

①

Name: Awaiz Ghaffar

Dep: BScs

Semester 15269 4th

ID ✓: 4

Assig. No. Computn Architecture

Subject

Submitted to: Amin Sir.

(i)

What is the general relationship among access time, memory cost and capacity?

As access time becomes faster the cost per bit increases, As memory size increases the cost per bit is smaller, Also with greater capacity, the access time becomes slower.

ii) Discuss different memory access methods in detail?

Another distinction among memory types is the method of accessing units of data. These include the following.

(2)

* Sequential Access:

Memory is organized into units of data called records. Access must be made in a specific linear sequence. Stored addressing information is used to separate records and assist in the retrieval process. A shared read-write mechanism is used and this must be moved from its current location to the desired location, passing and refetching each intermediate record.

Thus, the time to access an arbitrary record is highly variable.

Direct Access.

As with sequential access, direct access involves a shared read-write mechanism. However, individual blocks or records have a unique address based on physical location. Access is accomplished by direct access to reach a general vicinity plus sequential searching, counting, or waiting to reach the final

(3)

the final location. Again
access time is variable.

Random access:

Each addressable location in memory has a unique physically wired address in addressing mechanism. The time to access a given location is independent of the sequence of prior accesses and is constant. Thus, any location can be selected at random and directly addressed and accessed main memory. And some cache system are random access.

Associative:

This is random access type of memory that enables one to make a comparison of desired bit location with in a word of specified match, and to do this for all words simultaneously.

Thus a word for a specified portion of its contents rather than its address.

(4)

As with ordinary random-access memory each location has its own addressing mechanism and retrieval time is constant independent of location or prior access patterns. Cache memories may apply associative access.

iii) Discuss the importance of memory hierarchy.
iv) How does the principle of locality relate to the use of multiple memory levels?

Ans) slower and less expensive memory is used in higher stages with the most expensive being the registers in the processor as well as a cache. Main memory is slower and less expensive and is outside of the processor.

v) How many main memory address is interpreted in direct, associative and set associative?

Mapping?

Direct mapping?

The simplest technique.

maps each block of main memory

(5)

into only one possible ^{cache} line.

Associative line mapping:

Permits each main memory block to be loaded into any line of cache.

The cache control logic interprets a memory address simply as a tag and a word field.

To determine whether a block is in the cache the cache control logic must simultaneously examine every line's tag for a match.

Set Associative mapping:

A compromise that exhibits the strengths of both direct and associative approaches while reducing their disadvantages.

Q: NO. 2

Note each of following)

i) Memory unit transfer?

It is the maximum number of bits that can be read or written into the memory.

(6)

at a time. In case of main memory, it is mostly equal to word size. In case of external memory, unit of transfer is not limited to the word size. It is often larger and is referred to as blocks.

Memory of Performance Parameter:
The two most important characteristics of memory are capacity and performance. Three performance parameters are used.

Access time:

For random access memory, this is the time it takes to perform a read or write operation. That is the time from the instant that an address is presented to the memory to the instant that data have been stored or made available for use. For non-random access memory, access time is the time it takes to position the read/write mechanism and the desired location.

(7)

Memory cycle time:

This concept is primarily applied to random access memory and consist of the access time plus any additional time required before a second access can commence.

This additional time may be required for transients to die on signal times or to regenerate data if they are read destructively.

Note that memory cycle is concerned with the system bus not the processor.

iii) Disk Cache.

A portion of main memory can be used as a buffer to hold data temporarily that is to be read out to disk.

A few large transfer of data can be used instead of many

small transfer of data. Data can be retrieved rapidly from the software cache rather than slowly from the disk.

(8)

v) Principle of locality :-

The principle of locality states that data in the vicinity of a referenced word are likely to be referenced in the near future.

vi) Logical cache and Physical cache.

A logical cache, also known as a virtual cache, stores data using addresses. The processor accesses the cache directly without going through the MMU. A physical cache stores data using main memory physical addresses.

One obvious advantage of the logical cache is that cache access speed is faster than for a physical cache. A disadvantage has to do with the fact that most virtual memory systems supply each application with the same virtual memory address space.

(9)

Replacement Algorithms:

Once the cache has been filled when a new block is brought into the cache. One of the existing blocks must be replaced.

For direct mapping there is only one possible line for any particular block and no choice is possible. For the associative and set associative techniques a replacement algorithm is needed. To achieve high speed such an algorithm must be implemented in hardware.

vii) Possible approaches to cache coherence.

Possible approaches to cache coherence include the following:
Bus watching with write through

Each cache controller monitors the address lines to detect write operation to memory by other bus master.

(10)

if another master writes to a location in shared memory that also resides in the cache memory. The cache controller invalidates that cache entry. This strategy depends on the use of a write through Policy by all cache controllers.

Hardware Transparency :

Additional hardware is used to ensure that all updates to main memory via cache are reflected in all caches. Thus, if one processor modifies a word in its cache this update is written to main memory.

Non-Cacheable memory :

only a portion of main memory is shared by more than one processor and this is designated as non-cacheable. In such a system all access to shared memory are cache misses because the shared

(11)

The shared memory is never copied into the cache. The non-cacheable memory can be identified using chip select logic or high address bits.

Q: NO: 3

Different each of the following:

i) Sequential.

Sequential access memory is organized into units of data called records. Access must be made in a specific linear sequence. Stored addressing information is used to separate records and assist in the retrieval process. A shared read-write mechanism is used and this must be moved from its current location to the desired location. Passing and reflecting, searching, counting or waiting to reach the final location. Again time is variable.

ii) Direct :

As with sequential access, direct access involve a shared read-write mechanism. However individual blocks or records have an unique address based on physical location. Access is accomplished by direct access to reach the final location. Again access time is variable.

Random Access ?

Each addressable location in memory has a unique physically wired in addressing mechanism.

The time to access a given is independent of the sequence of prior accesses and is constant.

Thus any location can be selected at random and directly addressed and accessed in main memory and some cache system or random.

ii) Direct associative and set associative mapping.

Direct

The direct mapping technique is simple and inexpensive to implement. Its main disadvantage is that there is a fixed cache location for any given block. Thus if a program happens to reference words repeatedly from two different blocks that map into the same line then the blocks will be continually swapped in the cache and the hit ratio will be low.

Associative

With associative mapping there is flexibility as to which block to replace when a new block read into the cache. Replacement algorithms discussed later in this section are designed to maximize the hit ratio.

The principal disadvantage of associative mapping is the complex circuitry required to tags of

(14)

of all cache lines is parallel.
Set associative

Set associative mapping is a compromise that exhibits the strengths of both the direct and associative approaches while reducing their disadvantages. In this case the cache

consist of a number sets each of which consists of a number of lines. The relationship

$$C = m \times l$$

where

i = Cache set number

j = main memory block number

m = number of lines in the cache

l = number of lines in each set

iii) Split Cache and Unified Cache.

Has become common to split cache

One dedicated to instruction

one dedicated to data

Both exist at the same level.

typically as two L1 caches.

Advantages of unified cache.

- Higher hit rate.
 - Balance load of instruction and data fetches automatically.
 - Only one cache needs to be designed and implemented.
- Trend is toward split caches at the L1 and unified caches for higher levels.

Advantages of split cache
Eliminates cache contention b/w instruction fetch/decode unit and execution unit

important in pipelining

iv) write through write back.

Write Through

Simplest technique.

and write operation are made to main memory as well as to the cache memory.

The main disadvantage of this technique is that it generates substantial memory traffic and may create a bottleneck.

(6)

- write back
- minimize memory writes.
- updates are made. Only in cache.
- Portion of main memory are invalid and hence accesses by I/O modules can be allowed only through the cache.

(Q: NO 4.)

Direct solve.

i) In example, suppose 95% of memory accesses are found in level 1 cache. The average time to access a word can be expressed as

$$(0.95)(0.01 \text{ ms}) + (0.05)(0.01 \text{ ms} + 0.1 \text{ ms})$$
$$= 0.0095 + 0.0055 = 0.015 \text{ ms}$$

The average access time is much closer to 0.01 ms than to 0.1 ms as desired.

(ii)

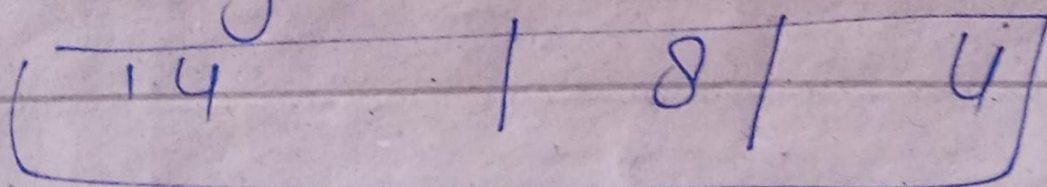
There are total of 8 k bytes / 16 bytes = 512 lines in the cache. Thus the cache consist of 256 sets of 2 lines each. Therefore 8 bits are needed to identify

(7)

set number. For the 64-Mbyte main memory a 26 bit address is needed. main memory consists of 64-M byte / 16 bytes = 2²² blocks. Therefore, the set plus tag length must be 22 bits. So the

tag length is 14 bits and the word field length is 4 bits main memory address.

TAG set word



END