

Documentation Ipanema

Baptiste Pires

May 31, 2024

1 Introduction

2 Structures

Dans cette section nous allons voir les différentes structures utilisées dans Ipanema et leurs utilités

2.1 struct sched_ipanema_entity

Nous allons commencer par la structure qui représente une entité d’ordonnancement. La structure est déclarée de cette façon :

```
1 struct sched_ipanema_entity {
2     union {
3         struct rb_node node_runqueue;
4         struct list_head node_list;
5     };
6
7     /* used for load balancing in policies */
8     struct list_head ipa_tasks;
9
10    int just_yielded;
11    int noproempt;
12
13    enum ipanema_state state;
14    struct ipanema_rq *rq;
15
16    /* Policy-specific metadata */
17    void *policy_metadata;
18
19    struct ipanema_policy *policy;
20 };
```

Les deux premiers champs, `struct rb_node node_runqueue` et `struct list_head node_list` sont utilisés pour maintenir une tâche dans une runqueue. Dû à l’union, on peut être dans une seule runqueue à un moment donné (soit dans une runqueue qui se sert d’un arbre rouge-noir comme CFS, ou une runqueue qui utilise une liste comme un algorithme FIFO par exemple).

Le champ `struct list_head ipa_tasks` est utilisé par les différentes politiques d’ordonnancement pendant l’équilibrage de charge. Par exemple, si on migre des tâches d’un cœur C_1 à un cœur C_2 , on a besoin de retirer les tâches de la runqueue de C_1 en ayant le verrou sur celle-ci. Quand on retire les tâches on peut les ajouter à la liste `ipa_tasks`. Une fois qu’on a retiré toutes les tâches qu’on souhaite migrer, on prend le verrou sur la runqueue de C_2 et on peut itérer sur la liste que l’on vient de construire.

`int just_yield` permet de savoir si on vient d’appeler `sched_yield()` ou non.

`int noproempt`

`enum ipanema_state state` représente l’état d’une entité Ipanema. Les états possibles sont les suivants :

```

1 enum ipanema_state {
2     IPANEMA_NOT_QUEUED, /* La tache n'est pas dans une runqueue */
3     IPANEMA_MIGRATING, /* La tache est en train de migrer d'un coeur a un autre */
4     IPANEMA_RUNNING, /* La tache est dans une runqueue et s'exécute */
5     IPANEMA_READY, /* La tache est dans une runqueue mais ne s'exécute pas */
6     IPANEMA_BLOCKED, /* La tache est bloqué (I/O, lock, etc...) */
7     IPANEMA_TERMINATED, /* La tache meurt (?) */
8     /*
9      * A special state that is equivalent to IPANEMA_READY, but makes it
10     * possible to figure out that the state change came from tick().
11     */
12     IPANEMA_READY_TICK
13 };

```

`struct ipanema_rq *rq` c'est un pointeur vers la `struct ipanema_rq` sur laquelle la tâche s'exécute (ou va s'exécuter).

2.1.1 struct ipanema_rq

Now we will take a look at the `struct ipanema_rq`. It's the structure that represents a runqueue in Ipanema, fairly simple struct made of the following fields :

```

1 struct ipanema_rq {
2     enum ipanema_rq_type type;
3     union {
4         struct rb_root root;
5         struct list_head head;
6     };
7
8     /*
9     * cpu to which this runqueue belongs (for per_cpu runqueues).
10    * For runqueues that are not per_cpu, rq->cpu > NR_CPUS.
11    */
12    unsigned int cpu;
13    enum ipanema_state state;
14
15    /* Number of tasks in this runqueue */
16    unsigned int nr_tasks;
17
18    int (*order_fn)(struct task_struct *a, struct task_struct *b);
19 };

```

`enum ipanema_rq_type` defines what kind of structure the runqueue will use. The enum is defined as follow : `enum ipanema_rq_type { RBTREE, LIST, FIFO };`. As we can see, three types of runqueues are supported now. One with an red-black tree structure, one with a list and with a fifo list.

The union of `struct rb_root` and `struct list_head` is the root of either the tree or list of tasks in this runqueue.

`int cpu` is the id of the cpu this runqueue belongs to.

`enum ipanema_state` demander à Redha pour être sûr. Defines the state in which the tasks must be when added to this runqueue.

`unsigned int nr_tasks` The number of tasks currently in the runqueue.

`int (*order_fn)(struct task_struct *a, struct task_struct *b)` is function pointer that will be called when a task is added to the runqueue or whenever an ordering is needed. It is not mandatory for the FIFO runqueues.

That concludes the part for the runqueue.

3 Ipanema architecture

Dans cette section nous allons voir les détails de l'architecture d'Ipanema.

3.1 Les variables globales d'Ipanema

Dans Ipanema, il y a deux variables globales importantes : `struct list_head ipanema_policies` et `per_cpu(struct task_struct *, ipanema_current)`. Nous allons voir comment ces deux variables sont utilisées pour comprendre leur utilité.

ipanema_policies La variable `struct list_head ipanema_policies` est une liste chaînée qui contient toutes les politiques d'ordonnancement disponibles. Cela correspond à tous les modules qui ont été insérés en tant que politiques. Cette liste est protégée par un `rwlock_t ipanema_rwlock` qui permet de la modifier/parcourir en toute sécurité.

Lorsque l'on souhaite ajouter une politique, on fait appel à la fonction `ipanema_add_policy(struct ipanema_policy *)` qui va insérer la nouvelle politique dans la liste.

Nous avons maintenant inséré une politique, mais aucun thread ne s'en sert pour l'instant. Pour ce faire, il va falloir se servir de l'appel système `SYS_sched_setattr` (voir `tools/ipanema/ipastart.c` pour un exemple d'utilisation). En somme, il appelle la fonction `__sched_setscheduler` qui va, si la politique demandée existe, modifier le champ de la `task_struct.ipanema.policy` pour le faire pointer vers la politique souhaitée. Ensuite, cette liste sera par exemple utilisée dès que l'on a besoin de trouver une tâche à exécuter.

ipanema_current Cette variable est l'une des plus importantes. Elle est définie `per_cpu`, ce qui signifie que pour tous les CPUs, il y a une variable d'allouée. Ces variables vont stocker le processus en cours d'exécution sur le cœur auquel elles appartiennent. Pour voir cela, il faut regarder la fonction `change_state` qui gère toutes les transitions d'états d'Ipanema. On peut voir que pour la transition `x -> RUNNING` (une tâche va commencer son exécution) : `per_cpu(ipanema_current, next_cpu) = p` où `next_cpu` est le CPU sur lequel la tâche va s'exécuter et `p` ladite tâche.

Il faut noter que juste avant, la fonction `change_rq` est appelée et la tâche `p` aura été retirée de sa runqueue où elle attendait à l'état `IPANEMA_READY`.

4 Ecrire un scheduler dans Ipanema

Dans cette section nous allons voir comment écrire un scheduler sous la forme d'un module en utilisant Ipanema.