

**SCHEDULING** -> NEL CONTESTO DEI SISTEMI OPERATIVI RIGUARDA IL COME LA CPU / LE CPU VENGANO ALLOCATE AI VARI PROCESSI

ESEMPIO: ABBIAMO UN SISTEMA CON UNA SOLA CPU E UN SOLO CORE, E ABBIAMO UN'UNITÀ DI ESECUZIONE E 10 PROCESSI CHE VOGLIONO ANDARE IN ESECUZIONE -> LO SCHEDULING FA IN MODO CHE IL NUCLEO DECIDA QUALI PROCESSI MANDARE IN ESECUZIONE E PER QUANTO

LO SCHEDULING IN GENERALE RIGUARDA IL COME ALLOCARE DEI LAVORI SU DELLE UNITÀ DI ESECUZIONE

PER QUANTO RIGUARDA IL SISTEMA OPERATIVO, CI SONO DUE PARTI CHE RIGUARDANO LA SCHEDULAZIONE DEI PROCESSI:

- IL MECCANISMO CHE PERMETTE AL KERNEL DI CAMBIARE PROCESSO -> CHIAMATO **CAMBIO DI CONTESTO (CONTEXT SWITCH)** -> SCRITTO IN LINGUAGGIO MACCHINA
- LO SCHEDULER È LA PARTE DEL KERNEL CHE DECIDE IL PROSSIMO PROCESSO DA MANDARE IN ESECUZIONE, SEGUENDO UNA CERTA **POLITICA (POLICY)** -> SCRITTO IN C

NELL'ESEMPIO CON UN SOLO PROCESSORE DA SINGOLO CORE, SE STA GIRANDO IL PROCESSO NON STA GIRANDO IL SISTEMA OPERATIVO

COME FA IL SISTEMA OPERATIVO A RIPRENDERE IL CONTROLLO ED EVENTUALMENTE UCCIDERE IL PROCESSO CHE SI STAVA ESEGUENDO E FARNE PARTIRE UN ALTRO?

CI SONO DIVERSI APPROCCI:

- **COOPERATIVO** -> IL SISTEMA OPERATIVO SI "FIDA" DEI PROCESSI
  - I PROCESSI FARANNO DEI SYSTEM CALL, FACENDO DELLE SYSTEM CALL CHIAMANO IL KERNEL E A QUEL PUNTO IL KERNEL PUÒ DECIDERE COSA FARE CON IL PROCESSO
  - SE UN PROCESSO NON FA SYSTEM CALL IL COMPUTER MUORE LOL, RIMANE PER SEMPRE NEL PROCESSO
- **non cooperativo** -> IL SISTEMA OPERATIVO RIPRENDE IL CONTROLLO A FORZA, TRAMITE UN TIMER INTERRUPT

NEL **CONTEXT-SWITCH** SI SALVANO I REGISTRI DEL PROCESSO CHE SI STAVA ESEGUENDO E SI VANNO A CARICARE I REGISTRI DEL PROCESSO CHE VOGLIAMO MANDARE IN ESECUZIONE

## OGNI PROCESSO HA UNO STACK UTENTE E UNO STACK KERNEL

OS @ run (kernel mode)	Hardware	Program (user mode)
		Process A
		...
	<b>timer interrupt</b> save regs(A) to k-stack(A) move to kernel mode jump to trap handler	
Handle the trap Call <code>switch()</code> routine save regs(A) to proc-struct(A) restore regs(B) from proc-struct(B) switch to k-stack(B) <b>return-from-trap (into B)</b>		
	restore regs(B) from k-stack(B) move to user mode jump to B's PC	
		Process B
		...

LO STACK POINTER È UN REGISTRO CHE PUNTA ALLA CIMA DELLO STACK, PERÒ UN PROCESSO PUÒ USARLO COME VUOLE -> LO STACK POINTER PUÒ AVERE UN VALORE SBAGLIATO, OPPURE IL PROCESSO POTREBBE DECIDERE DI UTILIZZARE LO STACK POINTER COME REGISTRO E USARLO PER CONTI, ECC.

QUANDO ARRIVA UN INTERRUPT AUTOMATICAMENTE LA CPU PASSA ALLO STACK KERNEL, CHE È UN VALORE CHE HA SETTATO IL KERNEL NELL'INIZIALIZZAZIONE DEL SISTEMA -> IN QUESTO MODO SAPPIAMO QUINDI CHE L'INTERRUPT HANDLER (IL PEZZO DI CODICE CHE GESTISCE L'INTERRUPT) AVRÀ UNO STACK VALIDO

### QUINDI SE VOGLIAMO PASSARE DA UN PROCESSO A AD UN PROCESSO B:

- SI STAVA ESEGUENDO IL PROCESSO A
- ARRIVA L'INTERRUPT
- VERRANNO SALVATI DEI REGISTRI SUL KERNEL STACK DEL PROCESSO A
- SI VA AD ESEGUIRE L'INTERRUPT HANDLER
- IL KERNEL DECIDE DI PASSARE AL PROCESSO B
- SI SALVERÀ NEL PCD (STRUTTURA CHE CONTIENE LE INFORMAZIONI DEI PROCESSI, IN QUESTO MOMENTO QUELLI DEL PROCESSO A) I VALORI ATTUALI DEL REGISTRO DELLA CPU E ANDRÀ A CARICARE QUELLI DEL PROCESSO B
- RITORNERÀ DALL'INTERRUPT

UN PROCESSO NELLA SUA VITA PUÒ ESSERE IN TANTI STATI:

- LO STATO **ready** -> IL PROCESSO È PRONTO AD ESSERE ESEGUITO MA NON SI STA ESEGUENDO
- LO STATO **running** -> IL PROCESSO SI STA ESEGUENDO
- LO STATO **blocked / waiting** -> IL PROCESSO NON HA RICEVUTO UN INPUT E STA ASPETTANDO (QUANDO RICEVUTO ENTRERÀ IN **ready**)

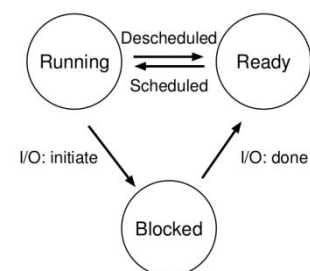


Figure 4.2: Process: State Transitions

LE POLITICHE DI SCHEDULING SONO QUEGLI ALGORITMI CHE DECIDONO QUALE PROCESSO MANDARE IN ESECUZIONE -> CI SONO DELLE METRICHE

PARTIAMO CON ASSUNZIONI IRREALISTICHE:

- CIASCUN JOB DURI LO STESSO TEMPO
- TUTTI I JOB ARRIVINO ALLO STESSO MOMENTO
- UNA VOLTA INIZIATO UN JOB SI ESEGUE FINO IN FONDO, SENZA INTERRUZIONI
- NON C'È INPUT/OUTPUT
- IL TEMPO DI CIASCUN JOB È NOTO A PRIORI

PRIMA METRICA -> **TEMPO DI TURN AROUND** =

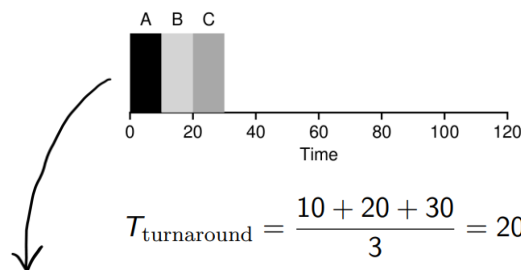
*tempo di completamento - tempo di arrivo*

Turn-around time

$$T_{\text{turnaround}} = T_{\text{completion}} - T_{\text{arrival}}$$

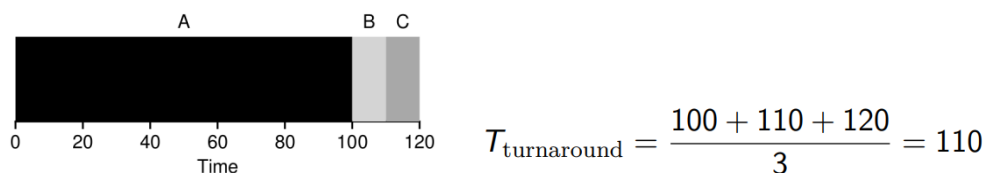
con le assunzioni iniziali,  $T_{\text{arrival}} = 0$  quindi:

$$T_{\text{turnaround}} = T_{\text{completion}}$$

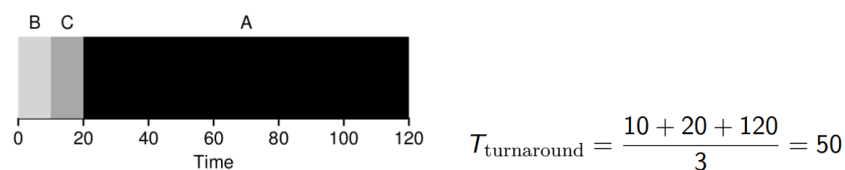


se a, B e C ARRIVANO ALLO STESSO TEMPO, L'ALGORITMO DEVE SCEGLIERE UN ORDINE

se a DURA 100, B e C 10



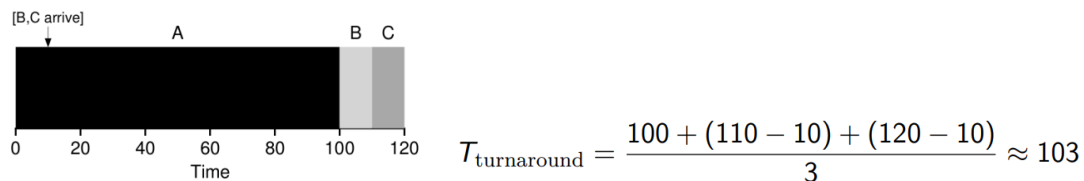
**EFFETTO CONVOGLIO:** GLI ALTRI DUE JOB SONO CORTI E IL PRIMO È LUNGO



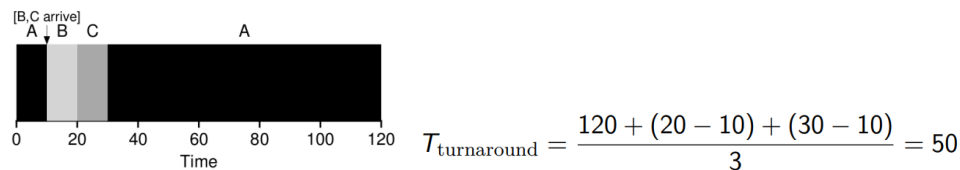
**SHORT JOB FIRST:** ALGORITMO PIÙ OTTIMALE

se non arrivano TUTTI ALLO STESSO MOMENTO

se arriva PRIMA a LUNGO, e POCO DOPO arriva B PIÙ CORTO, ORMAI SI STA ESEGUENDO a QUINDI NON POSSIAMO MANDARE IN ESECUZIONE GLI ALTRI



QUINDI COSA SI POTREBBE FARE? QUANDO ARRIVANO B e C METTO IN PAUSA L'ESECUZIONE DI a, COMPLETO B e C, e POI RIPRENDO a



### SHORTEST TIME-TO-COMPLETION FIRST

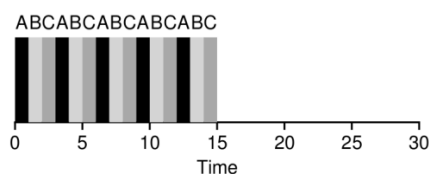
PUÒ PERÒ SOFFRIRE DI STARVATION PERCHÉ SE CONTINUANO AD ARRIVARE JOB CORTI, IL JOB a LUNGO NON FINIRÀ MAI

MA NON PUÒ ESSERE DAVVERO USATO, IN QUANTO IN UN SISTEMA OPERATIVO NON SAPPIAMO QUANTO SARÀ LUNGO UN PROCESSO

INOLTRE, MENTRE NEI SISTEMI NON INTERATTIVI IL TEMPO DI TURNAROUND MEDIO PUÒ ESSERE UNA METRICA SENSATA, NEI SISTEMI INTERATTIVI C'È UN'ALTRA METRICA MOLTO IMPORTANTE, OSSIA IL TEMPO DI RISPOSTA, CHE DEVE ESSERE MOLTO BASSO

*tempo di risposta = tempo di inizio esecuzione – tempo di arrivo*

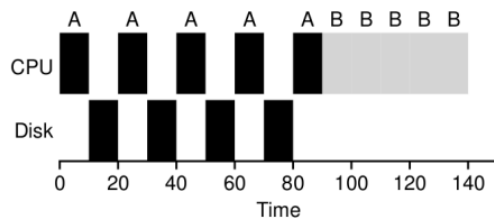
UN MODO PER AVERE UN TEMPO DI RISPOSTA BUONO È IL **ROUND-ROBIN** (GIROTONDO) -> MANDIAMO IN ESECUZIONE TUTTI I JOB CHE ABBIAMO PER UNA PICCOLA FETTINA DI TEMPO, ALTERNANDOLI



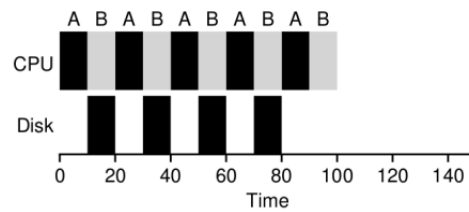
IL response-time con IL round-robin SI PUÒ OTTIMIZZARE RESTRINGENDO LA FETTINA DI TEMPO -> NON TROPPO PICCOLA PERCHÉ SE NO LA CPU FA SOLO CONTEXT-SWITCH E NON FA ALTRO DI UTILE

- evita la starvation ma penalizza i job corti

SAREBBE STUPIDO TENERE ALLOCATA LA CPU PER UN PROCESSO CHE ATTENDE INPUT/OUTPUT, QUINDI LO BUTTIAMO FUORI E POSSIAMO USARE LA CPU PER ALTRO



QUINDI



## ESEMPI DI ALGORITMI "REAL WORLD"

- MULTI-LEVEL FEEDBACK QUEUE (USATO DA WINDOWS E TEMPO FA DA LINUX)
- PROPORTIONAL SHARE E LINUX CFS (COMPLETELY FAIR SCHEDULER)

## MULTI-LEVEL FEEDBACK QUEUE (MLFQ):

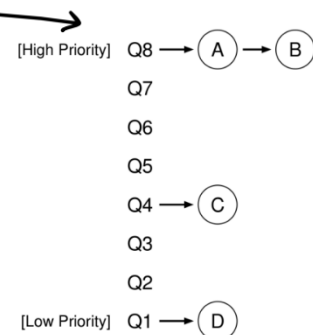
- CERCA DI UTILIZZARE SIA IL TURNAROUND TIME CHE IL RESPONSE TIME -> CERCA DI ESEGUIRE PRIMA I JOB CORTI, MA ANCHE OTTIMIZZANDO RESPONSE TIME
- L'IDEA È QUINDI CHE I PROCESSI INTERATTIVI AVRANNO UNA PRIORITÀ TENDENZIALMENTE PIÙ ALTA, MENTRE INVECE I PROCESSI CPU BOUND (CHE USANO PIÙ LA CPU) AVRANNO UNA PRIORITÀ PIÙ BASSA
- L'IDEA DEL MLFQ È DI AVERE PIÙ CODE A PRIORITÀ DIVERSA; E PER TUTTI I PROCESSI CHE STANNO SULLA STESSA CODA, VANNO IN ROUND-ROBIN
- LA PROPRIETÀ DI UN PROCESSO VERRÀ VARIATA IN BASE A COME SI COMPORTA IL PROCESSO; QUELLO CHE FA L'ALGORITMO È CERCARE DI PREVEDERE COME SI COMPORTERÀ UN PROCESSO BASANDOSI SU QUELLO CHE HA FATTO FIN'ORA

- NELL'ESEMPIO IN FOTO, ABBIAMO 8 CODE, DA Q1 A Q8, CON Q1 A BASSA PRIORITÀ E Q8 AD ALTA PRIORITÀ; LE REGOLE DICONO

- 1 Se  $p(A) > p(B)$ , gira A (e non gira B)
  - 2 Se  $p(A) = p(B)$ , A e B in RR
- CON  $p$  = PRIORITÀ
- CON QUESTE DUE REGOLE, STARVATION PER C E B

## QUINDI MODIFICHIAMO LE REGOLE FINO A QUELLE FINALI

- 1 Se  $p(A) > p(B)$ , gira A (e non gira B)
- 2 Se  $p(A) = p(B)$ , A e B in RR
- 3 Un nuovo job entra con la priorità massima
- 4 Quando un job usa un tempo fissato  $t$  a una certa priorità  $x$  (considerando la somma dei tempi usati nella coda  $x$ ), allora la sua priorità viene ridotta
- 5 Ogni  $s$  secondi, spostiamo tutti i job alla priorità più alta



IL PROBLEMA PIÙ GROSSO DELL'ALGORITMO È DEFINIRE I PARAMETRI (T, X, ECC.)

## COMPLETELY FAIR SCHEDULER (CFS):

- se io ho 3 processi che vogliono usare la CPU, do  $\frac{1}{3}$  di CPU a testa
- in generale manderò in esecuzione il processo che ha usato la CPU per meno tempo → perché se ne ha usata meno "se ne merita ancora" per essere pari agli altri
- il CFS misura il tempo utilizzato sulla CPU tramite **VIRTUAL RUNTIME** → **vruntime**
- se tutti i processi avessero la stessa identica priorità, il virtual runtime sarebbe semplicemente il runtime (quanto tempo un processo ha speso usando la CPU); il virtual runtime ci permette di avere processi con più priorità, in cui facciamo degli sconti, in cui lasciamo usare più CPU perché ha più priorità → gli viene conteggiato meno il tempo in cui sta sulla CPU
- in generale, quando c'è da scegliere quale processo mandare in esecuzione, il CFS sceglie quello che ha il virtual runtime più piccolo
- CFS usa vari parametri, uno è *sched\_latency* → il tempo assegnato ad ogni processo
- se ci sono  $n$  processi che sono pronti ad essere eseguiti, quindi nello stato ready, la fetta di tempo che viene data a ogni processo è  $\frac{\text{sched\_latency}}{n}$  → se però è troppo piccolo, vado a prendere il valore di un altro parametro, *min\_granularity*
- vengono mandati in esecuzione questi processi, e poi il vruntime di ogni processo verrà aggiornato in base a quanto tempo sono stati effettivamente in esecuzione

$$\text{slice}_k = \max\left(\frac{\text{sched\_latency}}{n}, \text{min\_granularity}\right)$$

$$\text{vruntime}_k += \text{runtime}_k$$

- nei sistemi reali, alcuni processi sono più importanti di altri, quindi c'è un valore *nice* che indica la priorità di un processo → più un processo è *nice*, più cede la CPU ad altri

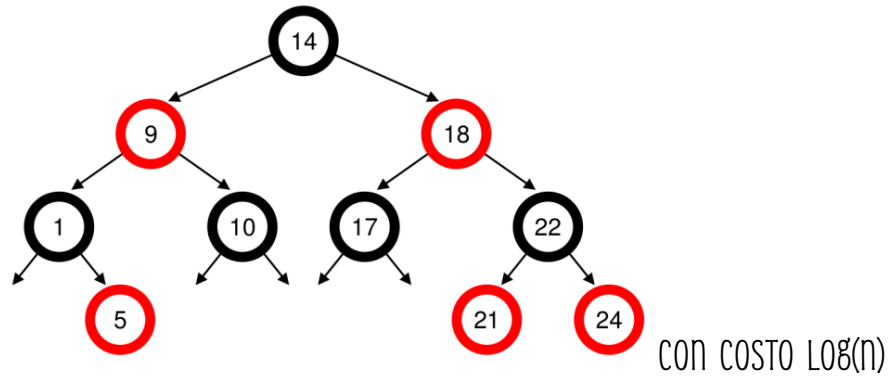
$$\text{slice}_k = \max\left(\text{sched\_latency} \cdot \frac{w_k}{\sum_{i=0}^{n-1} w_i}, \text{min\_granularity}\right)$$

$$\text{vruntime}_k += \frac{\text{runtime}_k}{w_k}$$

DOVE

```
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 */ 335, 272, 215, 172, 137,
    /* 10 */ 110, 87, 70, 56, 45,
    /* 15 */ 36, 29, 23, 18, 15,
};
```

- Per trovare efficientemente il minimo/aggiornare il vruntime dei processi pronti, essi vengono tenuti in un albero binario bilanciato di tipo rosso/nero:



- Che vruntime devo dare ad un processo nuovo?  
IL vruntime minimo che c'è nel sistema  $\rightarrow$  quindi non è "completely fair" con i processi che fanno frequentemente I/O