

Ordonnancement

Programmation concurrente

Grégoire Pichon

Univ. Claude Bernard Lyon 1

séance 3

Sylvain Brandel, Yves Caniou, Guillaume Damiaud, Meriem Ghali, Laurent Lefèvre, Thibaut Modrzyk, Grégoire Pichon, Alec Sadler, Florence Zara, Jerry Lacmou Zeutouo

- 1 Introduction
- 2 Algorithmes d'ordonnancement
 - Sans préemption
 - Avec préemption
- 3 Améliorations
- 4 Ordonnancement dans les systèmes
 - Linux
 - Windows
 - Evolutions
- 5 Conclusion

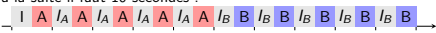
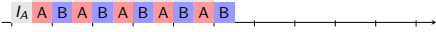
Rappels sur le cours

- Page du cours : <https://asr-lyon1.gitlabpages.inria.fr/prog-concurrente/>
- QCM, obligatoires : répondez « oui » à « Êtes-vous prêt à répondre aux questions sur l'ordonnancement ? »
(Rappel : pénalité sur la note de CCF si non fait une semaine après le CM)

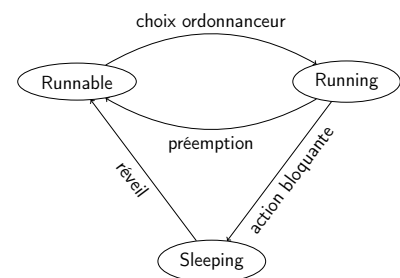
Temps partagé

- N processus/threads, M cœurs
- Cas simple : $N \leq M \rightsquigarrow$ chaque cœur exécute un processus
- Et si $N > M$? Il faut :
 - ▶ Décider quel processus lancer sur un cœur
 - ▶ Décider quand reprendre la main pour la passer à un autre processus
- C'est le rôle de l'**ordonnanceur** (scheduler)

Importance de la politique d'ordonnancement : cas d'école

- À un instant donné 2 processus (A et B) doivent s'exécuter.
- Chacun alterne $\frac{1}{2}$ s de lecture sur le disque et $\frac{1}{2}$ s de calcul.
- Chacun dure en tout 5 secondes.
- Si on les exécute
 - ▶ à la suite il faut 10 secondes :

 - ▶ en même temps il faut 5,5 secondes

- Si un processus est bloqué, il faut exécuter un autre processus

État de chaque processus



Un peu plus de réalisme

Idee simplifiée : « Si un processus est bloqué, il faut exécuter un autre processus ». Pas si simple :

- Les tâches arrivent aléatoirement (sauf cas particulier ex. temps réel).
- On ne connaît pas a priori leur durée.
- Plusieurs tâches peuvent être disponibles et assez souvent aucune ne l'est.
- Elles ont des importances différentes (p. ex. swap).
- Elles peuvent avoir des dépendances.
- Il y a parfois des contraintes externes (pilote).
- Le comportement des tâches est différent :
 - ▶ Processus interactifs (court, temps de réponse important)
 - ▶ Processus avec beaucoup d'entrée/sortie (analyse de disque)
 - ▶ Long calculs (rare)

Ordonnancement

Définition (Ordonnancement)

La sélection dans le temps des processus ou threads qui peuvent accéder à un processeur est l'**ordonnancement**. Le but est de maximiser :

- Le débit (nombre de processus traités par unité de temps)
- Le taux utile (le taux du processeur effectivement utilisé pour les tâches utilisateurs)

Il y a plusieurs types d'ordonnancement en fonction de la possibilité d'interrompre une tâche :

- **Ordonnancement collaboratif** : les tâches ne sont pas interruptibles.
- **Ordonnancement préemptif** : le système peut interrompre une tâche à tout moment.

Préemptif vs collaboratif

- L'ordonnancement collaboratif est plus simple
 - ▶ Les tâches rendent la main lorsqu'elles sont finies ou lors de certains *appels systèmes particuliers* (ex `read()`, `sleep()`, `pthread_yield()` ou `thread::yield()`, `SwitchToFiber()`)
 - ▶ Les tâches critiques sont protégées (écriture, périphérique)
 - ▶ Mais il est difficile d'assurer que le système reste réactif : Si une tâche non interruptible entre dans une boucle infinie...
 - ▶ ⇒ non-adapté aux systèmes d'exploitation grand public
- L'ordonnancement préemptif évite ce problème
 - ▶ Une boucle infinie dans un programme ne compromet pas le système
 - ▶ Le système n'est pas bloqué par certaines tâches trop longues
 - ▶ Mais cela pose des difficultés pour les tâches critiques (driver, écritures)
- Par exemple :
 - ▶ MsDos et Windows 3.1 utilisaient un ordonnancement collaboratif;
 - ▶ Windows 95 Windows 98 utilisaient un mode collaboratif pour certaines tâches.
 - ▶ Windows Nt, Xp et les suivants, MacOS, Unix, Linux... utilisent le préemptif.



Critère de performance

Un algorithme d'ordonnancement peut être jugé selon plusieurs critères :

- taux utile : utilisation du processeur ;
- débit : nombre de travaux réalisés ;
- temps de réponse (temps entre l'arrivée et la fin d'une tâche, parfois appelé « temps de traitement » ou « latence », *response-time* en anglais) : surtout pour les processus interactifs ;
- temps d'exécution moyen d'un ensemble de tâches ;
- temps de réponse total d'un ensemble de tâches ;
- équité entre les tâches (au moins assurer l'absence de famine).

Ces critères sont mutuellement contradictoires, par exemple pour minimiser le temps de réponse moyen il faut favoriser les petites tâches.



Ordonnancement à plusieurs niveaux

- Un processus élu n'utilise le processeur que pendant un court laps de temps (p. ex. 10ms).
- Si le choix du processus demande trop de temps (p. ex. 1ms) le taux utile diminue.
- On peut utiliser plusieurs niveaux d'ordonnancement :
 - ▶ Ordonnancement à court terme rapide pour choisir un nouveau processus (p. ex. FIFO ou liste de priorité)
 - ▶ Ordonnancement à long terme (calcul des priorités pour les processus, swap)



Temps-réel ou pas ?

- PC Bureautique (web, traitement de texte, ...)
 - ▶ But de l'ordonnancement : donner une sensation de fluidité à l'utilisateur
 - ▶ Difficilement quantifiable
 - ▶ Non « temps-réel »
- Multimédia (télévision numérique)
 - ▶ But de l'ordonnancement : *éviter* que les images n'arrivent trop tard (sinon, vidéo non-fluide), que le son arrive trop tard (sinon, craquements dans le haut parleur)
 - ▶ Notion de *date limite* (deadline) bien définie
 - ▶ *Temps-réel mou* (soft real-time) : pas de conséquences dramatiques si on manque une deadline
- Freinage ABS, injection moteur, commandes de vol d'un avion
 - ▶ But de l'ordonnancement : *garantir* que les traitements sont faits dans les temps
 - ▶ Ensemble de calculs à réaliser et deadlines bien définies
 - ▶ *Temps-réel dur/strict* (hard real-time) : conséquences dramatiques en cas de deadline manquée.



Ordonnancement généraliste ou temps-réel

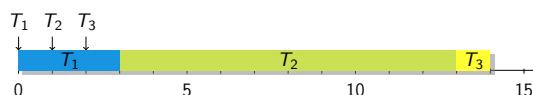
- Ordonnancement généraliste (par défaut sous Linux, Windows, Mac OS) :
 - ▶ Heuristiques complexes pour choisir les processus
 - ▶ Ne jamais donner trop longtemps le processeur au même processus
 - ▶ Fonctionne bien en pratique, mais aucune garantie
- Ordonnancement temps-réel (OS temps-réels dédiés, ou politiques spécifiques sous Linux, Windows, Mac OS) :
 - ▶ Doit être simple et prédictible
 - ▶ On connaît souvent les pires temps d'exécution (WCET)
 - ▶ Un processus prioritaire doit prendre la main tout de suite (exemple : frein ABS contre autoradio ? son contre image ?)



FIFO

Le plus simple, on exécute les tâches dans l'ordre d'arrivée : Par exemple :

	T_1	T_2	T_3
Durée	3	10	1
Arrivée	0	1	2



$$\text{temps de réponse moyen} \simeq \frac{(3-0)+(11-1)+(3-2)}{3} = 9$$

Cette méthode a les mêmes problèmes que toutes les queues.



Le plus court d'abord (Short Job First)

C'est la tâche disponible la plus courte qui est exécutée en premier sans tenir compte de l'ordre d'arrivée. Avec le même exemple :

	T_1	T_2	T_3
Durée	3	10	1
Arrivée	0	1	2



$$\text{temps de réponse moyen} \simeq \frac{(3-0)+(3-2)+(13-3)}{3} = 6$$

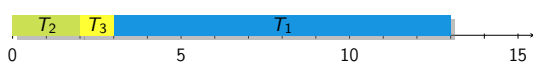
Le temps de réponse moyen est plus faible. Cet algorithme est efficace, il est optimal pour le temps de réponse moyen si les tâches arrivent en même temps. Mais il demande de connaître ou estimer la durée d'une tâche.



Avec priorité

Selon l'importance de la tâche, le système ou l'utilisateur attribue une priorité. C'est toujours la tâche la plus prioritaire à un instant donné qui est élue :

	T_1	T_2	T_3
Durée	10	2	1
Arrivée	0	0	0
Priorité	10	3	4



Problème

Les 2 derniers ordonnancements (SJF et à priorité) ne sont pas équitables :

- Si un processus est peu prioritaire (ou long pour SJF), il est bloqué par tous les autres

	T_1	T_2	T_3	T_4	T'_2
Durée	10	2	1	3	2
Arrivée	0	0	0	2	5
Priorité	10	3	4	6	3



- Il y a famine
- Idée de solution* : on peut augmenter la priorité des processus qui ne sont pas élus
- Nouveau problème* : une tâche longue devient prioritaire puis bloque le processeur (mauvais temps de réponse).
- ⇒ il faut une préemption pour découper la tâche.

1 Introduction

2 Algorithmes d'ordonnancement

- Sans préemption
- Avec préemption

3 Améliorations

4 Ordonnancement dans les systèmes

- Linux
- Windows
- Evolutions

5 Conclusion

Algorithmes précédents

- FIFO ne peut pas être préemptif (la tâche qui serait élue en cas de préemption est déjà celle exécutée).
- Les algorithmes SJF et à priorité peuvent l'être. Par ex :

	T_1	T_2	T_3
Durée	10	2	1
Arrivée	0	1	2
Priorité	10	2	5

► SJF



► Avec priorités

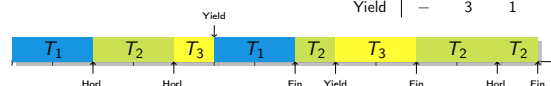


Tourniquet (Round Robin) - I

- On définit un *Quantum de temps* : une durée maximum d'exécution.
- Les tâches sont élues selon leur ordre d'arrivée (FIFO).
- Elles s'exécutent jusqu'à :
 - expiration du quantum
 - ou suspension (yield) par la tâche elle-même (appel système, défaut de page, terminaison...).

Par exemple avec un quantum de $q = 2$:

	T_1	T_2	T_3
Durée	4	6	3
Arrivée	0	1	2
Yield	—	3	1



Tourniquet (Round Robin) - II

- Toutes les tâches ont accès au processeur.
- Le temps d'attente est borné par $(n-1) \times q$ (n le nb. de tâches, q le quantum).
- Si le quantum est trop grand, c'est équivalent à l'algorithme FIFO.
- Si le quantum est trop court, il y a baisse du taux utile (trop de changements de contexte).
- Par exemple, on fixe le quantum de manière à ce que 80% des tâches terminent d'elles-mêmes. Ou bien (Linux Completely Fair Scheduler (CFS)) : quantum fixé tel que toutes les tâches aient une tranche de temps 50 fois par secondes (toutes les 20 ms).

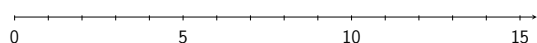
Ordonnancement à queues multiples

- Le système gère plusieurs listes de tâches.
- Chaque liste correspond à une priorité.
- Le scheduler choisit une tâche dans la première liste non vide.
- Par exemple :
 - Processus critique (swap, I/O...).
 - Processus système.
 - Processus interactif.
 - Processus de calcul.
 - Tâche de fond.
- L'ordre utilisé est très important. Par exemple, si le swap est moins prioritaire qu'une tâche qui a besoin de mémoire...
- On peut utiliser des algorithmes différents selon la liste.

Exemple

Avec un quantum de $q = 2$, ordonnez (en RR) :

	T_1	T_2	T_3	T_4
Durée	6	3	4	3
Arrivée	0	2	3	8
Priorité	4	2	2	4



Évolution de la priorité (non temps-réel) - I

- Problème classique :
 - Processus de calcul long en tâche de fond
 - Va de toutes façons prendre beaucoup de temps, inutile de le rendre prioritaire
 - Gourmand en CPU
 - Processus interactif (éditeur de texte, ...)
 - Temps de réponse court important pour l'utilisateur
 - Peu gourmand en CPU
- Principe de la solution :
 - Donner la main en priorité aux processus qui ont peu utilisé le CPU récemment.

Évolution de la priorité (non temps-réel) - II

Les algorithmes utilisant une priorité fixe sont susceptibles de conduire à la famine des processus peu prioritaires.

- On fait évoluer la priorité en fonction du temps passé hors du processeur.
- Par exemple (BSD), tous les n tops d'horloge, une nouvelle priorité est recalculée :

$$P_{\text{réelle}} = P_{\text{min}} + \frac{T_{\text{cpu}}}{n} + 2 \times P_{\text{nice}}$$

Avec :

- $P_{\text{réelle}}$ la priorité calculée (elle est ramenée à une valeur entre P_{min} et P_{max}).
- P_{min} la priorité minimum (la meilleure).
- T_{cpu} une estimation du temps passé dans le processeur.
- P_{nice} une valeur de « Bonne conduite » (cf. commande nice)



Évolution de la priorité - III

- T_{cpu} est calculée :
 - En augmentant de 1 à chaque top d'horloge la valeur du processus exécuté
 - En diminuant la valeur au cours du temps :

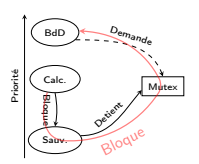
$$\text{Toutes les secondes } T_{\text{cpu}} = T_{\text{cpu}} \frac{2\text{load}}{2\text{load} + 1}$$

où *load* est la moyenne du nombre de threads en attente.



Inversion de priorité - I

Supposons 2 processus, une base de données et sa sauvegarde, qui partagent une zone protégée par un *mutex*.



- Comme la base a un rôle important, sa priorité est élevée.
- La sauvegarde a une priorité faible.
- À un instant donné la sauvegarde acquiert le *mutex*.
- Si une tâche de calcul de priorité moyenne arrive à ce moment,
- le calcul bloque la sauvegarde (car elle est plus prioritaire).
- La sauvegarde bloque la base de données car elle détient le *mutex*.

En pratique le calcul bloque la base de données qui est une tâche plus importante que lui.



Inversion de priorité - II

- Pour résoudre ce problème, certains systèmes utilisent *l'héritage de priorité*
 - Un processus qui bloque un processus plus prioritaire obtient sa priorité.
 - Ce n'est utilisé que dans des systèmes spécifiques (par ex. temps réel).
- Les autres utilisent des méthodes moins sûres :
 - Pour linux, l'évolution de la priorité fait en sorte que la tâche bloquante finira par être élue.
 - Pour windows, une tâche qui n'est pas élue pendant « longtemps » obtient la priorité maximale pendant 2 quantum de temps.



Linux

Le système linux ordonnance directement les threads (pas les processus). Il définit plusieurs classes de threads :

- 2 classes dites « temps réel »
 - Les FIFO (SCHED_FIFO)
 - Les Round Robin (SCHED_RR)qui ne sont pas réellement temps réel, car elles ne permettent pas d'obtenir de garantie sur les échéances.
- Les threads normaux (non temps-réels) :
 - Les threads standards (SCHED_OTHER) ← valeur par défaut
 - Les threads de calculs (SCHED_BATCH)
 - Les threads de tâches de fond (SCHED_IDLE)



Threads Temps Réel

- Ils sont toujours plus prioritaires que des threads non-temps réels.
- Ils ont une priorité statique comprise entre 1 et 99.
- Il faut être root pour utiliser ces priorités.
- La différence entre les deux est :
 - Les threads FIFO ne peuvent être interrompus que par un thread plus prioritaire.
 - Les threads RR sont préemptés lorsque leur quantum de temps est atteint (ou qu'un thread plus prioritaire arrive).



Threads non temps-réels - I

- Les threads normaux SCHED_OTHER :
 - Leur priorité statique est nulle (ils sont toujours plus lents que les threads temps réel).
 - Leur quantum de temps est plus faible.
 - Ils ont une priorité dynamique qui repose sur une valeur de *courtoisie* (*nice*) et qui évolue en fonction de leur utilisation du processeur.
 - Le système applique un bonus aux threads qui bloquent (sur des entrées sortie) et un malus aux processus qui utilisent le processeur.
 - ⇒ les processus interactifs sont favorisés.
 - En tout il y a 40 priorités pour les threads en temps partagé (donc 140 en tout).



Threads non temps-réels - II

- Les threads de traitements par lot SCHED_BATCH :
 - identiques aux précédents,
 - mais considérés forcément comme threads utilisant le processeur, ils sont pénalisés.
- Les tâches de fond SCHED_IDLE :
 - sans priorité statique ni courtoisie ;
 - ordonnées en dernier.



Windows

(hors-programme)

- Le système propose 32 priorités
 - de 0 à 15 pour les threads utilisateurs
 - de 16 à 31 pour le système (Temps Réel)
- Entre les threads de même priorité, le système utilise *Round Robin*
- La priorité se calcule en 2 parties :
 - Priorité de base du processus : `IDLE_PRIORITY_CLASS`, `BELOW_NORMAL_PRIORITY_CLASS`, `NORMAL_PRIORITY_CLASS`, `ABOVE_NORMAL_PRIORITY_CLASS`, `HIGH_PRIORITY_CLASS`, `REALTIME_PRIORITY_CLASS`.
 - La priorité du thread dans le processus : `THREAD_PRIORITY_IDLE`, `THREAD_PRIORITY_LOWEST`, `THREAD_PRIORITY_BELOW_NORMAL`, `THREAD_PRIORITY_NORMAL`, `THREAD_PRIORITY_ABOVE_NORMAL`, `THREAD_PRIORITY_HIGHEST`, `THREAD_PRIORITY_TIME_CRITICAL`,



Résumé

(hors-programme)

		Processus					
		Idle	Below n.	Normal	Above n.	High	R. T.
Threads	Idle	1	1	1	1	1	16
	Lowest	2	4	6	8	11	22
	Below n.	3	5	7	9	12	23
	Normal	4	6	8	10	13	24
	Above n.	5	7	9	11	14	25
	Highest	6	8	10	12	15	26
Time Crit.		15	15	15	15	15	31



En plus

(hors-programme)

- Des bonus/malus
 - bonus à tout processus dont la fenêtre arrive au premier plan,
 - bonus pour ceux qui se réveillent suite à une entrée-sortie,
 - malus pour un thread qui finit son quantum de temps,
 - bonus pour un thread qui n'a pas eu la main pendant une période « trop longue ».
- Des priorités limites
 - `THREAD_PRIORITY_IDLE` : 16 pour les temps réels et 1 pour les normaux.
 - `THREAD_PRIORITY_TIME_CRITICAL` : 31 pour les temps réels et 15 pour les normaux.
- Des threads particuliers de priorité 0 pour remplir de 0 la mémoire (sécurité).



Cas des multiprocesseurs

Dorénavant, les ordinateurs ont plusieurs cœurs de calculs. Cela permet d'exécuter plusieurs tâches à la fois mais pose de nouveaux problèmes d'ordonnancement.

- Les caches font que l'accès aux données n'est pas identique pour tous les processeurs.
- Deux tâches échangeant des données sur des processeurs différents vont générer des défauts de caches.
- Dans les modèles avec plusieurs processeurs, il y a des cas où chaque processeur a un banc de mémoire dédié (NUMA).
- Un processeur ayant besoin d'avoir accès à une mémoire sur un autre banc doit faire appel à la MMU de l'autre processeur.

Si l'ordonnanceur ne tient pas compte de cela, le résultat ne sera pas efficace.

Par exemple, pour un hyperviseur, il est préférable que les tâches d'une même machine virtuelle soit ordonnancées sur le même processeur : celui qui dispose de sa mémoire.



Isolations/jails

(hors-programme)

Les systèmes permettent maintenant de regrouper des processus pour les gérer ensembles et les isoler. Par exemple, sous linux les *cgroups* permettent de :

- limiter leur utilisation de ressources, arbitrer entre les différents groupes ;
- stopper, sauvegarder et redémarrer des groupes ;
- isoler les processus via des *espaces de nommage* :
 - identifiant de processus (PID), ces processus ne voient que ceux du même groupe ;
 - interface réseaux, IPC : ils ne voient que les outils de communications affectés au groupes ;
 - systèmes de fichiers : ils peuvent avoir leur propre point de montage ;
 - utilisateurs ;
 - ...



Virtualisation au niveau du système

(hors-programme)

L'isolation autorise à créer un système proche de la virtualisation :

- des systèmes quasi séparés ;
- mais qui partagent le même noyau ;
- demandent beaucoup moins de ressources ;
- Sous Linux ce sont les conteneurs, par exemple lxc, docker.

Cela permet :

- d'isoler des applications ;
- de faire cohabiter plusieurs versions d'une même bibliothèque ;
- d'implémenter des microservices.



Conclusion

- Le but de l'ordonnancement est de choisir la tâche qui va s'exécuter.
- Un mauvais ordonnancement peut être catastrophique.
- Sur les systèmes simples :
 - Repose souvent sur une priorité qui évolue au cours du temps.
 - 2 classes de processus : système (RT) et utilisateur.
 - Utilisation du *Round Robin*
- Système plus complexes :
 - Temps réel (besoin de garantie)
 - Automatisme (contraintes)



À retenir

- Différence entre préemptif et collaboratif
- Algorithme d'ordonnancement
 - FIFO
 - Round Robin
 - Priorité
- Évolution de la priorité dans le temps
- Inversion de priorité

