

Giuseppe Pelagatti

Programmazione e Struttura del sistema operativo Linux

Appunti del corso di
Architettura dei Calcolatori e Sistemi Operativi (AXO)

versione 2016

Parte N: Il Nucleo del Sistema Operativo

cap. N5 – Lo Scheduler

N.5 Lo Scheduler

1. Caratteristiche generali dello Scheduler

Il comportamento del SO è fortemente caratterizzato dalle politiche adottate per decidere quali processi eseguire e per quanto tempo eseguirli (*politiche di scheduling*).

Il componente del SO che realizza le politiche di scheduling è detto Scheduler. Il comportamento dello Scheduler è orientato a garantire le seguenti condizioni:

- che i processi più importanti vengano eseguiti prima dei processi meno importanti
- che i processi di pari importanza vengano eseguiti in maniera equa; in particolare ciò significa che nessun processo dovrebbe attendere il proprio turno di esecuzione per un tempo molto superiore agli altri.

Lo Scheduler è un componente critico nel funzionamento di un Sistema Operativo e molta ricerca viene svolta per migliorarne le prestazioni.

2. Concetti introduttivi

Dati n processi che devono essere eseguiti, la politica **Round Robin** consiste nell'assegnare ad ogni processo uno stesso *quanto* di tempo o *timeslice* in ordine circolare. Tutti i processi vengono quindi schedati per lo stesso tempo a turno. Questa politica è equa e garantisce che nessun processo rimanga bloccato per sempre (starvation).

Lo scheduler interviene in certi momenti per determinare quale processo mandare in esecuzione. La scelta del processo da mandare in esecuzione avviene nell'ambito di tutti i processi in stato di PRONTO esistenti nel sistema.

Dati 2 processi P e Q, lo scheduler deve calcolare una grandezza che determini quale dei 2 scegliere; dato che il termine priorità è usato con diversi significati particolari, noi useremo per questa grandezza il termine *diritto di esecuzione*. Lo scheduler sceglie quindi nell'ambito dei processi pronti quello che possiede il maggiore diritto di esecuzione.

I momenti in cui è significativo eseguire questa scelta sono:

- ogni volta che un processo si autosospende, per scegliere il prossimo processo da eseguire
- ogni volta che un processo viene risvegliato, perchè si estende l'insieme dei processi pronti; in questo caso se il nuovo processo diventato pronto ha un diritto di esecuzione superiore a quello corrente, il diritto di esecuzione diventa "diritto di preemption", cioè causa la sospensione del processo corrente (si ricorda però che la preemption è dilazionata fino al prossimo momento di ritorno a modo U)
- ogni volta che il processo in esecuzione è gestito con politica Round Robin e il suo timeslice è esaurito

3. Requisiti dei processi

Non tutti i processi hanno gli stessi requisiti relativamente allo scheduling. Ad un primo livello possiamo distinguere i processi nelle seguenti categorie:

1. Processi **Real-time (in senso stretto)**. Questi processi devono soddisfare dei vincoli di tempo estremamente stringenti e devono quindi essere schedati con estrema rapidità, *garantendo in ogni caso di non superare un preciso vincolo di ritardo massimo*; un esempio di processo di questo tipo è il controllo di un aereo o di un impianto.
2. Processi **semi-Real-time (soft Real time)**. Questi processi, pur richiedendo una relativa rapidità di risposta, non richiedono la garanzia assoluta di non superare un certo ritardo massimo. Un esempio di questo tipo è la scrittura di un CD o la esecuzione di un file Audio, che deve essere svolta producendo i dati con una certa continuità; se talvolta questa continuità non viene garantita si verificano dei disturbi, evento relativamente accettabile (almeno in confronto a una caduta di aereo o a un'esplosione di impianto). Tuttavia, il processo che scrive un CD deve avere un certo livello di precedenza che gli permetta normalmente di produrre i dati in tempo quando servono.
3. Processi **Normali**. Sono tutti gli altri processi. Questi processi possono avere diversi comportamenti:
 - **processi I/O bound** sono i processi che si sospendono spesso perchè hanno bisogno di I/O (ad esempio, un Text editor)
 - **processi CPU bound** sono invece i processi che tendono ad usare molto la CPU, perchè si autosospendono raramente (ad esempio, un Compilatore)

4. Classi di Scheduling

Per supportare le diverse categorie di processi lo Scheduler realizza diverse politiche di scheduling; ogni politica di scheduling è realizzata da una **Scheduler Class**. Nel descrittore di un processo il campo

```
constant struct sched_class * sched_class
```

contiene un puntatore alla struttura della scheduler class deputata a gestirlo.

Lo Scheduler è l'unico gestore delle runqueue – per questo motivo le altre funzioni devono chiedere allo scheduler di eseguire operazioni sulla runqueue. Questa scelta permette di organizzare le runqueue in maniera adatta alle diverse classi di scheduling senza dover modificare il resto del sistema.

Quando viene invocata una funzione dello scheduler (cfr. capitolo precedente, “interazione con lo scheduler”) tale funzione svolge poche funzioni preliminari e poi invoca la corrispondente funzione della scheduler class alla quale il task appartiene. Ad esempio, per i processi normali attualmente la scheduler class ha la struttura (semplificata) di figura 1, definita nel file `sched_fair.c` (che contiene l'implementazione del Completely Fair Scheduler – CFS – descritto più avanti).

```
1106 static const struct sched_class fair_sched_class = {
1107     .next                = &idle_sched_class,
1108     .enqueue_task        = enqueue_task_fair,
1109     .dequeue_task        = dequeue_task_fair,
1112     .check_preempt_curr  = check_preempt_wakeup,
1114     .pick_next_task      = pick_next_task_fair,
1115     .put_prev_task       = put_prev_task_fair,
1122     .set_curr_task       = set_curr_task_fair,
1123     .task_tick           = task_tick_fair,
1124     .task_new            = task_new_fair,
1125 };
```

Figura 4.1

Se quindi una funzione del nucleo invoca `enqueue_task`, all'interno di questa funzione verrà invocata la funzione corrispondente della `sched_class` del processo in esecuzione `p` nel modo seguente:

```
p->sched_class->enqueue_task
```

e, se, ad esempio, la `sched_class` del processo è `fair_sched_class`, la funzione effettivamente eseguita sarà `enqueue_task_fair`. In questo modo è possibile aggiungere nuove classi di scheduling al sistema senza doverlo modificare eccessivamente.

Ogni classe può a sua volta implementare più di una politica.

Attualmente le 3 classi più importanti supportate

1. `SCHED_FIFO` (First IN First Out)
2. `SCHED_RR` (Round Robin)
3. `SCHED_NORMAL`

L'ordine di queste 3 classi è un ordine anche dei diritti di esecuzione; un processo della classe `SCHED_FIFO` ha un diritto di esecuzione sempre superiore a quello di un processo di una delle due altre classi e un processo della classe `SCHED_RR` ha un diritto di esecuzione sempre superiore a un processo della classe `SCHED_NORMAL`.

La funzione `schedule()` invoca la funzione `pick_next_task()`, che a sua volta invoca le funzioni `pick_next_task()` specifiche di ogni singola classe in ordine di importanza delle classi per selezionare un nuovo task; `schedule()` poi procede al context switch utilizzando le funzioni della classe appropriata per togliere dall'esecuzione il task corrente (`prev`) e mettere in esecuzione il nuovo (`next`).

```

schedule( ){
    ...
    struct tsk_struct * prev;
    prev = CURR;
    if (prev->stato == ATTESA) {
        // toglì il task corrente prev dalla runqueue rq
        dequeue (prev);
    } /* if */

    // in modo analogo, se il task corrente (prev) è terminato dequeue (prev);

    //invoca la funzione di scelta del prossimo task
    next = pick_next_task(rq, prev);
    //se next è diverso da prev, esegui il context switch
    if (next != prev) {
        context_switch(prev, next);
        CURR->START = NOW // istante corrente salvato per prox tick
    }
    TIF_NEED_RESCHED = 0;
}

pick_next_task(rq, prev){
    for(ciascuna classe di scheduling, in ordine di importanza decrescente)
    {
        //invoca funzione di scelta del prossimo task per la classe in esame
        next = class->pick_next_task(rq, prev);
        if (next != NULL) return next; //appena trovi un task, ritorna
    }
    //pick_next_task restituisce sempre un puntatore valido:
    //a un task PRONTO con diritto di esecuzione massimo, se esiste
    //a prev, se non ci sono altri task pronti e il suo stato non è ATTESA
    //a IDLE, se nessuno dei 2 casi precedenti è praticabile
}

```

5. Scheduling dei Processi soft Realtime

Le classi SCHED_FIFO e SCHED_RR sono utilizzate per supportare i processi soft RT (Linux non supporta i processi RT in senso stretto, perchè non è in grado di garantire il non superamento di un ritardo massimo). In queste due classi il concetto fondamentale è quello di priorità statica.

A ogni processo di queste classi viene attribuita una priorità detta statica, perchè è attribuita all'inizio e non varia mai. I valori delle priorità statiche appartengono all'intervallo da 1 a 99 (99 è la più alta).

La priorità statica è memorizzata in `task_struct` nel campo `static_prio`.

Scheduling SCHED_FIFO

Un task di questa categoria viene eseguito senza alcun limite di tempo. Se esistono diversi task di questa categoria, sono schedulati in base alla priorità statica – un task con maggiore priorità ha un diritto di esecuzione maggiore di uno a priorità più bassa e può quindi causarne la preemption.

Scheduling SCHED_RR

I task in questa categoria sono simile ai precedenti, però nell'ambito della stessa fascia di priorità sono schedulati con una politica Round Robin; pertanto se esistono diversi task in questa categoria allo stesso livello di priorità, ognuno di loro viene eseguito per un timeslice.

6. Scheduling dei processi normali (SCHED_NORMAL)

6.1. Aspetti generali

I processi normali vengono presi in considerazione dallo Scheduler solo se non ci sono processi delle classi FIFO e RR in stato di PRONTO per l'esecuzione. L'attuale scheduler dei processi normali è chiamato ambiziosamente **Completely Fair Scheduler (CFS)**, perchè ambisce a simulare il seguente meccanismo ideale: dati N task, dedicare una CPU di potenza $1/N$ ad ogni task. Naturalmente il meccanismo ideale non è realizzabile con un numero di CPU inferiore a N , quindi in pratica la CPU (che qui supporremo unica per semplificare il discorso) deve essere assegnata per un **quanto** di tempo ad ogni task.

I problemi fondamentali che lo scheduler deve risolvere nel far questo sono:

1. *determinare ragionevolmente la durata dei quanti* (quanti troppo lunghi abbassano la responsività del sistema, quanti troppo corti causano troppo sovraccarico per i numerosi context switch)
2. permettere di *assegnare dei pesi ai processi*, in modo che ai processi più importanti sia assegnato più tempo che a quelli meno importanti
3. permettere ai processi che stanno a lungo in stato di ATTESA di tornare rapidamente in esecuzione quando vengono risvegliati

L'ultimo requisito è legato al seguente problema generato dallo scheduling basato sulla politica Round Robin: in assenza di meccanismi correttivi i processi I/O bound tenderebbero a funzionare molto male; ad esempio, si considerino un Editor di testo e un Compilatore. Ambedue sono processi normali. Tuttavia, dato che l'Editor è molto interattivo (I/O bound) tenderebbe ad eseguire per poco tempo e poi a sospendersi, mentre il compilatore (CPU bound) ogni volta che va in esecuzione tenderebbe ad eseguire per l'intero tempo a sua disposizione

6.2 Meccanismo fondamentale

Per semplificare la presentazione iniziale del meccanismo del CFS assumiamo che valgano le seguenti ipotesi:

1. tutti i processi hanno peso unitario: $t.LOAD = 1$ per tutti i task t presenti nella runqueue
2. i task non si autosospendano mai (non eseguono wait)

Sotto queste ipotesi il meccanismo base ha una logica molto semplice. Sia **NRT** il numero di task presenti nella runqueue a un certo istante:

- viene determinato un periodo **PER** di schedulazione durante il quale tutti i task della runqueue devono essere eseguiti
- ad ogni task viene assegnato un quanto di tempo $Q = PER/NRT$

I task vengono mantenuti in una coda ordinata **RB** e il funzionamento dello scheduler può essere schematizzato nel modo seguente:

1. viene estratto da RB (e reso CURR) il primo task della coda
2. CURR viene eseguito fino a quando scade il suo quanto di tempo Q
3. CURR viene sospeso e reinserito in RB in fondo alla coda
4. torna al passo 1

In pratica i task sono eseguiti a turno per esattamente PER/NRT ms; si osservi che il periodo di schedulazione deve essere considerato come una finestra che scorre nel tempo; non c'è una suddivisione rigida del tempo in periodi, ma in ogni istante si può fare riferimento all'inizio di un nuovo periodo di schedulazione. In altre parole, in un momento qualsiasi si può asserire che entro i prossimi PER ms tutti i task verranno eseguiti.

ATTENZIONE: da un punto di vista sintattico nelle formule si è spesso utilizzata la notazione $p.c$ (selezione di un campo c di una struttura p) anche quando sarebbe necessario utilizzare $p->v$ (perché p è un puntatore alla struttura che contiene c); questo semplifica la lettura, ma è scorretto se interpretato rigorosamente come codice C.

Determinazione del periodo di schedulazione PER

Il periodo di schedulazione varia dinamicamente con l'aumento o la diminuzione del numero di task presenti nella runqueue, non può quindi essere fissato rigidamente, ma è necessario cercare una mediazione tra i seguenti aspetti:

- un periodo troppo lungo può ritardare eccessivamente l'esecuzione di un processo
- un periodo troppo corto può produrre quanti eccessivamente corti al crescere di NRT, portando a commutazioni di contesto troppo frequenti

L'impostazione attualmente adottata da Linux consiste nel basare il calcolo di PER su due parametri di controllo (**SYSCTL** parameter) modificabili dall'amministratore del sistema:

- **LT** (latenza) – default: 6ms

- **GR** (granularità) – default: 0,75ms

Il periodo è calcolato con la seguente formula:

$$\text{PER} = \text{MAX}(\text{LT}, \text{NRT} * \text{GR})$$

In questo modo il quanto di un processo è LT/NRT fino a quando $\text{NRT} * \text{GR}$ è minore della latenza; quando $\text{NRT} * \text{GR}$ supera il valore LT (per i valori di default questo accade con $\text{NRT} > 8$) il periodo viene allungato in modo da evitare che il quanto scenda sotto il valore di granularità GR .

Adattamento del quanto ai pesi dei task

La formula $Q = \text{PER}/\text{NRT}$ vista sopra non tiene conto del peso dei task; in realtà il valore del quanto di tempo $t.Q$ assegnato a un task è proporzionale al suo peso rispetto al peso di tutti i task. Il calcolo del quanto utilizza le seguenti due variabili:

- RQL (rqload) è la somma dei pesi di tutti i task presenti sulla runqueue
- LC (load_coeff) il rapporto tra il peso di un task e RQL ;

$$\text{LC} = t.\text{LOAD}/\text{RQL}$$

Il quanto assegnato a un task t è calcolato nel modo seguente:

$$t.Q = \text{PER} * t.LC = (\text{PER}/\text{RQL}) * t.LOAD$$

Nell'ipotesi che tutti i task abbiano peso unitario queste formule si riducono alla $Q = \text{PER}/\text{NRT}$ indicata all'inizio; infatti $Q = \text{PER} * \text{LC} = \text{PER} * t.LOAD/\text{RQL} = \text{PER} * 1/(\text{NRT} * 1)$

Il Virtual Time

Per mantenere l'ordinamento dei task nella coda RB il CFS usa il concetto di **Virtual Runtime (VRT)**. Il VRT è una misura del tempo di esecuzione consumato da un processo, basato sulla modificazione del tempo reale in base a opportuni coefficienti, in modo che la decisione su quale sia il prossimo task da eseguire possa basarsi semplicemente sulla scelta del processo con VRT minimo; dato che RB è ordinato in base ai VRT dei task, il prossimo task da mettere in esecuzione sarà il primo di RB e viene indicato con **LFT** (leftmost) (si ricorda che il task in esecuzione non è contenuto in RB ma puntato dalla variabile CURR). Quando un task termina l'esecuzione viene reinserito nella coda RB nella posizione che gli compete in base al nuovo valore del VRT che possiede.

Consideriamo un task che ha eseguito per **DELTA** nsec e che abbandona l'esecuzione; lo scheduler incrementa in quel momento il suo tempo di esecuzione **SUM** e il suo **VRT**.

Il tempo di esecuzione viene semplicemente incrementato di **DELTA**:

$$\text{SUM} = \text{SUM} + \text{DELTA}$$

invece l'incremento del VRT viene corretto con un coefficiente **VRTC** (vrt_coeff) nel modo seguente:

$$t.VRTC = 1/t.LOAD$$

$$t.VRT = t.VRT + \text{DELTA} * t.VRTC$$

Si noti che mentre LC è proporzionale al peso di un processo VRTC è inversamente proporzionale a tale peso.

L'effetto del coefficiente VRTC è di far crescere il VRT dei processi più pesanti più lentamente del VRT dei processi più leggeri. In condizioni di funzionamento stabile (cioè un numero di processi costante, con tutti i processi che consumano tutto il loro quanto) la crescita del VRT di tutti i processi in un periodo è la stessa; infatti l'effetto di un quanto più lungo è esattamente compensato da una crescita più lenta del VRT; infatti dopo un quanto di tempo l'incremento di VRT è

$$Q * \text{VRTC} = (\text{PER} / \text{RQL}) * t.LOAD * (1 / t.LOAD) = \text{PER} / \text{RQL}$$

Come si vede, la variazione di VRT non dipende dal peso del processo.

Infine, per motivi che emergeranno più avanti, lo scheduler mantiene nella runqueue una variabile **VMIN** che rappresenta il

VRT minimo tra tutti i task presenti nella runqueue; questa variabile è aggiornata continuamente in base alla regola:

VMIN = MIN(CURR.VRT, LFT.VRT) //provvisoria

dove CURR.VRT è il VRT del task in esecuzione che viene anch'esso aggiornato continuamente. *Questa formula dovrà essere modificata* per rispondere ad alcuni problemi discussi più avanti.

Possiamo adesso interpretare lo pseudocodice seguente, relativo alla funzione `tick` invocata periodicamente in base agli interrupt del clock, dove:

- NOW è l'istante corrente
- START è l'istante in cui un task viene messo in esecuzione
- PREV è il valore di SUM al momento in cui un task viene messo in esecuzione

```
tick() {
    //aggiornamento dei parametri di CURR:
    DELTA = NOW - CURR->START;
    CURR->SUM = CURR->SUM + DELTA;
    CURR->VRT = CURR->VRT + DELTA * CURR.VRTC;
    CURR->START = NOW;
    //aggiornamento di VMIN della RQ
    VMIN = MIN(CURR->VRT,LFT->VRT)); //questa formula verrà corretta
    //controllo se è scaduto il quanto di tempo
    if ((CURR->SUM - CURR->PREV) >= Q) resched( );
}
```

La funzione `schedule()` invoca `pick_next_task`, che eventualmente arriva ad invocare `pick_next_task_fair()` se non esistono task pronti nelle classi di priorità superiore.

La funzione `pick_next_task_fair()`, che sceglie il nuovo task di classe `SCHED_NORMAL` e che viene invocata solo se le altre classi di scheduling non hanno restituito un task da mettere in esecuzione, svolge essenzialmente l'assegnazione a CURR del primo task (LFT) della coda RB; nel caso particolare in cui RB è vuota tenta di far proseguire CURR, ma, se CURR è stato eliminato dalla RUNQUEUE (ad esempio, perché ha eseguito una EXIT oppure una WAIT), mette in esecuzione IDLE.

Lo pseudocodice è il seguente:

```
pick_next_task_fair(rq, previous_task){
    if (LFT != NULL) {
        //RB non è vuoto
        CURR = LFT;
        //elimina LFT da RB ristrutturando la coda opportunamente
        CURR->PREV = CURR->SUM; //salva in PREV il valore di SUM
    }
    else {
        //il RB è vuoto
        if (CURR == NULL) // CURR è stato eliminato perchè in ATTESA
            CURR = IDLE;
        //altrimenti CURR non viene modificato
    }
    //a questo punto CURR può essere uguale a LFT precedente,
    //oppure a IDLE,
    //oppure al precedente CURR (cioè non è stato modificato)
    return CURR;
}
```


Esempio/Esercizio 1

In Figura 6.1 è presentato il formato in cui rappresenteremo la simulazione dello scheduling. La figura contiene 2 Tabelle: la prima serve a rappresentare le condizioni iniziali (da completare) e la seconda serve a rappresentare un evento.

Si deve completare la Tabella delle condizioni iniziali e compilare una Tabella di evento per ogni evento che si verifica *fino (incluso) al primo evento che avviene dopo un dato Limite di Simulazione*.

Per ogni evento significativo sono rappresentate:

1. Le caratteristiche dell'evento, cioè
 - a. il momento in cui si è verificato, come tempo trascorso dall'istante iniziale della simulazione
 - b. il tipo di evento (al momento consideriamo solamente Q_SCADE, cioè la scadenza del quanto di tempo)
 - c. il task nel cui contesto l'evento si è verificato
 - d. una variabile booleana RESCHEDULE che indica la necessità di rischedulare dopo aver servito l'evento (se il tipo di evento è Q_SCADE tale necessità è sempre vera)
2. I parametri globali della RUNQUEUE
3. I parametri dei diversi task, nell'ordine in cui sono memorizzati nel RB preceduti dal task corrente

In questo esempio la condizione iniziale è caratterizzata dai seguenti aspetti:

- esistono 3 task t1, t2 e t3, tutti di peso unitario e con coefficienti LC, VRTC identici
- pertanto NRT = 3; RQL = 3;
- quindi per ogni task LC = 0,33 e VRTC = 1
- anche i quanti di tempo sono identici ($PER * LC = 2 \text{ ms}$)
- i 3 task sono in esecuzione da tempi diversi (diversi SUM) ma i loro VRT sono abbastanza simili
- il task t1 è quello corrente;
- VMIN è il valore di t1.VRT, cioè di CURRENT

CONDIZIONI INIZIALI *****										
RUNQUEUE - NRT			PER		RQL		CURR		VMIN	
3			6,00		3,00		t1		100,00	
TASKS:	ID		LOAD		LC		Q		VRTC	SUM VRT
CURRENT	t1		1							10,00 100,00
RB TREE	t2		1							20,00 100,50
	t3		1							30,00 101,00

EVENT *****	TIME		TYPE		CONTEXT		RESCHEDULE	
RUNQUEUE - NRT			PER		RQL		CURR	VMIN
TASKS:	ID		LOAD		LC		Q	VRTC SUM VRT
CURRENT								
RB TREE								

Figura 6.1

In Figura 6.2 è mostrata la soluzione completa con Limite di Simulazione = 7 ms. Gli unici eventi che si verificano sono quelli di scadenza dei quanti (Q_SCADE). Tali eventi si verificano ogni 2ms e, come si vede in Figura 6.2, i task accumulano ad ogni esecuzione 2ms sia nella variabile SUM che nel VRT.

Dopo 6 ms la situazione è identica a quella iniziale, a parte ovviamente l'aumento dei tempi di esecuzione e dei VRT.

La simulazione si arresta con il primo evento che si verifica oltre i 7ms.

Nello svolgimento della simulazione è necessario considerare la seguente sequenza di operazioni:

1. determinare il prossimo evento che si deve verificare, il suo tipo e il suo tempo
2. eseguire gli aggiornamenti delle variabili modificate dalla funzione tick, cioè
 - a. curr.SUM
 - b. curr.VRT
 - c. VMIN
3. svolgere le operazioni richieste dal tipo di evento

Ad esempio, per compilare il primo evento di Figura 6.2 si procede nel modo seguente:

1. prossimo evento è la scadenza del quanto di t1 dopo 2ms dall'evento precedente (situazione iniziale, in questo caso), quindi tipo = Q_SCADE, tempo = 0 + 2 = 2ms
2. esecuzione di tick:
 - a. t1.SUM = 12
 - b. t1.VRT = 102
 - c. VMIN = min(102, 100.5) = 100.5
3. operazioni richieste dall'evento
 - a. t2 diventa CURRENT
 - b. t1 viene inserito nella posizione che gli compete nel RB (in questo caso alla fine) con i nuovi valori

CONDIZIONI INIZIALI *****								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t1 100,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t1 1.0 0,33 2,00 1,00 10,00 100,00								
RB TREE t2 1.0 0,33 2,00 1,00 20,00 100,50								
t3 1.0 0,33 2,00 1,00 30,00 101,00								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
2,00 Q_SCADE t1 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t2 100,50								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t2 1.0 0,33 2,00 1,00 20,00 100,50								
RB TREE t3 1.0 0,33 2,00 1,00 30,00 101,00								
t1 1.0 0,33 2,00 1,00 12,00 102,00								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
4,00 Q_SCADE t2 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t3 101,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t3 1.0 0,33 2,00 1,00 30,00 101,00								
RB TREE t1 1.0 0,33 2,00 1,00 12,00 102,00								
t2 1.0 0,33 2,00 1,00 22,00 102,50								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
6,00 Q_SCADE t3 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t1 102,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t1 1.0 0,33 2,00 1,00 12,00 102,00								
RB TREE t2 1.0 0,33 2,00 1,00 22,00 102,50								
t3 1.0 0,33 2,00 1,00 32,00 103,00								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
8,00 Q_SCADE t1 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t2 102,50								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t2 1.0 0,33 2,00 1,00 22,00 102,50								
RB TREE t3 1.0 0,33 2,00 1,00 32,00 103,00								
t1 1.0 0,33 2,00 1,00 14,00 104,00								

Figura 6.2

Esercizio 2.

Consideriamo ora il caso in cui i task hanno pesi diversi, come indicato nelle condizioni iniziali di Figura 6.3
L'esempio è simile al precedente, ma con pesi diversi per i 3 task.

CONDIZIONI INIZIALI *****									
RUNQUEUE - NRT PER RQL CURR VMIN									
3 6,00 3,00 t1 100,00									
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT		
CURRENT	t1	1.0				10,00	100,00		
RB TREE	t2	1.5				20,00	100,50		
	t3	0.5				30,00	101,00		

Figura 6.3

La soluzione è presentata in Figura 6.4

Si nota:

- i pesi indicati producono un forte squilibrio nell'esecuzione dei task; il quanto di t2 è il triplo del quanto di t1
- il VRT dei 3 task cresce invece di quantità identiche per ogni periodo di schedulazione (2 ms), quindi l'ordinamento di esecuzione si mantiene inalterato, consolidando il vantaggio dei task più pesanti, che ad ogni ciclo beneficiano di un quanto più lungo
- viene comunque garantito al task più leggero di eseguire almeno una volta per ogni periodo

CONDIZIONI INIZIALI *****								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t1 100,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t1 1.0 0,33 2,00 1,00 10,00 100,00								
RB TREE t2 1.5 0,50 3,00 0,67 20,00 100,50								
t3 0.5 0,17 1,00 2,00 30,00 101,00								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
2,00 Q_SCADE t1 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t2 100,50								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t2 1.5 0,50 3,00 0,67 20,00 100,50								
RB TREE t3 0.5 0,17 1,00 2,00 30,00 101,00								
t1 1.0 0,33 2,00 1,00 12,00 102,00								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
5,00 Q_SCADE t2 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t3 101,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t3 0.5 0,17 1,00 2,00 30,00 101,00								
RB TREE t1 1.0 0,33 2,00 1,00 12,00 102,00								
t2 1.5 0,50 3,00 0,67 23,00 102,50								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
6,00 Q_SCADE t3 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t1 102,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t1 1.0 0,33 2,00 1,00 12,00 102,00								
RB TREE t2 1.5 0,50 3,00 0,67 23,00 102,50								
t3 0.5 0,17 1,00 2,00 31,00 103,00								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
8,00 Q_SCADE t1 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t2 102,50								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t2 1.5 0,50 3,00 0,67 23,00 102,50								
RB TREE t3 0.5 0,17 1,00 2,00 31,00 103,00								
t1 1.0 0,33 2,00 1,00 14,00 104,00								

Figura 6.4

6.3 Gestione di wait e wakeup

L'idea di base del CFS per soddisfare il requisito 3, cioè garantire una risposta tempestiva dopo un wakeup si basa sul meccanismo fondamentale appena visto: dato che un processo in ATTESA non esegue, il suo VRT resta indietro rispetto agli altri, quindi al risveglio andrà in una posizione favorevole nel RB.

L'esecuzione di un evento di WAIT o di WAKEUP modifica il numero di processi presenti nella runqueue e quindi *richiede di ricalcolare* NRT e RQL e, per ogni task della runqueue il coefficiente LC e il quanto Q.

wakeup di un processo tw

Quando un processo viene risvegliato è possibile che il suo VRT sia molto basso (perchè è in ATTESA da lungo tempo) oppure sia ancora relativamente alto (ATTESA breve).

Il nuovo VRT che gli viene assegnato è:

$$tw.VRT = \text{MAX}(tw.VRT, (VMIN - LT/2))$$

In base a questa assegnazione il processo risvegliato parte con un valore di VRT che lo candida ad essere eseguito nel prossimo futuro, ma non gli dà un credito tale da distorcere completamente il sistema (come accadrebbe assegnando direttamente $tw.VRT$ a un task che è stato in attesa per un lungo periodo).

Tipicamente, se il processo ha fatto un'attesa molto breve gli viene lasciato il suo VRT.

Inoltre, la formula precedente mostra che un task risvegliato può avere un VRT inferiore a VMIN e quindi al prossimo aggiornamento di VMIN da parte della funzione `tick`, VMIN scenderebbe. Dato che è opportuno che *VMIN cresca in maniera monotona*, la formula vista prima per l'assegnamento di VMIN nella funzione `tick` deve essere modificata:

$$VMIN = \text{MAX}(VMIN, \text{MIN}(CURR.VRT, LFT.VRT));$$

In questo modo se esiste un task con VRT inferiore a VMIN, il suo valore non produce un assegnamento di VMIN.

Necessità di rescheduling

Gli eventi di WAIT richiedono di eseguire sempre un rescheduling; invece nel caso di WAKEUP per decidere se il nuovo processo deve causare un rescheduling vengono valutati due aspetti:

1. se il processo appartiene a una classe di scheduling superiore,
2. se il suo VRT è inferiore al VRT del processo corrente

La seconda valutazione viene però modificata con un correttivo che serve ad evitare che un processo che esegue attese brevissime causi commutazioni di contesto troppo frequenti; la formula applicata è contenuta nel seguente statement della funzione `check_preempt_curr(tw ...)`, invocata da `wakeup()`:

```
if ( (tw->schedule_class == classe con diritto > NORMAL) or
    ((tw->vrt + WGR * tw->load_coeff) < CURR->vrt ) ) resched();
```

dove WGR (wakeup granularity) è un parametro di configurazione (SYSCTL) con default 1ms.

Esercizio 3

Per simulare anche gli eventi di WAIT e di WAKEUP aggiungiamo alle condizioni iniziali l'indicazione degli eventi che si verificheranno. Per ogni coppia di eventi di WAIT/WAKEUP viene indicato:

- il task al quale si riferiscono
- il tipo dell'evento
- il tempo al quale si verificherà; l'indicazione del tempo è di 2 tipi:
 - se l'evento è WAIT il tempo indicato è il **tempo di esecuzione** dopo il quale il Task causa l'evento
 - se l'evento è WAKEUP il tempo è quello che intercorre tra l'evento di WAIT e il successivo risveglio tramite evento di WAKEUP (indipendentemente dal task in esecuzione)

In Figura 6.5a sono mostrate le stesse condizioni iniziali dell'esercizio 1, ma con l'aggiunta dei seguenti eventi:

- il task t1 si pone in ATTESA dopo 1 ms di esecuzione e verrà risvegliato dopo 5 ms
- il task t3 si pone in ATTESA dopo 3 ms di esecuzione e verrà risvegliato dopo 1 ms

CONDIZIONI INIZIALI *****									
RUNQUEUE - NRT PER RQL CURR VMIN									
3 6,00 3,00 t1 100,00									
TASKS: ID LOAD LC Q VRTC SUM VRT									
CURRENT t1 1 10,00 100,00									
RB TREE t2 1 20,00 100,50									
t3 1 30,00 101,00									
Events of task t1: WAIT at 1.0; WAKEUP after 5.0;									
Events of task t3: WAIT at 3.0; WAKEUP after 1.0;									
Figura 6.5									

L'esecuzione fino al limite di 11ms è mostrata nelle seguenti figure, nelle quali sono evidenziati gli aspetti nuovi.

Dopo 1ms di esecuzione il task t1 genera un evento di wait con i seguenti effetti (vedi figura 6.6a):

- i nuovi valori calcolati da tick sono: t1.SUM = 11; t1.VRT = 101; VMIN = 100.5
- il task t1 è spostato in fondo alla tabella che l'indicazione "WAITING" (si ricorda che il task non appartiene più alla runqueue ma è inserito in una waitqueue – lo riportiamo per comodità, perché dovremo recuperarlo al momento del WAKEUP)
- i task presenti nella runqueue sono diminuiti (si ricorda che i task in attesa non appartengono alla runqueue) e quindi vengono ricalcolati i parametri NRT e RQL della runqueue e i parametri LC e Q di ogni task
- viene messo in esecuzione il primo task di RB (t2)

CONDIZIONI INIZIALI *****									
RUNQUEUE - NRT PER RQL CURR VMIN									
3 6,00 3,00 t1 100,00									
TASKS: ID LOAD LC Q VRTC SUM VRT									
CURRENT t1 1.0 0,33 2,00 1,00 10,00 100,00									
RB TREE t2 1.0 0,33 2,00 1,00 20,00 100,50									
t3 1.0 0,33 2,00 1,00 30,00 101,00									
Events of task t1: WAIT at 1.0; WAKEUP after 5.0;									
Events of task t3: WAIT at 3.0; WAKEUP after 1.0;									

EVENT *****									
TIME TYPE CONTEXT RESCHEDULE									
1,00 WAIT t1 true									
RUNQUEUE - NRT PER RQL CURR VMIN									
2 6,00 2,00 t2 100,50									
TASKS: ID LOAD LC Q VRTC SUM VRT									
CURRENT t2 1.0 0,50 3,00 1,00 20,00 100,50									
RB TREE t3 1.0 0,50 3,00 1,00 30,00 101,00									
WAITING t1 1.0 0,33 2,00 1,00 11,00 101,00									

Figura 6.6a

Il task t2 esegue normalmente per 3ms fino allo scadere del proprio Q (t globale = 4ms), poi viene messo in esecuzione t3; dopo 2ms di esecuzione di t3 si verifica il WAKEUP di t1 (figura 6.6b). Infatti sono passati 5ms dal WAIT (3ms di esecuzione di t2 e 2ms di esecuzione di t3).

I valori calcolati da tick sono: $t3.SUM = 32$; $t3.VRT = 103$; $V_{MIN} = \max(101, \min(103, 103.5)) = 103$

Questo evento reinserisce t1 nella runqueue con VRT determinato dalla formula

$$tw.VRT = \max(tw.VRT, (V_{MIN} - LT/2))$$

quindi, dato che V_{MIN} vale 103, $LT/2$ vale 3 e $t1.VRT$ vale 101, $t1.VRT$ viene lasciato a 101.

A questo punto si deve decidere se rischedulare; la risposta è affermativa, perché la condizione riportata e valorizzata nella figura è vera. Lo scheduler trova che il task t1 ha VRT minimo e lo pone in esecuzione.

Si noti che t3 viene inserito nel RB prima di t2, perché il suo VRT è inferiore.

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		4,00		Q_SCADE		t2		true
RUNQUEUE -	NRT	PER	RQL	CURR	VMIN			
	2	6,00	2,00	t3	101,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t3	1.0	0,50	3,00	1,00	30,00	101,00	
RB TREE	t2	1.0	0,50	3,00	1,00	23,00	103,50	
WAITING	t1	1.0	0,33	2,00	1,00	11,00	101,00	

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		6,00		WAKEUP		t3		true
				$tw.vrt+WGR*tw.LC=101,00+1,00*0,33=101,33 < curr.vrt=103,00$				
RUNQUEUE -	NRT	PER	RQL	CURR	VMIN			
	3	6,00	3,00	t1	103,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t1	1.0	0,33	2,00	1,00	11,00	101,00	
RB TREE	t3	1.0	0,33	2,00	1,00	32,00	103,00	
	t2	1.0	0,33	2,00	1,00	23,00	103,50	

Figura 6.6b

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		8,00		Q_SCADE		t1		true
RUNQUEUE -	NRT	PER	RQL	CURR	VMIN			
	3	6,00	3,00	t3	103,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t3	1.0	0,33	2,00	1,00	32,00	103,00	
RB TREE	t1	1.0	0,33	2,00	1,00	13,00	103,00	
	t2	1.0	0,33	2,00	1,00	23,00	103,50	

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		9,00		WAIT		t3		true
RUNQUEUE -	NRT	PER	RQL	CURR	VMIN			
	2	6,00	2,00	t1	103,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t1	1.0	0,50	3,00	1,00	13,00	103,00	
RB TREE	t2	1.0	0,50	3,00	1,00	23,00	103,50	
WAITING	t3	1.0	0,33	2,00	1,00	33,00	104,00	

Figura 6.6c

Il task t1 esegue per tutto il suo quanto di 2 ms, fino al tempo 8, poi subentra t3 che genera un WAIT dopo 1 ms (Figura 6.6c); il task t3 infatti ha eseguito per 3ms dall'inizio (SUM era 30 inizialmente ed è diventato 33). Dopo 1ms ulteriore si verifica il WAKEUP di t3 (figura 6.6d); da notare che la condizione di rescheduling è falsa e quindi il task corrente (t1) prosegue la sua esecuzione fino allo scadere del quanto. T3 invece, pur essendo stato appena risvegliato, viene inserito in fondo alla coda RB.

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE					
		10,00		WAKEUP		t1		false					
tw.vrt+WGR*tw.LC = 104,00+1,00*0,33=104,33 < curr.vrt=104,00													
RUNQUEUE	-	NRT		PER		RQL		CURR		VMIN			
		3		6,00		3,00		t1		103,50			
TASKS:	ID		LOAD		LC		Q		VRTC		SUM		VRT
CURRENT	t1		1.0		0,33		2,00		1,00		14,00		104,00
RB TREE	t2		1.0		0,33		2,00		1,00		23,00		103,50
	t3		1.0		0,33		2,00		1,00		33,00		104,00

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE					
		11,00		Q_SCADE		t1		true					
RUNQUEUE - NRT		PER		RQL		CURR		VMIN					
3		6,00		3,00		t2		103,50					
TASKS:	ID		LOAD		LC		Q		VRTC		SUM		VRT
CURRENT	t2		1.0		0,33		2,00		1,00		23,00		103,50
RB TREE	t3		1.0		0,33		2,00		1,00		33,00		104,00
	t1		1.0		0,33		2,00		1,00		15,00		105,00

Figura 6.6d

Esercizio 4

Consideriamo la stessa situazione dell'esempio precedente, con il task t1 che si pone in ATTESA dopo 1 ms di esecuzione. In questo caso però il task t1 viene risvegliato 0,6 ms dopo che è andato in ATTESA. Mostriamo solamente la simulazione fino al risveglio di t1 (Figura 6.7)

Viene calcolato il valore di VRT da assegnare al task risvegliato t1:

$$t1.VRT = \text{MAX}(101, (101 - 3)) = 101$$

Infine viene valutata la necessità di rischedulare che risulta falsa, dato che

$$t1.VRT + WGR * t1.LC = 101 + (1*0,33) = 101,33 > t2.VRT = 101,1$$

In questo esempio la condizione di rescheduling è falsa grazie all'effetto di WGR, perchè senza WGR avremmo avuto

$$t1.VRT = 101 < t2.VRT = 101,1$$

cioè una condizione vera.

CONDIZIONI INIZIALI *****								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t1 100,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t1 1.0 0,33 2,00 1,00 10,00 100,00								
RB TREE t2 1.0 0,33 2,00 1,00 20,00 100,50								
t3 1.0 0,33 2,00 1,00 30,00 101,00								
Events of task t1: WAIT at 1.0; WAKEUP after 0.6;								
Events of task t3: WAIT at 3.0; WAKEUP after 1.0;								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
1,00 WAIT t1 true								
RUNQUEUE - NRT PER RQL CURR VMIN								
2 6,00 2,00 t2 100,50								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t2 1.0 0,50 3,00 1,00 20,00 100,50								
RB TREE t3 1.0 0,50 3,00 1,00 30,00 101,00								
WAITING t1 1.0 0,33 2,00 1,00 11,00 101,00								

EVENT ***** TIME TYPE CONTEXT RESCHEDULE								
1,60 WAKEUP t2 false								
$tw.vrt+WGR*tw.LC=101,00+1,00*0,33=101,33 < curr.vrt=101,10$								
RUNQUEUE - NRT PER RQL CURR VMIN								
3 6,00 3,00 t2 101,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t2 1.0 0,33 2,00 1,00 20,60 101,10								
RB TREE t3 1.0 0,33 2,00 1,00 30,00 101,00								
t1 1.0 0,33 2,00 1,00 11,00 101,00								

Figura 6.7

Esercizio 5

Nel seguente esercizio (figura 6.8) viene messo in esecuzione IDLE. Per IDLE i parametri non sono significativi e possono essere omissi.

Soluzione con limite di simulazione 12ms in figura 6.9

CONDIZIONI INIZIALI *****								
RUNQUEUE - NRT PER RQL CURR VMIN								
2 6,00 2,00 t1 100,00								
TASKS: ID LOAD LC Q VRTC SUM VRT								
CURRENT t1 1.0 0,50 3,00 1,00 10,00 100,00								
RB TREE t2 1.0 0,50 3,00 1,00 20,00 101,00								
Events of task t1: WAIT at 1.0; WAKEUP after 10.0;								
Events of task t2: WAIT at 1.0; WAKEUP after 6.0;								

Figura 6.8

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		1,00		WAIT		t1		true
RUNQUEUE -	NRT	PER	RQL	CURR	VMIN			
	1	6,00	1,00	t2	101,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t2	1.0	1,00	6,00	1,00	20,00	101,00	
RB	VUOTO							
WAITING	t1	1.0	0,50	3,00	1,00	11,00	101,00	

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		2,00		WAIT		t2		true
RUNQUEUE -	NRT	PER	RQL	CURR	VMIN			
	0	6,00	0,00	IDLE	102,00			
CURRENT	is	IDLE						
RB	VUOTO							
WAITING	t2	1.0	1,00	6,00	1,00	21,00	102,00	
WAITING	t1	1.0	0,50	3,00	1,00	11,00	101,00	

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		8,00		WAKEUP		IDLE		true
RUNQUEUE -	NRT	PER	RQL	CURR	VMIN			
	2	6,00	1,00	t2	102,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t2	1.0	1,00	6,00	1,00	21,00	102,00	
RB	VUOTO							
WAITING	t1	1.0	0,50	3,00	1,00	11,00	101,00	

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		11,00		WAKEUP		t2		true
		$tw.vrt+WGR*tw.LC=102,00+1,00*0,50=102,50 < curr.vrt=105,00$						
RUNQUEUE -	NRT	PER	RQL	CURR	VMIN			
	2	6,00	2,00	t1	105,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t1	1.0	0,50	3,00	1,00	11,00	102,00	
RB	TREE	t2	1.0	0,50	3,00	1,00	24,00	105,00

Figura 6.9

6.4 Creazione e cancellazione di task

Sia la cancellazione (evento EXIT) che la creazione di nuovi task (evento CLONE) comportano il ricalcolo dei parametri della runqueue e dei task dovuto al cambiamento del numero di processi.

L'evento EXIT comporta sempre la necessità di reschedulare.

Nel caso della creazione è necessario determinare il VRT iniziale da assegnare al processo e poi valutare la necessità di rescheduling; la regola adottata per determinare il VRT è la seguente:

$$\text{tnew.VRT} = \text{VMIN} + \text{tnew.Q} * \text{tnew.VRTC}$$

Attenzione alla sequenza di operazioni:

1. tnew viene creato
2. i parametri della runqueue e dei task vengono ricalcolati (incluso tnew.Q)
3. infine viene calcolato tnew.VRT con la formula esposta sopra

In base a questa assegnazione un processo nuovo (che ha tempo di esecuzione SUM=0) parte con un valore di VRT omogeneo agli altri processi. A differenza del caso di wakeup, al nuovo processo viene assegnato un valore di VRT che non lo posizionerà all'inizio della coda, ma comunque in modo da entrare nel periodo di scheduling che inizia con la sua creazione.

Questa assegnazione spiega perchè i valori di VRT sono generalmente molto più grandi di quelli di SUM nei processi, considerando che VMIN cresce monotonamente.

La necessità di rescheduling è valutata esattamente nello stesso modo del caso di wakeup:

```
if ( (tnew->schedule_class == classe con diritto > NORMAL) or
    ((tnew->vrt + WGR * tnew->load_coeff) < CURR->vrt ) ) resched();
```

Esercizio 6

Negli esercizi ai nuovi task creati si assegnano identificatori incrementando progressivamente rispetto a quelli già esistenti. Eseguire una simulazione partendo dalle seguenti condizioni iniziali (limite di simulazione 4ms.)

CONDIZIONI INIZIALI *****									
RUNQUEUE - NRT			PER		RQL		CURR		VMIN
3			6,00		3,00		t1		100,00
TASKS:	ID		LOAD		LC		Q		VRTC SUM VRT
CURRENT	t1		1.0		0,33		2,00		1,00 10,00 100,00
RB TREE	t2		1.0		0,33		2,00		1,00 20,00 100,50
	t3		1.0		0,33		2,00		1,00 30,00 101,00
Events of task t1: EXIT at 1.0;									
Events of task t3: CLONE at 1.0;									

Figura 6.10

Soluzione in Figura 6.11

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		1,00		EXIT		t1		true
RUNQUEUE -	NRT		PER	RQL	CURR	VMIN		
	2	6,00	2,00	t2	100,50			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t2	1.0	0,50	3,00	1,00	20,00	100,50	
RB TREE	t3	1.0	0,50	3,00	1,00	30,00	101,00	

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		4,00		Q_SCADE		t2		true
RUNQUEUE -	NRT		PER	RQL	CURR	VMIN		
	2	6,00	2,00	t3	101,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t3	1.0	0,50	3,00	1,00	30,00	101,00	
RB TREE	t2	1.0	0,50	3,00	1,00	23,00	103,50	

EVENT	*****	TIME		TYPE		CONTEXT		RESCHEDULE
		5,00		CLONE		t3		false
tnew.vrt+WGR*tnew.LC=104.00+1.00*0.33=104.33 < curr.vrt=102.00								
RUNQUEUE -	NRT		PER	RQL	CURR	VMIN		
	3	6,00	3,00	t3	102,00			
TASKS:	ID	LOAD	LC	Q	VRTC	SUM	VRT	
CURRENT	t3	1.0	0,33	2,00	1,00	31,00	102,00	
RB TREE	t2	1.0	0,33	2,00	1,00	23,00	103,50	
	t4	1.0	0,33	2,00	1,00	0,00	104,00	

Figura 6.11

6.5 Riassunto delle notazioni e delle regole del CFS

Notazione

Parametri **SYSCTL**:

LT (latency) – default: 6ms

GR (granularity) – default: 0,75ms

WGR (wakeup_granularity) - default: 1ms;

Altri Parametri globali:

TICK dello scheduler periodico

NRT = numero di task in stato di PRONTO + 1

NOW = istante corrente

Elementi della runqueue **RQ**:

CURR

VMIN

RB (Red Black Tree – è la coda ordinata in base al VRT dei task)

LFT (è il primo della RB, o leftmost)

IDLE (puntatore al descrittore di idle, che non è inserito in RB)

RQL = somma di tutti i t.LOAD, per tutti i task presenti sulla RQ

Elementi di ogni task **t**

VRT (virtual time)

SUM: è il tempo reale di esecuzione del task (alla creazione sum=0)

PREV: è il valore di SUM al momento in cui il task è diventato CURR

LOAD è il peso del task t

Coefficienti di peso

t.VRTC = vrt_coeff

t.LC = load_coeff

Regole

Calcolo dei coefficienti di peso

$t.LC = t.LOAD / RQL$

$t.VRTC = 1 / t.LOAD$

Calcolo del periodo di schedulazione PER

$PER = MAX(LT, NRT * GR)$

Calcolo del quanto di un task

$t.Q = PER * t.LC$

Aggiornamento periodico

$DELTA = NOW - START$

$SUM = SUM + DELTA$

$VRT = VRT + DELTA * VRTC$

$VMIN = MAX(VMIN, MIN(VRT, LFT.VRT))$

if (SUM - PREV) > Q --> resched

Wakeup di un processo tw

$tw.VRT = MAX(tw.VRT, (VMIN - LT/2))$

if ((tw.schedule_class == classe con diritto > NORMAL) or

((tw.VRT + WGR * tw.LC) < CURR.VRT)) resched

Creazione di un processo tnew

$tnew.VRT = VMIN + tnew.Q * tnew.VRTC$

if ((tnew.schedule_class == classe con diritto > NORMAL) or

((tnew.VRT + WGR * tnew.LC) < CURR.VRT)) resched

6.6 Il meccanismo di assegnazione dei pesi ai processi

L'assegnazione dei pesi ai processi si basa sul `nice_value`.

Nice value (priorità) e peso dei processi

L'utente può assegnare a un processo un `nice_value`. I `nice_value` vanno da -20 (massima priorità, bassa gentilezza – assegnabili solo dall'amministratore) a +19 (minima priorità, grande gentilezza). Il valore iniziale assegnato per default a un processo è 0. A parità di priorità e di politica di schedulazione, i processi che hanno `nice_value` maggiori ottengono in proporzione meno tempo di CPU rispetto a processi che hanno `nice_value` minori.

I `nice_value` sono trasformati in pesi dei task (`t.LOAD`); la regola di trasformazione, indicata in Tabella 6 estratta dal file `kernel/sched/sched.h`, corrisponde all'incirca alla seguente formula esponenziale:

$$t.LOAD = (1024 / 1.25^{\text{nice_value}})$$

```
kernel/sched/sched.h
/*
1120 * Nice levels are multiplicative, with a gentle 10% change for every
1121 * nice level changed. I.e. when a CPU-bound task goes from nice 0 to
1122 * nice 1, it will get ~10% less CPU time than another CPU-bound task
1123 * that remained on nice 0.
1124 *
1125 * The "10% effect" is relative and cumulative: from _any_ nice level,
1126 * if you go up 1 level, it's -10% CPU usage, if you go down 1 level
1127 * it's +10% CPU usage. (to achieve that we use a multiplier of 1.25.
1128 * If a task goes up by ~10% and another task goes down by ~10% then
1129 * the relative distance between them is ~25%.)
1130 */
1131 static const int prio_to_weight[40] = {
1132 /* -20 */      88761,      71755,      56483,      46273,      36291,
1133 /* -15 */      29154,      23254,      18705,      14949,      11916,
1134 /* -10 */      9548,       7620,       6100,       4904,       3906,
1135 /* -5  */      3121,       2501,       1991,       1586,       1277,
1136 /* 0   */      1024,       820,        655,        526,        423,
1137 /* 5   */      335,        272,        215,        172,        137,
1138 /* 10  */      110,        87,         70,         56,         45,
1139 /* 15  */      36,         29,         23,         18,         15,
1140 };
```

Tabella 6

La costante `NICE_0_LOAD` vale quindi 1024, perché è il peso associato al `nice_value` 0.

6.6 Esempi sul sistema reale (Approfondimento)

Utilizzando i `nice_value` si può osservare il comportamento del sistema reale.

Esempio 1: effetto del nicevalue su 2 processi CPU bound

Il programma di figura 6.10a crea 2 thread che eseguono ambedue per 10 volte la funzione `cpu_load`, che impiega 100ms. Il primo thread `tf1` modifica il proprio `nice_value` a 5, mentre il thread `tf2` mantiene il `nice_value` di default (0).

Il risultato dell'esecuzione, riportato in figura 6.10.b. Come si vede, il thread `tf2` riesce ad eseguire la funzione `cpu_load` il triplo delle volte rispetto a `tf1`.

Il valore dei pesi associati ai `nice_value` 0 e 5, come indicato in tabella 6, sono 1024 e 335, e il loro rapporto è 3.05; la prova conferma quindi che il rapporto tra i pesi dei task determina il rapporto tra i relativi tempi di esecuzione


```

void cpu_load(int iterations) {
    long long i,j, k;
    for (k=0; k<iterations; k++){
        //circa 100ms ad ogni iterazione
        for (j=0; j<MAX; j++){ i=j; }
    }
}

void * tfcpu(void * tfarg){
    int k;
    int arg = (int) tfarg;
    if (arg == 1) {          //thread 1 – esegue con nice = 5
        nice(5);    printf("start tf 1 \n");
    }
    else if (arg == 2) {    //thread 2 – esegue con nice di default = 0
        printf("start tf 2 \n");
    }
    //ciclo eseguito da ambedue i thread – durata 100 ms
    for (k=0; k<10; k++){
        printf("eseguito ciclo %d per tf %d \n", k, arg);
        cpu_load(1);
    }
    printf("tfcpu TERMINATA numero: %d\n", arg);
    return NULL;
}

int main(long argc, char * argv[], char * envp[]) {
    pthread_t t1, t2;
    int pid;
    pthread_create(&t1, NULL, &tfcpu, (void *) 1);
    pthread_create(&t2, NULL, &tfcpu, (void *) 2);
    printf("MAIN - ATTESA \n");
    pthread_join(t1, NULL);
    pthread_join(t2, NULL);
    printf("MAIN PROCESS TERMINATO \n");
    return printed_lines;
}

```

a) codice del programma

MAIN - ATTESA	eseguito ciclo 2 per tf 1
start tf 2	eseguito ciclo 7 per tf 2
eseguito ciclo 0 per tf 2	eseguito ciclo 8 per tf 2
start tf 1	eseguito ciclo 9 per tf 2
eseguito ciclo 0 per tf 1	eseguito ciclo 3 per tf 1
eseguito ciclo 1 per tf 2	tfcpu TERMINATA numero: 2
eseguito ciclo 2 per tf 2	eseguito ciclo 4 per tf 1
eseguito ciclo 3 per tf 2
eseguito ciclo 1 per tf 1	eseguito ciclo 9 per tf 1
eseguito ciclo 4 per tf 2	tfcpu TERMINATA numero: 1
eseguito ciclo 5 per tf 2	MAIN PROCESS TERMINATO
eseguito ciclo 6 per tf 2	

b) risultato dell'esecuzione

Figura 6.10 – sperimentazione dell'effetto dei nice_value

Esempio 2: lettura dei valori interni al sistema

Nella seguente figura sono mostrati i valori delle variabili VRT, START, SUM e PREV di un task, ottenuti tramite un kernel module apposito.

I valori sono stati campionati facendo eseguire un programma che esegue 3 volte un ciclo che richiede circa 100ms; la prima riga mostra i valori a inizio programma, le seconda a inizio ciclo, e le successive alla fine di ogni iterazione del ciclo.

Il programma è stato rieseguito 3 volte con nice_value 0, 5, 10, ai quali corrispondono i pesi 1024, 335, 110 rispettivamente. Nei risultati sono indicati anche i valori dei relativi vrt_coeff (VRTC) e la differenza tra i valori iniziale e finale del VRT e di SUM.

Tutti i valori temporali stampati sono in nanosecondi.

Si osserva che il rapporto tra SUM e VRT riflette esattamente il valore di VRTC; altri valori devono essere considerati con una certa tolleranza, per diversi motivi:

- la calibratura della durata del ciclo a 100 ms non è perfetta e non è del tutto stabile; il tempo effettivamente impiegato oscilla di alcuni percento
- nel sistema sono attivi altri task di sistema che interferiscono

Esempio – effetto di nice sui parametri dello scheduler

a) nice = 0 load 1024

```
[ 2453.427265] vrt: 29755 043792, start: 2453427 251962, sum:      586721, prev_sum:      586721
[ 2453.537675]
[ 2453.537679] vrt: 29852 647989, start: 2453536 559531, sum:  98 190918, prev_sum:  91 031363
[ 2453.636694]
[ 2453.636698] vrt: 29949 536024, start: 2453636 018675, sum: 195 078953, prev_sum: 191 411257
[ 2453.737772]
[ 2453.737776] vrt: 30047 380845, start: 2453737 174017, sum: 292 923774, prev_sum: 292 923774
[ 2453.846347]
[ 2453.846352] vrt: 30148 521595, start: 2453844 359992, sum: 394 064524, prev_sum: 390 213700
                delta vrt: 393 ms                delta sum: 394
```

b) nice = 5 load 335 (vrt_coeff: 3,1)

```
[ 2506.967770] vrt: 30221 851996, start: 2506967 762715, sum:      576577, prev_sum:      576577
[ 2507.076766]
[ 2507.076770] vrt: 30521 901471, start: 2507076 334538, sum:  98 737299, prev_sum:  79 506029
[ 2507.176417]
[ 2507.176417] vrt: 30818 620481, start: 2507176 416755, sum: 195 808464, prev_sum: 188 394180
[ 2507.280752]
[ 2507.280756] vrt: 31121 588611, start: 2507280 012015, sum: 294 924020, prev_sum: 288 610822
[ 2507.406717]
[ 2507.406721] vrt: 31424 604156, start: 2507404 016266, sum: 394 055088, prev_sum: 391 630771
                delta vrt: 1203 ≈ delta sum * 3,1    delta sum: 394
```

c) nice = 10 load 110 (vrt_coeff: 9,3)

```
[ 2935.301661] vrt: 32757 003504, start: 2935301 652169, sum:      604004, prev_sum:      604004
[ 2935.403602]
[ 2935.403606] vrt: 33618 675839, start: 2935400 010181, sum:  93 166464, prev_sum:  29 850013
[ 2935.499936]
[ 2935.499941] vrt: 34519 023316, start: 2935496 849093, sum: 189 883480, prev_sum: 133 064418
[ 2935.594520]
[ 2935.594523] vrt: 35406 688692, start: 2935592 208872, sum: 285 238161, prev_sum: 193 355737
[ 2935.691773]
[ 2935.691776] vrt: 36275 207653, start: 2935688 021581, sum: 378 536098, prev_sum: 294 564557
                delta vrt: 3518 ≈ delta sum * 9,3    delta sum: 378
```

Figura 6.11