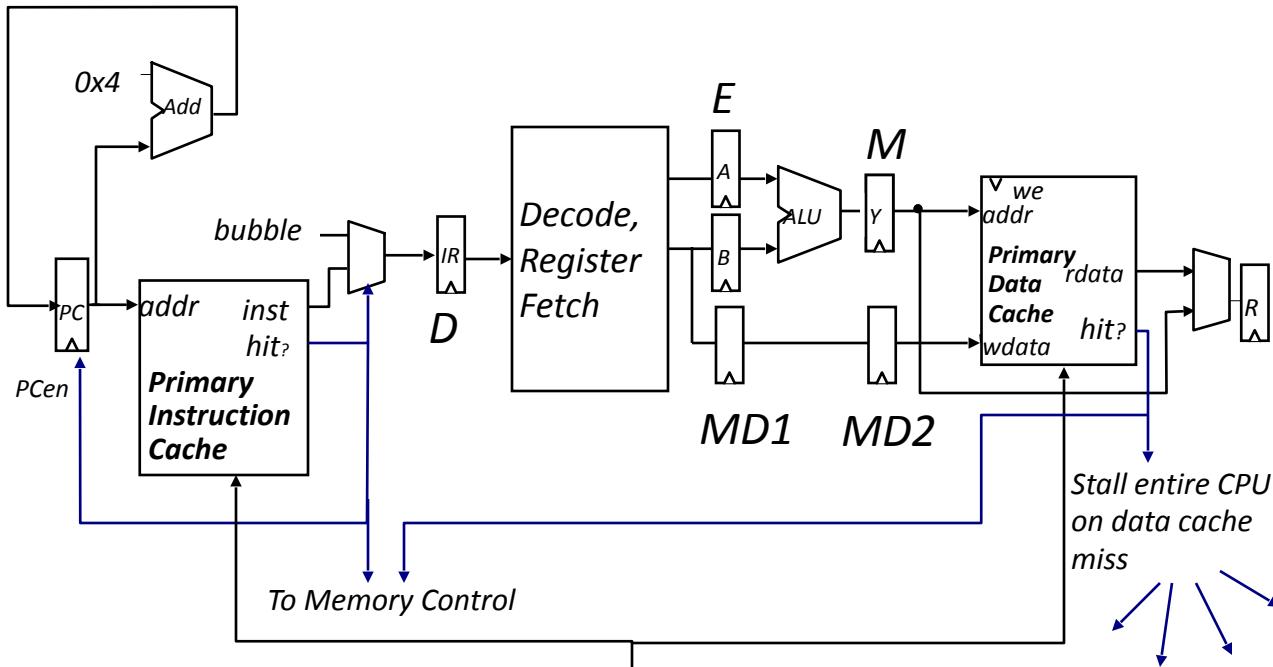


# Review from 61C

- **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.

# CPU-Cache Interaction

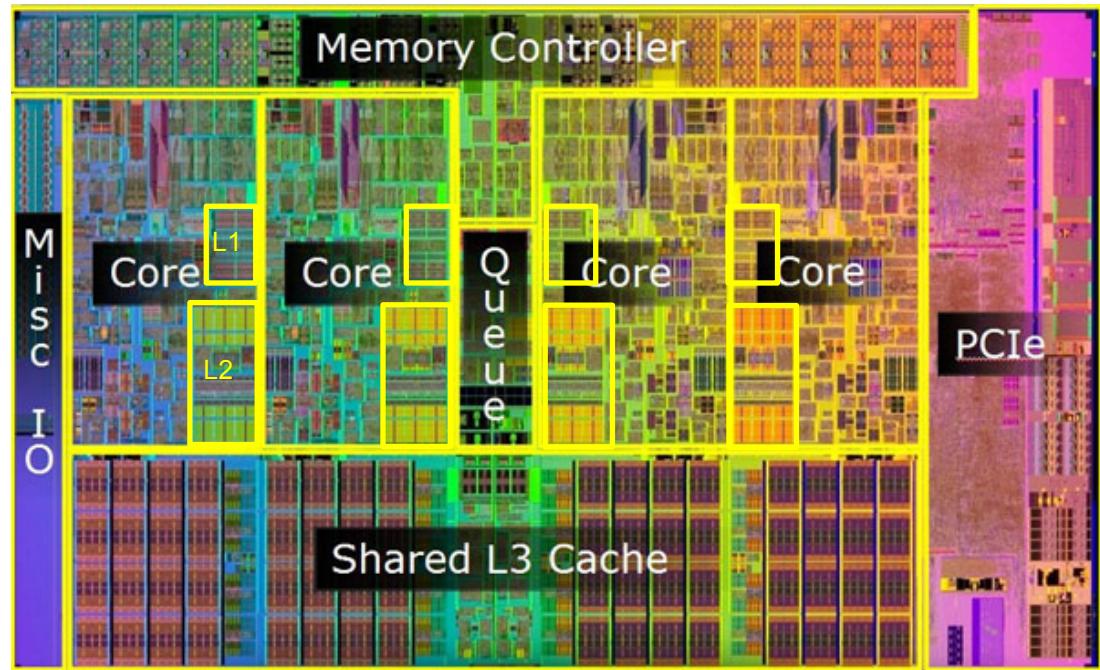
## (5-stage pipeline)



*Cache Refill Data from Lower Levels of  
Memory Hierarchy*

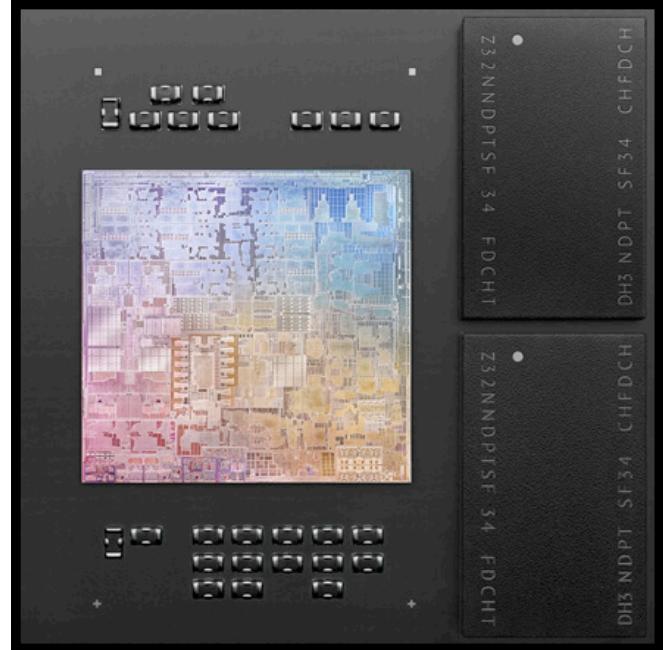
# Nahalem Die Photo (i7, i5)

- Per core:
  - 32KB L1 I-Cache (4-way set associative (SA))
  - 32KB L1 D-Cache (8-way SA)
  - 256KB unified L2 (8-way SA, 64B blocks)
  - Common L3 8MB cache
- Common L3 8MB cache

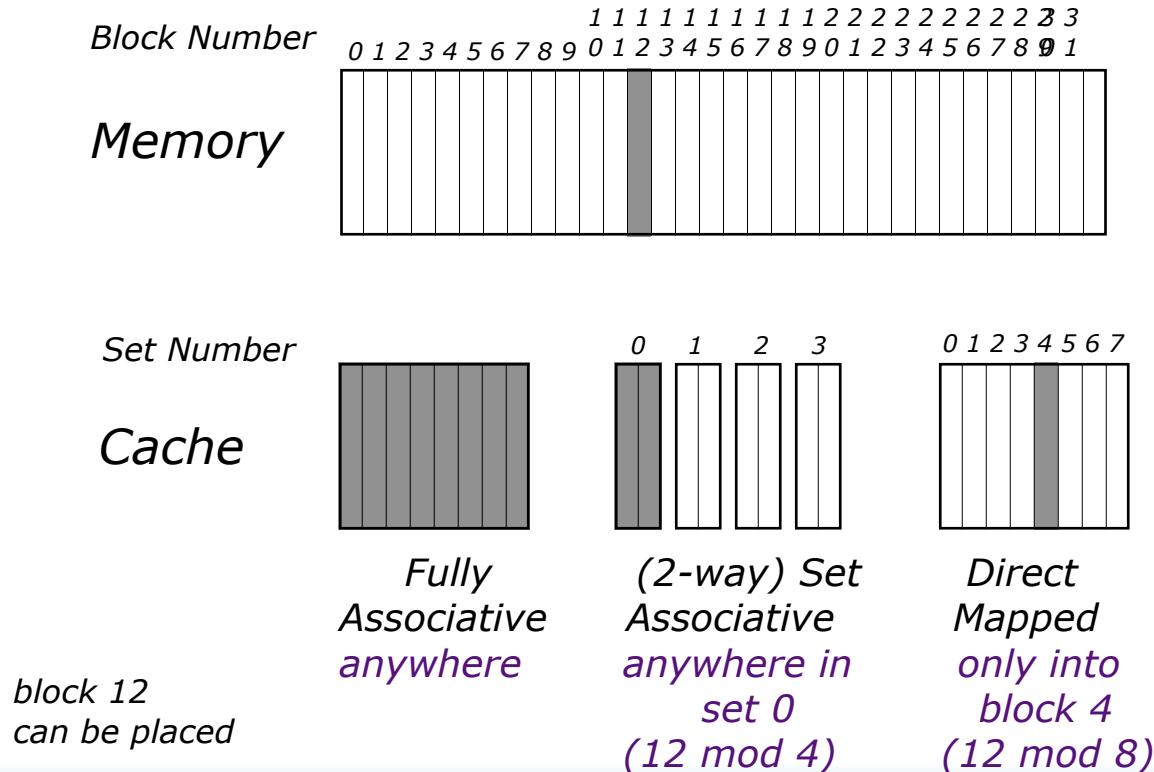


# Apple M1

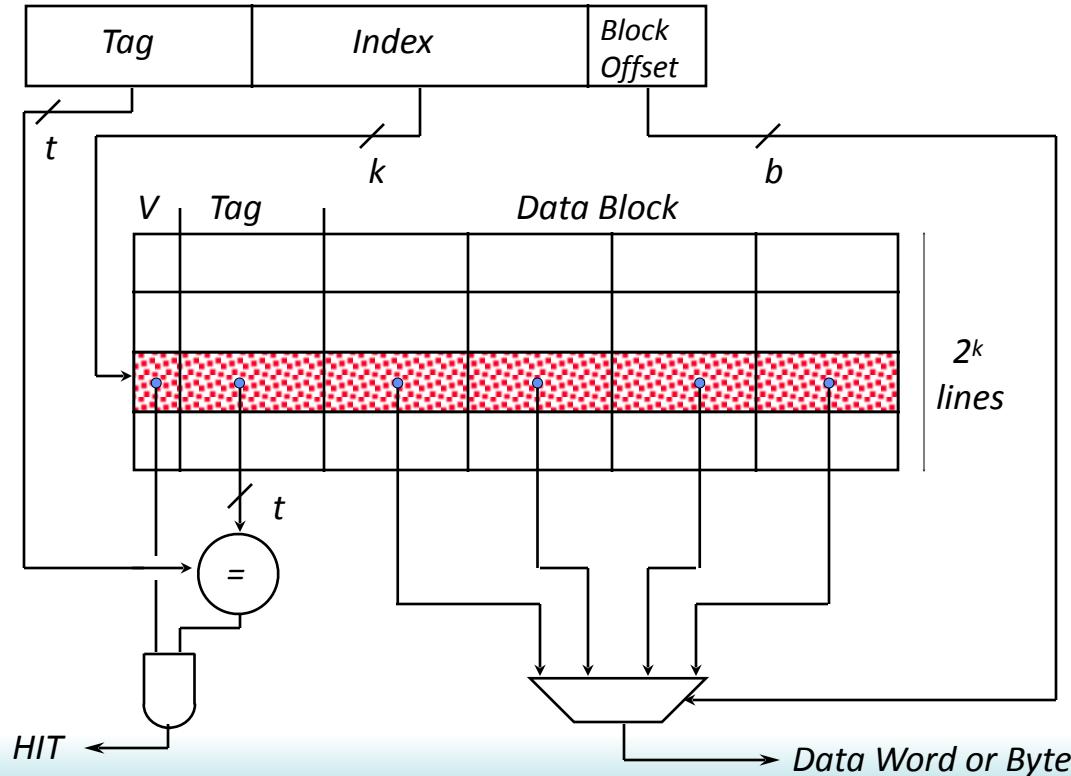
- ❑ 4 High-performance cores:
  - 192 KB of L1 I-cache
  - 128 KB of L1 data cache
  - Shared 12 MB L2 cache
- ❑ 4 Energy-efficient cores:
  - 128 KB L1 instruction cache,
  - 64 KB L1 data cache,
  - Shared 4 MB L2 cache.
- ❑ The SoC also has a 16MB System Level Cache shared by the GPU.



# Placement Policy



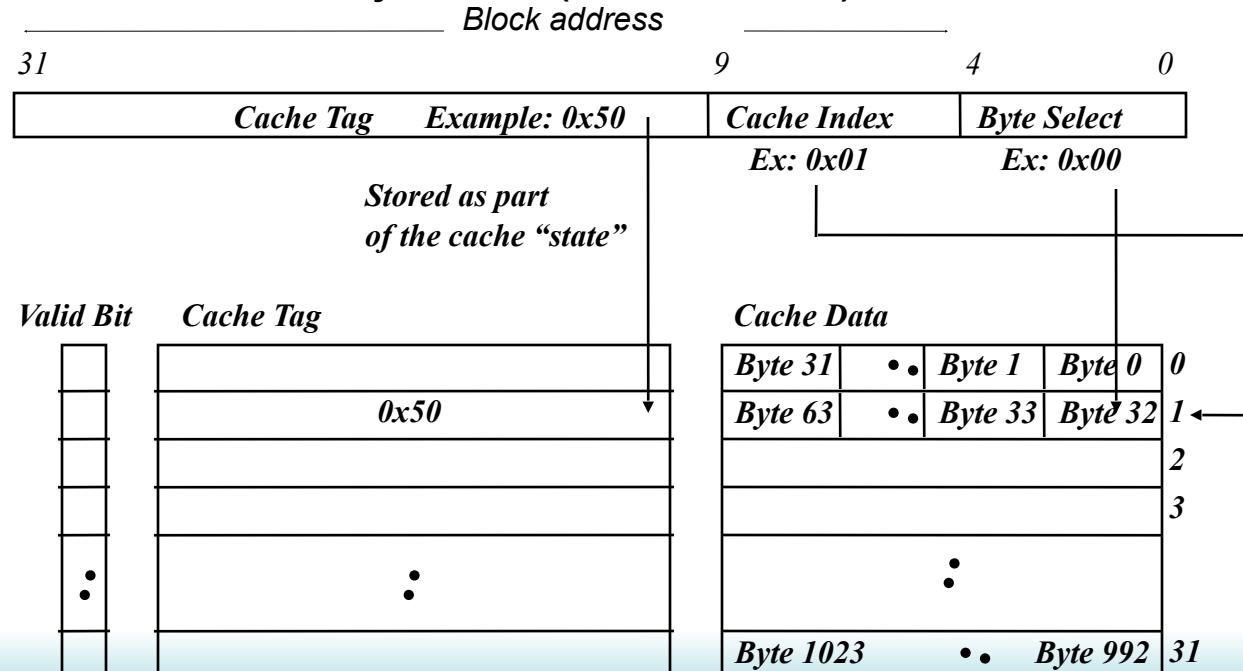
# Direct-Mapped Cache



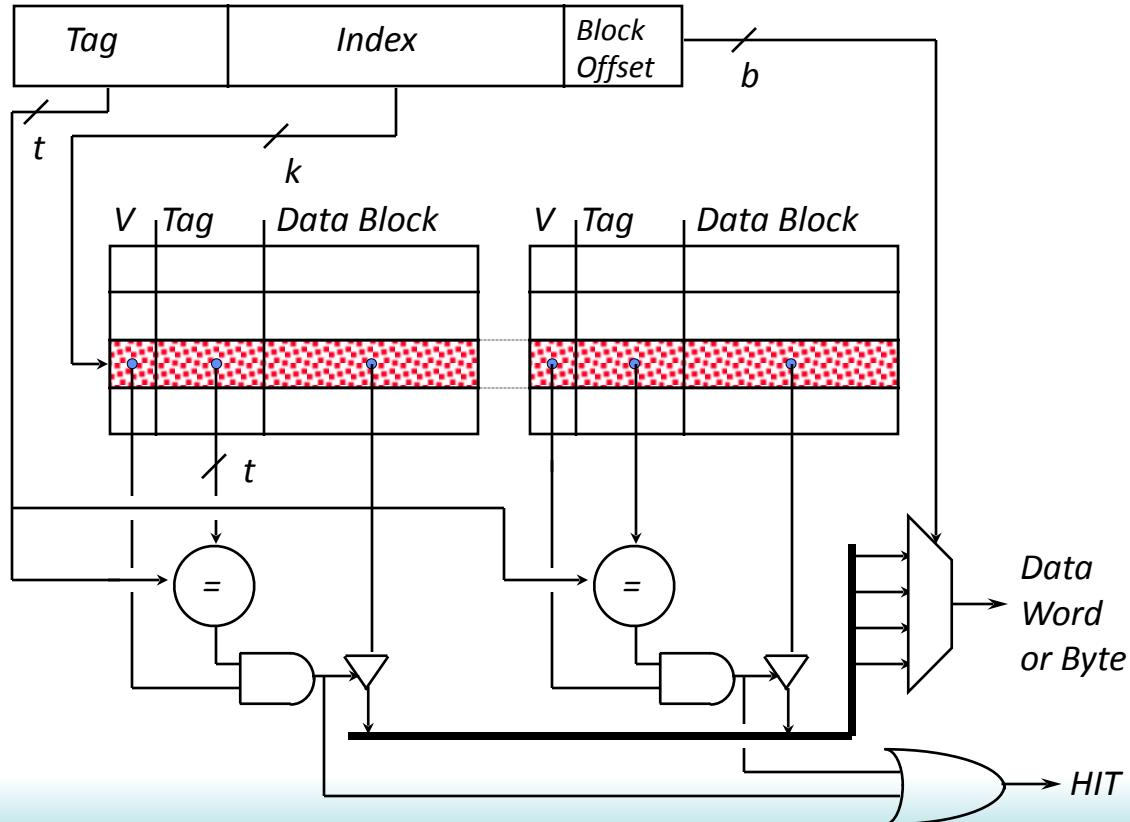
## Example: 1 KB Direct Mapped Cache with 32 B Blocks

For a  $2^N$  byte cache:

- The uppermost  $(32 - N)$  bits are always the Cache Tag
- The lowest  $M$  bits are the Byte Select (Block Size =  $2^M$ )



# 2-Way Set-Associative Cache



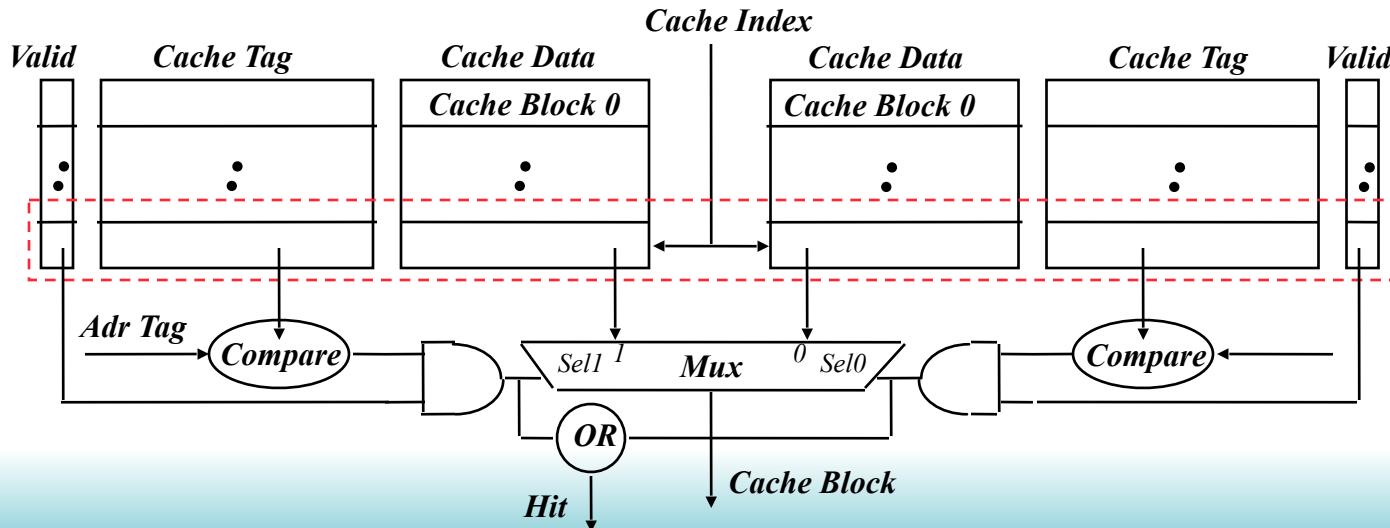
# Set Associative Cache

**N-way set associative:** N entries for each Cache Index

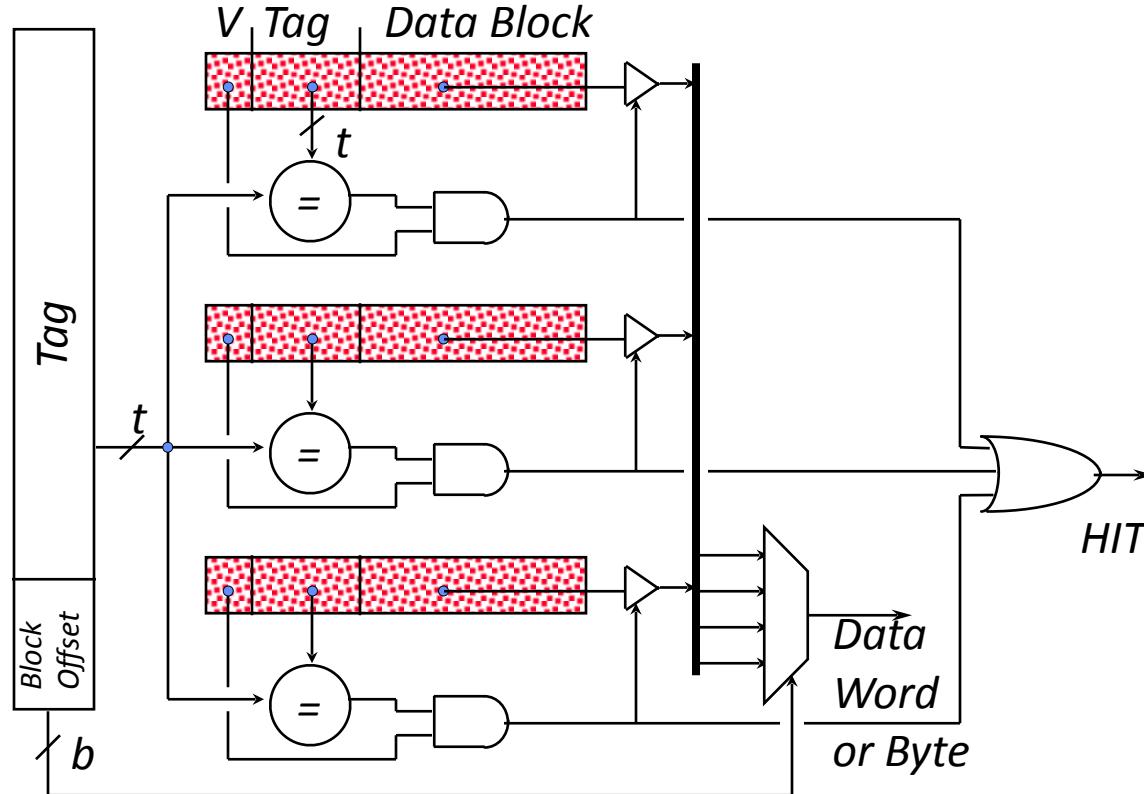
- (N direct mapped caches operates in parallel)

**Example: Two-way set associative cache**

- Cache Index selects a “set” from the cache
- The two tags in the set are compared to the input in parallel
- Data is selected based on the tag result



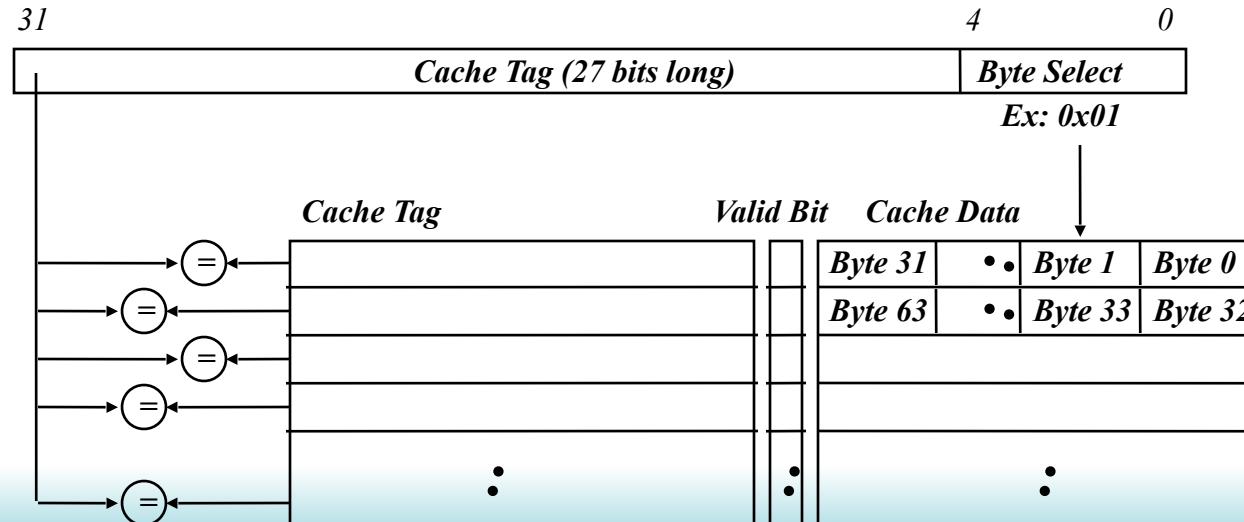
# Fully Associative Cache



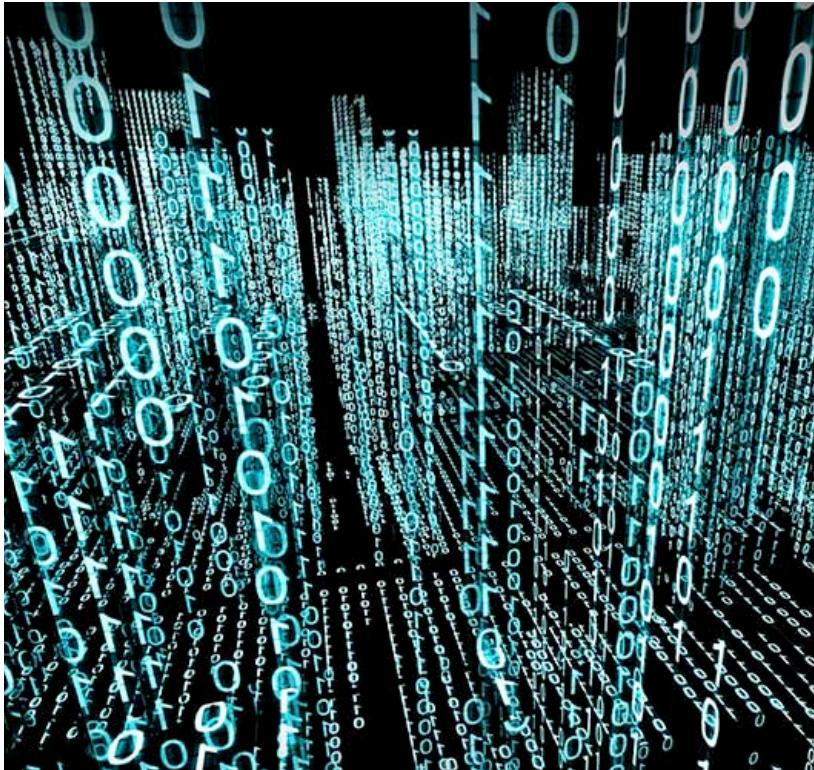
# Fully Associative

## Fully Associative Cache

- No Cache Index
- For read, compare the Cache Tags of all cache entries in parallel
- Example: Block Size = 32 B blocks, we need N 27-bit comparators



## *RAM Blocks and the Project*



# *Processor Design Considerations (FPGA Version)*

## **□ Register File: Consider distributed RAM (LUT RAM)**

- Size is close to what is needed: distributed RAM primitive configurations are 32 or 64 bits deep. Extra width is easily achieved by parallel arrangements.
- LUT-RAM configurations offer multi-porting options - useful for register files.
- Asynchronous read, might be useful by providing flexibility on where to put register read in the pipeline.

## **□ Instruction / Data Memories : Consider Block RAM**

- Higher density, lower cost for large number of bits
- A single 36kbit Block RAM implements 1K 32-bit words.
- Configuration stream based initialization, permits a simple “boot strap” procedure.

# *Processor Design Considerations (ASIC Version)*

## **Register File: use synthesized RAM**

- At this size (1k bits) synthesized is competitive with dense RAM block
- Latch-based instead of flip-flop-based would save on area.
- Asynchronous read, might be useful by providing flexibility on where to put register read in the pipeline.

## **Instruction / Data Caches : Use generated dense Block RAM**

- Higher density, lower cost for large number of bits
- We will provide for you