

Batch: B2 Roll No.: 110

Experiment / assignment / tutorial

No. 8

Grade: AA / AB / BB / BC / CC / CD / DD

Signature of the Staff In-charge with date

TITLE : Implementation of Cache Mapping Techniques.

AIM: To study and implement concept of various mapping techniques designed for cache memory.

Expected OUTCOME of Experiment: (Mention CO/CO's attained here)

CO3- Learn and evaluate memory organization and cache structure

Books/ Journals/ Websites referred:

1. Carl Hamacher, Zvonko Vranesic and Safwat Zaky, "Computer Organization", Fifth Edition, TataMcGraw-Hill.
2. Dr. M. Usha, T. S. Srikanth, "Computer System Architecture and Organization", First Edition, Wiley-India.

Pre Lab/ Prior Concepts:

Cache memory: The cache is a smaller, faster memory which stores copies of the data from the most frequently used main memory locations. As long as most memory

accesses are cached memory locations, the average latency of memory accesses will be closer to the cache latency than to the latency of main memory.

2. Hit Ratio: You want to increase as much as possible the likelihood of the cache containing the memory addresses that the processor wants.

$$\text{Hit Ratio} = \frac{\text{No. of hits}}{\text{No. of hits} + \text{No. of misses}}$$

There are only fewer cache lines than the main memory blocks, an algorithm is needed for mapping main memory blocks into cache lines. Further a means is needed for determining which main memory block currently occupies in a cache line. The choice of cache function dictates how the cache is organized. Three techniques can be used.

1. Direct mapping.
2. Associative mapping.
3. Set Associative mapping.

Direct Mapped Cache: The direct mapped cache is the simplest form of cache and the easiest to check for a hit. Since there is only one possible place that any memory location can be cached, there is nothing to search; the line either contains the memory information we are looking for, or it doesn't.

Unfortunately, the direct mapped cache also has the worst performance, because again there is only one place that any address can be stored. Let's look again at our 512 KB level 2 cache and 64 MB of system memory. As you recall this cache has 16,384 lines (assuming 32-byte cache lines) and so each one is shared by 4,096 memory addresses. In the absolute worst case, imagine that the processor needs 2 different addresses (call them X and Y) that both map to the same cache line, in alternating sequence (X, Y, X, Y). This could happen in a small loop if you were unlucky. The processor will load X from memory and store it in cache. Then it will look in the cache for Y, but Y uses the same cache line as X, so it won't be there. So Y is loaded from memory, and stored in the cache for future use. But then the processor requests X, and looks in the cache only to find Y. This conflict repeats over and over. The net result is that the hit ratio here is

0%. This is a worst case scenario, but in general the performance is worst for this type of mapping.

Fully Associative Cache: The fully associative cache has the best hit ratio because any line in the cache can hold any address that needs to be cached. This means the problem seen in the direct mapped cache disappears, because there is no dedicated single line that an address must use. However (you knew it was coming), this cache suffers from problems involving searching the cache. If a given address can be stored in any of 16,384 lines, how do you know where it is? Even with specialized hardware to do the searching, a performance penalty is incurred. And this penalty occurs for all accesses to memory, whether a cache hit occurs or not, because it is part of searching the cache to determine a hit. In addition, more logic must be added to determine which of the various lines to use when a new entry must be added (usually some form of a "least recently used" algorithm is employed to decide which cache line to use next). All this overhead adds cost, complexity and execution time.

Set Associative Cache (To be filled in by students)

This form of mapping is an enhanced form of direct mapping where the drawbacks of direct mapping are removed. Set associative addresses the problem of possible thrashing in the direct mapping method. It does this by saying that instead of having exactly one line that a block can map to in the cache, we will group a few lines together creating a set. Then a block in memory can map to any one of the lines of a specific set. Set-associative mapping allows that each word that is present in the cache can have two or more words in the main memory for the same index address. Set associative cache mapping combines the best of direct and associative cache mapping techniques.

Direct Mapping Implementation:

The mapping is expressed as

$$i = j \text{ modulo } m$$

i = cache line number

j = main memory block number

m = number of lines in the cache

- Address length = (s+w) bits
- Number of addressable units = 2^{s+w} words or bytes
- Block size = line size = 2^w words or bytes
- Number of blocks in main memory = $2^{s+w} / 2^w = 2^s$
- Number of lines in cache = $m = 2^r$
- Size of tag = (s-r) tags

Associative Mapping Implementation: (To be filled in by students)

Size of Tag = $\log_2(\text{total memory} / \text{word size})$

Address length = Tag + word size

Set Associative Mapping Implementation:

$$m = v * k$$

$$i = j \text{ mod } v$$

where

i = cache set number

j=main memory block number

v=number of sets

m=number of lines in the cache number of sets

k=number of lines in each set

ref - <https://www.geeksforgeeks.org/cache-memory-in-computer-organization/>

No of sets = length of memory / k

Implementation code: (direct mapping)

```
Memorysize=int(input('please enter Memory size in k words '))*1024

cachesize=int(input('please enter cache size in k words '))*1024

import math

tag=math.log2(Memorysize/cachesize)

block=math.log2(Memorysize)

word=Memorysize/cachesize /block

print("tag ",tag)

print("block ",block)

print("word ",word)
```

Post Lab Descriptive Questions

1. For a direct mapped cache, a main memory is viewed as consisting of 3 fields. List and define 3 fields.

Tag; holds the tag address of the memory location.

Block: holds the block number:

Word: holds the word number in the block

The memory is divided into tags, then tags are divided into blocks which are further divided into words. The block length is the length of the main memory.

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

Access time is increased, cost reduces, capacity increases.

Conclusion:

Thus we have understood the cache mapping techniques. We understood the set associative mapping, associative mapping and the direct mapping technique. All of these techniques have their own drawbacks and advantages. Set associative mapping sits midway between direct and associative mapping.

Date: 24 nov 22

Signature of faculty in-charge