

Proportional Share (PS)

It is based around a simple concept instead of optimizing for TAT or RT a scheduler might instead tried to guarantee that each job or process obtain a certain percentage of CPU time. A good example of PS scheduling is known as lottery scheduling.

Basic Concept

- Tickets: It is used to represent the share of a resource that a process should received. The percent of tickets that a process has represents its share of the system resource in problem.
Ex: Consider two processes A and B and further that A has 75 tickets while B has only 25. Thus what we would like is for A to receive 75% of the CPU and B the remaining 25%.
Holding a lottery is straight forward the schedule must know how many total tickets there are (in eg. 100). Therefore the scheduler then picks the winning ticket which is a number from 0-99 assuming A holds tickets 0-74 and B holds 75-99, the winning ticket simply determines whether A or B wins the .

63	85	70	39	76	17	29	41	36	39	15	85	68
A	B	A	A	B	A	A	A	A	A	A	B	A
B	A	A	A	A	A	A	A	A	A	A	B	A
83	63	62	43	0	49	12						

No. on the wheel.

in lottery scheduling.

The use of randomness leads to a probabilistic correctness in meeting the desired proportion. But no guarantee. In this example we only gets to run for 4 out 20 times slides instead of the desired 25% of allocation.

Ticket Mechanism:

Lottery scheduling also provides a no. of mechanisms to manipulate tickets in different & some time useful ways.

Concept of ticket currency: Currency allows a user with a set of tickets to allocate among their own jobs in whatever currency they would like. The system then automatically converts said currency into the correct global value.

Ticket transfer:

With transfer a process can temporarily hand off its tickets to another process. This availability is useful in a client-server setting where a process (client process) sends a message to a server asking it to do some work on the client's behalf. To speed up the work the client can pass the tickets to the server and thus tried to maximise the performance of the server while the server is handling the client's request. When finished, the server then transfers the tickets back to the client and all ^{is as} before.

Lottery Scheduling decision code

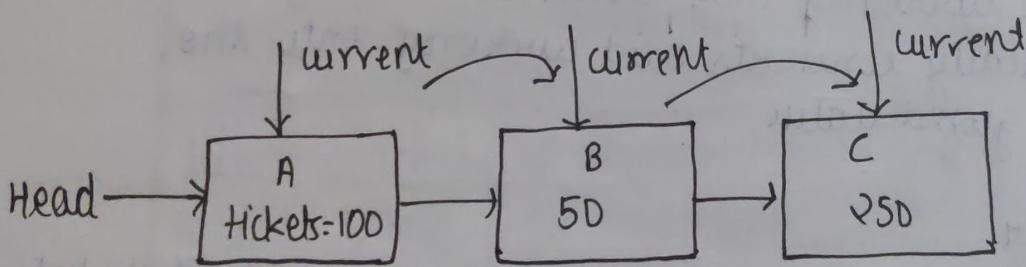
- ① int counter = 0
- ② int winner = getrandom(0, Total tickets);

// winner use same call to a random number generation to get value, between 0 and the total tickets.

- ③ node * current = head;

- ④ while (current)

```
{  
    counter = counter + current->tickets;  
    if (counter > winner)  
        break;  
    current = current->next;  
}
```



$$\text{counter} = 0$$

$$\text{winner} = 300$$

$$\text{counter} = 0 + 100$$

$$= 100$$

$$\text{counter} =$$

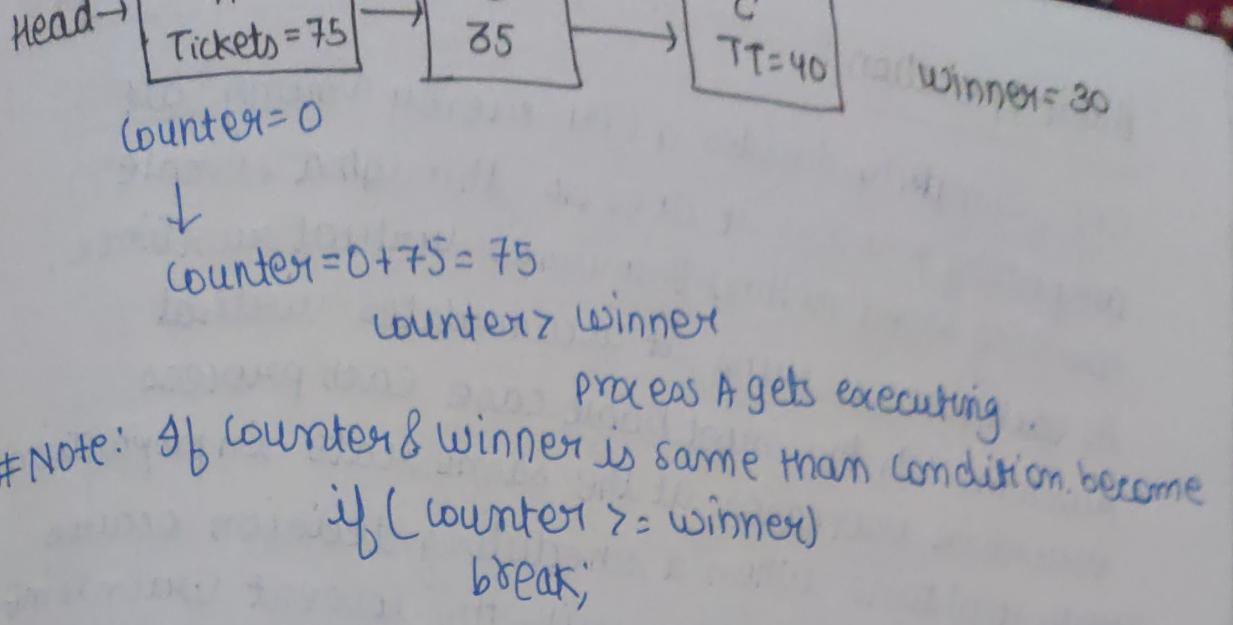
$$100 + 50$$

$$\text{counter} = 150 + 250$$

$$= 400$$

greater than
winner

: process C
start
executing



Stride Scheduling. A deterministic FSS

A	B	C	who runs
100	80	40	
0	0	0	A
100	0	0	B
100	80	0	C
100	80	40	C
100	80	80	

Completely fair Scheduling

The main focus of CFS is highly efficient and scalable manner. To achieve this goal CFS aims to spend very little time making scheduling decisions. Through both its inherent design and its clever use of datastructure is well suited to the task.

Basic Operation:

CFS simply divides a CPU evenly among all computing processes. It does so through a simple counting-based method known as virtual runtime. As each process runs, it accumulates virtual runtime. In the most basic case, each process' virtual runtime increases at the same rate in proportion with real time. When a scheduling decision occurs, CFS will pick the process with the lowest virtual runtime to run next.

WIM
The main issue is how does the scheduler know when to stop the currently running process and run the next one?

- (1) If CFS switches too often, fairness is increased as CFS will ensure that each process receives its share of CPU even over a single time window but at the cost of performance.
- (2) If CFS switches less often, performance is increased but at the cost of mean term fairness.

control Parameters

↳ Sched-Latency:

CFS uses this value to determine how long one process should run before considering a switch.

Typically Sched-Latency value is 46ms

CFS divides this value by the number (N) of processes running on the CPU to determine

the time slice for a process and thus ensure that over this period of time CFS will be completely fair.

for example: if there are N=4 processes running then per process slice = $\frac{46}{4} = 12\text{ ms}$.

CFS uses this value to determine how long one process should run before considering a switch

CFS then schedules the first process and runs it until

it has used 12ms of runtime and then checks to see if there is a process with lower runtime to run instead. In this case there is, and CFS is switched to one of the other three jobs and so forth.

e.g: Consider 4 jobs A, B, C, D each run for two time slices in this fashion two of them (C & D) then complete, leaving just to remaining which each run for 24ms in Round Robin fashion.

(Q) What if there are too many processes running?
→ Here CFS has another parameter min-granularity which is usually set to a value like 6ms. CFS will never set the time slice of a process to less than this value, ensuring that not too much time is spent in scheduling overhead.

For example:

If there are (too many) 10 processes are running over original calc. would divide sched latency by 10 to determine the time slice (4.8ms).

However because of min-granularity CFS will set the time slice of each process to 6ms instead.

Note: CFS utilizes a periodic timer interrupt which means it can only make decisions at fixed time intervals this interrupt goes off frequently (1ms), giving CFS a chance to wakeup and determine if the current job has reached the end of its run.

If a job/process has a time slice that is not a perfect multiple of the time interrupt interval i.e. CFS tracks runtime precisely which means that over the long interval it will eventually approximate ideal sharing of the CPU.

Weighing or Niceness

CFS also enables controls over process priority from enabling users or administrators to give some processes a higher share of the CPU. It does this through a classic UNIX mechanism known as the **Niceness** of a process.

The nice parameter can be set anywhere from -20 to +19 for a process with a default of zero.

+ve nice values imply lower priority and -ve nice values imply higher priority. CFS maps the nice value of each process to a weight as shown here.

```
static const int prio_to_weight[40] = {
```

```
/* -20 */ 88761, 71755, 56483, 46273, 36291
```

```
/* -15 */ 29154, 23254, 18705, 14949, 11916.
```

```
/* -10 */ 9548, 7620, 6100, 4904, 3906
```

```
/* -5 */ 3121, 2501, 1991, 1586, 1277
```

```
/* 0 */ 1024, 820, 655, 526, 423
```

```
/* 5 */ 355, 272, 215, 172, 137
```

```
/* 10 */ 110, 87, 70, 56, 45
```

```
/* 15 */ 36, 29, 23, 18, 15
```

```
/* 20 */ 10, 8, 7, 6, 5
```

```
}
```

These weights allow us to compute the effective time slice of each process. But accounting for their priority differences the formula we to do so is as follows

assuming N processes.

$$\text{Time-slice}_k = \frac{\text{Weight}_k}{\sum_{i=0}^{n-1} \text{Weight}_i} \text{ sched-latency.}$$

example

Consider there are 2 processes A and B. A, because its more precious job is given a higher priority by assigning it a nice value of -5. B, has default priority. (0).

weight of A = 3121 because it has priority -5

weight of B = 1024 0.

$$\text{Time-slice of A} = \frac{3121 \times 48}{\text{Weight}_A + \text{Weight}_B} = \frac{3121}{3121 + 1024} \times 48 = \frac{3121}{4145} \times 48 = 36$$

$$\text{Time-slice of B} = \frac{1024 \times 48}{4145} = 11.85 \text{ ms} \approx \underline{\underline{12 \text{ ms}}}$$

In addition to the generalize the time slice calculation, the way (F) calc. Vruntime must also be adopted. for new formula vruntime - The actual runtime that process i has accrued (runtime_i) and scales it inversely by the weight of the process by dividing the default weight of 1024 by its weight of i. As A's runtime will accumulates at one third rate of B's

$$V\text{runtime}_i = V\text{runtime}_i + \frac{\text{Weight}_i}{\text{Weight}_j} \times \text{runtime}_i$$

Note: One unique aspect of the construction of the table of weights above is that the table preserves CPU proportionality ratios when the diff. in nicevalue is constant

The main focus of CFS for a scheduler there are many facets of efficiency but one of them is as simple as this. When the scheduler has to find the next job to run. It should do so as quickly as possible. Simple data structure like list do not scale modern systems sometimes are comprised of many processes and thus searching through a long list every 80 milliseconds is wasteful.

CFS addresses this by keeping process in a red black tree. A red black tree is one of many types of balanced trees. In contrast to a simple binary tree (which can degenerate to list like performance under worst case insertion patterns) Balance trees do a little extra work to maintain low depth. and thus ensure that operations are logarithmic in time.

CFS doesn't keep all processes in this structure rather only running processes are kept there in. If a process goes to sleep (waiting for an I/O to complete or for a network packet to arrive) it is removed from the tree and keep track of elsewhere.

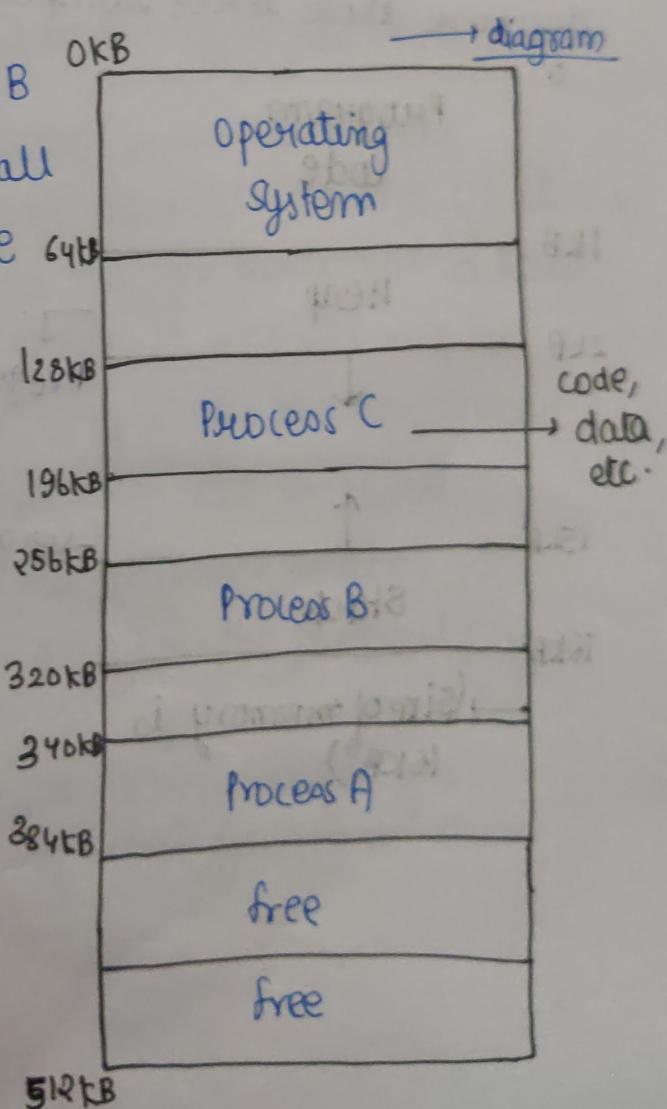
Memory Management

Multiprogramming and Time sharing.

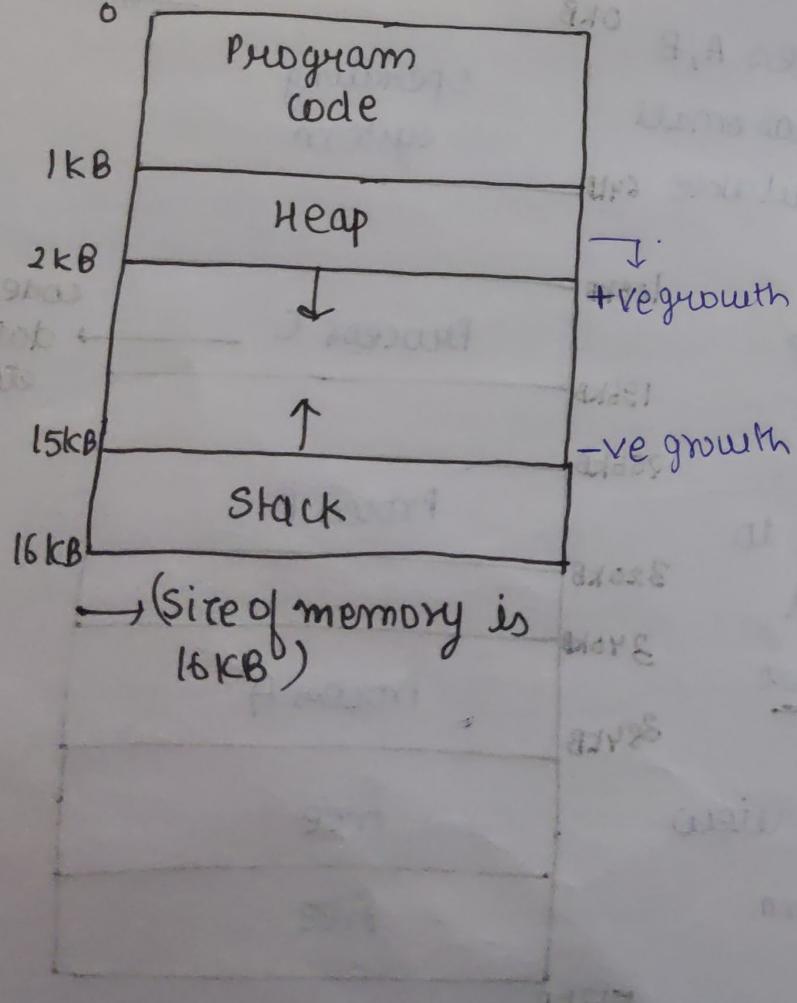
MP in which multiple processes were ready to run at a given time and the OS would switch b/w them. One way to implement Time sharing would be to run one process for a short time slice giving it full access to all memory than stop it. Save all of its state to some kind of disk (Physical memory) load some other processes states run it for a while and thus implement some kind of root sharing of the machine. Limitations of this approach is too slow when memory request increases.

In the diagram, 3 processes A, B and C. And each of them has small parts of 512 KB memory available for them.

However, we have to keep all processes demand and requires the operating system to create an easy to use abstraction of physical memory known as address space. Address space is the running programs view of memory in the system.



The address space of a process contains all of the memory state of the running program. For example the code of the program have to live in memory somewhere and thus they are in the address space. The program while it is running uses a stack to keep track of where it is in the function call chain as well as to allocate local variables and pass parameters and returns values to and from routines. Finally the heap is used for dynamically allocated user managed memory, such as that you might receive from a call to malloc() in C or new in CPP or (OOP) Java. So let us just assume those three components - code, stack and heap.



The main goal of OS for ease of use:

Transparency: OS virtual memory को ऐसे implement करता है जिससे running program अंदर निवार हो।

Protection:

Goal of OS:

1. Transparency

thus, fact prog

physical efficiency

vi term much in gm will include

protect

problem from a stor able conter (that

protected isol the ma

► window

Transparency: The OS should implement the virtual memory that is invisible to running program. Thus, the program should not be aware of the fact that memory is virtualized. Rather the program behaves as if it has its own private physical memory.

Efficiency: The OS should strive to make the virtualization as efficient as possible, both in terms of time (i.e. not making programs run much more slowly) and space.

In implementation time efficient virtualization, the OS will have to rely on hardware support, including hardware features such as TLB's.

Protection: The OS should make sure to protect processes from one another as well as the OS itself - from processes. When one process performs a load, a store, or an instruction fetch, it should not be able to access or affect in any way the memory contents of any other process or the OS itself. (that is, anything outside its address space).

Protection thus enables us to deliver the property of isolation among processes. Each process should be running in its own isolated world, safe from the # managers of (issues) of other faulty or even malicious processes.

Hardware based Translation (Address Translation)
The hardware transforms each memory access where the desired location or information is actually loaded. Thus on each and every memory reference, an address translation is performed by the Hardware to redirect application memory references to their actual locations in memory. Of course the hardware alone can't virtualize memory as it just provides the low level mechanism for doing so efficiently. The OS must get involved takes place, it must thus manage memory, keeping track of which locations are free and which are in use.

Assumptions:

- ① Specifically we will assume for now that the users address space must be placed contiguously in physical memory.
- ② The size of the address space is not too big specifically that it is less than the size of physical memory.
- ③ Each address space is exactly the same size.

128: move 0XD C1. ebx) / . eax;
132: addl \$ 0X03\$, +eax;
135: movl +eax, 0XD C1. ebx)

loads a value from memory, it increments it by 3 and then stores the value back into memory.

void func()

int x=3000; push x on stack
x=x+3;
printf("1.d", x) using set of

into assembly, which look something like this (3 lines written previously).

When these instructions run from the perspective of the process, the following memory access takes place.

Step 1: fetch instructions at address 128

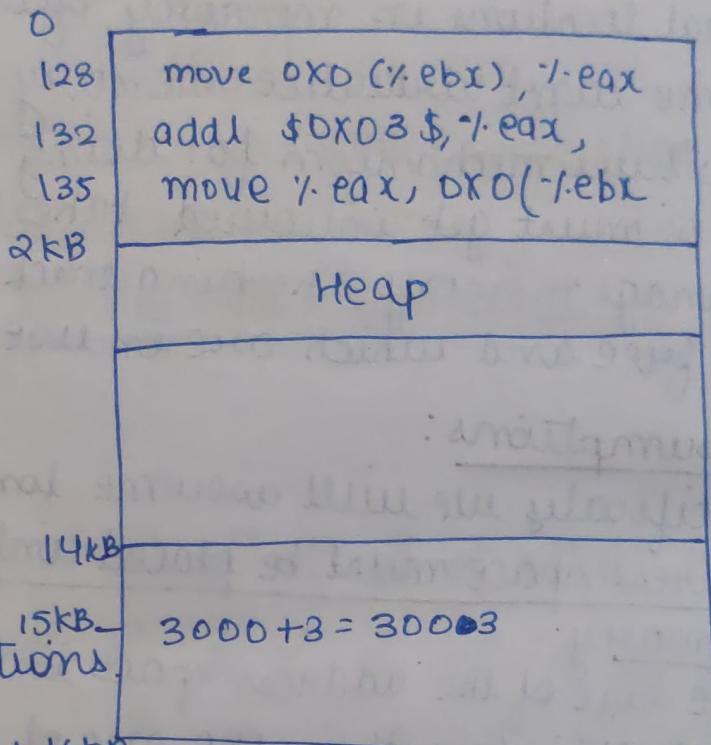
Step 2: execute the ins. and load the value from ebx to eax.

Step 3: fetch instructions at address 132

Step 4: execute this instruction.

Step 5: fetch the instruction at 16KB address 135

Step 6: execute this instruction and stored the value to address by 15KB



* Interposition [3]

↳ transparency: To virtualize memory OS wants to place the process somewhere else in physical memory, not at add. 0.

Thus we have problem: How can be relocate this process in memory in way that is transparent to the process?

Dynamic Relocation

Base and
Bound (limit)

First time sharing machines is a simple idea referred to as base and bounds also referred as dynamic relocation specifically we will need two hardware registers within each CPU one is called the base register and the other is bound register. This base and bound pair is going to allow us to place the address space anywhere we did like in physical memory and do so while ensuring that the process can only access its own address space. However when a program starts running the OS decides wherein physical memory, it should be loaded and sets the base register to that value.

Now, when any memory reference is generated by the process it is translated by the processor in the following manner.

$$\text{physical memory} = \frac{\text{virtual}}{\text{address}} + \text{base}$$

Suppose a process of address space of size 4KB has been loaded at physical address at 16KB. so what are the result of the address translation for 1KB, for 3000B, for 4400B.

$$\text{for } 1\text{KB} \rightarrow 16\text{KB} + 1\text{KB} = 17\text{KB}$$

$$\text{for } 3000\text{B} \rightarrow 16\text{KB} + 3\text{KB} = 19\text{KB}$$

$$\text{for } 4400\text{B} \rightarrow \text{overflow}$$

~~#~~ segmentation

So we have been putting the entire address space of each process in memory with the base and bound registers the OS can easily relocate processes to different parts of physical memory.

We have noticed there is a big chunk of free space right in the middle between the stack and the heap.

Although the space between the stack and heap is not being used by the process, it is still taking up physical memory when we relocate the entire address space somewhere in physical memory. Thus

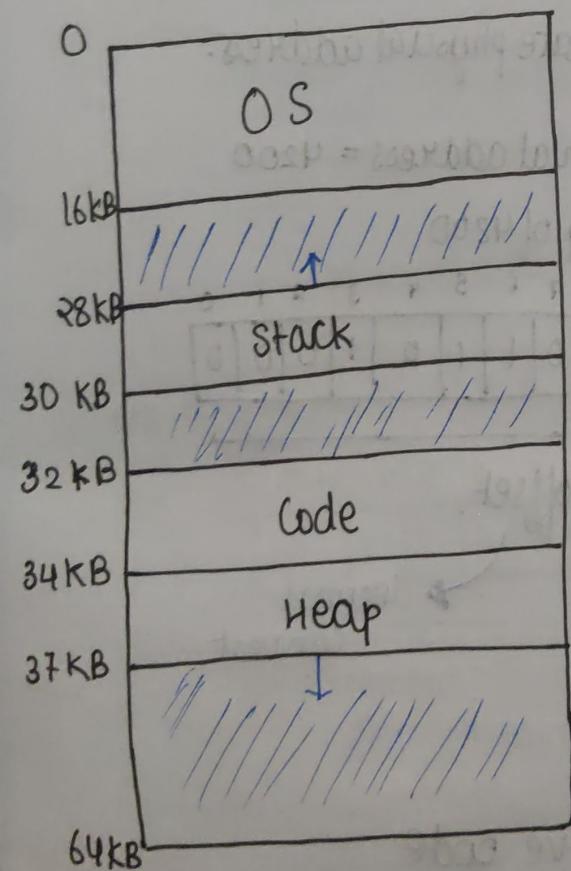
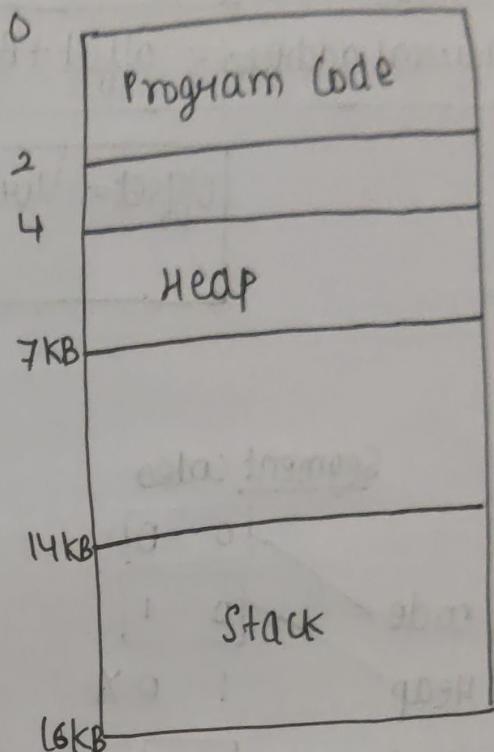
The simple approach of using a base and bounds register pair to virtualize memory is wasteful. It also makes it quite hard to run a program when the entire address space doesn't fit into memory.

To solve this problem new idea adopted which is known as segmentation or generalized based/bound method. The idea is simple instead of having just one base & bound pair in our memory management unit (MMU).

A segment is just a contiguous portion of the address space of a particular length and in our address space we have three logically different segments (Code, stack and heap).

what segmentation allows the OS to do is to place each one of these segments in different parts of physical memory and thus avoid filling physical memory with unused virtual address space.

segment	Base	Size
Code	32 KB	2 KB
Heap	34 KB	2 KB
Stack	28 KB	2 KB

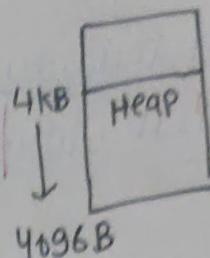


* Defragmentation

Imp page:

Virtual address = 4200

Heap = 34KB



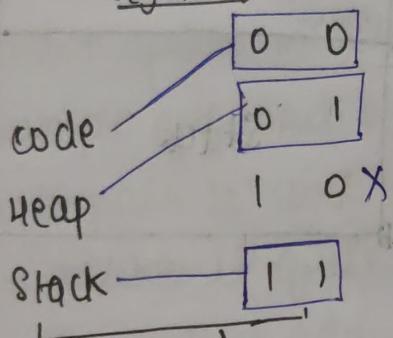
physical address = offset + base

offset = Virtual address - start heap address

$$\begin{aligned} \text{offset} &= 4200 - 4096 \\ &= 104 \end{aligned}$$

$$\begin{aligned} \text{physical address} &= 104 + 3 \\ &= 138 \end{aligned}$$

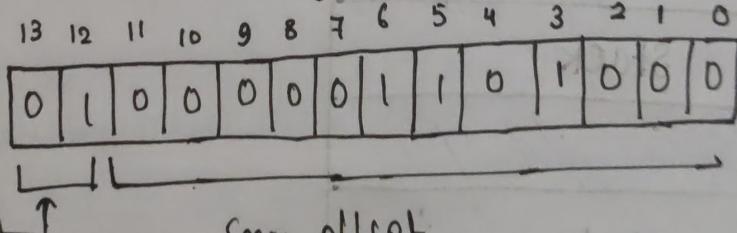
Segment codes



calculate physical address.

virtual address = 4200

Binary of 4200



Segment

codes



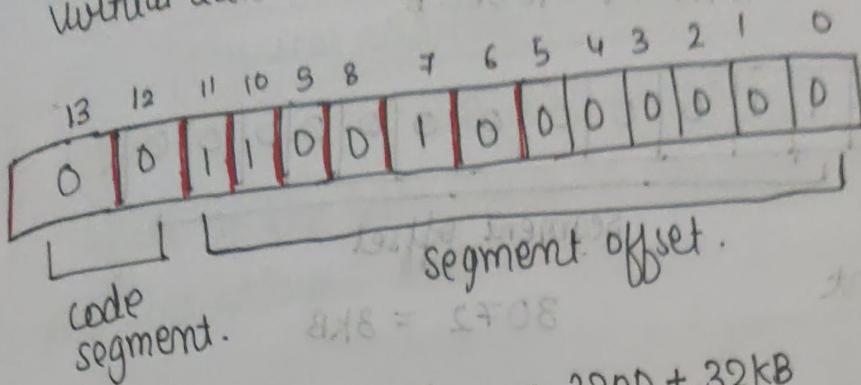
This will tell
the base

Segm. offset.

decimal
convert

for heap and code it is +ve code.

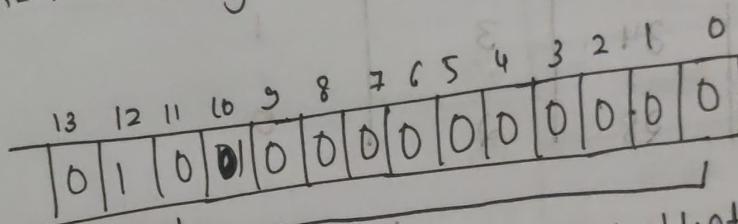
virtual address = 3200



$$\text{physical address} = 3200 + 32\text{KB}$$

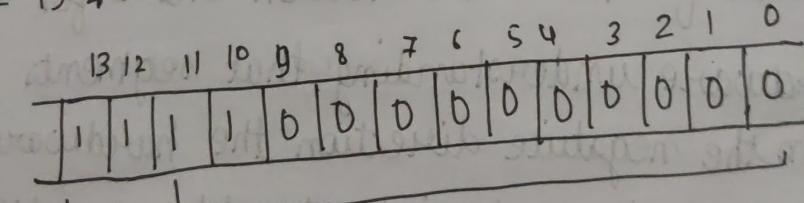
$$= \underline{\underline{3232\text{KB}}}$$

$$V.A = 5\text{KB} = 5120\text{Bytes}$$



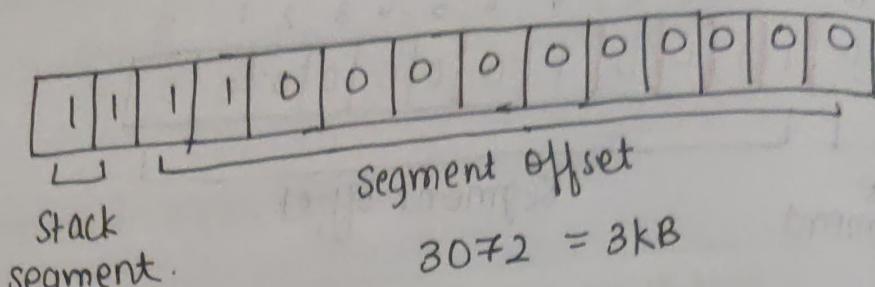
$$\text{physical add} = 1024 + 34 = \underline{\underline{35\text{KB}}}$$

$$VA = 15\text{KB} = 15,360$$



$$\text{virtual address} = 15\text{KB}$$

$$= 15 \times 1024 = 15360\text{ B.}$$



Stack have -ve growth

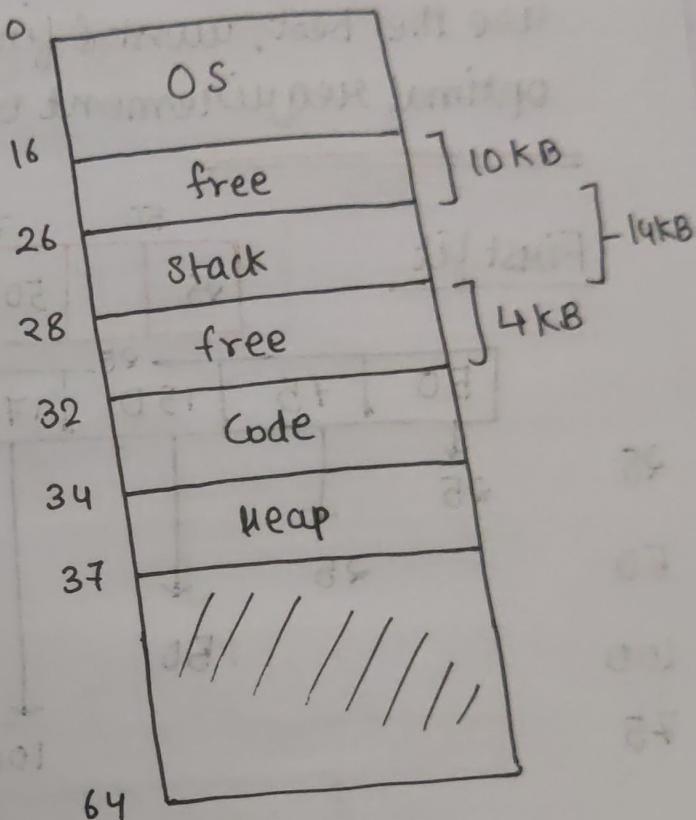
<u>segment</u>	<u>Base</u>	<u>size (max. 4KB) grow</u>
code	32	2
Heap	34	3
stack	28	2

The first thing we need a little extra hardware support. Instead of just base and bound's value the hardware also needs to know which way the segment grows (a bit that is set to 1 when the segment grows for the +ve direction and 0 for -ve) with the hardware understanding that segments can grow in the negative direction the hardware must now translate such virtual address slightly differently. For example assume virtual address to be 15KB, which should map to physical address 27KB so in this segment bit is 1,1 and offset is 3KB. A segment can be maximum size is 4KB

And thus the correct -ve offset is $3KB - 4KB = -1KB$.

$$\text{physical address} = 28KB - 1KB = 27KB.$$

which equals to -1KB we simply add the -ve offset to the base 28KB to arrive at the correct physical address 27KB.



Suppose 14KB, 10KB, 250, 500, 600, 360.

200, 50, 300, 100, 40. → chunks segments.

First fit: for 200 → 250 left out 50

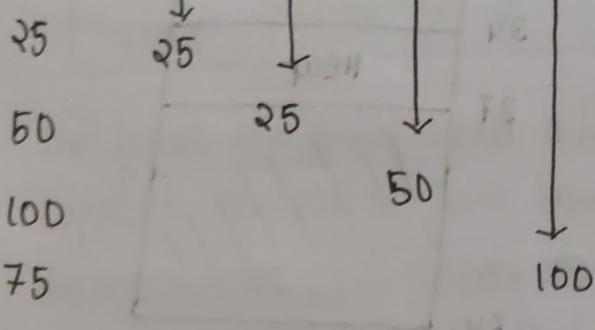
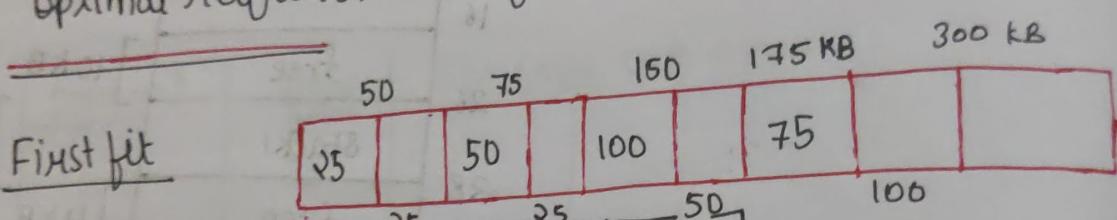
Q) Process Request are given as

25KB, 50KB, 100KB, 75KB

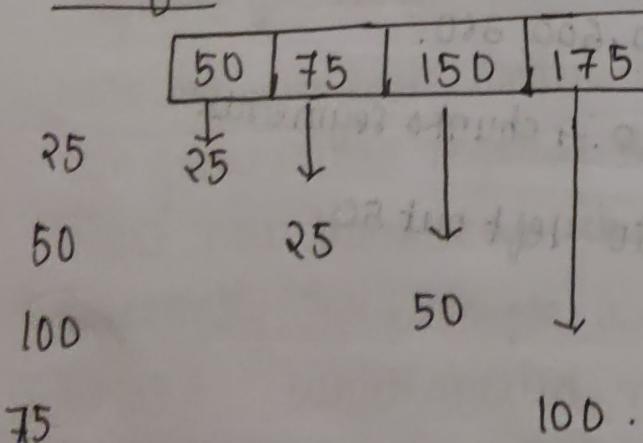
and free list of space is

50KB, 75KB, 150KB, 175KB, 300KB.

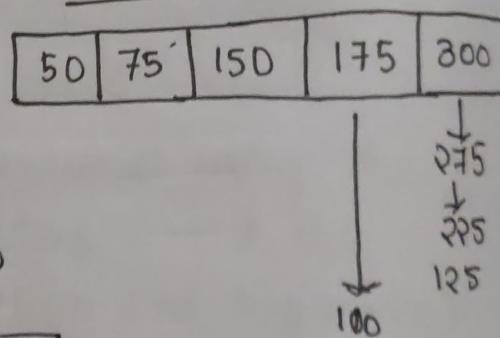
use the best, worst & first fit algorithm to find the optimal requirement of the process



Best fit



worst fit



Consider 5 memory partitions of size

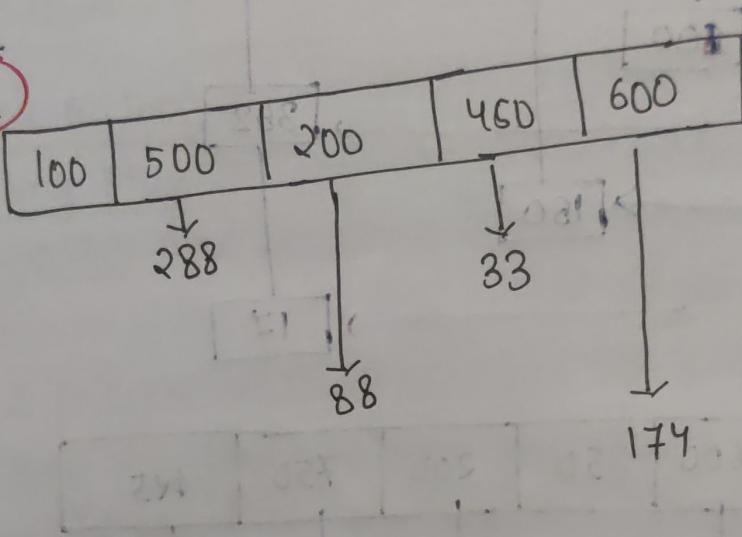
100 KB, 500 KB, 200 KB, 450 KB, 600 KB in same order

96 seq. of 4 q for blocks of size

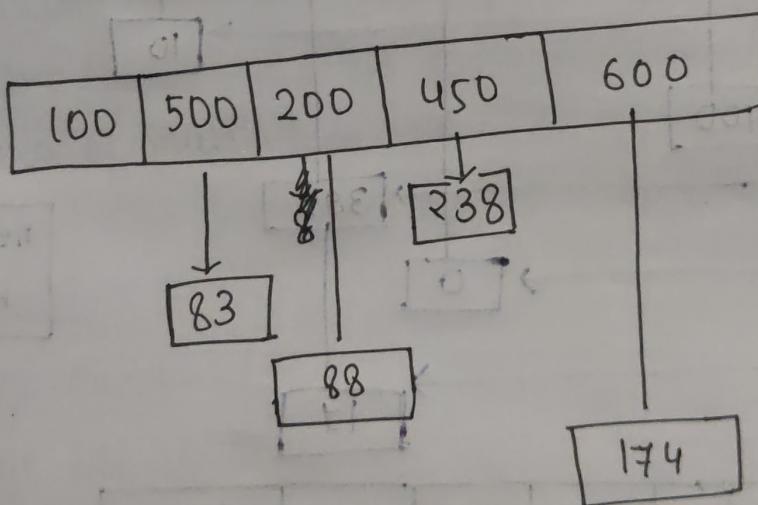
212, 417, 112 KB, 426 KB in same order come,

then which of the following algos makes the efficient use of memory.

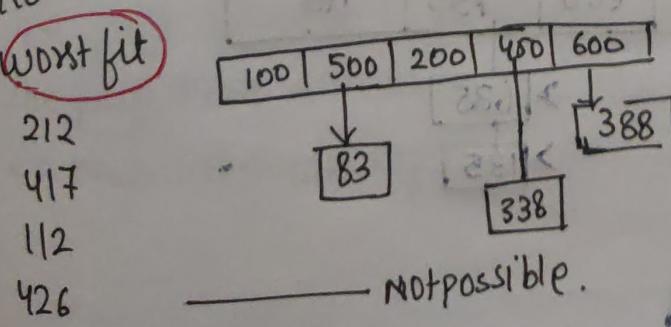
first fit



Best fit



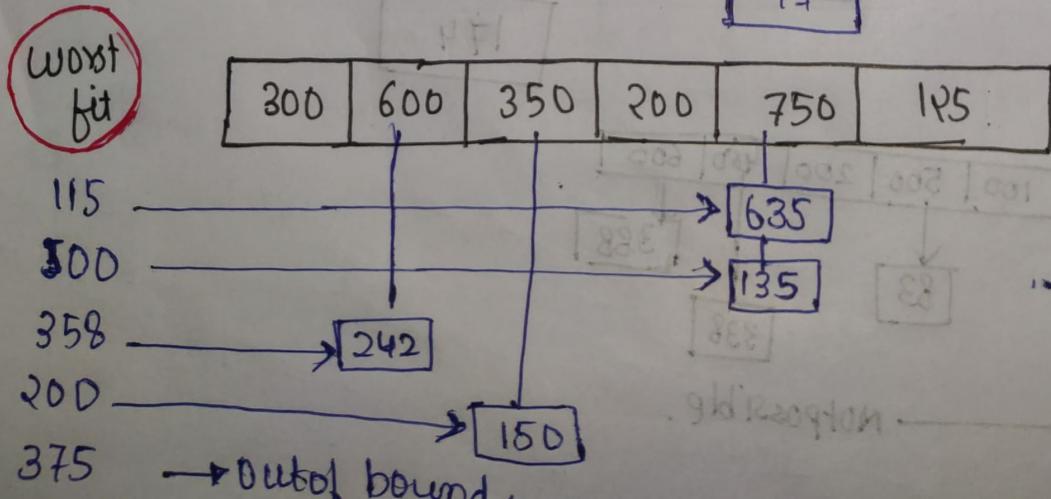
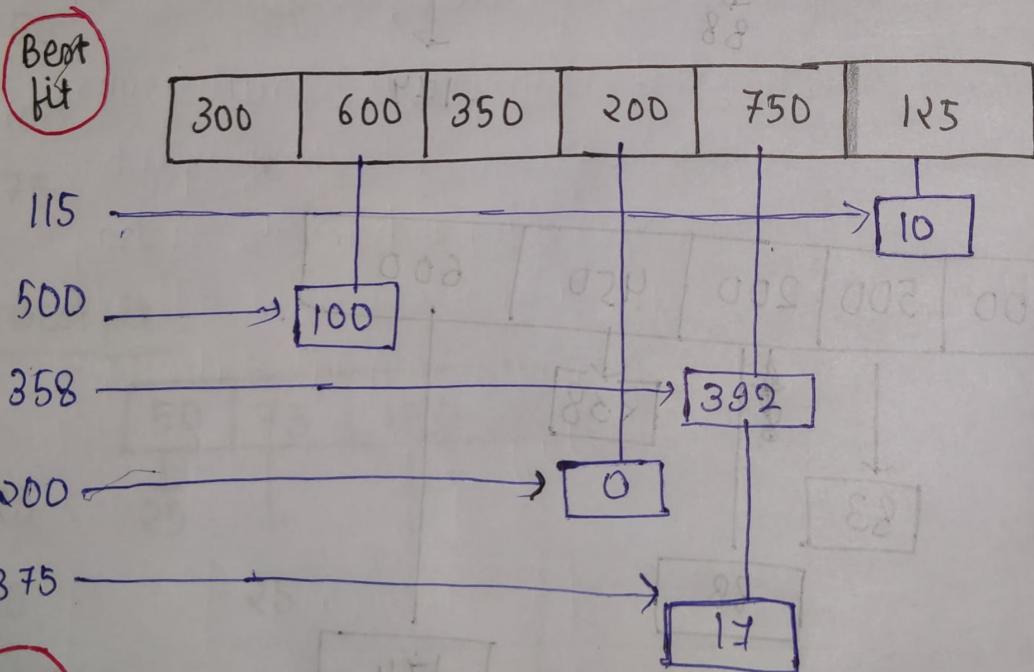
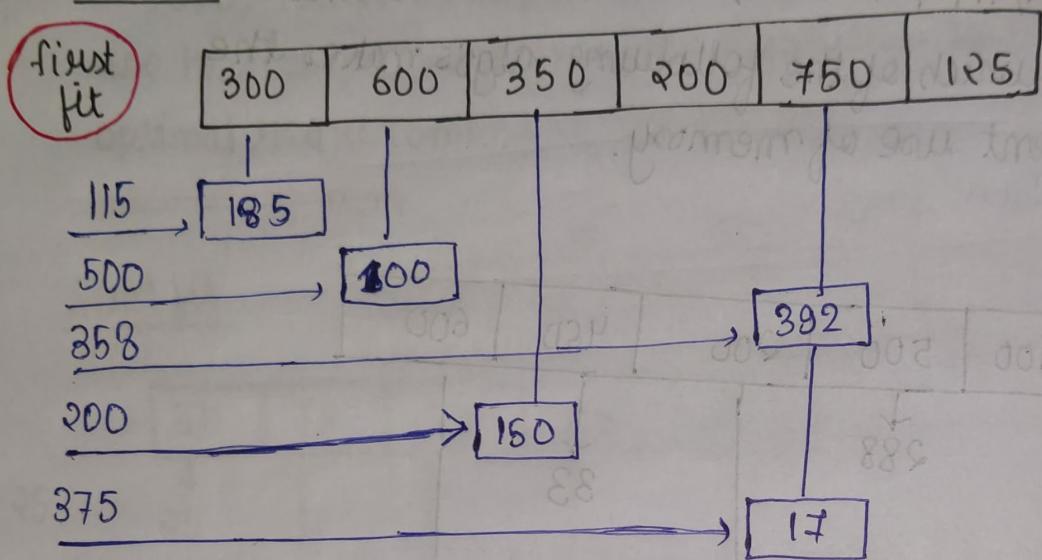
worst fit



NOT POSSIBLE.

Q) Given six memory partitions of size 300KB, 600KB, 350KB, 200KB, 750KB, 125KB in order.

How would the first fit, best fit and worst fit algo's place processes of size 115KB, 500KB, 358KB, 200KB, 375KB



Paging. [Because segmentation has problems in managing free space as memory becomes fragmented]

It is another memory management scheme that avoids external fragmentation and the need for compaction whereas segmentation does not. It also solves the considerable problem of fitting memory chunks of varying size onto the backing store.

Paging: → address space into fixed size units we call a page

Basic Method:

The basic method for implementing paging involves breaking physical memory into fixed size blocks called frames and breaking logical memory (address space) into blocks of the same size called pages. When a process is to be executed its pages are loaded into any available frames from some source (from store). The backing store is divided into fixed size blocks that are the same size as the memory frames or clusters of multiple frames. Every address generated by the CPU is divided into two parts:

- ↳ page number (P)
- ↳ page offset (d)

The page number is used as an index into a page table. The page table contains the base address of each page in physical memory.

$$\text{addr} = \text{PA} - \text{RA}$$

This base address is combined with the page offset to define the physical memory address that is sent to the memory unit.

The page size is defined by the hardware. The size of the page is a power of 2 varying b/w 512 Bytes and 1 MB per page depending on the computer architecture.

$$\text{page size} = ? \rightarrow 512 \text{ bytes} - 16 \text{B}$$

(a) what is the difference b/w 32 bit and 64 bit architecture?

The selection of power of 2 as a page size makes the translation of a logical address into a page number and a page offset. If the size of the logical address space is 2^m and a page size is 2^n , then the high order $m-n$ bits of a logical address represents the page number and the m lower order bits represents the offset.

Page number	offset
$m-n$	n

$$\text{logical address} = 16 \text{KB} \quad \text{page size} = 4 \text{KB}$$

$$= 2^4 \times 2^{10} \text{B}$$

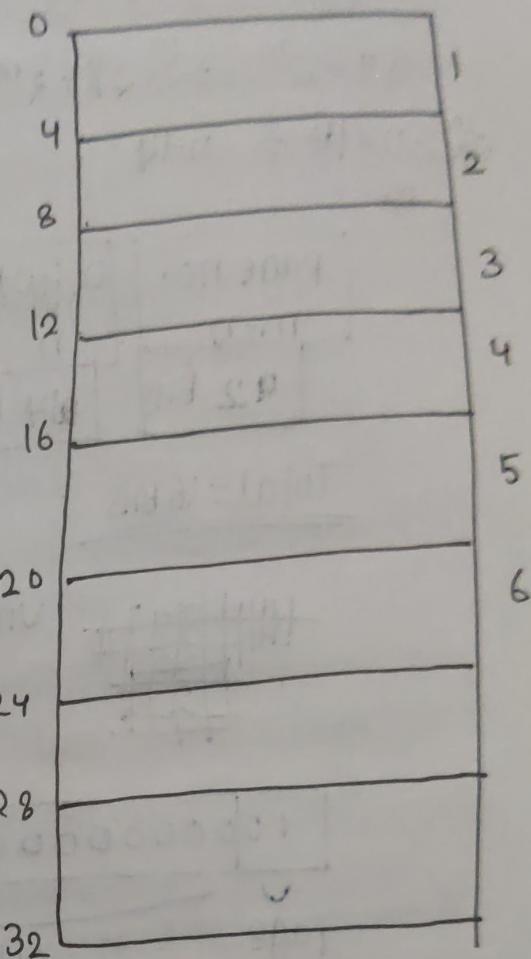
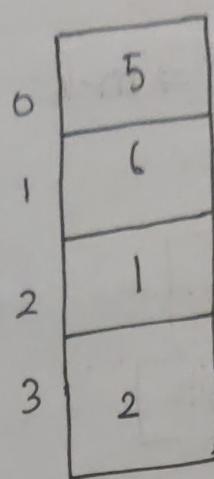
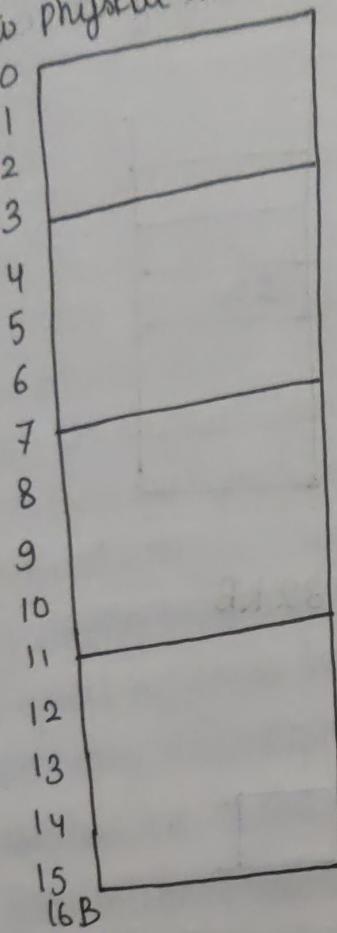
$$= 2^{14} \text{B}^m$$

$$= 2^2 \times 1,024 \text{B}$$

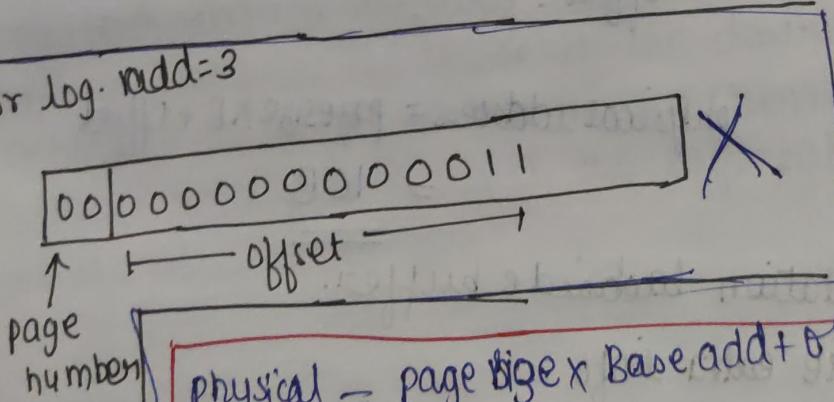
$$= 2^{12} \text{B}^n$$

$$\underline{m-n = 2 \text{ bits}}$$

Consider the memory of 16KB size 4KB each, $n=2$, $m=4$
 using a page size of 4 Bytes and physical memory of 32 bytes.
 So how the programmers view of memory can be mapped
 into physical memory.



for log. vadd = 3



$$\text{physical} = \text{page size} \times \text{Base add} + \text{offset}$$

$$4 \times 5 + 3 = \underline{\underline{23}}$$

Virtual address = 8 bit.

Page size = 16 Bytes

$$4 \cdot 2^4 = 16 \quad 4=n$$

Bits = 4

Total Pages = 16

32x1024

Q) logical address size = 64 kB

Page size = 16 kB

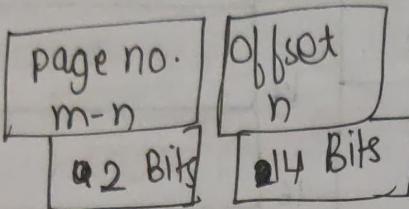
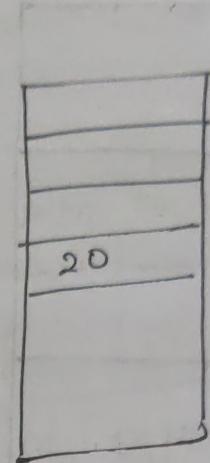
physical memory size = 128.

virtual address = 32 kB,

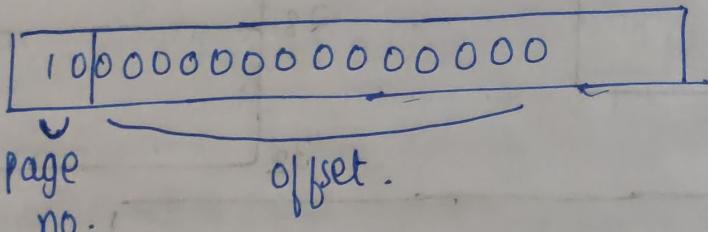
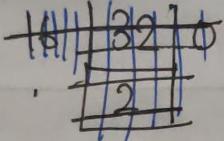
find out the physical
memory address

$$64 = 2^m \Rightarrow 2^10 \times 2^6 = 2^m \Rightarrow m = 16$$

$$2^{10} \times n = 16 \Rightarrow n = 4$$

Total = 16 Bits

virtual address = 32 kB

physical address = pagesize \times B + offset

$$= 16 \text{ B}$$

Interviewed
pair of views

TLB: Translation lookaside buffer.

- ↳ Share extra info - looking up
- ↳ And with help of cache it boost up the speed of the process.

TLB-hit
TLB-miss

[Best eg: Browser]

address = 10000000000000000000000000000000

00000000000000000000000000000000

00000000000000000000000000000000

00000000000000000000000000000000

00000000000000000000000000000000

Using paging as a core mechanism to support virtual memory can lead to high performance overheads by chopping the address space into small fix size units (pages), paging requires a large amount of mapping information. Because that mapping information is generally stored in physical memory, paging logically requires an extra memory lookup for each virtual address generated by the program. Going to memory for translation information before every instruction fetch or explicit load or store is prohibitively slow. And thus our problem is how to speed up memory translation.

By using hardware translation lookaside buffer (TLB) to speed address translation a TLB is part of the memory management unit and it is simply a hardware cache of popular virtual to physical address translation. (Address translation cache). Upon each virtual memory reference the hardware first checks the TLB to see if the desired translation is held there in, if so the translation is performed without having to consult the page table they may get the physical address.

TLB Algorithm

↳ Assumption:

- Simple linear page table:
- Hardware managed TLB

The following steps are as follows:

- (1) Extract the virtual page number from the virtual address and check if the TLB holds the translation for this VPN (virtual page no.)
- (2) If it does a TLB hit is success then we can now extract the page frame number from the relevant TLB entry, concatenate that onto the offset from the original virtual address, and form the desired physical address, and access memory, assuming protection checks do not fail.
- (3) If the CPU doesn't find the translation in the TLB then the hardware access the page table to find the translation and assuming that the virtual memory reference generated by the process is valid and accessible, updates the TLB with the translation.

Example: Accessing an Array.

We have an array of 10 four byte integers in memory starting at virtual address 100.

Assume further that we have a small 8-bit virtual address space, with 16 byte pages; Thus

	00	04	08	12	16
00					
01					
02					
03					
04	miss a[3]	a[0]	hit a[1]	a[2]	hit a[6]
05	miss a[7]	a[4]	a[5]	a[9]	
06		a[6]			
07					
08					
09					
10					
11					
12					
13					
14					
15					

Some TLB stores address space identifier in each TLB entry. And ASID uniquely identifies each process and is used to provide address space protection for that process. When the TLB attempts to resolve virtual page numbers it ensures that the ASID for the currently running process matches the ASID associated with the virtual page. If the ASID do not match the attempt is treated as TLB miss. In addition to providing address space protection and an ASID allows the TLB to contain entries for several different processes simultaneously. If the TLB does not support separate ASID's then every time a new page table is selected

, the TLB must be flushed/ erased to ensure that the next executing process doesn't use the wrong Translation information, otherwise the TLB could include old entries that contain valid virtual addresses but have incorrect/ invalid physical addresses left over from the previous process.

The percentage of times that the page no. of interest is found in the TLB is called Hit Ratio.

For eg: A 80% Hit ratio means that we find the desired page number in the TLB 80% of the time. If it takes 100ns to access memory then a mapped memory access takes 100ns when the page no. is in the TLB. If we failed to find the page no. in the TLB then, we must first access memory for the page table and frame number, and 100 nanoseconds, and then access the desired byte in memory (100 nanoseconds) for a total of 200ns.

To find the effective memory access time, we must weigh each case by its probability

V.IMP

$$\text{Effective access time} = P \times \text{hit memory time} + (1-P) \times \text{miss memory time}$$

80. ~~Q1~~ → Hit Ratio

$$P = \frac{80}{100} = \frac{80}{100} \times 100 + \frac{20}{100} \times 200 \\ = 120 \text{ ns}$$

Hit Ratio 99%.

$$P = \frac{99}{100} \times 100 + \frac{1}{100} \times 200 \\ P = 101$$