

# Cluster-wide Power-up/power-down race avoidance algorithm

This file documents the algorithm which is used to coordinate CPU and cluster setup and teardown operations and to manage hardware coherency controls safely.

The section "Rationale" explains what the algorithm is for and why it is needed. "Basic model" explains general concepts using a simplified view of the system. The other sections explain the actual details of the algorithm in use.

## Rationale

In a system containing multiple CPUs, it is desirable to have the ability to turn off individual CPUs when the system is idle, reducing power consumption and thermal dissipation.

In a system containing multiple clusters of CPUs, it is also desirable to have the ability to turn off entire clusters.

Turning entire clusters off and on is a risky business, because it involves performing potentially destructive operations affecting a group of independently running CPUs, while the OS continues to run. This means that we need some coordination in order to ensure that critical cluster-level operations are only performed when it is truly safe to do so.

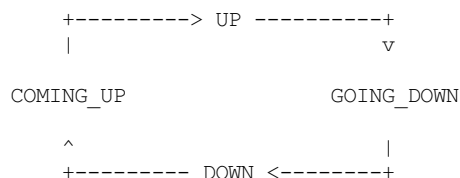
Simple locking may not be sufficient to solve this problem, because mechanisms like Linux spinlocks may rely on coherency mechanisms which are not immediately enabled when a cluster powers up. Since enabling or disabling those mechanisms may itself be a non-atomic operation (such as writing some hardware registers and invalidating large caches), other methods of coordination are required in order to guarantee safe power-down and power-up at the cluster level.

The mechanism presented in this document describes a coherent memory based protocol for performing the needed coordination. It aims to be as lightweight as possible, while providing the required safety properties.

## Basic model

Each cluster and CPU is assigned a state, as follows:

- DOWN
- COMING\_UP
- UP
- GOING\_DOWN



DOWN:

The CPU or cluster is not coherent, and is either powered off or suspended, or is ready to be powered off or suspended.

COMING\_UP:

The CPU or cluster has committed to moving to the UP state. It may be part way through the process of initialisation and enabling coherency.

UP:

The CPU or cluster is active and coherent at the hardware level. A CPU in this state is not necessarily being used actively by the kernel.

GOING\_DOWN:

The CPU or cluster has committed to moving to the DOWN state. It may be part way through the process of teardown and coherency exit.

Each CPU has one of these states assigned to it at any point in time. The CPU states are described in the "CPU state" section, below.

Each cluster is also assigned a state, but it is necessary to split the state value into two parts (the "cluster" state and "inbound" state) and to introduce additional states in order to avoid races between different CPUs in the cluster simultaneously modifying the state. The cluster-level states are described in the "Cluster state" section.

To help distinguish the CPU states from cluster states in this discussion, the state names are given a *CPU\_* prefix for the CPU states, and a *CLUSTER\_* or *INBOUND\_* prefix for the cluster states.

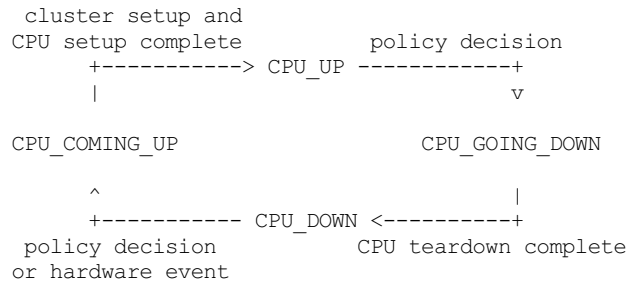
## CPU state

In this algorithm, each individual core in a multi-core processor is referred to as a "CPU". CPUs are assumed to be single-threaded: therefore, a CPU can only be doing one thing at a single point in time.

This means that CPUs fit the basic model closely.

The algorithm defines the following states for each CPU in the system:

- CPU\_DOWN
- CPU\_COMING\_UP
- CPU\_UP
- CPU\_GOING\_DOWN



The definitions of the four states correspond closely to the states of the basic model.

Transitions between states occur as follows.

A trigger event (spontaneous) means that the CPU can transition to the next state as a result of making local progress only, with no requirement for any external event to happen.

**CPU\_DOWN:**

A CPU reaches the CPU\_DOWN state when it is ready for power-down. On reaching this state, the CPU will typically power itself down or suspend itself, via a WFI instruction or a firmware call.

Next state:

CPU\_COMING\_UP

Conditions:

none

Trigger events:

- an explicit hardware power-up operation, resulting from a policy decision on another CPU;
- a hardware event, such as an interrupt.

**CPU\_COMING\_UP:**

A CPU cannot start participating in hardware coherency until the cluster is set up and coherent. If the cluster is not ready, then the CPU will wait in the CPU\_COMING\_UP state until the cluster has been set up.

Next state:

CPU\_UP

Conditions:

The CPU's parent cluster must be in CLUSTER\_UP.

Trigger events:

Transition of the parent cluster to CLUSTER\_UP.

Refer to the "Cluster state" section for a description of the CLUSTER\_UP state.

**CPU\_UP:**

When a CPU reaches the CPU\_UP state, it is safe for the CPU to start participating in local coherency.

This is done by jumping to the kernel's CPU resume code.

Note that the definition of this state is slightly different from the basic model definition: CPU\_UP does not mean that the CPU is coherent yet, but it does mean that it is safe to resume the kernel. The kernel handles the rest of the resume procedure, so the remaining steps are not visible as part of the race avoidance algorithm.

The CPU remains in this state until an explicit policy decision is made to shut down or suspend the CPU.

Next state:

CPU\_GOING\_DOWN

Conditions:

none

Trigger events:

explicit policy decision

**CPU\_GOING\_DOWN:**

While in this state, the CPU exits coherency, including any operations required to achieve this (such as cleaning data caches).

Next state:

CPU\_DOWN

Conditions:

local CPU teardown complete

Trigger events:

(spontaneous)

## Cluster state

A cluster is a group of connected CPUs with some common resources. Because a cluster contains multiple CPUs, it can be doing multiple things at the same time. This has some implications. In particular, a CPU can start up while another CPU is tearing the cluster down.

In this discussion, the "outbound side" is the view of the cluster state as seen by a CPU tearing the cluster down. The "inbound side" is the view of the cluster state as seen by a CPU setting the CPU up.

In order to enable safe coordination in such situations, it is important that a CPU which is setting up the cluster can advertise its state independently of the CPU which is tearing down the cluster. For this reason, the cluster state is split into two parts:

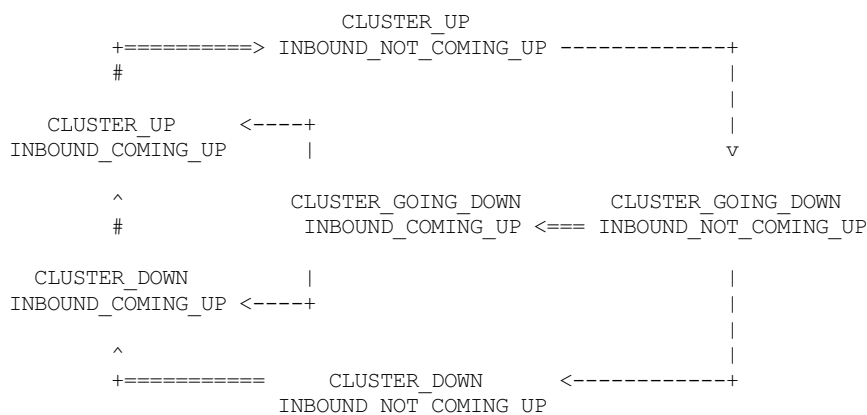
"cluster" state: The global state of the cluster; or the state on the outbound side:

- CLUSTER\_DOWN
- CLUSTER\_UP
- CLUSTER\_GOING\_DOWN

"inbound" state: The state of the cluster on the inbound side.

- INBOUND\_NOT\_COMING\_UP
- INBOUND\_COMING\_UP

The different pairings of these states results in six possible states for the cluster as a whole:



Transitions -----> can only be made by the outbound CPU, and only involve changes to the "cluster" state.

Transitions ===##> can only be made by the inbound CPU, and only involve changes to the "inbound" state, except where there is no further transition possible on the outbound side (i.e., the outbound CPU has put the cluster into the CLUSTER\_DOWN state).

The race avoidance algorithm does not provide a way to determine which exact CPUs within the cluster play these roles. This must be decided in advance by some other means. Refer to the section "Last man and first man selection" for more explanation.

CLUSTER\_DOWN/INBOUND\_NOT\_COMING\_UP is the only state where the cluster can actually be powered down.

The parallelism of the inbound and outbound CPUs is observed by the existence of two different paths from CLUSTER\_GOING\_DOWN/ INBOUND\_NOT\_COMING\_UP (corresponding to GOING\_DOWN in the basic model) to CLUSTER\_DOWN/INBOUND\_COMING\_UP (corresponding to COMING\_UP in the basic model). The second path avoids cluster teardown completely.

CLUSTER\_UP/INBOUND\_COMING\_UP is equivalent to UP in the basic model. The final transition to CLUSTER\_UP/INBOUND\_NOT\_COMING\_UP is trivial and merely resets the state machine ready for the next cycle.

Details of the allowable transitions follow.

The next state in each case is notated

<cluster state>/<inbound state> (<transitioner>)

where the <transitioner> is the side on which the transition can occur; either the inbound or the outbound side.

#### CLUSTER\_DOWN/INBOUND\_NOT\_COMING\_UP:

Next state:

CLUSTER\_DOWN/INBOUND\_COMING\_UP (inbound)

Conditions:

none

Trigger events:

- a. an explicit hardware power-up operation, resulting from a policy decision on another CPU;
- b. a hardware event, such as an interrupt.

#### CLUSTER\_DOWN/INBOUND\_COMING\_UP:

In this state, an inbound CPU sets up the cluster, including enabling of hardware coherency at the cluster level and any other operations (such as cache invalidation) which are required in order to achieve this.

The purpose of this state is to do sufficient cluster-level setup to enable other CPUs in the cluster to enter coherency safely.

Next state:

CLUSTER\_UP/INBOUND\_COMING\_UP (inbound)

Conditions:

cluster-level setup and hardware coherency complete

Trigger events:

(spontaneous)

#### CLUSTER\_UP/INBOUND\_COMING\_UP:

Cluster-level setup is complete and hardware coherency is enabled for the cluster. Other CPUs in the cluster can safely enter coherency.

This is a transient state, leading immediately to CLUSTER\_UP/INBOUND\_NOT\_COMING\_UP. All other CPUs on the cluster should consider treat these two states as equivalent.

Next state:

CLUSTER\_UP/INBOUND\_NOT\_COMING\_UP (inbound)

Conditions:

none

Trigger events:

(spontaneous)

#### CLUSTER\_UP/INBOUND\_NOT\_COMING\_UP:

Cluster-level setup is complete and hardware coherency is enabled for the cluster. Other CPUs in the cluster can safely enter coherency.

The cluster will remain in this state until a policy decision is made to power the cluster down.

Next state:

CLUSTER\_GOING\_DOWN/INBOUND\_NOT\_COMING\_UP (outbound)

Conditions:

none

Trigger events:

policy decision to power down the cluster

#### CLUSTER\_GOING\_DOWN/INBOUND\_NOT\_COMING\_UP:

An outbound CPU is tearing the cluster down. The selected CPU must wait in this state until all CPUs in the cluster are in the CPU\_DOWN state.

When all CPUs are in the CPU\_DOWN state, the cluster can be torn down, for example by cleaning data caches and exiting cluster-level coherency.

To avoid wasteful unnecessary teardown operations, the outbound should check the inbound cluster state for asynchronous transitions to INBOUND\_COMING\_UP. Alternatively, individual CPUs can be checked for entry into CPU\_COMING\_UP or CPU\_UP.

Next states:

CLUSTER\_DOWN/INBOUND\_NOT\_COMING\_UP (outbound)

Conditions:

cluster torn down and ready to power off

Trigger events:

(spontaneous)

CLUSTER\_GOING\_DOWN/INBOUND\_COMING\_UP (inbound)

Conditions:

none

Trigger events:

- a. an explicit hardware power-up operation, resulting from a policy decision on another CPU;
- b. a hardware event, such as an interrupt.

CLUSTER\_GOING\_DOWN/INBOUND\_COMING\_UP:

The cluster is (or was) being torn down, but another CPU has come online in the meantime and is trying to set up the cluster again.

If the outbound CPU observes this state, it has two choices:

- a. back out of teardown, restoring the cluster to the CLUSTER\_UP state;
- b. finish tearing the cluster down and put the cluster in the CLUSTER\_DOWN state; the inbound CPU will set up the cluster again from there.

Choice (a) permits the removal of some latency by avoiding unnecessary teardown and setup operations in situations where the cluster is not really going to be powered down.

Next states:

CLUSTER\_UP/INBOUND\_COMING\_UP (outbound)

Conditions:

cluster-level setup and hardware coherency complete

Trigger events:

(spontaneous)

CLUSTER\_DOWN/INBOUND\_COMING\_UP (outbound)

Conditions:

cluster torn down and ready to power off

Trigger events:

(spontaneous)

## Last man and First man selection

The CPU which performs cluster tear-down operations on the outbound side is commonly referred to as the "last man".

The CPU which performs cluster setup on the inbound side is commonly referred to as the "first man".

The race avoidance algorithm documented above does not provide a mechanism to choose which CPUs should play these roles.

Last man:

When shutting down the cluster, all the CPUs involved are initially executing Linux and hence coherent. Therefore, ordinary spinlocks can be used to select a last man safely, before the CPUs become non-coherent.

First man:

Because CPUs may power up asynchronously in response to external wake-up events, a dynamic mechanism is needed to make sure that only one CPU attempts to play the first man role and do the cluster-level initialisation: any other CPUs must wait for this to complete before proceeding.

Cluster-level initialisation may involve actions such as configuring coherency controls in the bus fabric.

The current implementation in `mcpm_head.S` uses a separate mutual exclusion mechanism to do this arbitration. This mechanism is documented in detail in `vlocks.txt`.

## Features and Limitations

Implementation:

The current ARM-based implementation is split between `arch/arm/common/mcpm_head.S` (low-level inbound CPU operations) and `arch/arm/common/mcpm_entry.c` (everything else):

`__mcpm_cpu_going_down()` signals the transition of a CPU to the CPU\_GOING\_DOWN state.

`__mcpm_cpu_down()` signals the transition of a CPU to the CPU\_DOWN state.

A CPU transitions to CPU\_COMING\_UP and then to CPU\_UP via the low-level power-up code in mcpm\_head.S. This could involve CPU-specific setup code, but in the current implementation it does not.

\_\_mcpm\_outbound\_enter\_critical() and \_\_mcpm\_outbound\_leave\_critical() handle transitions from CLUSTER\_UP to CLUSTER\_GOING\_DOWN and from there to CLUSTER\_DOWN or back to CLUSTER\_UP (in the case of an aborted cluster power-down).

These functions are more complex than the \_\_mcpm\_cpu\_\*(\*) functions due to the extra inter-CPU coordination which is needed for safe transitions at the cluster level.

A cluster transitions from CLUSTER\_DOWN back to CLUSTER\_UP via the low-level power-up code in mcpm\_head.S. This typically involves platform-specific setup code, provided by the platform-specific power\_up\_setup function registered via mcpm\_sync\_init.

Deep topologies:

As currently described and implemented, the algorithm does not support CPU topologies involving more than two levels (i.e., clusters of clusters are not supported). The algorithm could be extended by replicating the cluster-level states for the additional topological levels, and modifying the transition rules for the intermediate (non-outermost) cluster levels.

## Colophon

Originally created and documented by Dave Martin for Linaro Limited, in collaboration with Nicolas Pitre and Achin Gupta.

Copyright (C) 2012-2013 Linaro Limited Distributed under the terms of Version 2 of the GNU General Public License, as defined in linux/COPYING.