

1.  $2^{12}$  bytes =  $\boxed{4\text{KiB}}$
2.  $2^{12}$  bytes / 2 bytes/word =  $\boxed{2^{11}$  words
3.  $2^{11}$  words  $\Rightarrow$   $\boxed{2^{11}$  blocks
4.  $\boxed{1 \text{ bit}}$  (since 2 bytes in word)
5.  $\boxed{11 \text{ bits}}$  (since  $2^{11}$  blocks)
6.  $\boxed{12 \text{ bits}}$  ( $24 - (1 + 11)$ )
7.  $\boxed{14 \text{ bits}}$  ( $12 + 1 + 1$ )  $\rightarrow$
7.  $\boxed{7 \times 2^9 \text{ bytes}}$   $7 \times 2^{12} \text{ bits} = 7 \times 2^9 \text{ bytes of metadata in cache}$
8.  $(7 \times 2^9) + 2^{12} \text{ bytes} = \boxed{7.5\text{KiB}}$   $\leftarrow$  bigger than our 6KiB budget.
9.  $\boxed{2^9 \text{ blocks}}$  ( $2^{11}/4$ )
10.  ~~$\boxed{3 \text{ bits}}$  (since  $8 = 2^3$  bytes in 4 words)~~
10.  ~~$\boxed{1 \text{ bit}}$  (since 2 bytes per word)~~
11.  ~~$\boxed{2 \text{ bits}}$  (since  $2^2$  bytes per word)~~
12.  ~~$\boxed{9 \text{ bits}}$  (since  $2^9$  blocks in cache)~~
13.  ~~$\boxed{12 \text{ bits}}$  ( $24 - (9 + 2 + 1)$ )~~
14.  ~~$\boxed{14 \text{ bits}}$  (same as previously)~~
14.  $\Rightarrow 7 \times 2^{10} \text{ bits} = 7 \times 2^7 \text{ bytes of metadata in cache}$
15.  $7 \times 2^7 + 2^{12} \text{ bytes} = 4.875 \text{ KiB}$
16. Because a cache block is  $4 \times 16$  bit words, a 16 bit bus means that 4 values will be loaded from memory on a cache read miss.
17.  $6 - 4.875 = 1.125 \text{ KiB} = 1152 \text{ bytes left in our budget for write-through buffer.}$   
 Buffer entry size = 20 bits (metadata & data) + 24 bits (address) = 44 bits (5.5 bytes)  
 $(1152 \times 8) / 44 \approx 209$  buffer entries.
18. The additional complexity of write-back would involve checking the dirty bit when writing to a word in the cache. If the bit was set, then the whole block (4 words) must be written to memory (over a restricted 16-bit bus). However write-back will put an updated value in the buffer on each write of one word only. Note that the write-back could be improved by widening the dirty metadata to 2-bits, and therefore selecting dirty words out of larger cache blocks.
19. You would likely use write-through, so that the device can get the most up to date data in the mapped region. (eg. a framebuffer, you want the screen to always reflect current data.

20: With some small supporting features in the hardware, like the addition of a bit of metadata to each cache block specifying that the address should use either write-through or write-back, it should be possible to specify particular memory regions that use different cache writing schemes.

21: There are two separate scenarios to worry about on cache reads, which depend on the nature of the device using the mmaped region: if the device writes to the mmaped region, or if the CPU is the only source of change to the region.

If the device is writing, then caching might be detrimental. It would mean that the CPU may never invalidate a cache block and need to go to memory, meaning it is isolated from the changes occurring in the memory region. In this case, one might not want to cache that address range at all. This might be an input device, like a keyboard. ~~For~~ For this, we'd treat all reads in the region as cache misses, and load from memory. However, if only the CPU is writing, then a cache miss can be ~~loaded~~ loaded into cache, and then treated as a normal cache block. Subsequent reads, while the block remains valid, can be treated as cache hits. This would be used in cases of output-only devices, like a screen using an mmaped frame buffer.