



# Memory Hierarchy

CS:APP 6.4 intro

These slides adapted from materials provided by the textbook authors.

# Memory Heirarchy

- Storage technologies and trends
- Locality of reference
- Caching in the memory hierarchy

# Example Memory

## Hierarchy

Smaller,  
faster,  
and  
costlier  
(per byte)  
storage  
devices

Larger,  
slower,  
and  
cheaper  
(per byte)  
storage  
devices

L6:

Remote secondary storage  
(e.g., Web servers)

L5:

Local secondary storage  
(local disks)

L4:

Main memory  
(DRAM)

L0:

Regs

L1:

L1 cache  
(SRAM)

L2:

L2 cache  
(SRAM)

L3:

L3 cache  
(SRAM)

CPU registers hold words  
retrieved from the L1 cache.

L1 cache holds cache lines  
retrieved from the L2 cache.

L2 cache holds cache lines  
retrieved from L3 cache

L3 cache holds cache lines  
retrieved from main memory.

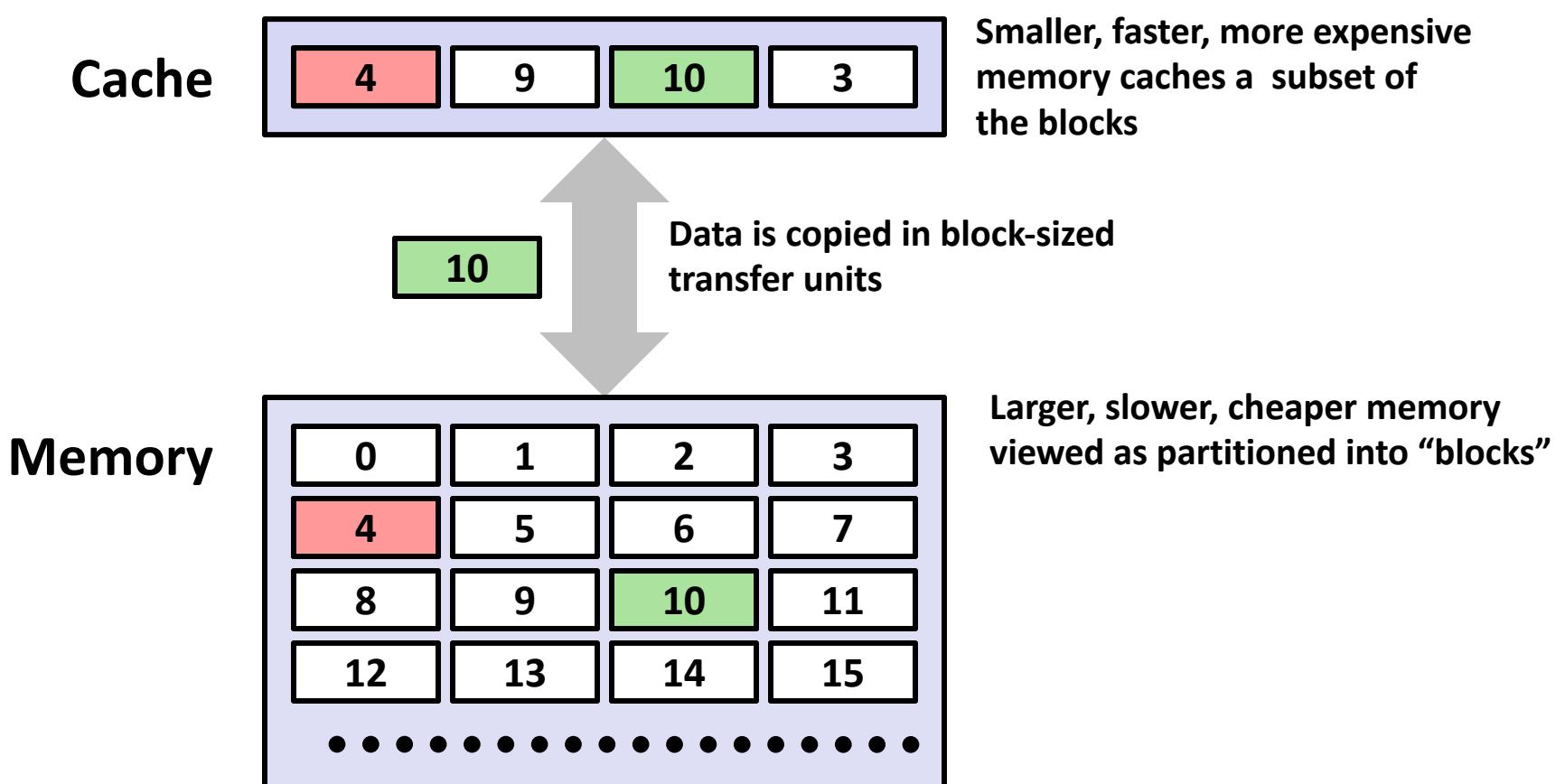
Main memory holds  
disk blocks retrieved  
from local disks.

Local disks hold files  
retrieved from disks  
on remote servers

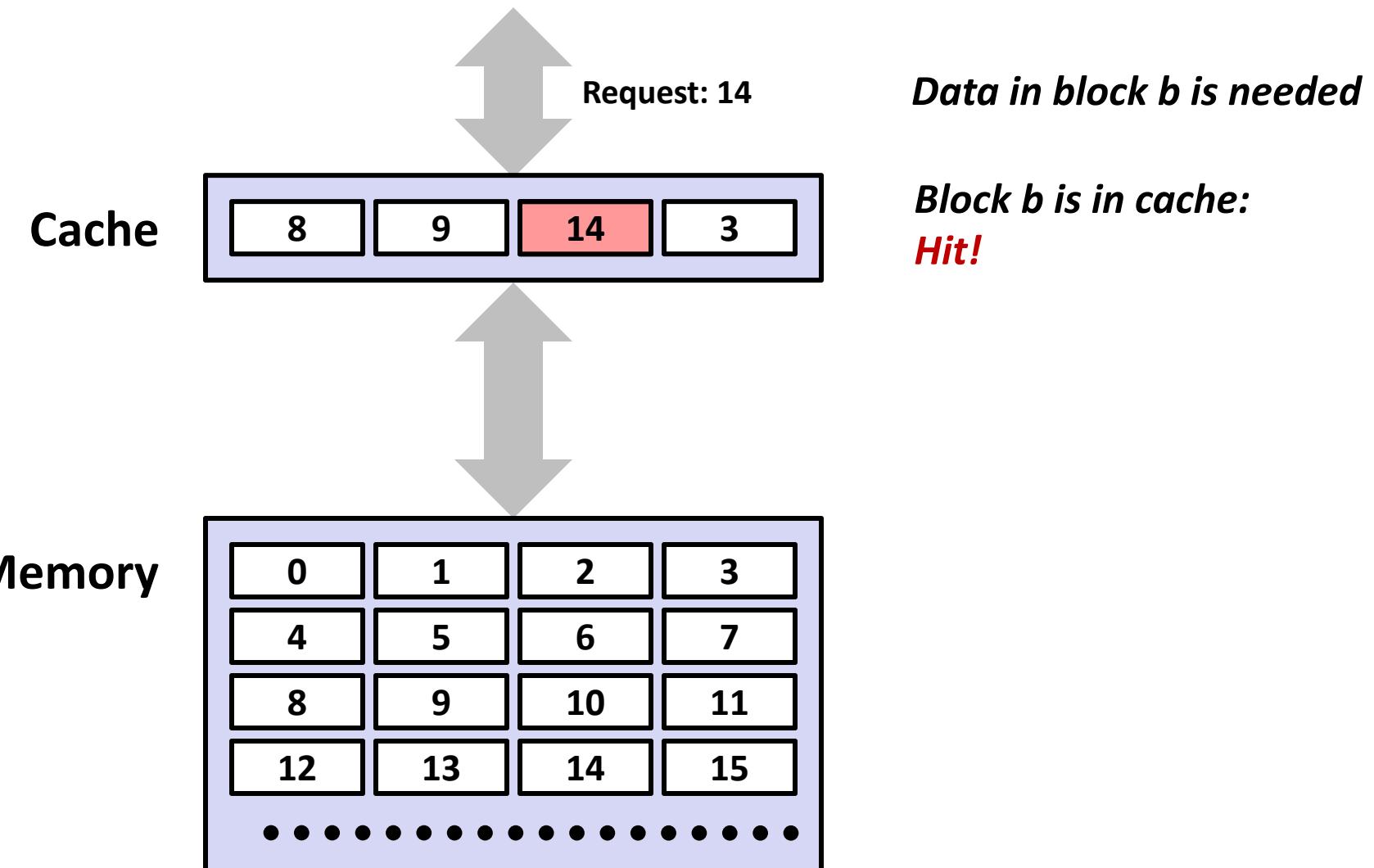
# Caches

- ***Cache:*** A smaller, faster storage device that acts as a staging area for a subset of the data in a larger, slower device.
- **Fundamental idea of a memory hierarchy:**
  - For each  $k$ , the faster, smaller device at level  $k$  serves as a cache for the larger, slower device at level  $k+1$ .
- **Why do memory hierarchies work?**
  - Because of locality, programs tend to access the data at level  $k$  more often than they access the data at level  $k+1$ .
  - Thus, the storage at level  $k+1$  can be slower, and thus larger and cheaper per bit.
- ***Big Idea:*** The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.

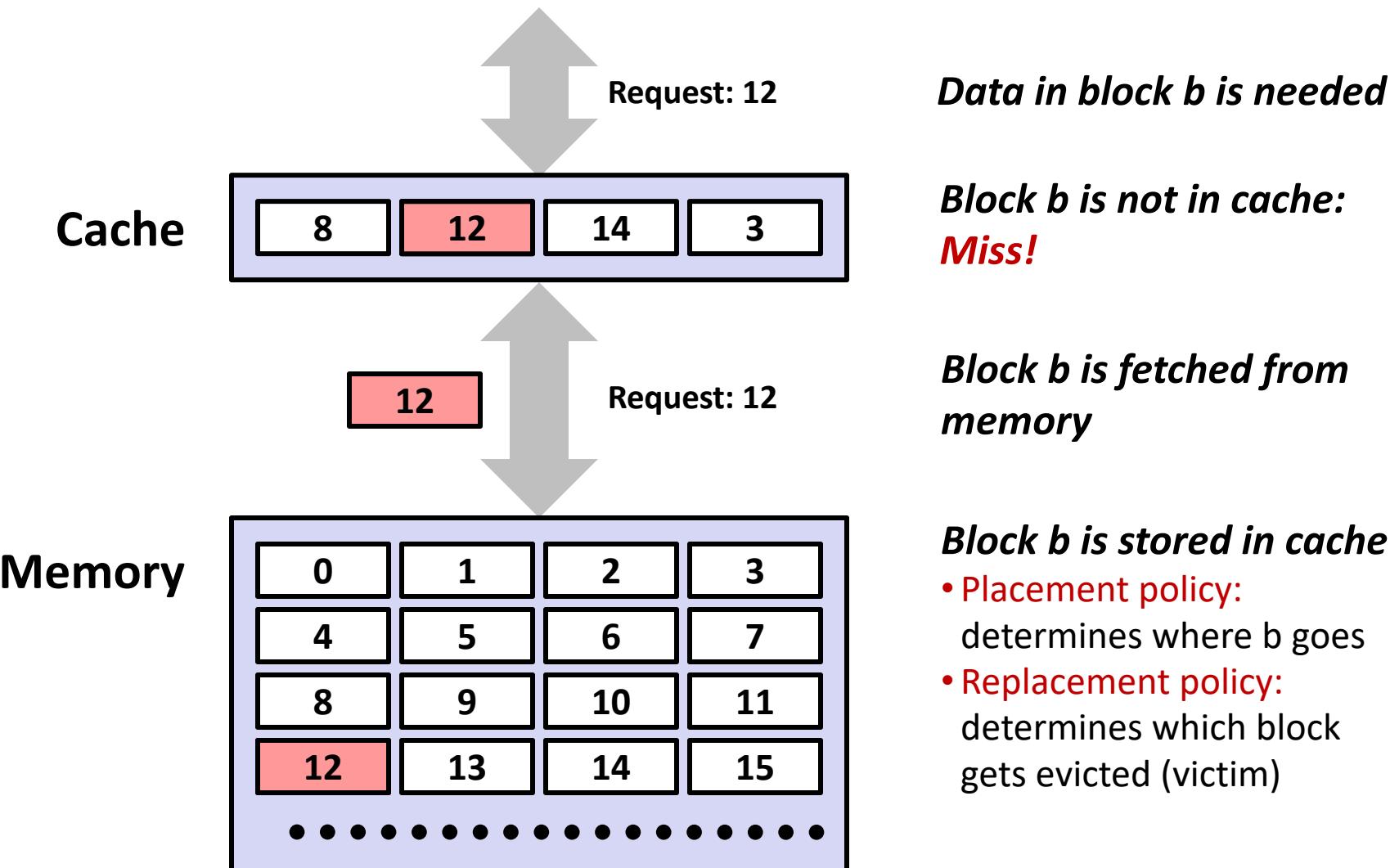
# General Cache Concepts



# General Cache Concepts: Hit



# General Cache Concepts: Miss



# General Caching Concepts:

## Types of Cache Misses

### ■ **Cold (compulsory) miss**

- Cold misses occur because the cache is empty.

### ■ **Conflict miss**

- Most caches limit blocks at level  $k+1$  to a small subset (sometimes a singleton) of the block positions at level  $k$ .
  - E.g. Block  $i$  at level  $k+1$  must be placed in block  $(i \bmod 4)$  at level  $k$ .
- Conflict misses occur when the level  $k$  cache is large enough, but multiple data objects all map to the same level  $k$  block.
  - E.g. Referencing blocks 0, 8, 0, 8, 0, 8, ... would miss every time.

### ■ **Capacity miss**

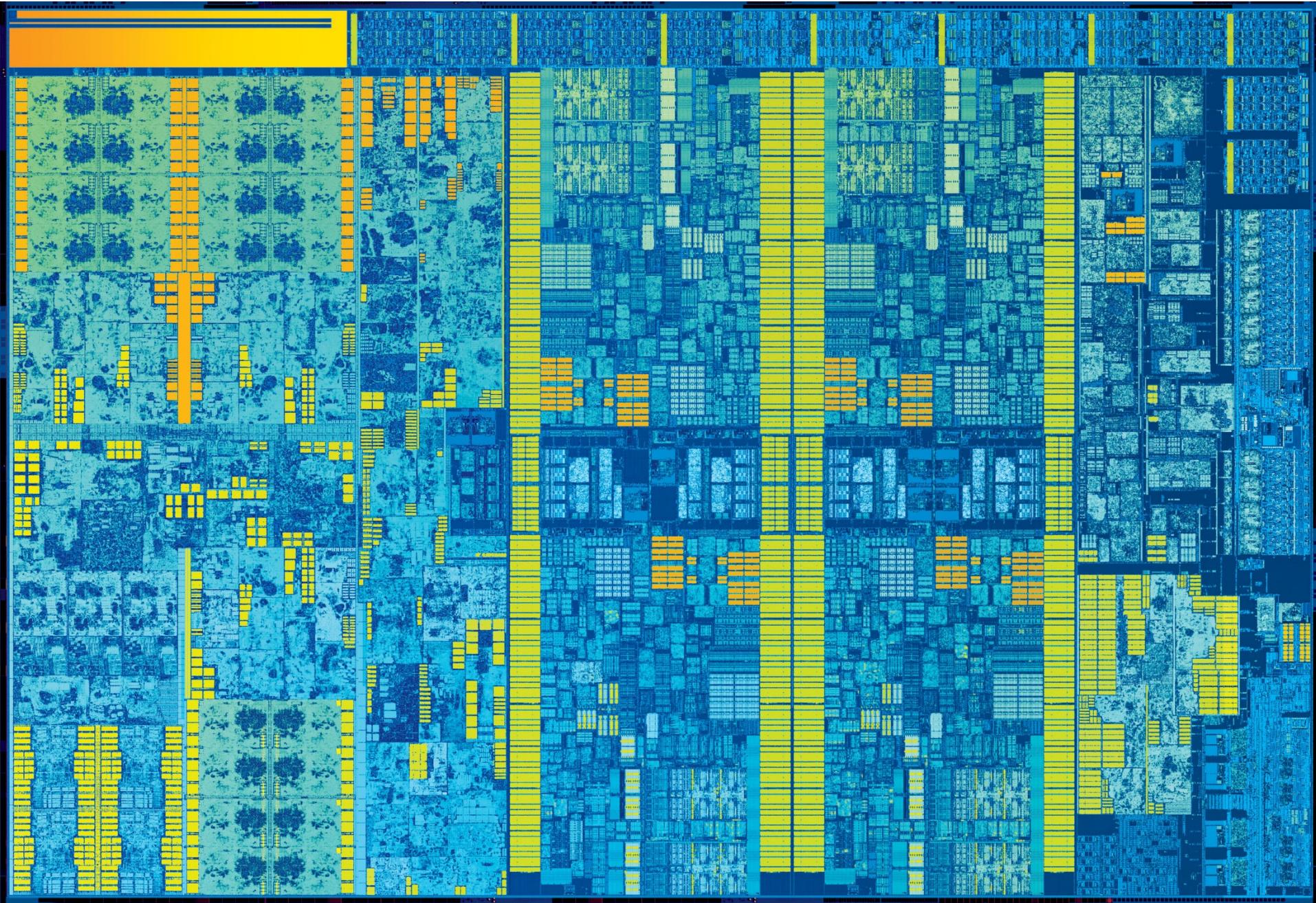
- Occurs when the set of active cache blocks (**working set**) is larger than the cache.

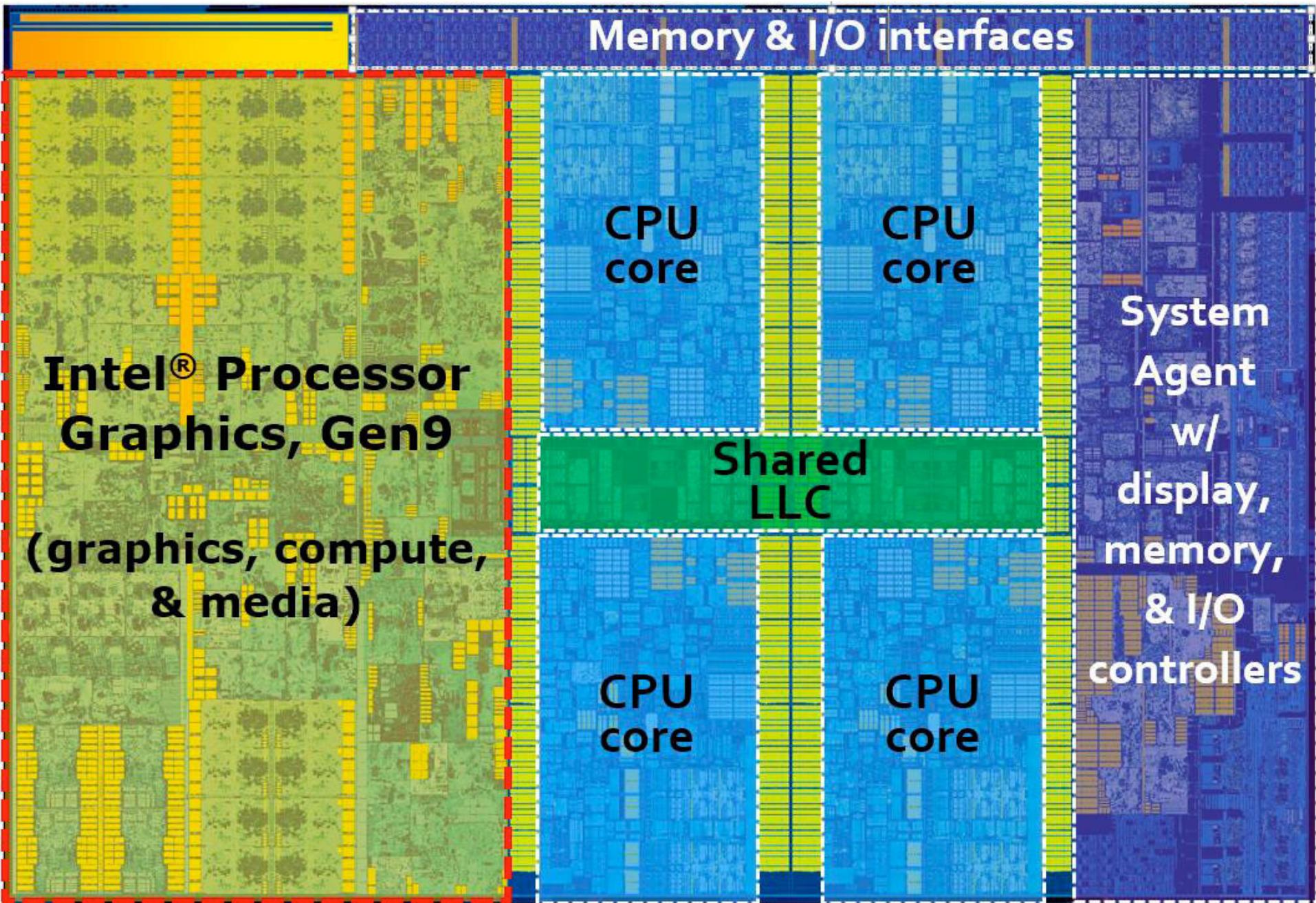
# Examples of Caching in the Mem. Hierarchy

Cache Type	What is Cached?	Where is it Cached?	Latency (cycles)	Managed By
Registers	4-8 bytes words	CPU core	0	Compiler
TLB	Address translations	On-Chip TLB	0	Hardware MMU
L1 cache	64-byte blocks	On-Chip L1	4	Hardware
L2 cache	64-byte blocks	On-Chip L2	10	Hardware
Virtual Memory	4-KB pages	Main memory	100	Hardware + OS
Buffer cache	Parts of files	Main memory	100	OS
Disk cache	Disk sectors	Disk controller	100,000	Disk firmware
Network buffer cache	Parts of files	Local disk	10,000,000	NFS client
Browser cache	Web pages	Local disk	10,000,000	Web browser
Web cache	Web pages	Remote server disks	1,000,000,000	Web proxy server

# Summary

- The speed gap between CPU, memory and mass storage continues to widen.
- Well-written programs exhibit a property called *locality*.
- Memory hierarchies based on *caching* close the gap by exploiting locality.

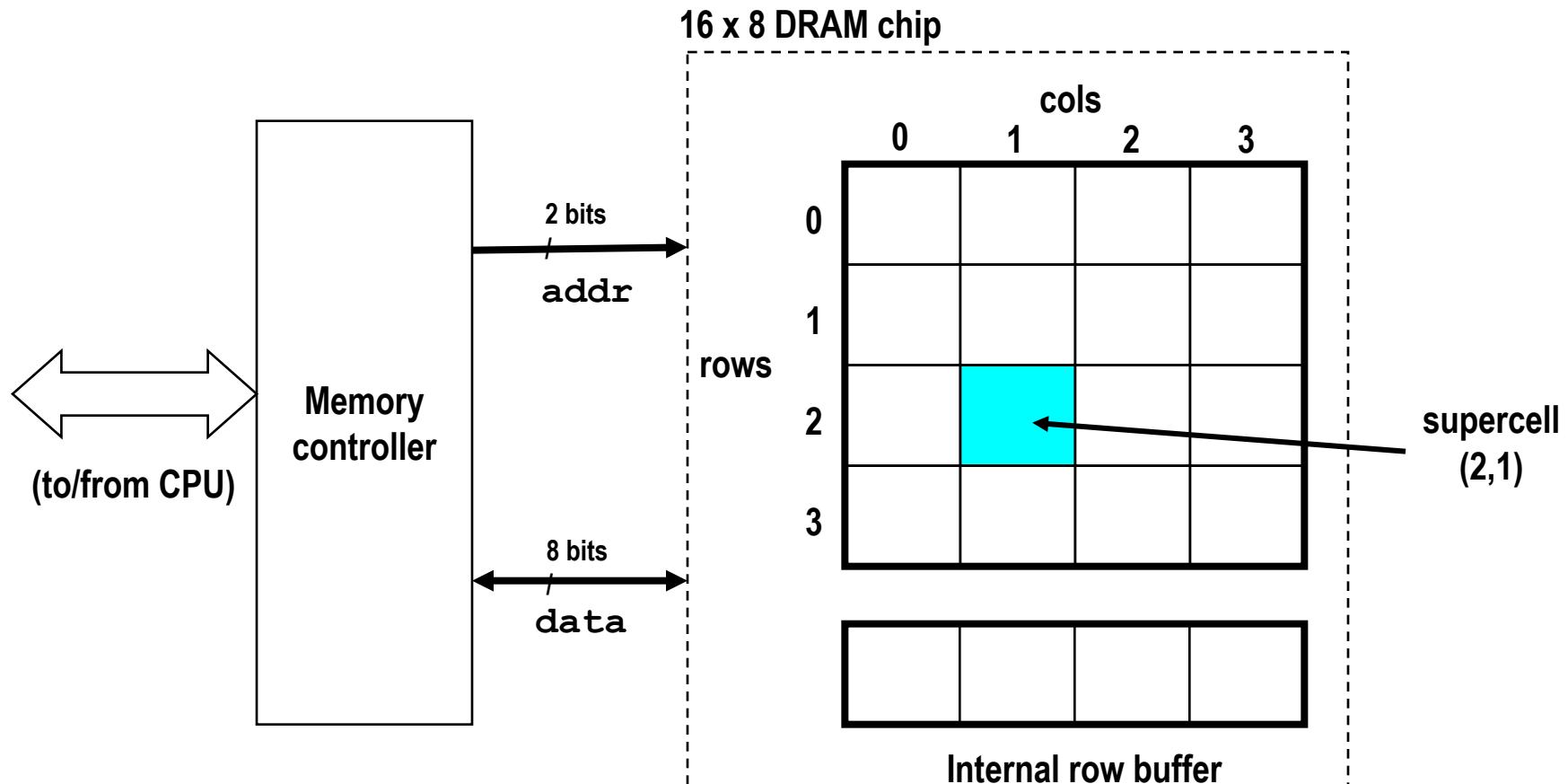




# Conventional DRAM Organization

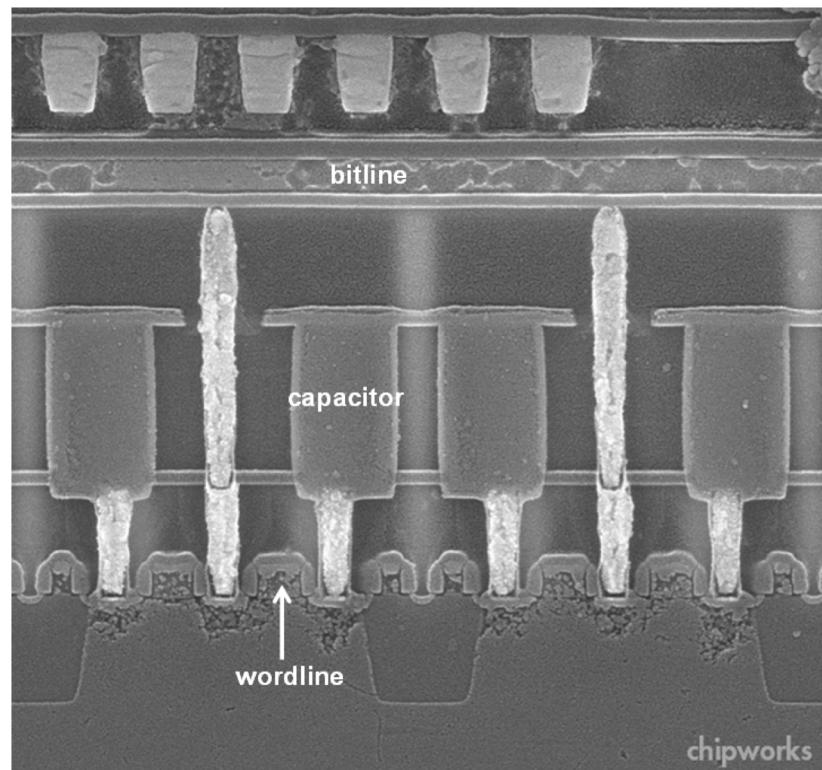
## ■ $d \times w$ DRAM:

- $dw$  total bits organized as  $d$  **supercells** of size  $w$  bits



# What's in a DRAM cell?

- Dynamic RAM stores a charge in a capacitor
- That charge depletes over time and when read
- Reading and writing DRAM requires both reading and writing to replenish charge
- “refresh” needed as charge depeletes

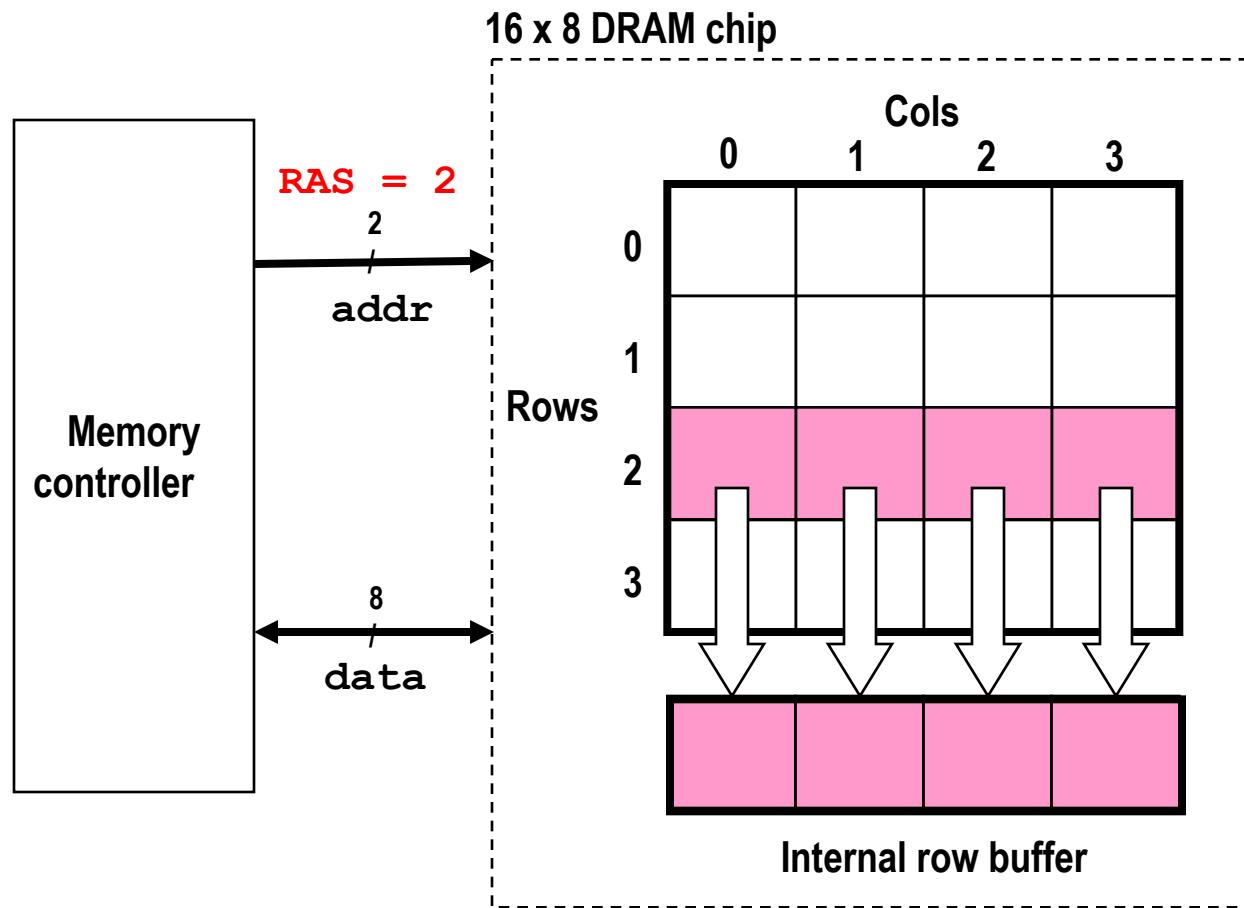


65nm Embedded DRAM in Xbox GPU

# Reading DRAM Supercell (2,1)

Step 1(a): Row access strobe (**RAS**) selects row 2.

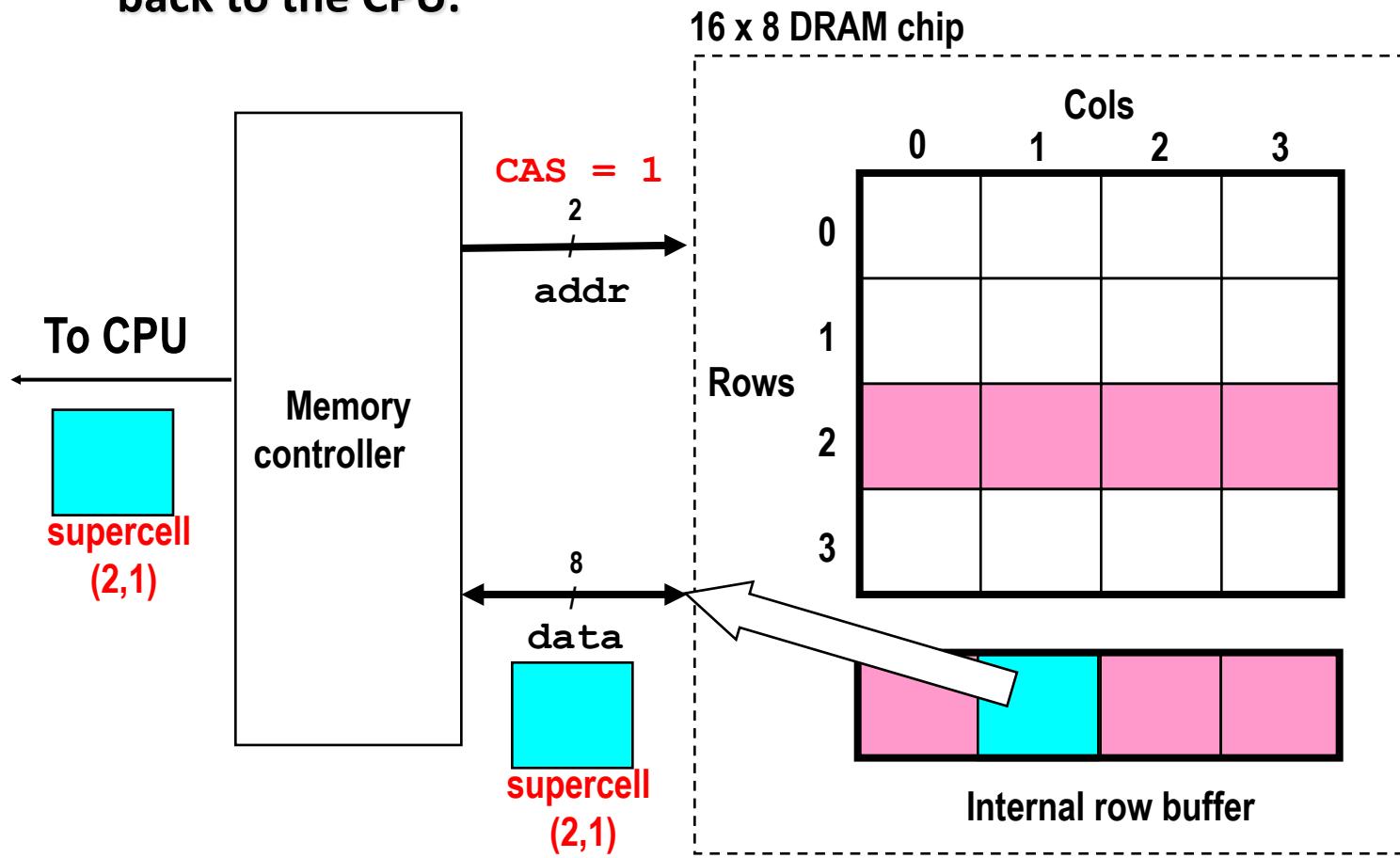
Step 1(b): Row 2 copied from DRAM array to row buffer.



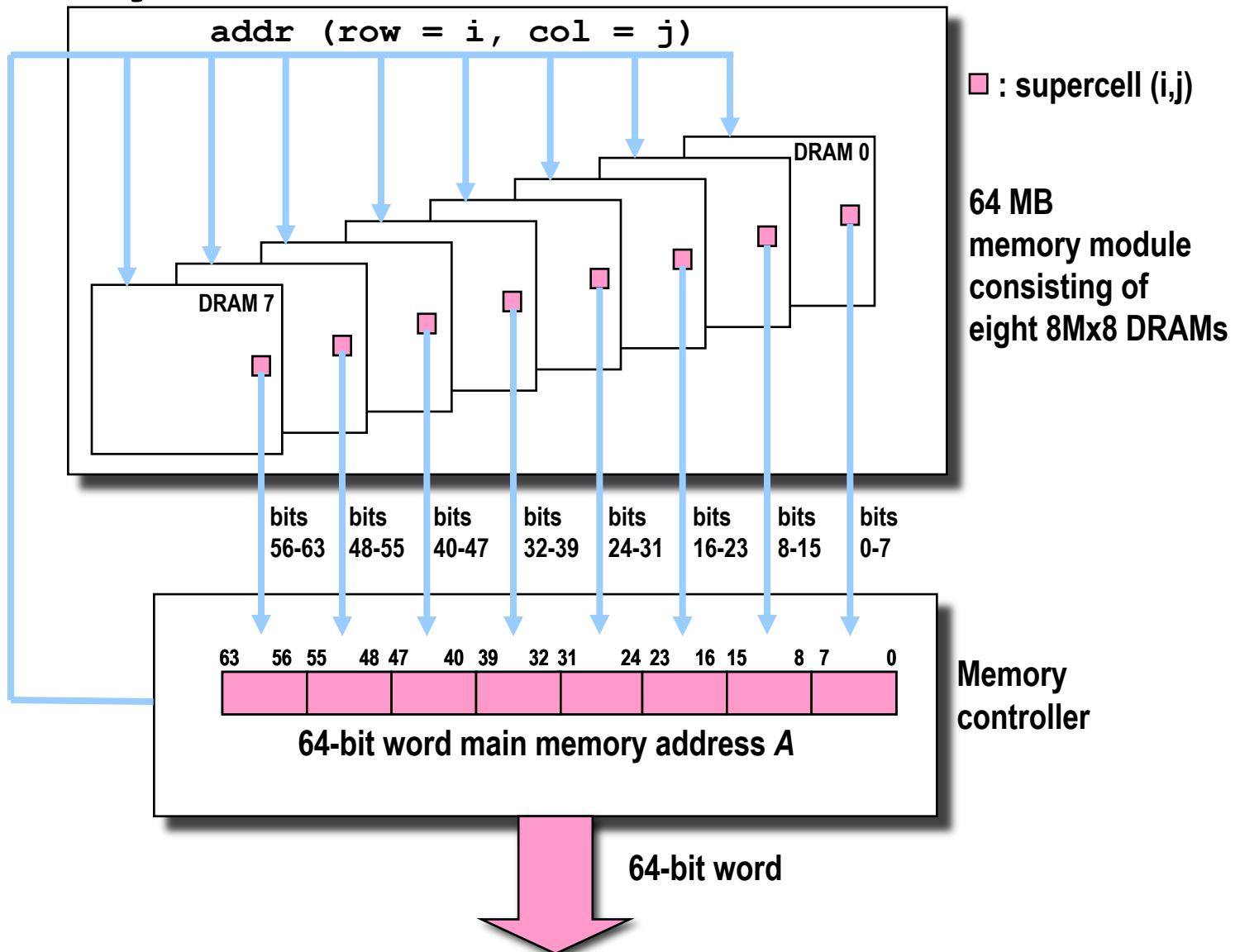
# Reading DRAM Supercell (2,1)

Step 2(a): Column access strobe (**CAS**) selects column 1.

Step 2(b): Supercell (2,1) copied from buffer to data lines, and eventually back to the CPU.



# Memory Modules



# Enhanced DRAMs

- **Basic DRAM cell has not changed since its invention in 1966.**
  - Commercialized by Intel in 1970.
- **DRAM cores with better interface logic and faster I/O :**
  - Synchronous DRAM (**SDRAM**)
    - Uses a conventional clock signal instead of asynchronous control
    - Allows reuse of the row addresses (e.g., RAS, CAS, CAS, CAS)
  - Double data-rate synchronous DRAM (**DDR SDRAM**)
    - Double edge clocking sends two bits per cycle per pin
    - Different types distinguished by size of small prefetch buffer:
      - **DDR** (2 bits), **DDR2** (4 bits), **DDR3** (8 bits)
    - By 2010, standard for most server and desktop systems
    - Intel Core i7 supports only DDR3 SDRAM

# Storage Trends

## SRAM

Metric	1985	1990	1995	2000	2005	2010	2015	2015:1985
\$/MB	2,900	320	256	100	75	60	320	116
access (ns)	150	35	15	3	2	1.5	200	115

## DRAM

Metric	1985	1990	1995	2000	2005	2010	2015	2015:1985
\$/MB	880	100	30	1	0.1	0.06	0.02	44,000
access (ns)	200	100	70	60	50	40	20	10
typical size (MB)	0.256	4	16	64	2,000	8,000	16,000	62,500

## Disk

Metric	1985	1990	1995	2000	2005	2010	2015	2015:1985
\$/GB	100,000	8,000	300	10	5	0.3	0.03	3,333,333
access (ms)	75	28	10	8	5	3	3	25
typical size (GB)	0.01	0.16	1	20	160	1,500	3,000	300,000

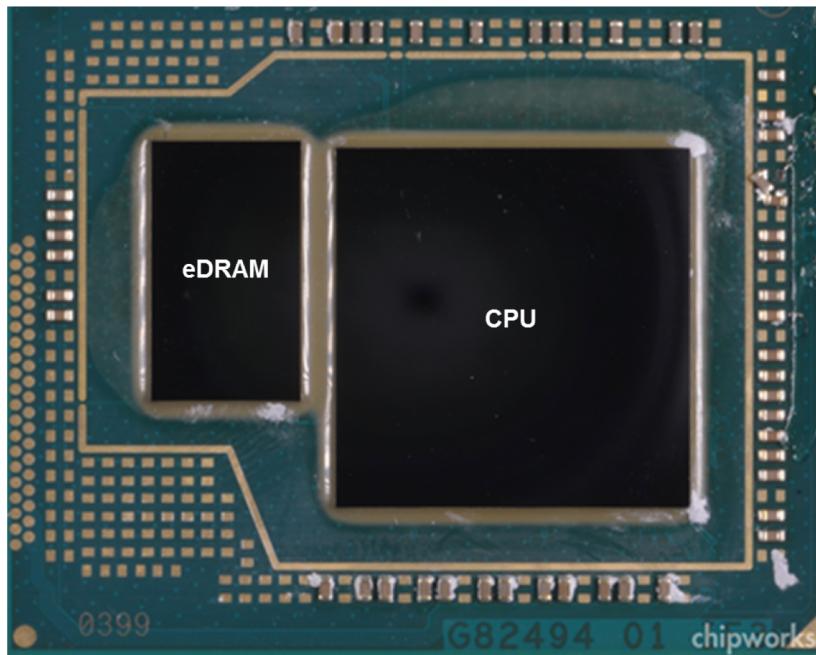
# CPU Clock Rates

Inflection point in computer history  
when designers hit the “Power Wall”



	1985	1990	1995	2003	2005	2010	2015	2015:1985
CPU	80286	80386	Pentium	P-4	Core 2	Core i7(n)	Core i7(h)	
Clock rate (MHz)	6	20	150	3,300	2,000	2,500	3,000	500
Cycle time (ns)	166	50	6	0.30	0.50	0.4	0.33	500
Cores	1	1	1	1	2	4	4	4
Effective cycle time (ns)	166	50	6	0.30	0.25	0.10	0.08	2,075
							(n) Nehalem processor	
							(h) Haswell processor	

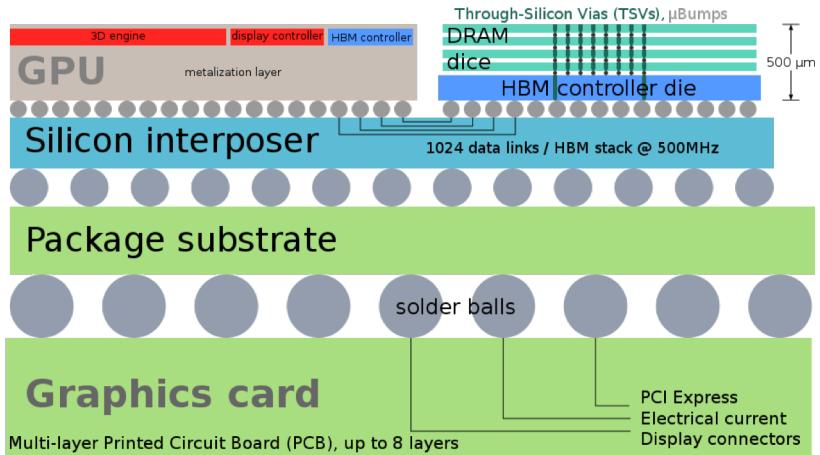
# Battling Memory Wall - Embedded DRAM



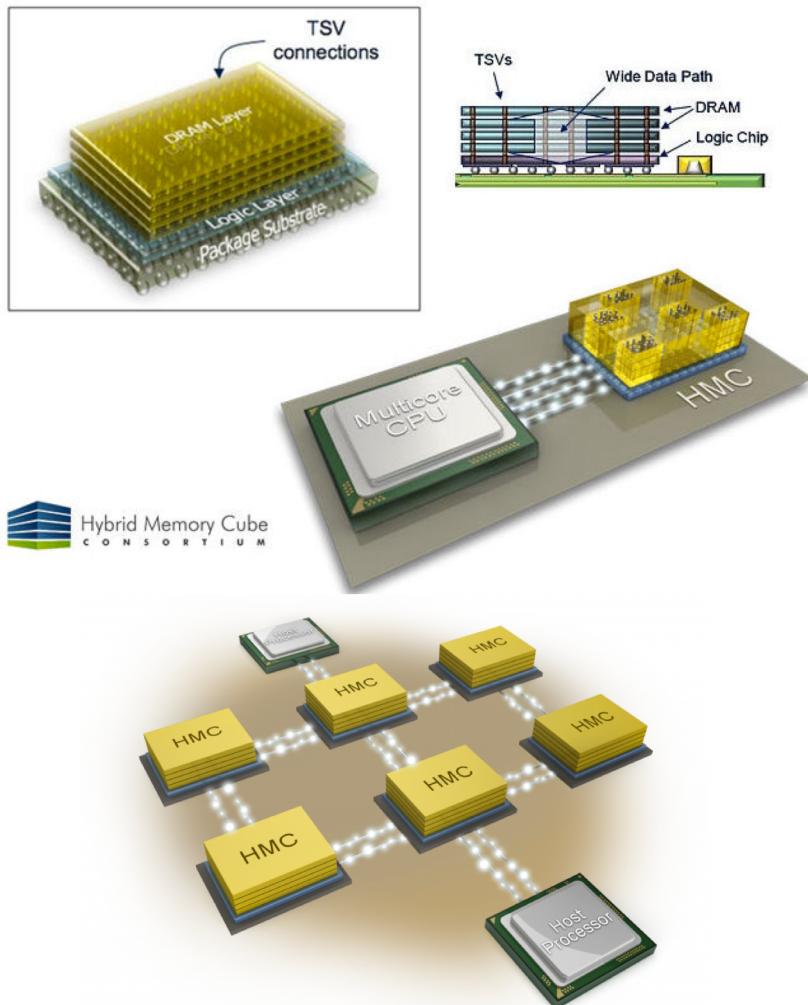
- DRAM density is much higher than SRAM
- Intel used “embedded DRAM” in Haswell CPU’s as L4 cache
- Large (128MB) co-packaged with CPU
- Future designs will have 16GB L4 cache

# Battling Memory Wall - More embedded DRAM

- HBM – High Bandwidth Memory
- Co-packaged with CPU / GPU
- Many short wires



# Battling Memory Wall - Hybrid Memory Cube



- Overcome memory wall by shorter wires
- Turn memory system into a network
- Not clear all will be adopted, but technology path is to higher integration
- Won't mix & match