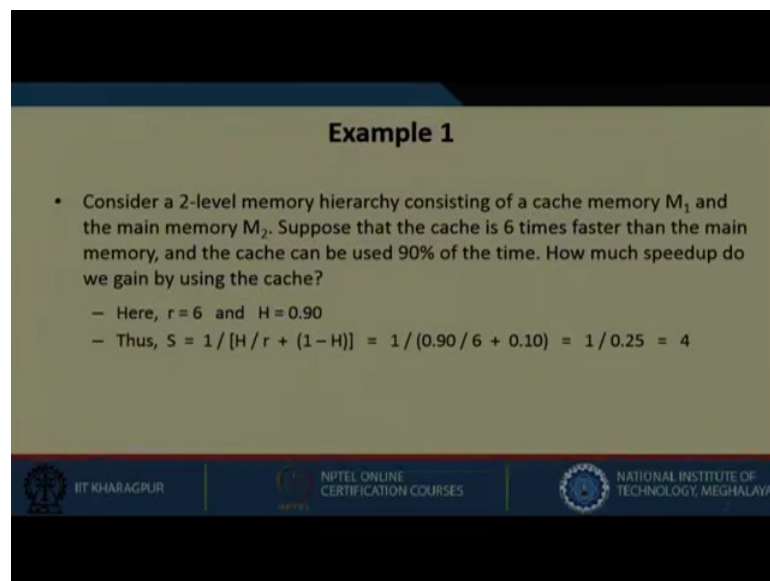**Computer Architecture and Organization**
**Prof. Kamalika Datta**
**Department of Computer Science and Engineering**
**National Institute of Technology, Meghalaya**

**Lecture - 29**
**Memory Hierarchy Design (PART II)**

Welcome to next lecture 29. Here we will be discussing about memory hierarchy design in detail.

(Refer Slide Time: 00:28)



Let us take this example here. Consider a two level memory hierarchy consisting of a cache memory M1 and main memory M2. Generally, the level which is closest to the processor we call it cache memory and then the next level of memory can be different levels of cache memory; or we can also have a main memory.

Suppose that cache is 6 times faster than main memory and the cache can be used 90% of the time. How much speed up do we gain by using cache?

So, here r will be 6 and H will be 0.90. So, we simply put in into this formula and we get 4. So, we are getting a speed up of 4.
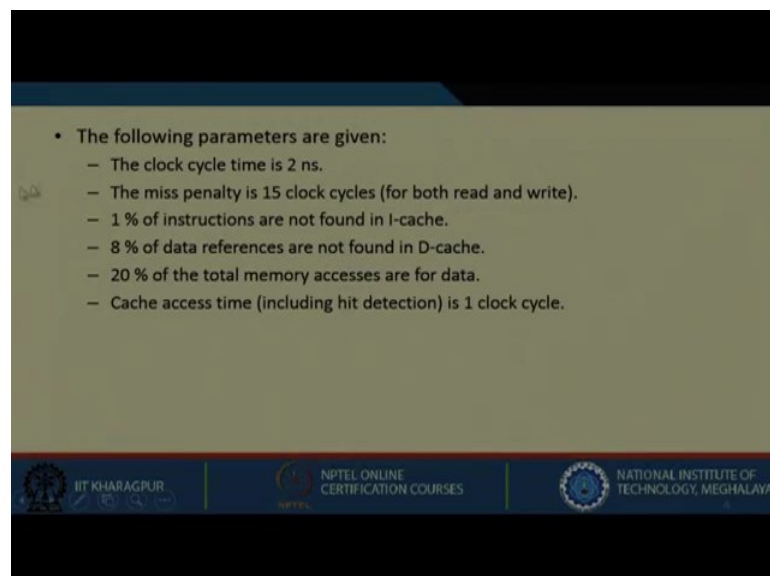
(Refer Slide Time: 02:25)



This is pretty similar to what we discussed way back in Amadahl's law. Now I take another example where we consider two-level memory hierarchy with separate instruction and data caches in level 1, and main memory in level 2. So, here the scenario is we have I-cache, and we have D-cache. The instructions are brought and kept in instruction cache the data are brought and kept in data cache.

The advantage would be if you want to access data and instructions simultaneously, you can do that.
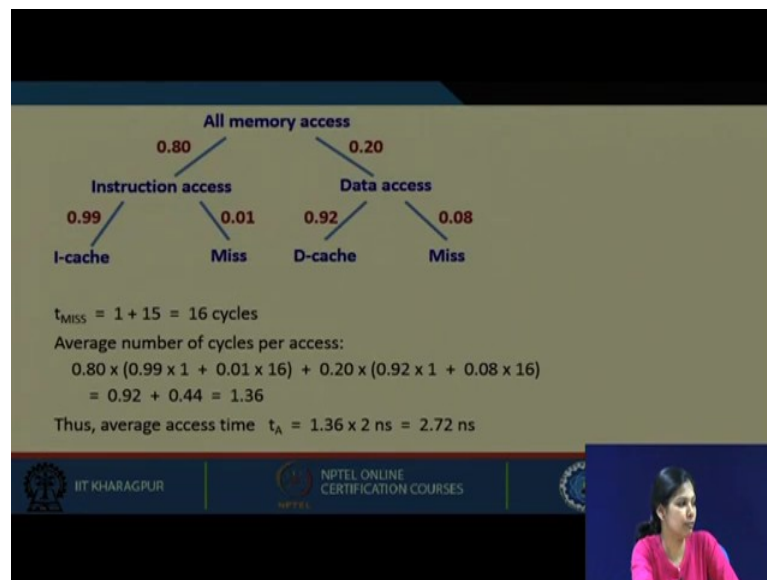
(Refer Slide Time: 03:13)

The following parameters are given.

For miss penalty total 15 clock cycles are required. Out of which it says that 1% of instructions are not found in I-cache; and 8% of data references are not found in D-cache. 20% of the total memory access are for data, and cache access time including hit detection is 1 clock cycle.

So, these are the information that are provided and we will be seeing how we can find the average access time.
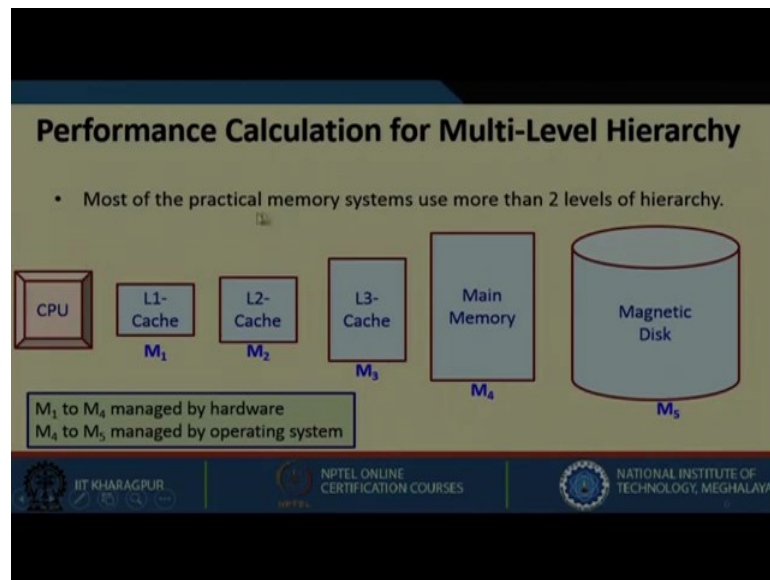
(Refer Slide Time: 04:32)



So, let us see, out of total memory access 80% are for instruction and 20% are for data. Out of this 80% of instruction, I-cache can be used 99%, and miss rate is 0.01 where the data is not found in I-cache.

Similarly, for data access 8% there is miss. So, 0.92 will be hit.

The average access time is calculated as shown.

So, in this example we have seen what happens when there is a miss, what happens when data is found in cache, and now the cache is divided into two different levels, one is I-cache and another is D-cache.

So, when you have two different caches, the calculation will be made in this fashion. Let us see performance calculation for multi level hierarchy.

In general, we have multi-level hierarchy rather than two-level hierarchy. In two-level hierarchy, we have only two levels M1 and M2, but in reality we have multiple levels. For multiple levels whatever we have done for two levels can be extended to multi levels.

So, if you see this particular diagram which is considered as the most practical memory system that; obviously, we use more than two levels of hierarchy. The first level is M1 which is L1 cache, the next level is M2 which is using L2 cache, the next level is M3 which is L3 cache, then M4 is the main memory, and finally M5 is the magnetic disk.

So, how do we calculate? As we can see the access time of this will be much faster than this, then this, then this, and data from main memory is brought through L3 to L2 and then to L1, and when it is in L1 because of locality of reference, we will get the data whenever CPU is asking for it. But initially remember there will be a miss because there will be no data in the cache.

Whenever the data is loaded in main memory, whenever it is required it has to be first brought from main memory through different levels to L1 cache, and finally it goes to CPU. So, for the first time there will be miss always, but for consecutive times there will

be hit because of locality of reference. So, M1 to M4 is managed by hardware; this is very important. But M4 to M5 is handled by operating system.

(Refer Slide Time: 10:31)



Now, let us see this we have one more level; three levels of memory hierarchy, where tL1 is the access time of M1, tL2 is access time of M2, HL1 is the hit ratio of M1; that means, the percentage time the data or instruction is found in L1 cache. HL2 is the percentage time the data or instruction is found in L2 cache, which is not found in L1 cache.

So, hit ratio of M2 is with respect to the residual access that try to access M2. Consider a three level hierarchy that consists of L1 cache, L2 cache and main memory. Whenever there is a miss in L1 we go to L2. Obviously, if we do not find the data or instruction in L1 cache we move to L2 cache. What will be the average access time now? Earlier we have done average access time calculation taking into consideration that CPU first access L1 cache.

The access time tA is given by this expression.

(Refer Slide Time: 13:03)



So, what will be the average access time? The expression is shown.

This is the equation for three levels of memory hierarchy.

(Refer Slide Time: 16:14)



So, what is the implication of the memory hierarchy to the CPU? We can say that the processors designed without memory hierarchy are simple because all memory access takes same amount of time. Misses in a memory hierarchy implies variable memory access time as seen by the CPU, because access time of L1 cache is different from access time of L2 cache, and it is different from main memory.

Some mechanism is required to determine whether or not the requested information is present in top level of memory hierarchy. So, how can we check this? Check happens on every memory access and affects the hit time and implemented in hardware to provide acceptable performance. So, what we are doing.
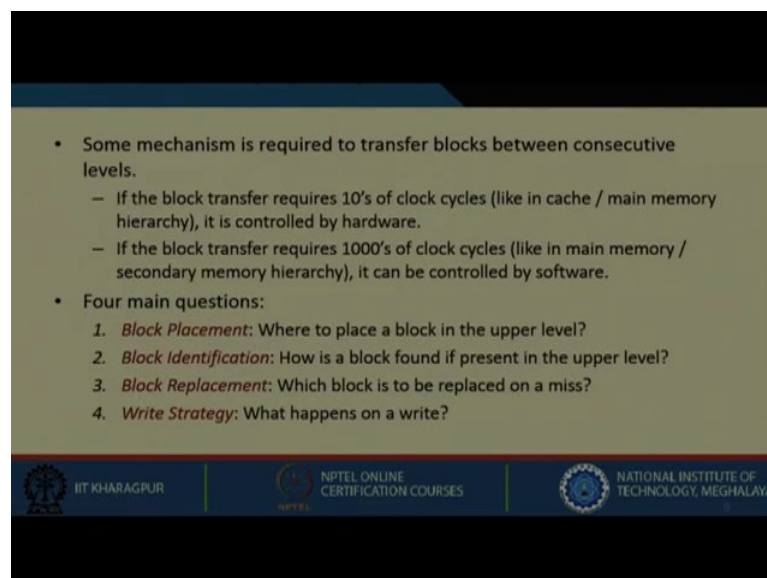
We are first checking into the top level L1. Is there any mechanism that we in advance know that this particular word is there or not, or the checking can it be made little fast so that hit time can made little faster. So, there are some kinds of things that can be done, check happens on every memory access and that is why it affects the hit times. If we can have some kind of hardware implementation for this matching, it can be performed in reasonable time.

(Refer Slide Time: 18:33)



Some mechanism is required to transfer blocks between consecutive levels; as we know that if there is a miss in both L1 and L2, you have to bring a block from main memory to L2 level and then from L2 to L1 level. So, whenever there is a miss in some level you have to transfer a block. If the block transfer requires 10s of clock cycles like cache or main memory, it is controlled by hardware; but if the block transfer requires 1000s of clock cycles like in main memory, secondary memory etc., it can be controlled by software.
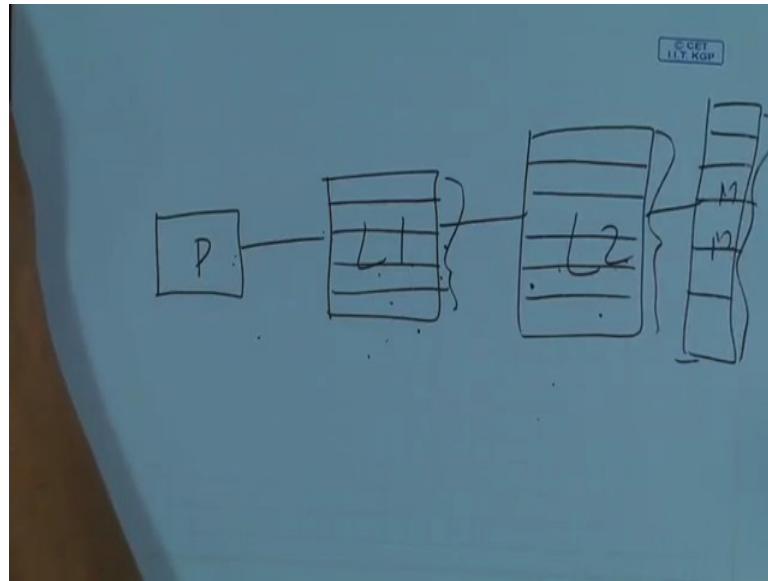
At this point of time few questions arise. One is block placement. Let us say we have different levels of memory.

So, again we have processor, we have L1 cache, we have L2 cache, and we have main memory. Now it has got some blocks.

Now, the point is where we will make the placement of block. From here we will bring the block and place it here. It needs to be understood next is block identification. What is blocked identification? How is a block found if present in a upper level? So, we must be matching something when the CPU generates a logical address and it says that this is the logical address I am looking for this particular word, then something much must be matched to know if that block is present in the upper level or not. So, block identification is also important.
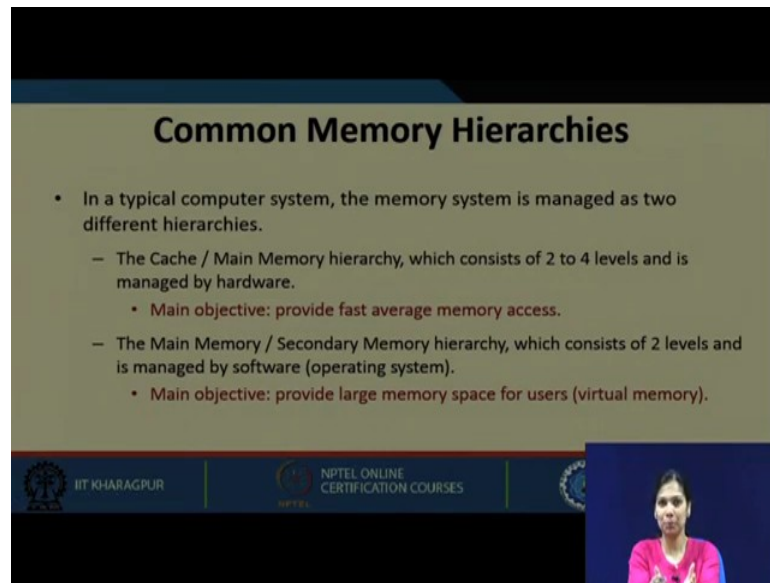
Similarly, block replacement … which block is to be replaced on a miss. What do you mean by that? Now you see in this particular diagram we have more number of blocks, here we have less number of blocks, and we are saying that every time the processor accesses the L1 cache it will be the fastest. So, all the programs or data are brought into L1 cache and then it is accessed.

Now, the size of this is small. So, at a time limited number of blocks can be present here. If you want to run your program that requires more number of blocks, in that case what you have to do? You have to replace some of the existing blocks from here and bring a new block. So, for that some block replacement strategies must be there; that means,

some blocks can be removed from this particular level and then some blocks from this level can be brought in. So, this is called block replacement.

Write strategy means what happens on a write whenever we are performing a write operation. What strategies must be used for this write operation?

(Refer Slide Time: 22:56)



In a typical computer system the memory system is managed as two different hierarchies, cache memory / main memory hierarchy which consists of 2 to 4 levels and is managed by hardware. The main objective of this particular hierarchy is to provide fast average access time for instructions and data.

Then we have main memory / secondary memory hierarchy that consists of two levels and is managed by software, that is by the operating system which is the interface between the user and the hardware part of our computer. The main objective is to provide large memory space for the user, called virtual memory.
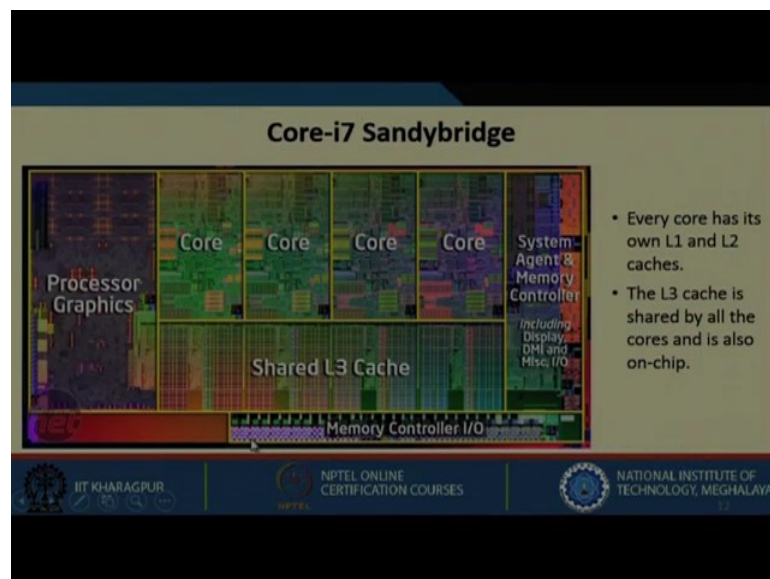
So, basically this virtual memory is a concept where it says that you have a large memory space at your disposal, but actually you are having the space as main memory only. Programs and data are brought from secondary storage to main memory as and when it is required. This is called virtual memory, which will not be taken up in this particular course in detail.

(Refer Slide Time: 24:51)



This is the typical memory hierarchy of Intel core i. There are four cores; core 0, core 1, core 2, and core 3. Inside each of the cores you have registers, you have separate L1 cache for data and instruction, we have an unified L2 cache.

(Refer Slide Time: 27:39)



Next is core-i7 SandyBridge. Here you have 4 cores, every core has its own L1 and L2 cache, and L3 is shared by all the cores and is also on chip. This is the processor, graphics processor, the memory controller, etc.

So, we come to the end of lecture 29 where we discussed about memory hierarchy in details. We have seen that how memory hierarchy actually affects the overall performance.

Thank you.