

Impact du cache important : presque tout est un fichier donc tout passe par lui.

Metadata fs en direct, seulem^t les datas dans le cache.

Cache très utile pr mmap!

Au moment du dead, page-fault → il voit que mapping de fichier

→ vm_operations_struct : *fault

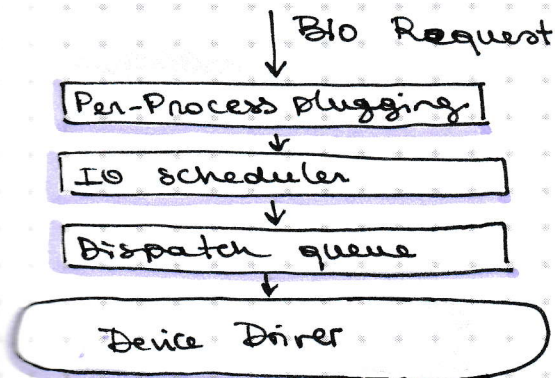
2 types de mapping :
• private read → page cache write → duplique
• shared changem^t propagé → écrit ds page cache

Block layer:

→ requestes IO : start, r/w, taille
• hints (SYNC), flags (FUA, FLUSH)

Vie d'une requête

- crée ds block layer quand IO
- peut ê delayed, merged (IO scheduler, multi queue handling)
- dispatch dans le bon driver
- driver signale qd IO finit

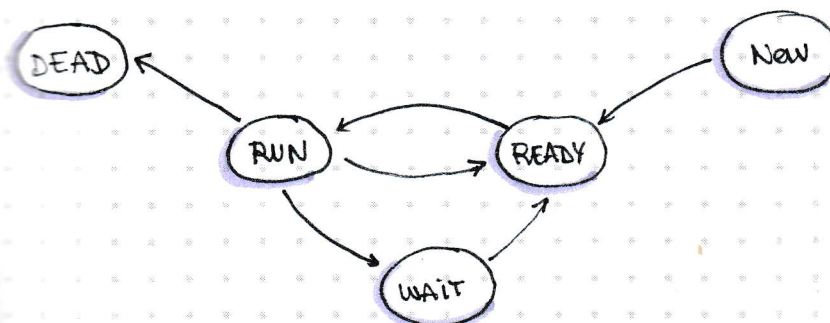


• NOOP : just pass requests into dispatch queue

• CFQ : Completely fair queueing
• acté par défaut
• sync > async
• support IO priorities, cgroups, sync request idling...

• Deadline :
• read > write
• reduce seeking : sort
• dispatch request at least after deadline expired.
→ le - tard possible

SCHEDULING



Contraintes:

- No starvat^o
- Fair
- Fast ⇒ $O(1)$ (du -, en amorti)

Ça qui nous intéresse : # le monde tourne en un tps, de manière équitable, # le tps qui tourne.

Choisir qui tourne \rightarrow prend du tps. \rightarrow - de tps pr process

Temps CPU : tps d'exec sur processeur
calcul en global ou par process

- tps user \rightarrow max
- tps $\$/$ kernel \rightarrow min

Tps d'exec : abnité de tps d'exec d'un process.

process qui finit : global

process qui finit pas : tps d'exec / tps \rightarrow ratio
(tps passé en running)

Tps turn around : tps entre new et dead \rightarrow tps que le process met à finir.

Débit (throughput) : diff entre exec et turn around.

But : tjs qqn qui tourne

Pb : on ne sait pas qd le process va faire un syscall
 \rightarrow prédire le futur?

On peut faire des stats sur les anciennes exec pour essayer de prédire un scheduling idéal.

Tps de réponse : tps qu'une appli met à demander

\rightarrow tps entre new et 1^{er} running

on aimerai que ce tps de réponse soit le + petit possible.

Tps d'attente : tps passé en ready : tps attendu alors qu'elle était prête.

on cherche à le minimiser \rightarrow il faut faire des choix.

ex : shell. Passe son tps en waiting. Mais on a envie qu'elle soit très rapide \rightarrow peu de tps en ready!

Appli interactive : interagit ac qqch (user, $\$/$, réseau...)

passer son tps à faire des syscall

Il y a des applis qui font peu de syscall (peu ou pas d'interact°)
 \rightarrow + calculatoire.

Goulot d'étranglement : limite le débit. Ralentis l'appli et l'empêche d'aller aussi vite qu'elle voudrait



ex : vitesse du disque pour accès au fichier, tps CPU pr compil, qt de mémoire, jeu limité par CPU & GPU, ...

Identifier les goulots d'étranglement permet de savoir où opti.

Goulot d'étranglement du scheduler : IO & tps CPU.

\rightarrow Tâches IO Bound (syscall bloquant)

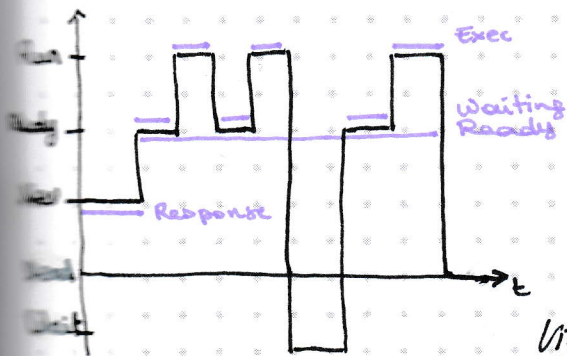
\rightarrow Tâches CPU Bound

\rightarrow interactif

\rightarrow calculatoire

- Tâches IO Bound : min waiting time
- Tâches CPU Bound : exec time ou turnaround time

Rq: en pratique le bound change pendant l'exéc



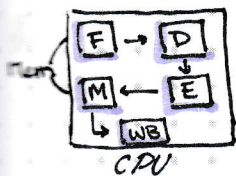
On veut min. exec time
 → mais le nbr d'instr ne change pas?
 l'exec time ne se resume pas à ça
 → access mem, gest° cache...

à chaque rescheduling → tps d'init
 (not. TLB (cache paginat°))
 → flush à chaque changement de process.



On perd trop de tps en exec à reemplir ts les caches.

Processeur : Fetch • Decode • Exec • Mem • Writeback



Pls calculs peuvent être fait en m. tps → parallèle.
 à chaque tick on avance (+ vite, ↑ fréquence)

1 seule instr qui nous intéresse à la fois.

À aller + vite :

Fetch (1)	Decode (2)	Exec (1)	Mem (4)	Writeback (1)
1				
2	1			
3	2	1		
4	3	2	1	
5	4	3	2	1

→ Pipeline

Sans pipeline : 1 instr 9 tick, 1 instr tous les 9 ticks

Avec pipeline : 1 instr 9 tick, 1 instr tous les 4 ticks

→ C'est l'étage le + lent qui ralentit.

À aller + vite, on peut split (Mem en 4 part, Decode en 2)

→ on monte artificiellement en fréquence.

Intel allait jusqu'à 35 étages sur Pentium 4

Pb : ça chauffe.

- jmp → solut° ?
- on vire tout (on perd potentiellement 100%)
- on exec 8^{uit} (c'est au code d'être bien fait)
- prédire le futur ? → predictor de branchement
 on charge les instr qu'on a le + de chance d'exé. → likely / unlikely

ex: `cmp` —
`nop`
`je` —
`mov 0, a`
 ↳ `mov 1, a`

Au final, augmenter autant la fréquence n'est pas la solut° → archi + petite, cœurs + petits, mais pls cœurs. (Archi Core)

Rq: Basic block → ~~block~~ entre instr et jmp.

Pour pas perdre tous les caches, predictor de branchement...

→ exec CPU Bound + lgtps, - sv+ } équitable
→ exec IO Bound + sv+, - lgtps

Au final :

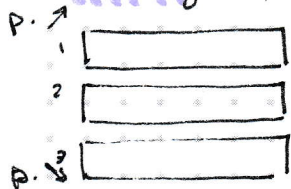
IO Bound → waiting → Quantum →
CPU Bound → waiting → Quantum →

Comment détecter quel bound ?

IO Bound → nb syscall / s, ne finit pas ses quantum.

CPU Bound → finit le quantum sv+

Ready: pls queues par priorité



On prend d'abord de la queue haute priorité
puis on descend.

Quand on devient IO Bound, priorité ↑
" " " CPU Bound, priorité ↓

(ça marche pas tout car on priorise trop les CPU Bound au final)

→ Windows (en vrai, il prend pas que IO / CPU Bound en compte.
graphisme in kernel → priorité qui change en fct)
(affichage priorité ↑) (trucs chelous avec installateurs)