

Process Scheduling [Chapter 2]

Multitasking OS comes in two flavours

- (1) Co-operative scheduling
- (2) Preemptive

All the Linux slice of comes in Preemptive Scheduling flavour.

In Preemptive Multitasking Schedule decide which Process to start Running and which Process to stop Running.

The time a process runs before it preempted is usually predetermined, and is called the time slice of the process. Process's time.

Managing the time slice enables the scheduler to make global scheduling decision.

Cooperative Scheduling In this Method

Process does not stop running until it voluntarily stops. the act of a process voluntarily suspending itself is called yielding.

Limitation of Cooperative Scheduling

- (1) Scheduled Can not take the global decision.
- (2) Any process can monopolize the system.

I/O Bound Virtual Processes Round Robin Processes

I/O Bound Process spends most of the time waiting for some input/output event to occur.

Process Round Robin Process spends most of this time executing the code. These kind of Process runs Continuously without waiting for any I/O event and stopped while preempted by the scheduler.

Priority Preemptive

Linux implements two varieties of Priority Preemptive values.

- (1) Nice Values: A number from -20 to +19 with a default value of 0. Large Nice Value Gets priority.

to lower Priority.

⑤ Real-time priority: this value are 0 - 99.

higher the Real time value higher the Priority.

⇒ Time Slice is a time that how long a task can run until it is preempted. If the time-slice will be large then system will be less interactive and when timeslice will be low then context switching will be extra overhead on the system thus can the performance low.

⇒ The Linux Scheduling Algorithm:

Scheduled classes, Linux scheduling is modular, enabling different Algo to schedule different types of processes this modularity is called scheduled classes.

The Complete Fair Scheduler (CFS) is the Registered scheduler class for normal process called SCHED-NORMA

In Linux,

CFS is defined in kernel/sched/fair.c

→ Process Scheduling in Unix Systems

→ Need to check

⇒ Fair Scheduling: This allocation is
which we assume that our process
is in ideal state and every process
gets an equal amount of time.
Suppose if we have n processes then
each process will get 1/n time
of process.

Perfect multi-tasking is the situation in
which every process gets 1/n time
for the whole time on processes

But CFS will run each process for some
amount of time, rounds robin, selecting
next the process that has been the
last. CFS calculate how long a
process should be run as a function
of the total number of runnable
processes.

In CFS we calculate the weight of each process by using the nice value of that process. Each process then runs for a "time slice" proportional to its weight divided by the total weight of all runnable threads.

If the number of Runnable Process will reach to infinity then the time slice for each process will approach to zero which increase the unnecessary overhead of switching. To avoid CFS impose a limit on timeslice which is the minimum timeslice to be allocated a process. By default it's 1 millisecond. More work to be done on this

Linux Scheduling Implementations

We will discuss here CFS (Actual implementation) which lives in kernel/sched/fair.c

- (1) Time Accounting
- (2) Process Selection
- (3) The Scheduling Entry Point
- (4) Sleeping and waking up

① Time Accounting

All Process Scheduler ~~is~~ must account for the time that a process runs.

Most Unix System maintains the account not how much a process Actual runs on Processor.

~~The Scheduled Entity / Structure~~

CFS does / not have the feature of using the time slice. But It has to maintain the account for time a process Actual Runs for that it uses scheduled entity structure

Struct sched_entity, defined in Linux Scheduler To keep track of process accounting.

this structure embed in the Process descriptor structure like struct

The VIRTUAL RUNTIME

The Virtual runtime variable stored the virtual runtime of a process which is the actual runtime (the amount of time

Spent Running) Normalized by the number of runnable processes.

CFP uses runtime to account for how long a process has been and thus how much longer it ought to run.

update_crr() calculates the execution time of the current process and stores that value in delta_exec.

update_crr() is invoked periodically by system timer and also whenever a process becomes runnable or becomes unrunnable.

① Process Selection

In Linux CFS selects the process which has smallest runtime to next run.

CFS uses a Red Black tree to manage the list of runnable processes and efficiently finds the process with smallest runtime.

Picking the next task: \rightarrow CFS search the leftmost child in a tree which have the smallest runtime. It is performed by an concurrent entity () . It has complexity $O(\log N)$ for N Nodes if the tree is balanced.

If the tree of Ready is Null it means there is no task to Run so scheduling the idle task.

Adding Processes to tree

It happens when a Process become Runnable or forced Call create a process. Adding of Process in tree is performed by enqueue entity ();

It add the Process in the same way as we add a node in BST. ()

Removing Process

It happens when a process blocks or terminates.

driving entity ()

③ The Scheduler Entity Point

The main entity point in the process schedule is the function `schedule()` defined in `kernel/sched.c`, this function is called by rest of kernel to schedule other process.

`schedule()` is a generic call & internally calls each scheduling class and checks the highest priority tasks from highest priority scheduling task.

Pending tasks ()

If we know that gfs and tasks are in the fail class we can call that function directly:

as Sleeping and Waking Up

Whenever a process goes into sleep state for whatever reason it is removed from ready ring and calls `schedule()` to select a new process to execute. Waking up will set the interrupt. The timer is set at

removable, removable from the wait queue
and added back to the tail-blocked list

WaitQueue: This is a simple list of

processes waiting for an event to occur.
Waitqueue are implemented in the kernel
by a WaitQueueHead. It, WaitQueue
are created statically via DECLARE-
WAITQUEUE() or dynamically via
init_waitqueue_head().

Process Put themselves on a waitqueues
and make themselves not removable.
When the event associated with
Waitqueue occurs, the process on
the queue are awakened.

/* q is the wait queue we wish to
sleep on */

DEFINE_WAIT(wait);

It will create a wait queue entry.

After waitqueue(q, &wait);

Add itself to a wait queue via
add_waitqueue(). This wait queue

Awakened the process when the condition for which it is waiting occurs.

while (Condition)

{

 prepare-to-wait (Aq, Await, TASK_INTERRUPTIBLE).

* It changes the Process State to either TASK_INTERRUPTIBLE or. This function also adds the task to the wait queue if necessary, which is needed on subsequent iteration of loop.

If (signal-handling (current))

/> handle signal

Schedule ()

finish-wait (Aq, Await);

If set the task to TASK_RUNNING and removes itself from the wait queue,

example: notify in filesystem in wait queue.

Waking up:

Waking up is handled via wake-up(), which wakes up all the tasks waiting on the given wait queue. It calls try-to-wake-up() which sets the task state to TASK_RUNNING and adds the task Red slack after.

For example when some data arrives in hard disk, the VFS calls wake-up() on the wait queue that holds the process waiting for the data.

Premotion and Context Switching:

Context switching is switching from one runnable process to another by context-switch() function, defined in Kernel/sched.c. If it is called by function schedule() when a new process has been selected to run, it does two basic jobs:

- (i) switch_mm(): To switch the virtual memory mapping of that process to next process.

(2) Switch_to(): To switch the Process State

from previous process to current's. It does the saving ~~info~~ of current state, Process Register etc.

User Preemption: It occurs when the kernel

is about to return to user-space. Need_Retire is set and therefore the scheduler is invoked.

No user-preemption can occur

- ① When Returning to User-Space from a System Call.
- ② Returning to user-space from an interrupt handled.

Kernel Preemption

Any task in the kernel which is not holding any lock can be preempted.

To enable the Kernel Preemption in kernel kernel embed a variable preempt_Grant in the thread_info structure.

Everytime when process acquired a lock
the value increase by 1 and vice versa.
Whenever this value is 0 the kernel
is preempted.

Upon Return from the interrupt of returning
to kernel space value of need reschedule
is set and preempt count is 0
then most important task is scheduled