# CS 354 - Machine Organization & Programming
## Tuesday Oct 17, and Thursday Oct 19, 2023

**Print paper copies of this outline for best use.**

**Week 07 Activity: Heap Practice Assignment**

**Project p3A:** DUE on or before Friday 10/20 & p3B on10/27

**Homework 3:** DUE on or before Monday 10/16

**Learning Objectives**

- Describe the relative difference in speed and size of various types of memory and storage.
- Identify and describe the units of transfer used by each storage type.
- Define, identify, and describe spatial locality and temporal locality
- Identify good locality in program code
- Explain why programs with good locality work better with caching.
- Compute stride (in words) of array memory accesses
- Determine if common algorithms produce good or bad locality for each type.
- Define and use basic cache terminology
- Convert hex to binary, use bits of an address to determine if the address is in a given cache
- Extract bits and compute the set index, tag bits, and byte, determine if byte is in the cache

**This Week: MEMORY MANAGEMENT via CACHING blocks of memory for fast access**

| | |
|---|---|
| Memory Hierarchy | Rethinking Addressing |
| Locality & Caching | Designing a Cache: Sets and Tags |
| Bad Locality | Basic Cache Lines |
| Caching: Basic Idea & Terms | Basic Cache Operation |
| Designing a Cache: Blocks | Basic Cache Practice |

**Next Week**: Vary cache set size and Cache Writes
B&O 6.4.3 Set Associative Caches
6.4.4 Fully Associative Caches
6.4.5 Issues with Writes
6.4.6 Anatomy of a Real Cache Hierarchy
6.4.7 Performance Impact of Cache Parameters

Note: p4A and p4B will be released next week
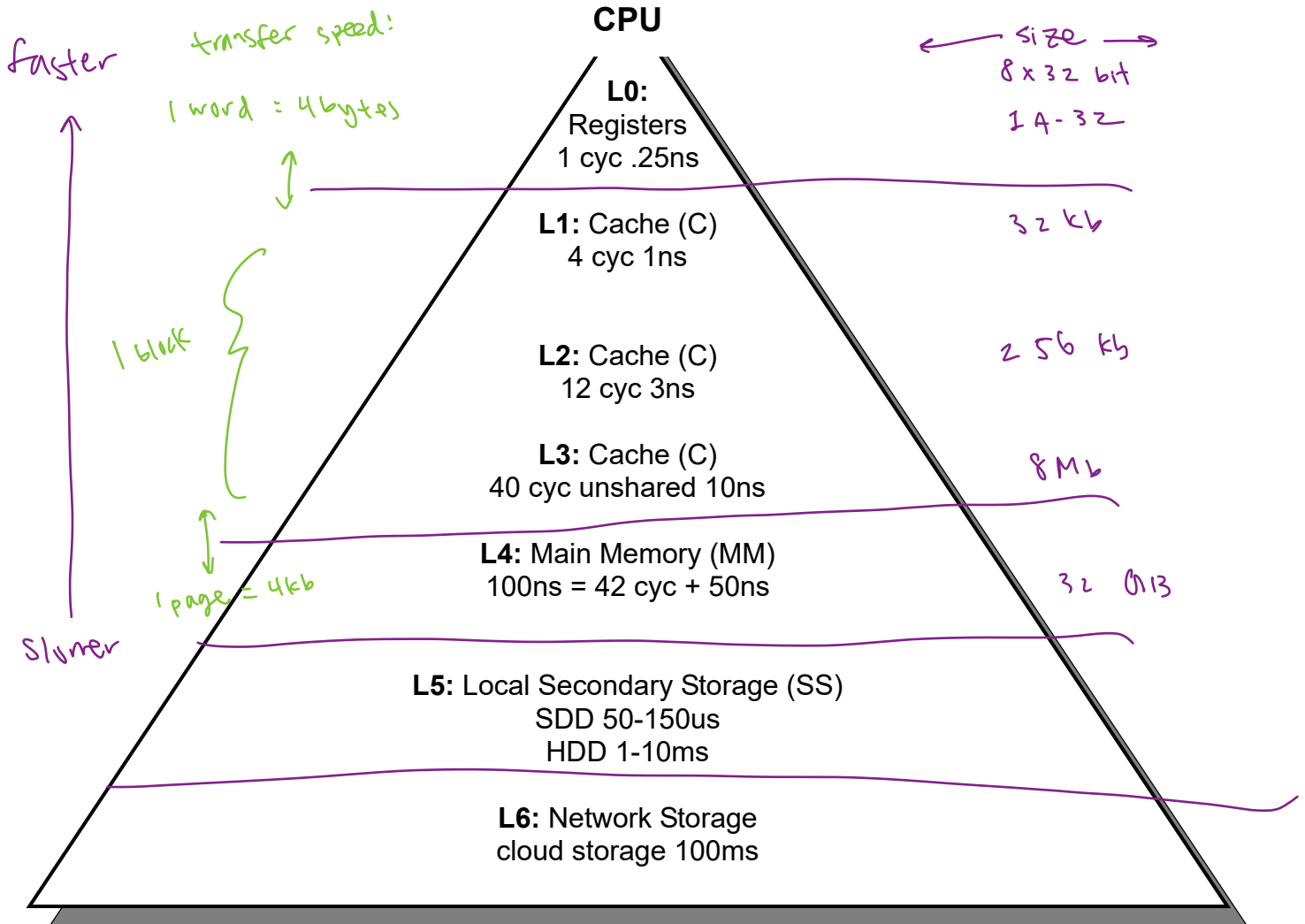
Get p3A and p3B done this week and avoid the rush!
p3 - implement and test alloc (partA) and free (partB) by Monday and submit progress
p3 - implement immediate coalescing by Wednesday and submit progress
p3 - complete testing and debugging by Friday and complete final submission

# Memory Hierarchy

❋ *The memory hierarchygives the illusion of having lots of fast memory.*

**CPU**

*faster*

*transfer speed!*

*1 word = 4bytes*

*1 block*

*1 page = 4kb*

*Slower*

*size*
*8 x 32 bit*
*1 A-32*

**L0:**
Registers
1 cyc .25ns

**L1:** Cache (C)
4 cyc 1ns

*32 kb*

**L2:** Cache (C)
12 cyc 3ns

*256 kb*

**L3:** Cache (C)
40 cyc unshared 10ns

*8 Mb*

**L4:** Main Memory (MM)
100ns = 42 cyc + 50ns

*32  Gb13*

**L5:** Local Secondary Storage (SS)
SDD 50-150us
HDD 1-10ms

**L6:** Network Storage
cloud storage 100ms

## Cache

**is a smaller faster mem that acts as a staging area for data stored in a larger slower mem**

## Memory Units

*word*: size used by CPU     transfer betweenL1 & CPU     *= 4bytes / word*

*block*: size used by C     transfer betweenC levels & MM     *= 32bytes / block*

*page*: size used by MM     transfer betweenMM & SS     *= 4 kb/ page*

## Memory Transfer Time: https://simple.wikipedia.org/wiki/Orders_of_magnitude_(time)

*cpu cycles*:**used to measure time**

*latency*:**memory access time (delay)**

# Locality & Caching

**What?**

*temporal locality*: when a recently accessed memory location    *int i in for*
is repeatedly accessed in the near future

*spatial locality*: when a recently accessed memory location
   is followed by nearby memory locations being accessed in the near future

locality is designed into    hardware
                              o.s.
                              application

**Example**

```
int sumArray(int a[], int size, int step) {
   int sum = 0;
   for (int i = 0; i < size; i += step)
      sum += a[i];
   return sum;
}
```

→ List the variables that clearly demonstrate temporal locality.   *i, sum, stop*

→ List the variables that clearly demonstrate spatial locality.   *a , a[i]*

*stride*:   step size in words   (4 bytes)
            good spacial locality when stride ~ 1 word

❇ *The caching system uses locality to predict what the cpu will need in the near future.*

**How?** The caching system

temporal: anticipates data will be reused so it copies   *value into cache*

spatial: anticipates nearby data will be used so it copies   *a block*

*cache block*:   unit of memory transfered between main mem.
                 and cache levels

❇ *Programs with good locality run faster since they work better with the caching system!*

**Why?** Programs with good locality   *maximize use of data @ top hierarchy*
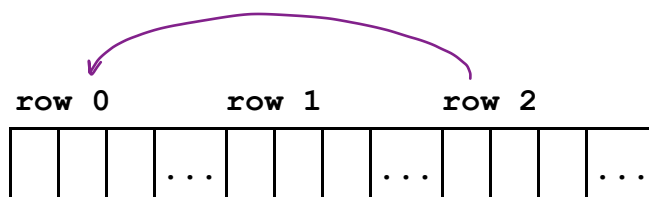
# Bad Locality

**Why is this code bad?**

```
int a[ROWS][COLS];

for (int c = 0; c < COLS; c++)
   for (int r = 0; r < ROWS; r++)
      a[r][c] = r * c;
```

access:

row 0      row 1      row 2

a    4 7    2 5 8    3 6 9

*not good spacial*

*stride length = # cols*

→ How would you improve the code to reduce stride?   *change for loop order*

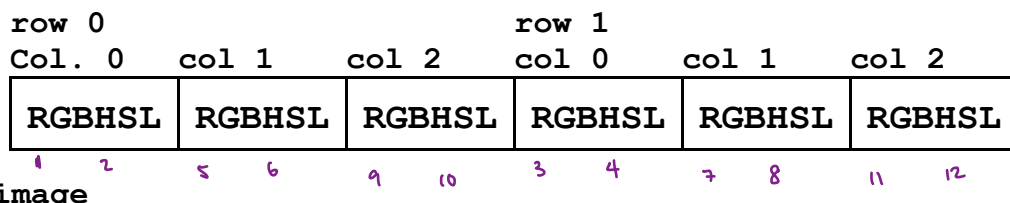**Key Questions for Determining Spatial Locality:**

1. What does the memory layout look like for the data?   *2D SAA are in row-major order*

2. What is the stride of the code across the data?

**Why is this code bad?**

*H=2 , W=3*

```
struct {
   float rgb[3];
   float hsl[3];
} image[HEIGHT][WIDTH];

for (int v = 0; v < 3; v++)
   for (int c = 0; c < WIDTH; c++)
      for (int r = 0; r < HEIGHT; r++) {
         image[r][c].rgb[v] = 0;
         image[r][c].hsl[v] = 0;
      }
```

| row 0 | | | row 1 | | |
|-------|-------|-------|-------|-------|-------|
| Col. 0 | col 1 | col 2 | col 0 | col 1 | col 2 |
| RGBHSL | RGBHSL | RGBHSL | RGBHSL | RGBHSL | RGBHSL |

image   1  2   5  6   9 10   3  4   7  8   11 12

*For each row*

*for each col*

*for each v*

*.rgb [v]*

*.hsl [v]* ~~.hsl [v]~~

*for each v*   *is better*

*.hsl [v]*

➢ How would you improve the code to reduce stride?

**Good or bad locality?**

◆ Instruction Flow:

sequencing?   *good spacial locality , poor temporal*

*no repetition*

selection?   *poor  " , poor  "*

repetition?   *dept. on loop , good temporal locality for data*

*repeated access*

◆ Searching Algorithms:

linear search   *array. good spacial , good temp for value being matched*

*linked list - poor spacial,  "  "  and temp val*

binary search   *array - poor spacial*

# Caching: Basic Idea & Terms

**Assume**: Memory is divided into 32 byte blocks and all blocks are already in main memory.
Cache L1 has 4 locations to store blocks and L2 has 16 locations to store blocks.

→ Update the memory hierarchy below given blocks are accessed in this sequence:
22,11,22,44,11,33,11,22,55,27,44

*working set* →

L1 Miss
L2 Miss
L3 Hit

*cache miss*

block not found

cold miss

when cache has room but block not there

capacity miss

when cache is too small for working set

conflict miss

when ≥ 2 blocks map to same location

*cache hit* = FASTER

when block is found in cache

*placement policies*

lvl 1. unrestricted

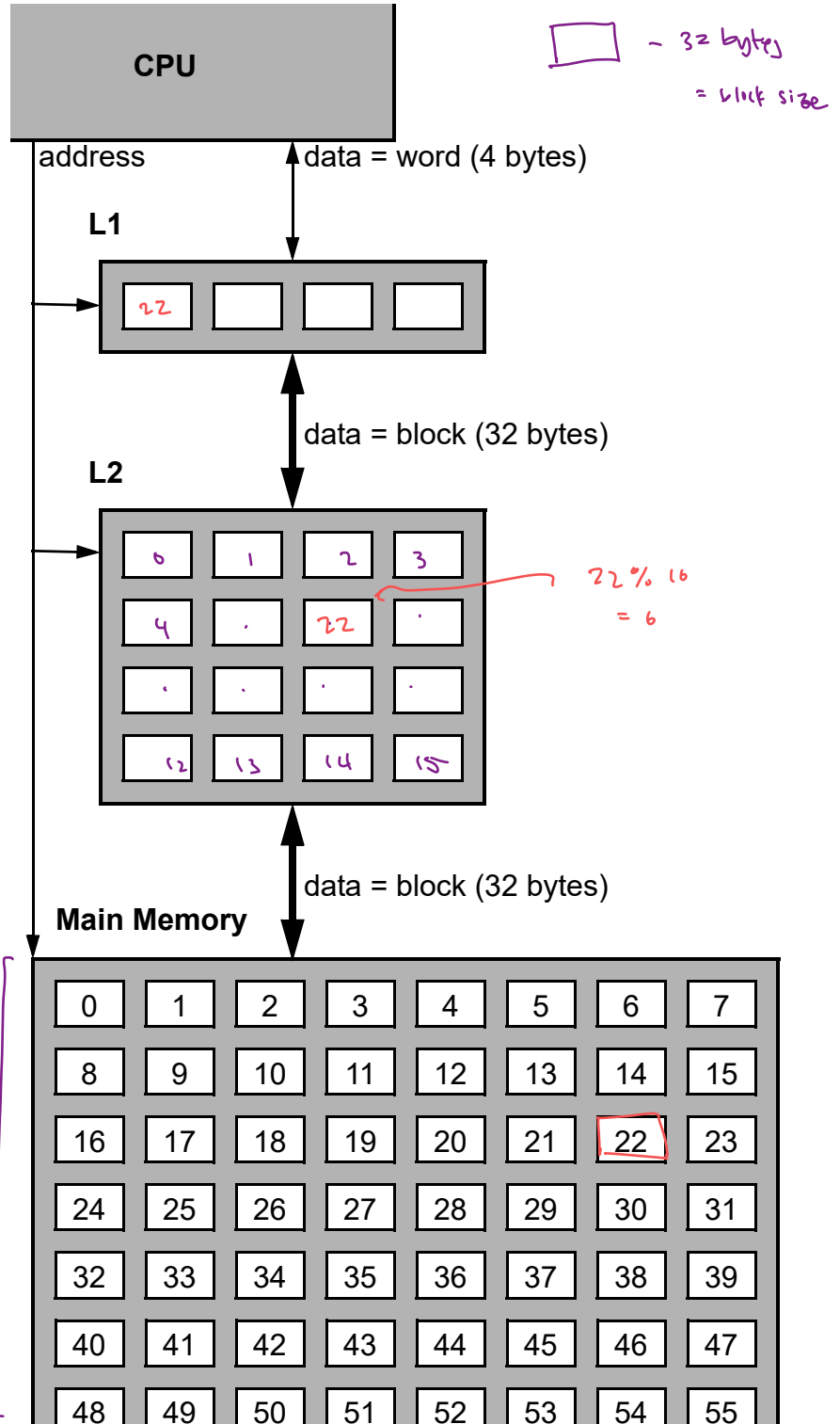lvl 2. restricted

*replacement policies*

lvl 1. choose any

lvl 2. no choice, replace existing block

*victim block*

cache block chosen to be replace

*working set*

set of blocks used during working time

**CPU**

address          data = word (4 bytes)

**L1**

| 22 | | | |

data = block (32 bytes)

**L2**

| 0 | 1 | 2 | 3 |
| 4 | . | 22 | . |
| . | . | . | . |
| 12 | 13 | 14 | 15 |

22 % 16 = 6

□ ~ 32 bytes = block size

data = block (32 bytes)

**Main Memory**

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 |
| 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 |
| 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 |

1 page everything fits

# Designing a Cache: Blocks

❋ *The bits of an address* are used to determine if blocks containing that address is in the cache

**How many bytes in an address space?**
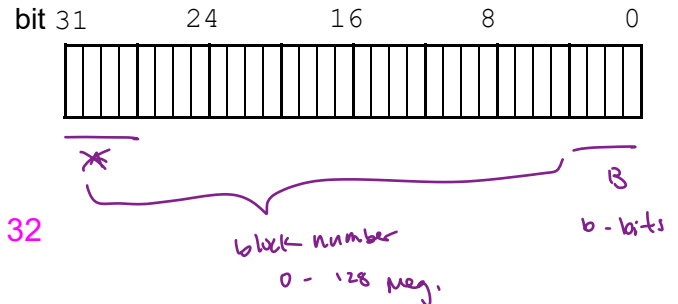
Let M be number of bytes in AS, IA-32 is 4GB

32-bit Address Breakdown
bit 31       24        16        8        0

$4GB = 2^{32}$

$M = 2^m$
$m = \log_2 M$

$\therefore 32 = \log_2(4GB)$

block number
0 - 128 meg.

B
b - bits

Thus m is number of bits in an address, IA-32 is 32

**How big is a block?**  # of bytes- 8 bits / byte

❋ *Cache blocks must be big enough* to capture spacial locality
*but small enough* to minimize latency

Let B be number of bytes per block, IA-32 is 32 bytes /block

$B = 2^b = 32 \text{ bytes}$
$b = \log_2 B = 5$

*b bits*: # of addr. bits to determine which byte in a block

   *word offset* identifies which word in a block ( 8 words /block)

   *byte offset* identifies which bytes in a word ( 4 bytes / word)

➤ What is the problem with using the most significant bits (left side) for the b bits?
we are picking bytes very far away from each other as a small change in
the most sig. bit skips many bits ⟶ lose spacial locality ∴ use least sig. bits

**How many 32-byte blocks of memory in a 32-bit address space?**

$$\frac{2^{32}}{2^5} = 2^{32-5} = 2^{27} = 2^7 \cdot 2^{20} = 128 \cdot 1MB = 128 M \text{ of blocks}$$

❋ *The remaining bits of an address encode the block number.* = 124,217,728 blocks

# Rethinking Adressing

✳ *An address identifies which byte in the VAS to access.*

✳ *An address is* divided into parts to
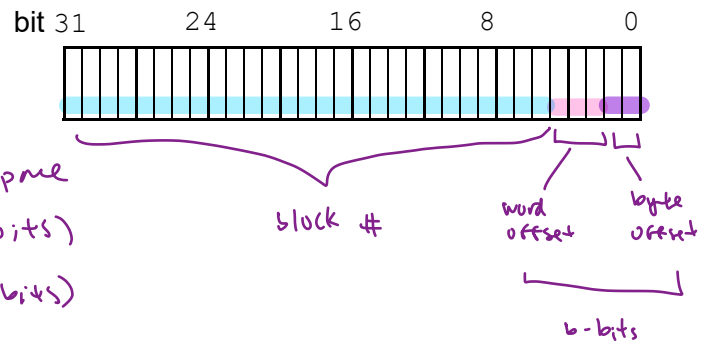
access mem. in steps

**Memory Access in Caching System**

32-bit Address Breakdown

bit 31    24    16    8    0



step 1. Identify which  block in virtual addr space

step 2. Identify which  word in block (3 bits)

step 3. Identify which  byte in word (2 bits)

block #    word offset    byte offset

6-bits

0xFFFFFFFF

**IA-32 4GB Virtual Address Space**

0x800000020

**Words in Block**

7

6

0x800000018

5

4

0x800000010

**Blocks in VAS**

**Bytes in Word**

3

3

67108865

2

1

0x8000000B →  **67108864**

2

0

0x800000008

67108863

1

67108862

0

0x800000000

① don't offset

| 8 | 0 | 0 | 0 | 0 | 0 | 0 | B | <- address in hex for a char |

| 1000 | 0000 | 0000 | 0000 | 0000 | 0000 | 000**0** | **1011** | **1.** byte #2147483659 in VAS |

| 1000 | 0000 | 0000 | 0000 | 0000 | 0000 | 000 | | **1.** MM and C access |

2 = 0 10      **2.** L1 passes to CPU

3 = 11      **3.** CPU accesses

# Designing a Cache: Sets & Tags

❋ *A cache must be searched* if unrestricted placement policy

→ Problem? slow - O(n) where n = # locations in cache (L2)

Improvement? limit (restrict) where each block can be stored (L2)

_set_: where block is uniquely mapped in a cache

— 27 most sig. bits

❋ *The block number bits of an address* are divided into 2 parts

1. set - maps blocks to specific set in cache
2. tag - uniquely identify blocks in the set

**How many sets in the cache?**

Let S be <span style="color:magenta">the number of sets in cache</span>
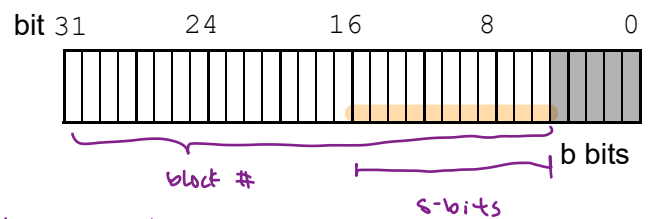
$S = 2^s$     let $S = 1024$     $S = 8192$
$s = \log_2 S$

$s = 10$ bits     $S = 13$

32-bit Address Breakdown
bit 31    24    16    8    0

block #

s-bits

b bits

⭐ _s bits_: bits that identify which set the block maps to.

(in next least sig. bits after (before) b-bits)

➤ What is the problem with using the most significant bits (left side) for the s bits?

lose spatial locality

→ How many blocks map to each set for a 32-bit AS and a cache with 1024 sets? 8192 sets?

$2^{32-5} = \dfrac{2^{27} \text{ blocks}}{2^{10} \text{ sets}} = 2^{17}$ blocks/set     $\dfrac{2^{27} \text{ blocks}}{2^{13} \text{ sets}} = 2^{14}$ blocks/set
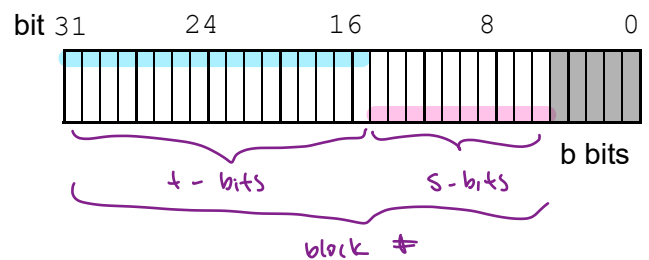
**Since different blocks map to the same set how do we know which block is in a set?**

use remaining bits as unique tag

32-bit Address Breakdown
bit 31    24    16    8    0

t - bits

s-bits

b bits

block #

_t bits_: bits of addr. that identify which blocks in the set

❋ *When a block is copied into a cache* <span style="color:magenta">*its t bits are also stored as its tag*</span>
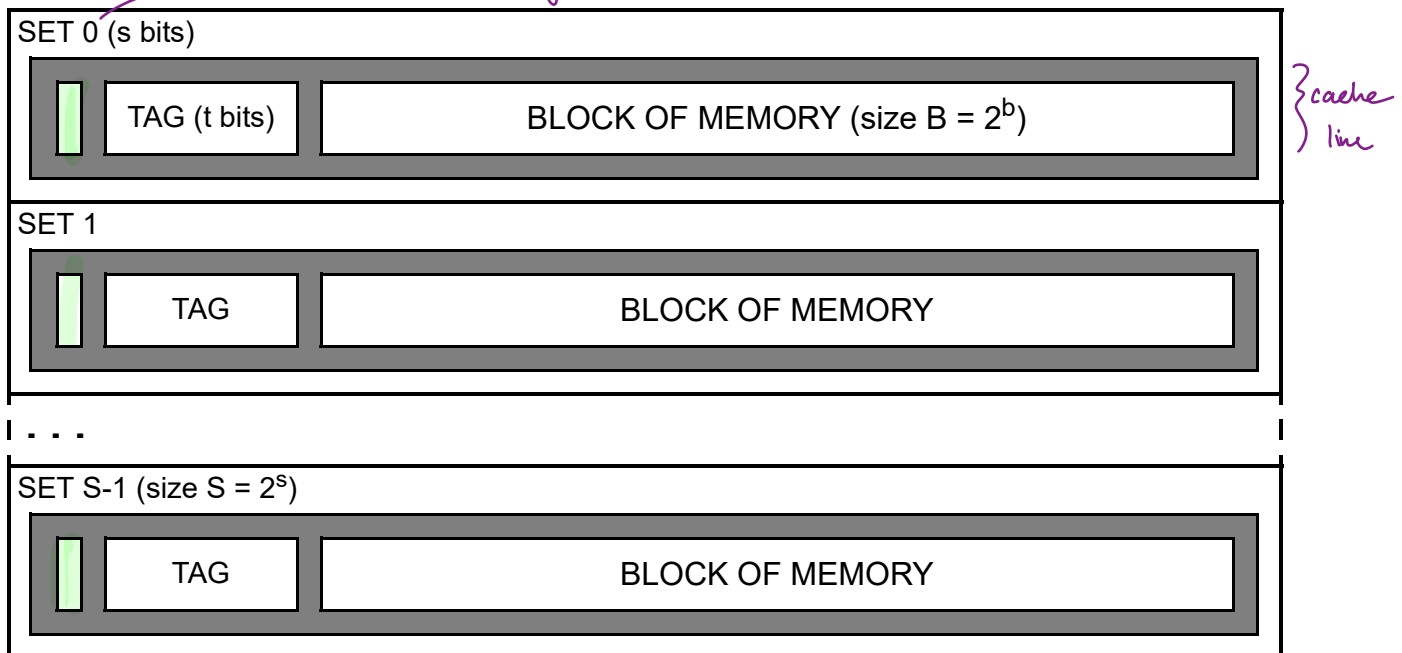
# Basic Cache Lines

**What?** A _line_ is

- locatun in the cache that can store one block of memory
- composed of storage for the block and info needed for cache operation

✽ *In our <u>basic cache</u> each cache set* has only one line

**Basic Cache Diagram** — treat like array index

SET 0 (s bits)

| TAG (t bits) | BLOCK OF MEMORY (size B = $2^b$) |

} cache line

SET 1

| TAG | BLOCK OF MEMORY |

. . .

SET S-1 (size S = $2^s$)

| TAG | BLOCK OF MEMORY |

→ How do you know if a line in the cache is used or not?

use a status bit (v bit)

if v = 1, cache block is copied to cache line

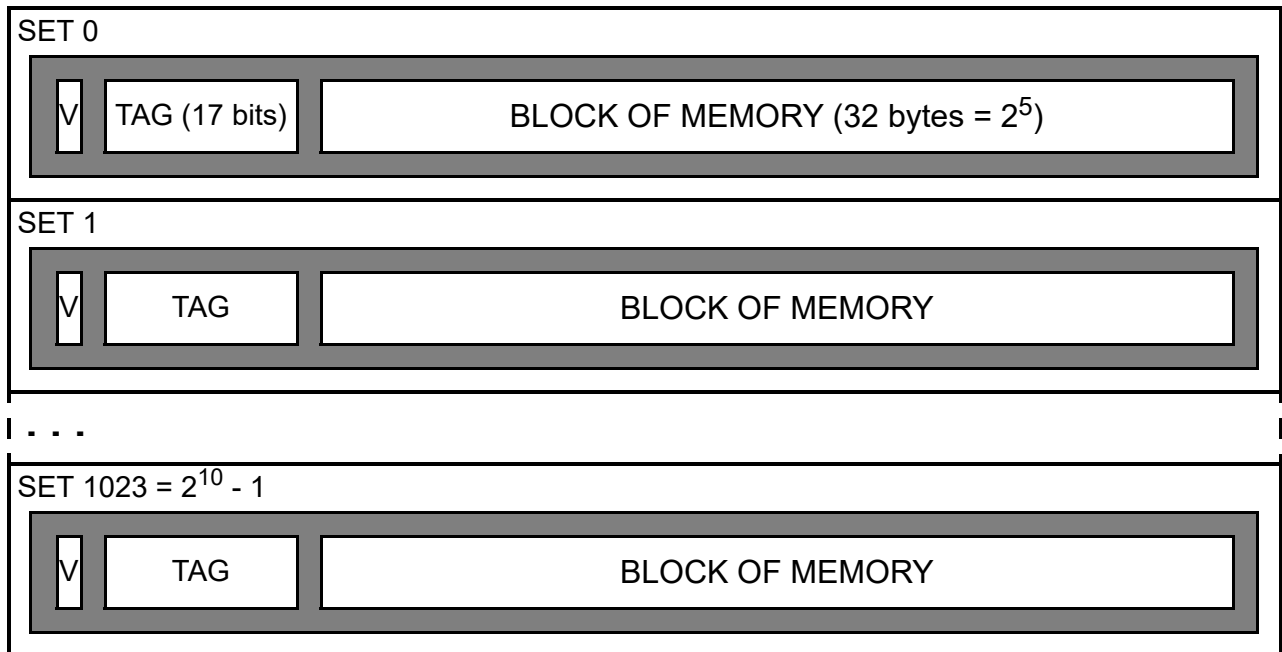→ How big is a basic cache given S sets with blocks having B bytes?

$$C = S \times B$$

\# of bytes   \# of sets   bytes/ block

b/c 1 block/ set

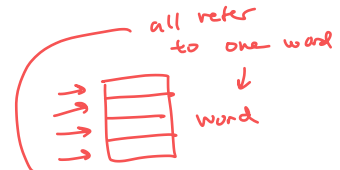tag and v-bit are not included in cache size
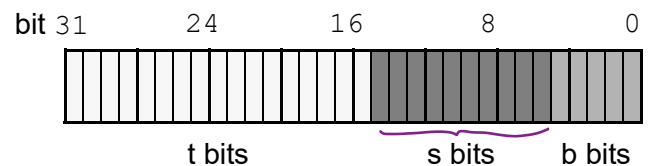
# Basic Cache Operation

**Basic Cache Diagram**

```
┌─────────────────────────────────────────────────────────────────────┐
│ SET 0                                                                 │
│  ┌──────────────────────────────────────────────────────────────┐    │
│  │ V │ TAG (17 bits) │     BLOCK OF MEMORY (32 bytes = 2⁵)       │    │
│  └──────────────────────────────────────────────────────────────┘    │
├─────────────────────────────────────────────────────────────────────┤
│ SET 1                                                                 │
│  ┌──────────────────────────────────────────────────────────────┐    │
│  │ V │    TAG    │          BLOCK OF MEMORY                      │    │
│  └──────────────────────────────────────────────────────────────┘    │
```

. . .

```
┌─────────────────────────────────────────────────────────────────────┐
│ SET 1023 = 2¹⁰ - 1                                                    │
│  ┌──────────────────────────────────────────────────────────────┐    │
│  │ V │    TAG    │          BLOCK OF MEMORY                      │    │
│  └──────────────────────────────────────────────────────────────┘    │
└─────────────────────────────────────────────────────────────────────┘
```

→ How big is this basic cache?    $C = S \times B = 1024 \times 2^5 = 2^{15}$ bytes

$2^{10} \qquad 32 \qquad = 32\ Kb$

*all refer to one word* → *word*

**How does a cache process a request for a <u>word</u> at a <u>particular address</u>?**    *4-bytes*

1. *Set Selection*  identify the set
   extract s-bit
   └→ use as index

2. *Line Matching*  extract t-bits
   compare t-bits w/ stored tag in line of the set

   if no match or valid bit is 0 , cache miss
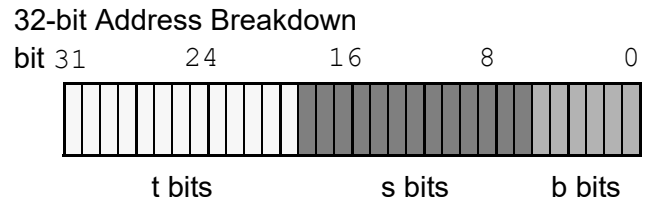
   if match and valid bit is 1 , cache hit

   For L1 cache only , must now extract
   word from block using word offset

**32-bit Address Breakdown**

bit 31        24        16        8        0

|          t bits          |  s bits  | b bits |

on HWs & Exams,
questions will be
"is this addr in the cache?"

# Basic Cache Practice

**You are given the following 32-bit address breakdown used by a cache:**

32-bit Address Breakdown

bit 31          24          16          8          0

t bits          s bits          b bits

→ How big are the blocks?

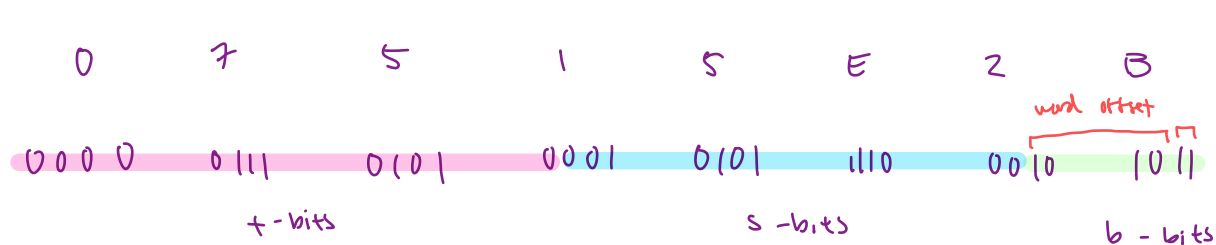$B = 2^b = 2^6 = 64$ bytes / block

→ How many sets?

$S = 2^s = 2^{13} = 8192$ sets / cache

→ How big is this basic cache?

1 block / cache line

$C = S \times B = 2^{13} \times 2^6 = 2^{19} = 2^9 \cdot 2^{10} = 512K$

**Assume the cache design above is given the following specific address: `0x07515E2B`**

memorize
hex → bin

| 0 | 7 | 5 | 1 | 5 | E | 2 | B |
|---|---|---|---|---|---|---|---|

word offset

0000  0111  0101  0001  0101  1110  0010  1011

t - bits          s - bits          b - bits

→ Which set should be checked given the address above?   s - bits

0010  1011  1100  0  = 8 + 16 + 32 + 64 + 256 + 1024

= 1400₁₀  ∴  set index is  1400

→ Which <u>word in the block</u> does the L1 cache access for the address?

1010

10  1011  →  = word 10 of 64 byte block

which word, byte

➤ Which byte in the word does the address specify?

11  → = byte 3 of word 10 of our block

**Assume address above maps to a set with its line having the following V status and tag.**

→ Does the address above produce a hit or miss?

V  tag

1.) 1  0x0750  MISS  - wrong tag

2.) 0  0x0750  "    "

3.) 1  0x00EA  HIT          0  0600  111  0  1010
                            0    0    E      A

4.) 0  0x00EA  MISS - Not valid

1) blocks, bytes, data

2) tag - which block

3) v - valid or not