

# **The Mach cpu\_server: An Implementation of Processor Allocation**

David L.Black

Department of Computer Science  
Carnegie-Mellon University  
Pittsburgh, PA 15213

**Version of:**  
14 August 1990

## **Abstract**

The CPU-Server is a user-mode server that performs processor allocation for the Mach operating system. This document describes the server and its user interfaces.

This research was sponsored by the Defense Advanced Research Projects Agency (DOD), ARPA Order No. 4864, monitored by the Space and Naval Warfare Systems Command under contract N00039-84-C-0467.

The views and conclusions contained in this document are those of the authors and should not be interpreted as representing official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or of the U.S. Government.

## 1. Introduction

This document describes the `cpu_server`, a user-mode server that performs processor allocation for the Mach operating system, and its library interfaces. The server implements processor allocation policy; the actual mechanisms are contained in the Mach kernel (see related document). The library interfaces hide the details of server interaction from applications with simple scheduling requirements. This server is an example implementation of processor allocation policy; many other policies can be implemented using the mechanisms provided by the kernel.

Not all releases of Mach support this server. It is mainly interesting on multi-processor machines, so no major effort has been made to support it on the standard uni-processor versions of Mach. However, the new releases of Mach support it on all machines. On systems that support it, this software (MiG interface to server and library interfaces) is available in the `cpu` library (`/usr/mach/lib/libcpu.a`, `-lcpu` command to `ld`).

## 2. Concepts

The server performs processor allocation by assigning processors to processor sets provided by its clients. This allows the clients to use and manage the processors without giving clients complete control over them; only the server has the port capabilities required to reassign processors. Processor sets are entities exported by the Mach kernel; threads assigned to a processor set run exclusively on processors assigned to that set and vice versa (with the exception of some Unix<sup>™</sup> system calls).

The server interface is designed around a class of objects called requests. A request consists of the following components:

- A run duration,
- A sequence of <processor set, number of processor> pairs.

A request is satisfied by assigning each processor set its corresponding number of processors for the run duration specified. The server enforces internal limits on the number of processors and the maximum run time. Current limits are 15 minutes and 75% or less of the processors on the system.

## 3. Implementation

The server satisfies requests in a greedy fashion with strict adherence to the order in which they are received. For example, if the server has 10 processors and receives requests for 4, 7, and 2 processors, it will satisfy the request for 4 first, and then the requests for 7 and 2 together. This algorithm was chosen for its simplicity and lack of starvation; more sophisticated algorithms that make better use of the processors by satisfying requests out of order could be used.

The server implementation was based on the `cthreads` library and the Mach Interface Generator(MiG). Two threads are used internally; one manages the assignment of processors to requests, the other manages all interactions with clients. Clients communicate with the server via `rpcs`; MiG-generated interfaces hide the details of the message formats. In addition the server can optionally generate notification messages to indicate that processors have been allocated to a request and that processors are about to be removed from a request; these messages can be used for internal synchronization of applications that require multiple processors for proper execution.

## 4. Interface

The server's interface to its clients uses remote procedure calls to implement the following primitives:

- **cpu\_request\_create**(server, total\_processors, run\_time, \*delay, \*request) -- create a request for *total\_processors* processors for *run\_time* seconds. The request object and a delay estimate are returned. The initial server port (*server*) is obtained by looking up the name "cpu\_server" with the local name service.
- **cpu\_request\_add**(request, processor\_set, processors, \*processors\_left) -- Add the tuple <*processor\_set*, *processors*> to the specified request. Return the number of processor remaining in the request (i.e. that can be used in future **cpu\_request\_add** calls).
- **cpu\_request\_set\_notify**(request, notify\_port) - Ask the server to send a notification message to *notify\_port* after the processors are allocated and 1 second before removing the processors. Applications that receive these messages can use them to make sure that execution only takes place while all of the processors are allocated; the end message can be used to initiate a barrier synchronization to stop cleanly, and the start message can be used to exit the barrier.
- **cpu\_request\_activate**(request, options, total\_time, \*delay) -- Activate the request (i.e. request the server to satisfy it). A maximum delay until the requested processors will be allocated is returned. No further **cpu\_request\_add** calls are permitted on the request. The following options are supported:
  - **Destroy** - Destroy the processor sets when the request is completed. This is intended to support naive users by preventing a program that overruns its request from going into suspended animation; destroying its processor sets forces the program back into the default processor set where it will continue to run.
  - **Repeat** - Repeat the request for *total\_time*. The time for each instance of the request is specified by the *run\_time* argument of the **cpu\_request\_create** operation that created the request.
- **cpu\_request\_activate\_task**(request, options, total\_time, task, \*delay) -- Identical to **cpu\_request\_activate**, but informs the server that the application using the processor sets is a task; this allows the server to optimize assignment and removal of processors by using **task\_suspend** and **task\_resume** -- processor assignment is much faster if the processor is idle.
- **cpu\_request\_destroy**(request) - Destroy the specified request, freeing any processors that were allocated to it. If the **Notify end** option applies, an end notification will be sent and the freeing of the processors will be delayed by 1 second.
- **cpu\_request\_status**(request, \*reserved\_processors, \*assigned\_processors, \*active, \*options, \*time) - Find out information about a request including whether it has processors assigned to it, and how long until the assigned processors will be removed, or the maximum delay until processors will be assigned to it.
- **cpu\_server\_info**(server, \*max\_time, \*max\_total\_time, \*max\_processors, \*delay) -- Obtain information about the server. *max\_time* is the maximum time for a single request in **cpu\_request\_create**. *max\_total\_time* is the maximum total time for a repeating request created by the **Repeat** option to **cpu\_request\_activate** or **cpu\_request\_activate\_task**. *max\_processors* is the maximum number of processors for any request. *delay* is the current maximum delay until a new request can be satisfied. The **cpuinfo** program uses this call to determine if the *cpu\_server* is available and what its current situation is.

## 5. Library Interfaces

The above server interface along with the kernel interface will be used directly by programs that require explicit control over which threads are executing on which processors at which time. For applications with less stringent processor allocation requirements, simple library interfaces which hide all of the internal details of interaction with the kernel and server may be appropriate. Four interfaces have been developed, the `allocate`, `task`, `hook`, and `task-hook` interfaces.

The `allocate` interface supports a single allocation of a pool of processors. It exports the following two calls:

- **allocate\_processors**(num\_processors, time, interactive) - Allocate *num\_processors* processors for the specified time. If the time is larger than the server's maximum slice time, then a repeating request is automatically submitted. If *interactive* is TRUE, errors generate printf's, and the user is asked whether the server's maximum delay is acceptable; no allocation is performed if the answer is no.
- **deallocate\_processors**() - Free the processors allocated by **allocate\_processors**. This must be called by a thread in the same task as the thread that did the allocation.

**allocate\_processors** does not return until the allocation of processors has started; it performs a **task\_assign** internally so that the initial thread and all threads and tasks subsequently created share the allocated processors. If a program overruns its time allocation, it will continue to run, but without dedicated processors.

The `task` interface is identical to the `allocate` interface, but is restricted to applications consisting of a single task so that the server can exploit efficiencies available in this case (suspending the task before removing processors). The server will perform its own **task\_suspend** and **task\_resume** calls on the task in this case. The task interface exports the following two calls:

- **task\_allocate\_processors**(num\_processors, num\_seconds, interactive) - Identical to **allocate\_processors**, but also promises the server that all or the threads using the processors will be in the current task.
- **task\_deallocate\_processors**() - Deallocate the processors allocated by **task\_allocate\_processors**.

The `hook` interface supports allocation of a pool of processors, with user scheduling hooks. It exports the following two calls:

- **allocate\_processors\_with\_hooks**(num\_processors, num\_seconds, start\_hook, end\_hook, interactive) - Allocate the specified number of processors for the specified number of seconds. If the number of seconds is greater than the server's maximum slice time, a repeating request is automatically submitted. *start\_hook* is called each time after the processors are allocated. Similarly, *end\_hook* is called approximately 1 second before any processor deallocation. *interactive* functions identically to the previous interfaces.
- **deallocate\_processors\_with\_hooks**() - Free the processors allocated by **allocate\_processors\_with\_hooks**. The *end\_hook* will not be called, but users should be aware of the race between this deallocate call and the server ending a time slice.

A thread must be dedicated to the **allocate\_processors\_with\_hooks** call; i.e. the calling thread does not return until the allocation request is finished or terminated. This dedicated thread is assigned to the allocated processors (this can be reversed by performing a **thread\_assign\_default** or **thread\_assign** as part of the first invocation of *start\_hook*. All threads within its task and subsequently created tasks are also assigned (c.f. **task\_assign**). Both *start\_hook* and *end\_hook* must return.

Finally, there is the task-hook interface, which combines the functionality of the hook interface with the server optimization of the task interface; the calls in this interface are the hook interface calls with task\_ prefixed. Users should note that the interfaces are independent; processors must be deallocated with the deallocate call from the same interface that was used to allocate them.

## 6. Extensions

The `cpu_server` could be extended in a number of ways:

- Time Sensitivity - Allow more of the machine to be allocated during off peak periods.
- Request Size Sensitivity - Use a declining priority system instead of absolute ordering to speed throughput. This is a standard technique from schedulers for physical memory machines; other similar techniques are also applicable (the current processor allocation problem for multiprocessors is essentially similar to the memory allocation problem for schedulers for physical memory machines).
- User Sensitivity - Reserve some portion of the machine for certain users based on administrative decisions (e.g. they paid for some portion of the machine, so they should always be able to get at least that portion).
- Non-Uniform Topologies - For NUMA (Non Uniform Memory Access), all topology knowledge and policy is located in the server. Techniques such as first-fit and best-fit may be applicable to the processor allocation server for such a machine (e.g. the clustered Gigamax machine being built by Encore).

In addition the server and the kernel interface (thread assignment) could be used directly in a language runtime that is sophisticated enough to control and allocate processors to its threads.

## Table of Contents

<b>1. Introduction</b>	<b>1</b>
<b>2. Concepts</b>	<b>1</b>
<b>3. Implementation</b>	<b>1</b>
<b>4. Interface</b>	<b>2</b>
<b>5. Library Interfaces</b>	<b>3</b>
<b>6. Extensions</b>	<b>4</b>