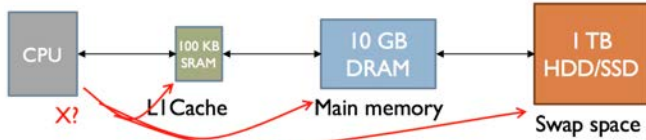


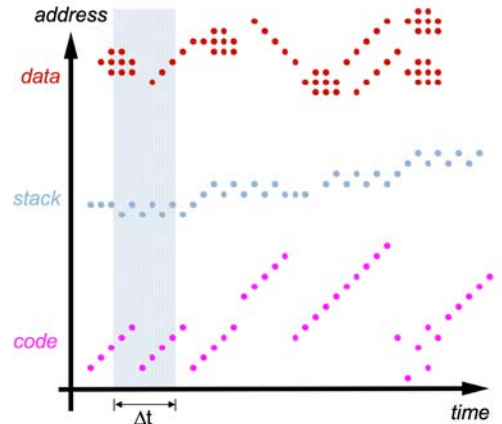
# Computation Structures

## Memory Hierarchy & Caches Worksheet

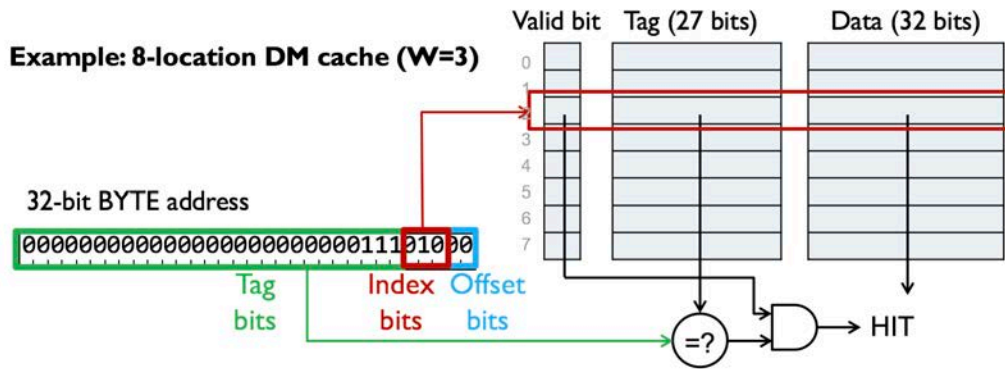


Keep the most often-used data in a small, fast SRAM (often local to CPU chip). The reason this strategy works: LOCALITY.

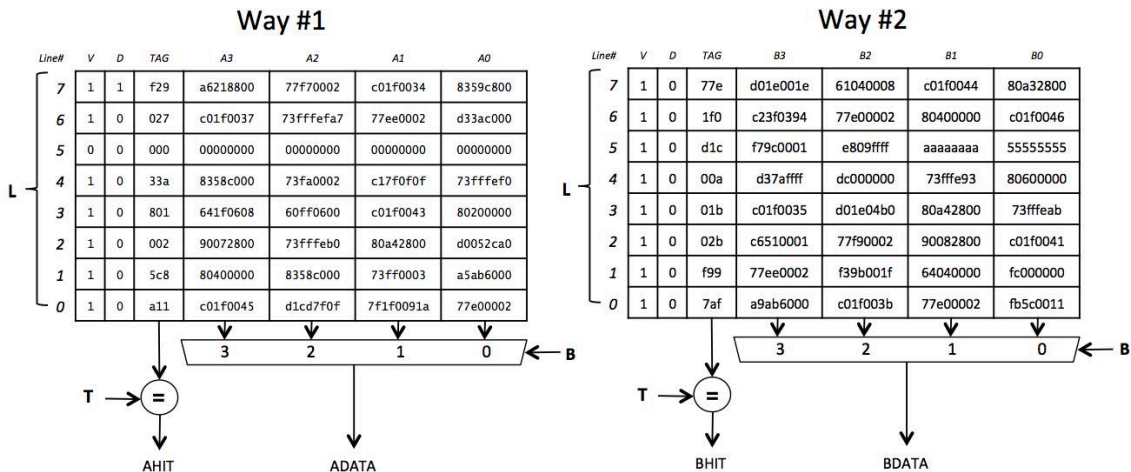
*Locality of reference:* Access to address  $X$  at time  $t$  implies that access to address  $X + \Delta X$  at time  $t + \Delta t$  becomes more probable as  $\Delta X$  and  $\Delta t$  approach zero.



$$AMAT = HitTime + MissRatio * MissPenalty$$



**Example: 2-way set-associative cache, 8 sets, 4-word block size, write-back**



*Replacement strategy choices:* least-recently used (LRU); first in, first out (FIFO); random  
*Write-policy choices:* write-through, write-behind, write-back

**Problem 1.**

- (A) The timing for a particular cache is as follows: checking the cache takes 1 cycle. If there's a hit the data is returned to the CPU at the end of the first cycle. If there's a miss, it takes 10 *additional* cycles to retrieve the word from main memory, store it in the cache, and return it to the CPU. If we want an average memory access time of 1.4 cycles, what is the minimum possible value for the cache's hit ratio?

$$1.4 = 1 + (1 - \alpha)10$$

$$\therefore 0.4 = 1 - \alpha$$

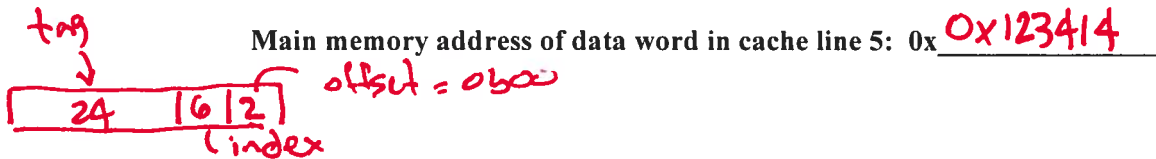
Minimum possible value of hit ratio: 0.96

- (B) If the cache block size, i.e., words/cache line, is doubled but the total number of data words in the cache is unchanged, how will the following cache parameters change? Please circle the best answer.

- # of offset bits: UNCHANGED ... **+1** ... -1 ... 2x ... 0.5x ... CAN'T TELL
- # of tag bits: **UNCHANGED** ... +1 ... -1 ... 2x ... 0.5x ... CAN'T TELL
- # of cache lines: UNCHANGED ... +1 ... -1 ... 2x ... **0.5x** ... CAN'T TELL

Consider a direct-mapped cache with 64 total data words with 1 word/cache line, which uses a LRU replacement strategy and a write-back write strategy. This cache architecture is used for parts (C) through (F).

- (C) If cache line number 5 is valid and its tag field has the value 0x1234, what is the address in main memory of the data word currently residing in cache line 5?



The program shown on the right repeatedly executes an inner loop that sums the 16 elements of an array that is stored starting in location 0x310.

The program is executed for many iterations, then a measurement of the cache statistics is made during one iteration through all the code, i.e., starting with the execution of the instruction labeled `outer_loop`: until just before the next time that instruction is executed.

```

    . = 0
outer_loop:
    CMOVE(16,R0) // initialize loop index J
    CMOVE(0,R1)

loop:
    // add up elements in array
    SUBC(R0,1,R0) // decrement index
    MULC(R0,4,R2) // convert to byte offset
    LD(R2,0x310,R3)// load value from A[J]
    ADD(R3,R1,R1) // add to sum
    BNE(R0,loop) // loop until all words are summed

BR(outer_loop) // perform test again!
    
```

(D) In total, how many instruction fetches occur during one complete iteration of the outer loop?  
How many data reads?

$i\text{fetch} = 2 + (5 \times 16) + 1$       Number of instruction fetches: 83  
 $d\text{read} = 1 \times 16$       Number of data reads: 16

(E) How many instruction fetch misses occur during one complete iteration of the outer loop?  
How many data read misses? Hint: remember that the array starts at address 0x310.

4 locations in array  
 conflict w/ 4 instructions  
 ⇒ other insts. in cache  
 ⇒ other data in cache  
 Number of instruction fetch misses: 4  
 Number of data read misses: 4

(F) What is the hit ratio measured after one complete iteration of the outer loop?

$\text{total accesses} = 83 + 16 = 99$   
 $\text{total hits} = 99 - 8$   
 Hit ratio:  $\frac{99 - 8}{99} = \frac{91}{99}$

**Problem 2.**

The Beta Engineering Team is working on the design of a cache. They've decided that the cache will have a total of  $2^{10} = 1024$  data words, but are still thinking about the other aspects of the cache architecture.

First assume the team chooses to build a direct-mapped write-back cache with a block size of 4 words. ⇒ 4 bit offset

(A) Please answer the following questions:

$\frac{1024 \text{ words}}{4 \text{ wds/line}} = 256 \text{ lines}$       Number of lines in the cache: 256  
 tag | 20 | 8 | 4 | offset      Number of bits in the tag field for each cache entry: 20

(B) This cache takes 2 clock cycles to determine if a memory access is a hit or a miss and, if it's a hit, return data to the Beta. If the access is a miss, the cache takes 20 additional clock cycles to fill the cache line and return the requested word to the Beta. If the hit rate is 90%, what is the Beta's average memory access time in clock cycles?

Average memory access time assuming 90% hit rate (clock cycles): 4

$AMAT = 2 + (1 - .9)(20) = 2 + 2 = 4$

Now assume the team chooses to build a 2-way set-associative write-back cache with a block size of 4 words. *The total number of data words in the entire cache is still 1024.* The cache uses a LRU replacement strategy.  $\Rightarrow$  128 lines in each way (7 bits of index)

(C) Please answer the following questions:

Address bits used as offset (including byte offset): A[ 3 : 0 ]

Address bits used as cache line index: A[ 10 : 4 ]

Address bits used for tag comparison: A[ 31 : 11 ]

(D) To implement the LRU replacement strategy this cache requires some additional state for each set. How many state bits are required for each set?

Number of state bits needed for each set for LRU: 1

To test this set-associative cache, the team runs the benchmark code shown on the right. The code sums the elements of a 16-element array. The first instruction of the code is at location 0x0 and the first element of the array is at location 0x10000. Assume that the cache is empty when execution starts and remember *the cache has a block size of 4 words.*

```

. = 0x0
  CMOVE(0, R0)
  CMOVE(0, R1)
L: LD(R0, A, R2)
  ADD(R2, R1, R1)
  ADDC(R0, 4, R0)
  CMLTCL(R0, 64, R2)
  BT(R2, L)
  HALT()

```

(E) How many instruction misses will occur when running the benchmark?

62 cache lines  
 Number of instruction misses when running the benchmark: 2

(F) How many data misses (i.e., misses caused by the memory access from the LD instruction) will occur when running the benchmark?

Number of data misses when running the benchmark: 4

16-element array  $\Rightarrow$  4 cache lines

```

. = 0x10000
A: LONG(1)
  LONG(2)
...
  LONG(15)
  LONG(16)

```

(g) What's the exact hit rate when the complete benchmark is executed?

Benchmark hit rate:  $\frac{99-6}{99} = \frac{93}{99} \approx 94\%$

$\# \text{ inst. fetches} = 2 + 16 * 5 + 1 = 83$

$\# \text{ data fetches} = 16 \Rightarrow \text{total fetches} = 99$

**Problem 3.**

The program from the Cache performance lab is shown at the right. Assume the program is being run on a Beta with a cache with the following parameters:

- 2-way set-associative
- block size of 2, i.e., 2 data words are stored in each cache line
- total number of data words in the cache is 32
- LRU replacement strategy

(A) The cache will divide the 32-bit address supplied by the Beta into three fields: B bits of block offset (including byte offset bits), L bits of cache line index, and T bits of tag field. Based on the cache parameters given above, what are the appropriate values for B, L, and T?

3 bytes/line value for B: 3  
 8 lines/way value for L: 3  
 32 - L - B value for T: 26

```

I = 0x240 // location of program
A = 0x420 // location of array A
N = 16 // size of array (in words)

. = I // start program here
test:
240 CMOVE(N,R0) // initialize loop index J
244 CMOVE(0,R1)

loop: // add up elements in array
248 SUBC(R0,1,R0) // decrement index
24C MULC(R0,4,R2) // convert to byte offset
LD(R2,A,R3) // load value from A[J]
ADD(R3,R1,R1) // add to sum
BNE(R0,loop) // loop N times

BR(test) // perform test again!
    
```

```

// allocate space to hold array
. = A
STORAGE(N) // N words
    
```

(B) If the MULC instruction is resident in a cache line, what will be its cache line index? the value of the tag field for the cache?

0x24C = 00100100 1100  
 offset: 0b 100  
 line: 0b001  
 tag: 0b1001  
 Cache line index for MULC when resident in cache: 1  
 Tag field for MULC when resident in cache: 0x 9

(C) With the values of I, A, and N as shown, list all the values j (0 ≤ j < N) where the location holding the value A[j] will map to the same cache line index as the MULC instruction in the program.

List all j where A[j] have the same cache line index as MULC: 10, 11  
 A[8] @ 0x420 ⇒ cache line 4. 2 words/line ⇒ A[8], A[9] are in cache line 4.

(D) If the outer loop is run many times, give the steady-state hit ratio for the cache, i.e., assume that the number of compulsory misses as the cache is first filled are insignificant compared to the number of hits and misses during execution.

everything fits in the cache!  
 Steady-state hit ratio (%): 100%



(D) This cache is used to run the following benchmark program. The code starts at memory address 0; the array referenced by the code has its first element at memory address 0x2000. First determine the number of memory accesses (both instruction and data) made during each iteration through the loop. Then estimate the steady-state average hit ratio for the program, i.e., the average hit ratio after many iterations through the loop.

```

. = 0
  CMOVE(0,R0)      // byte index into array
  CMOVE(0,R1)      // initialize checksum accumulator
loop:
  LD(R0,array,R2)  // load next element of array
  SHLC(R1,1,R1)    // shift checksum
  ADDC(1,R1,R1)    // increment checksum
  ADD(R2,R1,R1)    // include data value in checksum
  ADDC(R0,4,R0)    // byte index of next array element
  CMLTTC(R0,1000,R2) // process 250 entries
  BT(R2,loop)
  HALT()

```

```

. = 0x2000
array:
... array contents here ...

```

*each iteration: 7 ifetches, 1 data access*  
*instructions occupy 4 cache lines (2 words/line)*  
*↳ 2-way ⇒ 100% hit on ifetch*  
*data accesses: new word every cycle*  
*↳ 2 words/line ⇒ 50% hit*

Number of memory accesses made during each iteration of the loop: 8

Estimated steady-state average hit ratio:  $\frac{15}{16} = 93.75\%$

*1 miss every two iterations*

**Problem 5.**

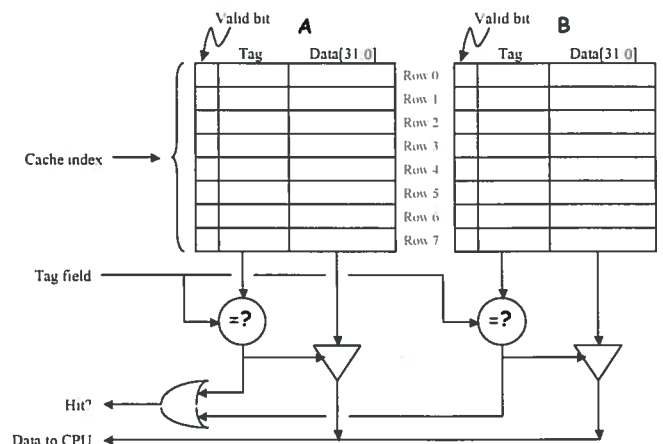
Consider the diagram to the right for a 2-way set associative cache to be used with our Beta design. Each cache line holds a single 32-bit word of data along with its associated tag and valid bit (0 when the cache line is invalid, 1 when the cache line is valid).

(A) The Beta produces 32-bit byte addresses, A[31:0]. To ensure the best cache performance, which address bits should be used for the cache index? For the tag field?

*8 cache lines*  
 ↳ address bits used for cache index: A[4:2]

address bits used for tag field: A[31:5]

*remaining bits are tag.*



(B) Suppose the Beta does a read of location 0x5678. Identify which cache location(s) would be checked to see if that location is in the cache. For each location specify the cache section (A or B) and row number (0 through 7). E.g., 3A for row 3, section A. If there is a cache hit on this access what would be the contents of the tag data for the cache line that holds the data for this location?

0x5678  
 0b101100111100 cache location(s) checked on access to 0x5678: 6A, 6B  
 2 B 3 | (offset) cache tag data on hit for location 0x5678 (hex): 0x 2B3  
 tag ← line index ← 0x6

(C) Assume that checking the cache on each read takes 1 cycle and that refilling the cache on a miss takes an *additional* 8 cycles. If we wanted the *average* access time over many reads to be 1.1 cycles, what is the minimum hit ratio the cache must achieve during that period of time? You needn't simplify your answer.

$$1.1 = 1 + (1 - HR) \cdot (8)$$

minimum hit ratio for 1.1 cycle average access time:  $\frac{7.9}{8}$

(D) Estimate the approximate cache hit ratio for the following program. Assume the cache is empty before execution begins (all the valid bits are 0) and that an LRU replacement strategy is used. Remember the cache is used for both instruction and data (LD) accesses.

LINE	ADDR
0	0
1	4
2	8
3	C
4	10
5	14
6	18
7	1C
0	20
1	24

```

. = 0
CMOVE(source, R0)
CMOVE(0, R1)
CMOVE(0x1000, R2)
loop: LD(R0, 0, R3)
      ADDC(R0, 4, R0)
      ADD(R3, R1, R1)
      SUBC(R2, 1, R2)
      BNE(R2, loop)
      ST(R1, source)
      HALT()
  
```

each iteration

- 5 inst. fetches  $\approx$  100% hit
  - 1 data fetch  $\approx$  0% hit  
 ↳ no data word fetched twice
- 6 total

```

. = 0x100
source:
  
```

```

. = . + 0x4000 // Set source to 0x100, reserve 1000 words
  
```

approximate hit ratio:  $\frac{5}{6}$

(E) After the program of part (D) has finished execution what information is stored in row 4 of the cache? Give the addresses for the two locations that are cached (one in each of the sections) or briefly explain why that information can't be determined.

Addresses whose data is cached in "Row 4": 0x10 and 0x40FB  
 ADDR 40E0 40E4 40E8 40EC 40F0 40F4 40F8 40FC  
 LINE 0 1 2 3 4 5 6 7  
 ^ data

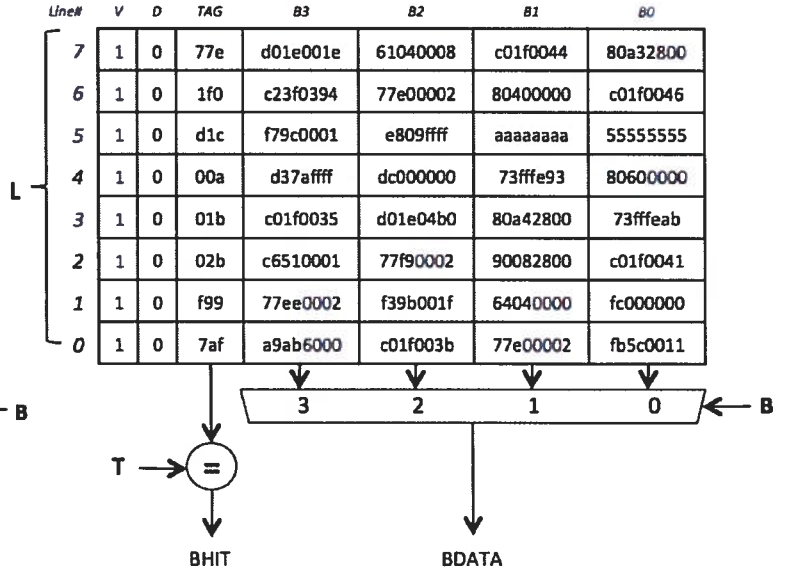
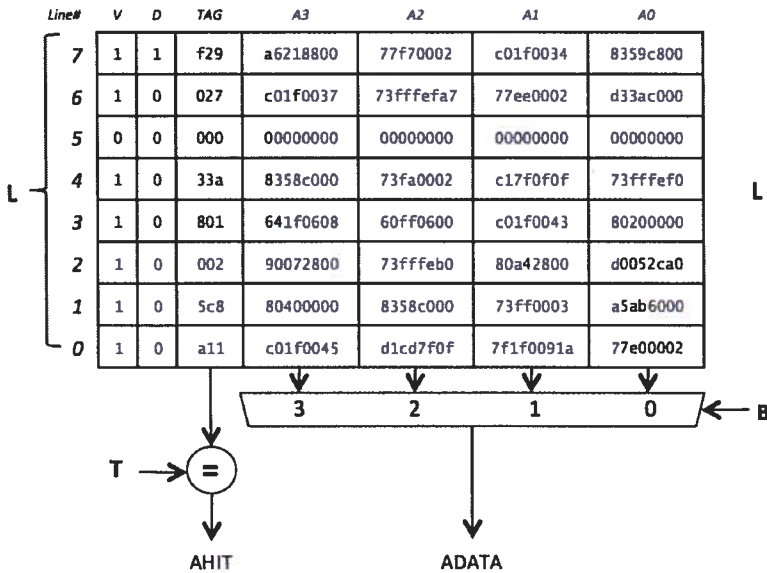


**Problem 6.**

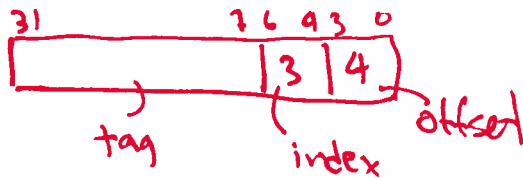
A standard unpipelined Beta is connected to a 2-way set-associative cache containing 8 sets, with a block size of 4 32-bit words. The cache uses a LRU replacement strategy. At a particular point during execution, a snapshot is taken of the cache contents, which are shown below. All values are in hex; assume that any hex digits not shown are 0.

**Way #1**

**Way #2**



- (A) The cache uses bits from the 32-bit byte address produced by the Beta to select the appropriate set (L), as input to the tag comparisons (T) and to select the appropriate word from the data block (B). For correct and optimal performance what are the appropriate portions of the address to use for L, T and B? Express your answer in the form “A[N:M]” for N and M in the range 0 to 31, or write “CAN’T TELL”.



Address bits to use for L: A[6:4]  
 Address bits to use for T: A[31:7]  
 Address bits to use for B: A[3:2]

- (B) For the following addresses, if the contents of the specified location appear in the cache, give the location’s 32-bit contents in hex (determined by using the appropriate value from the cache). If the contents of the specified location are NOT in the cache, write “MISS”.

$B=0, L=0, T=0x1422$  Contents of location 0xA1100 (in hex) or “MISS”: 0x miss

$B=2, L=4, T=0xA$  Contents of location 0x548 (in hex) or “MISS”: 0x D0000000

- (C) Ignoring the current contents of the cache, is it possible for the contents of locations 0x0 and 0x1000 to both be present in the cache simultaneously?

Locations 0x0 and 0x1000 present simultaneously (circle one): YES ... NO

*both map to cache line 0, but it's a 2-way cache*

(D) (Give a one-sentence explanation of how the D bit got set to 1 for Line #7 of Way #1.

**One sentence explanation**

ST changed a value of a word on that cache line but it hasn't yet been written back to memory.

(E) The following code snippet sums the elements of the 32-element integer array X. Assume this code is executing on a Beta with a cache architecture as described above and that, initially, the cache is empty, i.e., all the V bits have been set to 0. Compute the hit ratio as this program runs until it executes the HALT() instruction, a total of  $2 + (6 \cdot 32) + 1 = 195$  instruction fetches and 32 data accesses.

Hit ratio:  $\frac{216}{227} = 95.29\%$

+ 2 {  
 . = 0  
 CMOVE(0, R0) // loop counter  
 CMOVE(0, R1) // accumulated sum

loop:

6 ifetch  
 1 dfetch  
 ↑  
 32 iterations

{  
 SHLC(R0, 2, R2) // convert loop counter to byte offset  
 LD(R2, X, R3) // load next value from array  
 ADD(R3, R1, R1) // add value to sum  
 ADDC(R0, 1, R0) // increment loop counter  
 CMLTTC(R0, 32, R2) // finished with all 32 elements?  
 BT(R2, loop) // nope, keep going

+ 1 HALT() // all done, sum in R1

X: LONG(1) // the 32-element integer array X  
 LONG(2)  
 ...  
 LONG(32)

11 misses total out of 227 accesses

- 2-way cache: 100% hit rate on ifetch
- block size = 4:
  - 3 misses to load instructions
  - 32 data accesses, but only 25% miss (each miss loads 4 words)
  - ⇒ 8 data misses total

**Problem 7.**

After his geek hit single *I Hit the Line*, renegade singer Johnny Cache has decided he'd better actually learn how a cache works. He bought three Beta processors, identical except for their cache architectures:

- **Beta1** has a 64-line direct-mapped cache
- **Beta2** has a 2-way set associative cache, LRU, with a total of 64 lines
- **Beta3** has a 4-way set associative cache, LRU, with a total of 64 lines

Note that each cache has the same total capacity: 64 lines, each holding a single 32-bit **word** of data or instruction. All three machines use the same cache for data and instructions fetched from main memory.

Johnny has written a simple test program:

```
// Try a little cache benchmark
I = 0x1000 // where program lives
A = 0x2000 // data region 1
B = 0x3000 // data region 2
N = 16 // size of data regions (BYTES!)

. = I // start program here
P: CMOVE (1000, R6) // outer loop count
Q: CMOVE (N, R0) // loop index I (array offset)
R: SUBC (R0, 4, R0) // I = I-1
LD (R0, A, R1) // read A[I]
LD (R0, B, R2) // read B[I]
BNE (R0, R)
SUBC (R6, 1, R6) // repeat many times
BNE (R6, Q)
HALT ()
```

Johnny runs his program on each Beta, and finds that one Beta model outperforms the other two.

(A) (2 points) Which Beta gets the highest hit ratio on the above benchmark?

3 regions: 0x1000+, 0x2000+, 0x3000+ Circle one: Beta1 Beta2 **Beta3**

(B) (2 points) Johnny changes the value of **B** in his program to **0x2000** (same as **A**), and finds a substantial improvement in the hit rate attained by one of the Beta models (approaching 100%). Which model shows this marked improvement?

now just 2 regions: 0x1000+, 0x2000+ Circle one: Beta1 **Beta2** Beta3

(C) (3 points) Finally, Johnny sets **I**, **A**, and **B** each to **0x0**, and sets **N** to **64**. What is the TOTAL number of cache misses that will occur executing this version of the program on each of the Beta models?

TOTAL cache misses running on Beta1: 16; Beta2: 16; Beta3: 16

only compulsory misses for all these caches

MIT OpenCourseWare  
<https://ocw.mit.edu/>

6.004 Computation Structures  
Spring 2017

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.