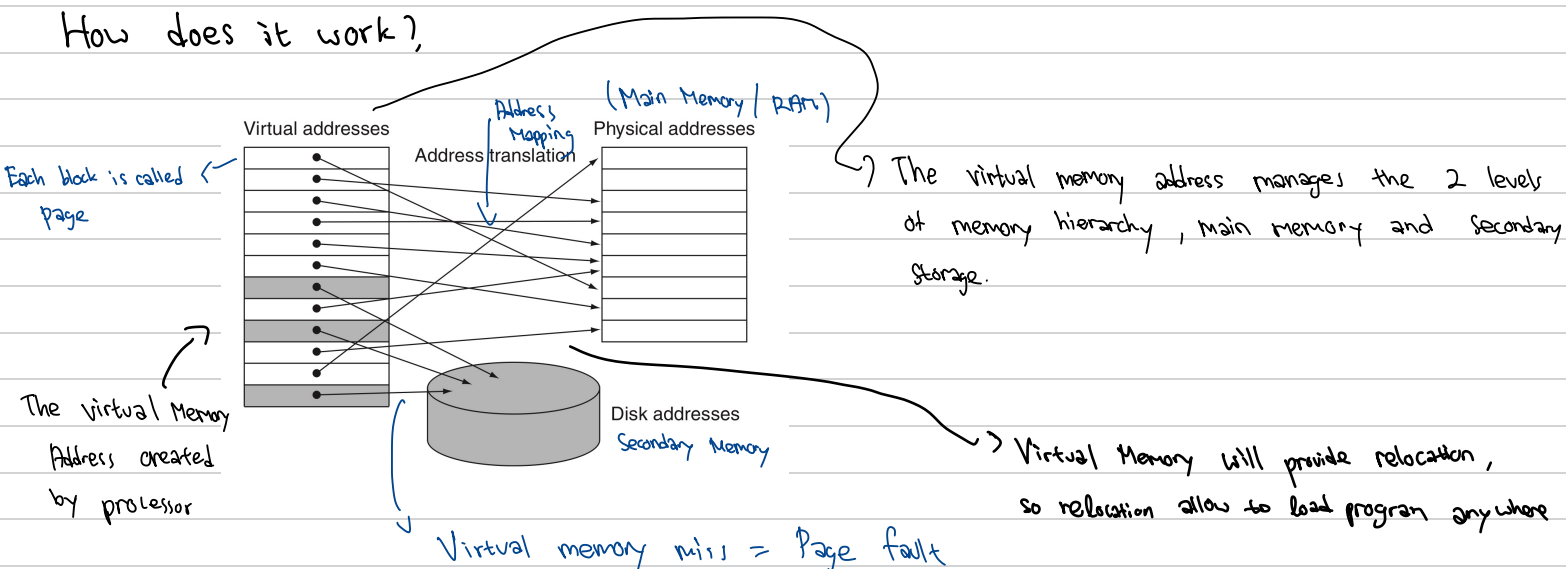


5.7 ~ 5.8

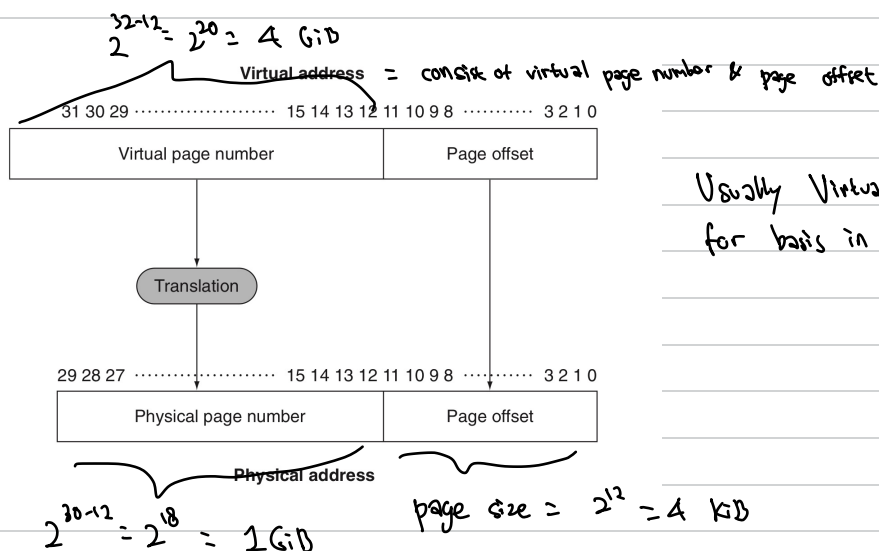
# Virtual Memory

Definition: A Memory Management technique that makes main memory (RAM) to have larger memory by using secondary memory.

How does it work?



How does Address Mapping work?



Usually Virtual pages > Physical Pages for basis in illusion of virtual Memory

## Key decisions in Designing Virtual Memory systems

Background: Page fault → takes million of clock cycles to process

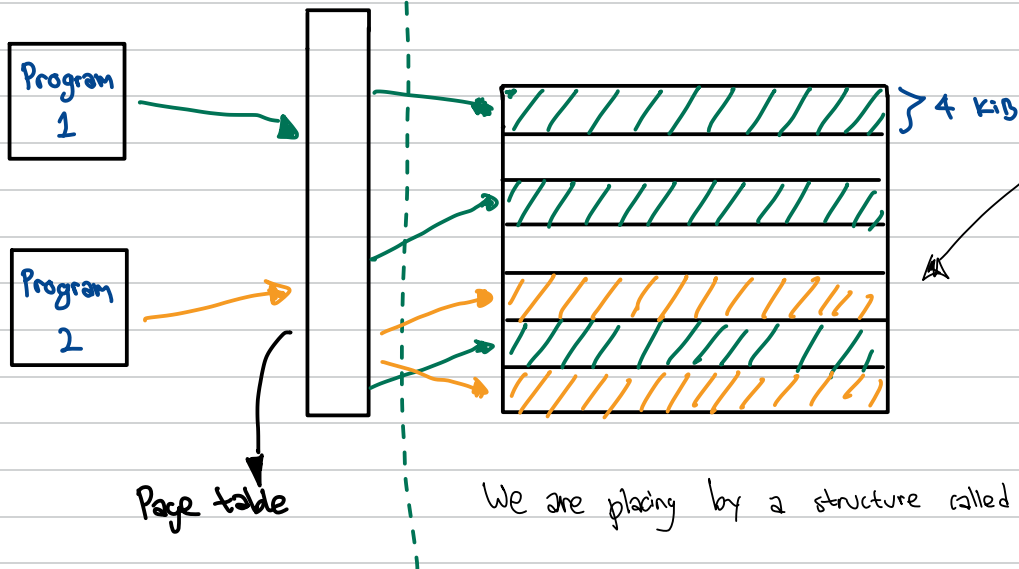
- Pages need to be large
- Pages fault can be handle by software
- Virtual Memory chooses write back instead of write-through

Since disk transfer time is small compared to its access time.

# Placing a Page and Finding it Again

Virtual Memory

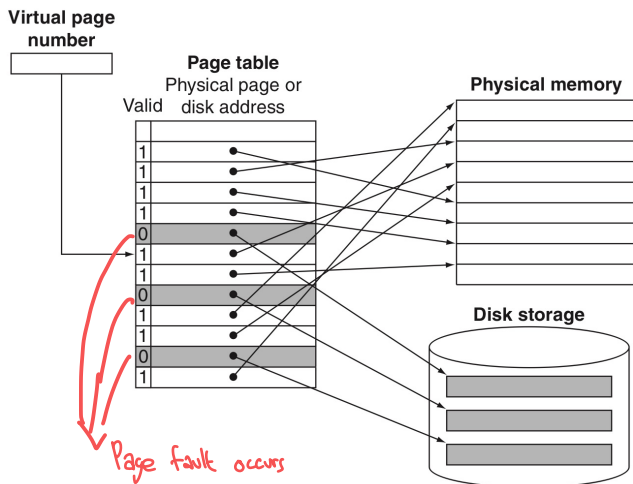
Physical Memory



We are placing by a structure called **page table**,

Page tables is indexed by the page number to discover the corresponding page number.

## Page faults



Page fault  $\rightarrow$  OS gets the control

Find Available pages

Need keep track of the location in secondary memory of each page in virtual address space. Using LRU - least recently used scheme

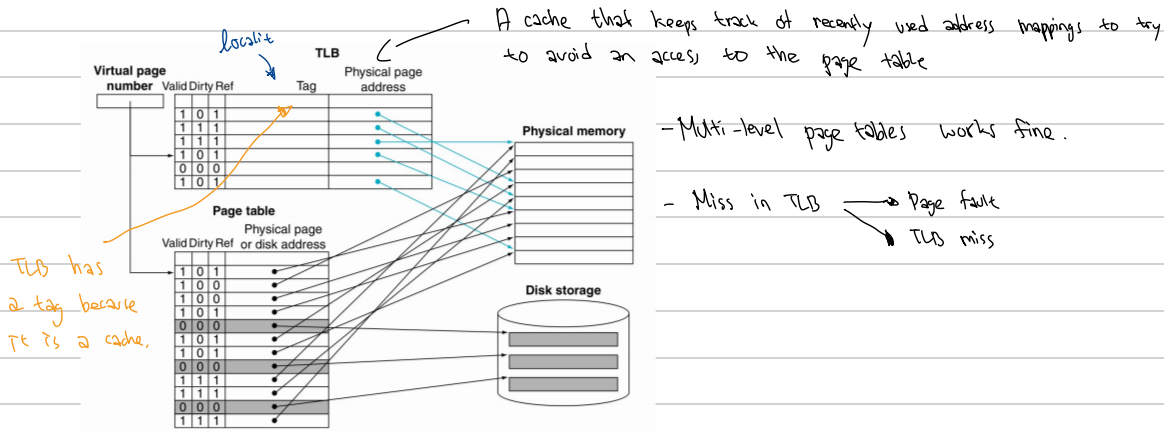
Also create space "**swap space**" for page replacement

Besides that, Also record each virtual page is stored on disk.

## 5 techniques for reducing the amount of storage in page table

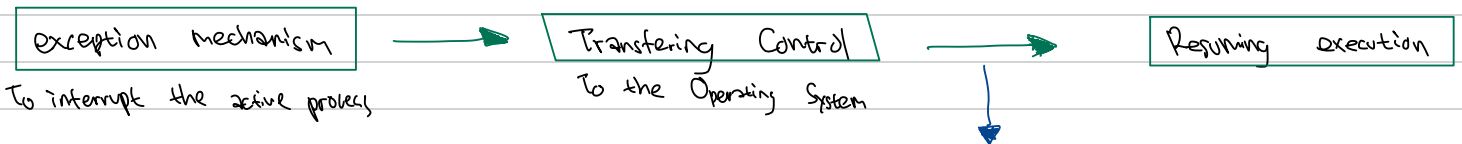
- Keep limit of the register  $\rightarrow$  Address space need in just 1 direction
- Create 2 separate page tables & 2 separate limit  $\rightarrow$  can grow from 2 direction
- Hashing the virtual address, making (page table  $\approx$  physical page) sized
- Page tables to paged  $\rightarrow$  page tables placed in virtual address space.
- Multiple level of page tables

## TLB: Translation - lookaside buffer



## How to handle TLB Misses & Page faults

- How can TLB Misses: no entry in the TLB Matches a virtual address



If it's the **Virtual Address** caused page fault.

3 steps will be taken:

1. Look up page table → find location of referenced page.
2. Choose a physical page to replace
3. Read referenced page to the physical page

Need to be declared if the TLB miss/Page fault happens.

## Common Framework for Memory Hierarchy

Block placement scheme we learn until now

Scheme name	Number of sets	Blocks per set
Direct mapped	Number of blocks in cache	1
Set associative	Number of blocks in the cache Associativity	Associativity (typically 2-16)
Fully associative	1	Number of blocks in the cache

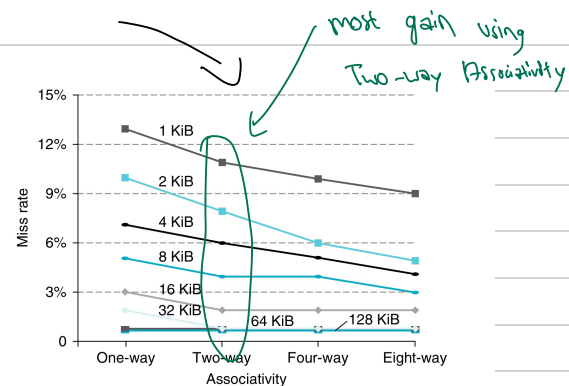
chose depend on the MISS vs COST of implementing

Ex: L2 cache → higher associativity (hit times x critical)

Virtual Memory → Full Associativity

- Misses are very expensive
- Can easily index, no extra hardware
- Allow software to use sophisticated schemes

↑ degree → ↓ miss rate



### 2 Replacement Blocks Techniques

- Random: Randomly choose, use some hardware
- Least Recently Used (LRU) = replace block that hasn't been used for the longer time.  
↳ used in larger associativity

## 2 techniques on Write

### Write-through used in cache



Definition: Information written to cache & lower level memory hierarchy (RAM for a cache)

Advantages:

- Easier to implement
- Miss are cheaper

### Write-back

used in Virtual memory



Definition: The information is written to cache only.

The information written to lower level memory hierarchy

only when it's replaced

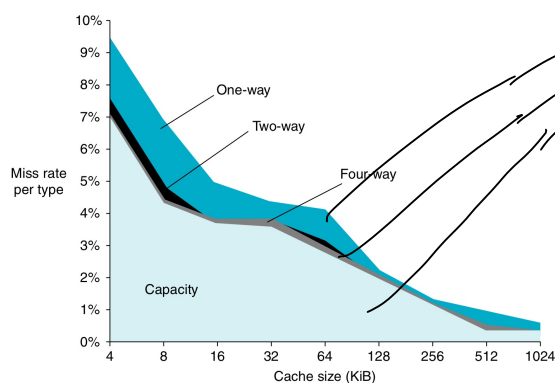
Advantages:

- Effective use of high-bandwidth transfer

- Multiple writes  $\rightarrow$  1 writes

- Individual words can be written by the processor

## Miss Rate Sources



### 3 Sources of Miss Rate

1. Compulsory misses (cold-start misses):

cause: first access to a block that isn't in the cache.

2. Capacity misses

cause: cache can't contain all the blocks

3. Conflict misses

cause: multiple blocks compete for the same set

## Design Consideration

Design change	Effect on miss rate	Possible negative performance effect
Increases cache size	Decreases capacity misses	May increase access time
Increases associativity	Decreases miss rate due to conflict misses	May increase access time
Increases block size	Decreases miss rate for a wide range of block sizes due to spatial locality	Increases miss penalty. Very large block could increase miss rate

Block placement scheme