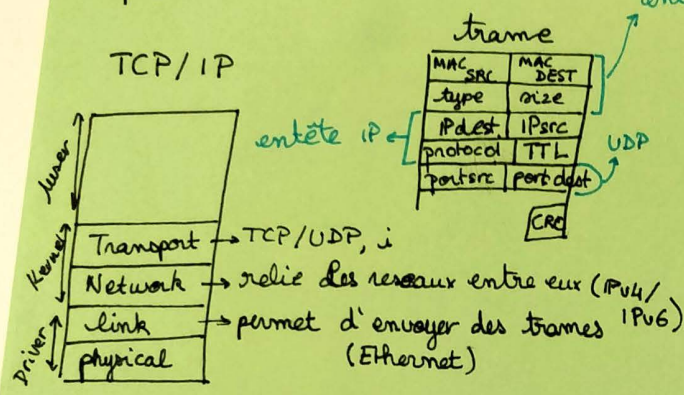


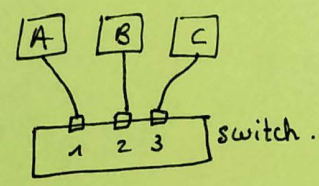
Programme :

- Process / Thread / Task
- Scheduling
- Mechanisme de synchro
- Réseau
- file system.

Reprendre les cours de Réseau.



hub : passif → broadcast.
capacité du réseau limitée car 1 seule personne peut parler en m^{ême} temps.
éviter domaine de collision.
Ce dernier doit être diminué au max.



Le switch tient une table d'apprentissage.

MAC	PORT
A	1
B	2

Si A envoie une trame le switch voit que A est sur le port 1.
Donc il le note dans la table.
Lorsque B répond le switch note B sur le port 2.

MTU: combien de temps je peux parler. (en mb de bytes).

Network:

Permet d'uniformiser les adresses. = adresses IP : XX.XX.XX.XX.

Routeur fait le lien entre 2 réseaux.

Netmask : met les bits du prefix à 1 les autres à 0. Cela permet de savoir qui est sur le même réseau.

Revoir les tables de routage.

→ table d'apprentissage entre routeur

On met la MAC adresse du routeur dans le paquet. Puis le routeur va remplacer sa MAC adresse par celle du prochain routeur.

Des packets sont fragmentés afin de s'adapter aux différentes MTU.

TTL: valeur qui est decrementée de 1 à chaque routeur, le packet est perdu lorsque la TTL atteint 0. exprimé en sec (car avant 1sec = 1routeur). Mais du coup permet de définir un timeout pour attendre le packet.

Transport:

DP: permet de dialoguer. DNS: on envoie un nom, ça renvoie une IP.
quette AXFR: permet de demander un nom.

TCP: permet de mettre les paquets ds l'ordre.

Reprendre le cours de Chewie s/ les ACK, Handshake.

Socket: API, implem de Burckley.

une
call initial: $\text{fd socket (INET, STREAM ou GRAM)}$
protocol \leftarrow type de connexion

initie la connexion connect (fd, sock addr) \rightarrow 3 way handshake.

s'accroche à un port bind (fd, sockaddr).
listen (fd, —)

send () } qd on sait a qui on
recv () } parle
sinon - sendto et recvfrom.

fd = accept (fd, &sockaddr)

Process vs Thread

Concurrence entre les process. plus file d'execut° ds le m process.

Process = State + file d'execut° (= thread \neq) if partagent les signaux, file descript° et l'adress space. \rightarrow bout de memoire alloué.

On peut avoir plus fois le m programme qui execute en // grâce à un fork.

Pour communiquer entre 2 process on utilise un pipe. Mais pas pratique car doit être prévu à la création. Du coup autre moyen: sockets. \neq types: \rightarrow INET

= ipv4, tcp/ip \rightarrow UNIX = pr la communicat° locale. On s'y connecte grace à connect (ip + port). Vu qu'on est en local on a un fichier special de type socket. du coup on donne à connect le chemin du fichier. Les sockets UNIX permettent de mettre des permissions s/ le fichier.

Pour écrire s/ l'ecran, on se connecte s/ le ^(serveur \rightarrow Xorg) X. grapha et on lui demande d'afficher le nécessaire. Le serveur grapha est lance en mode user et le gestionnaire de login permet de recuperer les fd necessaires (ds lesquels la data à afficher est contenue). Le gestionnaire de login / logind permet de creer les sessions.

Une autre methode:

Memorie (SHM) permet de partager en memoire via un fichier.
mkfifo.

Pour executer 2 \neq thread en m temps il faut 2 stacks \neq et un set de registre \neq , on partage les variables globales et une zone de memoire mappable.

Du coup pour lancer un thread on alloue une stack et un nouveau set de registre.

Les threads partagent les registres de il faut enregistrer les valeurs des registres à un autre.

AGAIN → bloque si jamais le fichier n'est pas prêt à être lu.
 fork est remplacé par clone qui permet de préciser les paramètres de la duplicat°.
 task → permet une meilleure relat° avec le kernel mais passage de l'un à l'autre + lent.
 pthread → implémenté ds la lib C. permet de faciliter l'appel à clone + d'autres détails.

SY52 (3)

8452

(2)

Scheduling

Rappel : time sharing = exécuter plus tâches en m temps.

Tache = temps d'exécuto + deadline + dependance.
(pas forcément lié au système).

batch : 1 programme à la fois qui tourne jusqu'à la fin.

→ Turnaround time = temps entre la ^(début?) fin de l'exécuto et l'obtent° du résultat.

→ Throughput (débit) = nb de tâche à la (heure/journée, ...).

Le but est de minimiser le turnaround et de maximiser le débit.

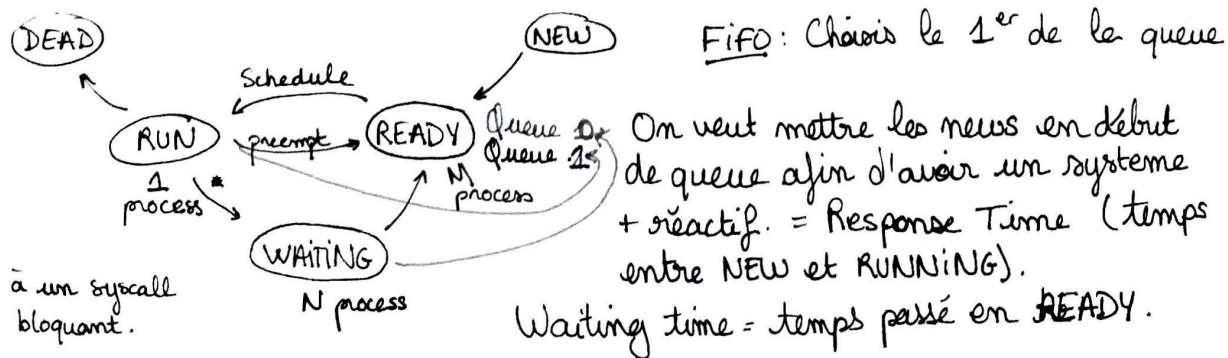
Pour ~~maximiser~~ minimiser le turnaround on fait du FCFS (First Come, First Served).

Pour maximiser le débit on fait du SJF (Shortest Job First).

Time Sharing

Il n'y a plus de temps d'exécution, ni de deadline, ni de dépendances.

d'algo qui détermine qui passe après (scheduler) doivent s'exécuter en temps constant sinon il prennent du temps sur celui d'exécuto. Starvation = tt les process doivent tourner. Fair = pas de favoritisme entre les process (donner le m temps à chacun).



Round Robin: ajoute des interruptions lié à un timer, qui peut envoyer un process de RUN à READY. On choisit un quantum de temps maximum pr le RUN. pb: si le process d'avant est passé en waiting on tourne sur la fin de sa clock. (la clock est longue à reset).

Système en temps réel : connaît le plan d'exécuto à l'avance. ex s/ un avion.

2 types de programmes :
 → Interactif : utilise des syscalls ex: shell, http, init
 → Non Interactif : ex: compilateur, program qui calcule.

Bottleneck (goulet d'étranglement) = ralentit le flux.

↳ CPU-Bound, Mem-Bound, IO-Bound (ex: téléchargement) → Svt pr les interactif.
(cpu lent). (ex: accès trop fréquent à la mémoire) tendance à ne pas finir

↳ Svt pr les non-interactif, tendance à ~~ne pas~~ finir le quantum le quantum.

On va prioriser les process interactifs, pour les repérer on compte le nb de fois qu'il passe en waiting.

Pour les non-interactif il est + rapide de faire tourner le process sur une durée de $2n$ que sur 2 durées de n car on a pas à changer les caches.

types de caches = L_1, L_2, TLB, \dots L_2 = accès mémoire L_1 = Core Cache.

de cache est utilisé lorsqu'on charge + de mémoire que demandé. ex si on travaille s/ un tableau il est plus rapide de charger le tableau lorsqu'on traite la 1^{ère} case que si on va chercher cases par case en mémoire.

Dans la pipeline d'exécution on a un cache L_1 au niveau du fetch pr charger les instruct° et un au niveau de memory où sera chargé la data. C'est un L_2 qui coordonne les deux L_1 . de TLB permet de sauvegarder les permissions et accès à la page. Il est vidé lors d'un changement de process.

On répare la queue en 2, une pour les CPU-Bound et une pour IO-Bound

On prend n queues allant de 0 à $n-1$, la queue 0 est la + prioritaire.

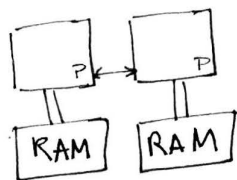
Une coup les programmes vont aller ds la queue correspondant à leur priorité.

RUN → READY = ↘ Priorité; RUN → WAITING = ↗ Priorité; READY → RUNNING = ↗ Priorité.
Plus la priorité est basse, + on alloue un quantum long.

• Ticket Scheduling: les programmes ont des tickets et ils utilisent 1 ticket par round on met une priorité à ceux qui ont le - de ticket

CFS (= Completely Fair Scheduling) = on unifie le temps d'exécution des programmes.

NUMA: Système à deux processeurs.



↔ = canal de dialogue

Si on veut accéder à qq chose qui est sur la rame voisine on met un peu + de temps.

File System.

Abstraction Storage

→ files

→ hiérarchie de nom

But: partager et nommer des b informations.

Block Device (Disk)

→ lire et écrire des blocks. (= suite de data continue) s'apparente à un buffer de taille fixe.

Sizes: - 2048 (cd-rom)

- 4K (disque moderne).

SYS2 (3)

You
fro
the

On utilise les BTree pour organiser les blocks ds le disque

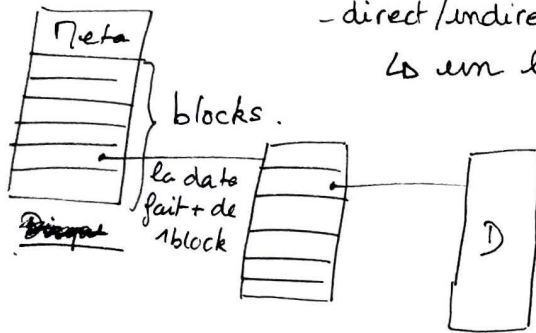
inode: - un en mémoire (utilisé par l'OS)

- un s/ le disque

- il contient de la data et metadata, il est référencé par un N°.

- direct/indirect block/double/triple indirect.

↳ un lien de chaque qui correspondent aux tailles de blocks = modèles classique.



Directory: fichier à un autre sauf qu'on change le type et que de la data on met un tableau d'entry (= nom + n°-inode + length du nom).

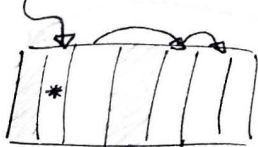
→ Superblock placé en début de fileSystem

on en met plus car si il est effacé c'est le merde. des superblocks secondaires ne contiennent que les stats.

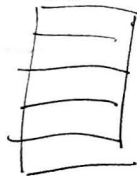
FAT (= File Allocat° Table)

root

*Dir



FAT



pb: on écrit tjrs au même endroit.

pratique car l'implem prend peu de place et est simple.

(regarder direntry) → roots?

Dans une carte SD on a une mini RAM de laquelle se trouve la FAT. (pr aller + vite).
Du coup elle est commit avant qu'on retire la carte de l'ordi.

- zfs

- Merkle Tree.

- btrfs.

arbre et chaque nœud est hashé et contient le nœud de son parent
ex: git.

1 seul FS visible s/ linux.

mount() → prend un disque et l'expose ds un dossier.

VFS → Virtual File System

table contient les pointeurs vers et les fonctions d'un objet, elle est placée
début de la structure de l'objet.

même principe est utilisé dans les fichiers. ils possèdent des pointeurs sur
et les f° de manipulatio de fichier (read, write, open, ...)

on a un compteur de read, ainsi si deux pointeurs pointent s/ le m° fichier
un des deux close, le fichier n'est pas détruit. principe que les shared ptr.

Quand `read` est appelé on appelle `vfs_read` qui lui même appelle le `read` du fichier.

↳ gère les buffers liés au kernel

le numero major indique la bonne struct file operation

le numero minor indique quel objet il faut traiter.

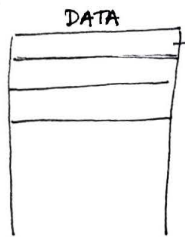
voir `ioctl` → syscall special pour gerer les IO s/ devices extérieurs.

Multithread.

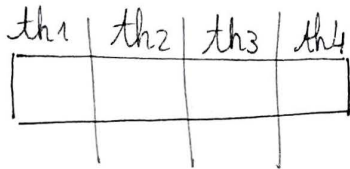
SIMD → Single Instruction Multiple Data.

exemple avec addition sur un vecteur:

il existe des registres qui prennent 4 flottant à la fois pour faire des opérat° 4x plus vite.



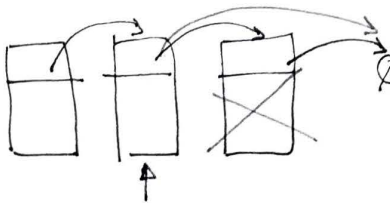
↳ quand on charge de la data on charge une cacheline entière.



Pour aller + vite on peut diviser le calcul s/ le vecteur en le confiant à 4 thread s'exécutant en mêm temps. Un thread

s'occupe des cases de 0 à N et non pas $\% \text{nb_thread} = \text{th_id}$ car l'écriture et lecture des la data se fait par cacheline. Si 2 east thread écrivent s/ la même data, (chacun dans leur cash) les cash se synchronisent ce que ↑ le temps d'exécution. = False sharing.

exemple sur une liste chaînée



Imaginons qu'un thread A s'apprête à lire un pointeur next et qu'un thread B supprime l'elt pointé au mêm moment. A va lire de la mémoire qui n'est plus dispo.

On a les sect° critiques et des endroits dangereux (ici ajout, suppression, ...) de thread va donc devoir ê le seul. On lock donc la partie critique qui va devoir attendre que le reste soit finis

Atomic Primitives

→ test and ~~test~~ set (TAS)

→ compare and swap (CAS)

volatile: dis au compilateur de fetch la variable

pb producteur / consommateur : un thread produit et l'autre consomme.

prod → Queue → conso.

```
lock
for (..)
{
  if (v==0)
  {
    v=1
    break
  }
}
```

unlock

↳ spin lock : le CPU attend activem^t

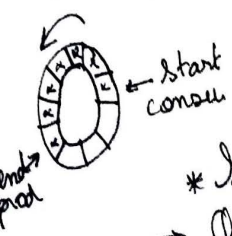
* else sleep(0)

↳ mutex : dans pendant un temps

mieux si
→ j'ai besoin d'aller vite
→ j'ai besoin d'aller vite

⇒ Tt les autres cas.

la queue est implémentée s/s la forme d'une ring buffer



Le pb c'est si \neq thread et consumer ou producteur.
On peut mettre des mutex quand on modifie la queue mais

* Semaphore (demander à Jeanne)

- Queue de process en attente et count.
- Quand on a besoin de compter.

→ attribut d'une ressource par un thread stocké de une res.

Deadlock: lock indéblocable :

T1	T2
lock(A)	lock(B)
lock(B)	lock(A)
unlock(A)	

ici T2 attend que T1 unlock A
mais T1 attend que T2 unlock B pour unlock A.

live lock: Deux personnes qui se gêne dans un couloir

T1 _____ lock(A)
T2 _____ lock(A) x unlock(A)

* ici T1 + prioritaire que T2 donc il est exécuté avant T2 qui a déjà lock A. Du coup on fait une inversion de priorité pr que T2 soit exécuté avant pr pouvoir libérer A.