

UNIT IV

MEMORY SYSTEM

Basic concepts of Semiconductor RAMs - ROMs – Speed, Size and Cost – Cache memories – Performance consideration – **Virtual memory** – Memory Management requirements – Secondary storage - Case Study: Memory Organization in Multiprocessors



Recap the previous Class



Virtual memories

- Recall that an important challenge in the design of a computer system is to **provide a large, fast memory system at an affordable cost.**
- Architectural solutions to increase the **effective speed and size** of the memory system.
- **Cache memories** were developed to increase the effective speed of the memory system.
- **Virtual memory** is an architectural solution to increase the effective size of the memory system.

Virtual memories (contd..)

- Recall that the addressable memory space depends on the number of address bits in a computer.
 - For example, if a computer issues 32-bit addresses, the addressable memory space is 4G bytes.
- **Physical main memory** in a computer is generally **not as large** as the entire possible addressable space.
 - Physical memory typically ranges from a few hundred megabytes to 1G bytes.
- Large programs that **cannot fit completely** into the main memory have their **parts stored on secondary storage** devices such as magnetic disks.
 - Pieces of programs must be transferred to the main memory from secondary storage before they can be executed.

Virtual memories (contd..)

- When a new piece of a program is to be transferred to the main memory, and the **main memory is full**, then some other piece in the main memory **must be replaced**.
- Operating system automatically transfers data between the main memory and secondary storage.
 - Application programmer need not be concerned with this transfer.
 - Also, application programmer **does not need to be aware of the limitations** imposed by the available physical memory.



Virtual memories (contd..)

- Techniques that automatically move program and data between main memory and secondary storage when they are required for execution are called virtual-memory techniques.
- Programs and processors reference an instruction or data independent of the size of the main memory.
- Processor issues binary addresses for instructions and data.
 - These binary addresses are called **logical or virtual addresses**.

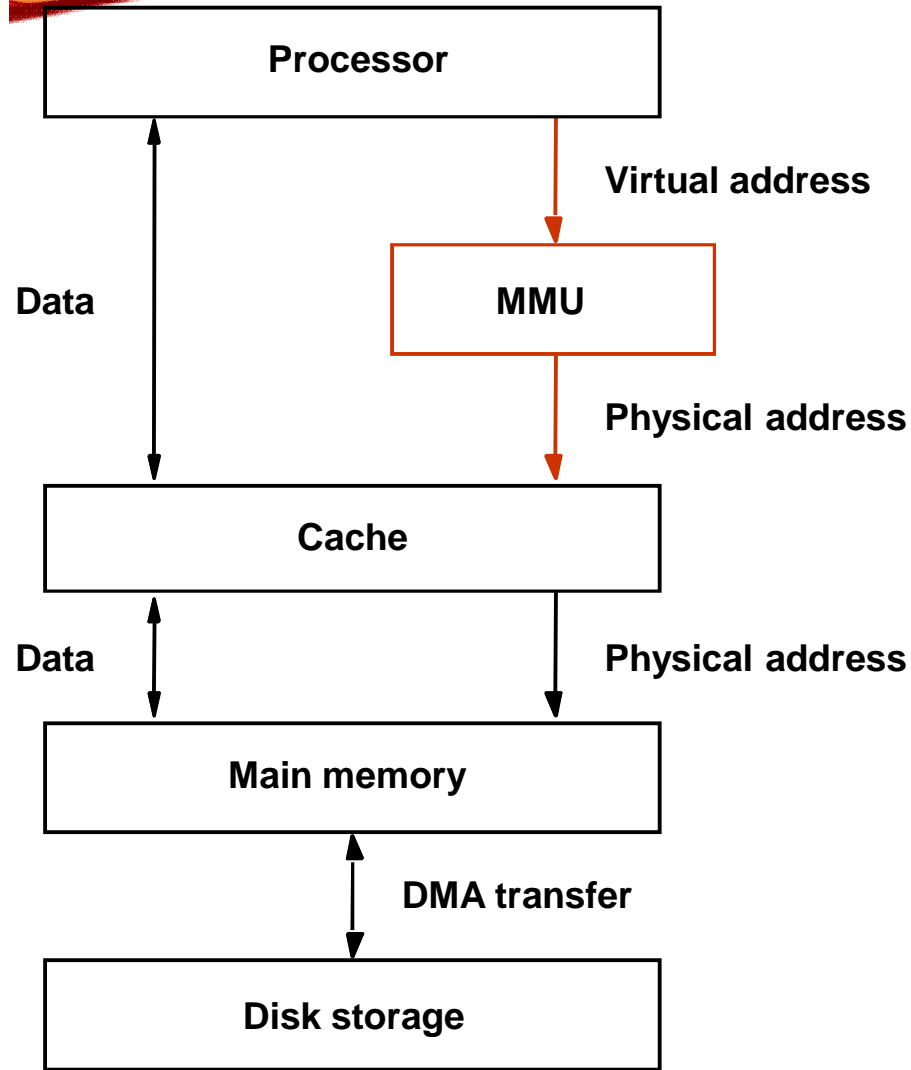


Virtual memories (contd..)

- Virtual addresses are translated into physical addresses by a **combination of hardware and software subsystems.**
 - If virtual address refers to a part of the program that is currently in the main memory, it is accessed immediately.
 - If the address refers to a part of the program that **is not currently in the main memory**, it is first transferred to the main memory before it can be used.



Virtual memory organization



- Memory management unit (**MMU**) translates virtual addresses into physical addresses.
- If the desired data or instructions **are in the main memory they are fetched** as described previously.
- If the desired data or instructions are **not in the main memory, they must be transferred** from secondary storage to the main memory.
- MMU causes the operating system to bring the data from the secondary storage into the main memory.



Address translation

- Assume that program and data are composed of fixed-length units called **pages**.
- A page consists of a **block of words** that occupy contiguous locations in the main memory.
- Page is a **basic unit of information** that is transferred between secondary storage and main memory.
- Size of a page commonly ranges from **2K to 16K bytes**.
 - Pages should not be too small, because the access time of a secondary storage device is much larger than the main memory.
 - Pages should not be too large, else a large portion of the page may not be used, and it will occupy valuable space in the main memory.



Address translation (contd..)

- Concepts of virtual memory are similar to the concepts of cache memory.
- **Cache memory:**
 - Introduced to bridge the speed gap between the processor and the main memory.
 - Implemented in hardware.
- **Virtual memory:**
 - Introduced to bridge the speed gap between the main memory and secondary storage.
 - Implemented in part by software.

Address translation (contd..)

- Each **virtual or logical address** generated by a processor is interpreted as a virtual page number (high-order bits) plus an offset (low-order bits) that specifies the location of a particular byte within that page.
- Information about the main memory location of each page is kept in the **page table**.
 - Main memory address where the page is stored.
 - Current status of the page.
- Area of the main memory that can hold a page is called as **page frame**.
- Starting address of the page table is kept in a **page table base register**.



Address translation (contd..)

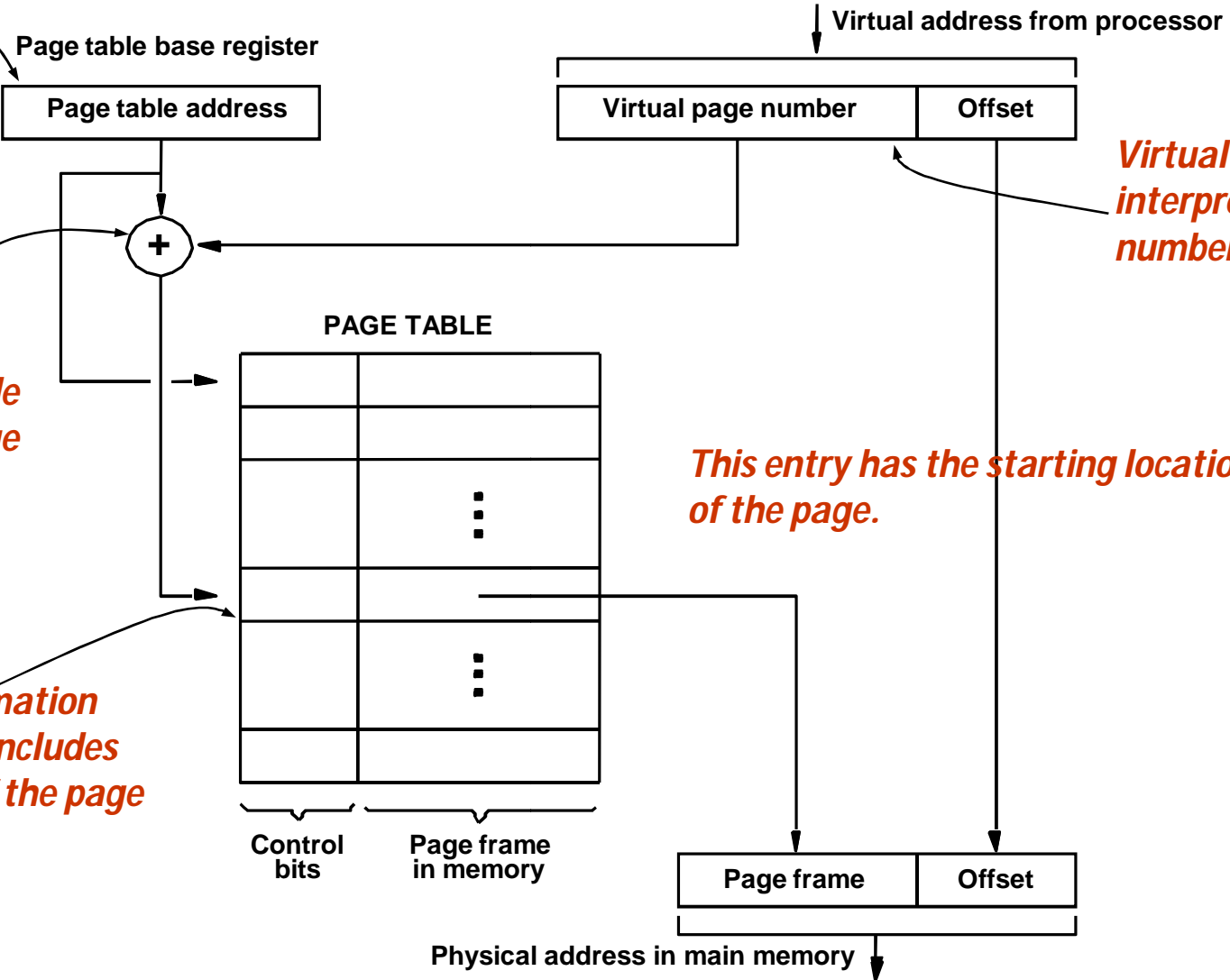
- Virtual page number generated by the processor **is added to the contents of the page table base register.**
 - This provides the address of the corresponding entry in the page table.
- The contents of this location in the page table give the starting address of the page if the page is currently in the main memory.



PTBR holds the address of the page table.

PTBR + virtual page number provide the entry of the page in the page table.

Page table holds information about each page. This includes the starting address of the page in the main memory.



Virtual address is interpreted as page number and offset.

This entry has the starting location of the page.



Address translation (contd..)

- Page table entry for a page also **includes some control bits** which describe the status of the page while it is in the main memory.
- One bit indicates the **validity of the page**.
 - Indicates whether the page is actually loaded into the main memory.
 - Allows the operating system to invalidate the page without actually removing it.
- One bit indicates whether the page has been **modified** during its residency in the main memory.
 - This bit determines whether the page should be written back to the disk when it is removed from the main memory.
 - Similar to the dirty or modified bit in case of cache memory.



Address translation (contd..)

- Where should the page table be located?
- Recall that the page table is used by the MMU for every read and write access to the memory.
 - Ideal location for the page table is within the MMU.
- Page table is quite large.
- MMU is implemented as part of the processor chip.
- Impossible to include a complete page table on the chip.
- Page table is kept in the main memory.
- A copy of a small portion of the page table can be accommodated within the MMU.
 - Portion consists of page table entries that correspond to the most recently accessed pages.



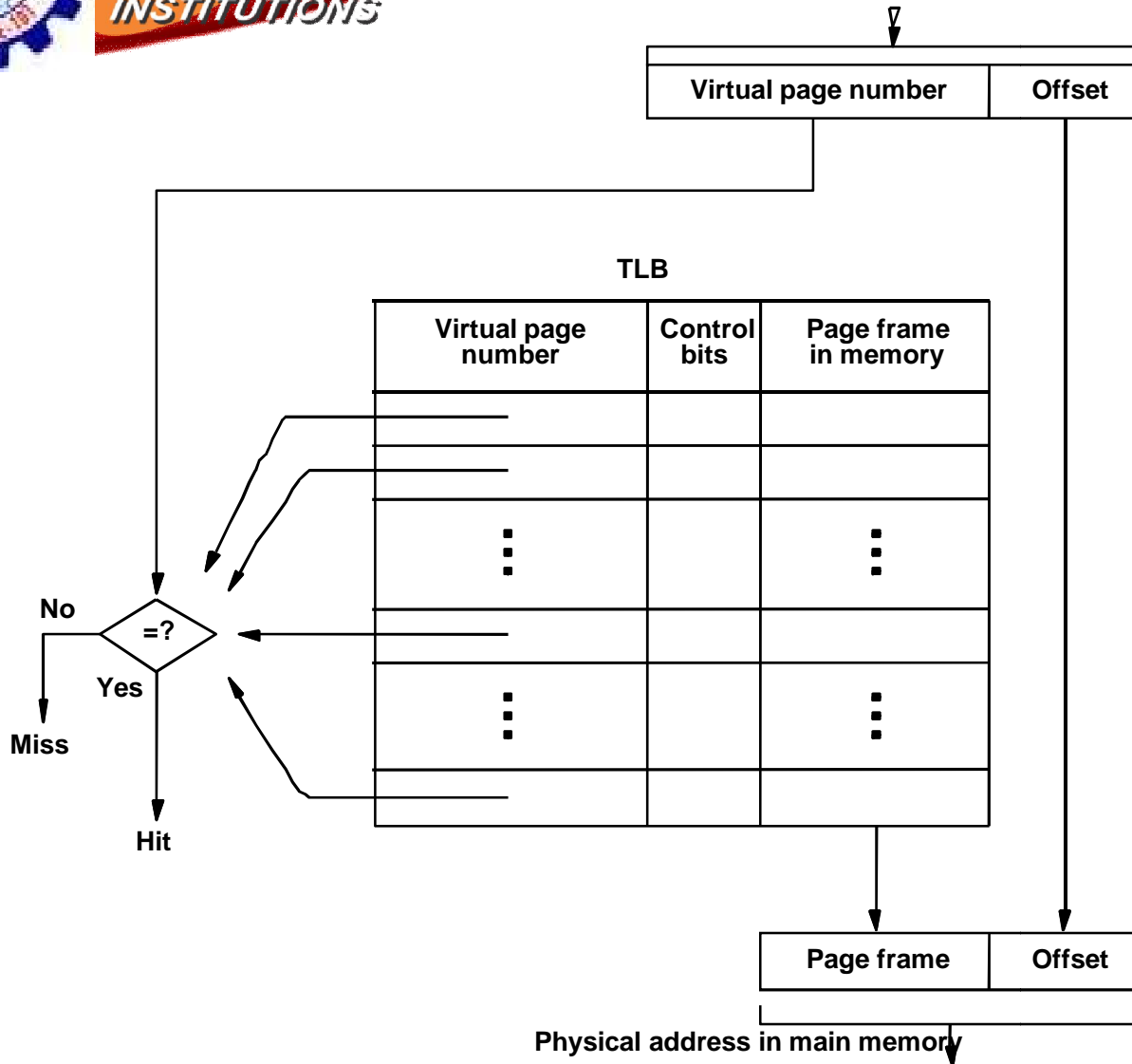
Address translation (contd..)

- A small cache called as **Translation Lookaside Buffer (TLB)** is included in the MMU.
 - TLB holds page table entries of the **most recently accessed pages**.
- Recall that cache memory holds most recently accessed blocks from the main memory.
 - **Operation of the TLB and page table** in the main memory is **similar to the operation of the cache and main memory**.
- **Page table entry for a page includes:**
 - Address of the page frame where the page resides in the main memory.
 - Some control bits.
- In addition to the above for each page, TLB must hold the virtual page number for each page.



Address translation (contd..)

Associative-mapped TLB



- **High-order bits** of the virtual address generated by the processor select the virtual page.
- These bits are compared to the virtual page numbers in the TLB.
- If there is a match, a **hit occurs** and the corresponding address of the page frame is read.
- If there is no match, a **miss occurs** and the page table within the main memory must be consulted.
- **Set-associative mapped TLBs** are found in commercial processors.



Address translation (contd..)

- How to keep the entries of the TLB coherent with the contents of the page table in the main memory?
- Operating system may change the contents of the page table in the main memory.
 - Simultaneously it must also invalidate the corresponding entries in the TLB.
- A **control bit** is provided in the TLB to invalidate an entry.
- If an entry is invalidated, then the TLB gets the information for that entry from the page table.
 - Follows the same process that it would follow if the entry is not found in the TLB or if a "miss" occurs.



Address translation (contd..)

- What happens if a program generates an access to **a page that is not in the main memory?**
- In this case, a **page fault** is said to occur.
 - Whole page must be brought into the main memory from the disk, before the execution can proceed.
- Upon **detecting a page fault by the MMU**, following actions occur:
 - MMU asks the operating system to intervene by raising an exception.
 - Processing of the active task which caused the page fault is interrupted.
 - Control is transferred to the operating system.
 - Operating system copies the requested page from secondary storage to the main memory.
 - Once the page is copied, control is returned to the task which was interrupted.



Address translation (contd..)

- Servicing of a page fault requires transferring the requested page from secondary storage to the main memory.
- This transfer may incur a long delay.
- While the page is being transferred the operating system may:
 - Suspend the execution of the task that caused the page fault.
 - Begin execution of another task whose pages are in the main memory.
- Enables efficient use of the processor.



TEXT BOOK

Carl Hamacher, Zvonko Vranesic and Safwat Zaky, "Computer Organization", McGraw-Hill, 6th Edition 2012.

REFERENCES

1. David A. Patterson and John L. Hennessey, "Computer organization and design", MorganKauffman ,Elsevier, 5th edition, 2014.
2. William Stallings, "Computer Organization and Architecture designing for Performance", Pearson Education 8th Edition, 2010
3. John P.Hayes, "Computer Architecture and Organization", McGraw Hill, 3rd Edition, 2002
4. M. Morris R. Mano "Computer System Architecture" 3rd Edition 2007
5. David A. Patterson "Computer Architecture: A Quantitative Approach", Morgan Kaufmann; 5th edition 2011

THANK YOU