

1. Cache principles
 - a. Cache holds a much smaller subset of RAM's contents
 - i. Bring in values from RAM while executing a program
 - ii. If we need a value, look in cache first before going to memory
 - iii. Benefit from having cache because it is faster than RAM, and closer to the CPU
 - iv. However, cache is smaller than RAM, and costs more per bit than RAM
 - b. Terminology
 - i. Memory *block* – unit of main memory stored in a cache line
 - ii. Cache *line* – basic unit of a cache
 1. Contains a block of memory
 2. Also contains a tag, and control bits to determine when line should be replaced
 - iii. Cache *tag* – number stored with a line
 1. Along with the line's position in the cache, determines the address of the block from main memory
 - iv. Cache *hit* – data that the CPU asks for is in the cache
 - v. Cache *miss* – data that the CPU asks for is *not* in the cache
 - vi. *Eviction* – removing memory from the cache and replacing with new memory
 - vii. *Dirty* – line contains updated data that differs from the corresponding main memory
 1. Before evicting the line from cache, must copy back to main memory
 2. Otherwise, would lose changes to that data
2. Cache addressing
 - a. What type of addresses does the cache use?
 - b. Cache addresses can be *physical*, using actual memory addresses
 - i. Slower, as the cache must wait for the MMU
 - ii. Memory management unit (MMU) converts virtual to physical addresses
 - iii. However, physically addressed caches don't have to deal with the aliasing problem below
 - c. Can also be *logical*, using the same virtual addresses that the process uses
 - i. Faster, as the cache doesn't need to wait for MMU
 1. Can reuse the same address that the process was using
 - ii. However, these caches encounter the *aliasing* problem
 1. Same virtual address can be used in multiple user processes
 2. Cache would say there's a hit
 3. However, that data doesn't correspond to process currently running
 4. Need to flush cache every time we switch processes
 5. Incur a performance penalty every time we flush cache
 - a. Lose all our built-up cache data
 - b. Must go back to RAM to fetch new values

- 3. Cache design
 - a. Multiple-level caches
 - i. See these in many modern CPUs
 - ii. L1 cache is the smallest and closest to the CPU
 - iii. L2 cache is larger than the L1, but is further away from the CPU
 - 1. Takes more time to access the L2 than the L1
 - iv. Same applies for L3
 - 1. Even larger than L2, but even slower
 - b. Unified or split
 - i. Split - have separate caches for data or instructions
 - ii. Unified – combine the two
 - c. Cache size
 - i. Larger cache means slower cache
 - ii. There is no optimum cache size
 - 1. Exact specifications will always depend on machine being run