OS : Coordinates everything
   ↳ Middleman b/w Hardware and user

OS handles interactions with the disk and performs
   storage management.

```c
#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>
int main(){
    printf("This is parent. PID: %d\n",(int)getpid());
    int rc=fork();
    if(rc<0){
        printf("Fork failed!\n");
        exit(1);
    }
    else if(rc==0){
        printf("This is child. PID: %d\n",(int)getpid());
    }
    else{
        printf("This is parent my child is %d\n",rc);
    }
    printf("Who am I? PID:%d\n",(int)getpid());
    return 0;
}
```



```c
#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>
#include <sys/wait.h>
int main(){
    printf("This is parent. PID: %d\n",(int)getpid());
    int rc=fork();
    if(rc<0){
        printf("Fork failed!\n");
        exit(1);
    }
    else if(rc==0){
        printf("This is child. PID: %d\n",(int)getpid());
    }
    else{
        int rc_wait=wait(   );
        printf("This is parent my child is %d\n",rc);
    }
    printf("Who am I? PID:%d\n",(int)getpid());
    return 0;
}
```

```
    }
    printf("Who am I? PID:%d\n",(int)...
    return 0.
}
// wait() makes the parent wait for any one of the child and returns the pid of chi
// we can use wait() in a loop to wait for all children (loop will it doesnt return
// waitpid(child_pid) waits for the specified child
// wait(), waitpid() only works for parent-child pair where parent waits for the ch
```

OS → divides CPU work , makes it feel like any process has infinite access to CPU.
Process → Takes resources

• Memory :

   ┌ Static Memory → Variables [Stack]
   └ Dynamic Memory → malloc, calloc [Heap]



What Constitutes a Process?
- Unique Identifier (Process ID)  getpid
- Memory Image
  - Code and data (static)
  - Stack and Heap (Dynamic)
- CPU Context: Registers
  - Program Counter
  - Current Operands
  - Stack Pointer
- File Descriptors
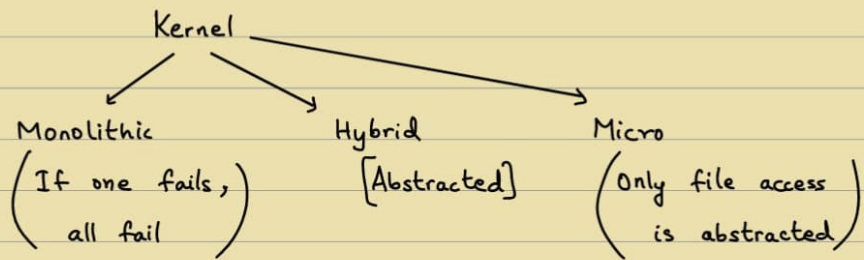  - Pointers to open files and devices

Memory Image of Process
Code
Data
Stack
Heap

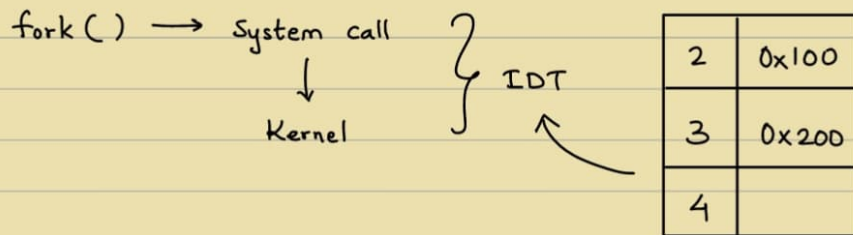Create memory image for the process
⇓
Allocate PC, SP

- **Kernel :** Handles direct instructions with hardware
  - ↳ Part of the code

Kernel
- Monolithic
  - ( If one fails, all fail )
- Hybrid
  - [Abstracted]
- Micro
  - ( Only file access is abstracted )

<u>Kernel</u> has stack
  - ↳ Process

- **API** ⟶ System calls (fork, exec, wait)
  - ↳ Resides in Kernel

- **Limited Direct Execution (LDE)** ⟶ CPU

fork ( ) ⟶ System call
              ↓
           Kernel

} IDT

| 2 | 0x100 |
|---|-------|
| 3 | 0x200 |
| 4 |       |

System call ⟹ Mode change ⟹ User mode to Kernel Mode
  - ↳ TRAP Instruction : Changes privilege

IDT ⟹ OS code gets executed
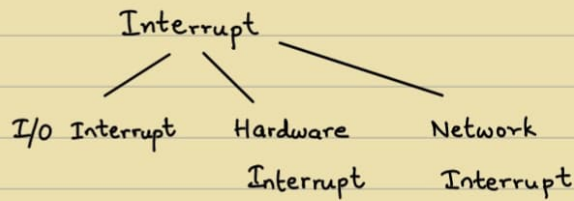  - ↳ Interrupt disrupt Table

Context ⟶ PC, SP ⟶ Saved

User stack ✗
<u>Kernel stack</u> ✓
  - ↳ Per process level
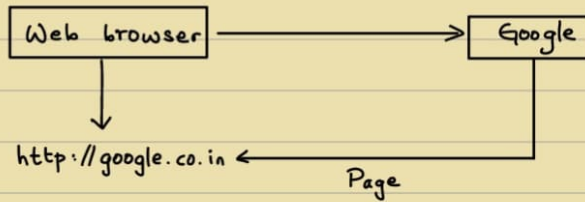
10

11  syscall ⟵ OS

12

Interrupt
```
        /        |        \
I/O Interrupt  Hardware   Network
               Interrupt   Interrupt
```

- Context Switch:
  How to decide which process to run next $\longrightarrow$ Scheduling

  In reality, Jobs can arrive at any time.

```
┌──────────────┐                    ┌─────────┐
│ Web browser  │ ─────────────────> │ Google  │
└──────────────┘                    └─────────┘
        │                                │
        v                                │
 http://google.co.in  <──────────────────┘
                         Page
```

Port $\rightarrow$ 80

| Process id | Port |
|------------|------|
|            |      |
|            |      |

$\leftarrow$ Mapping

- Socket $\longrightarrow$ Abstraction $\longrightarrow$ Allows communication
    $\hookrightarrow$ API

Stream $\longrightarrow$ continuous (TCP)
     $\hookrightarrow$ Mutual connection (strict) + Connection oriented
Datagram $\longrightarrow$ No memory + No persistent connection (Stateless connection)
     $\hookrightarrow$ UDP
     $\longrightarrow$ Unrealiable (by default)



| TCP vs UDP | |
|---|---|
| **TCP** | **UDP** |
| Connection Oriented | Not Connection Oriented |
| Reliability (order is maintained and retransmission) | No reliability at L4 |
| Higher overhead - reliability, error checking, etc | Low overhead |
| Flow control (based on network) | No implicit flow control |
| Error detection - retransmit erroneous packets | Has some error checking - Erroneous packets are discarded without notification |
| Congestion Control | No Congestion Control |
| Use cases: HTTP/HTTPS, File transfer, Mail | Use cases: Streaming data, VoIP, DNS queries, .. |

Flow control : Control speed of data transmission.

- UDP $\longrightarrow$ Unreliable at L4 ∵ Fire and Forget, i.e. responsible for giving data but not ensuring recieving.

  Headers differentiate TCP and UDP.
  
  $\hookrightarrow$ 32 bits $\searrow$ lightweight (8 bytes)

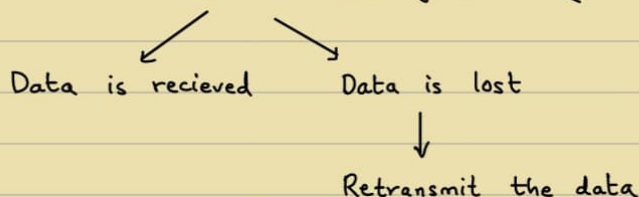  TCP $\longrightarrow$ 3-way Handshake (Agreement)
  L4 ensures packets getting delivered in-order.
  $\hookrightarrow$ Transmission Rate

- <u>TCP Header</u> :
  ① Sequence number : Allows sequencing of data
  ② Acknowledgement number : Sent back from reciever, to let the sender know that data was recieved and is waiting for next data.

  Data is sent $\longrightarrow$ Wait to recieve acknowledgement (OS starts a timer)

  Data is recieved    Data is lost

  Retransmit the data

  OS does not want to waste resources
  $\Rightarrow$ ∴ Starts a timer and if it elapses $\Rightarrow$ Timed Out

  Data can be lost even when acknowledgement no. > sequence no.

  Acknowledgement cache $\longrightarrow$ Prevents <u>Deadlock</u>
  
  $\swarrow$
  
  Sender sends & waits for response (ACK)
  Reciever sends back & still on the way

  " IIIT Timer $\longrightarrow$ 1½ hour delay "

  - RTT : Round Trip Time
    Sample RTT : How much time
    Estimated RTT : Keep updating estimation on every sample
    $(1-\alpha)$ ERRT + $(\alpha)$ SRTT
    $\hookrightarrow$ If very high, ∴ Works only on current & not past
    Default : 0.125

DeviationRTT : Worst case scenario , can happen but rarely

$$DRTT = (1-\beta) DRTT + \beta \times | SRTT - ERTT |$$

↳ Weighted moving range          $(\beta : 0.75)$

Timeout Interval = ERTT + (4 × DRTT)

↳
$\begin{cases} 3 \longrightarrow 85\% \text{ confidence} \\ 2 \longrightarrow 65\% \text{ confidence} \\ 4 \longrightarrow 99\% \text{ confidence} \end{cases}$

[Based on experiments]

Q) Check acknowledgement constantly for every packet ?

    Disadvantage (if check) ⟹ High overhead & Latency

    Solution ⟹ Delayed Acknowledgement

            ↳ Wait after sending multiple packets
                  and reciever sends ACK

To prevent deadlocks :

    Timer on the reciever end

③ Window : Flow Control

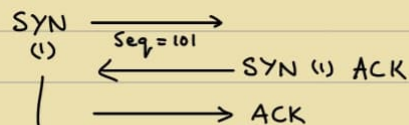    Reciever has a limit to recieve data & then waiting time.

    Send a 0 byte data to know whether reciever can recieve data.

Reciever & Sender ⟶ Both send data , but have seperate values of ACK.

Sequence number ⟶ Any random number within the bounds.

      ↳ Decided after agreement

3 - way handshake ⟶ " Can we start connecting ? "

          SYN    ⟶
          (1)   Seq = 101
             ⟵ SYN (1) ACK
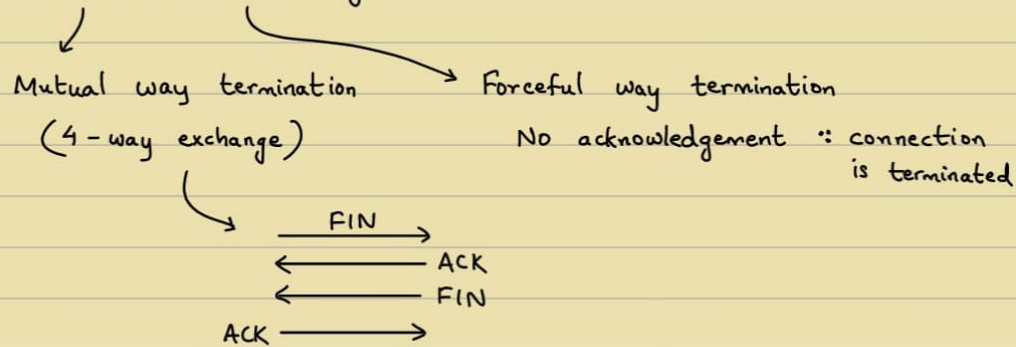             ⟶ ACK

        ↳ 1 : Setting the SYN bit , i.e. sending the acknowledgement
                Agreement on Seq. no. from both sides

        ACK → 1 : Accepts the mutual seq. no.

        SYN ACK : Mutual Agreement & that both recieve data

    SYN : 1 ⟹ Connection establishes and sender initiates sending

Closing connection : FIN and RST flags

Mutual way termination       → Forceful way termination
(4 - way exchange)           No acknowledgement ∴ connection
                                                        is terminated

                  FIN  →
                 ←     ACK
                 ←     FIN
       ACK  →

Note : Size of TCP Header : 30 bits to 60 bits (Option Speed)

P bit : [Push bit]
     Send data ⟹ Reciever gets and puts in buffer (Generally)
     P = 1 ⟹ Data is directly sent to process , No buffer used

U bit : [Urgent bit]
     Some packets sent need to be processed urgently and rest can be
     processed later

- Memory :
     Every Process requires memory.
     Memory Virtualization ⟶ Every process feels like its getting entire memory.

         Some of the RAM is used up by OS , code and data . Rest of
         the memory is usable RAM , i.e. our code and data.

     Multiple Processes running at the same time ⟹ Fast

     Ways to achieve process virtualisation :
     (1) Partitioning OS into slots
         Disadvantages : - Other slots data can be accessed.
                     - Each process may not be fixed amount of data.

{ Saturday : 8:30 Tutorial ⟶ Class }

Flow control : Just b/w sender and reciever (L4)

Congestion control : C-bit , In the network layer

- **Memory virtualization** :

    e.g. Supermarket vs e-shopping

    Every process ⟵ Address space

    ⟶ Virtual Address Space , mapped to Physical memory
    (VA)

    $*p = 3 \implies \&p$ ⟵ Virtual memory

    Note : Only 1 Physical Memory ⟵ Physical Address (PA)

    Processes might not always stay in Physical memory , it can be
    in disk as well from overflowing . (Note : Swap Space)

    * Translation : Memory Management Unit (MMU)
        ⟶ Fast
        ⟶ Conversion + Fetch

    - Goals of Virtualization :
        (1) Transparancy ⟵ Each process gets an illusion of infinite space
        (2) Efficiency ⟵ Use hardware
        (3) Protection

    - For Process Virtualization :
        Mechanism : Limited Directed Execution (LDE)
        Policies : Scheduling Algorithms

    - For Memory Virtualization :
        Assumptions :
        (1) User address space is contiguous in Physical Memory.
            Process ⟶ Block
        (2) Size of Address space is too big, less than size of Physical memory
            Process ⟶ No overflow
        (3) Every Address space is of equal size.

    Memory Virtualisation :
        Q) How do I divide Physical Address space.?
        g) How do I map Virtual Address to Physical Address ?

Reference the variable ← Stack

Note : Base and Bound approach :

→ Values are stored in MMU

Then, Add Base Address

VA → Process starts at 0
But PA : Process = 0 + Base value
= VA + Base

Accessing after bound → Fault : Address out-of-bound error
Allocation : Dynamic ⟹ Dynamic Relocation
VA to PA : Address Translation

→ MMU not OS

\* MMU → @ Context-Switch
→ Has only 1 pair of Base and Bound

\* Segmentation : (Memory Management System)
Code, Stack & Heap — Generalised Base and Bounds
∀ Segment : Has a base and bound, Not Fixed
→ Dividing PA space

Segment Registers in MMU.
First → Identify which segment

\* Address Translation :
Stack → goes up ∴ Subtract
Real Address (PA) = VA + base address

Code : No offset
Heap : Offset

Calculation

VA ⟶ Every Process

Map : VA ⟶ PA : Base and Bounds ⟶ Disadvantages :
        ↘                         (1) Fragmentation
        0 to Max               (2) Unused Memory

- Identify segment :
  - (1) Implicit :
          Stack
  - (2) Explicit : (VA → 14 bits)
          First 2 bits ⟶ Segment (12-13)
                   Rest ⟶ offset   (0-11)

          00 : Code
          01 : Heap
          11 : Stack

    ∴ Bits are unused ⇒ Code and Heap have same bit in some OS.

- Simple Address Translation :

      Many OS : Segmentation ✗
                  Paging ✓

     PA : Base + Offset
              ↘
                  Data − Max stack data

-               Segmentation
                    /          \
        Fine - grained          Course - grained
        ├ Small size segments    ├ Large size segments
        ├ High overhead
  When different sized segments , Disadvantage :
  External Fragmentation ← Extra memory is available
                          But not contiguous , i.e. has holes

      Defragmentation : High overhead

      Note : Free space management algorithm , Closest Fit, First fit are algorithms

Fixed size segments $\Rightarrow$ Code, Heap & Stack ×

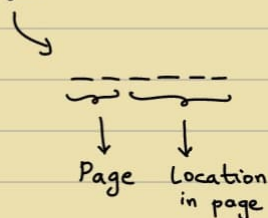$\quad\searrow$ Disadvantage : Internal Fragmentation

$\qquad\qquad\searrow$ Better than external fragmentation

Page : Fixed size segment in a process
Paging : VA and PA are both divided

$\quad\therefore$ If 4 pages $\rightarrow 2^{\textcircled{2}}$ : 2 bits
$\qquad$ 64 bytes VA $\rightarrow 2^{\textcircled{6}}$ : 6 bits

$\qquad\qquad\qquad\qquad\qquad\qquad$ Page $\quad$ Location
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ in page

VPN : Virtual Page Number $\left.\right\}\leftarrow$ Every page has a number
PFN : Page Frame Number

In PA : Page Frame
In VA : Page

Advantage :
$\quad$ (1) Flexibility : OS finds a free page
$\quad$ (2) Simplicity : ∀ Process , ∃ Page table

$\qquad\qquad\qquad\qquad\searrow$ Mapping b/w VA and PA
$\qquad\qquad\qquad\qquad\qquad$ Array can be used , e.g. VPN[i] = PFN$_j$

Note : Multiple VPNs can be mapped to a PFN (not same process)

$\quad$ If Mem$_{Process}$ > Mem$_{RAM}$ $\Rightarrow$ Swap in & Swap out from Hard disk

– Process of Translation :

$\qquad$ VPN $\qquad\qquad$ Offset
$\qquad$ VA$_5$ VA$_4$ $\quad$ VA$_3$ VA$_2$ VA$_1$ VA$_0$

$\quad$ e.g. mov 21 %eax
$\qquad (21)_{10} = (10101)_2 = (010101)_2$

01 : VA space ⎱
0101 : offset  ⎰ → & if $VPN_1$ = $PFN_7$ (from page table)

Then  $\underbrace{010101}_{VA}$ ⟹ $\underbrace{110101}_{PA}$

$(110101)_2 = (117)_{10}$

Page Table Base Register ⟶ store base of the address of the page table.

e.g. 32 bit address space with 4 KB pages
   no. of bits for offset?

   4 KB = $2^{12}$ ⟹ 12 bits to represent
            32 - 12 = 20 bits for mapping, i.e. VPN
                ∴ $2^{20}$ mappings per process
      Each mapping ⟶ 4 bytes
            ∴ $2^{20} \times 4$ = 4 MB per process per page

   If 100 processes ⟹ $\underbrace{400 \text{ MB}}$ for address translation
                        ⤵
                         → Page Table
            ∴ Use pages to store page tables

- <u>Page Table</u> :
      Page Table Entry (PTE)
         VPN , PFN , V, P, $\underbrace{\text{Present}}$ , $\underbrace{\text{Dirty}}$ , Ref
                    ↓    ↓      ↓           ↓              ↘ If page is recently used
              Valid Bit                               ↘ If page is modified (from memory) after
                                          ↘ If page in Physical memory or secondary
                         ↘ Protection Bit

      ∴ Disadvantage : Fetch twice for every $\underbrace{\text{translation}}$
                    (Extract memory)                        ⤵
- Efficient Translation : Caching                          ↘ PTBR
   Translation Lookaside Buffer (TLB) – Register              ↓
   ⤷ Address Translation cache                            Page Table
      ↘ In MMU                                                ↓
                                                            PTE
                                                             ↓
                                                           Find

TLB → Hit rate of cache ↑

Cache : VPN ⟷ PFN

Size of cache ↑ ⟹ Hit rate ↑
Size of page ↑ ⟹ Hit rate ↑

- Spatial locality : Space dimension, i.e. same space accessed
- Temporal locality : Time dimension, i.e. recently accessed.

• Hardware TLB handling  — Hardware goes through the page in $O(1)$.
  Software TLB handling  — Hardware raises exception (TRAP)
                          Context Switch from Hardware to Software

@ TLB miss handlers  ⟶ stored in OS Kernel, avoiding translation
                          (Physical memory)

TLB ⟹ Fully Associative cache
        (No direct mapping, random)

• TLB :
  - Valid bit :
      Translation is valid or not ⟶ Page Table Entry case
      Process A to B ⟹ B should not use A's mapping ⟶ TLB case
                ∴ A ⟶ invalid

  - ASID :
      Every process has a Address Space Identifier.
      Mapping : Process id ⟷ ASID
                        ↳ Small (∵ Cache)
                    ↳ Disadvantage : Large amount of bits (32 to 64)

• (i) Segmentation + Paging  }⟶ Reduce Page Table size
  (ii) Encoding

        Page size ↑ ⟹ Mapping ↓ but Internal Fragmentation,
                          lot of empty sizes will be left

        (i) Segmentation + Paging ⟶ Getting best of both ?
              ↳ External      ↳ Internal      } Issues
                Fragmentation    Fragmentation

Page Table ⟶ stored in RAM
⟶ divided into pages
↳ Stored using Pages

Page Table contains mappings, entries.
Each mapping is not important, only valid entries are
Invalid entries ⟶ don't store, store only when it becomes valid.
Valid entries ⟶ store
∴ Multi - Page Table (Tree -like)
i.e. Meta-index, multiple pointers heirarchy.

✳
[ Explore 2-level Page Table
  Study 3-level Page Table ]
✳

Page Directory ⟶ Data structure for multi-level page tables.

PDBR ⟶ less space (∵ Only valid PTEs are stored)

Disadvantage : Time overhead

12/9
- Inverted Page Table : 1 page table per process
  disadvantage : Linear scan is expensive

- OS creates Memory Heirarchy.
  Physical Memory ⟶ Faster access
                    Less capacity
  Disk ⟶ Slower access
         More capacity

On-demand Pages
- Swap Space :
  ↳ Allocated in disk
                divided into blocks ≈ Page size

- Present bit :
  1 : page is in the page table
  0 : page is in the swap space

- Valid bit = 0 : Page is not in Physical memory ⟹ In disk

    AKA Page Fault ⟶ Handler

Process moved to
block state

⟶ What to swap, when & from where
AKA Page Replacement

- AMAT
- Optimal Replacement Policy (Theoretical Optimal)

    Belady Replacement Policy (But Not Practical)

      ↳ Possible future access

Reduce any policy to Belady for comparison

(i) FIFO ⟶ First in First Out

    Cache size ↑ ⟹ no. of hits ↑   (Belady's Anamoly)

    ⟶ Not always, depends on ordering of stream of data.

(ii) LRU ⟶ Least Recently used
    - Valid bit = 1
    - Access bit = 1

      ↓ Recently modified

    Access bit = 0 ← Evict a page

- Thrashing ⟶ Excessive Page fault Handlers

$$\left\{ \begin{array}{l} \text{Sat} \rightarrow \text{Tut (IMP)} \begin{bmatrix} \text{Bonus} \\ \text{Quiz} \end{bmatrix} \\ \text{Tues} \rightarrow \text{Tut (Attendance) [MP2]} \\ \text{Fri} \rightarrow \text{None} \end{array} \right\}$$

SMTP uses TCP

$\left(\text{Spotify} - \text{Web series}\right)$

Install Apache / ngnix webserver

→ <u>Peer 2 Peer Network</u> :

Spotify : Midserver + Peer2Peer
(Client server)

→ Application <u>level</u> <u>Protocol</u> :

HTTP : 80

HTTPS : 443 (HTTP over secure network)

∀ connection → 3 way Handshake (∵ <u>HTTP</u> runs on TCP)

↳ Stateless

• 1.0 HTTP ← For fetching every connection, connection is established.

1.1 HTTP ← Maintain connection to fetch all objects and then terminates.

2.0 HTTP & 3.0 HTTP

→ <u>Cookies</u> :

∵ HTTP is stateless, Data stored in browser cache.

Metadata is stored in server.

Web cache :

Browser cache ⎫
ISP cache ⎬ If recently changed, use cache
Regional cache ⎭
Main server

• Traditionally,

Client ⟷ Server : Takes time
(Delay)

∴ Add cache at server.

Content distribution Networks (CDN) :

↳ Servers distributed among diff. locations.

(Logics : Potential viewership + location)

Enter Deep : Build smaller clusters in more sites
Bring Home : Build large cluster in lesser sites


Cache does not contain originals
↳ Conditional GET : Data from server would only be header, no body


HTTP 2.0 → Current
HTTP 3.0 → Future, but has support now.


→ <u>Domain Name System</u> (DNS) :
Translation b/w name and I/P address. Allows aliasing.
Uses UDP.
- UDP call to get I/P address & HTTP call happens, i.e. then client uses HTTP to send request to server.
- DNS record
- BIND → Turns machine to DNS Server → Port 53
i.e. machine acts like a directory, returns IP address
- Any request (HTTP, SMTP, etc) sent translates the hostname to IP address using DNS


- 1. <u>Root DNS</u>                    → Report to IANA
        ↳     12 different Organizations are responsible for DNS creation, etc
2. Top-level Domain (TLD) : .com, .org, .ai domains
3. Authoritative Domain


- How it works : (Recursive Process), When you send a request to a host,
1. Request is sent through Port 53
2. Local DNS (ISP Provides) ← Cache kind-of
        If mapping exists, fetch. Else goes to next level (Recursion)
3. Root DNS (if it has, it will give back, else send to next level)
4. TLD DNS (if it has, it will give back, else send to next level)
5. Authoritative DNS (Then, get I/P address and return to client)
6. Then, you send a HTTP request to the server.


- DNS servers stores Resource Records (RR) : (name, value, type, ttl)
                                                        ↙           ↳ Each DNS record
                                                    tuple


ttl : time after which it needs to be updated from the root DNS.
    ↳ Time To Live

(i) A type DNS record : (Authoritative)
   Name : Hostname
   Value : I/P address
   $\begin{pmatrix} AAA \longrightarrow IPV6 \\ A \longrightarrow IPV4 \end{pmatrix}$

(ii) CNAME type $\longrightarrow$ Mapping one domain to another domain
        $\hookrightarrow$ Canonical name record (alias)

(iii) NS $\longrightarrow$ Named Server
        Main server itself (Authoritative)

(iv) MX type $\longrightarrow$ Mail
            (mail.outlook.com)

27/9

- Concurrency :
  Multiple execution points sharing the same memory space
  Same process — Multiple threads sharing same memory space.
          - The threads synchronize with each other.
          - Should not overlap ideally.


      Thread : Smallest unit of execution.
              Seperate P.C
              Seperate stack of local variables
              Creates child processes
              Same address space / memory


      Concurrency in Single core machine $\longrightarrow$ Context switching

        fork $\longrightarrow$ child is created, PID diff. from that of parent.
                Memory copy is created, code is shared.

        exec $\longrightarrow$ Child is created with same PID.
                Replace the memory, $\therefore$ diff. memory
                Inter-process communication $\longrightarrow$ reading / writing from shared files

            Multiprocessing : Multiple processes executing in diff. cores of the CPU
                            in the same time on the same data.

Scheduling → Thread level

Not Process level

- Kernel Space → Kernel level threads → Scheduled by Kernel
  User level threads → Scheduled by user, i.e. user level libraries ← pthreads

- Concurrency : Interleaving
  - Context - switching in case of single core
  - Dealing with lot of things at once.

  Parallelism : At same point of time → Multiple execution
           ∴ Multiple cores are required.
  - Doing lot of things at once
  - Subset of concurrency

- TCB : For every thread, ∃ TCB
  Kernel level threads : Scheduled by OS, Handles system calls.
  User threads : What process a developer writes, e.g. processing a file.

- Start routine



```c
intro > C concurrency_sample_alpha.c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <pthread.h>
4
5   int counter = 0;
6   int max_index;
7
8   void *worker_thread(void *arg)
9   {
10      printf("%s\n", (char *) arg);
11      return NULL;
12  }
13
14  int main (int argc, char *argv[])
15  {
16      pthread_t thread_p1, thread_p2;
17      printf("Starting the threading demo\n");
18      pthread_create(&thread_p1, NULL, worker_thread,"thre
19      pthread_create(&thread_p2, NULL, worker_thread,"thre
20      pthread_join(thread_p1, NULL);
21      pthread_join(thread_p2,NULL);
22      printf("end\n");
23      return 0;
24  }
```

```
thread 1
thread 2
end
karthikvaidhyanathan@MacBook-Pro-73 Intro % ./a.out
Starting the threading demo
thread 1
thread 2
end
karthikvaidhyanathan@MacBook-Pro-73 Intro % ./a.out
Starting the threading demo
thread 2
thread 1
end
karthikvaidhyanathan@MacBook-Pro-73 Intro %
```

Non - deterministic + No wait()
Depends on scheduler

Reason : Race condition
           ↳ access b/w shared variables
      No synchronisation