

---

# **Pacemaker Explained**

*Release 2.1.4*

**the Pacemaker project contributors**

**Jun 15, 2022**



# CONTENTS

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Abstract</b>   | <b>3</b> |
| <b>2</b> | <b>Table of Contents</b>  | <b>5</b> |
| 2.1      | Introduction . . . . .  | 5        |
| 2.1.1    | The Scope of this Document . . . . .                                | 5        |
| 2.1.2    | What Is Pacemaker? . . . . .  | 5        |
| 2.2      | Cluster-Wide Configuration . . . . .                                | 12       |
| 2.2.1    | Configuration Layout . . . . .                                      | 12       |
| 2.2.2    | CIB Properties . . . . .  | 13       |
| 2.2.3    | Cluster Options . . . . .   | 14       |
| 2.3      | Cluster Nodes . . . . .   | 19       |
| 2.3.1    | Defining a Cluster Node . . . . .                                   | 19       |
| 2.3.2    | Node Attributes . . . . .   | 19       |
| 2.3.3    | Tracking Node Health . . . . .                                      | 21       |
| 2.4      | Cluster Resources . . . . .   | 24       |
| 2.4.1    | What is a Cluster Resource? . . . . .                               | 24       |
| 2.4.2    | Resource Classes . . . . .  | 24       |
| 2.4.3    | Resource Properties . . . . .                                       | 27       |
| 2.4.4    | Resource Options . . . . .  | 28       |
| 2.4.5    | Resource Operations . . . . .                                       | 34       |
| 2.5      | Resource Constraints . . . . .                                      | 39       |
| 2.5.1    | Scores . . . . .  | 39       |
| 2.5.2    | Deciding Which Nodes a Resource Can Run On . . . . .                | 39       |
| 2.5.3    | Specifying the Order in which Resources Should Start/Stop . . . . . | 42       |
| 2.5.4    | Placing Resources Relative to other Resources . . . . .             | 43       |
| 2.5.5    | Resource Sets . . . . .   | 46       |
| 2.5.6    | Ordering Sets of Resources . . . . .                                | 46       |
| 2.5.7    | Colocating Sets of Resources . . . . .                              | 50       |
| 2.6      | Fencing . . . . .   | 53       |
| 2.6.1    | What Is Fencing? . . . . .  | 53       |
| 2.6.2    | Why Is Fencing Necessary? . . . . .                                 | 53       |
| 2.6.3    | Fence Devices . . . . .   | 54       |
| 2.6.4    | Fence Agents . . . . .  | 54       |
| 2.6.5    | When a Fence Device Can Be Used . . . . .                           | 54       |
| 2.6.6    | Limitations of Fencing Resources . . . . .                          | 55       |
| 2.6.7    | Special Meta-Attributes for Fencing Resources . . . . .             | 55       |
| 2.6.8    | Special Instance Attributes for Fencing Resources . . . . .         | 55       |
| 2.6.9    | Default Check Type . . . . .  | 59       |
| 2.6.10   | Unfencing . . . . .   | 59       |
| 2.6.11   | Fencing and Quorum . . . . .  | 59       |

|        |  |     |
|--------|--|-----|
| 2.6.12 | Fencing Timeouts . . . . .   | 60  |
| 2.6.13 | Fence Devices Dependent on Other Resources . . . . .                 | 61  |
| 2.6.14 | Configuring Fencing . . . . .  | 61  |
| 2.6.15 | Fencing Topologies . . . . .   | 67  |
| 2.6.16 | Remapping Reboots . . . . .  | 70  |
| 2.7    | Alerts . . . . .   | 71  |
| 2.7.1  | Alert Agents . . . . .   | 71  |
| 2.7.2  | Alert Recipients . . . . .   | 71  |
| 2.7.3  | Alert Meta-Attributes . . . . .                                      | 72  |
| 2.7.4  | Alert Instance Attributes . . . . .                                  | 73  |
| 2.7.5  | Alert Filters . . . . .  | 73  |
| 2.7.6  | Using the Sample Alert Agents . . . . .                              | 74  |
| 2.7.7  | Writing an Alert Agent . . . . .                                     | 75  |
| 2.8    | Rules . . . . .  | 76  |
| 2.8.1  | Rule Properties . . . . .  | 77  |
| 2.8.2  | Node Attribute Expressions . . . . .                                 | 77  |
| 2.8.3  | Date/Time Expressions . . . . .                                      | 79  |
| 2.8.4  | Resource Expressions . . . . .                                       | 83  |
| 2.8.5  | Operation Expressions . . . . .                                      | 83  |
| 2.8.6  | Using Rules to Determine Resource Location . . . . .                 | 84  |
| 2.8.7  | Using Rules to Define Options . . . . .                              | 85  |
| 2.9    | Advanced Configuration . . . . .                                     | 89  |
| 2.9.1  | Specifying When Recurring Actions are Performed . . . . .            | 89  |
| 2.9.2  | Handling Resource Failure . . . . .                                  | 89  |
| 2.9.3  | Moving Resources . . . . .   | 91  |
| 2.9.4  | Reloading an Agent After a Definition Change . . . . .               | 96  |
| 2.10   | Advanced Resource Types . . . . .                                    | 96  |
| 2.10.1 | Groups - A Syntactic Shortcut . . . . .                              | 96  |
| 2.10.2 | Clones - Resources That Can Have Multiple Active Instances . . . . . | 98  |
| 2.10.3 | Bundles - Containerized Resources . . . . .                          | 110 |
| 2.11   | Reusing Parts of the Configuration . . . . .                         | 116 |
| 2.11.1 | Reusing Resource Definitions . . . . .                               | 116 |
| 2.11.2 | Reusing Rules, Options and Sets of Operations . . . . .              | 119 |
| 2.11.3 | Tagging Configuration Elements . . . . .                             | 120 |
| 2.12   | Utilization and Placement Strategy . . . . .                         | 122 |
| 2.12.1 | Utilization attributes . . . . .                                     | 122 |
| 2.12.2 | Placement Strategy . . . . .   | 123 |
| 2.12.3 | Allocation Details . . . . .   | 124 |
| 2.12.4 | Limitations and Workarounds . . . . .                                | 125 |
| 2.13   | Access Control Lists (ACLs) . . . . .                                | 125 |
| 2.13.1 | ACL Prerequisites . . . . .  | 125 |
| 2.13.2 | ACL Configuration . . . . .  | 126 |
| 2.13.3 | ACL Roles . . . . .  | 126 |
| 2.13.4 | ACL Targets and Groups . . . . .                                     | 127 |
| 2.13.5 | ACL Examples . . . . .   | 127 |
| 2.13.6 | ACL Limitations . . . . .  | 130 |
| 2.14   | Status – Here be dragons . . . . .                                   | 130 |
| 2.14.1 | Node Status . . . . .  | 130 |
| 2.14.2 | Transient Node Attributes . . . . .                                  | 131 |
| 2.14.3 | Operation History . . . . .  | 132 |
| 2.15   | Multi-Site Clusters and Tickets . . . . .                            | 134 |
| 2.15.1 | Challenges for Multi-Site Clusters . . . . .                         | 135 |
| 2.15.2 | Conceptual Overview . . . . .  | 135 |
| 2.15.3 | Configuring Ticket Dependencies . . . . .                            | 136 |

|          |  |            |
|----------|--|------------|
| 2.15.4   | Managing Multi-Site Clusters . . . . . | 137        |
| 2.15.5   | For more information . . . . .         | 139        |
| 2.16     | Sample Configurations . . . . .        | 139        |
| 2.16.1   | Empty . . . . .                        | 139        |
| 2.16.2   | Simple . . . . .                       | 139        |
| 2.16.3   | Advanced Configuration . . . . .       | 140        |
| <b>3</b> | <b>Index</b>                           | <b>143</b> |
|          | <b>Index</b>                           | <b>145</b> |



*Configuring Pacemaker Clusters*





## ABSTRACT

This document definitively explains Pacemaker's features and capabilities, particularly the XML syntax used in Pacemaker's Cluster Information Base (CIB).



## TABLE OF CONTENTS

### 2.1 Introduction

#### 2.1.1 The Scope of this Document

This document is intended to be an exhaustive reference for configuring Pacemaker. To achieve this, it focuses on the XML syntax used to configure the CIB.

For those that are allergic to XML, multiple higher-level front-ends (both command-line and GUI) are available. These tools will not be covered in this document, though the concepts explained here should make the functionality of these tools more easily understood.

Users may be interested in other parts of the [Pacemaker documentation set](#), such as *Clusters from Scratch*, a step-by-step guide to setting up an example cluster, and *Pacemaker Administration*, a guide to maintaining a cluster.

#### 2.1.2 What Is Pacemaker?

Pacemaker is a high-availability *cluster resource manager* – software that runs on a set of hosts (a *cluster of nodes*) in order to preserve integrity and minimize downtime of desired services (*resources*).<sup>1</sup> It is maintained by the [ClusterLabs](#) community.

Pacemaker's key features include:

- Detection of and recovery from node- and service-level failures
- Ability to ensure data integrity by fencing faulty nodes
- Support for one or more nodes per cluster
- Support for multiple resource interface standards (anything that can be scripted can be clustered)
- Support (but no requirement) for shared storage
- Support for practically any redundancy configuration (active/passive, N+1, etc.)
- Automatically replicated configuration that can be updated from any node
- Ability to specify cluster-wide relationships between services, such as ordering, colocation and anti-colocation
- Support for advanced service types, such as *clones* (services that need to be active on multiple nodes), *promotable clones* (clones that can run in one of two roles), and containerized services

---

<sup>1</sup> *Cluster* is sometimes used in other contexts to refer to hosts grouped together for other purposes, such as high-performance computing (HPC), but Pacemaker is not intended for those purposes.

- Unified, scriptable cluster management tools

---

### Note: Fencing

*Fencing*, also known as *STONITH* (an acronym for Shoot The Other Node In The Head), is the ability to ensure that it is not possible for a node to be running a service. This is accomplished via *fence devices* such as intelligent power switches that cut power to the target, or intelligent network switches that cut the target's access to the local network.

Pacemaker represents fence devices as a special class of resource.

A cluster cannot safely recover from certain failure conditions, such as an unresponsive node, without fencing.

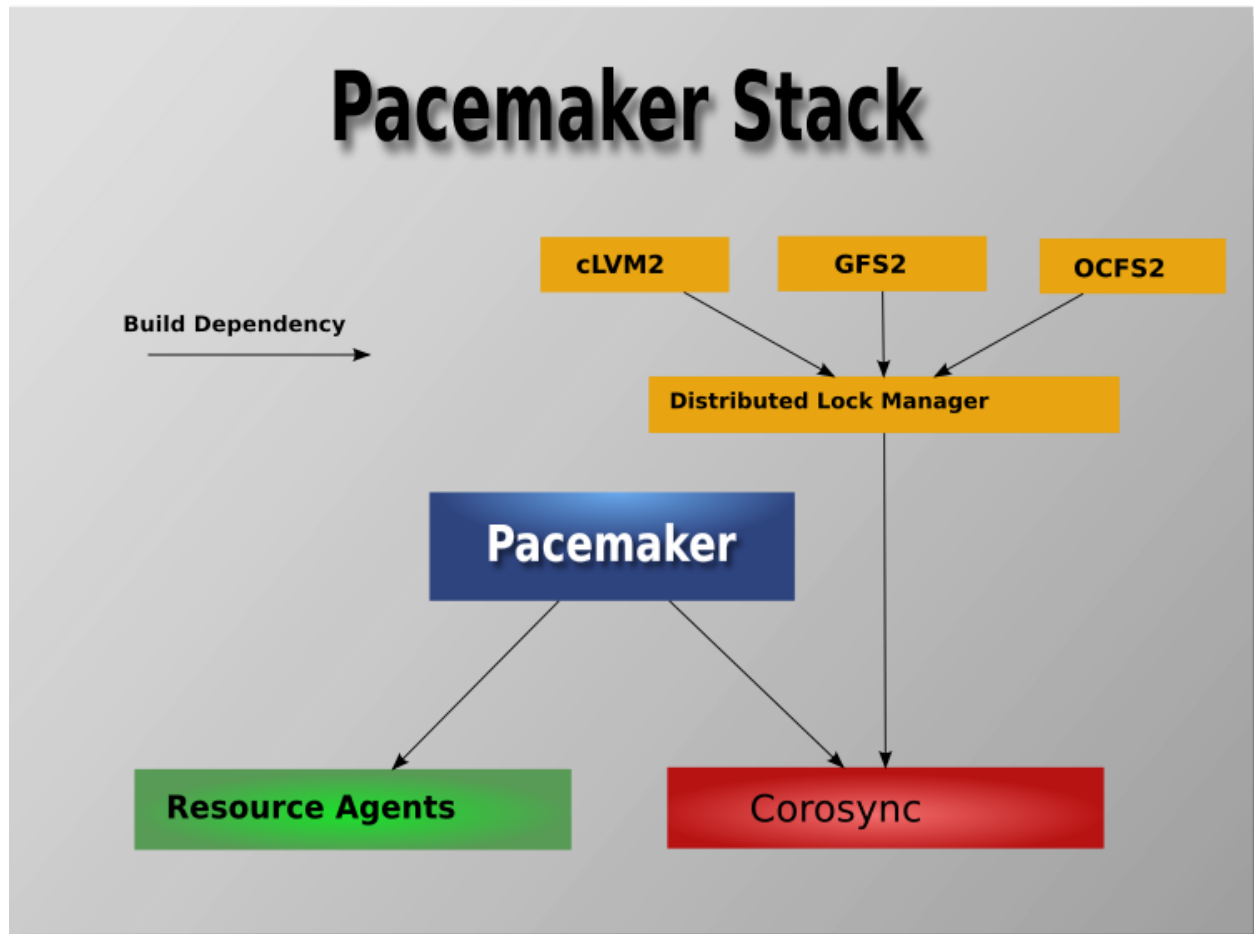
---

### Cluster Architecture

At a high level, a cluster can be viewed as having these parts (which together are often referred to as the *cluster stack*):

- **Resources:** These are the reason for the cluster's being – the services that need to be kept highly available.
- **Resource agents:** These are scripts or operating system components that start, stop, and monitor resources, given a set of resource parameters. These provide a uniform interface between Pacemaker and the managed services.
- **Fence agents:** These are scripts that execute node fencing actions, given a target and fence device parameters.
- **Cluster membership layer:** This component provides reliable messaging, membership, and quorum information about the cluster. Currently, Pacemaker supports [Corosync](#) as this layer.
- **Cluster resource manager:** Pacemaker provides the brain that processes and reacts to events that occur in the cluster. These events may include nodes joining or leaving the cluster; resource events caused by failures, maintenance, or scheduled activities; and other administrative actions. To achieve the desired availability, Pacemaker may start and stop resources and fence nodes.
- **Cluster tools:** These provide an interface for users to interact with the cluster. Various command-line and graphical (GUI) interfaces are available.

Most managed services are not, themselves, cluster-aware. However, many popular open-source cluster filesystems make use of a common *Distributed Lock Manager* (DLM), which makes direct use of Corosync for its messaging and membership capabilities and Pacemaker for the ability to fence nodes.

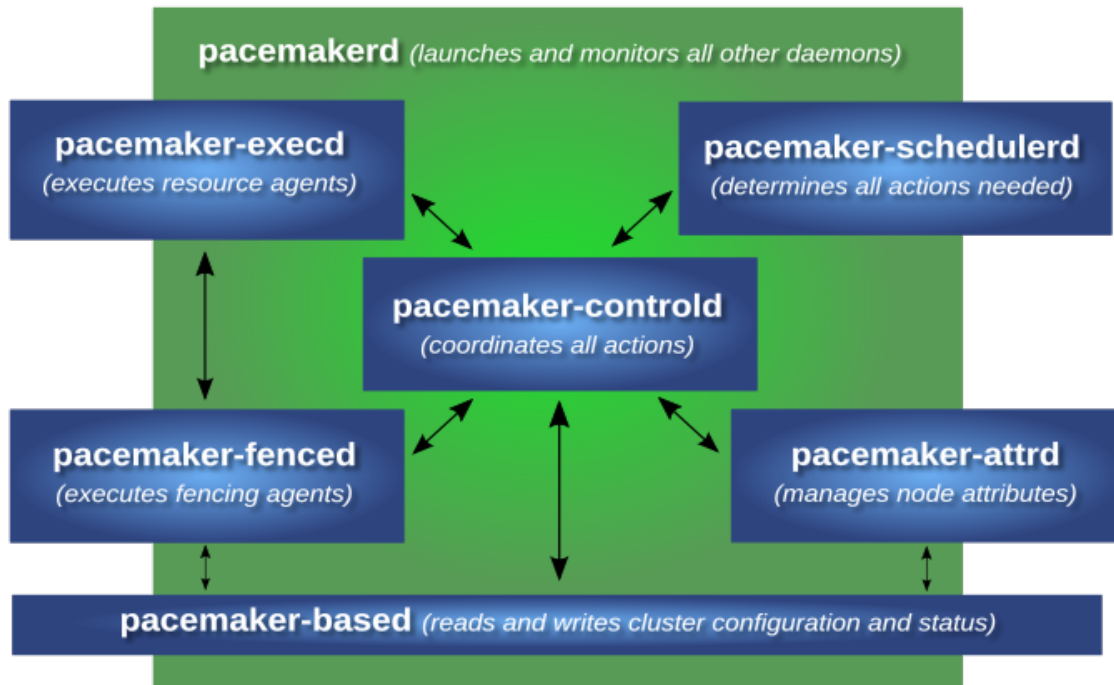


## Pacemaker Architecture

Pacemaker itself is composed of multiple daemons that work together:

- `pacemakerd`
- `pacemaker-attrd`
- `pacemaker-based`
- `pacemaker-controld`
- `pacemaker-execd`
- `pacemaker-fenced`
- `pacemaker-schedulerd`

# Pacemaker internals



## ClusterLabs

Pacemaker's main process (**pacemakerd**) spawns all the other daemons, and respawns them if they unexpectedly exit.

The *Cluster Information Base* (CIB) is an [XML](#) representation of the cluster's configuration and the state of all nodes and resources. The *CIB manager* (**pacemaker-based**) keeps the CIB synchronized across the cluster, and handles requests to modify it.

The *attribute manager* (**pacemaker-attribd**) maintains a database of attributes for all nodes, keeps it synchronized across the cluster, and handles requests to modify them. These attributes are usually recorded in the CIB.

Given a snapshot of the CIB as input, the *scheduler* (**pacemaker-schedulerd**) determines what actions are necessary to achieve the desired state of the cluster.

The *local executor* (**pacemaker-execd**) handles requests to execute resource agents on the local cluster node, and returns the result.

The *fencer* (**pacemaker-fenced**) handles requests to fence nodes. Given a target node, the fencer decides which cluster node(s) should execute which fencing device(s), and calls the necessary fencing agents (either directly, or via requests to the fencer peers on other nodes), and returns the result.

The *controller* (**pacemaker-controld**) is Pacemaker's coordinator, maintaining a consistent view of the cluster membership and orchestrating all the other components.

Pacemaker centralizes cluster decision-making by electing one of the controller instances as the *Designated Controller* (DC). Should the elected DC process (or the node it is on) fail, a new one is quickly established. The DC responds to cluster events by taking a current snapshot of the CIB, feeding it to the scheduler, then

asking the executors (either directly on the local node, or via requests to controller peers on other nodes) and the fencer to execute any necessary actions.

---

**Note: Old daemon names**

The Pacemaker daemons were renamed in version 2.0. You may still find references to the old names, especially in documentation targeted to version 1.1.

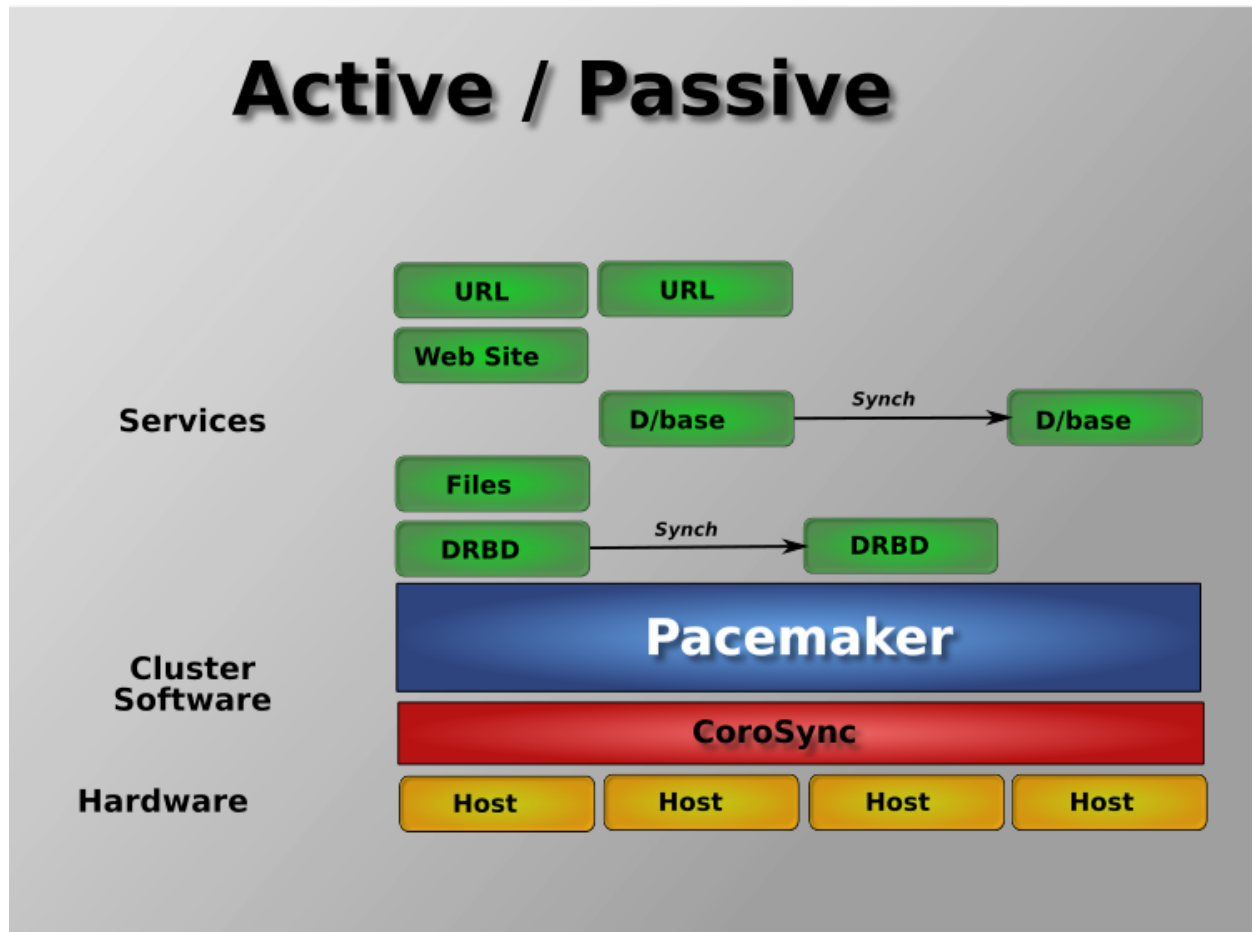
| Old name          | New name           |
|-------------------|--------------------|
| attrd             | pacemaker-attrd    |
| cib               | pacemaker-based    |
| crmd              | pacemaker-controld |
| lrmd              | pacemaker-execd    |
| stonithd          | pacemaker-fenced   |
| pacemaker_remoted | pacemaker-remoted  |

---

**Node Redundancy Designs**

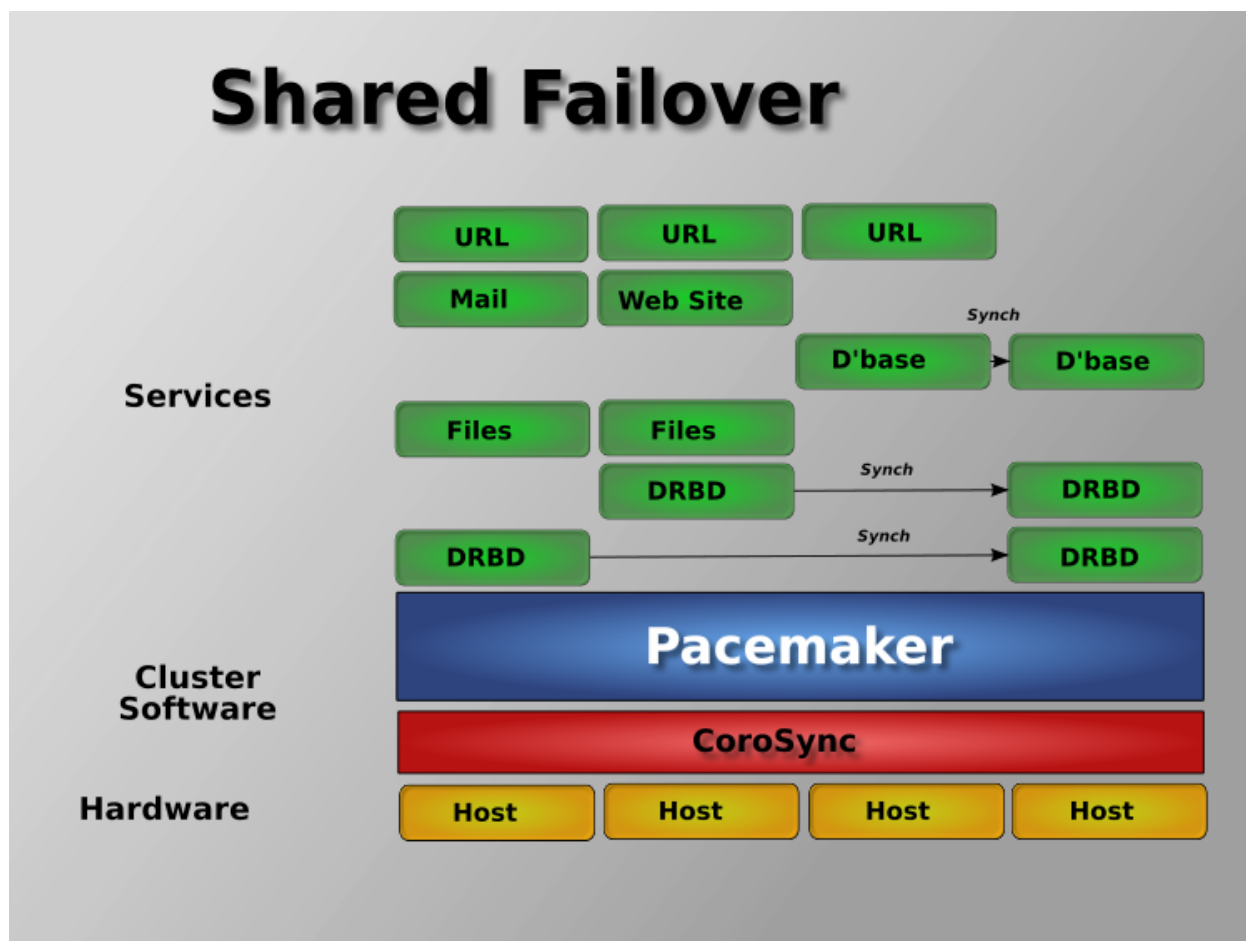
Pacemaker supports practically any [node redundancy configuration](#) including *Active/Active*, *Active/Passive*, *N+1*, *N+M*, *N-to-1* and *N-to-N*.

Active/passive clusters with two (or more) nodes using Pacemaker and [DRBD](#) are a cost-effective high-availability solution for many situations. One of the nodes provides the desired services, and if it fails, the other node takes over.

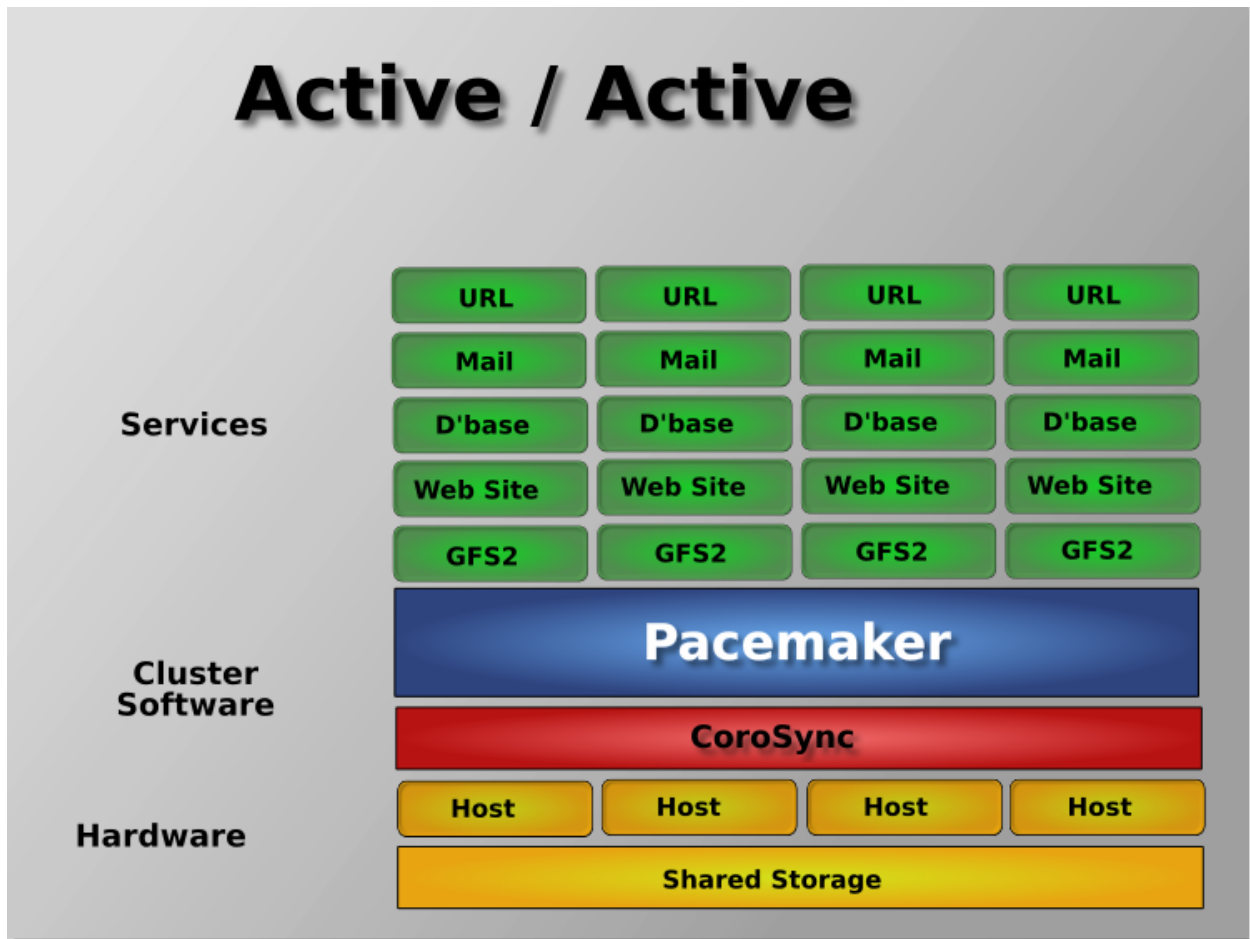


Pacemaker also supports multiple nodes in a shared-failover design, reducing hardware costs by allowing several active/passive clusters to be combined and share a common backup node.





When shared storage is available, every node can potentially be used for failover. Pacemaker can even run multiple copies of services to spread out the workload. This is sometimes called N to N Redundancy.



## 2.2 Cluster-Wide Configuration

### 2.2.1 Configuration Layout

The cluster is defined by the Cluster Information Base (CIB), which uses XML notation. The simplest CIB, an empty one, looks like this:

#### An empty configuration

```
<cib crm_feature_set="3.6.0" validate-with="pacemaker-3.5" epoch="1" num_updates="0" admin_epoch=
  ↪"0">
  <configuration>
    <crm_config/>
    <nodes/>
    <resources/>
    <constraints/>
  </configuration>
  <status/>
</cib>
```

The empty configuration above contains the major sections that make up a CIB:

- **cib**: The entire CIB is enclosed with a **cib** element. Certain fundamental settings are defined as attributes of this element.
  - **configuration**: This section – the primary focus of this document – contains traditional configuration information such as what resources the cluster serves and the relationships among them.
    - \* **crm\_config**: cluster-wide configuration options
    - \* **nodes**: the machines that host the cluster
    - \* **resources**: the services run by the cluster
    - \* **constraints**: indications of how resources should be placed
  - **status**: This section contains the history of each resource on each node. Based on this data, the cluster can construct the complete current state of the cluster. The authoritative source for this section is the local executor (pacemaker-execd process) on each cluster node, and the cluster will occasionally repopulate the entire section. For this reason, it is never written to disk, and administrators are advised against modifying it in any way.

In this document, configuration settings will be described as properties or options based on how they are defined in the CIB:

- Properties are XML attributes of an XML element.
- Options are name-value pairs expressed as **nvpair** child elements of an XML element.

Normally, you will use command-line tools that abstract the XML, so the distinction will be unimportant; both properties and options are cluster settings you can tweak.

## 2.2.2 CIB Properties

Certain settings are defined by CIB properties (that is, attributes of the **cib** tag) rather than with the rest of the cluster configuration in the **configuration** section.

The reason is simply a matter of parsing. These options are used by the configuration database which is, by design, mostly ignorant of the content it holds. So the decision was made to place them in an easy-to-find location.

Table 1: **CIB Properties**

| Attribute     | Description   |
|---------------|---|
| admin_epoch   | When a node joins the cluster, the cluster performs a check to see which node has the best configuration. It asks the node with the highest ( <b>admin_epoch</b> , <b>epoch</b> , <b>num_updates</b> ) tuple to replace the configuration on all the nodes – which makes setting them, and setting them correctly, very important. <b>admin_epoch</b> is never modified by the cluster; you can use this to make the configurations on any inactive nodes obsolete.<br><b>Warning:</b> Never set this value to zero. In such cases, the cluster cannot tell the difference between your configuration and the “empty” one used when nothing is found on disk. |
| epoch         | The cluster increments this every time the configuration is updated (usually by the administrator).   |
| num_updates   | The cluster increments this every time the configuration or status is updated (usually by the cluster) and resets it to 0 when epoch changes.   |
| validate-with | Determines the type of XML validation that will be done on the configuration. If set to <b>none</b> , the cluster will not verify that updates conform to the DTD (nor reject ones that don’t).   |

Continued on next page

Table 1 – continued from previous page

| Attribute        | Description   |
|------------------|---|
| cib-last-written | Indicates when the configuration was last written to disk. Maintained by the cluster; for informational purposes only.  |
| have-quorum      | Indicates if the cluster has quorum. If false, this may mean that the cluster cannot start resources or fence other nodes (see <code>no-quorum-policy</code> below). Maintained by the cluster. |
| dc-uuid          | Indicates which cluster node is the current leader. Used by the cluster when placing resources and determining the order of some events. Maintained by the cluster.                             |

### 2.2.3 Cluster Options

Cluster options, as you might expect, control how the cluster behaves when confronted with various situations.

They are grouped into sets within the `crm_config` section. In advanced configurations, there may be more than one set. (This will be described later in the chapter on *Rules* where we will show how to have the cluster use different sets of options during working hours than during weekends.) For now, we will describe the simple case where each option is present at most once.

You can obtain an up-to-date list of cluster options, including their default values, by running the `man pacemaker-schedulerd` and `man pacemaker-controld` commands.

Table 2: Cluster Options

| Option                 | Default | Description  |
|------------------------|---------|--|
| cluster-name           |         | An (optional) name for the cluster as a whole. This is mostly for users' convenience for use as desired in administration, but this can be used in the Pacemaker configuration in <i>Rules</i> (as the <code>#cluster-name</code> <i>node attribute</i> ). It may also be used by higher-level tools when displaying cluster information, and by certain resource agents (for example, the <code>ocf:heartbeat:GFS2</code> agent stores the cluster name in filesystem meta-data). |
| dc-version             |         | Version of Pacemaker on the cluster's DC. Determined automatically by the cluster. Often includes the hash which identifies the exact Git changeset it was built from. Used for diagnostic purposes.   |
| cluster-infrastructure |         | The messaging stack on which Pacemaker is currently running. Determined automatically by the cluster. Used for informational and diagnostic purposes.  |

Continued on next page

Table 2 – continued from previous page

| Option                 | Default | Description  |
|------------------------|---------|--|
| no-quorum-policy       | stop    | What to do when the cluster does not have quorum. Allowed values: <ul style="list-style-type: none"> <li>• <b>ignore</b>: continue all resource management</li> <li>• <b>freeze</b>: continue resource management, but don't recover resources from nodes not in the affected partition</li> <li>• <b>stop</b>: stop all resources in the affected cluster partition</li> <li>• <b>demote</b>: demote promotable resources and stop all other resources in the affected cluster partition (<i>since 2.0.5</i>)</li> <li>• <b>suicide</b>: fence all nodes in the affected cluster partition</li> </ul> |
| batch-limit            | 0       | The maximum number of actions that the cluster may execute in parallel across all nodes. The “correct” value will depend on the speed and load of your network and cluster nodes. If zero, the cluster will impose a dynamically calculated limit only when any node has high load. If -1, the cluster will not impose any limit.  |
| migration-limit        | -1      | The number of <i>live migration</i> actions that the cluster is allowed to execute in parallel on a node. A value of -1 means unlimited.   |
| symmetric-cluster      | true    | Whether resources can run on any node by default (if false, a resource is allowed to run on a node only if a <i>location constraint</i> enables it)  |
| stop-all-resources     | false   | Whether all resources should be disallowed from running (can be useful during maintenance)   |
| stop-orphan-resources  | true    | Whether resources that have been deleted from the configuration should be stopped. This value takes precedence over <b>is-managed</b> (that is, even unmanaged resources will be stopped when orphaned if this value is <b>true</b> )  |
| stop-orphan-actions    | true    | Whether recurring <i>operations</i> that have been deleted from the configuration should be cancelled  |
| start-failure-is-fatal | true    | Whether a failure to start a resource on a particular node prevents further start attempts on that node? If <b>false</b> , the cluster will decide whether the node is still eligible based on the resource's current failure count and <i>migration-threshold</i> .   |
| enable-startup-probes  | true    | Whether the cluster should check the pre-existing state of resources when the cluster starts   |
| maintenance-mode       | false   | Whether the cluster should refrain from monitoring, starting and stopping resources  |

Continued on next page

Table 2 – continued from previous page

| Option                   | Default | Description  |
|--------------------------|---------|--|
| stonith-enabled          | true    | Whether the cluster is allowed to fence nodes (for example, failed nodes and nodes with resources that can't be stopped). If true, at least one fence device must be configured before resources are allowed to run.<br>If false, unresponsive nodes are immediately assumed to be running no resources, and resource recovery on online nodes starts without any further protection (which can mean <i>data loss</i> if the unresponsive node still accesses shared storage, for example). See also the <i>requires</i> resource meta-attribute.  |
| stonith-action           | reboot  | Action the cluster should send to the fence agent when a node must be fenced. Allowed values are <b>reboot</b> , <b>off</b> , and (for legacy agents only) <b>poweroff</b> .   |
| stonith-timeout          | 60s     | How long to wait for <b>on</b> , <b>off</b> , and <b>reboot</b> fence actions to complete by default.  |
| stonith-max-attempts     | 10      | How many times fencing can fail for a target before the cluster will no longer immediately re-attempt it.  |
| stonith-watchdog-timeout | 0       | If nonzero, and the cluster detects <b>have-watchdog</b> as <b>true</b> , then watchdog-based self-fencing will be performed via SBD when fencing is required, without requiring a fencing resource explicitly configured.<br>If this is set to a positive value, unseen nodes are assumed to self-fence within this much time.<br><b>Warning:</b> It must be ensured that this value is larger than the <b>SBD_WATCHDOG_TIMEOUT</b> environment variable on all nodes. Pacemaker verifies the settings individually on all nodes and prevents startup or shuts down if configured wrongly on the fly. It is strongly recommended that <b>SBD_WATCHDOG_TIMEOUT</b> be set to the same value on all nodes.<br>If this is set to a negative value, and <b>SBD_WATCHDOG_TIMEOUT</b> is set, twice that value will be used.<br><b>Warning:</b> In this case, it is essential (and currently not verified by pacemaker) that <b>SBD_WATCHDOG_TIMEOUT</b> is set to the same value on all nodes. |
| concurrent-fencing       | false   | Whether the cluster is allowed to initiate multiple fence actions concurrently. Fence actions initiated externally, such as via the <b>stonith_admin</b> tool or an application such as DLM, or by the fencer itself such as recurring device monitors and <b>status</b> and <b>list</b> commands, are not limited by this option.   |
| fence-reaction           | stop    | How should a cluster node react if notified of its own fencing? A cluster node may receive notification of its own fencing if fencing is misconfigured, or if fabric fencing is in use that doesn't cut cluster communication. Allowed values are <b>stop</b> to attempt to immediately stop pacemaker and stay stopped, or <b>panic</b> to attempt to immediately reboot the local node, falling back to stop on failure. The default is likely to be changed to <b>panic</b> in a future release. ( <i>since 2.0.3</i> )   |

Continued on next page

Table 2 – continued from previous page

| Option                 | Default | Description   |
|------------------------|---------|---|
| priority-fencing-delay | 0       | Apply this delay to any fencing targeting the lost nodes with the highest total resource priority in case we don't have the majority of the nodes in our cluster partition, so that the more significant nodes potentially win any fencing match (especially meaningful in a split-brain of a 2-node cluster). A promoted resource instance takes the resource's priority plus 1 if the resource's priority is not 0. Any static or random delays introduced by <code>pcmk_delay_base</code> and <code>pcmk_delay_max</code> configured for the corresponding fencing resources will be added to this delay. This delay should be significantly greater than (safely twice) the maximum delay from those parameters. ( <i>since 2.0.4</i> ) |
| cluster-delay          | 60s     | Estimated maximum round-trip delay over the network (excluding action execution). If the DC requires an action to be executed on another node, it will consider the action failed if it does not get a response from the other node in this time (after considering the action's own timeout). The "correct" value will depend on the speed and load of your network and cluster nodes.   |
| dc-deadtime            | 20s     | How long to wait for a response from other nodes during startup. The "correct" value will depend on the speed/load of your network and the type of switches used.   |
| cluster-ipc-limit      | 500     | The maximum IPC message backlog before one cluster daemon will disconnect another. This is of use in large clusters, for which a good value is the number of resources in the cluster multiplied by the number of nodes. The default of 500 is also the minimum. Raise this if you see "Evicting client" messages for cluster daemon PIDs in the logs.  |
| pe-error-series-max    | -1      | The number of scheduler inputs resulting in errors to save. Used when reporting problems. A value of -1 means unlimited (report all), and 0 means none.   |
| pe-warn-series-max     | 5000    | The number of scheduler inputs resulting in warnings to save. Used when reporting problems. A value of -1 means unlimited (report all), and 0 means none.   |
| pe-input-series-max    | 4000    | The number of "normal" scheduler inputs to save. Used when reporting problems. A value of -1 means unlimited (report all), and 0 means none.  |
| enable-acl             | false   | Whether <i>Access Control Lists (ACLs)</i> should be used to authorize modifications to the CIB   |
| placement-strategy     | default | How the cluster should allocate resources to nodes (see <i>Utilization and Placement Strategy</i> ). Allowed values are <code>default</code> , <code>utilization</code> , <code>balanced</code> , and <code>minimal</code> .  |
| node-health-strategy   | none    | How the cluster should react to node health attributes (see <i>Tracking Node Health</i> ). Allowed values are <code>none</code> , <code>migrate-on-red</code> , <code>only-green</code> , <code>progressive</code> , and <code>custom</code> .  |
| node-health-base       | 0       | The base health score assigned to a node. Only used when <code>node-health-strategy</code> is <code>progressive</code> .  |
| node-health-green      | 0       | The score to use for a node health attribute whose value is <code>green</code> . Only used when <code>node-health-strategy</code> is <code>progressive</code> or <code>custom</code> .  |

Continued on next page

Table 2 – continued from previous page

| Option                   | Default | Description   |
|--------------------------|---------|---|
| node-health-yellow       | 0       | The score to use for a node health attribute whose value is <b>yellow</b> . Only used when <b>node-health-strategy</b> is <b>progressive</b> or <b>custom</b> .   |
| node-health-red          | 0       | The score to use for a node health attribute whose value is <b>red</b> . Only used when <b>node-health-strategy</b> is <b>progressive</b> or <b>custom</b> .  |
| cluster-recheck-interval | 15min   | Pacemaker is primarily event-driven, and looks ahead to know when to recheck the cluster for failure timeouts and most time-based rules ( <i>since 2.0.3</i> ). However, it will also recheck the cluster after this amount of inactivity. This has two goals: rules with <b>date_spec</b> are only guaranteed to be checked this often, and it also serves as a fail-safe for some kinds of scheduler bugs. A value of 0 disables this polling; positive values are a time interval.   |
| shutdown-lock            | false   | The default of false allows active resources to be recovered elsewhere when their node is cleanly shut down, which is what the vast majority of users will want. However, some users prefer to make resources highly available only for failures, with no recovery for clean shutdowns. If this option is true, resources active on a node when it is cleanly shut down are kept “locked” to that node (not allowed to run elsewhere) until they start again on that node after it rejoins (or for at most <b>shutdown-lock-limit</b> , if set). Stonith resources and Pacemaker Remote connections are never locked. Clone and bundle instances and the promoted role of promotable clones are currently never locked, though support could be added in a future release. Locks may be manually cleared using the <b>--refresh</b> option of <b>crm_resource</b> (both the resource and node must be specified; this works with remote nodes if their connection resource’s <b>target-role</b> is set to <b>Stopped</b> , but not if Pacemaker Remote is stopped on the remote node without disabling the connection resource). ( <i>since 2.0.4</i> ) |
| shutdown-lock-limit      | 0       | If <b>shutdown-lock</b> is true, and this is set to a nonzero time duration, locked resources will be allowed to start after this much time has passed since the node shutdown was initiated, even if the node has not rejoined. (This works with remote nodes only if their connection resource’s <b>target-role</b> is set to <b>Stopped</b> .) ( <i>since 2.0.4</i> )  |
| remove-after-stop        | false   | <i>Deprecated</i> Should the cluster remove resources from Pacemaker’s executor after they are stopped? Values other than the default are, at best, poorly tested and potentially dangerous. This option is deprecated and will be removed in a future release.   |
| startup-fencing          | true    | <i>Advanced Use Only:</i> Should the cluster fence unseen nodes at start-up? Setting this to false is unsafe, because the unseen nodes could be active and running resources but unreachable.   |
| election-timeout         | 2min    | <i>Advanced Use Only:</i> If you need to adjust this value, it probably indicates the presence of a bug.  |

Continued on next page



Table 2 – continued from previous page

| Option                    | Default | Description  |
|---------------------------|---------|--|
| shutdown-escalation       | 20min   | <i>Advanced Use Only:</i> If you need to adjust this value, it probably indicates the presence of a bug.   |
| join-integration-timeout  | 3min    | <i>Advanced Use Only:</i> If you need to adjust this value, it probably indicates the presence of a bug.   |
| join-finalization-timeout | 30min   | <i>Advanced Use Only:</i> If you need to adjust this value, it probably indicates the presence of a bug.   |
| transition-delay          | 0s      | <i>Advanced Use Only:</i> Delay cluster recovery for the configured interval to allow for additional or related events to occur. This can be useful if your configuration is sensitive to the order in which ping updates arrive. Enabling this option will slow down cluster recovery under all conditions. |

## 2.3 Cluster Nodes

### 2.3.1 Defining a Cluster Node

Each cluster node will have an entry in the **nodes** section containing at least an ID and a name. A cluster node's ID is defined by the cluster layer (Corosync).

#### Example Corosync cluster node entry

```
<node id="101" uname="pcmk-1"/>
```

In normal circumstances, the admin should let the cluster populate this information automatically from the cluster layer.

#### Where Pacemaker Gets the Node Name

The name that Pacemaker uses for a node in the configuration does not have to be the same as its local hostname. Pacemaker uses the following for a Corosync node's name, in order of most preferred first:

- The value of **name** in the **nodelist** section of **corosync.conf**
- The value of **ring0\_addr** in the **nodelist** section of **corosync.conf**
- The local hostname (value of **uname -n**)

If the cluster is running, the **crm\_node -n** command will display the local node's name as used by the cluster.

If a Corosync **nodelist** is used, **crm\_node --name-for-id** with a Corosync node ID will display the name used by the node with the given Corosync **nodeid**, for example:

```
crm_node --name-for-id 2
```

### 2.3.2 Node Attributes

Pacemaker allows node-specific values to be specified using *node attributes*. A node attribute has a name, and may have a distinct value for each node.

Node attributes come in two types, *permanent* and *transient*. Permanent node attributes are kept within the **node** entry, and keep their values even if the cluster restarts on a node. Transient node attributes are kept in the CIB's **status** section, and go away when the cluster stops on the node.

While certain node attributes have specific meanings to the cluster, they are mainly intended to allow administrators and resource agents to track any information desired.

For example, an administrator might choose to define node attributes for how much RAM and disk space each node has, which OS each uses, or which server room rack each node is in.

Users can configure *Rules* that use node attributes to affect where resources are placed.

### Setting and querying node attributes

Node attributes can be set and queried using the **crm\_attribute** and **attrd\_updater** commands, so that the user does not have to deal with XML configuration directly.

Here is an example command to set a permanent node attribute, and the XML configuration that would be generated:

#### Result of using **crm\_attribute** to specify which kernel **pcmk-1** is running

```
# crm_attribute --type nodes --node pcmk-1 --name kernel --update $(uname -r)

<node id="1" uname="pcmk-1">
  <instance_attributes id="nodes-1-attributes">
    <nvpair id="nodes-1-kernel" name="kernel" value="3.10.0-862.14.4.el7.x86_64"/>
  </instance_attributes>
</node>
```

To read back the value that was just set:

```
# crm_attribute --type nodes --node pcmk-1 --name kernel --query
scope=nodes name=kernel value=3.10.0-862.14.4.el7.x86_64
```

The **--type nodes** indicates that this is a permanent node attribute; **--type status** would indicate a transient node attribute.

### Special node attributes

Certain node attributes have special meaning to the cluster.

Node attribute names beginning with **#** are considered reserved for these special attributes. Some special attributes do not start with **#**, for historical reasons.

Certain special attributes are set automatically by the cluster, should never be modified directly, and can be used only within *Rules*; these are listed under *built-in node attributes*.

For true/false values, the cluster considers a value of “1”, “y”, “yes”, “on”, or “true” (case-insensitively) to be true, “0”, “n”, “no”, “off”, “false”, or unset to be false, and anything else to be an error.

Table 3: Node attributes with special significance

| Name                       | Description   |
|----------------------------|---|
| fail-count-*               | Attributes whose names start with <code>fail-count-</code> are managed by the cluster to track how many times particular resource operations have failed on this node. These should be queried and cleared via the <code>crm_failcount</code> or <code>crm_resource --cleanup</code> commands rather than directly.   |
| last-failure-*             | Attributes whose names start with <code>last-failure-</code> are managed by the cluster to track when particular resource operations have most recently failed on this node. These should be cleared via the <code>crm_failcount</code> or <code>crm_resource --cleanup</code> commands rather than directly.   |
| maintenance                | Similar to the <code>maintenance-mode</code> <i>cluster option</i> , but for a single node. If true, resources will not be started or stopped on the node, resources and individual clone instances running on the node will become unmanaged, and any recurring operations for those will be cancelled.<br><b>Warning:</b> Restarting pacemaker on a node that is in single-node maintenance mode will likely lead to undesirable effects. If <code>maintenance</code> is set as a transient attribute, it will be erased when Pacemaker is stopped, which will immediately take the node out of maintenance mode and likely get it fenced. Even if permanent, if Pacemaker is restarted, any resources active on the node will have their local history erased when the node rejoins, so the cluster will no longer consider them running on the node and thus will consider them managed again, leading them to be started elsewhere. This behavior might be improved in a future release. |
| probe_complete             | This is managed by the cluster to detect when nodes need to be reprobbed, and should never be used directly.  |
| resource-discovery-enabled | If the node is a remote node, fencing is enabled, and this attribute is explicitly set to false (unset means true in this case), resource discovery (probes) will not be done on this node. This is highly discouraged; the <code>resource-discovery</code> location constraint property is preferred for this purpose.   |
| shutdown                   | This is managed by the cluster to orchestrate the shutdown of a node, and should never be used directly.  |
| site-name                  | If set, this will be used as the value of the <code>#site-name</code> node attribute used in rules. (If not set, the value of the <code>cluster-name</code> cluster option will be used as <code>#site-name</code> instead.)  |
| standby                    | If true, the node is in standby mode. This is typically set and queried via the <code>crm_standby</code> command rather than directly.  |
| terminate                  | If the value is true or begins with any nonzero number, the node will be fenced. This is typically set by tools rather than directly.   |
| #digests-*                 | Attributes whose names start with <code>#digests-</code> are managed by the cluster to detect when <i>Unfencing</i> needs to be redone, and should never be used directly.  |
| #node-unfenced             | When the node was last unfenced (as seconds since the epoch). This is managed by the cluster and should never be used directly.   |

### 2.3.3 Tracking Node Health

A node may be functioning adequately as far as cluster membership is concerned, and yet be “unhealthy” in some respect that makes it an undesirable location for resources. For example, a disk drive may be reporting

SMART errors, or the CPU may be highly loaded.

Pacemaker offers a way to automatically move resources off unhealthy nodes.

### Node Health Attributes

Pacemaker will treat any node attribute whose name starts with `#health` as an indicator of node health. Node health attributes may have one of the following values:

Table 4: Allowed Values for Node Health Attributes

| Value               | Intended significance  |
|---------------------|--|
| <code>red</code>    | This indicator is unhealthy  |
| <code>yellow</code> | This indicator is becoming unhealthy   |
| <code>green</code>  | This indicator is healthy  |
| <i>integer</i>      | A numeric score to apply to all resources on this node (0 or positive is healthy, negative is unhealthy) |

### Node Health Strategy

Pacemaker assigns a node health score to each node, as the sum of the values of all its node health attributes. This score will be used as a location constraint applied to this node for all resources.

The `node-health-strategy` cluster option controls how Pacemaker responds to changes in node health attributes, and how it translates `red`, `yellow`, and `green` to scores.

Allowed values are:

Table 5: Node Health Strategies

| Value                       | Effect   |
|-----------------------------|--|
| <code>none</code>           | Do not track node health attributes at all.  |
| <code>migrate-on-red</code> | Assign the value of <code>-INFINITY</code> to <code>red</code> , and 0 to <code>yellow</code> and <code>green</code> . This will cause all resources to move off the node if any attribute is <code>red</code> .   |
| <code>only-green</code>     | Assign the value of <code>-INFINITY</code> to <code>red</code> and <code>yellow</code> , and 0 to <code>green</code> . This will cause all resources to move off the node if any attribute is <code>red</code> or <code>yellow</code> .  |
| <code>progressive</code>    | Assign the value of the <code>node-health-red</code> cluster option to <code>red</code> , the value of <code>node-health-yellow</code> to <code>yellow</code> , and the value of <code>node-health-green</code> to <code>green</code> . Each node is additionally assigned a score of <code>node-health-base</code> (this allows resources to start even if some attributes are <code>yellow</code> ). This strategy gives the administrator finer control over how important each value is. |
| <code>custom</code>         | Track node health attributes using the same values as <code>progressive</code> for <code>red</code> , <code>yellow</code> , and <code>green</code> , but do not take them into account. The administrator is expected to implement a policy by defining <i>Rules</i> referencing node health attributes.   |

## Exempting a Resource from Health Restrictions

If you want a resource to be able to run on a node even if its health score would otherwise prevent it, set the resource's `allow-unhealthy-nodes` meta-attribute to `true` (*available since 2.1.3*).

This is particularly useful for node health agents, to allow them to detect when the node becomes healthy again. If you configure a health agent without this setting, then the health agent will be banned from an unhealthy node, and you will have to investigate and clear the health attribute manually once it is healthy to allow resources on the node again.

If you want the meta-attribute to apply to a clone, it must be set on the clone itself, not on the resource being cloned.

## Configuring Node Health Agents

Since Pacemaker calculates node health based on node attributes, any method that sets node attributes may be used to measure node health. The most common are resource agents and custom daemons.

Pacemaker provides examples that can be used directly or as a basis for custom code. The `ocf:pacemaker:HealthCPU`, `ocf:pacemaker:HealthIOWait`, and `ocf:pacemaker:HealthSMART` resource agents set node health attributes based on CPU and disk status.

To take advantage of this feature, add the resource to your cluster (generally as a cloned resource with a recurring monitor action, to continually check the health of all nodes). For example:

### Example HealthIOWait resource configuration

```
<clone id="resHealthIOWait-clone">
  <primitive class="ocf" id="HealthIOWait" provider="pacemaker" type="HealthIOWait">
    <instance_attributes id="resHealthIOWait-instance_attributes">
      <nvpair id="resHealthIOWait-instance_attributes-red_limit" name="red_limit" value="30"/>
      <nvpair id="resHealthIOWait-instance_attributes-yellow_limit" name="yellow_limit" value="10"/>
    </instance_attributes>
    <operations>
      <op id="resHealthIOWait-monitor-interval-5" interval="5" name="monitor" timeout="5"/>
      <op id="resHealthIOWait-start-interval-0s" interval="0s" name="start" timeout="10s"/>
      <op id="resHealthIOWait-stop-interval-0s" interval="0s" name="stop" timeout="10s"/>
    </operations>
  </primitive>
</clone>
```

The resource agents use `attrd_updater` to set proper status for each node running this resource, as a node attribute whose name starts with `#health` (for `HealthIOWait`, the node attribute is named `#health-iowait`).

When a node is no longer faulty, you can force the cluster to make it available to take resources without waiting for the next monitor, by setting the node health attribute to green. For example:

### Force node1 to be marked as healthy

```
# attrd_updater --name "#health-iowait" --update "green" --node "node1"
```

## 2.4 Cluster Resources

### 2.4.1 What is a Cluster Resource?

A resource is a service made highly available by a cluster. The simplest type of resource, a *primitive* resource, is described in this chapter. More complex forms, such as groups and clones, are described in later chapters.

Every primitive resource has a *resource agent*. A resource agent is an external program that abstracts the service it provides and present a consistent view to the cluster.

This allows the cluster to be agnostic about the resources it manages. The cluster doesn't need to understand how the resource works because it relies on the resource agent to do the right thing when given a **start**, **stop** or **monitor** command. For this reason, it is crucial that resource agents are well-tested.

Typically, resource agents come in the form of shell scripts. However, they can be written using any technology (such as C, Python or Perl) that the author is comfortable with.

### 2.4.2 Resource Classes

Pacemaker supports several classes of agents:

- OCF
- LSB
- Systemd
- Upstart (deprecated)
- Service
- Fencing
- Nagios Plugins

#### Open Cluster Framework

The OCF standard<sup>1</sup> is basically an extension of the Linux Standard Base conventions for init scripts to:

- support parameters,
- make them self-describing, and
- make them extensible

OCF specs have strict definitions of the exit codes that actions must return<sup>2</sup>.

The cluster follows these specifications exactly, and giving the wrong exit code will cause the cluster to behave in ways you will likely find puzzling and annoying. In particular, the cluster needs to distinguish a completely stopped resource from one which is in some erroneous and indeterminate state.

Parameters are passed to the resource agent as environment variables, with the special prefix `OCF_RESKEY_`. So, a parameter which the user thinks of as `ip` will be passed to the resource agent as `OCF_RESKEY_ip`. The number and purpose of the parameters is left to the resource agent; however, the resource agent should use the **meta-data** command to advertise any that it supports.

---

<sup>1</sup> See <https://github.com/ClusterLabs/OCF-spec/tree/main/ra>. The Pacemaker implementation has been somewhat extended from the OCF specs.

<sup>2</sup> The resource-agents source code includes the **ocf-tester** script, which can be useful in this regard.

The OCF class is the most preferred as it is an industry standard, highly flexible (allowing parameters to be passed to agents in a non-positional manner) and self-describing.

For more information, see the [reference](#) and the *Resource Agents* chapter of *Pacemaker Administration*.

## Linux Standard Base

*LSB* resource agents are more commonly known as *init scripts*. If a full path is not given, they are assumed to be located in `/etc/init.d`.

Commonly, they are provided by the OS distribution. In order to be used with a Pacemaker cluster, they must conform to the LSB specification<sup>3</sup>.

**Warning:** Many distributions or particular software packages claim LSB compliance but ship with broken init scripts. For details on how to check whether your init script is LSB-compatible, see the *Resource Agents* chapter of *Pacemaker Administration*. Common problematic violations of the LSB standard include:

- Not implementing the **status** operation at all
- Not observing the correct exit status codes for **start/stop/status** actions
- Starting a started resource returns an error
- Stopping a stopped resource returns an error

---

**Important:** Remember to make sure the computer is *not* configured to start any services at boot time – that should be controlled by the cluster.

---

## Systemd

Most Linux distributions have replaced the old *SysV* style of initialization daemons and scripts with *Systemd*.

Pacemaker is able to manage these services *if they are present*.

Instead of init scripts, systemd has *unit files*. Generally, the services (unit files) are provided by the OS distribution, but there are online guides for converting from init scripts<sup>4</sup>.

---

**Important:** Remember to make sure the computer is *not* configured to start any services at boot time – that should be controlled by the cluster.

---

## Upstart

Some distributions replaced the old *SysV* style of initialization daemons (and scripts) with *Upstart*.

Pacemaker is able to manage these services *if they are present*.

Instead of init scripts, Upstart has *jobs*. Generally, the services (jobs) are provided by the OS distribution.

---

<sup>3</sup> See [http://refspecs.linux-foundation.org/LSB\\_3.0.0/LSB-Core-generic/LSB-Core-generic/iniscrptact.html](http://refspecs.linux-foundation.org/LSB_3.0.0/LSB-Core-generic/LSB-Core-generic/iniscrptact.html) for the LSB Spec as it relates to init scripts.

<sup>4</sup> For example, <http://0pointer.de/blog/projects/systemd-for-admins-3.html>

---

**Important:** Remember to make sure the computer is *not* configured to start any services at boot time – that should be controlled by the cluster.

---

**Warning:** Upstart support is deprecated in Pacemaker. Upstart is no longer an actively maintained project, and test platforms for it are no longer readily usable. Support will likely be dropped entirely at the next major release of Pacemaker.

### System Services

Since there are various types of system services (**systemd**, **upstart**, and **lsb**), Pacemaker supports a special **service** alias which intelligently figures out which one applies to a given cluster node.

This is particularly useful when the cluster contains a mix of **systemd**, **upstart**, and **lsb**.

In order, Pacemaker will try to find the named service as:

- an LSB init script
- a Systemd unit file
- an Upstart job

### STONITH

The STONITH class is used exclusively for fencing-related resources. This is discussed later in *Fencing*.

### Nagios Plugins

Nagios Plugins<sup>5</sup> are a way to monitor services. Pacemaker can use these as resources, to react to a change in the service's status.

To use plugins as resources, Pacemaker must have been built with support, and OCF-style meta-data for the plugins must be installed on nodes that can run them. Meta-data for several common plugins is provided by the [nagios-agents-metadata](#) project.

The supported parameters for such a resource are same as the long options of the plugin.

Start and monitor actions for plugin resources are implemented as invoking the plugin. A plugin result of “OK” (0) is treated as success, a result of “WARN” (1) is treated as a successful but degraded service, and any other result is considered a failure.

A plugin resource is not going to change its status after recovery by restarting the plugin, so using them alone does not make sense with **on-fail** set (or left to default) to **restart**. Another value could make sense, for example, if you want to fence or standby nodes that cannot reach some external service.

A more common use case for plugin resources is to configure them with a **container** meta-attribute set to the name of another resource that actually makes the service available, such as a virtual machine or container.

With **container** set, the plugin resource will automatically be colocated with the containing resource and ordered after it, and the containing resource will be considered failed if the plugin resource fails. This

---

<sup>5</sup> The project has two independent forks, hosted at <https://www.nagios-plugins.org/> and <https://www.monitoring-plugins.org/>. Output from both projects' plugins is similar, so plugins from either project can be used with pacemaker.



allows monitoring of a service inside a virtual machine or container, with recovery of the virtual machine or container if the service fails.

Configuring a virtual machine as a guest node, or a container as a *bundle*, is the preferred way of monitoring a service inside, but plugin resources can be useful when it is not practical to modify the virtual machine or container image for this purpose.

### 2.4.3 Resource Properties

These values tell the cluster which resource agent to use for the resource, where to find that resource agent and what standards it conforms to.

Table 6: **Properties of a Primitive Resource**

| Field    | Description  |
|----------|--|
| id       | Your name for the resource   |
| class    | The standard the resource agent conforms to. Allowed values: <code>lsb</code> , <code>nagios</code> , <code>ocf</code> , <code>service</code> , <code>stonith</code> , <code>systemd</code> , <code>upstart</code> |
| type     | The name of the Resource Agent you wish to use. E.g. <code>IPaddr</code> or <code>Filesystem</code>  |
| provider | The OCF spec allows multiple vendors to supply the same resource agent. To use the OCF resource agents supplied by the Heartbeat project, you would specify <code>heartbeat</code> here.                           |

The XML definition of a resource can be queried with the `crm_resource` tool. For example:

```
# crm_resource --resource Email --query-xml
```

might produce:

#### A system resource definition

```
<primitive id="Email" class="service" type="exim"/>
```

**Note:** One of the main drawbacks to system services (LSB, systemd or Upstart) resources is that they do not allow any parameters!

#### An OCF resource definition

```
<primitive id="Public-IP" class="ocf" type="IPaddr" provider="heartbeat">
  <instance_attributes id="Public-IP-params">
    <nvpair id="Public-IP-ip" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

## 2.4.4 Resource Options

Resources have two types of options: *meta-attributes* and *instance attributes*. Meta-attributes apply to any type of resource, while instance attributes are specific to each resource agent.

### Resource Meta-Attributes

Meta-attributes are used by the cluster to decide how a resource should behave and can be easily set using the `--meta` option of the `crm_resource` command.

Table 7: Meta-attributes of a Primitive Resource

| Field               | Default   | Description  |
|---------------------|---|--|
| priority            | 0   | If not all resources can be active, the cluster will stop lower priority resources in order to keep higher priority ones active.   |
| critical            | true  | Use this value as the default for <b>influence</b> in all <i>colocation constraints</i> involving this resource, as well as the implicit colocation constraints created if this resource is in a <i>group</i> . For details, see <i>Colocation Influence</i> . (since 2.1.0)   |
| target-role         | Started   | What state should the cluster attempt to keep this resource in? Allowed values: <ul style="list-style-type: none"> <li>• <b>Stopped</b>: Force the resource to be stopped</li> <li>• <b>Started</b>: Allow the resource to be started (and in the case of <i>promotable clone resources</i>, promoted if appropriate)</li> <li>• <b>Unpromoted</b>: Allow the resource to be started, but only in the unpromoted role if the resource is <i>promotable</i></li> <li>• <b>Promoted</b>: Equivalent to <b>Started</b></li> </ul> |
| is-managed          | TRUE  | Is the cluster allowed to start and stop the resource? Allowed values: <b>true</b> , <b>false</b>  |
| maintenance         | FALSE   | Similar to the <b>maintenance-mode cluster option</b> , but for a single resource. If true, the resource will not be started, stopped, or monitored on any node. This differs from <b>is-managed</b> in that monitors will not be run. Allowed values: <b>true</b> , <b>false</b>  |
| resource-stickiness | 1 for individual clone instances, 0 for all other resources | A score that will be added to the current node when a resource is already active. This allows running resources to stay where they are, even if they would be placed elsewhere if they were being started from a stopped state.  |

Continued on next page

Table 7 – continued from previous page

| Field               | Default   | Description   |
|---------------------|---|---|
| requires            | <code>quorum</code> for resources with a class of <code>stonith</code> , otherwise <code>unfencing</code> if <code>unfencing</code> is active in the cluster, otherwise <code>fencing</code> if <code>stonith-enabled</code> is <code>true</code> , otherwise <code>quorum</code> | <p>Conditions under which the resource can be started. Allowed values:</p> <ul style="list-style-type: none"> <li>• <b>nothing</b>: can always be started</li> <li>• <b>quorum</b>: The cluster can only start this resource if a majority of the configured nodes are active</li> <li>• <b>fencing</b>: The cluster can only start this resource if a majority of the configured nodes are active <i>and</i> any failed or unknown nodes have been <i>fenced</i></li> <li>• <b>unfencing</b>: The cluster can only start this resource if a majority of the configured nodes are active <i>and</i> any failed or unknown nodes have been fenced <i>and</i> only on nodes that have been <i>unfenced</i></li> </ul> |
| migration-threshold | INFINITY  | How many failures may occur for this resource on a node, before this node is marked ineligible to host this resource. A value of 0 indicates that this feature is disabled (the node will never be marked ineligible); by contrast, the cluster treats INFINITY (the default) as a very large but finite number. This option has an effect only if the failed operation specifies <code>on-fail</code> as <code>restart</code> (the default), and additionally for failed <code>start</code> operations, if the cluster property <code>start-failure-is-fatal</code> is <code>false</code> .  |
| failure-timeout     | 0   | How many seconds to wait before acting as if the failure had not occurred, and potentially allowing the resource back to the node on which it failed. A value of 0 indicates that this feature is disabled.   |

Continued on next page

Table 7 – continued from previous page

| Field                      | Default  | Description  |
|----------------------------|--|--|
| multiple-active            | stop_start   | What should the cluster do if it ever finds the resource active on more than one node?<br>Allowed values: <ul style="list-style-type: none"> <li>• <b>block</b>: mark the resource as unmanaged</li> <li>• <b>stop_only</b>: stop all active instances and leave them that way</li> <li>• <b>stop_start</b>: stop all active instances and start the resource in one location only</li> <li>• <b>stop_unexpected</b>: stop all active instances except where the resource should be active (this should be used only when extra instances are not expected to disrupt existing instances, and the resource agent’s monitor of an existing instance is capable of detecting any problems that could be caused; note that any resources ordered after this will still need to be restarted)</li> </ul> |
| allow-migrate              | TRUE for ocf:pacemaker:remote resources, FALSE otherwise | Whether the cluster should try to “live migrate” this resource when it needs to be moved (see <i>Migrating Resources</i> )   |
| allow-unhealthy-nodes      | FALSE  | Whether the resource should be able to run on a node even if the node’s health score would otherwise prevent it (see <i>Tracking Node Health</i> ) (since 2.1.3)   |
| container-attribute-target |  | Specific to bundle resources; see <i>Bundle Node Attributes</i>  |
| remote-node                |  | The name of the Pacemaker Remote guest node this resource is associated with, if any. If specified, this both enables the resource as a guest node and defines the unique name used to identify the guest node. The guest must be configured to run the Pacemaker Remote daemon when it is started.<br><b>WARNING:</b> This value cannot overlap with any resource or node IDs.  |
| remote-port                | 3121   | If <b>remote-node</b> is specified, the port on the guest used for its Pacemaker Remote connection. The Pacemaker Remote daemon on the guest must be configured to listen on this port.  |
| remote-addr                | value of <b>remote-node</b>                              | If <b>remote-node</b> is specified, the IP address or hostname used to connect to the guest via Pacemaker Remote. The Pacemaker Remote daemon on the guest must be configured to accept connections on this address.   |

Continued on next page

Table 7 – continued from previous page

| Field                  | Default | Description   |
|------------------------|---------|---|
| remote-connect-timeout | 60s     | If <code>remote-node</code> is specified, how long before a pending guest connection will time out. |

As an example of setting resource options, if you performed the following commands on an LSB Email resource:

```
# crm_resource --meta --resource Email --set-parameter priority --parameter-value 100
# crm_resource -m -r Email -p multiple-active -v block
```

the resulting resource definition might be:

#### An LSB resource with cluster options

```
<primitive id="Email" class="lsb" type="exim">
  <meta_attributes id="Email-meta_attributes">
    <nvpair id="Email-meta_attributes-priority" name="priority" value="100"/>
    <nvpair id="Email-meta_attributes-multiple-active" name="multiple-active" value="block"/>
  </meta_attributes>
</primitive>
```

In addition to the cluster-defined meta-attributes described above, you may also configure arbitrary meta-attributes of your own choosing. Most commonly, this would be done for use in *rules*. For example, an IT department might define a custom meta-attribute to indicate which company department each resource is intended for. To reduce the chance of name collisions with cluster-defined meta-attributes added in the future, it is recommended to use a unique, organization-specific prefix for such attributes.

### Setting Global Defaults for Resource Meta-Attributes

To set a default value for a resource option, add it to the `rscl_defaults` section with `crm_attribute`. For example,

```
# crm_attribute --type rscl_defaults --name is-managed --update false
```

would prevent the cluster from starting or stopping any of the resources in the configuration (unless of course the individual resources were specifically enabled by having their `is-managed` set to `true`).

### Resource Instance Attributes

The resource agents of some resource classes (`lsb`, `systemd` and `upstart` *not* among them) can be given parameters which determine how they behave and which instance of a service they control.

If your resource agent supports parameters, you can add them with the `crm_resource` command. For example,

```
# crm_resource --resource Public-IP --set-parameter ip --parameter-value 192.0.2.2
```

would create an entry in the resource like this:

#### An example OCF resource with instance attributes

```
<primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

For an OCF resource, the result would be an environment variable called `OCF_RESKEY_ip` with a value of `192.0.2.2`.

The list of instance attributes supported by an OCF resource agent can be found by calling the resource agent with the `meta-data` command. The output contains an XML description of all the supported attributes, their purpose and default values.

### Displaying the metadata for the Dummy resource agent template

```
# export OCF_ROOT=/usr/lib/ocf
# $OCF_ROOT/resource.d/pacemaker/Dummy meta-data
```

```

<?xml version="1.0"?>
<!DOCTYPE resource-agent SYSTEM "ra-api-1.dtd">
<resource-agent name="Dummy" version="2.0">
<version>1.1</version>

<longdesc lang="en">
This is a dummy OCF resource agent. It does absolutely nothing except keep track
of whether it is running or not, and can be configured so that actions fail or
take a long time. Its purpose is primarily for testing, and to serve as a
template for resource agent writers.
</longdesc>
<shortdesc lang="en">Example stateless resource agent</shortdesc>

<parameters>
<parameter name="state" unique-group="state">
<longdesc lang="en">
Location to store the resource state in.
</longdesc>
<shortdesc lang="en">State file</shortdesc>
<content type="string" default="/var/run/Dummy-RESOURCE_ID.state" />
</parameter>

<parameter name="passwd" reloadable="1">
<longdesc lang="en">
Fake password field
</longdesc>
<shortdesc lang="en">Password</shortdesc>
<content type="string" default="" />
</parameter>

<parameter name="fake" reloadable="1">
<longdesc lang="en">
Fake attribute that can be changed to cause a reload
</longdesc>
<shortdesc lang="en">Fake attribute that can be changed to cause a reload</shortdesc>
<content type="string" default="dummy" />
</parameter>

<parameter name="op_sleep" reloadable="1">
<longdesc lang="en">
Number of seconds to sleep during operations. This can be used to test how
the cluster reacts to operation timeouts.
</longdesc>
<shortdesc lang="en">Operation sleep duration in seconds.</shortdesc>
<content type="string" default="0" />
</parameter>

<parameter name="fail_start_on" reloadable="1">
<longdesc lang="en">
Start, migrate_from, and reload-agent actions will return failure if running on
the host specified here, but the resource will run successfully anyway (future
monitor calls will find it running). This can be used to test on-fail=ignore.
</longdesc>
<shortdesc lang="en">Report bogus start failure on specified host</shortdesc>
<content type="string" default="" />
</parameter>
<parameter name="envfile" reloadable="1">
<longdesc lang="en">

```

```

If this is set, the environment will be dumped to this file for every call.

```

```

</longdesc>
<shortdesc lang="en">Environment dump file</shortdesc>
<content type="string" default="" />
</parameter>

```

```

</parameters>

```

## 2.4.5 Resource Operations

*Operations* are actions the cluster can perform on a resource by calling the resource agent. Resource agents must support certain common operations such as start, stop, and monitor, and may implement any others.

Operations may be explicitly configured for two purposes: to override defaults for options (such as timeout) that the cluster will use whenever it initiates the operation, and to run an operation on a recurring basis (for example, to monitor the resource for failure).

### An OCF resource with a non-default start timeout

```
<primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
  <operations>
    <op id="Public-IP-start" name="start" timeout="60s"/>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

Pacemaker identifies operations by a combination of name and interval, so this combination must be unique for each resource. That is, you should not configure two operations for the same resource with the same name and interval.

### Operation Properties

Operation properties may be specified directly in the `op` element as XML attributes, or in a separate `meta_attributes` block as `nvpair` elements. XML attributes take precedence over `nvpair` elements if both are specified.

Table 8: Properties of an Operation

| Field    | Default | Description  |
|----------|---------|--|
| id       |         | A unique name for the operation.   |
| name     |         | The action to perform. This can be any action supported by the agent; common values include <code>monitor</code> , <code>start</code> , and <code>stop</code> .                            |
| interval | 0       | How frequently (in seconds) to perform the operation. A value of 0 means “when needed”. A positive value defines a <i>recurring action</i> , which is typically used with <i>monitor</i> . |
| timeout  |         | How long to wait before declaring the action has failed  |

Continued on next page



Table 8 – continued from previous page

| Field          | Default  | Description  |
|----------------|--|--|
| on-fail        | Varies by action: <ul style="list-style-type: none"> <li>• <b>stop:</b> <code>fence</code> if <code>stonith-enabled</code> is true or <code>block</code> otherwise</li> <li>• <b>demote:</b> <code>on-fail</code> of the <code>monitor</code> action with <code>role</code> set to <code>Promoted</code>, if present, enabled, and configured to a value other than <code>demote</code>, or <code>restart</code> otherwise</li> <li>• all other actions: <code>restart</code></li> </ul> | The action to take if this action ever fails. Allowed values: <ul style="list-style-type: none"> <li>• <b>ignore:</b> Pretend the resource did not fail.</li> <li>• <b>block:</b> Don't perform any further operations on the resource.</li> <li>• <b>stop:</b> Stop the resource and do not start it elsewhere.</li> <li>• <b>demote:</b> Demote the resource, without a full restart. This is valid only for <code>promote</code> actions, and for <code>monitor</code> actions with both a nonzero <code>interval</code> and <code>role</code> set to <code>Promoted</code>; for any other action, a configuration error will be logged, and the default behavior will be used. (<i>since 2.0.5</i>)</li> <li>• <b>restart:</b> Stop the resource and start it again (possibly on a different node).</li> <li>• <b>fence:</b> STONITH the node on which the resource failed.</li> <li>• <b>standby:</b> Move <i>all</i> resources away from the node on which the resource failed.</li> </ul> |
| enabled        | TRUE   | If <b>false</b> , ignore this operation definition. This is typically used to pause a particular recurring <code>monitor</code> operation; for instance, it can complement the respective resource being unmanaged ( <code>is-managed=false</code> ), as this alone will <i>not block any configured monitoring</i> . Disabling the operation does not suppress all actions of the given type. Allowed values: <b>true</b> , <b>false</b> .  |
| record-pending | TRUE   | If <b>true</b> , the intention to perform the operation is recorded so that GUIs and CLI tools can indicate that an operation is in progress. This is best set as an <i>operation default</i> (see <i>Setting Global Defaults for Operations</i> ). Allowed values: <b>true</b> , <b>false</b> .   |
| role           |  | Run the operation only on node(s) that the cluster thinks should be in the specified role. This only makes sense for recurring <code>monitor</code> operations. Allowed (case-sensitive) values: <b>Stopped</b> , <b>Started</b> , and in the case of <i>promotable clone resources</i> , <b>Unpromoted</b> and <b>Promoted</b> .  |

**Note:** When `on-fail` is set to `demote`, recovery from failure by a successful demote causes the cluster to recalculate whether and where a new instance should be promoted. The node with the failure is eligible, so if promotion scores have not changed, it will be promoted again.

There is no direct equivalent of `migration-threshold` for the promoted role, but the same effect can be achieved with a location constraint using a *rule* with a node attribute expression for the resource's fail count.

For example, to immediately ban the promoted role from a node with any failed promote or promoted instance monitor:

```
<rsc_location id="loc1" rsc="my_primitive">
  <rule id="rule1" score="-INFINITY" role="Promoted" boolean-op="or">
    <expression id="expr1" attribute="fail-count-my_primitive#promote_0"
      operation="gte" value="1"/>
    <expression id="expr2" attribute="fail-count-my_primitive#monitor_10000"
      operation="gte" value="1"/>
  </rule>
</rsc_location>
```

This example assumes that there is a promotable clone of the `my_primitive` resource (note that the primitive name, not the clone name, is used in the rule), and that there is a recurring 10-second-interval monitor configured for the promoted role (fail count attributes specify the interval in milliseconds).

---

## Monitoring Resources for Failure

When Pacemaker first starts a resource, it runs one-time `monitor` operations (referred to as *probes*) to ensure the resource is running where it's supposed to be, and not running where it's not supposed to be. (This behavior can be affected by the `resource-discovery` location constraint property.)

Other than those initial probes, Pacemaker will *not* (by default) check that the resource continues to stay healthy<sup>6</sup>. You must configure `monitor` operations explicitly to perform these checks.

### An OCF resource with a recurring health check

```
<primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
  <operations>
    <op id="Public-IP-start" name="start" timeout="60s"/>
    <op id="Public-IP-monitor" name="monitor" interval="60s"/>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

By default, a `monitor` operation will ensure that the resource is running where it is supposed to. The `target-role` property can be used for further checking.

For example, if a resource has one `monitor` operation with `interval=10` `role=Started` and a second `monitor` operation with `interval=11` `role=Stopped`, the cluster will run the first monitor on any nodes it thinks *should* be running the resource, and the second monitor on any nodes that it thinks *should not* be running the resource (for the truly paranoid, who want to know when an administrator manually starts a service by mistake).

---

**Note:** Currently, monitors with `role=Stopped` are not implemented for *clone* resources.

---

## Monitoring Resources When Administration is Disabled

Recurring `monitor` operations behave differently under various administrative settings:

---

<sup>6</sup> Currently, anyway. Automatic monitoring operations may be added in a future version of Pacemaker.

- When a resource is unmanaged (by setting `is-managed=false`): No monitors will be stopped.

If the unmanaged resource is stopped on a node where the cluster thinks it should be running, the cluster will detect and report that it is not, but it will not consider the monitor failed, and will not try to start the resource until it is managed again.

Starting the unmanaged resource on a different node is strongly discouraged and will at least cause the cluster to consider the resource failed, and may require the resource's `target-role` to be set to **Stopped** then **Started** to be recovered.

- When a node is put into standby: All resources will be moved away from the node, and all **monitor** operations will be stopped on the node, except those specifying `role` as **Stopped** (which will be newly initiated if appropriate).
- When the cluster is put into maintenance mode: All resources will be marked as unmanaged. All monitor operations will be stopped, except those specifying `role` as **Stopped** (which will be newly initiated if appropriate). As with single unmanaged resources, starting a resource on a node other than where the cluster expects it to be will cause problems.

## Setting Global Defaults for Operations

You can change the global default values for operation properties in a given cluster. These are defined in an `op_defaults` section of the CIB's `configuration` section, and can be set with `crm_attribute`. For example,

```
# crm_attribute --type op_defaults --name timeout --update 20s
```

would default each operation's `timeout` to 20 seconds. If an operation's definition also includes a value for `timeout`, then that value would be used for that operation instead.

## When Implicit Operations Take a Long Time

The cluster will always perform a number of implicit operations: **start**, **stop** and a non-recurring **monitor** operation used at startup to check whether the resource is already active. If one of these is taking too long, then you can create an entry for them and specify a longer timeout.

### An OCF resource with custom timeouts for its implicit actions

```
<primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
  <operations>
    <op id="public-ip-startup" name="monitor" interval="0" timeout="90s"/>
    <op id="public-ip-start" name="start" interval="0" timeout="180s"/>
    <op id="public-ip-stop" name="stop" interval="0" timeout="15min"/>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

## Multiple Monitor Operations

Provided no two operations (for a single resource) have the same name and interval, you can have as many **monitor** operations as you like. In this way, you can do a superficial health check every minute and

progressively more intense ones at higher intervals.

To tell the resource agent what kind of check to perform, you need to provide each monitor with a different value for a common parameter. The OCF standard creates a special parameter called `OCF_CHECK_LEVEL` for this purpose and dictates that it is “made available to the resource agent without the normal `OCF_RESKEY` prefix”.

Whatever name you choose, you can specify it by adding an `instance_attributes` block to the `op` tag. It is up to each resource agent to look for the parameter and decide how to use it.

**An OCF resource with two recurring health checks, performing different levels of checks specified via `OCF_CHECK_LEVEL`.**

```
<primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
  <operations>
    <op id="public-ip-health-60" name="monitor" interval="60">
      <instance_attributes id="params-public-ip-depth-60">
        <nvpair id="public-ip-depth-60" name="OCF_CHECK_LEVEL" value="10"/>
      </instance_attributes>
    </op>
    <op id="public-ip-health-300" name="monitor" interval="300">
      <instance_attributes id="params-public-ip-depth-300">
        <nvpair id="public-ip-depth-300" name="OCF_CHECK_LEVEL" value="20"/>
      </instance_attributes>
    </op>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-level" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

### Disabling a Monitor Operation

The easiest way to stop a recurring monitor is to just delete it. However, there can be times when you only want to disable it temporarily. In such cases, simply add `enabled=false` to the operation’s definition.

**Example of an OCF resource with a disabled health check**

```
<primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
  <operations>
    <op id="public-ip-check" name="monitor" interval="60s" enabled="false"/>
  </operations>
  <instance_attributes id="params-public-ip">
    <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
  </instance_attributes>
</primitive>
```

This can be achieved from the command line by executing:

```
# cibadmin --modify --xml-text '<op id="public-ip-check" enabled="false"/>'
```

Once you’ve done whatever you needed to do, you can then re-enable it with

```
# cibadmin --modify --xml-text '<op id="public-ip-check" enabled="true"/>'
```

## 2.5 Resource Constraints

### 2.5.1 Scores

Scores of all kinds are integral to how the cluster works. Practically everything from moving a resource to deciding which resource to stop in a degraded cluster is achieved by manipulating scores in some way.

Scores are calculated per resource and node. Any node with a negative score for a resource can't run that resource. The cluster places a resource on the node with the highest score for it.

#### Infinity Math

Pacemaker implements **INFINITY** (or equivalently, **+INFINITY**) internally as a score of 1,000,000. Addition and subtraction with it follow these three basic rules:

- Any value + **INFINITY** = **INFINITY**
- Any value - **INFINITY** = **-INFINITY**
- **INFINITY** - **INFINITY** = **-INFINITY**

---

**Note:** What if you want to use a score higher than 1,000,000? Typically this possibility arises when someone wants to base the score on some external metric that might go above 1,000,000.

The short answer is you can't.

The long answer is it is sometimes possible work around this limitation creatively. You may be able to set the score to some computed value based on the external metric rather than use the metric directly. For nodes, you can store the metric as a node attribute, and query the attribute when computing the score (possibly as part of a custom resource agent).

---

### 2.5.2 Deciding Which Nodes a Resource Can Run On

*Location constraints* tell the cluster which nodes a resource can run on.

There are two alternative strategies. One way is to say that, by default, resources can run anywhere, and then the location constraints specify nodes that are not allowed (an *opt-out* cluster). The other way is to start with nothing able to run anywhere, and use location constraints to selectively enable allowed nodes (an *opt-in* cluster).

Whether you should choose opt-in or opt-out depends on your personal preference and the make-up of your cluster. If most of your resources can run on most of the nodes, then an opt-out arrangement is likely to result in a simpler configuration. On the other-hand, if most resources can only run on a small subset of nodes, an opt-in configuration might be simpler.

#### Location Properties

Table 9: Attributes of a `rsc_location` Element

| Attribute                       | Default             | Description  |
|---------------------------------|---------------------|--|
| <code>id</code>                 |                     | A unique name for the constraint (required)  |
| <code>rsc</code>                |                     | The name of the resource to which this constraint applies. A location constraint must either have a <code>rsc</code> , have a <code>rsc-pattern</code> , or contain at least one resource set.   |
| <code>rsc-pattern</code>        |                     | A pattern matching the names of resources to which this constraint applies. The syntax is the same as <a href="#">POSIX</a> extended regular expressions, with the addition of an initial <code>!</code> indicating that resources <i>not</i> matching the pattern are selected. If the regular expression contains submatches, and the constraint is governed by a <i>rule</i> , the submatches can be referenced as <code>%1</code> through <code>%9</code> in the rule's <code>score-attribute</code> or a rule expression's <code>attribute</code> . A location constraint must either have a <code>rsc</code> , have a <code>rsc-pattern</code> , or contain at least one resource set.   |
| <code>node</code>               |                     | The name of the node to which this constraint applies. A location constraint must either have a <code>node</code> and <code>score</code> , or contain at least one rule.   |
| <code>score</code>              |                     | Positive values indicate a preference for running the affected resource(s) on <code>node</code> – the higher the value, the stronger the preference. Negative values indicate the resource(s) should avoid this node (a value of <b>-INFINITY</b> changes “should” to “must”). A location constraint must either have a <code>node</code> and <code>score</code> , or contain at least one rule.   |
| <code>resource-discovery</code> | <code>always</code> | <p>Whether Pacemaker should perform resource discovery (that is, check whether the resource is already running) for this resource on this node. This should normally be left as the default, so that rogue instances of a service can be stopped when they are running where they are not supposed to be. However, there are two situations where disabling resource discovery is a good idea: when a service is not installed on a node, discovery might return an error (properly written OCF agents will not, so this is usually only seen with other agent types); and when Pacemaker Remote is used to scale a cluster to hundreds of nodes, limiting resource discovery to allowed nodes can significantly boost performance.</p> <ul style="list-style-type: none"> <li>• <b>always</b>: Always perform resource discovery for the specified resource on this node.</li> <li>• <b>never</b>: Never perform resource discovery for the specified resource on this node. This option should generally be used with a <b>-INFINITY</b> score, although that is not strictly required.</li> <li>• <b>exclusive</b>: Perform resource discovery for the specified resource only on this node (and other nodes similarly marked as <b>exclusive</b>). Multiple location constraints using <b>exclusive</b> discovery for the same resource across different nodes creates a subset of nodes resource-discovery is exclusive to. If a resource is marked for <b>exclusive</b> discovery on one or more nodes, that resource is only allowed to be placed within that subset of nodes.</li> </ul> |

**Warning:** Setting `resource-discovery` to `never` or `exclusive` removes Pacemaker's ability to detect and stop unwanted instances of a service running where it's not supposed to be. It is up to the system administrator (you!) to make sure that the service can *never* be active on nodes without

resource-discovery (such as by leaving the relevant software uninstalled).

### Asymmetrical “Opt-In” Clusters

To create an opt-in cluster, start by preventing resources from running anywhere by default:

```
# crm_attribute --name symmetric-cluster --update false
```

Then start enabling nodes. The following fragment says that the web server prefers **sles-1**, the database prefers **sles-2** and both can fail over to **sles-3** if their most preferred node fails.

#### Opt-in location constraints for two resources

```
<constraints>
  <rsc_location id="loc-1" rsc="Webserver" node="sles-1" score="200"/>
  <rsc_location id="loc-2" rsc="Webserver" node="sles-3" score="0"/>
  <rsc_location id="loc-3" rsc="Database" node="sles-2" score="200"/>
  <rsc_location id="loc-4" rsc="Database" node="sles-3" score="0"/>
</constraints>
```

### Symmetrical “Opt-Out” Clusters

To create an opt-out cluster, start by allowing resources to run anywhere by default:

```
# crm_attribute --name symmetric-cluster --update true
```

Then start disabling nodes. The following fragment is the equivalent of the above opt-in configuration.

#### Opt-out location constraints for two resources

```
<constraints>
  <rsc_location id="loc-1" rsc="Webserver" node="sles-1" score="200"/>
  <rsc_location id="loc-2-do-not-run" rsc="Webserver" node="sles-2" score="-INFINITY"/>
  <rsc_location id="loc-3-do-not-run" rsc="Database" node="sles-1" score="-INFINITY"/>
  <rsc_location id="loc-4" rsc="Database" node="sles-2" score="200"/>
</constraints>
```

### What if Two Nodes Have the Same Score

If two nodes have the same score, then the cluster will choose one. This choice may seem random and may not be what was intended, however the cluster was not given enough information to know any better.

#### Constraints where a resource prefers two nodes equally

```
<constraints>
  <rsc_location id="loc-1" rsc="Webserver" node="sles-1" score="INFINITY"/>
  <rsc_location id="loc-2" rsc="Webserver" node="sles-2" score="INFINITY"/>
  <rsc_location id="loc-3" rsc="Database" node="sles-1" score="500"/>
  <rsc_location id="loc-4" rsc="Database" node="sles-2" score="300"/>
  <rsc_location id="loc-5" rsc="Database" node="sles-2" score="200"/>
</constraints>
```

In the example above, assuming no other constraints and an inactive cluster, **Webserver** would probably be placed on **sles-1** and **Database** on **sles-2**. It would likely have placed **Webserver** based on the node's uname and **Database** based on the desire to spread the resource load evenly across the cluster. However other factors can also be involved in more complex configurations.

### 2.5.3 Specifying the Order in which Resources Should Start/Stop

*Ordering constraints* tell the cluster the order in which certain resource actions should occur.

**Important:** Ordering constraints affect *only* the ordering of resource actions; they do *not* require that the resources be placed on the same node. If you want resources to be started on the same node *and* in a specific order, you need both an ordering constraint *and* a colocation constraint (see *Placing Resources Relative to other Resources*), or alternatively, a group (see *Groups - A Syntactic Shortcut*).

#### Ordering Properties

Table 10: Attributes of a `rsc_order` Element

| Field        | Default                      | Description  |
|--------------|------------------------------|--|
| id           |                              | A unique name for the constraint   |
| first        |                              | Name of the resource that the <b>then</b> resource depends on  |
| then         |                              | Name of the dependent resource   |
| first-action | start                        | The action that the <b>first</b> resource must complete before <b>then-action</b> can be initiated for the <b>then</b> resource. Allowed values: <b>start</b> , <b>stop</b> , <b>promote</b> , <b>demote</b> .   |
| then-action  | value of <b>first-action</b> | The action that the <b>then</b> resource can execute only after the <b>first-action</b> on the <b>first</b> resource has completed. Allowed values: <b>start</b> , <b>stop</b> , <b>promote</b> , <b>demote</b> .  |
| kind         | Mandatory                    | How to enforce the constraint. Allowed values: <ul style="list-style-type: none"> <li>• <b>Mandatory:</b> <b>then-action</b> will never be initiated for the <b>then</b> resource unless and until <b>first-action</b> successfully completes for the <b>first</b> resource.</li> <li>• <b>Optional:</b> The constraint applies only if both specified resource actions are scheduled in the same transition (that is, in response to the same cluster state). This means that <b>then-action</b> is allowed on the <b>then</b> resource regardless of the state of the <b>first</b> resource, but if both actions happen to be scheduled at the same time, they will be ordered.</li> <li>• <b>Serialize:</b> Ensure that the specified actions are never performed concurrently for the specified resources. <b>First-action</b> and <b>then-action</b> can be executed in either order, but one must complete before the other can be initiated. An example use case is when resource start-up puts a high load on the host.</li> </ul> |

Continued on next page



Table 10 – continued from previous page

| Field       | Default   | Description   |
|-------------|---|---|
| symmetrical | TRUE for <b>Mandatory</b> and <b>Optional</b> kinds. FALSE for <b>Serialize</b> kind. | If true, the reverse of the constraint applies for the opposite action (for example, if B starts after A starts, then B stops before A stops). <b>Serialize</b> orders cannot be symmetrical. |

Promote and demote apply to *promotable* clone resources.

### Optional and mandatory ordering

Here is an example of ordering constraints where **Database** *must* start before **Webserver**, and **IP** *should* start before **Webserver** if they both need to be started:

#### Optional and mandatory ordering constraints

```
<constraints>
  <rsc_order id="order-1" first="IP" then="Webserver" kind="Optional"/>
  <rsc_order id="order-2" first="Database" then="Webserver" kind="Mandatory" />
</constraints>
```

Because the above example lets `symmetrical` default to TRUE, **Webserver** must be stopped before **Database** can be stopped, and **Webserver** should be stopped before **IP** if they both need to be stopped.

## 2.5.4 Placing Resources Relative to other Resources

*Colocation constraints* tell the cluster that the location of one resource depends on the location of another one.

Colocation has an important side-effect: it affects the order in which resources are assigned to a node. Think about it: You can't place A relative to B unless you know where B is<sup>1</sup>.

So when you are creating colocation constraints, it is important to consider whether you should colocate A with B, or B with A.

**Important:** Colocation constraints affect *only* the placement of resources; they do *not* require that the resources be started in a particular order. If you want resources to be started on the same node *and* in a specific order, you need both an ordering constraint (see *Specifying the Order in which Resources Should Start/Stop*) and a colocation constraint, or alternatively, a group (see *Groups - A Syntactic Shortcut*).

### Colocation Properties

Table 11: Attributes of a `rsc_colocation` Constraint

| Field | Default | Description                                  |
|-------|---------|--|
| id    |         | A unique name for the constraint (required). |

Continued on next page

<sup>1</sup> While the human brain is sophisticated enough to read the constraint in any order and choose the correct one depending on the situation, the cluster is not quite so smart. Yet.

Table 11 – continued from previous page

| Field          | Default   | Description  |
|----------------|---|--|
| rsc            |   | The name of a resource that should be located relative to <b>with-rsc</b> . A colocation constraint must either contain at least one <i>resource set</i> , or specify both <b>rsc</b> and <b>with-rsc</b> .  |
| with-rsc       |   | The name of the resource used as the colocation target. The cluster will decide where to put this resource first and then decide where to put <b>rsc</b> . A colocation constraint must either contain at least one <i>resource set</i> , or specify both <b>rsc</b> and <b>with-rsc</b> .               |
| node-attribute | #uname  | If <b>rsc</b> and <b>with-rsc</b> are specified, this node attribute must be the same on the node running <b>rsc</b> and the node running <b>with-rsc</b> for the constraint to be satisfied. (For details, see <i>Colocation by Node Attribute</i> .)   |
| score          | 0   | Positive values indicate the resources should run on the same node. Negative values indicate the resources should run on different nodes. Values of +/- INFINITY change “should” to “must”.  |
| rsc-role       | Started   | If <b>rsc</b> and <b>with-rsc</b> are specified, and <b>rsc</b> is a <i>promotable clone</i> , the constraint applies only to <b>rsc</b> instances in this role. Allowed values: <b>Started</b> , <b>Promoted</b> , <b>Unpromoted</b> . For details, see <i>Promotable Clone Constraints</i> .           |
| with-rsc-role  | Started   | If <b>rsc</b> and <b>with-rsc</b> are specified, and <b>with-rsc</b> is a <i>promotable clone</i> , the constraint applies only to <b>with-rsc</b> instances in this role. Allowed values: <b>Started</b> , <b>Promoted</b> , <b>Unpromoted</b> . For details, see <i>Promotable Clone Constraints</i> . |
| influence      | value of critical meta-attribute for <b>rsc</b> | Whether to consider the location preferences of <b>rsc</b> when <b>with-rsc</b> is already active. Allowed values: <b>true</b> , <b>false</b> . For details, see <i>Colocation Influence</i> . (since 2.1.0)   |

## Mandatory Placement

Mandatory placement occurs when the constraint’s score is **+INFINITY** or **-INFINITY**. In such cases, if the constraint can’t be satisfied, then the **rsc** resource is not permitted to run. For **score=INFINITY**, this includes cases where the **with-rsc** resource is not active.

If you need resource **A** to always run on the same machine as resource **B**, you would add the following constraint:

### Mandatory colocation constraint for two resources

```
<rsc_colocation id="colocate" rsc="A" with-rsc="B" score="INFINITY"/>
```

Remember, because **INFINITY** was used, if **B** can’t run on any of the cluster nodes (for whatever reason) then **A** will not be allowed to run. Whether **A** is running or not has no effect on **B**.

Alternatively, you may want the opposite – that **A** *cannot* run on the same machine as **B**. In this case, use **score="-INFINITY"**.

**Mandatory anti-colocation constraint for two resources**

```
<rsc_colocation id="anti-colocate" rsc="A" with-rsc="B" score="-INFINITY"/>
```

Again, by specifying **-INFINITY**, the constraint is binding. So if the only place left to run is where **B** already is, then **A** may not run anywhere.

As with **INFINITY**, **B** can run even if **A** is stopped. However, in this case **A** also can run if **B** is stopped, because it still meets the constraint of **A** and **B** not running on the same node.

**Advisory Placement**

If mandatory placement is about “must” and “must not”, then advisory placement is the “I’d prefer if” alternative. For constraints with scores greater than **-INFINITY** and less than **INFINITY**, the cluster will try to accommodate your wishes but may ignore them if the alternative is to stop some of the cluster resources.

As in life, where if enough people prefer something it effectively becomes mandatory, advisory colocation constraints can combine with other elements of the configuration to behave as if they were mandatory.

**Advisory colocation constraint for two resources**

```
<rsc_colocation id="colocate-maybe" rsc="A" with-rsc="B" score="500"/>
```

**Colocation by Node Attribute**

The `node-attribute` property of a colocation constraints allows you to express the requirement, “these resources must be on similar nodes”.

As an example, imagine that you have two Storage Area Networks (SANs) that are not controlled by the cluster, and each node is connected to one or the other. You may have two resources **r1** and **r2** such that **r2** needs to use the same SAN as **r1**, but doesn’t necessarily have to be on the same exact node. In such a case, you could define a *node attribute* named **san**, with the value **san1** or **san2** on each node as appropriate. Then, you could colocate **r2** with **r1** using `node-attribute` set to **san**.

**Colocation Influence**

By default, if A is colocated with B, the cluster will take into account A’s preferences when deciding where to place B, to maximize the chance that both resources can run.

For a detailed look at exactly how this occurs, see [Colocation Explained](#).

However, if `influence` is set to **false** in the colocation constraint, this will happen only if B is inactive and needing to be started. If B is already active, A’s preferences will have no effect on placing B.

An example of what effect this would have and when it would be desirable would be a nonessential reporting tool colocated with a resource-intensive service that takes a long time to start. If the reporting tool fails enough times to reach its migration threshold, by default the cluster will want to move both resources to another node if possible. Setting `influence` to **false** on the colocation constraint would mean that the reporting tool would be stopped in this situation instead, to avoid forcing the service to move.

The **critical** resource meta-attribute is a convenient way to specify the default for all colocation constraints and groups involving a particular resource.

**Note:** If a noncritical resource is a member of a group, all later members of the group will be treated as noncritical, even if they are marked as (or left to default to) critical.

## 2.5.5 Resource Sets

*Resource sets* allow multiple resources to be affected by a single constraint.

### A set of 3 resources

```
<resource_set id="resource-set-example">
  <resource_ref id="A"/>
  <resource_ref id="B"/>
  <resource_ref id="C"/>
</resource_set>
```

Resource sets are valid inside `rsc_location`, `rsc_order` (see *Ordering Sets of Resources*), `rsc_colocation` (see *Colocating Sets of Resources*), and `rsc_ticket` (see *Configuring Ticket Dependencies*) constraints.

A resource set has a number of properties that can be set, though not all have an effect in all contexts.

Table 12: Attributes of a `resource_set` Element

| Field                    | Default   | Description   |
|--------------------------|---|---|
| <code>id</code>          |   | A unique name for the set (required)  |
| <code>sequential</code>  | true  | Whether the members of the set must be acted on in order. Meaningful within <code>rsc_order</code> and <code>rsc_colocation</code> .  |
| <code>require-all</code> | true  | Whether all members of the set must be active before continuing. With the current implementation, the cluster may continue even if only one member of the set is started, but if more than one member of the set is starting at the same time, the cluster will still wait until all of those have started before continuing (this may change in future versions). Meaningful within <code>rsc_order</code> . |
| <code>role</code>        |   | The constraint applies only to resource set members that are <i>Promotable clones</i> in this role. Meaningful within <code>rsc_location</code> , <code>rsc_colocation</code> and <code>rsc_ticket</code> . Allowed values: <code>Started</code> , <code>Promoted</code> , <code>Unpromoted</code> . For details, see <i>Promotable Clone Constraints</i> .   |
| <code>action</code>      | value of <code>first-action</code> in the enclosing ordering constraint | The action that applies to <i>all members</i> of the set. Meaningful within <code>rsc_order</code> . Allowed values: <code>start</code> , <code>stop</code> , <code>promote</code> , <code>demote</code> .  |
| <code>score</code>       |   | <i>Advanced use only.</i> Use a specific score for this set within the constraint.  |

## 2.5.6 Ordering Sets of Resources

A common situation is for an administrator to create a chain of ordered resources, such as:

### A chain of ordered resources

```
<constraints>
  <rsc_order id="order-1" first="A" then="B" />
  <rsc_order id="order-2" first="B" then="C" />
  <rsc_order id="order-3" first="C" then="D" />
</constraints>
```

### Visual representation of the four resources' start order for the above constraints



## Ordered Set

To simplify this situation, *Resource Sets* can be used within ordering constraints:

### A chain of ordered resources expressed as a set

```
<constraints>
  <rsc_order id="order-1">
    <resource_set id="ordered-set-example" sequential="true">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_order>
</constraints>
```

While the set-based format is not less verbose, it is significantly easier to get right and maintain.

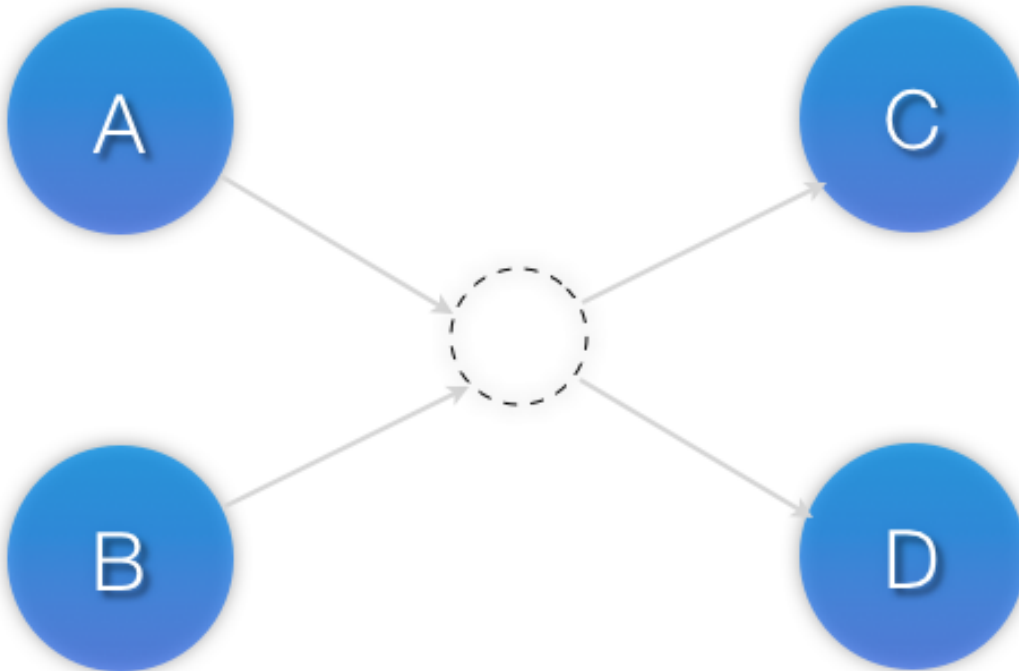
**Important:** If you use a higher-level tool, pay attention to how it exposes this functionality. Depending on the tool, creating a set **A B** may be equivalent to **A then B**, or **B then A**.

## Ordering Multiple Sets

The syntax can be expanded to allow sets of resources to be ordered relative to each other, where the members of each individual set may be ordered or unordered (controlled by the `sequential` property). In the example below, **A** and **B** can both start in parallel, as can **C** and **D**, however **C** and **D** can only start once *both* **A** and **B** are active.

**Ordered sets of unordered resources**

```
<constraints>
  <rsc_order id="order-1">
    <resource_set id="ordered-set-1" sequential="false">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="ordered-set-2" sequential="false">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_order>
</constraints>
```

**Visual representation of the start order for two ordered sets of unordered resources**

Of course either set – or both sets – of resources can also be internally ordered (by setting `sequential="true"`) and there is no limit to the number of sets that can be specified.

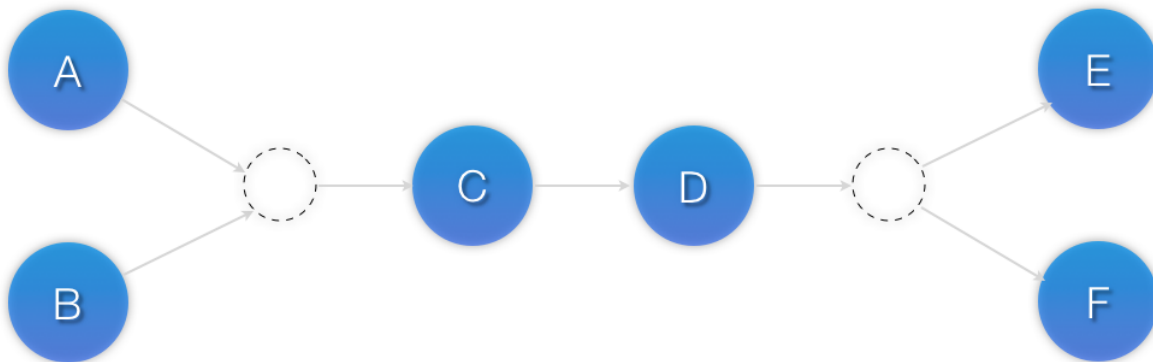
**Advanced use of set ordering - Three ordered sets, two of which are internally unordered**

```

<constraints>
  <rsc_order id="order-1">
    <resource_set id="ordered-set-1" sequential="false">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="ordered-set-2" sequential="true">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
    <resource_set id="ordered-set-3" sequential="false">
      <resource_ref id="E"/>
      <resource_ref id="F"/>
    </resource_set>
  </rsc_order>
</constraints>

```

Visual representation of the start order for the three sets defined above



**Important:** An ordered set with `sequential=false` makes sense only if there is another set in the constraint. Otherwise, the constraint has no effect.

### Resource Set OR Logic

The unordered set logic discussed so far has all been “AND” logic. To illustrate this take the 3 resource set figure in the previous section. Those sets can be expressed, **(A and B) then (C) then (D) then (E and F)**.

Say for example we want to change the first set, **(A and B)**, to use “OR” logic so the sets look like this: **(A or B) then (C) then (D) then (E and F)**. This functionality can be achieved through the use of the `require-all` option. This option defaults to `TRUE` which is why the “AND” logic is used by default. Setting `require-all=false` means only one resource in the set needs to be started before continuing on to the next set.

Resource Set “OR” logic: Three ordered sets, where the first set is internally unordered with “OR” logic

```
<constraints>
  <rsc_order id="order-1">
    <resource_set id="ordered-set-1" sequential="false" require-all="false">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="ordered-set-2" sequential="true">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
    <resource_set id="ordered-set-3" sequential="false">
      <resource_ref id="E"/>
      <resource_ref id="F"/>
    </resource_set>
  </rsc_order>
</constraints>
```

---

**Important:** An ordered set with `require-all=false` makes sense only in conjunction with `sequential=false`. Think of it like this: `sequential=false` modifies the set to be an unordered set using “AND” logic by default, and adding `require-all=false` flips the unordered set’s “AND” logic to “OR” logic.

---

## 2.5.7 Colocating Sets of Resources

Another common situation is for an administrator to create a set of colocated resources.

The simplest way to do this is to define a resource group (see *Groups - A Syntactic Shortcut*), but that cannot always accurately express the desired relationships. For example, maybe the resources do not need to be ordered.

Another way would be to define each relationship as an individual constraint, but that causes a difficult-to-follow constraint explosion as the number of resources and combinations grow.

Colocation chain as individual constraints, where A is placed first, then B, then C, then D

```
<constraints>
  <rsc_colocation id="coloc-1" rsc="D" with-rsc="C" score="INFINITY"/>
  <rsc_colocation id="coloc-2" rsc="C" with-rsc="B" score="INFINITY"/>
  <rsc_colocation id="coloc-3" rsc="B" with-rsc="A" score="INFINITY"/>
</constraints>
```

To express complicated relationships with a simplified syntax<sup>2</sup>, *resource sets* can be used within colocation constraints.

---

<sup>2</sup> which is not the same as saying easy to follow



### Equivalent colocation chain expressed using resource\_set

```
<constraints>
  <rsc_colocation id="coloc-1" score="INFINITY" >
    <resource_set id="colocated-set-example" sequential="true">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_colocation>
</constraints>
```

**Note:** Within a `resource_set`, the resources are listed in the order they are *placed*, which is the reverse of the order in which they are *colocated*. In the above example, resource **A** is placed before resource **B**, which is the same as saying resource **B** is colocated with resource **A**.

As with individual constraints, a resource that can't be active prevents any resource that must be colocated with it from being active. In both of the two previous examples, if **B** is unable to run, then both **C** and by inference **D** must remain stopped.

**Important:** If you use a higher-level tool, pay attention to how it exposes this functionality. Depending on the tool, creating a set **A B** may be equivalent to **A with B**, or **B with A**.

Resource sets can also be used to tell the cluster that entire *sets* of resources must be colocated relative to each other, while the individual members within any one set may or may not be colocated relative to each other (determined by the set's `sequential` property).

In the following example, resources **B**, **C**, and **D** will each be colocated with **A** (which will be placed first). **A** must be able to run in order for any of the resources to run, but any of **B**, **C**, or **D** may be stopped without affecting any of the others.

### Using colocated sets to specify a shared dependency

```
<constraints>
  <rsc_colocation id="coloc-1" score="INFINITY" >
    <resource_set id="colocated-set-2" sequential="false">
      <resource_ref id="B"/>
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
    <resource_set id="colocated-set-1" sequential="true">
      <resource_ref id="A"/>
    </resource_set>
  </rsc_colocation>
</constraints>
```

**Note:** Pay close attention to the order in which resources and sets are listed. While the members of any one sequential set are placed first to last (i.e., the colocation dependency is last with first), multiple sets are

placed last to first (i.e. the colocation dependency is first with last).

---

**Important:** A colocated set with `sequential="false"` makes sense only if there is another set in the constraint. Otherwise, the constraint has no effect.

---

There is no inherent limit to the number and size of the sets used. The only thing that matters is that in order for any member of one set in the constraint to be active, all members of sets listed after it must also be active (and naturally on the same node); and if a set has `sequential="true"`, then in order for one member of that set to be active, all members listed before it must also be active.

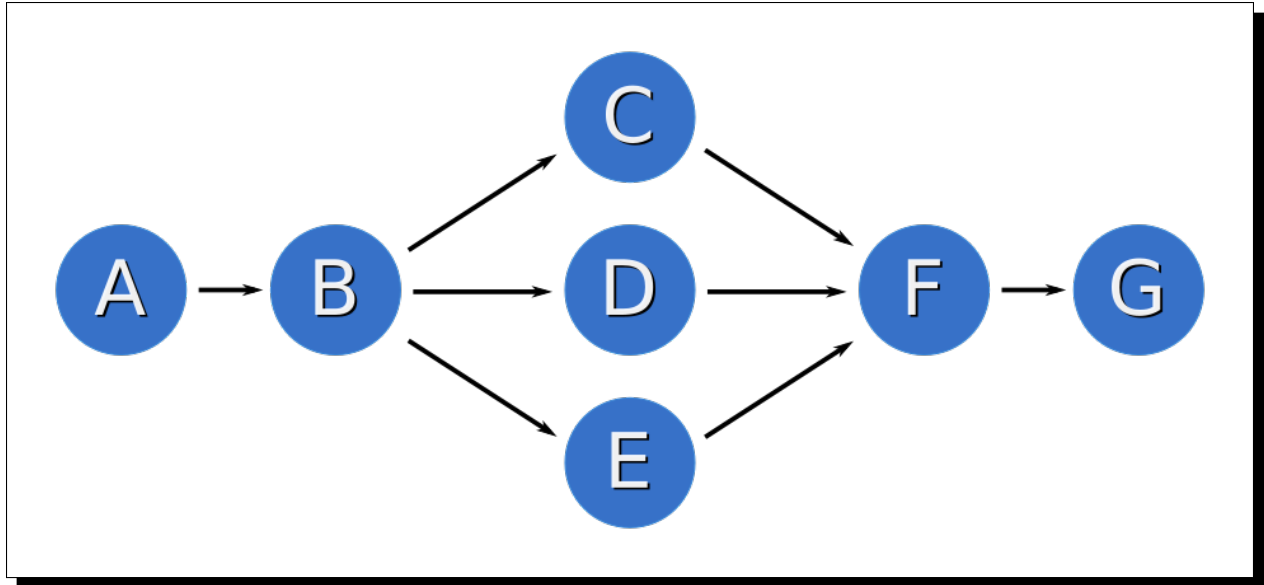
If desired, you can restrict the dependency to instances of promotable clone resources that are in a specific role, using the set's `role` property.

**Colocation in which the members of the middle set have no interdependencies, and the last set listed applies only to promoted instances**

```
<constraints>
  <rsc_colocation id="coloc-1" score="INFINITY" >
    <resource_set id="colocated-set-1" sequential="true">
      <resource_ref id="F"/>
      <resource_ref id="G"/>
    </resource_set>
    <resource_set id="colocated-set-2" sequential="false">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
      <resource_ref id="E"/>
    </resource_set>
    <resource_set id="colocated-set-3" sequential="true" role="Promoted">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
  </rsc_colocation>
</constraints>
```

**Visual representation of the above example (resources are placed from left to right)**

---



**Note:** Unlike ordered sets, colocated sets do not use the `require-all` option.

## 2.6 Fencing

### 2.6.1 What Is Fencing?

*Fencing* is the ability to make a node unable to run resources, even when that node is unresponsive to cluster commands.

Fencing is also known as *STONITH*, an acronym for “Shoot The Other Node In The Head”, since the most common fencing method is cutting power to the node. Another method is “fabric fencing”, cutting the node’s access to some capability required to run resources (such as network access or a shared disk).

### 2.6.2 Why Is Fencing Necessary?

Fencing protects your data from being corrupted by malfunctioning nodes or unintentional concurrent access to shared resources.

Fencing protects against the “split brain” failure scenario, where cluster nodes have lost the ability to reliably communicate with each other but are still able to run resources. If the cluster just assumed that uncommunicative nodes were down, then multiple instances of a resource could be started on different nodes.

The effect of split brain depends on the resource type. For example, an IP address brought up on two hosts on a network will cause packets to randomly be sent to one or the other host, rendering the IP useless. For a database or clustered file system, the effect could be much more severe, causing data corruption or divergence.

Fencing is also used when a resource cannot otherwise be stopped. If a resource fails to stop on a node, it cannot be started on a different node without risking the same type of conflict as split-brain. Fencing the original node ensures the resource can be safely started elsewhere.

Users may also configure the `on-fail` property of *Resource Operations* or the `loss-policy` property of *ticket constraints* to `fence`, in which case the cluster will fence the resource's node if the operation fails or the ticket is lost.

### 2.6.3 Fence Devices

A *fence device* or *fencing device* is a special type of resource that provides the means to fence a node.

Examples of fencing devices include intelligent power switches and IPMI devices that accept SNMP commands to cut power to a node, and iSCSI controllers that allow SCSI reservations to be used to cut a node's access to a shared disk.

Since fencing devices will be used to recover from loss of networking connectivity to other nodes, it is essential that they do not rely on the same network as the cluster itself, otherwise that network becomes a single point of failure.

Since loss of a node due to power outage is indistinguishable from loss of network connectivity to that node, it is also essential that at least one fence device for a node does not share power with that node. For example, an on-board IPMI controller that shares power with its host should not be used as the sole fencing device for that host.

Since fencing is used to isolate malfunctioning nodes, no fence device should rely on its target functioning properly. This includes, for example, devices that ssh into a node and issue a shutdown command (such devices might be suitable for testing, but never for production).

### 2.6.4 Fence Agents

A *fence agent* or *fencing agent* is a `stonith`-class resource agent.

The fence agent standard provides commands (such as `off` and `reboot`) that the cluster can use to fence nodes. As with other resource agent classes, this allows a layer of abstraction so that Pacemaker doesn't need any knowledge about specific fencing technologies – that knowledge is isolated in the agent.

Pacemaker supports two fence agent standards, both inherited from no-longer-active projects:

- Red Hat Cluster Suite (RHCS) style: These are typically installed in `/usr/sbin` with names starting with `fence_`.
- Linux-HA style: These typically have names starting with `external/`. Pacemaker can support these agents using the `fence_legacy` RHCS-style agent as a wrapper, *if* support was enabled when Pacemaker was built, which requires the `cluster-glue` library.

### 2.6.5 When a Fence Device Can Be Used

Fencing devices do not actually “run” like most services. Typically, they just provide an interface for sending commands to an external device.

Additionally, fencing may be initiated by Pacemaker, by other cluster-aware software such as DRBD or DLM, or manually by an administrator, at any point in the cluster life cycle, including before any resources have been started.

To accommodate this, Pacemaker does not require the fence device resource to be “started” in order to be used. Whether a fence device is started or not determines whether a node runs any recurring monitor for the device, and gives the node a slight preference for being chosen to execute fencing using that device.

By default, any node can execute any fencing device. If a fence device is disabled by setting its `target-role` to `Stopped`, then no node can use that device. If a location constraint with a negative score prevents a

specific node from “running” a fence device, then that node will never be chosen to execute fencing using the device. A node may fence itself, but the cluster will choose that only if no other nodes can do the fencing.

A common configuration scenario is to have one fence device per target node. In such a case, users often configure anti-location constraints so that the target node does not monitor its own device.

## 2.6.6 Limitations of Fencing Resources

Fencing resources have certain limitations that other resource classes don’t:

- They may have only one set of meta-attributes and one set of instance attributes.
- If *Rules* are used to determine fencing resource options, these might be evaluated only when first read, meaning that later changes to the rules will have no effect. Therefore, it is better to avoid confusion and not use rules at all with fencing resources.

These limitations could be revisited if there is sufficient user demand.

## 2.6.7 Special Meta-Attributes for Fencing Resources

The table below lists special resource meta-attributes that may be set for any fencing resource.

Table 13: Additional Properties of Fencing Resources

| Field    | Type   | Default | Description  |
|----------|--------|---------|--|
| provides | string |         | Any special capability provided by the fence device. Currently, only one such capability is meaningful: <i>unfencing</i> . |

## 2.6.8 Special Instance Attributes for Fencing Resources

The table below lists special instance attributes that may be set for any fencing resource (*not* meta-attributes, even though they are interpreted by Pacemaker rather than the fence agent). These are also listed in the man page for `pacemaker-fenced`.

Table 14: Additional Properties of Fencing Resources

| Field           | Type   | Default | Description  |
|-----------------|--------|---------|--|
| stonith-timeout | time   |         | This is not used by Pacemaker (see the <code>pcmk_reboot_timeout</code> , <code>pcmk_off_timeout</code> , etc. properties instead), but it may be used by Linux-HA fence agents.   |
| pcmk_host_map   | string |         | A mapping of node names to ports for devices that do not understand the node names. Example: <code>node1:1;node2:2,3</code> tells the cluster to use port 1 for <code>node1</code> and ports 2 and 3 for <code>node2</code> . If <code>pcmk_host_check</code> is explicitly set to <code>static-list</code> , either this or <code>pcmk_host_list</code> must be set. The port portion of the map may contain special characters such as spaces if preceded by a backslash ( <i>since 2.1.2</i> ). |

Continued on next page

Table 14 – continued from previous page

| Field                          | Type    | Default  | Description  |
|--------------------------------|---------|--|--|
| <code>pcmk_host_list</code>    | string  |  | A list of machines controlled by this device. If <code>pcmk_host_check</code> is explicitly set to <code>static-list</code> , either this or <code>pcmk_host_map</code> must be set.   |
| <code>pcmk_host_check</code>   | string  | Value appropriate to other parameters (see “Default Check Type” below) | The method Pacemaker should use to determine which nodes can be targeted by this device. Allowed values: <ul style="list-style-type: none"> <li>• <b>static-list</b>: targets are listed in the <code>pcmk_host_list</code> or <code>pcmk_host_map</code> attribute</li> <li>• <b>dynamic-list</b>: query the device via the agent’s <code>list</code> action</li> <li>• <b>status</b>: query the device via the agent’s <code>status</code> action</li> <li>• <b>none</b>: assume the device can fence any node</li> </ul>  |
| <code>pcmk_delay_max</code>    | time    | 0s   | Enable a delay of no more than the time specified before executing fencing actions. Pacemaker derives the overall delay by taking the value of <code>pcmk_delay_base</code> and adding a random delay value such that the sum is kept below this maximum. This is sometimes used in two-node clusters to ensure that the nodes don’t fence each other at the same time.  |
| <code>pcmk_delay_base</code>   | time    | 0s   | Enable a static delay before executing fencing actions. This can be used, for example, in two-node clusters to ensure that the nodes don’t fence each other, by having separate fencing resources with different values. The node that is fenced with the shorter delay will lose a fencing race. The overall delay introduced by pacemaker is derived from this value plus a random delay such that the sum is kept below the maximum delay. A single device can have different delays per node using a host map ( <i>since 2.1.2</i> ), for example <code>node1:0s;node2:5s</code> . |
| <code>pcmk_action_limit</code> | integer | 1  | The maximum number of actions that can be performed in parallel on this device. A value of -1 means unlimited. Node fencing actions initiated by the cluster (as opposed to an administrator running the <code>stonith_admin</code> tool or the fencer running recurring device monitors and <code>status</code> and <code>list</code> commands) are additionally subject to the <code>concurrent-fencing</code> cluster property.   |

Continued on next page

Table 14 – continued from previous page

| Field                            | Type    | Default   | Description   |
|----------------------------------|---------|---|---|
| <code>pcmk_host_argument</code>  | string  | <code>port</code> otherwise <code>plug</code> if supported according to the metadata of the fence agent | <i>Advanced use only.</i> Which parameter should be supplied to the fence agent to identify the node to be fenced. Some devices support neither the standard <code>plug</code> nor the deprecated <code>port</code> parameter, or may provide additional ones. Use this to specify an alternate, device-specific parameter. A value of <code>none</code> tells the cluster not to supply any additional parameters. |
| <code>pcmk_reboot_action</code>  | string  | <code>reboot</code>   | <i>Advanced use only.</i> The command to send to the resource agent in order to reboot a node. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.   |
| <code>pcmk_reboot_timeout</code> | time    | 60s   | <i>Advanced use only.</i> Specify an alternate timeout to use for <code>reboot</code> actions instead of the value of <code>stonith-timeout</code> . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.   |
| <code>pcmk_reboot_retries</code> | integer | 2   | <i>Advanced use only.</i> The maximum number of times to retry the <code>reboot</code> command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up.      |
| <code>pcmk_off_action</code>     | string  | <code>off</code>  | <i>Advanced use only.</i> The command to send to the resource agent in order to shut down a node. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.  |
| <code>pcmk_off_timeout</code>    | time    | 60s   | <i>Advanced use only.</i> Specify an alternate timeout to use for <code>off</code> actions instead of the value of <code>stonith-timeout</code> . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.  |
| <code>pcmk_off_retries</code>    | integer | 2   | <i>Advanced use only.</i> The maximum number of times to retry the <code>off</code> command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up.         |

Continued on next page

Table 14 – continued from previous page

| Field                             | Type    | Default | Description   |
|-----------------------------------|---------|---------|---|
| <code>pcmk_list_action</code>     | string  | list    | <i>Advanced use only.</i> The command to send to the resource agent in order to list nodes. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.  |
| <code>pcmk_list_timeout</code>    | time    | 60s     | <i>Advanced use only.</i> Specify an alternate timeout to use for <code>list</code> actions instead of the value of <code>stonith-timeout</code> . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.   |
| <code>pcmk_list_retries</code>    | integer | 2       | <i>Advanced use only.</i> The maximum number of times to retry the <code>list</code> command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up.    |
| <code>pcmk_monitor_action</code>  | string  | monitor | <i>Advanced use only.</i> The command to send to the resource agent in order to report extended status. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.  |
| <code>pcmk_monitor_timeout</code> | time    | 60s     | <i>Advanced use only.</i> Specify an alternate timeout to use for <code>monitor</code> actions instead of the value of <code>stonith-timeout</code> . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.  |
| <code>pcmk_monitor_retries</code> | integer | 2       | <i>Advanced use only.</i> The maximum number of times to retry the <code>monitor</code> command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up. |
| <code>pcmk_status_action</code>   | string  | status  | <i>Advanced use only.</i> The command to send to the resource agent in order to report status. Some devices do not support the standard commands or may provide additional ones. Use this to specify an alternate, device-specific command.   |

Continued on next page



Table 14 – continued from previous page

| Field                            | Type    | Default | Description  |
|----------------------------------|---------|---------|--|
| <code>pcmk_status_timeout</code> | time    | 60s     | <i>Advanced use only.</i> Specify an alternate timeout to use for <code>status</code> actions instead of the value of <code>stonith-timeout</code> . Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.  |
| <code>pcmk_status_retries</code> | integer | 2       | <i>Advanced use only.</i> The maximum number of times to retry the <code>status</code> command within the timeout period. Some devices do not support multiple connections, and operations may fail if the device is busy with another task, so Pacemaker will automatically retry the operation, if there is time remaining. Use this option to alter the number of times Pacemaker retries before giving up. |

### 2.6.9 Default Check Type

If the user does not explicitly configure `pcmk_host_check` for a fence device, a default value appropriate to other configured parameters will be used:

- If either `pcmk_host_list` or `pcmk_host_map` is configured, `static-list` will be used;
- otherwise, if the fence device supports the `list` action, and the first attempt at using `list` succeeds, `dynamic-list` will be used;
- otherwise, if the fence device supports the `status` action, `status` will be used;
- otherwise, `none` will be used.

### 2.6.10 Unfencing

With fabric fencing (such as cutting network or shared disk access rather than power), it is expected that the cluster will fence the node, and then a system administrator must manually investigate what went wrong, correct any issues found, then reboot (or restart the cluster services on) the node.

Once the node reboots and rejoins the cluster, some fabric fencing devices require an explicit command to restore the node's access. This capability is called *unfencing* and is typically implemented as the fence agent's `on` command.

If any cluster resource has `requires` set to `unfencing`, then that resource will not be probed or started on a node until that node has been unfenced.

### 2.6.11 Fencing and Quorum

In general, a cluster partition may execute fencing only if the partition has quorum, and the `stonith-enabled` cluster property is set to true. However, there are exceptions:

- The requirements apply only to fencing initiated by Pacemaker. If an administrator initiates fencing using the `stonith_admin` command, or an external application such as DLM initiates fencing using Pacemaker's C API, the requirements do not apply.

- A cluster partition without quorum is allowed to fence any active member of that partition. As a corollary, this allows a `no-quorum-policy` of `suicide` to work.
- If the `no-quorum-policy` cluster property is set to `ignore`, then quorum is not required to execute fencing of any node.

## 2.6.12 Fencing Timeouts

Fencing timeouts are complicated, since a single fencing operation can involve many steps, each of which may have a separate timeout.

Fencing may be initiated in one of several ways:

- An administrator may initiate fencing using the `stonith_admin` tool, which has a `--timeout` option (defaulting to 2 minutes) that will be used as the fence operation timeout.
- An external application such as DLM may initiate fencing using the Pacemaker C API. The application will specify the fence operation timeout in this case, which might or might not be configurable by the user.
- The cluster may initiate fencing itself. In this case, the `stonith-timeout` cluster property (defaulting to 1 minute) will be used as the fence operation timeout.

However fencing is initiated, the initiator contacts Pacemaker's fencer (`pacemaker-fenced`) to request fencing. This connection and request has its own timeout, separate from the fencing operation timeout, but usually happens very quickly.

The fencer will contact all fencers in the cluster to ask what devices they have available to fence the target node. The fence operation timeout will be used as the timeout for each of these queries.

Once a fencing device has been selected, the fencer will check whether any action-specific timeout has been configured for the device, to use instead of the fence operation timeout. For example, if `stonith-timeout` is 60 seconds, but the fencing device has `pcmk_reboot_timeout` configured as 90 seconds, then a timeout of 90 seconds will be used for reboot actions using that device.

A device may have retries configured, in which case the timeout applies across all attempts. For example, if a device has `pcmk_reboot_retries` configured as 2, and the first reboot attempt fails, the second attempt will only have whatever time is remaining in the action timeout after subtracting how much time the first attempt used. This means that if the first attempt fails due to using the entire timeout, no further attempts will be made. There is currently no way to configure a per-attempt timeout.

If more than one device is required to fence a target, whether due to failure of the first device or a fencing topology with multiple devices configured for the target, each device will have its own separate action timeout.

For all of the above timeouts, the fencer will generally multiply the configured value by 1.2 to get an actual value to use, to account for time needed by the fencer's own processing.

Separate from the fencer's timeouts, some fence agents have internal timeouts for individual steps of their fencing process. These agents often have parameters to configure these timeouts, such as `login-timeout`, `shell-timeout`, or `power-timeout`. Many such agents also have a `disable-timeout` parameter to ignore their internal timeouts and just let Pacemaker handle the timeout. This causes a difference in retry behavior. If `disable-timeout` is not set, and the agent hits one of its internal timeouts, it will report that as a failure to Pacemaker, which can then retry. If `disable-timeout` is set, and Pacemaker hits a timeout for the agent, then there will be no time remaining, and no retry will be done.

### 2.6.13 Fence Devices Dependent on Other Resources

In some cases, a fence device may require some other cluster resource (such as an IP address) to be active in order to function properly.

This is obviously undesirable in general: fencing may be required when the depended-on resource is not active, or fencing may be required because the node running the depended-on resource is no longer responding.

However, this may be acceptable under certain conditions:

- The dependent fence device should not be able to target any node that is allowed to run the depended-on resource.
- The depended-on resource should not be disabled during production operation.
- The `concurrent-fencing` cluster property should be set to `true`. Otherwise, if both the node running the depended-on resource and some node targeted by the dependent fence device need to be fenced, the fencing of the node running the depended-on resource might be ordered first, making the second fencing impossible and blocking further recovery. With concurrent fencing, the dependent fence device might fail at first due to the depended-on resource being unavailable, but it will be retried and eventually succeed once the resource is brought back up.

Even under those conditions, there is one unlikely problem scenario. The DC always schedules fencing of itself after any other fencing needed, to avoid unnecessary repeated DC elections. If the dependent fence device targets the DC, and both the DC and a different node running the depended-on resource need to be fenced, the DC fencing will always fail and block further recovery. Note, however, that losing a DC node entirely causes some other node to become DC and schedule the fencing, so this is only a risk when a stop or other operation with `on-fail` set to `fencing` fails on the DC.

### 2.6.14 Configuring Fencing

Higher-level tools can provide simpler interfaces to this process, but using Pacemaker command-line tools, this is how you could configure a fence device.

1. Find the correct driver:

```
# stonith_admin --list-installed
```

**Note:** You may have to install packages to make fence agents available on your host. Searching your available packages for `fence-` is usually helpful. Ensure the packages providing the fence agents you require are installed on every cluster node.

2. Find the required parameters associated with the device (replacing `$AGENT_NAME` with the name obtained from the previous step):

```
# stonith_admin --metadata --agent $AGENT_NAME
```

3. Create a file called `stonith.xml` containing a primitive resource with a class of `stonith`, a type equal to the agent name obtained earlier, and a parameter for each of the values returned in the previous step.
4. If the device does not know how to fence nodes based on their uname, you may also need to set the special `pcmk_host_map` parameter. See *Special Meta-Attributes for Fencing Resources* for details.
5. If the device does not support the `list` command, you may also need to set the special `pcmk_host_list` and/or `pcmk_host_check` parameters. See *Special Meta-Attributes for Fencing Resources* for details.

- If the device does not expect the victim to be specified with the `port` parameter, you may also need to set the special `pcmk_host_argument` parameter. See *Special Meta-Attributes for Fencing Resources* for details.

- Upload it into the CIB using `cibadmin`:

```
# cibadmin --create --scope resources --xml-file stonith.xml
```

- Set `stonith-enabled` to true:

```
# crm_attribute --type crm_config --name stonith-enabled --update true
```

- Once the stonith resource is running, you can test it by executing the following, replacing `$NODE_NAME` with the name of the node to fence (although you might want to stop the cluster on that machine first):

```
# stonith_admin --reboot $NODE_NAME
```

## Example Fencing Configuration

For this example, we assume we have a cluster node, `pcmk-1`, whose IPMI controller is reachable at the IP address 192.0.2.1. The IPMI controller uses the username `testuser` and the password `abc123`.

- Looking at what's installed, we may see a variety of available agents:

```
# stonith_admin --list-installed
```

```
(... some output omitted ...)
fence_idrac
fence_ilo3
fence_ilo4
fence_ilo5
fence_imm
fence_ipmilan
(... some output omitted ...)
```

Perhaps after some reading some man pages and doing some Internet searches, we might decide `fence_ipmilan` is our best choice.

- Next, we would check what parameters `fence_ipmilan` provides:

```
# stonith_admin --metadata -a fence_ipmilan
```

```
<resource-agent name="fence_ipmilan" shortdesc="Fence agent for IPMI">
  <symlink name="fence_ilo3" shortdesc="Fence agent for HP iLO3"/>
  <symlink name="fence_ilo4" shortdesc="Fence agent for HP iLO4"/>
  <symlink name="fence_ilo5" shortdesc="Fence agent for HP iLO5"/>
  <symlink name="fence_imm" shortdesc="Fence agent for IBM Integrated Management Module"/>
  <symlink name="fence_idrac" shortdesc="Fence agent for Dell iDRAC"/>
  <longdesc>fence_ipmilan is an I/O Fencing agent which can be used with machines controlled
  ↳ by IPMI. This agent calls support software ipmitool (http://ipmitool.sf.net/). WARNING! This
  ↳ fence agent might report success before the node is powered off. You should use -m/method
  ↳ onoff if your fence device works correctly with that option.</longdesc>
  <vendor-url/>
  <parameters>
    <parameter name="action" unique="0" required="0">
```

(continues on next page)

(continued from previous page)

```

    <getopt mixed="-o, --action=[action]"/>
    <content type="string" default="reboot"/>
    <shortdesc lang="en">Fencing action</shortdesc>
  </parameter>
  <parameter name="auth" unique="0" required="0">
    <getopt mixed="-A, --auth=[auth]"/>
    <content type="select">
      <option value="md5"/>
      <option value="password"/>
      <option value="none"/>
    </content>
    <shortdesc lang="en">IPMI Lan Auth type.</shortdesc>
  </parameter>
  <parameter name="cipher" unique="0" required="0">
    <getopt mixed="-C, --cipher=[cipher]"/>
    <content type="string"/>
    <shortdesc lang="en">Ciphersuite to use (same as ipmitool -C parameter)</shortdesc>
  </parameter>
  <parameter name="hexadecimal_kg" unique="0" required="0">
    <getopt mixed="--hexadecimal-kg=[key]"/>
    <content type="string"/>
    <shortdesc lang="en">Hexadecimal-encoded Kg key for IPMIv2 authentication</shortdesc>
  </parameter>
  <parameter name="ip" unique="0" required="0" obsoletes="ipaddr">
    <getopt mixed="-a, --ip=[ip]"/>
    <content type="string"/>
    <shortdesc lang="en">IP address or hostname of fencing device</shortdesc>
  </parameter>
  <parameter name="ipaddr" unique="0" required="0" deprecated="1">
    <getopt mixed="-a, --ip=[ip]"/>
    <content type="string"/>
    <shortdesc lang="en">IP address or hostname of fencing device</shortdesc>
  </parameter>
  <parameter name="ipport" unique="0" required="0">
    <getopt mixed="-u, --ipport=[port]"/>
    <content type="integer" default="623"/>
    <shortdesc lang="en">TCP/UDP port to use for connection with device</shortdesc>
  </parameter>
  <parameter name="lanplus" unique="0" required="0">
    <getopt mixed="-P, --lanplus"/>
    <content type="boolean" default="0"/>
    <shortdesc lang="en">Use Lanplus to improve security of connection</shortdesc>
  </parameter>
  <parameter name="login" unique="0" required="0" deprecated="1">
    <getopt mixed="-l, --username=[name]"/>
    <content type="string"/>
    <shortdesc lang="en">Login name</shortdesc>
  </parameter>
  <parameter name="method" unique="0" required="0">
    <getopt mixed="-m, --method=[method]"/>
    <content type="select" default="onoff">
      <option value="onoff"/>
      <option value="cycle"/>
    </content>
    <shortdesc lang="en">Method to fence</shortdesc>
  </parameter>

```

(continues on next page)

(continued from previous page)

```

<parameter name="passwd" unique="0" required="0" deprecated="1">
  <getopt mixed="-p, --password=[password]"/>
  <content type="string"/>
  <shortdesc lang="en">Login password or passphrase</shortdesc>
</parameter>
<parameter name="passwd_script" unique="0" required="0" deprecated="1">
  <getopt mixed="-S, --password-script=[script]"/>
  <content type="string"/>
  <shortdesc lang="en">Script to run to retrieve password</shortdesc>
</parameter>
<parameter name="password" unique="0" required="0" obsoletes="passwd">
  <getopt mixed="-p, --password=[password]"/>
  <content type="string"/>
  <shortdesc lang="en">Login password or passphrase</shortdesc>
</parameter>
<parameter name="password_script" unique="0" required="0" obsoletes="passwd_script">
  <getopt mixed="-S, --password-script=[script]"/>
  <content type="string"/>
  <shortdesc lang="en">Script to run to retrieve password</shortdesc>
</parameter>
<parameter name="plug" unique="0" required="0" obsoletes="port">
  <getopt mixed="-n, --plug=[ip]"/>
  <content type="string"/>
  <shortdesc lang="en">IP address or hostname of fencing device (together with --port-as-
↪ip)</shortdesc>
</parameter>
<parameter name="port" unique="0" required="0" deprecated="1">
  <getopt mixed="-n, --plug=[ip]"/>
  <content type="string"/>
  <shortdesc lang="en">IP address or hostname of fencing device (together with --port-as-
↪ip)</shortdesc>
</parameter>
<parameter name="privlvl" unique="0" required="0">
  <getopt mixed="-L, --privlvl=[level]"/>
  <content type="select" default="administrator">
    <option value="callback"/>
    <option value="user"/>
    <option value="operator"/>
    <option value="administrator"/>
  </content>
  <shortdesc lang="en">Privilege level on IPMI device</shortdesc>
</parameter>
<parameter name="target" unique="0" required="0">
  <getopt mixed="--target=[targetaddress]"/>
  <content type="string"/>
  <shortdesc lang="en">Bridge IPMI requests to the remote target address</shortdesc>
</parameter>
<parameter name="username" unique="0" required="0" obsoletes="login">
  <getopt mixed="-l, --username=[name]"/>
  <content type="string"/>
  <shortdesc lang="en">Login name</shortdesc>
</parameter>
<parameter name="quiet" unique="0" required="0">
  <getopt mixed="-q, --quiet"/>
  <content type="boolean"/>
  <shortdesc lang="en">Disable logging to stderr. Does not affect --verbose or --debug-
↪file or logging to syslog.</shortdesc>

```

(continues on next page)

(continued from previous page)

```

</parameter>
<parameter name="verbose" unique="0" required="0">
  <getopt mixed="-v, --verbose"/>
  <content type="boolean"/>
  <shortdesc lang="en">Verbose mode</shortdesc>
</parameter>
<parameter name="debug" unique="0" required="0" deprecated="1">
  <getopt mixed="-D, --debug-file=[debugfile]"/>
  <content type="string"/>
  <shortdesc lang="en">Write debug information to given file</shortdesc>
</parameter>
<parameter name="debug_file" unique="0" required="0" obsoletes="debug">
  <getopt mixed="-D, --debug-file=[debugfile]"/>
  <content type="string"/>
  <shortdesc lang="en">Write debug information to given file</shortdesc>
</parameter>
<parameter name="version" unique="0" required="0">
  <getopt mixed="-V, --version"/>
  <content type="boolean"/>
  <shortdesc lang="en">Display version information and exit</shortdesc>
</parameter>
<parameter name="help" unique="0" required="0">
  <getopt mixed="-h, --help"/>
  <content type="boolean"/>
  <shortdesc lang="en">Display help and exit</shortdesc>
</parameter>
<parameter name="delay" unique="0" required="0">
  <getopt mixed="--delay=[seconds]"/>
  <content type="second" default="0"/>
  <shortdesc lang="en">Wait X seconds before fencing is started</shortdesc>
</parameter>
<parameter name="ipmitool_path" unique="0" required="0">
  <getopt mixed="--ipmitool-path=[path]"/>
  <content type="string" default="/usr/bin/ipmitool"/>
  <shortdesc lang="en">Path to ipmitool binary</shortdesc>
</parameter>
<parameter name="login_timeout" unique="0" required="0">
  <getopt mixed="--login-timeout=[seconds]"/>
  <content type="second" default="5"/>
  <shortdesc lang="en">Wait X seconds for cmd prompt after login</shortdesc>
</parameter>
<parameter name="port_as_ip" unique="0" required="0">
  <getopt mixed="--port-as-ip"/>
  <content type="boolean"/>
  <shortdesc lang="en">Make "port/plug" to be an alias to IP address</shortdesc>
</parameter>
<parameter name="power_timeout" unique="0" required="0">
  <getopt mixed="--power-timeout=[seconds]"/>
  <content type="second" default="20"/>
  <shortdesc lang="en">Test X seconds for status change after ON/OFF</shortdesc>
</parameter>
<parameter name="power_wait" unique="0" required="0">
  <getopt mixed="--power-wait=[seconds]"/>
  <content type="second" default="2"/>
  <shortdesc lang="en">Wait X seconds after issuing ON/OFF</shortdesc>
</parameter>

```

(continues on next page)

(continued from previous page)

```

<parameter name="shell_timeout" unique="0" required="0">
  <getopt mixed="--shell-timeout=[seconds]" />
  <content type="second" default="3" />
  <shortdesc lang="en">Wait X seconds for cmd prompt after issuing command</shortdesc>
</parameter>
<parameter name="retry_on" unique="0" required="0">
  <getopt mixed="--retry-on=[attempts]" />
  <content type="integer" default="1" />
  <shortdesc lang="en">Count of attempts to retry power on</shortdesc>
</parameter>
<parameter name="sudo" unique="0" required="0" deprecated="1">
  <getopt mixed="--use-sudo" />
  <content type="boolean" />
  <shortdesc lang="en">Use sudo (without password) when calling 3rd party software</
↪shortdesc>
</parameter>
<parameter name="use_sudo" unique="0" required="0" obsoletes="sudo">
  <getopt mixed="--use-sudo" />
  <content type="boolean" />
  <shortdesc lang="en">Use sudo (without password) when calling 3rd party software</
↪shortdesc>
</parameter>
<parameter name="sudo_path" unique="0" required="0">
  <getopt mixed="--sudo-path=[path]" />
  <content type="string" default="/usr/bin/sudo" />
  <shortdesc lang="en">Path to sudo binary</shortdesc>
</parameter>
</parameters>
<actions>
  <action name="on" automatic="0" />
  <action name="off" />
  <action name="reboot" />
  <action name="status" />
  <action name="monitor" />
  <action name="metadata" />
  <action name="manpage" />
  <action name="validate-all" />
  <action name="diag" />
  <action name="stop" timeout="20s" />
  <action name="start" timeout="20s" />
</actions>
</resource-agent>

```

Once we've decided what parameter values we think we need, it is a good idea to run the fence agent's status action manually, to verify that our values work correctly:

```
# fence_ipmilan --lanplus -a 192.0.2.1 -l testuser -p abc123 -o status
```

```
Chassis Power is on
```

- Based on that, we might create a fencing resource configuration like this in `stonith.xml` (or any file name, just use the same name with `cibadmin` later):

```

<primitive id="Fencing-pcmk-1" class="stonith" type="fence_ipmilan" >
  <instance_attributes id="Fencing-params" >
    <nvpair id="Fencing-lanplus" name="lanplus" value="1" />

```

(continues on next page)



(continued from previous page)

```

<nvpair id="Fencing-ip" name="ip" value="192.0.2.1" />
<nvpair id="Fencing-password" name="password" value="testuser" />
<nvpair id="Fencing-username" name="username" value="abc123" />
</instance_attributes>
<operations >
  <op id="Fencing-monitor-10m" interval="10m" name="monitor" timeout="300s" />
</operations>
</primitive>

```

**Note:** Even though the man page shows that the `action` parameter is supported, we do not provide that in the resource configuration. Pacemaker will supply an appropriate action whenever the fence device must be used.

4. In this case, we don't need to configure `pcmk_host_map` because `fence_ipmilan` ignores the target node name and instead uses its `ip` parameter to know how to contact the IPMI controller.
5. We do need to let Pacemaker know which cluster node can be fenced by this device, since `fence_ipmilan` doesn't support the `list` action. Add a line like this to the agent's instance attributes:

```
<nvpair id="Fencing-pcmk_host_list" name="pcmk_host_list" value="pcmk-1" />
```

6. We don't need to configure `pcmk_host_argument` since `ip` is all the fence agent needs (it ignores the target name).
7. Make the configuration active:

```
# cibadmin --create --scope resources --xml-file stonith.xml
```

8. Set `stonith-enabled` to true (this only has to be done once):

```
# crm_attribute --type crm_config --name stonith-enabled --update true
```

9. Since our cluster is still in testing, we can reboot `pcmk-1` without bothering anyone, so we'll test our fencing configuration by running this from one of the other cluster nodes:

```
# stonith_admin --reboot pcmk-1
```

Then we will verify that the node did, in fact, reboot.

We can repeat that process to create a separate fencing resource for each node.

With some other fence device types, a single fencing resource is able to be used for all nodes. In fact, we could do that with `fence_ipmilan`, using the `port-as-ip` parameter along with `pcmk_host_map`. Either approach is fine.

## 2.6.15 Fencing Topologies

Pacemaker supports fencing nodes with multiple devices through a feature called *fencing topologies*. Fencing topologies may be used to provide alternative devices in case one fails, or to require multiple devices to all be executed successfully in order to consider the node successfully fenced, or even a combination of the two.

Create the individual devices as you normally would, then define one or more `fencing-level` entries in the `fencing-topology` section of the configuration.

- Each fencing level is attempted in order of ascending `index`. Allowed values are 1 through 9.

- If a device fails, processing terminates for the current level. No further devices in that level are exercised, and the next level is attempted instead.
- If the operation succeeds for all the listed devices in a level, the level is deemed to have passed.
- The operation is finished when a level has passed (success), or all levels have been attempted (failed).
- If the operation failed, the next step is determined by the scheduler and/or the controller.

Some possible uses of topologies include:

- Try on-board IPMI, then an intelligent power switch if that fails
- Try fabric fencing of both disk and network, then fall back to power fencing if either fails
- Wait up to a certain time for a kernel dump to complete, then cut power to the node

Table 15: Attributes of a fencing-level Element

| Attribute        | Description   |
|------------------|---|
| id               | A unique name for this element (required)   |
| target           | The name of a single node to which this level applies   |
| target-pattern   | An extended regular expression (as defined in <a href="#">POSIX</a> ) matching the names of nodes to which this level applies               |
| target-attribute | The name of a node attribute that is set (to <b>target-value</b> ) for nodes to which this level applies                                    |
| target-value     | The node attribute value (of <b>target-attribute</b> ) that is set for nodes to which this level applies                                    |
| index            | The order in which to attempt the levels. Levels are attempted in ascending order <i>until one succeeds</i> . Valid values are 1 through 9. |
| devices          | A comma-separated list of devices that must all be tried for this level   |

**Note:** Fencing topology with different devices for different nodes

```
<cib crm_feature_set="3.6.0" validate-with="pacemaker-3.5" admin_epoch="1" epoch="0" num_updates="0"
↪">
  <configuration>
    ...
    <fencing-topology>
      <!-- For pcmk-1, try poison-pill and fail back to power -->
      <fencing-level id="f-p1.1" target="pcmk-1" index="1" devices="poison-pill"/>
      <fencing-level id="f-p1.2" target="pcmk-1" index="2" devices="power"/>

      <!-- For pcmk-2, try disk and network, and fail back to power -->
      <fencing-level id="f-p2.1" target="pcmk-2" index="1" devices="disk,network"/>
      <fencing-level id="f-p2.2" target="pcmk-2" index="2" devices="power"/>
    </fencing-topology>
    ...
  </configuration>
  <status/>
</cib>
```

### Example Dual-Layer, Dual-Device Fencing Topologies

The following example illustrates an advanced use of **fencing-topology** in a cluster with the following properties:

- 2 nodes (prod-mysql1 and prod-mysql2)
- the nodes have IPMI controllers reachable at 192.0.2.1 and 192.0.2.2
- the nodes each have two independent Power Supply Units (PSUs) connected to two independent Power Distribution Units (PDUs) reachable at 198.51.100.1 (port 10 and port 11) and 203.0.113.1 (port 10 and port 11)
- fencing via the IPMI controller uses the `fence_ipmilan` agent (1 fence device per controller, with each device targeting a separate node)
- fencing via the PDUs uses the `fence_apc_snmp` agent (1 fence device per PDU, with both devices targeting both nodes)
- a random delay is used to lessen the chance of a “death match”
- fencing topology is set to try IPMI fencing first then dual PDU fencing if that fails

In a node failure scenario, Pacemaker will first select `fence_ipmilan` to try to kill the faulty node. Using the fencing topology, if that method fails, it will then move on to selecting `fence_apc_snmp` twice (once for the first PDU, then again for the second PDU).

The fence action is considered successful only if both PDUs report the required status. If any of them fails, fencing loops back to the first fencing method, `fence_ipmilan`, and so on, until the node is fenced or the fencing action is cancelled.

---

**Note: First fencing method: single IPMI device per target**

Each cluster node has its own dedicated IPMI controller that can be contacted for fencing using the following primitives:

```
<primitive class="stonith" id="fence_prod-mysql1_ipmi" type="fence_ipmilan">
  <instance_attributes id="fence_prod-mysql1_ipmi-instance_attributes">
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-ipaddr" name="ipaddr" value="192.0.2.1"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-login" name="login" value="fencing"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-passwd" name="passwd" value="finishme"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-lanplus" name="lanplus" value="true"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-pcmk_host_list" name="pcm_k_host_list"
    value="prod-mysql1"/>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-pcmk_delay_max" name="pcm_k_delay_max"
    value="8s"/>
  </instance_attributes>
</primitive>
<primitive class="stonith" id="fence_prod-mysql2_ipmi" type="fence_ipmilan">
  <instance_attributes id="fence_prod-mysql2_ipmi-instance_attributes">
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-ipaddr" name="ipaddr" value="192.0.2.2"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-login" name="login" value="fencing"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-passwd" name="passwd" value="finishme"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-lanplus" name="lanplus" value="true"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-pcmk_host_list" name="pcm_k_host_list"
    value="prod-mysql2"/>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-pcmk_delay_max" name="pcm_k_delay_max"
    value="8s"/>
  </instance_attributes>
</primitive>
```

**Note: Second fencing method: dual PDU devices**

Each cluster node also has 2 distinct power supplies controlled by 2 distinct PDUs:

- Node 1: PDU 1 port 10 and PDU 2 port 10
- Node 2: PDU 1 port 11 and PDU 2 port 11

The matching fencing agents are configured as follows:

```
<primitive class="stonith" id="fence_apc1" type="fence_apc_snmp">
  <instance_attributes id="fence_apc1-instance_attributes">
    <nvpair id="fence_apc1-instance_attributes-ipaddr" name="ipaddr" value="198.51.100.1"/>
    <nvpair id="fence_apc1-instance_attributes-login" name="login" value="fencing"/>
    <nvpair id="fence_apc1-instance_attributes-passwd" name="passwd" value="fencing"/>
    <nvpair id="fence_apc1-instance_attributes-pcmk_host_list"
      name="pcm_k_host_map" value="prod-mysql1:10;prod-mysql2:11"/>
    <nvpair id="fence_apc1-instance_attributes-pcmk_delay_max" name="pcm_k_delay_max" value="8s"/>
  </instance_attributes>
</primitive>
<primitive class="stonith" id="fence_apc2" type="fence_apc_snmp">
  <instance_attributes id="fence_apc2-instance_attributes">
    <nvpair id="fence_apc2-instance_attributes-ipaddr" name="ipaddr" value="203.0.113.1"/>
    <nvpair id="fence_apc2-instance_attributes-login" name="login" value="fencing"/>
    <nvpair id="fence_apc2-instance_attributes-passwd" name="passwd" value="fencing"/>
    <nvpair id="fence_apc2-instance_attributes-pcmk_host_list"
      name="pcm_k_host_map" value="prod-mysql1:10;prod-mysql2:11"/>
    <nvpair id="fence_apc2-instance_attributes-pcmk_delay_max" name="pcm_k_delay_max" value="8s"/>
  </instance_attributes>
</primitive>
```

---

**Note: Fencing topology**

Now that all the fencing resources are defined, it's time to create the right topology. We want to first fence using IPMI and if that does not work, fence both PDUs to effectively and surely kill the node.

```
<fencing-topology>
  <fencing-level id="level-1-1" target="prod-mysql1" index="1" devices="fence_prod-mysql1_ipmi" />
  <fencing-level id="level-1-2" target="prod-mysql1" index="2" devices="fence_apc1,fence_apc2" />
  <fencing-level id="level-2-1" target="prod-mysql2" index="1" devices="fence_prod-mysql2_ipmi" />
  <fencing-level id="level-2-2" target="prod-mysql2" index="2" devices="fence_apc1,fence_apc2" />
</fencing-topology>
```

In `fencing-topology`, the lowest `index` value for a target determines its first fencing method.

---

## 2.6.16 Remapping Reboots

When the cluster needs to reboot a node, whether because `stonith-action` is `reboot` or because a reboot was requested externally (such as by `stonith_admin --reboot`), it will remap that to other commands in two cases:

- If the chosen fencing device does not support the `reboot` command, the cluster will ask it to perform `off` instead.

- If a fencing topology level with multiple devices must be executed, the cluster will ask all the devices to perform **off**, then ask the devices to perform **on**.

To understand the second case, consider the example of a node with redundant power supplies connected to intelligent power switches. Rebooting one switch and then the other would have no effect on the node. Turning both switches off, and then on, actually reboots the node.

In such a case, the fencing operation will be treated as successful as long as the **off** commands succeed, because then it is safe for the cluster to recover any resources that were on the node. Timeouts and errors in the **on** phase will be logged but ignored.

When a reboot operation is remapped, any action-specific timeout for the remapped action will be used (for example, `pcmk_off_timeout` will be used when executing the **off** command, not `pcmk_reboot_timeout`).

## 2.7 Alerts

*Alerts* may be configured to take some external action when a cluster event occurs (node failure, resource starting or stopping, etc.).

### 2.7.1 Alert Agents

As with resource agents, the cluster calls an external program (an *alert agent*) to handle alerts. The cluster passes information about the event to the agent via environment variables. Agents can do anything desired with this information (send an e-mail, log to a file, update a monitoring system, etc.).

#### Simple alert configuration

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh" />
  </alerts>
</configuration>
```

In the example above, the cluster will call `my-script.sh` for each event.

Multiple alert agents may be configured; the cluster will call all of them for each event.

Alert agents will be called only on cluster nodes. They will be called for events involving Pacemaker Remote nodes, but they will never be called *on* those nodes.

### 2.7.2 Alert Recipients

Usually, alerts are directed towards a recipient. Thus, each alert may be additionally configured with one or more recipients. The cluster will call the agent separately for each recipient.

#### Alert configuration with recipient

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <recipient id="my-alert-recipient" value="some-address"/>
    </alert>
  </alerts>
</configuration>
```

In the above example, the cluster will call `my-script.sh` for each event, passing the recipient `some-address` as an environment variable.

The recipient may be anything the alert agent can recognize – an IP address, an e-mail address, a file name, whatever the particular agent supports.

### 2.7.3 Alert Meta-Attributes

As with resource agents, meta-attributes can be configured for alert agents to affect how Pacemaker calls them.

Table 16: Meta-Attributes of an Alert

| Meta-Attribute   | Default       | Description  |
|------------------|---------------|--|
| timestamp-format | %H:%M:%S.%06N | Format the cluster will use when sending the event's timestamp to the agent. This is a string as used with the <code>date(1)</code> command. |
| timeout          | 30s           | If the alert agent does not complete within this amount of time, it will be terminated.  |

Meta-attributes can be configured per alert agent and/or per recipient.

#### Alert configuration with meta-attributes

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <meta_attributes id="my-alert-attributes">
        <nvpair id="my-alert-attributes-timeout" name="timeout"
          value="15s"/>
      </meta_attributes>
      <recipient id="my-alert-recipient1" value="someuser@example.com">
        <meta_attributes id="my-alert-recipient1-attributes">
          <nvpair id="my-alert-recipient1-timestamp-format"
            name="timestamp-format" value="%D %H:%M"/>
        </meta_attributes>
      </recipient>
      <recipient id="my-alert-recipient2" value="otheruser@example.com">
        <meta_attributes id="my-alert-recipient2-attributes">
          <nvpair id="my-alert-recipient2-timestamp-format"
            name="timestamp-format" value="%c"/>
        </meta_attributes>
      </recipient>
    </alert>
  </alerts>
</configuration>
```

In the above example, the `my-script.sh` will get called twice for each event, with each call using a 15-second timeout. One call will be passed the recipient `someuser@example.com` and a timestamp in the format `%D %H:%M`, while the other call will be passed the recipient `otheruser@example.com` and a timestamp in the format `%c`.

## 2.7.4 Alert Instance Attributes

As with resource agents, agent-specific configuration values may be configured as instance attributes. These will be passed to the agent as additional environment variables. The number, names and allowed values of these instance attributes are completely up to the particular agent.

### Alert configuration with instance attributes

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <meta_attributes id="my-alert-attributes">
        <nvpair id="my-alert-attributes-timeout" name="timeout"
          value="15s"/>
      </meta_attributes>
      <instance_attributes id="my-alert-options">
        <nvpair id="my-alert-options-debug" name="debug"
          value="false"/>
      </instance_attributes>
      <recipient id="my-alert-recipient1"
        value="someuser@example.com"/>
    </alert>
  </alerts>
</configuration>
```

## 2.7.5 Alert Filters

By default, an alert agent will be called for node events, fencing events, and resource events. An agent may choose to ignore certain types of events, but there is still the overhead of calling it for those events. To eliminate that overhead, you may select which types of events the agent should receive.

### Alert configuration to receive only node events and fencing events

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <select>
        <select_nodes />
        <select_fencing />
      </select>
      <recipient id="my-alert-recipient1"
        value="someuser@example.com"/>
    </alert>
  </alerts>
</configuration>
```

The possible options within `<select>` are `<select_nodes>`, `<select_fencing>`, `<select_resources>`, and `<select_attributes>`.

With `<select_attributes>` (the only event type not enabled by default), the agent will receive alerts when a node attribute changes. If you wish the agent to be called only when certain attributes change, you can configure that as well.

## Alert configuration to be called when certain node attributes change

```
<configuration>
  <alerts>
    <alert id="my-alert" path="/path/to/my-script.sh">
      <select>
        <select_attributes>
          <attribute id="alert-standby" name="standby" />
          <attribute id="alert-shutdown" name="shutdown" />
        </select_attributes>
      </select>
      <recipient id="my-alert-recipient1" value="someuser@example.com"/>
    </alert>
  </alerts>
</configuration>
```

Node attribute alerts are currently considered experimental. Alerts may be limited to attributes set via `attrd_updater`, and agents may be called multiple times with the same attribute value.

## 2.7.6 Using the Sample Alert Agents

Pacemaker provides several sample alert agents, installed in `/usr/share/pacemaker/alerts` by default.

While these sample scripts may be copied and used as-is, they are provided mainly as templates to be edited to suit your purposes. See their source code for the full set of instance attributes they support.

## Sending cluster events as SNMP traps

```
<configuration>
  <alerts>
    <alert id="snmp_alert" path="/path/to/alert_snmp.sh">
      <instance_attributes id="config_for_alert_snmp">
        <nvpair id="trap_node_states" name="trap_node_states"
          value="all"/>
      </instance_attributes>
      <meta_attributes id="config_for_timestamp">
        <nvpair id="ts_fmt" name="timestamp-format"
          value="%Y-%m-%d,%H:%M:%S.%01N"/>
      </meta_attributes>
      <recipient id="snmp_destination" value="192.168.1.2"/>
    </alert>
  </alerts>
</configuration>
```

## Sending cluster events as e-mails

```
<configuration>
  <alerts>
    <alert id="smtp_alert" path="/path/to/alert_smtp.sh">
      <instance_attributes id="config_for_alert_smtp">
        <nvpair id="email_sender" name="email_sender"
          value="donotreply@example.com"/>
      </instance_attributes>
      <recipient id="smtp_destination" value="admin@example.com"/>
    </alert>
  </alerts>
</configuration>
```



## 2.7.7 Writing an Alert Agent

Table 17: **Environment variables passed to alert agents**

| Environment Variable      | Description   |
|---------------------------|---|
| CRM_alert_kind            | The type of alert ( <b>node</b> , <b>fencing</b> , <b>resource</b> , or <b>attribute</b> )  |
| CRM_alert_node            | Name of affected node   |
| CRM_alert_node_sequence   | A sequence number increased whenever an alert is being issued on the local node, which can be used to reference the order in which alerts have been issued by Pacemaker. An alert for an event that happened later in time reliably has a higher sequence number than alerts for earlier events. Be aware that this number has no cluster-wide meaning.           |
| CRM_alert_recipient       | The configured recipient  |
| CRM_alert_timestamp       | A timestamp created prior to executing the agent, in the format specified by the <b>timestamp-format</b> meta-attribute. This allows the agent to have a reliable, high-precision time of when the event occurred, regardless of when the agent itself was invoked (which could potentially be delayed due to system load, etc.).                                 |
| CRM_alert_timestamp_epoch | The same time as <b>CRM_alert_timestamp</b> , expressed as the integer number of seconds since January 1, 1970. This (along with <b>CRM_alert_timestamp_usec</b> ) can be useful for alert agents that need to format time in a specific way rather than let the user configure it.   |
| CRM_alert_timestamp_usec  | The same time as <b>CRM_alert_timestamp</b> , expressed as the integer number of microseconds since <b>CRM_alert_timestamp_epoch</b> .  |
| CRM_alert_version         | The version of Pacemaker sending the alert  |
| CRM_alert_desc            | Detail about event. For <b>node</b> alerts, this is the node's current state ( <b>member</b> or <b>lost</b> ). For <b>fencing</b> alerts, this is a summary of the requested fencing operation, including origin, target, and fencing operation error code, if any. For <b>resource</b> alerts, this is a readable string equivalent of <b>CRM_alert_status</b> . |
| CRM_alert_nodeid          | ID of node whose status changed (provided with <b>node</b> alerts only)   |
| CRM_alert_rc              | The numerical return code of the fencing or resource operation (provided with <b>fencing</b> and <b>resource</b> alerts only)   |
| CRM_alert_task            | The requested fencing or resource operation (provided with <b>fencing</b> and <b>resource</b> alerts only)  |
| CRM_alert_exec_time       | The (wall-clock) time, in milliseconds, that it took to execute the action. If the action timed out, <b>CRM_alert_status</b> will be 2, <b>CRM_alert_desc</b> will be "Timed Out", and this value will be the action timeout. May not be supported on all platforms. ( <b>resource</b> alerts only) ( <i>since 2.0.1</i> )  |
| CRM_alert_interval        | The interval of the resource operation ( <b>resource</b> alerts only)   |
| CRM_alert_rsc             | The name of the affected resource ( <b>resource</b> alerts only)  |
| CRM_alert_status          | A numerical code used by Pacemaker to represent the operation result ( <b>resource</b> alerts only)   |
| CRM_alert_target_rc       | The expected numerical return code of the operation ( <b>resource</b> alerts only)  |
| CRM_alert_attribute_name  | The name of the node attribute that changed ( <b>attribute</b> alerts only)   |
| CRM_alert_attribute_value | The new value of the node attribute that changed ( <b>attribute</b> alerts only)  |

Special concerns when writing alert agents:

- Alert agents may be called with no recipient (if none is configured), so the agent must be able to handle this situation, even if it only exits in that case. (Users may modify the configuration in stages, and

add a recipient later.)

- If more than one recipient is configured for an alert, the alert agent will be called once per recipient. If an agent is not able to run concurrently, it should be configured with only a single recipient. The agent is free, however, to interpret the recipient as a list.
- When a cluster event occurs, all alerts are fired off at the same time as separate processes. Depending on how many alerts and recipients are configured, and on what is done within the alert agents, a significant load burst may occur. The agent could be written to take this into consideration, for example by queueing resource-intensive actions into some other instance, instead of directly executing them.
- Alert agents are run as the `hacluster` user, which has a minimal set of permissions. If an agent requires additional privileges, it is recommended to configure `sudo` to allow the agent to run the necessary commands as another user with the appropriate privileges.
- As always, take care to validate and sanitize user-configured parameters, such as `CRM_alert_timestamp` (whose content is specified by the user-configured `timestamp-format`), `CRM_alert_recipient`, and all instance attributes. Mostly this is needed simply to protect against configuration errors, but if some user can modify the CIB without having `hacluster`-level access to the cluster nodes, it is a potential security concern as well, to avoid the possibility of code injection.

---

**Note: `ocf:pacemaker:ClusterMon` compatibility**

The alerts interface is designed to be backward compatible with the external scripts interface used by the `ocf:pacemaker:ClusterMon` resource, which is now deprecated. To preserve this compatibility, the environment variables passed to alert agents are available prepended with `CRM_notify_` as well as `CRM_alert_`. One break in compatibility is that `ClusterMon` ran external scripts as the `root` user, while alert agents are run as the `hacluster` user.

---

## 2.8 Rules

Rules can be used to make your configuration more dynamic, allowing values to change depending on the time or the value of a node attribute. Examples of things rules are useful for:

- Set a higher value for *resource-stickiness* during working hours, to minimize downtime, and a lower value on weekends, to allow resources to move to their most preferred locations when people aren't around to notice.
- Automatically place the cluster into maintenance mode during a scheduled maintenance window.
- Assign certain nodes and resources to a particular department via custom node attributes and meta-attributes, and add a single location constraint that restricts the department's resources to run only on those nodes.

Each constraint type or property set that supports rules may contain one or more **rule** elements specifying conditions under which the constraint or properties take effect. Examples later in this chapter will make this clearer.

## 2.8.1 Rule Properties

Table 18: **Attributes of a rule Element**

| Attribute       | Default        | Description  |
|-----------------|----------------|--|
| id              |                | A unique name for this element (required)  |
| role            | <b>Started</b> | The rule is in effect only when the resource is in the specified role. Allowed values are <b>Started</b> , <b>Unpromoted</b> , and <b>Promoted</b> . A rule with a <b>role</b> of <b>Promoted</b> cannot determine the initial location of a clone instance and will only affect which of the active instances will be promoted. |
| score           |                | If this rule is used in a location constraint and evaluates to true, apply this score to the constraint. Only one of <b>score</b> and <b>score-attribute</b> may be used.  |
| score-attribute |                | If this rule is used in a location constraint and evaluates to true, use the value of this node attribute as the score to apply to the constraint. Only one of <b>score</b> and <b>score-attribute</b> may be used.  |
| boolean-op      | <b>and</b>     | If this rule contains more than one condition, a value of <b>and</b> specifies that the rule evaluates to true only if all conditions are true, and a value of <b>or</b> specifies that the rule evaluates to true if any condition is true.   |

A **rule** element must contain one or more conditions. A condition may be an **expression** element, a **date\_expression** element, or another **rule** element.

## 2.8.2 Node Attribute Expressions

Expressions are rule conditions based on the values of node attributes.

Table 19: **Attributes of an expression Element**

| Attribute | Default  | Description   |
|-----------|--|---|
| id        |  | A unique name for this element (required)   |
| attribute |  | The node attribute to test (required)   |
| type      | The default type for <b>lt</b> , <b>gt</b> , <b>lte</b> , and <b>gte</b> operations is <b>number</b> if either value contains a decimal point character, or <b>integer</b> otherwise. The default type for all other operations is <b>string</b> . If a numeric parse fails for either value, then the values are compared as type <b>string</b> . | How the node attributes should be compared. Allowed values are <b>string</b> , <b>integer</b> ( <i>since 2.0.5</i> ), <b>number</b> , and <b>version</b> . <b>integer</b> truncates floating-point values if necessary before performing a 64-bit integer comparison. <b>number</b> performs a double-precision floating-point comparison ( <i>32-bit integer before 2.0.5</i> ). |

Continued on next page

Table 19 – continued from previous page

| Attribute    | Default        | Description   |
|--------------|----------------|---|
| operation    |                | <p>The comparison to perform (required). Allowed values:</p> <ul style="list-style-type: none"> <li>• <b>lt</b>: <b>True if the node attribute value</b> is less than the comparison value</li> <li>• <b>gt</b>: <b>True if the node attribute value</b> is greater than the comparison value</li> <li>• <b>lte</b>: <b>True if the node attribute value</b> is less than or equal to the comparison value</li> <li>• <b>gte</b>: <b>True if the node attribute value</b> is greater than or equal to the comparison value</li> <li>• <b>eq</b>: <b>True if the node attribute value</b> is equal to the comparison value</li> <li>• <b>ne</b>: <b>True if the node attribute value</b> is not equal to the comparison value</li> <li>• <b>defined</b>: <b>True if the node has the</b> named attribute</li> <li>• <b>not_defined</b>: <b>True if the node does not</b> have the named attribute</li> </ul> |
| value        |                | User-supplied value for comparison (required for operations other than <b>defined</b> and <b>not_defined</b> )  |
| value-source | <b>literal</b> | <p>How the <b>value</b> is derived. Allowed values:</p> <ul style="list-style-type: none"> <li>• <b>literal</b>: <b>value</b> is a literal string to compare against</li> <li>• <b>param</b>: <b>value</b> is the name of a resource parameter to compare against (only valid in location constraints)</li> <li>• <b>meta</b>: <b>value</b> is the name of a resource meta-attribute to compare against (only valid in location constraints)</li> </ul>   |

In addition to custom node attributes defined by the administrator, the cluster defines special, built-in node attributes for each node that can also be used in rule expressions.

Table 20: **Built-in Node Attributes**

| Name          | Value   |
|---------------|---|
| #uname        | <i>Node name</i>  |
| #id           | Node ID   |
| #kind         | Node type. Possible values are <b>cluster</b> , <b>remote</b> , and <b>container</b> . Kind is <b>remote</b> for Pacemaker Remote nodes created with the <code>ocf:pacemaker:remote</code> resource, and <b>container</b> for Pacemaker Remote guest nodes and bundle nodes |
| #is_dc        | <b>true</b> if this node is the cluster's Designated Controller (DC), <b>false</b> otherwise  |
| #cluster-name | The value of the <b>cluster-name</b> cluster property, if set   |
| #site-name    | The value of the <b>site-name</b> node attribute, if set, otherwise identical to <b>#cluster-name</b>   |
| #role         | The role the relevant promotable clone resource has on this node. Valid only within a rule for a location constraint for a promotable clone resource.   |

### 2.8.3 Date/Time Expressions

Date/time expressions are rule conditions based (as the name suggests) on the current date and time.

A `date_expression` element may optionally contain a `date_spec` or `duration` element depending on the context.

Table 21: Attributes of a `date_expression` Element

| Attribute              | Description  |
|------------------------|--|
| <code>id</code>        | A unique name for this element (required)  |
| <code>start</code>     | A date/time conforming to the <a href="#">ISO8601</a> specification. May be used when <code>operation</code> is <code>in_range</code> (in which case at least one of <code>start</code> or <code>end</code> must be specified) or <code>gt</code> (in which case <code>start</code> is required).  |
| <code>end</code>       | A date/time conforming to the <a href="#">ISO8601</a> specification. May be used when <code>operation</code> is <code>in_range</code> (in which case at least one of <code>start</code> or <code>end</code> must be specified) or <code>lt</code> (in which case <code>end</code> is required).  |
| <code>operation</code> | Compares the current date/time with the start and/or end date, depending on the context. Allowed values: <ul style="list-style-type: none"> <li>• <code>gt</code>: True if the current date/time is after <code>start</code></li> <li>• <code>lt</code>: True if the current date/time is before <code>end</code></li> <li>• <code>in_range</code>: True if the current date/time is after <code>start</code> (if specified) and before either <code>end</code> (if specified) or <code>start</code> plus the value of the <code>duration</code> element (if one is contained in the <code>date_expression</code>)</li> <li>• <code>date_spec</code>: True if the current date/time matches the specification given in the contained <code>date_spec</code> element (described below)</li> </ul> |

---

**Note:** There is no `eq`, `neq`, `gte`, or `lte` operation, since they would be valid only for a single second.

---

#### Date Specifications

A `date_spec` element is used to create a cron-like expression relating to time. Each field can contain a single number or range. Any field not supplied is ignored.

Table 22: Attributes of a `date_spec` Element

| Attribute              | Description   |
|------------------------|---|
| <code>id</code>        | A unique name for this element (required)   |
| <code>hours</code>     | Allowed values: 0-23 (where 0 is midnight and 23 is 11 p.m.)  |
| <code>monthdays</code> | Allowed values: 1-31 (depending on month and year)  |
| <code>weekdays</code>  | Allowed values: 1-7 (where 1 is Monday and 7 is Sunday)   |
| <code>yeardays</code>  | Allowed values: 1-366 (depending on the year)   |
| <code>months</code>    | Allowed values: 1-12  |
| <code>weeks</code>     | Allowed values: 1-53 (depending on weekyear)  |
| <code>years</code>     | Year according to the Gregorian calendar  |
| <code>weekyears</code> | Year in which the week started; for example, 1 January 2005 can be specified in ISO 8601 as “2005-001 Ordinal”, “2005-01-01 Gregorian” or “2004-W53-6 Weekly” and thus would match <code>years="2005"</code> or <code>weekyears="2004"</code> |
| <code>moon</code>      | Allowed values are 0-7 (where 0 is the new moon and 4 is full moon). Seriously, you can use this. This was implemented to demonstrate the ease with which new comparisons could be added.   |

For example, `monthdays="1"` matches the first day of every month, and `hours="09-17"` matches the hours between 9 a.m. and 5 p.m. (inclusive).

At this time, multiple ranges (e.g. `weekdays="1,2"` or `weekdays="1-2,5-6"`) are not supported.

---

**Note:** Pacemaker can calculate when evaluation of a `date_expression` with an operation of `gt`, `lt`, or `in_range` will next change, and schedule a cluster re-check for that time. However, it does not do this for `date_spec`. Instead, it evaluates the `date_spec` whenever a cluster re-check naturally happens via a cluster event or the `cluster-recheck-interval` cluster option.

For example, if you have a `date_spec` enabling a resource from 9 a.m. to 5 p.m., and `cluster-recheck-interval` has been set to 5 minutes, then sometime between 9 a.m. and 9:05 a.m. the cluster would notice that it needs to start the resource, and sometime between 5 p.m. and 5:05 p.m. it would realize that it needs to stop the resource. The timing of the actual start and stop actions will further depend on factors such as any other actions the cluster may need to perform first, and the load of the machine.

---

## Durations

A `duration` is used to calculate a value for `end` when one is not supplied to `in_range` operations. It contains one or more attributes each containing a single number. Any attribute not supplied is ignored.

Table 23: Attributes of a duration Element

| Attribute | Description   |
|-----------|---|
| id        | A unique name for this element (required)             |
| seconds   | This many seconds will be added to the total duration |
| minutes   | This many minutes will be added to the total duration |
| hours     | This many hours will be added to the total duration   |
| weeks     | This many weeks will be added to the total duration   |
| months    | This many months will be added to the total duration  |
| years     | This many years will be added to the total duration   |

### Example Time-Based Expressions

A small sample of how time-based expressions can be used:

**True if now is any time in the year 2005**

```
<rule id="rule1" score="INFINITY">
  <date_expression id="date_expr1" start="2005-001" operation="in_range">
    <duration id="duration1" years="1"/>
  </date_expression>
</rule>
```

or equivalently:

```
<rule id="rule2" score="INFINITY">
  <date_expression id="date_expr2" operation="date_spec">
    <date_spec id="date_spec2" years="2005"/>
  </date_expression>
</rule>
```

**9 a.m. to 5 p.m. Monday through Friday**

```
<rule id="rule3" score="INFINITY">
  <date_expression id="date_expr3" operation="date_spec">
    <date_spec id="date_spec3" hours="9-16" weekdays="1-5"/>
  </date_expression>
</rule>
```

Note that the 16 matches all the way through 16:59:59, because the numeric value of the hour still matches.

**9 a.m. to 6 p.m. Monday through Friday or anytime Saturday**

```
<rule id="rule4" score="INFINITY" boolean-op="or">
  <date_expression id="date_expr4-1" operation="date_spec">
    <date_spec id="date_spec4-1" hours="9-16" weekdays="1-5"/>
  </date_expression>
  <date_expression id="date_expr4-2" operation="date_spec">
    <date_spec id="date_spec4-2" weekdays="6"/>
  </date_expression>
</rule>
```

### 9 a.m. to 5 p.m. or 9 p.m. to 12 a.m. Monday through Friday

```
<rule id="rule5" score="INFINITY" boolean-op="and">
  <rule id="rule5-nested1" score="INFINITY" boolean-op="or">
    <date_expression id="date_expr5-1" operation="date_spec">
      <date_spec id="date_spec5-1" hours="9-16"/>
    </date_expression>
    <date_expression id="date_expr5-2" operation="date_spec">
      <date_spec id="date_spec5-2" hours="21-23"/>
    </date_expression>
  </rule>
  <date_expression id="date_expr5-3" operation="date_spec">
    <date_spec id="date_spec5-3" weekdays="1-5"/>
  </date_expression>
</rule>
```

### Mondays in March 2005

```
<rule id="rule6" score="INFINITY" boolean-op="and">
  <date_expression id="date_expr6-1" operation="date_spec">
    <date_spec id="date_spec6" weekdays="1"/>
  </date_expression>
  <date_expression id="date_expr6-2" operation="in_range">
    start="2005-03-01" end="2005-04-01"/>
  </date_expression>
</rule>
```

**Note:** Because no time is specified with the above dates, 00:00:00 is implied. This means that the range includes all of 2005-03-01 but none of 2005-04-01. You may wish to write end as "2005-03-31T23:59:59" to avoid confusion.

### A full moon on Friday the 13th

```
<rule id="rule7" score="INFINITY" boolean-op="and">
  <date_expression id="date_expr7" operation="date_spec">
    <date_spec id="date_spec7" weekdays="5" monthdays="13" moon="4"/>
  </date_expression>
</rule>
```



## 2.8.4 Resource Expressions

An `rsc_expression` (since 2.0.5) is a rule condition based on a resource agent's properties. This rule is only valid within an `rsc_defaults` or `op_defaults` context. None of the matching attributes of `class`, `provider`, and `type` are required. If one is omitted, all values of that attribute will match. For instance, omitting `type` means every type will match.

Table 24: Attributes of a `rsc_expression` Element

| Attribute | Description   |
|-----------|---|
| id        | A unique name for this element (required)   |
| class     | The standard name to be matched against resource agents   |
| provider  | If given, the vendor to be matched against resource agents (only relevant when <code>class</code> is <code>ocf</code> ) |
| type      | The name of the resource agent to be matched  |

### Example Resource-Based Expressions

A small sample of how resource-based expressions can be used:

**True for all `ocf:heartbeat:IPaddr2` resources**

```
<rule id="rule1" score="INFINITY">
  <rsc_expression id="rule_expr1" class="ocf" provider="heartbeat" type="IPaddr2"/>
</rule>
```

**Provider doesn't apply to non-OCF resources**

```
<rule id="rule2" score="INFINITY">
  <rsc_expression id="rule_expr2" class="stonith" type="fence_xvm"/>
</rule>
```

## 2.8.5 Operation Expressions

An `op_expression` (since 2.0.5) is a rule condition based on an action of some resource agent. This rule is only valid within an `op_defaults` context.

Table 25: Attributes of an `op_expression` Element

| Attribute | Description   |
|-----------|---|
| id        | A unique name for this element (required)   |
| name      | The action name to match against. This can be any action supported by the resource agent; common values include <code>monitor</code> , <code>start</code> , and <code>stop</code> (required). |
| interval  | The interval of the action to match against. If not given, only the name attribute will be used to match.   |

### Example Operation-Based Expressions

A small sample of how operation-based expressions can be used:

**True for all monitor actions**

```
<rule id="rule1" score="INFINITY">
  <op_expression id="rule_expr1" name="monitor"/>
</rule>
```

**True for all monitor actions with a 10 second interval**

```
<rule id="rule2" score="INFINITY">
  <op_expression id="rule_expr2" name="monitor" interval="10s"/>
</rule>
```

## 2.8.6 Using Rules to Determine Resource Location

A location constraint may contain one or more top-level rules. The cluster will act as if there is a separate location constraint for each rule that evaluates as true.

Consider the following simple location constraint:

**Prevent resource `webserver` from running on node `node3`**

```
<rsc_location id="ban-apache-on-node3" rsc="webserver"
  score="-INFINITY" node="node3"/>
```

The same constraint can be more verbosely written using a rule:

**Prevent resource `webserver` from running on node `node3` using a rule**

```
<rsc_location id="ban-apache-on-node3" rsc="webserver">
  <rule id="ban-apache-rule" score="-INFINITY">
    <expression id="ban-apache-expr" attribute="#uname"
      operation="eq" value="node3"/>
  </rule>
</rsc_location>
```

The advantage of using the expanded form is that one could add more expressions (for example, limiting the constraint to certain days of the week), or activate the constraint by some node attribute other than node name.

### Location Rules Based on Other Node Properties

The expanded form allows us to match on node properties other than its name. If we rated each machine's CPU power such that the cluster had the following nodes section:

**Sample node section with node attributes**

```
<nodes>
  <node id="uuid1" uname="c001n01" type="normal">
    <instance_attributes id="uuid1-custom_attrs">
      <nvpair id="uuid1-cpu_mips" name="cpu_mips" value="1234"/>
    </instance_attributes>
  </node>
  <node id="uuid2" uname="c001n02" type="normal">
    <instance_attributes id="uuid2-custom_attrs">
      <nvpair id="uuid2-cpu_mips" name="cpu_mips" value="5678"/>
    </instance_attributes>
  </node>
</nodes>
```

then we could prevent resources from running on underpowered machines with this rule:

**Rule using a node attribute (to be used inside a location constraint)**

```
<rule id="need-more-power-rule" score="-INFINITY">
  <expression id="need-more-power-expr" attribute="cpu_mips"
    operation="lt" value="3000"/>
</rule>
```

**Using score-attribute Instead of score**

When using `score-attribute` instead of `score`, each node matched by the rule has its score adjusted differently, according to its value for the named node attribute. Thus, in the previous example, if a rule inside a location constraint for a resource used `score-attribute="cpu_mips"`, `c001n01` would have its preference to run the resource increased by 1234 whereas `c001n02` would have its preference increased by 5678.

**2.8.7 Using Rules to Define Options**

Rules may be used to control a variety of options:

- *Cluster options* (`cluster_property_set` elements)
- *Node attributes* (`instance_attributes` or `utilization` elements inside a `node` element)
- *Resource options* (`utilization`, `meta_attributes`, or `instance_attributes` elements inside a resource definition element or `op`, `rsc_defaults`, `op_defaults`, or `template` element)
- *Operation properties* (`meta_attributes` elements inside an `op` or `op_defaults` element)

**Note:** Attribute-based expressions for meta-attributes can only be used within `operations` and `op_defaults`. They will not work with resource configuration or `rsc_defaults`. Additionally, attribute-based expressions cannot be used with cluster options.

## Using Rules to Control Resource Options

Often some cluster nodes will be different from their peers. Sometimes, these differences – e.g. the location of a binary or the names of network interfaces – require resources to be configured differently depending on the machine they’re hosted on.

By defining multiple `instance_attributes` objects for the resource and adding a rule to each, we can easily handle these special cases.

In the example below, `mySpecialRsc` will use `eth1` and port 9999 when run on `node1`, `eth2` and port 8888 on `node2` and default to `eth0` and port 9999 for all other nodes.

### Defining different resource options based on the node name

```
<primitive id="mySpecialRsc" class="ocf" type="Special" provider="me">
  <instance_attributes id="special-node1" score="3">
    <rule id="node1-special-case" score="INFINITY" >
      <expression id="node1-special-case-expr" attribute="#uname"
        operation="eq" value="node1"/>
    </rule>
    <nvpair id="node1-interface" name="interface" value="eth1"/>
  </instance_attributes>
  <instance_attributes id="special-node2" score="2" >
    <rule id="node2-special-case" score="INFINITY">
      <expression id="node2-special-case-expr" attribute="#uname"
        operation="eq" value="node2"/>
    </rule>
    <nvpair id="node2-interface" name="interface" value="eth2"/>
    <nvpair id="node2-port" name="port" value="8888"/>
  </instance_attributes>
  <instance_attributes id="defaults" score="1" >
    <nvpair id="default-interface" name="interface" value="eth0"/>
    <nvpair id="default-port" name="port" value="9999"/>
  </instance_attributes>
</primitive>
```

The order in which `instance_attributes` objects are evaluated is determined by their score (highest to lowest). If not supplied, the score defaults to zero. Objects with an equal score are processed in their listed order. If the `instance_attributes` object has no rule, or a rule that evaluates to `true`, then for any parameter the resource does not yet have a value for, the resource will use the parameter values defined by the `instance_attributes`.

For example, given the configuration above, if the resource is placed on `node1`:

- `special-node1` has the highest score (3) and so is evaluated first; its rule evaluates to `true`, so `interface` is set to `eth1`.
- `special-node2` is evaluated next with score 2, but its rule evaluates to `false`, so it is ignored.
- `defaults` is evaluated last with score 1, and has no rule, so its values are examined; `interface` is already defined, so the value here is not used, but `port` is not yet defined, so `port` is set to 9999.

## Using Rules to Control Resource Defaults

Rules can be used for resource and operation defaults. The following example illustrates how to set a different `resource-stickiness` value during and outside work hours. This allows resources to automatically move back to their most preferred hosts, but at a time that (in theory) does not interfere with business activities.

**Change resource-stickiness during working hours**

```

<rsc_defaults>
  <meta_attributes id="core-hours" score="2">
    <rule id="core-hour-rule" score="0">
      <date_expression id="nine-to-five-Mon-to-Fri" operation="date_spec">
        <date_spec id="nine-to-five-Mon-to-Fri-spec" hours="9-16" weekdays="1-5"/>
      </date_expression>
    </rule>
    <nvpair id="core-stickiness" name="resource-stickiness" value="INFINITY"/>
  </meta_attributes>
  <meta_attributes id="after-hours" score="1" >
    <nvpair id="after-stickiness" name="resource-stickiness" value="0"/>
  </meta_attributes>
</rsc_defaults>

```

Rules may be used similarly in `instance_attributes` or `utilization` blocks.

Any single block may directly contain only a single rule, but that rule may itself contain any number of rules.

`rsc_expression` and `op_expression` blocks may additionally be used to set defaults on either a single resource or across an entire class of resources with a single rule. `rsc_expression` may be used to select resource agents within both `rsc_defaults` and `op_defaults`, while `op_expression` may only be used within `op_defaults`. If multiple rules succeed for a given resource agent, the last one specified will be the one that takes effect. As with any other rule, boolean operations may be used to make more complicated expressions.

**Default all IPAddr2 resources to stopped**

```

<rsc_defaults>
  <meta_attributes id="op-target-role">
    <rule id="op-target-role-rule" score="INFINITY">
      <rsc_expression id="op-target-role-expr" class="ocf" provider="heartbeat"
        type="IPAddr2"/>
    </rule>
    <nvpair id="op-target-role-nvpair" name="target-role" value="Stopped"/>
  </meta_attributes>
</rsc_defaults>

```

**Default all monitor action timeouts to 7 seconds**

```

<op_defaults>
  <meta_attributes id="op-monitor-defaults">
    <rule id="op-monitor-default-rule" score="INFINITY">
      <op_expression id="op-monitor-default-expr" name="monitor"/>
    </rule>
    <nvpair id="op-monitor-timeout" name="timeout" value="7s"/>
  </meta_attributes>
</op_defaults>

```

Default the timeout on all 10-second-interval monitor actions on IPAddr2 resources to 8 seconds

```
<op_defaults>
  <meta_attributes id="op-monitor-and">
    <rule id="op-monitor-and-rule" score="INFINITY">
      <rsc_expression id="op-monitor-and-rsc-expr" class="ocf" provider="heartbeat"
        type="IPAddr2"/>
      <op_expression id="op-monitor-and-op-expr" name="monitor" interval="10s"/>
    </rule>
    <nvpair id="op-monitor-and-timeout" name="timeout" value="8s"/>
  </meta_attributes>
</op_defaults>
```

## Using Rules to Control Cluster Options

Controlling cluster options is achieved in much the same manner as specifying different resource options on different nodes.

The following example illustrates how to set `maintenance_mode` during a scheduled maintenance window. This will keep the cluster running but not monitor, start, or stop resources during this time.

Schedule a maintenance window for 9 to 11 p.m. CDT Sept. 20, 2019

```
<crm_config>
  <cluster_property_set id="cib-bootstrap-options">
    <nvpair id="bootstrap-stonith-enabled" name="stonith-enabled" value="1"/>
  </cluster_property_set>
  <cluster_property_set id="normal-set" score="10">
    <nvpair id="normal-maintenance-mode" name="maintenance-mode" value="false"/>
  </cluster_property_set>
  <cluster_property_set id="maintenance-window-set" score="1000">
    <nvpair id="maintenance-nvpair1" name="maintenance-mode" value="true"/>
    <rule id="maintenance-rule1" score="INFINITY">
      <date_expression id="maintenance-date1" operation="in_range"
        start="2019-09-20 21:00:00 -05:00" end="2019-09-20 23:00:00 -05:00"/>
    </rule>
  </cluster_property_set>
</crm_config>
```

---

**Important:** The `cluster_property_set` with an id set to “cib-bootstrap-options” will *always* have the highest priority, regardless of any scores. Therefore, rules in another `cluster_property_set` can never take effect for any properties listed in the bootstrap set.

---

## 2.9 Advanced Configuration

### 2.9.1 Specifying When Recurring Actions are Performed

By default, recurring actions are scheduled relative to when the resource started. In some cases, you might prefer that a recurring action start relative to a specific date and time. For example, you might schedule an in-depth monitor to run once every 24 hours, and want it to run outside business hours.

To do this, set the operation's `interval-origin`. The cluster uses this point to calculate the correct `start-delay` such that the operation will occur at `interval-origin` plus a multiple of the operation interval.

For example, if the recurring operation's interval is 24h, its `interval-origin` is set to 02:00, and it is currently 14:32, then the cluster would initiate the operation after 11 hours and 28 minutes.

The value specified for `interval` and `interval-origin` can be any date/time conforming to the [ISO8601 standard](#). By way of example, to specify an operation that would run on the first Monday of 2021 and every Monday after that, you would add:

**Example recurring action that runs relative to base date/time**

```
<op id="intensive-monitor" name="monitor" interval="P7D" interval-origin="2021-W01-1"/>
```

### 2.9.2 Handling Resource Failure

By default, Pacemaker will attempt to recover failed resources by restarting them. However, failure recovery is highly configurable.

#### Failure Counts

Pacemaker tracks resource failures for each combination of node, resource, and operation (start, stop, monitor, etc.).

You can query the fail count for a particular node, resource, and/or operation using the `crm_failcount` command. For example, to see how many times the 10-second monitor for `myrsc` has failed on `node1`, run:

```
# crm_failcount --query -r myrsc -N node1 -n monitor -I 10s
```

If you omit the node, `crm_failcount` will use the local node. If you omit the operation and interval, `crm_failcount` will display the sum of the fail counts for all operations on the resource.

You can use `crm_resource --cleanup` or `crm_failcount --delete` to clear fail counts. For example, to clear the above monitor failures, run:

```
# crm_resource --cleanup -r myrsc -N node1 -n monitor -I 10s
```

If you omit the resource, `crm_resource --cleanup` will clear failures for all resources. If you omit the node, it will clear failures on all nodes. If you omit the operation and interval, it will clear the failures for all operations on the resource.

**Note:** Even when cleaning up only a single operation, all failed operations will disappear from the status display. This allows us to trigger a re-check of the resource's current status.

Higher-level tools may provide other commands for querying and clearing fail counts.

The `crm_mon` tool shows the current cluster status, including any failed operations. To see the current fail counts for any failed resources, call `crm_mon` with the `--failcounts` option. This shows the fail counts per resource (that is, the sum of any operation fail counts for the resource).

### Failure Response

Normally, if a running resource fails, pacemaker will try to stop it and start it again. Pacemaker will choose the best location to start it each time, which may be the same node that it failed on.

However, if a resource fails repeatedly, it is possible that there is an underlying problem on that node, and you might desire trying a different node in such a case. Pacemaker allows you to set your preference via the `migration-threshold` resource meta-attribute.<sup>1</sup>

If you define `migration-threshold` to  $N$  for a resource, it will be banned from the original node after  $N$  failures there.

---

**Note:** The `migration-threshold` is per *resource*, even though fail counts are tracked per *operation*. The operation fail counts are added together to compare against the `migration-threshold`.

---

By default, fail counts remain until manually cleared by an administrator using `crm_resource --cleanup` or `crm_failcount --delete` (hopefully after first fixing the failure's cause). It is possible to have fail counts expire automatically by setting the `failure-timeout` resource meta-attribute.

---

**Important:** A successful operation does not clear past failures. If a recurring monitor operation fails once, succeeds many times, then fails again days later, its fail count is 2. Fail counts are cleared only by manual intervention or failure timeout.

---

For example, setting `migration-threshold` to 2 and `failure-timeout` to 60s would cause the resource to move to a new node after 2 failures, and allow it to move back (depending on stickiness and constraint scores) after one minute.

---

**Note:** `failure-timeout` is measured since the most recent failure. That is, older failures do not individually time out and lower the fail count. Instead, all failures are timed out simultaneously (and the fail count is reset to 0) if there is no new failure for the timeout period.

---

There are two exceptions to the migration threshold: when a resource either fails to start or fails to stop.

If the cluster property `start-failure-is-fatal` is set to `true` (which is the default), start failures cause the fail count to be set to `INFINITY` and thus always cause the resource to move immediately.

Stop failures are slightly different and crucial. If a resource fails to stop and fencing is enabled, then the cluster will fence the node in order to be able to start the resource elsewhere. If fencing is disabled, then the cluster has no way to continue and will not try to start the resource elsewhere, but will try to stop it again after any failure timeout or clearing.

---

<sup>1</sup> The naming of this option was perhaps unfortunate as it is easily confused with live migration, the process of moving a resource from one node to another without stopping it. Xen virtual guests are the most common example of resources that can be migrated in this manner.



## 2.9.3 Moving Resources

### Moving Resources Manually

There are primarily two occasions when you would want to move a resource from its current location: when the whole node is under maintenance, and when a single resource needs to be moved.

### Standby Mode

Since everything eventually comes down to a score, you could create constraints for every resource to prevent them from running on one node. While Pacemaker configuration can seem convoluted at times, not even we would require this of administrators.

Instead, you can set a special node attribute which tells the cluster “don’t let anything run here”. There is even a helpful tool to help query and set it, called `crm_standby`. To check the standby status of the current machine, run:

```
# crm_standby -G
```

A value of `on` indicates that the node is *not* able to host any resources, while a value of `off` says that it *can*. You can also check the status of other nodes in the cluster by specifying the `-node` option:

```
# crm_standby -G --node sles-2
```

To change the current node’s standby status, use `-v` instead of `-G`:

```
# crm_standby -v on
```

Again, you can change another host’s value by supplying a hostname with `--node`.

A cluster node in standby mode will not run resources, but still contributes to quorum, and may fence or be fenced by nodes.

### Moving One Resource

When only one resource is required to move, we could do this by creating location constraints. However, once again we provide a user-friendly shortcut as part of the `crm_resource` command, which creates and modifies the extra constraints for you. If `Email` were running on `sles-1` and you wanted it moved to a specific location, the command would look something like:

```
# crm_resource -M -r Email -H sles-2
```

Behind the scenes, the tool will create the following location constraint:

```
<rsc_location id="cli-prefer-Email" rsc="Email" node="sles-2" score="INFINITY"/>
```

It is important to note that subsequent invocations of `crm_resource -M` are not cumulative. So, if you ran these commands:

```
# crm_resource -M -r Email -H sles-2
# crm_resource -M -r Email -H sles-3
```

then it is as if you had never performed the first command.

To allow the resource to move back again, use:

```
# crm_resource -U -r Email
```

Note the use of the word *allow*. The resource *can* move back to its original location, but depending on `resource-stickiness`, location constraints, and so forth, it might stay where it is.

To be absolutely certain that it moves back to `sles-1`, move it there before issuing the call to `crm_resource -U`:

```
# crm_resource -M -r Email -H sles-1
# crm_resource -U -r Email
```

Alternatively, if you only care that the resource should be moved from its current location, try:

```
# crm_resource -B -r Email
```

which will instead create a negative constraint, like:

```
<rsc_location id="cli-ban-Email-on-sles-1" rsc="Email" node="sles-1" score="-INFINITY"/>
```

This will achieve the desired effect, but will also have long-term consequences. As the tool will warn you, the creation of a `-INFINITY` constraint will prevent the resource from running on that node until `crm_resource -U` is used. This includes the situation where every other cluster node is no longer available!

In some cases, such as when `resource-stickiness` is set to `INFINITY`, it is possible that you will end up with the problem described in *What if Two Nodes Have the Same Score*. The tool can detect some of these cases and deals with them by creating both positive and negative constraints. For example:

```
<rsc_location id="cli-ban-Email-on-sles-1" rsc="Email" node="sles-1" score="-INFINITY"/>
<rsc_location id="cli-prefer-Email" rsc="Email" node="sles-2" score="INFINITY"/>
```

which has the same long-term consequences as discussed earlier.

### Moving Resources Due to Connectivity Changes

You can configure the cluster to move resources when external connectivity is lost in two steps.

#### Tell Pacemaker to Monitor Connectivity

First, add an `ocf:pacemaker:ping` resource to the cluster. The `ping` resource uses the system utility of the same name to test whether a list of machines (specified by DNS hostname or IP address) are reachable, and uses the results to maintain a node attribute.

The node attribute is called `pingd` by default, but is customizable in order to allow multiple ping groups to be defined.

Normally, the `ping` resource should run on all cluster nodes, which means that you'll need to create a clone. A template for this can be found below, along with a description of the most interesting parameters.

Table 26: Commonly Used `ocf:pacemaker:ping` Resource Parameters

| Resource Parameter | Description  |
|--------------------|--|
| dampen             | The time to wait (dampening) for further changes to occur. Use this to prevent a resource from bouncing around the cluster when cluster nodes notice the loss of connectivity at slightly different times. |
| multiplier         | The number of connected ping nodes gets multiplied by this value to get a score. Useful when there are multiple ping nodes configured.   |
| host_list          | The machines to contact in order to determine the current connectivity status. Allowed values include resolvable DNS connectivity host names, IPv4 addresses, and IPv6 addresses.                          |

**Example ping resource that checks node connectivity once every minute**

```
<clone id="Connected">
  <primitive id="ping" class="ocf" provider="pacemaker" type="ping">
    <instance_attributes id="ping-attrs">
      <nvpair id="ping-dampen" name="dampen" value="5s"/>
      <nvpair id="ping-multiplier" name="multiplier" value="1000"/>
      <nvpair id="ping-hosts" name="host_list" value="my.gateway.com www.bigcorp.com"/>
    </instance_attributes>
    <operations>
      <op id="ping-monitor-60s" interval="60s" name="monitor"/>
    </operations>
  </primitive>
</clone>
```

**Important:** You're only half done. The next section deals with telling Pacemaker how to deal with the connectivity status that `ocf:pacemaker:ping` is recording.

**Tell Pacemaker How to Interpret the Connectivity Data**

**Important:** Before attempting the following, make sure you understand *Rules*.

There are a number of ways to use the connectivity data.

The most common setup is for people to have a single ping target (for example, the service network's default gateway), to prevent the cluster from running a resource on any unconnected node.

**Don't run a resource on unconnected nodes**

```
<rule id="WebServer-no-connectivity" rsc="Webserver">
  <rule id="ping-exclude-rule" score="-INFINITY" >
    <expression id="ping-exclude" attribute="pingd" operation="not_defined"/>
  </rule>
```

A more complex setup is to have a number of ping targets configured. You can require the cluster to only run resources on nodes that can connect to all (or a minimum subset) of them.

#### Run only on nodes connected to three or more ping targets

```
<primitive id="ping" provider="pacemaker" class="ocf" type="ping">
... <!-- omitting some configuration to highlight important parts -->
  <nvpair id="ping-multiplier" name="multiplier" value="1000"/>
...
</primitive>
...
<rsc_location id="WebServer-connectivity" rsc="Webserver">
  <rule id="ping-prefer-rule" score="-INFINITY" >
    <expression id="ping-prefer" attribute="pingd" operation="lt" value="3000"/>
  </rule>
</rsc_location>
```

Alternatively, you can tell the cluster only to *prefer* nodes with the best connectivity, by using `score-attribute` in the rule. Just be sure to set `multiplier` to a value higher than that of resource-stickiness (and don't set either of them to `INFINITY`).

#### Prefer node with most connected ping nodes

```
<rsc_location id="WebServer-connectivity" rsc="Webserver">
  <rule id="ping-prefer-rule" score-attribute="pingd" >
    <expression id="ping-prefer" attribute="pingd" operation="defined"/>
  </rule>
</rsc_location>
```

It is perhaps easier to think of this in terms of the simple constraints that the cluster translates it into. For example, if `sles-1` is connected to all five ping nodes but `sles-2` is only connected to two, then it would be as if you instead had the following constraints in your configuration:

#### How the cluster translates the above location constraint

```
<rsc_location id="ping-1" rsc="Webserver" node="sles-1" score="5000"/>
<rsc_location id="ping-2" rsc="Webserver" node="sles-2" score="2000"/>
```

The advantage is that you don't have to manually update any constraints whenever your network connectivity changes.

You can also combine the concepts above into something even more complex. The example below shows how you can prefer the node with the most connected ping nodes provided they have connectivity to at least three (again assuming that `multiplier` is set to 1000).

#### More complex example of choosing location based on connectivity

```
<rsc_location id="WebServer-connectivity" rsc="Webserver">
  <rule id="ping-exclude-rule" score="-INFINITY" >
    <expression id="ping-exclude" attribute="pingd" operation="lt" value="3000"/>
  </rule>
  <rule id="ping-prefer-rule" score-attribute="pingd" >
    <expression id="ping-prefer" attribute="pingd" operation="defined"/>
  </rule>
</rsc_location>
```

## Migrating Resources

Normally, when the cluster needs to move a resource, it fully restarts the resource (that is, it stops the resource on the current node and starts it on the new node).

However, some types of resources, such as many virtual machines, are able to move to another location without loss of state (often referred to as live migration or hot migration). In pacemaker, this is called resource migration. Pacemaker can be configured to migrate a resource when moving it, rather than restarting it.

Not all resources are able to migrate; see the [migration checklist](#) below. Even those that can, won't do so in all situations. Conceptually, there are two requirements from which the other prerequisites follow:

- The resource must be active and healthy at the old location; and
- everything required for the resource to run must be available on both the old and new locations.

The cluster is able to accommodate both *push* and *pull* migration models by requiring the resource agent to support two special actions: `migrate_to` (performed on the current location) and `migrate_from` (performed on the destination).

In push migration, the process on the current location transfers the resource to the new location where it is later activated. In this scenario, most of the work would be done in the `migrate_to` action and, if anything, the activation would occur during `migrate_from`.

Conversely for pull, the `migrate_to` action is practically empty and `migrate_from` does most of the work, extracting the relevant resource state from the old location and activating it.

There is no wrong or right way for a resource agent to implement migration, as long as it works.

### Migration Checklist

- The resource may not be a clone.
- The resource agent standard must be OCF.
- The resource must not be in a failed or degraded state.
- The resource agent must support `migrate_to` and `migrate_from` actions, and advertise them in its meta-data.
- The resource must have the `allow-migrate` meta-attribute set to `true` (which is not the default).

If an otherwise migratable resource depends on another resource via an ordering constraint, there are special situations in which it will be restarted rather than migrated.

For example, if the resource depends on a clone, and at the time the resource needs to be moved, the clone has instances that are stopping and instances that are starting, then the resource will be restarted. The scheduler is not yet able to model this situation correctly and so takes the safer (if less optimal) path.

Also, if a migratable resource depends on a non-migratable resource, and both need to be moved, the migratable resource will be restarted.

## 2.9.4 Reloading an Agent After a Definition Change

The cluster automatically detects changes to the configuration of active resources. The cluster's normal response is to stop the service (using the old definition) and start it again (with the new definition). This works, but some resource agents are smarter and can be told to use a new set of options without restarting.

To take advantage of this capability, the resource agent must:

- Implement the **reload-agent** action. What it should do depends completely on your application!

---

**Note:** Resource agents may also implement a **reload** action to make the managed service reload its own *native* configuration. This is different from **reload-agent**, which makes effective changes in the resource's *Pacemaker* configuration (specifically, the values of the agent's reloadable parameters).

---

- Advertise the **reload-agent** operation in the **actions** section of its meta-data.
- Set the **reloadable** attribute to 1 in the **parameters** section of its meta-data for any parameters eligible to be reloaded after a change.

Once these requirements are satisfied, the cluster will automatically know to reload the resource (instead of restarting) when a reloadable parameter changes.

---

**Note:** Metadata will not be re-read unless the resource needs to be started. If you edit the agent of an already active resource to set a parameter reloadable, the resource may restart the first time the parameter value changes.

---

---

**Note:** If both a reloadable and non-reloadable parameter are changed simultaneously, the resource will be restarted.

---

## 2.10 Advanced Resource Types

### 2.10.1 Groups - A Syntactic Shortcut

One of the most common elements of a cluster is a set of resources that need to be located together, start sequentially, and stop in the reverse order. To simplify this configuration, we support the concept of groups.

#### A group of two primitive resources

```
<group id="shortcut">
  <primitive id="Public-IP" class="ocf" type="IPAddr" provider="heartbeat">
    <instance_attributes id="params-public-ip">
      <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
    </instance_attributes>
  </primitive>
  <primitive id="Email" class="lsb" type="exim"/>
</group>
```

Although the example above contains only two resources, there is no limit to the number of resources a group can contain. The example is also sufficient to explain the fundamental properties of a group:

- Resources are started in the order they appear in (**Public-IP** first, then **Email**)
- Resources are stopped in the reverse order to which they appear in (**Email** first, then **Public-IP**)

If a resource in the group can't run anywhere, then nothing after that is allowed to run, too.

- If **Public-IP** can't run anywhere, neither can **Email**;
- but if **Email** can't run anywhere, this does not affect **Public-IP** in any way

The group above is logically equivalent to writing:

#### How the cluster sees a group resource

```
<configuration>
  <resources>
    <primitive id="Public-IP" class="ocf" type="IPaddr" provider="heartbeat">
      <instance_attributes id="params-public-ip">
        <nvpair id="public-ip-addr" name="ip" value="192.0.2.2"/>
      </instance_attributes>
    </primitive>
    <primitive id="Email" class="lsb" type="exim"/>
  </resources>
  <constraints>
    <rsc_colocation id="xxx" rsc="Email" with-rsc="Public-IP" score="INFINITY"/>
    <rsc_order id="yyy" first="Public-IP" then="Email"/>
  </constraints>
</configuration>
```

Obviously as the group grows bigger, the reduced configuration effort can become significant.

Another (typical) example of a group is a DRBD volume, the filesystem mount, an IP address, and an application that uses them.

## Group Properties

Table 27: Properties of a Group Resource

| Field | Description                 |
|-------|-----------------------------|
| id    | A unique name for the group |

## Group Options

Groups inherit the `priority`, `target-role`, and `is-managed` properties from primitive resources. See *Resource Options* for information about those properties.

Table 28: Group-specific configuration options

| Meta-Attribute | Default | Description  |
|----------------|---------|--|
| ordered        | true    | If <b>true</b> , group members will be started in the order they are listed in the configuration (and stopped in the reverse order). |

### Group Instance Attributes

Groups have no instance attributes. However, any that are set for the group object will be inherited by the group's children.

### Group Contents

Groups may only contain a collection of cluster resources (see *Resource Properties*). To refer to a child of a group resource, just use the child's id instead of the group's.

### Group Constraints

Although it is possible to reference a group's children in constraints, it is usually preferable to reference the group itself.

#### Some constraints involving groups

```
<constraints>
  <rsc_location id="group-prefers-node1" rsc="shortcut" node="node1" score="500"/>
  <rsc_colocation id="webserver-with-group" rsc="Webserver" with-rsc="shortcut"/>
  <rsc_order id="start-group-then-webserver" first="Webserver" then="shortcut"/>
</constraints>
```

### Group Stickiness

Stickiness, the measure of how much a resource wants to stay where it is, is additive in groups. Every active resource of the group will contribute its stickiness value to the group's total. So if the default `resource-stickiness` is 100, and a group has seven members, five of which are active, then the group as a whole will prefer its current location with a score of 500.

## 2.10.2 Clones - Resources That Can Have Multiple Active Instances

*Clone* resources are resources that can have more than one copy active at the same time. This allows you, for example, to run a copy of a daemon on every node. You can clone any primitive or group resource<sup>1</sup>.

<sup>1</sup> Of course, the service must support running multiple instances.



## Anonymous versus Unique Clones

A clone resource is configured to be either *anonymous* or *globally unique*.

Anonymous clones are the simplest. These behave completely identically everywhere they are running. Because of this, there can be only one instance of an anonymous clone active per node.

The instances of globally unique clones are distinct entities. All instances are launched identically, but one instance of the clone is not identical to any other instance, whether running on the same node or a different node. As an example, a cloned IP address can use special kernel functionality such that each instance handles a subset of requests for the same IP address.

## Promotable clones

If a clone is *promotable*, its instances can perform a special role that Pacemaker will manage via the **promote** and **demote** actions of the resource agent.

Services that support such a special role have various terms for the special role and the default role: primary and secondary, master and replica, controller and worker, etc. Pacemaker uses the terms *promoted* and *unpromoted* to be agnostic to what the service calls them or what they do.

All that Pacemaker cares about is that an instance comes up in the unpromoted role when started, and the resource agent supports the **promote** and **demote** actions to manage entering and exiting the promoted role.

## Clone Properties

Table 29: Properties of a Clone Resource

| Field | Description                 |
|-------|-----------------------------|
| id    | A unique name for the clone |

## Clone Options

*Options* inherited from primitive resources: **priority**, **target-role**, **is-managed**

Table 30: Clone-specific configuration options

| Field           | Default | Description  |
|-----------------|---------|--|
| globally-unique | false   | If <b>true</b> , each clone instance performs a distinct function  |
| clone-max       | 0       | The maximum number of clone instances that can be started across the entire cluster. If 0, the number of nodes in the cluster will be used.  |
| clone-node-max  | 1       | If <b>globally-unique</b> is <b>true</b> , the maximum number of clone instances that can be started on a single node  |
| clone-min       | 0       | Require at least this number of clone instances to be runnable before allowing resources depending on the clone to be runnable. A value of 0 means require all clone instances to be runnable. |

Continued on next page

Table 30 – continued from previous page

| Field             | Default | Description  |
|-------------------|---------|--|
| notify            | false   | Call the resource agent's <b>notify</b> action for all active instances, before and after starting or stopping any clone instance. The resource agent must support this action. Allowed values: <b>false</b> , <b>true</b>   |
| ordered           | false   | If <b>true</b> , clone instances must be started sequentially instead of in parallel. Allowed values: <b>false</b> , <b>true</b>   |
| interleave        | false   | When this clone is ordered relative to another clone, if this option is <b>false</b> (the default), the ordering is relative to <i>all</i> instances of the other clone, whereas if this option is <b>true</b> , the ordering is relative only to instances on the same node. Allowed values: <b>false</b> , <b>true</b> |
| promotable        | false   | If <b>true</b> , clone instances can perform a special role that Pacemaker will manage via the resource agent's <b>promote</b> and <b>demote</b> actions. The resource agent must support these actions. Allowed values: <b>false</b> , <b>true</b>  |
| promoted-max      | 1       | If <b>promotable</b> is <b>true</b> , the number of instances that can be promoted at one time across the entire cluster   |
| promoted-node-max | 1       | If <b>promotable</b> is <b>true</b> and <b>globally-unique</b> is <b>false</b> , the number of clone instances can be promoted at one time on a single node  |

**Note: Deprecated Terminology**

In older documentation and online examples, you may see promotable clones referred to as *multi-state*, *stateful*, or *master/slave*; these mean the same thing as *promotable*. Certain syntax is supported for backward compatibility, but is deprecated and will be removed in a future version:

- Using a **master** tag, instead of a **clone** tag with the **promotable** meta-attribute set to **true**
- Using the **master-max** meta-attribute instead of **promoted-max**
- Using the **master-node-max** meta-attribute instead of **promoted-node-max**
- Using **Master** as a role name instead of **Promoted**
- Using **Slave** as a role name instead of **Unpromoted**

**Clone Contents**

Clones must contain exactly one primitive or group resource.

**A clone that runs a web server on all nodes**

```
<clone id="apache-clone">
  <primitive id="apache" class="lsb" type="apache">
    <operations>
      <op id="apache-monitor" name="monitor" interval="30"/>
    </operations>
  </primitive>
</clone>
```

**Warning:** You should never reference the name of a clone's child (the primitive or group resource being cloned). If you think you need to do this, you probably need to re-evaluate your design.

## Clone Instance Attribute

Clones have no instance attributes; however, any that are set here will be inherited by the clone's child.

## Clone Constraints

In most cases, a clone will have a single instance on each active cluster node. If this is not the case, you can indicate which nodes the cluster should preferentially assign copies to with resource location constraints. These constraints are written no differently from those for primitive resources except that the clone's **id** is used.

### Some constraints involving clones

```
<constraints>
  <rsc_location id="clone-prefers-node1" rsc="apache-clone" node="node1" score="500"/>
  <rsc_colocation id="stats-with-clone" rsc="apache-stats" with="apache-clone"/>
  <rsc_order id="start-clone-then-stats" first="apache-clone" then="apache-stats"/>
</constraints>
```

Ordering constraints behave slightly differently for clones. In the example above, **apache-stats** will wait until all copies of **apache-clone** that need to be started have done so before being started itself. Only if *no* copies can be started will **apache-stats** be prevented from being active. Additionally, the clone will wait for **apache-stats** to be stopped before stopping itself.

Colocation of a primitive or group resource with a clone means that the resource can run on any node with an active instance of the clone. The cluster will choose an instance based on where the clone is running and the resource's own location preferences.

Colocation between clones is also possible. If one clone **A** is colocated with another clone **B**, the set of allowed locations for **A** is limited to nodes on which **B** is (or will be) active. Placement is then performed normally.

## Promotable Clone Constraints

For promotable clone resources, the **first-action** and/or **then-action** fields for ordering constraints may be set to **promote** or **demote** to constrain the promoted role, and colocation constraints may contain **rsc-role** and/or **with-rsc-role** fields.

### Constraints involving promotable clone resources

```
<constraints>
  <rsc_location id="db-prefers-node1" rsc="database" node="node1" score="500"/>
  <rsc_colocation id="backup-with-db-unpromoted" rsc="backup"
    with-rsc="database" with-rsc-role="Unpromoted"/>
  <rsc_colocation id="myapp-with-db-promoted" rsc="myApp"
    with-rsc="database" with-rsc-role="Promoted"/>
  <rsc_order id="start-db-before-backup" first="database" then="backup"/>
  <rsc_order id="start-myapp-after-backup" first="backup" then="myApp"
    then-action="start"/>
</constraints>
```

In the example above, **myApp** will wait until one of the database copies has been started and promoted before being started itself on the same node. Only if no copies can be promoted will **myApp** be prevented from being active. Additionally, the cluster will wait for **myApp** to be stopped before demoting the database.

Colocation of a primitive or group resource with a promotable clone resource means that it can run on any node with an active instance of the promotable clone resource that has the specified role (**Promoted** or **Unpromoted**). In the example above, the cluster will choose a location based on where database is running in the promoted role, and if there are multiple promoted instances it will also factor in **myApp**'s own location preferences when deciding which location to choose.

Colocation with regular clones and other promotable clone resources is also possible. In such cases, the set of allowed locations for the **rsc** clone is (after role filtering) limited to nodes on which the **with-rsc** promotable clone resource is (or will be) in the specified role. Placement is then performed as normal.

### Using Promotable Clone Resources in Colocation Sets

When a promotable clone is used in a *resource set* inside a colocation constraint, the resource set may take a **role** attribute.

In the following example, an instance of **B** may be promoted only on a node where **A** is in the promoted role. Additionally, resources **C** and **D** must be located on a node where both **A** and **B** are promoted.

#### Colocate C and D with A's and B's promoted instances

```
<constraints>
  <rsc_colocation id="coloc-1" score="INFINITY" >
    <resource_set id="colocated-set-example-1" sequential="true" role="Promoted">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="colocated-set-example-2" sequential="true">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_colocation>
</constraints>
```

### Using Promotable Clone Resources in Ordered Sets

When a promotable clone is used in a *resource set* inside an ordering constraint, the resource set may take an **action** attribute.

#### Start C and D after first promoting A and B

```
<constraints>
  <rsc_order id="order-1" score="INFINITY" >
    <resource_set id="ordered-set-1" sequential="true" action="promote">
      <resource_ref id="A"/>
      <resource_ref id="B"/>
    </resource_set>
    <resource_set id="ordered-set-2" sequential="true" action="start">
      <resource_ref id="C"/>
      <resource_ref id="D"/>
    </resource_set>
  </rsc_order>
</constraints>
```

---

In the above example, **B** cannot be promoted until **A** has been promoted. Additionally, resources **C** and **D** must wait until **A** and **B** have been promoted before they can start.

### Clone Stickiness

To achieve a stable allocation pattern, clones are slightly sticky by default. If no value for `resource-stickiness` is provided, the clone will use a value of 1. Being a small value, it causes minimal disturbance to the score calculations of other resources but is enough to prevent Pacemaker from needlessly moving copies around the cluster.

---

**Note:** For globally unique clones, this may result in multiple instances of the clone staying on a single node, even after another eligible node becomes active (for example, after being put into standby mode then made active again). If you do not want this behavior, specify a `resource-stickiness` of 0 for the clone temporarily and let the cluster adjust, then set it back to 1 if you want the default behavior to apply again.

---



---

**Important:** If `resource-stickiness` is set in the `rsc_defaults` section, it will apply to clone instances as well. This means an explicit `resource-stickiness` of 0 in `rsc_defaults` works differently from the implicit default used when `resource-stickiness` is not specified.

---

### Clone Resource Agent Requirements

Any resource can be used as an anonymous clone, as it requires no additional support from the resource agent. Whether it makes sense to do so depends on your resource and its resource agent.

#### Resource Agent Requirements for Globally Unique Clones

Globally unique clones require additional support in the resource agent. In particular, it must only respond with `${OCF_SUCCESS}` if the node has that exact instance active. All other probes for instances of the clone should result in `${OCF_NOT_RUNNING}` (or one of the other OCF error codes if they are failed).

Individual instances of a clone are identified by appending a colon and a numerical offset, e.g. **apache:2**.

Resource agents can find out how many copies there are by examining the `OCF_RESKEY_CRM_meta_clone_max` environment variable and which instance it is by examining `OCF_RESKEY_CRM_meta_clone`.

The resource agent must not make any assumptions (based on `OCF_RESKEY_CRM_meta_clone`) about which numerical instances are active. In particular, the list of active copies will not always be an unbroken sequence, nor always start at 0.

#### Resource Agent Requirements for Promotable Clones

Promotable clone resources require two extra actions, **demote** and **promote**, which are responsible for changing the state of the resource. Like **start** and **stop**, they should return `${OCF_SUCCESS}` if they completed successfully or a relevant error code if they did not.

The states can mean whatever you wish, but when the resource is started, it must come up in the unpromoted role. From there, the cluster will decide which instances to promote.

In addition to the clone requirements for monitor actions, agents must also *accurately* report which state they are in. The cluster relies on the agent to report its status (including role) accurately and does not indicate to the agent what role it currently believes it to be in.

Table 31: Role implications of OCF return codes

| Monitor Return Code  | Description          |
|----------------------|----------------------|
| OCF_NOT_RUNNING      | Stopped              |
| OCF_SUCCESS          | Running (Unpromoted) |
| OCF_RUNNING_PROMOTED | Running (Promoted)   |
| OCF_FAILED_PROMOTED  | Failed (Promoted)    |
| Other                | Failed (Unpromoted)  |

## Clone Notifications

If the clone has the `notify` meta-attribute set to **true**, and the resource agent supports the `notify` action, Pacemaker will call the action when appropriate, passing a number of extra variables which, when combined with additional context, can be used to calculate the current state of the cluster and what is about to happen to it.

Table 32: Environment variables supplied with Clone notify actions

| Variable                                     | Description                                |
|--|--|
| OCF_RESKEY_CRM_meta_notify_type              | Allowed values: <b>pre</b> , <b>post</b>   |
| OCF_RESKEY_CRM_meta_notify_operation         | Allowed values: <b>start</b> , <b>stop</b> |
| OCF_RESKEY_CRM_meta_notify_start_resource    | Resources to be started                    |
| OCF_RESKEY_CRM_meta_notify_stop_resource     | Resources to be stopped                    |
| OCF_RESKEY_CRM_meta_notify_active_resource   | Resources that are running                 |
| OCF_RESKEY_CRM_meta_notify_inactive_resource | Resources that are not running             |
| OCF_RESKEY_CRM_meta_notify_start_uname       | Nodes on which resources will be started   |
| OCF_RESKEY_CRM_meta_notify_stop_uname        | Nodes on which resources will be stopped   |
| OCF_RESKEY_CRM_meta_notify_active_uname      | Nodes on which resources are running       |

The variables come in pairs, such as `OCF_RESKEY_CRM_meta_notify_start_resource` and `OCF_RESKEY_CRM_meta_notify_start_uname`, and should be treated as an array of whitespace-separated elements.

`OCF_RESKEY_CRM_meta_notify_inactive_resource` is an exception, as the matching **uname** variable does not exist since inactive resources are not running on any node.

Thus, in order to indicate that **clone:0** will be started on **sles-1**, **clone:2** will be started on **sles-3**, and **clone:3** will be started on **sles-2**, the cluster would set:

**Notification variables**

```
OCF_RESKEY_CRM_meta_notify_start_resource="clone:0 clone:2 clone:3"
OCF_RESKEY_CRM_meta_notify_start_uname="sles-1 sles-3 sles-2"
```

**Note:** Pacemaker will log but otherwise ignore failures of notify actions.

**Interpretation of Notification Variables****Pre-notification (stop):**

- Active resources: `$OCF_RESKEY_CRM_meta_notify_active_resource`
- Inactive resources: `$OCF_RESKEY_CRM_meta_notify_inactive_resource`
- Resources to be started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources to be stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`

**Post-notification (stop) / Pre-notification (start):**

- Active resources
  - `$OCF_RESKEY_CRM_meta_notify_active_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Inactive resources
  - `$OCF_RESKEY_CRM_meta_notify_inactive_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Resources that were started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources that were stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`

**Post-notification (start):**

- Active resources:
  - `$OCF_RESKEY_CRM_meta_notify_active_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Inactive resources:
  - `$OCF_RESKEY_CRM_meta_notify_inactive_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources that were started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources that were stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`

## Extra Notifications for Promotable Clones

Table 33: Extra environment variables supplied for promotable clones

| Variable                                       | Description   |
|--|---|
| OCF_RESKEY_CRM_meta_notify_promoted_resource   | Resources that are running in the promoted role             |
| OCF_RESKEY_CRM_meta_notify_unpromoted_resource | Resources that are running in the unpromoted role           |
| OCF_RESKEY_CRM_meta_notify_promote_resource    | Resources to be promoted                                    |
| OCF_RESKEY_CRM_meta_notify_demote_resource     | Resources to be demoted                                     |
| OCF_RESKEY_CRM_meta_notify_promote_uname       | Nodes on which resources will be promoted                   |
| OCF_RESKEY_CRM_meta_notify_demote_uname        | Nodes on which resources will be demoted                    |
| OCF_RESKEY_CRM_meta_notify_promoted_uname      | Nodes on which resources are running in the promoted role   |
| OCF_RESKEY_CRM_meta_notify_unpromoted_uname    | Nodes on which resources are running in the unpromoted role |

## Interpretation of Promotable Notification Variables

### Pre-notification (demote):

- Active resources: `$OCF_RESKEY_CRM_meta_notify_active_resource`
- Promoted resources: `$OCF_RESKEY_CRM_meta_notify_promoted_resource`
- Unpromoted resources: `$OCF_RESKEY_CRM_meta_notify_unpromoted_resource`
- Inactive resources: `$OCF_RESKEY_CRM_meta_notify_inactive_resource`
- Resources to be started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources to be promoted: `$OCF_RESKEY_CRM_meta_notify_promote_resource`
- Resources to be demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Resources to be stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`

### Post-notification (demote) / Pre-notification (stop):

- Active resources: `$OCF_RESKEY_CRM_meta_notify_active_resource`
- Promoted resources:
  - `$OCF_RESKEY_CRM_meta_notify_promoted_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Unpromoted resources: `$OCF_RESKEY_CRM_meta_notify_unpromoted_resource`
- Inactive resources: `$OCF_RESKEY_CRM_meta_notify_inactive_resource`
- Resources to be started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources to be promoted: `$OCF_RESKEY_CRM_meta_notify_promote_resource`



- Resources to be demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Resources to be stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Resources that were demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`

#### Post-notification (stop) / Pre-notification (start)

- Active resources:
  - `$OCF_RESKEY_CRM_meta_notify_active_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Promoted resources:
  - `$OCF_RESKEY_CRM_meta_notify_promoted_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Unpromoted resources:
  - `$OCF_RESKEY_CRM_meta_notify_unpromoted_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Inactive resources:
  - `$OCF_RESKEY_CRM_meta_notify_inactive_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Resources to be started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources to be promoted: `$OCF_RESKEY_CRM_meta_notify_promote_resource`
- Resources to be demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Resources to be stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Resources that were demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Resources that were stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`

#### Post-notification (start) / Pre-notification (promote)

- Active resources:
  - `$OCF_RESKEY_CRM_meta_notify_active_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Promoted resources:
  - `$OCF_RESKEY_CRM_meta_notify_promoted_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Unpromoted resources:
  - `$OCF_RESKEY_CRM_meta_notify_unpromoted_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Inactive resources:
  - `$OCF_RESKEY_CRM_meta_notify_inactive_resource`

- plus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources to be started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources to be promoted: `$OCF_RESKEY_CRM_meta_notify_promote_resource`
- Resources to be demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Resources to be stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Resources that were started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources that were demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Resources that were stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`

#### Post-notification (promote)

- Active resources:
  - `$OCF_RESKEY_CRM_meta_notify_active_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Promoted resources:
  - `$OCF_RESKEY_CRM_meta_notify_promoted_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_demote_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_promote_resource`
- Unpromoted resources:
  - `$OCF_RESKEY_CRM_meta_notify_unpromoted_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_start_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_promote_resource`
- Inactive resources:
  - `$OCF_RESKEY_CRM_meta_notify_inactive_resource`
  - plus `$OCF_RESKEY_CRM_meta_notify_stop_resource`
  - minus `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources to be started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources to be promoted: `$OCF_RESKEY_CRM_meta_notify_promote_resource`
- Resources to be demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Resources to be stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`
- Resources that were started: `$OCF_RESKEY_CRM_meta_notify_start_resource`
- Resources that were promoted: `$OCF_RESKEY_CRM_meta_notify_promote_resource`
- Resources that were demoted: `$OCF_RESKEY_CRM_meta_notify_demote_resource`
- Resources that were stopped: `$OCF_RESKEY_CRM_meta_notify_stop_resource`

## Monitoring Promotable Clone Resources

The usual monitor actions are insufficient to monitor a promotable clone resource, because Pacemaker needs to verify not only that the resource is active, but also that its actual role matches its intended one.

Define two monitoring actions: the usual one will cover the unpromoted role, and an additional one with `role="Promoted"` will cover the promoted role.

### Monitoring both states of a promotable clone resource

```
<clone id="myPromotableRsc">
  <meta_attributes id="myPromotableRsc-meta">
    <nvpair name="promotable" value="true"/>
  </meta_attributes>
  <primitive id="myRsc" class="ocf" type="myApp" provider="myCorp">
    <operations>
      <op id="public-ip-unpromoted-check" name="monitor" interval="60"/>
      <op id="public-ip-promoted-check" name="monitor" interval="61" role="Promoted"/>
    </operations>
  </primitive>
</clone>
```

**Important:** It is crucial that *every* monitor operation has a different interval! Pacemaker currently differentiates between operations only by resource and interval; so if (for example) a promotable clone resource had the same monitor interval for both roles, Pacemaker would ignore the role when checking the status – which would cause unexpected return codes, and therefore unnecessary complications.

## Determining Which Instance is Promoted

Pacemaker can choose a promotable clone instance to be promoted in one of two ways:

- Promotion scores: These are node attributes set via the `crm_attribute` command using the `--promotion` option, which generally would be called by the resource agent's start action if it supports promotable clones. This tool automatically detects both the resource and host, and should be used to set a preference for being promoted. Based on this, `promoted-max`, and `promoted-node-max`, the instance(s) with the highest preference will be promoted.
- Constraints: Location constraints can indicate which nodes are most preferred to be promoted.

### Explicitly preferring node1 to be promoted

```
<rsc_location id="promoted-location" rsc="myPromotableRsc">
  <rule id="promoted-rule" score="100" role="Promoted">
    <expression id="promoted-exp" attribute="#uname" operation="eq" value="node1"/>
  </rule>
</rsc_location>
```

### 2.10.3 Bundles - Containerized Resources

Pacemaker supports a special syntax for launching a service inside a `container` with any infrastructure it requires: the *bundle*.

Pacemaker bundles support `Docker`, `podman` (*since 2.0.1*), and `rkt` container technologies.<sup>2</sup>

#### A bundle for a containerized web server

```
<bundle id="httpd-bundle">
  <podman image="pcmk:http" replicas="3"/>
  <network ip-range-start="192.168.122.131"
    host-netmask="24"
    host-interface="eth0">
    <port-mapping id="httpd-port" port="80"/>
  </network>
  <storage>
    <storage-mapping id="httpd-syslog"
      source-dir="/dev/log"
      target-dir="/dev/log"
      options="rw"/>
    <storage-mapping id="httpd-root"
      source-dir="/srv/html"
      target-dir="/var/www/html"
      options="rw,Z"/>
    <storage-mapping id="httpd-logs"
      source-dir-root="/var/log/pacemaker/bundles"
      target-dir="/etc/httpd/logs"
      options="rw,Z"/>
  </storage>
  <primitive class="ocf" id="httpd" provider="heartbeat" type="apache"/>
</bundle>
```

#### Bundle Prerequisites

Before configuring a bundle in Pacemaker, the user must install the appropriate container launch technology (Docker, podman, or rkt), and supply a fully configured container image, on every node allowed to run the bundle.

Pacemaker will create an implicit resource of type **ocf:heartbeat:docker**, **ocf:heartbeat:podman**, or **ocf:heartbeat:rkt** to manage a bundle's container. The user must ensure that the appropriate resource agent is installed on every node allowed to run the bundle.

#### Bundle Properties

Table 34: XML Attributes of a bundle Element

| Attribute   | Description                             |
|-------------|---|
| id          | A unique name for the bundle (required) |
| description | Arbitrary text (not used by Pacemaker)  |

<sup>2</sup> Docker is a trademark of Docker, Inc. No endorsement by or association with Docker, Inc. is implied.

A bundle must contain exactly one `docker`, `podman`, or `rkt` element.

## Bundle Container Properties

Table 35: XML attributes of a `docker`, `podman`, or `rkt` Element

| Attribute                      | Default   | Description  |
|--------------------------------|---|--|
| <code>image</code>             |   | Container image tag (required)   |
| <code>replicas</code>          | Value of <code>promoted-max</code> if that is positive, else 1                                  | A positive integer specifying the number of container instances to launch  |
| <code>replicas-per-host</code> | 1   | A positive integer specifying the number of container instances allowed to run on a single node  |
| <code>promoted-max</code>      | 0   | A non-negative integer that, if positive, indicates that the containerized service should be treated as a promotable service, with this many replicas allowed to run the service in the promoted role  |
| <code>network</code>           |   | If specified, this will be passed to the <code>docker run</code> , <code>podman run</code> , or <code>rkt run</code> command as the network setting for the container.   |
| <code>run-command</code>       | <code>/usr/sbin/pacemaker-remoted</code> if bundle contains a <b>primitive</b> , otherwise none | This command will be run inside the container when launching it (“PID 1”). If the bundle contains a <b>primitive</b> , this command <i>must</i> start <code>pacemaker-remoted</code> (but could, for example, be a script that does other stuff, too). |
| <code>options</code>           |   | Extra command-line options to pass to the <code>docker run</code> , <code>podman run</code> , or <code>rkt run</code> command  |

**Note:** Considerations when using cluster configurations or container images from Pacemaker 1.1:

- If the container image has a pre-2.0.0 version of Pacemaker, set `run-command` to `/usr/sbin/pacemaker_remoted` (note the underbar instead of dash).
- `masters` is accepted as an alias for `promoted-max`, but is deprecated since 2.0.0, and support for it will be removed in a future version.

## Bundle Network Properties

A bundle may optionally contain one `<network>` element.

Table 36: XML attributes of a network Element

| Attribute      | Default | Description   |
|----------------|---------|---|
| add-host       | TRUE    | If TRUE, and <code>ip-range-start</code> is used, Pacemaker will automatically ensure that <code>/etc/hosts</code> inside the containers has entries for each <i>replica name</i> and its assigned IP.  |
| ip-range-start |         | If specified, Pacemaker will create an implicit <code>ocf:heartbeat:IPaddr2</code> resource for each container instance, starting with this IP address, using up to <code>replicas</code> sequential addresses. These addresses can be used from the host's network to reach the service inside the container, though it is not visible within the container itself. Only IPv4 addresses are currently supported.   |
| host-netmask   | 32      | If <code>ip-range-start</code> is specified, the IP addresses are created with this CIDR netmask (as a number of bits).   |
| host-interface |         | If <code>ip-range-start</code> is specified, the IP addresses are created on this host interface (by default, it will be determined from the IP address).   |
| control-port   | 3121    | If the bundle contains a <code>primitive</code> , the cluster will use this integer TCP port for communication with Pacemaker Remote inside the container. Changing this is useful when the container is unable to listen on the default port, for example, when the container uses the host's network rather than <code>ip-range-start</code> (in which case <code>replicas-per-host</code> must be 1), or when the bundle may run on a Pacemaker Remote node that is already listening on the default port. Any <code>PCMK_remote_port</code> environment variable set on the host or in the container is ignored for bundle connections. |

---

**Note:** Replicas are named by the bundle id plus a dash and an integer counter starting with zero. For example, if a bundle named **httpd-bundle** has `replicas=2`, its containers will be named **httpd-bundle-0** and **httpd-bundle-1**.

---

Additionally, a `network` element may optionally contain one or more `port-mapping` elements.

Table 37: Attributes of a port-mapping Element

| Attribute     | Default              | Description   |
|---------------|----------------------|---|
| id            |                      | A unique name for the port mapping (required)   |
| port          |                      | If this is specified, connections to this TCP port number on the host network (on the container's assigned IP address, if <b>ip-range-start</b> is specified) will be forwarded to the container network. Exactly one of <b>port</b> or <b>range</b> must be specified in a <b>port-mapping</b> .   |
| internal-port | value of <b>port</b> | If <b>port</b> and this are specified, connections to <b>port</b> on the host's network will be forwarded to this port on the container network.  |
| range         |                      | If this is specified, connections to these TCP port numbers (expressed as <i>first_port-last_port</i> ) on the host network (on the container's assigned IP address, if <b>ip-range-start</b> is specified) will be forwarded to the same ports in the container network. Exactly one of <b>port</b> or <b>range</b> must be specified in a <b>port-mapping</b> . |

---

**Note:** If the bundle contains a **primitive**, Pacemaker will automatically map the **control-port**, so it is not necessary to specify that port in a **port-mapping**.

---

## Bundle Storage Properties

A bundle may optionally contain one **storage** element. A **storage** element has no properties of its own, but may contain one or more **storage-mapping** elements.

Table 38: Attributes of a storage-mapping Element

| Attribute       | Default | Description   |
|-----------------|---------|---|
| id              |         | A unique name for the storage mapping (required)  |
| source-dir      |         | The absolute path on the host's filesystem that will be mapped into the container. Exactly one of <b>source-dir</b> and <b>source-dir-root</b> must be specified in a <b>storage-mapping</b> .  |
| source-dir-root |         | The start of a path on the host's filesystem that will be mapped into the container, using a different subdirectory on the host for each container instance. The subdirectory will be named the same as the <i>replica name</i> . Exactly one of <b>source-dir</b> and <b>source-dir-root</b> must be specified in a <b>storage-mapping</b> . |
| target-dir      |         | The path name within the container where the host storage will be mapped (required)   |
| options         |         | A comma-separated list of file system mount options to use when mapping the storage   |

---

**Note:** Pacemaker does not define the behavior if the source directory does not already exist on the host. However, it is expected that the container technology and/or its resource agent will create the source

---

directory in that case.

---

**Note:** If the bundle contains a `primitive`, Pacemaker will automatically map the equivalent of `source-dir=/etc/pacemaker/authkey` `target-dir=/etc/pacemaker/authkey` and `source-dir-root=/var/log/pacemaker/bundles` `target-dir=/var/log` into the container, so it is not necessary to specify those paths in a `storage-mapping`.

---

**Important:** The `PCMK_authkey_location` environment variable must not be set to anything other than the default of `/etc/pacemaker/authkey` on any node in the cluster.

---

**Important:** If SELinux is used in enforcing mode on the host, you must ensure the container is allowed to use any storage you mount into it. For Docker and podman bundles, adding “Z” to the mount options will create a container-specific label for the mount that allows the container access.

---

### Bundle Primitive

A bundle may optionally contain one *primitive* resource. The primitive may have operations, instance attributes, and meta-attributes defined, as usual.

If a bundle contains a primitive resource, the container image must include the Pacemaker Remote daemon, and at least one of `ip-range-start` or `control-port` must be configured in the bundle. Pacemaker will create an implicit `ocf:pacemaker:remote` resource for the connection, launch Pacemaker Remote within the container, and monitor and manage the primitive resource via Pacemaker Remote.

If the bundle has more than one container instance (replica), the primitive resource will function as an implicit *clone* – a *promotable clone* if the bundle has `promoted-max` greater than zero.

**Note:** If you want to pass environment variables to a bundle’s Pacemaker Remote connection or primitive, you have two options:

- Environment variables whose value is the same regardless of the underlying host may be set using the container element’s `options` attribute.
  - If you want variables to have host-specific values, you can use the *storage-mapping* element to map a file on the host as `/etc/pacemaker/pcmk-init.env` in the container (*since 2.0.3*). Pacemaker Remote will parse this file as a shell-like format, with variables set as `NAME=VALUE`, ignoring blank lines and comments starting with “#”.
- 

**Important:** When a bundle has a `primitive`, Pacemaker on all cluster nodes must be able to contact Pacemaker Remote inside the bundle’s containers.

- The containers must have an accessible network (for example, `network` should not be set to “none” with a `primitive`).
- The default, using a distinct network space inside the container, works in combination with `ip-range-start`. Any firewall must allow access from all cluster nodes to the `control-port` on the container IPs.



- If the container shares the host's network space (for example, by setting `network` to "host"), a unique `control-port` should be specified for each bundle. Any firewall must allow access from all cluster nodes to the `control-port` on all cluster and remote node IPs.
- 

## Bundle Node Attributes

If the bundle has a `primitive`, the primitive's resource agent may want to set node attributes such as *promotion scores*. However, with containers, it is not apparent which node should get the attribute.

If the container uses shared storage that is the same no matter which node the container is hosted on, then it is appropriate to use the promotion score on the bundle node itself.

On the other hand, if the container uses storage exported from the underlying host, then it may be more appropriate to use the promotion score on the underlying host.

Since this depends on the particular situation, the `container-attribute-target` resource meta-attribute allows the user to specify which approach to use. If it is set to `host`, then user-defined node attributes will be checked on the underlying host. If it is anything else, the local node (in this case the bundle node) is used as usual.

This only applies to user-defined attributes; the cluster will always check the local node for cluster-defined attributes such as `#uname`.

If `container-attribute-target` is `host`, the cluster will pass additional environment variables to the primitive's resource agent that allow it to set node attributes appropriately: `CRM_meta_container_attribute_target` (identical to the meta-attribute value) and `CRM_meta_physical_host` (the name of the underlying host).

---

**Note:** When called by a resource agent, the `attrd_updater` and `crm_attribute` commands will automatically check those environment variables and set attributes appropriately.

---

## Bundle Meta-Attributes

Any meta-attribute set on a bundle will be inherited by the bundle's primitive and any resources implicitly created by Pacemaker for the bundle.

This includes options such as `priority`, `target-role`, and `is-managed`. See *Resource Options* for more information.

## Limitations of Bundles

Restarting pacemaker while a bundle is unmanaged or the cluster is in maintenance mode may cause the bundle to fail.

Bundles may not be explicitly cloned or included in groups. This includes the bundle's primitive and any resources implicitly created by Pacemaker for the bundle. (If `replicas` is greater than 1, the bundle will behave like a clone implicitly.)

Bundles do not have instance attributes, utilization attributes, or operations, though a bundle's primitive may have them.

A bundle with a primitive can run on a Pacemaker Remote node only if the bundle uses a distinct `control-port`.

## 2.11 Reusing Parts of the Configuration

Pacemaker provides multiple ways to simplify the configuration XML by reusing parts of it in multiple places.

Besides simplifying the XML, this also allows you to manipulate multiple configuration elements with a single reference.

### 2.11.1 Reusing Resource Definitions

If you want to create lots of resources with similar configurations, defining a *resource template* simplifies the task. Once defined, it can be referenced in primitives or in certain types of constraints.

#### Configuring Resources with Templates

The primitives referencing the template will inherit all meta-attributes, instance attributes, utilization attributes and operations defined in the template. And you can define specific attributes and operations for any of the primitives. If any of these are defined in both the template and the primitive, the values defined in the primitive will take precedence over the ones defined in the template.

Hence, resource templates help to reduce the amount of configuration work. If any changes are needed, they can be done to the template definition and will take effect globally in all resource definitions referencing that template.

Resource templates have a syntax similar to that of primitives.

#### Resource template for a migratable Xen virtual machine

```
<template id="vm-template" class="ocf" provider="heartbeat" type="Xen">
  <meta_attributes id="vm-template-meta_attributes">
    <nvpair id="vm-template-meta_attributes-allow-migrate" name="allow-migrate" value="true"/>
  </meta_attributes>
  <utilization id="vm-template-utilization">
    <nvpair id="vm-template-utilization-memory" name="memory" value="512"/>
  </utilization>
  <operations>
    <op id="vm-template-monitor-15s" interval="15s" name="monitor" timeout="60s"/>
    <op id="vm-template-start-0" interval="0" name="start" timeout="60s"/>
  </operations>
</template>
```

Once you define a resource template, you can use it in primitives by specifying the `template` property.

#### Xen primitive resource using a resource template

```
<primitive id="vm1" template="vm-template">
  <instance_attributes id="vm1-instance_attributes">
    <nvpair id="vm1-instance_attributes-name" name="name" value="vm1"/>
    <nvpair id="vm1-instance_attributes-xmfile" name="xmfile" value="/etc/xen/shared-vm/vm1"/>
  </instance_attributes>
</primitive>
```

In the example above, the new primitive `vm1` will inherit everything from `vm-template`. For example, the equivalent of the above two examples would be:

#### Equivalent Xen primitive resource not using a resource template

```
<primitive id="vm1" class="ocf" provider="heartbeat" type="Xen">
  <meta_attributes id="vm-template-meta_attributes">
    <nvpair id="vm-template-meta_attributes-allow-migrate" name="allow-migrate" value="true"/>
  </meta_attributes>
  <utilization id="vm-template-utilization">
    <nvpair id="vm-template-utilization-memory" name="memory" value="512"/>
  </utilization>
  <operations>
    <op id="vm-template-monitor-15s" interval="15s" name="monitor" timeout="60s"/>
    <op id="vm-template-start-0" interval="0" name="start" timeout="60s"/>
  </operations>
  <instance_attributes id="vm1-instance_attributes">
    <nvpair id="vm1-instance_attributes-name" name="name" value="vm1"/>
    <nvpair id="vm1-instance_attributes-xmfile" name="xmfile" value="/etc/xen/shared-vm/vm1"/>
  </instance_attributes>
</primitive>
```

If you want to overwrite some attributes or operations, add them to the particular primitive's definition.

#### Xen resource overriding template values

```
<primitive id="vm2" template="vm-template">
  <meta_attributes id="vm2-meta_attributes">
    <nvpair id="vm2-meta_attributes-allow-migrate" name="allow-migrate" value="false"/>
  </meta_attributes>
  <utilization id="vm2-utilization">
    <nvpair id="vm2-utilization-memory" name="memory" value="1024"/>
  </utilization>
  <instance_attributes id="vm2-instance_attributes">
    <nvpair id="vm2-instance_attributes-name" name="name" value="vm2"/>
    <nvpair id="vm2-instance_attributes-xmfile" name="xmfile" value="/etc/xen/shared-vm/vm2"/>
  </instance_attributes>
  <operations>
    <op id="vm2-monitor-30s" interval="30s" name="monitor" timeout="120s"/>
    <op id="vm2-stop-0" interval="0" name="stop" timeout="60s"/>
  </operations>
</primitive>
```

In the example above, the new primitive `vm2` has special attribute values. Its `monitor` operation has a longer timeout and interval, and the primitive has an additional `stop` operation.

To see the resulting definition of a resource, run:

```
# crm_resource --query-xml --resource vm2
```

To see the raw definition of a resource in the CIB, run:

```
# crm_resource --query-xml-raw --resource vm2
```

## Using Templates in Constraints

A resource template can be referenced in the following types of constraints:

- **order** constraints (see *Specifying the Order in which Resources Should Start/Stop*)
- **colocation** constraints (see *Placing Resources Relative to other Resources*)
- **rsc\_ticket** constraints (for multi-site clusters as described in *Configuring Ticket Dependencies*)

Resource templates referenced in constraints stand for all primitives which are derived from that template. This means, the constraint applies to all primitive resources referencing the resource template. Referencing resource templates in constraints is an alternative to resource sets and can simplify the cluster configuration considerably.

For example, given the example templates earlier in this chapter:

```
<rsc_colocation id="vm-template-colo-base-rsc" rsc="vm-template" rsc-role="Started" with-rsc="base-
↳rsc" score="INFINITY"/>
```

would colocate all VMs with **base-rsc** and is the equivalent of the following constraint configuration:

```
<rsc_colocation id="vm-colo-base-rsc" score="INFINITY">
  <resource_set id="vm-colo-base-rsc-0" sequential="false" role="Started">
    <resource_ref id="vm1"/>
    <resource_ref id="vm2"/>
  </resource_set>
  <resource_set id="vm-colo-base-rsc-1">
    <resource_ref id="base-rsc"/>
  </resource_set>
</rsc_colocation>
```

---

**Note:** In a colocation constraint, only one template may be referenced from either **rsc** or **with-rsc**; the other reference must be a regular resource.

---

## Using Templates in Resource Sets

Resource templates can also be referenced in resource sets.

For example, given the example templates earlier in this section, then:

```
<rsc_order id="order1" score="INFINITY">
  <resource_set id="order1-0">
    <resource_ref id="base-rsc"/>
    <resource_ref id="vm-template"/>
    <resource_ref id="top-rsc"/>
  </resource_set>
</rsc_order>
```

is the equivalent of the following constraint using a sequential resource set:

```
<rsc_order id="order1" score="INFINITY">
  <resource_set id="order1-0">
    <resource_ref id="base-rsc"/>
    <resource_ref id="vm1"/>
    <resource_ref id="vm2"/>
  </resource_set>
</rsc_order>
```

(continues on next page)

(continued from previous page)

```

    <resource_ref id="top-rsc"/>
  </resource_set>
</rsc_order>

```

Or, if the resources referencing the template can run in parallel, then:

```

<rsc_order id="order2" score="INFINITY">
  <resource_set id="order2-0">
    <resource_ref id="base-rsc"/>
  </resource_set>
  <resource_set id="order2-1" sequential="false">
    <resource_ref id="vm-template"/>
  </resource_set>
  <resource_set id="order2-2">
    <resource_ref id="top-rsc"/>
  </resource_set>
</rsc_order>

```

is the equivalent of the following constraint configuration:

```

<rsc_order id="order2" score="INFINITY">
  <resource_set id="order2-0">
    <resource_ref id="base-rsc"/>
  </resource_set>
  <resource_set id="order2-1" sequential="false">
    <resource_ref id="vm1"/>
    <resource_ref id="vm2"/>
  </resource_set>
  <resource_set id="order2-2">
    <resource_ref id="top-rsc"/>
  </resource_set>
</rsc_order>

```

### 2.11.2 Reusing Rules, Options and Sets of Operations

Sometimes a number of constraints need to use the same set of rules, and resources need to set the same options and parameters. To simplify this situation, you can refer to an existing object using an `id-ref` instead of an `id`.

So if for one resource you have

```

<rsc_location id="WebServer-connectivity" rsc="Webserver">
  <rule id="ping-prefer-rule" score-attribute="pingd" >
    <expression id="ping-prefer" attribute="pingd" operation="defined"/>
  </rule>
</rsc_location>

```

Then instead of duplicating the rule for all your other resources, you can instead specify:

#### Referencing rules from other constraints

```

<rsc_location id="WebDB-connectivity" rsc="WebDB">
  <rule id-ref="ping-prefer-rule"/>
</rsc_location>

```

**Important:** The cluster will insist that the `rule` exists somewhere. Attempting to add a reference to a non-existing rule will cause a validation failure, as will attempting to remove a `rule` that is referenced elsewhere.

The same principle applies for `meta_attributes` and `instance_attributes` as illustrated in the example below:

#### Referencing attributes, options, and operations from other resources

```
<primitive id="mySpecialRsc" class="ocf" type="Special" provider="me">
  <instance_attributes id="mySpecialRsc-attrs" score="1" >
    <nvpair id="default-interface" name="interface" value="eth0"/>
    <nvpair id="default-port" name="port" value="9999"/>
  </instance_attributes>
  <meta_attributes id="mySpecialRsc-options">
    <nvpair id="failure-timeout" name="failure-timeout" value="5m"/>
    <nvpair id="migration-threshold" name="migration-threshold" value="1"/>
    <nvpair id="stickiness" name="resource-stickiness" value="0"/>
  </meta_attributes>
  <operations id="health-checks">
    <op id="health-check" name="monitor" interval="60s"/>
    <op id="health-check" name="monitor" interval="30min"/>
  </operations>
</primitive>
<primitive id="myOther1Rsc" class="ocf" type="Other" provider="me">
  <instance_attributes id-ref="mySpecialRsc-attrs"/>
  <meta_attributes id-ref="mySpecialRsc-options"/>
  <operations id-ref="health-checks"/>
</primitive>
```

`id-ref` can similarly be used with `resource_set` (in any constraint type), `nvpair`, and `operations`.

### 2.11.3 Tagging Configuration Elements

Pacemaker allows you to *tag* any configuration element that has an XML ID.

The main purpose of tagging is to support higher-level user interface tools; Pacemaker itself only uses tags within constraints. Therefore, what you can do with tags mostly depends on the tools you use.

#### Configuring Tags

A tag is simply a named list of XML IDs.

#### Tag referencing three resources

```
<tags>
  <tag id="all-vms">
    <obj_ref id="vm1"/>
    <obj_ref id="vm2"/>
    <obj_ref id="vm3"/>
  </tag>
</tags>
```

What you can do with this new tag depends on what your higher-level tools support. For example, a tool might allow you to enable or disable all of the tagged resources at once, or show the status of just the tagged resources.

A single configuration element can be listed in any number of tags.

## Using Tags in Constraints and Resource Sets

Pacemaker itself only uses tags in constraints. If you supply a tag name instead of a resource name in any constraint, the constraint will apply to all resources listed in that tag.

### Constraint using a tag

```
<rsc_order id="order1" first="storage" then="all-vms" kind="Mandatory" />
```

In the example above, assuming the `all-vms` tag is defined as in the previous example, the constraint will behave the same as:

### Equivalent constraints without tags

```
<rsc_order id="order1-1" first="storage" then="vm1" kind="Mandatory" />
<rsc_order id="order1-2" first="storage" then="vm2" kind="Mandatory" />
<rsc_order id="order1-3" first="storage" then="vm3" kind="Mandatory" />
```

A tag may be used directly in the constraint, or indirectly by being listed in a *resource set* used in the constraint. When used in a resource set, an expanded tag will honor the set's `sequential` property.

## Filtering With Tags

The `crm_mon` tool can be used to display lots of information about the state of the cluster. On large or complicated clusters, this can include a lot of information, which makes it difficult to find the one thing you are interested in. The `--resource=` and `--node=` command line options can be used to filter results. In their most basic usage, these options take a single resource or node name. However, they can also be supplied with a tag name to display several objects at once.

For instance, given the following CIB section:

```
<resources>
  <primitive class="stonith" id="Fencing" type="fence_xvm"/>
  <primitive class="ocf" id="dummy" provider="pacemaker" type="Dummy"/>
  <group id="inactive-group">
    <primitive class="ocf" id="inactive-dummy-1" provider="pacemaker" type="Dummy"/>
    <primitive class="ocf" id="inactive-dummy-2" provider="pacemaker" type="Dummy"/>
  </group>
  <clone id="inactive-clone">
    <primitive id="inactive-dhcpd" class="lsb" type="dhcpd"/>
  </clone>
</resources>
<tags>
  <tag id="inactive-rscs">
```

(continues on next page)

(continued from previous page)

```
<obj_ref id="inactive-group"/>
<obj_ref id="inactive-clone"/>
</tag>
</tags>
```

The following would be output for `crm_mon --resource=inactive-rscs -r`:

```
Cluster Summary:
* Stack: corosync
* Current DC: cluster02 (version 2.0.4-1.e97f9675f.git.el7-e97f9675f) - partition with quorum
* Last updated: Tue Oct 20 16:09:01 2020
* Last change: Tue May 5 12:04:36 2020 by hacluster via crmd on cluster01
* 5 nodes configured
* 27 resource instances configured (4 DISABLED)

Node List:
* Online: [ cluster01 cluster02 ]

Full List of Resources:
* Clone Set: inactive-clone [inactive-dhcpd] (disabled):
  * Stopped (disabled): [ cluster01 cluster02 ]
* Resource Group: inactive-group (disabled):
  * inactive-dummy-1 (ocf::pacemaker:Dummy): Stopped (disabled)
  * inactive-dummy-2 (ocf::pacemaker:Dummy): Stopped (disabled)
```

## 2.12 Utilization and Placement Strategy

Pacemaker decides where to place a resource according to the resource allocation scores on every node. The resource will be allocated to the node where the resource has the highest score.

If the resource allocation scores on all the nodes are equal, by the default placement strategy, Pacemaker will choose a node with the least number of allocated resources for balancing the load. If the number of resources on each node is equal, the first eligible node listed in the CIB will be chosen to run the resource.

Often, in real-world situations, different resources use significantly different proportions of a node's capacities (memory, I/O, etc.). We cannot balance the load ideally just according to the number of resources allocated to a node. Besides, if resources are placed such that their combined requirements exceed the provided capacity, they may fail to start completely or run with degraded performance.

To take these factors into account, Pacemaker allows you to configure:

1. The capacity a certain node provides.
2. The capacity a certain resource requires.
3. An overall strategy for placement of resources.

### 2.12.1 Utilization attributes

To configure the capacity that a node provides or a resource requires, you can use *utilization attributes* in **node** and **resource** objects. You can name utilization attributes according to your preferences and define as many name/value pairs as your configuration needs. However, the attributes' values must be integers.



### Specifying CPU and RAM consumed by several resources

```
<primitive id="rsc-small" class="ocf" provider="pacemaker" type="Dummy">
  <utilization id="rsc-small-utilization">
    <nvpair id="rsc-small-utilization-cpu" name="cpu" value="1"/>
    <nvpair id="rsc-small-utilization-memory" name="memory" value="1024"/>
  </utilization>
</primitive>
<primitive id="rsc-medium" class="ocf" provider="pacemaker" type="Dummy">
  <utilization id="rsc-medium-utilization">
    <nvpair id="rsc-medium-utilization-cpu" name="cpu" value="2"/>
    <nvpair id="rsc-medium-utilization-memory" name="memory" value="2048"/>
  </utilization>
</primitive>
<primitive id="rsc-large" class="ocf" provider="pacemaker" type="Dummy">
  <utilization id="rsc-large-utilization">
    <nvpair id="rsc-large-utilization-cpu" name="cpu" value="3"/>
    <nvpair id="rsc-large-utilization-memory" name="memory" value="3072"/>
  </utilization>
</primitive>
```

A node is considered eligible for a resource if it has sufficient free capacity to satisfy the resource's requirements. The nature of the required or provided capacities is completely irrelevant to Pacemaker – it just makes sure that all capacity requirements of a resource are satisfied before placing a resource to a node.

**Note:** Utilization is supported for bundles (*since 2.1.3*), but only for bundles with an inner primitive. Any resource utilization values should be specified for the inner primitive, but any priority meta-attribute should be specified for the outer bundle.

## 2.12.2 Placement Strategy

After you have configured the capacities your nodes provide and the capacities your resources require, you need to set the `placement-strategy` in the global cluster options, otherwise the capacity configurations have *no effect*.

Four values are available for the `placement-strategy`:

- **default**

Utilization values are not taken into account at all. Resources are allocated according to allocation scores. If scores are equal, resources are evenly distributed across nodes.

- **utilization**

Utilization values are taken into account *only* when deciding whether a node is considered eligible (i.e. whether it has sufficient free capacity to satisfy the resource's requirements). Load-balancing is still done based on the number of resources allocated to a node.

- **balanced**

Utilization values are taken into account when deciding whether a node is eligible to serve a resource *and* when load-balancing, so an attempt is made to spread the resources in a way that optimizes resource performance.

- **minimal**

Utilization values are taken into account *only* when deciding whether a node is eligible to serve a resource. For load-balancing, an attempt is made to concentrate the resources on as few nodes as possible, thereby enabling possible power savings on the remaining nodes.

Set `placement-strategy` with `crm_attribute`:

```
# crm_attribute --name placement-strategy --update balanced
```

Now Pacemaker will ensure the load from your resources will be distributed evenly throughout the cluster, without the need for convoluted sets of colocation constraints.

### 2.12.3 Allocation Details

#### Which node is preferred to get consumed first when allocating resources?

- The node with the highest node weight gets consumed first. Node weight is a score maintained by the cluster to represent node health.
- If multiple nodes have the same node weight:
- If `placement-strategy` is `default` or `utilization`, the node that has the least number of allocated resources gets consumed first.
  - If their numbers of allocated resources are equal, the first eligible node listed in the CIB gets consumed first.
- If `placement-strategy` is `balanced`, the node that has the most free capacity gets consumed first.
  - If the free capacities of the nodes are equal, the node that has the least number of allocated resources gets consumed first.
    - \* If their numbers of allocated resources are equal, the first eligible node listed in the CIB gets consumed first.
- If `placement-strategy` is `minimal`, the first eligible node listed in the CIB gets consumed first.

#### Which node has more free capacity?

If only one type of utilization attribute has been defined, free capacity is a simple numeric comparison.

If multiple types of utilization attributes have been defined, then the node that is numerically highest in the the most attribute types has the most free capacity. For example:

- If `nodeA` has more free `cpus`, and `nodeB` has more free `memory`, then their free capacities are equal.
- If `nodeA` has more free `cpus`, while `nodeB` has more free `memory` and `storage`, then `nodeB` has more free capacity.

#### Which resource is preferred to be assigned first?

- The resource that has the highest `priority` (see *Resource Options*) gets allocated first.
- If their priorities are equal, check whether they are already running. The resource that has the highest score on the node where it's running gets allocated first, to prevent resource shuffling.
- If the scores above are equal or the resources are not running, the resource has the highest score on the preferred node gets allocated first.
- If the scores above are equal, the first runnable resource listed in the CIB gets allocated first.

### 2.12.4 Limitations and Workarounds

The type of problem Pacemaker is dealing with here is known as the [knapsack problem](#) and falls into the [NP-complete](#) category of computer science problems – a fancy way of saying “it takes a really long time to solve”.

Clearly in a HA cluster, it’s not acceptable to spend minutes, let alone hours or days, finding an optimal solution while services remain unavailable.

So instead of trying to solve the problem completely, Pacemaker uses a *best effort* algorithm for determining which node should host a particular service. This means it arrives at a solution much faster than traditional linear programming algorithms, but by doing so at the price of leaving some services stopped.

In the contrived example at the start of this chapter:

- `rsc-small` would be allocated to `node1`
- `rsc-medium` would be allocated to `node2`
- `rsc-large` would remain inactive

Which is not ideal.

There are various approaches to dealing with the limitations of pacemaker’s placement strategy:

- **Ensure you have sufficient physical capacity.**

It might sound obvious, but if the physical capacity of your nodes is (close to) maxed out by the cluster under normal conditions, then failover isn’t going to go well. Even without the utilization feature, you’ll start hitting timeouts and getting secondary failures.

- **Build some buffer into the capabilities advertised by the nodes.**

Advertise slightly more resources than we physically have, on the (usually valid) assumption that a resource will not use 100% of the configured amount of CPU, memory and so forth *all* the time. This practice is sometimes called *overcommit*.

- **Specify resource priorities.**

If the cluster is going to sacrifice services, it should be the ones you care about (comparatively) the least. Ensure that resource priorities are properly set so that your most important resources are scheduled first.

## 2.13 Access Control Lists (ACLs)

By default, the `root` user or any user in the `haclient` group can modify Pacemaker’s CIB without restriction. Pacemaker offers *access control lists (ACLs)* to provide more fine-grained authorization.

---

**Important:** Being able to modify the CIB’s resource section allows a user to run any executable file as root, by configuring it as an LSB resource with a full path.

---

### 2.13.1 ACL Prerequisites

In order to use ACLs:

- The `enable-acl` *cluster option* must be set to true.
- Desired users must have user accounts in the `haclient` group on all cluster nodes in the cluster.

- If your CIB was created before Pacemaker 1.1.12, it might need to be updated to the current schema (using `cibadmin --upgrade` or a higher-level tool equivalent) in order to use the syntax documented here.
- Prior to the 2.1.0 release, the Pacemaker software had to have been built with ACL support. If you are using an older release, your installation supports ACLs only if the output of the command `pacemakerd --features` contains `acls`. In newer versions, ACLs are always enabled.

### 2.13.2 ACL Configuration

ACLs are specified within an `acls` element of the CIB. The `acls` element may contain any number of `acl_role`, `acl_target`, and `acl_group` elements.

### 2.13.3 ACL Roles

An ACL *role* is a collection of permissions allowing or denying access to particular portions of the CIB. A role is configured with an `acl_role` element in the CIB `acls` section.

Table 39: Properties of an `acl_role` element

| Attribute   | Description                                    |
|-------------|--|
| id          | A unique name for the role ( <i>required</i> ) |
| description | Arbitrary text (not used by Pacemaker)         |

An `acl_role` element may contain any number of `acl_permission` elements.

Table 40: Properties of an `acl_permission` element

| Attribute   | Description   |
|-------------|---|
| id          | A unique name for the permission ( <i>required</i> )  |
| description | Arbitrary text (not used by Pacemaker)  |
| kind        | The access being granted. Allowed values are <code>read</code> , <code>write</code> , and <code>deny</code> . A value of <code>write</code> grants both read and write access.  |
| object-type | The name of an XML element in the CIB to which the permission applies. (Exactly one of <code>object-type</code> , <code>xpath</code> , and <code>reference</code> must be specified for a permission.)  |
| attribute   | If specified, the permission applies only to <code>object-type</code> elements that have this attribute set (to any value). If not specified, the permission applies to all <code>object-type</code> elements. May only be used with <code>object-type</code> .   |
| reference   | The ID of an XML element in the CIB to which the permission applies. (Exactly one of <code>object-type</code> , <code>xpath</code> , and <code>reference</code> must be specified for a permission.)  |
| xpath       | An <code>XPath</code> specification selecting an XML element in the CIB to which the permission applies. Attributes may be specified in the XPath to select particular elements, but the permissions apply to the entire element. (Exactly one of <code>object-type</code> , <code>xpath</code> , and <code>reference</code> must be specified for a permission.) |

**Important:**

- Permissions are applied to the selected XML element’s entire XML subtree (all elements enclosed within it).
- Write permission grants the ability to create, modify, or remove the element and its subtree, and also the ability to create any “scaffolding” elements (enclosing elements that do not have attributes other than an ID).
- Permissions for more specific matches (more deeply nested elements) take precedence over more general ones.
- If multiple permissions are configured for the same match (for example, in different roles applied to the same user), any **deny** permission takes precedence, then **write**, then lastly **read**.

### 2.13.4 ACL Targets and Groups

ACL targets correspond to user accounts on the system.

Table 41: **Properties of an `acl_target` element**

| Attribute | Description  |
|-----------|--|
| id        | The name of a user on the system ( <i>required</i> ) |

ACL groups may be specified, but are not currently used by Pacemaker. This is expected to change in a future version.

Table 42: **Properties of an `acl_group` element**

| Attribute | Description   |
|-----------|---|
| id        | The name of a group on the system ( <i>required</i> ) |

Each `acl_target` and `acl_group` element may contain any number of `role` elements.

Table 43: **Properties of a `role` element**

| Attribute | Description  |
|-----------|--|
| id        | The <code>id</code> of an <code>acl_role</code> element that specifies permissions granted to the enclosing target or group. |

**Important:** The `root` and `hacluster` user accounts always have full access to the CIB, regardless of ACLs. For all other user accounts, when `enable-acl` is true, permission to all parts of the CIB is denied by default (permissions must be explicitly granted).

### 2.13.5 ACL Examples

```
<acls>

<acl_role id="read_all">
  <acl_permission id="read_all-cib" kind="read" xpath="/cib" />
</acl_role>

<acl_role id="operator">

  <acl_permission id="operator-maintenance-mode" kind="write"
    xpath="//crm_config//nvpair[@name='maintenance-mode']" />

  <acl_permission id="operator-maintenance-attr" kind="write"
    xpath="//nvpair[@name='maintenance']" />

  <acl_permission id="operator-target-role" kind="write"
    xpath="//resources//meta_attributes/nvpair[@name='target-role']" />

  <acl_permission id="operator-is-managed" kind="write"
    xpath="//resources//nvpair[@name='is-managed']" />

  <acl_permission id="operator-rsc_location" kind="write"
    object-type="rsc_location" />

</acl_role>

<acl_role id="administrator">
  <acl_permission id="administrator-cib" kind="write" xpath="/cib" />
</acl_role>

<acl_role id="minimal">

  <acl_permission id="minimal-standby" kind="read"
    description="allow reading standby node attribute (permanent or transient)"
    xpath="//instance_attributes/nvpair[@name='standby']"/>

  <acl_permission id="minimal-maintenance" kind="read"
    description="allow reading maintenance node attribute (permanent or transient)"
    xpath="//nvpair[@name='maintenance']"/>

  <acl_permission id="minimal-target-role" kind="read"
    description="allow reading resource target roles"
    xpath="//resources//meta_attributes/nvpair[@name='target-role']"/>

  <acl_permission id="minimal-is-managed" kind="read"
    description="allow reading resource managed status"
    xpath="//resources//meta_attributes/nvpair[@name='is-managed']"/>

  <acl_permission id="minimal-deny-instance-attributes" kind="deny"
    xpath="//instance_attributes"/>

  <acl_permission id="minimal-deny-meta-attributes" kind="deny"
    xpath="//meta_attributes"/>

  <acl_permission id="minimal-deny-operations" kind="deny"
    xpath="//operations"/>

  <acl_permission id="minimal-deny-utilization" kind="deny"
```

(continues on next page)

(continued from previous page)

```

    xpath="//utilization"/>

    <acl_permission id="minimal-nodes" kind="read"
      description="allow reading node names/IDs (attributes are denied separately)"
      xpath="/cib/configuration/nodes"/>

    <acl_permission id="minimal-resources" kind="read"
      description="allow reading resource names/agents (parameters are denied separately)"
      xpath="/cib/configuration/resources"/>

    <acl_permission id="minimal-deny-constraints" kind="deny"
      xpath="/cib/configuration/constraints"/>

    <acl_permission id="minimal-deny-topology" kind="deny"
      xpath="/cib/configuration/fencing-topology"/>

    <acl_permission id="minimal-deny-op_defaults" kind="deny"
      xpath="/cib/configuration/op_defaults"/>

    <acl_permission id="minimal-deny-rsc_defaults" kind="deny"
      xpath="/cib/configuration/rsc_defaults"/>

    <acl_permission id="minimal-deny-alerts" kind="deny"
      xpath="/cib/configuration/alerts"/>

    <acl_permission id="minimal-deny-acls" kind="deny"
      xpath="/cib/configuration/acls"/>

    <acl_permission id="minimal-cib" kind="read"
      description="allow reading cib element and crm_config/status sections"
      xpath="/cib"/>

  </acl_role>

  <acl_target id="alice">
    <role id="minimal"/>
  </acl_target>

  <acl_target id="bob">
    <role id="read_all"/>
  </acl_target>

  <acl_target id="carol">
    <role id="read_all"/>
    <role id="operator"/>
  </acl_target>

  <acl_target id="dave">
    <role id="administrator"/>
  </acl_target>

</acls>

```

In the above example, the user **alice** has the minimal permissions necessary to run basic Pacemaker CLI tools, including using **crm\_mon** to view the cluster status, without being able to modify anything. The user **bob** can view the entire configuration and status of the cluster, but not make any changes. The user **carol** can read everything, and change selected cluster properties as well as resource roles and location constraints.

Finally, `dave` has full read and write access to the entire CIB.

Looking at the `minimal` role in more depth, it is designed to allow read access to the `cib` tag itself, while denying access to particular portions of its subtree (which is the entire CIB).

This is because the DC node is indicated in the `cib` tag, so `crm_mon` will not be able to report the DC otherwise. However, this does change the security model to allow by default, since any portions of the CIB not explicitly denied will be readable. The `cib` read access could be removed and replaced with read access to just the `crm_config` and `status` sections, for a safer approach at the cost of not seeing the DC in status output.

For a simpler configuration, the `minimal` role allows read access to the entire `crm_config` section, which contains cluster properties. It would be possible to allow read access to specific properties instead (such as `stonith-enabled`, `dc-uuid`, `have-quorum`, and `cluster-name`) to restrict access further while still allowing status output, but cluster properties are unlikely to be considered sensitive.

## 2.13.6 ACL Limitations

### Actions performed via IPC rather than the CIB

ACLs apply *only* to the CIB.

That means ACLs apply to command-line tools that operate by reading or writing the CIB, such as `crm_attribute` when managing permanent node attributes, `crm_mon`, and `cibadmin`.

However, command-line tools that communicate directly with Pacemaker daemons via IPC are not affected by ACLs. For example, users in the `haclient` group may still do the following, regardless of ACLs:

- Query transient node attribute values using `crm_attribute` and `attrd_updater`.
- Query basic node information using `crm_node`.
- Erase resource operation history using `crm_resource`.
- Query fencing configuration information, and execute fencing against nodes, using `stonith_admin`.

### ACLs and Pacemaker Remote

ACLs apply to commands run on Pacemaker Remote nodes using the Pacemaker Remote node's name as the ACL user name.

The idea is that Pacemaker Remote nodes (especially virtual machines and containers) are likely to be purpose-built and have different user accounts from full cluster nodes.

## 2.14 Status – Here be dragons

Most users never need to understand the contents of the status section and can be happy with the output from `crm_mon`.

However for those with a curious inclination, this section attempts to provide an overview of its contents.

### 2.14.1 Node Status

In addition to the cluster's configuration, the CIB holds an up-to-date representation of each cluster node in the `status` section.



**A bare-bones status entry for a healthy node cl-virt-1**

```
<node_state id="1" uname="cl-virt-1" in_ccm="true" crmd="online" crm-debug-origin="do_update_
→resource" join="member" expected="member">
  <transient_attributes id="1"/>
  <lrmd id="1"/>
</node_state>
```

Users are highly recommended *not* to modify any part of a node's state *directly*. The cluster will periodically regenerate the entire section from authoritative sources, so any changes should be done with the tools appropriate to those sources.

Table 44: Authoritative Sources for State Information

| CIB Object           | Authoritative Source |
|----------------------|----------------------|
| node_state           | pacemaker-controld   |
| transient_attributes | pacemaker-attd       |
| lrmd                 | pacemaker-execd      |

The fields used in the `node_state` objects are named as they are largely for historical reasons and are rooted in Pacemaker's origins as the resource manager for the older Heartbeat project. They have remained unchanged to preserve compatibility with older versions.

Table 45: Node Status Fields

| Field            | Description  |
|------------------|--|
| id               | Unique identifier for the node. Corosync-based clusters use a numeric counter.   |
| uname            | The node's name as known by the cluster  |
| in_ccm           | Is the node a member at the cluster communication layer? Allowed values: <b>true</b> , <b>false</b> .                          |
| crmd             | Is the node a member at the pacemaker layer? Allowed values: <b>online</b> , <b>offline</b> .                                  |
| crm-debug-origin | The name of the source function that made the most recent change (for debugging purposes).                                     |
| join             | Does the node participate in hosting resources? Allowed values: <b>down</b> , <b>pending</b> , <b>member</b> , <b>banned</b> . |
| expected         | Expected value for <code>join</code> .   |

The cluster uses these fields to determine whether, at the node level, the node is healthy or is in a failed state and needs to be fenced.

## 2.14.2 Transient Node Attributes

Like regular *Node Attributes*, the name/value pairs listed in the `transient_attributes` section help to describe the node. However they are forgotten by the cluster when the node goes offline. This can be useful, for instance, when you want a node to be in standby mode (not able to run resources) just until the next reboot.

In addition to any values the administrator sets, the cluster will also store information about failed resources here.

#### A set of transient node attributes for node cl-virt-1

```
<transient_attributes id="cl-virt-1">
  <instance_attributes id="status-cl-virt-1">
    <nvpair id="status-cl-virt-1-pingd" name="pingd" value="3"/>
    <nvpair id="status-cl-virt-1-probe_complete" name="probe_complete" value="true"/>
    <nvpair id="status-cl-virt-1-fail-count-pingd:0.monitor_30000" name="fail-count-pingd:0
    ↪#monitor_30000" value="1"/>
    <nvpair id="status-cl-virt-1-last-failure-pingd:0" name="last-failure-pingd:0" value=
    ↪"1239009742"/>
  </instance_attributes>
</transient_attributes>
```

In the above example, we can see that a monitor on the `pingd:0` resource has failed once, at 09:22:22 UTC 6 April 2009.<sup>1</sup>

We also see that the node is connected to three `pingd` peers and that all known resources have been checked for on this machine (`probe_complete`).

### 2.14.3 Operation History

A node's resource history is held in the `lrm_resources` tag (a child of the `lrm` tag). The information stored here includes enough information for the cluster to stop the resource safely if it is removed from the configuration section. Specifically, the resource's `id`, `class`, `type` and `provider` are stored.

#### A record of the `apcstonith` resource

```
<lrm_resource id="apcstonith" type="fence_apc_snmp" class="stonith"/>
```

Additionally, we store the last job for every combination of `resource`, `action` and `interval`. The concatenation of the values in this tuple are used to create the id of the `lrm_rsc_op` object.

Table 46: Contents of an `lrm_rsc_op` job

| Field                  | Description   |
|------------------------|---|
| <code>id</code>        | Identifier for the job constructed from the resource's <code>operation</code> and <code>interval</code> .   |
| <code>call-id</code>   | The job's ticket number. Used as a sort key to determine the order in which the jobs were executed.   |
| <code>operation</code> | The action the resource agent was invoked with.   |
| <code>interval</code>  | The frequency, in milliseconds, at which the operation will be repeated. A one-off job is indicated by 0.   |
| <code>op-status</code> | The job's status. Generally this will be either 0 (done) or -1 (pending). Rarely used in favor of <code>rc-code</code> .  |
| <code>rc-code</code>   | The job's result. Refer to the <i>Resource Agents</i> chapter of <i>Pacemaker Administration</i> for details on what the values here mean and how they are interpreted. |

Continued on next page

<sup>1</sup> You can use the standard `date` command to print a human-readable version of any seconds-since-epoch value, for example `date -d @1239009742`.

Table 46 – continued from previous page

| Field             | Description   |
|-------------------|---|
| last-rc-change    | Machine-local date/time, in seconds since epoch, at which the job first returned the current value of <code>rc-code</code> . For diagnostic purposes.   |
| exec-time         | Time, in milliseconds, that the job was running for. For diagnostic purposes.   |
| queue-time        | Time, in seconds, that the job was queued for in the local executor. For diagnostic purposes.   |
| crm_feature_set   | The version which this job description conforms to. Used when processing <code>op-digest</code> .   |
| transition-key    | A concatenation of the job's graph action number, the graph number, the expected result and the UUID of the controller instance that scheduled it. This is used to construct <code>transition-magic</code> (below).   |
| transition-magic  | A concatenation of the job's <code>op-status</code> , <code>rc-code</code> and <code>transition-key</code> . Guaranteed to be unique for the life of the cluster (which ensures it is part of CIB update notifications) and contains all the information needed for the controller to correctly analyze and process the completed job. Most importantly, the decomposed elements tell the controller if the job entry was expected and whether it failed. |
| op-digest         | An MD5 sum representing the parameters passed to the job. Used to detect changes to the configuration, to restart resources if necessary.   |
| crm-debug-origins | The origin of the current values. For diagnostic purposes.  |

### Simple Operation History Example

A monitor operation (determines current state of the `apcstonith` resource)

```
<lrms_resource id="apcstonith" type="fence_apc_snmp" class="stonith">
  <lrms_rsc_op id="apcstonith_monitor_0" operation="monitor" call-id="2"
    rc-code="7" op-status="0" interval="0"
    crm-debug-origins="do_update_resource" crm_feature_set="3.0.1"
    op-digest="2e3da9274d3550dc6526fb24bfcba0"
    transition-key="22:2:7:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    transition-magic="0:7;22:2:7:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239008085" exec-time="10" queue-time="0"/>
</lrms_resource>
```

In the above example, the job is a non-recurring monitor operation often referred to as a “probe” for the `apcstonith` resource.

The cluster schedules probes for every configured resource on a node when the node first starts, in order to determine the resource's current state before it takes any further action.

From the `transition-key`, we can see that this was the 22nd action of the 2nd graph produced by this instance of the controller (2668bbeb-06d5-40f9-936d-24cb7f87006a).

The third field of the `transition-key` contains a 7, which indicates that the job expects to find the resource inactive. By looking at the `rc-code` property, we see that this was the case.

As that is the only job recorded for this node, we can conclude that the cluster started the resource elsewhere.

### Complex Operation History Example

## Resource history of a pingd clone with multiple jobs

```
<lrms_resource id="pingd:0" type="pingd" class="ocf" provider="pacemaker">
  <lrms_rsc_op id="pingd:0_monitor_30000" operation="monitor" call-id="34"
    rc-code="0" op-status="0" interval="30000"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1"
    transition-key="10:11:0:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239009741" exec-time="10" queue-time="0"/>
  <lrms_rsc_op id="pingd:0_stop_0" operation="stop"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1" call-id="32"
    rc-code="0" op-status="0" interval="0"
    transition-key="11:11:0:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239009741" exec-time="10" queue-time="0"/>
  <lrms_rsc_op id="pingd:0_start_0" operation="start" call-id="33"
    rc-code="0" op-status="0" interval="0"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1"
    transition-key="31:11:0:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239009741" exec-time="10" queue-time="0" />
  <lrms_rsc_op id="pingd:0_monitor_0" operation="monitor" call-id="3"
    rc-code="0" op-status="0" interval="0"
    crm-debug-orig="do_update_resource" crm_feature_set="3.0.1"
    transition-key="23:2:7:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-rc-change="1239008085" exec-time="20" queue-time="0"/>
</lrms_resource>
```

When more than one job record exists, it is important to first sort them by `call-id` before interpreting them.

Once sorted, the above example can be summarized as:

1. A non-recurring monitor operation returning 7 (not running), with a `call-id` of 3
2. A stop operation returning 0 (success), with a `call-id` of 32
3. A start operation returning 0 (success), with a `call-id` of 33
4. A recurring monitor returning 0 (success), with a `call-id` of 34

The cluster processes each job record to build up a picture of the resource's state. After the first and second entries, it is considered stopped, and after the third it considered active.

Based on the last operation, we can tell that the resource is currently active.

Additionally, from the presence of a `stop` operation with a lower `call-id` than that of the `start` operation, we can conclude that the resource has been restarted. Specifically this occurred as part of actions 11 and 31 of transition 11 from the controller instance with the key `2668bbeb...`. This information can be helpful for locating the relevant section of the logs when looking for the source of a failure.

## 2.15 Multi-Site Clusters and Tickets

Apart from local clusters, Pacemaker also supports multi-site clusters. That means you can have multiple, geographically dispersed sites, each with a local cluster. Failover between these clusters can be coordinated manually by the administrator, or automatically by a higher-level entity called a *Cluster Ticket Registry* (*CTR*).

### 2.15.1 Challenges for Multi-Site Clusters

Typically, multi-site environments are too far apart to support synchronous communication and data replication between the sites. That leads to significant challenges:

- How do we make sure that a cluster site is up and running?
- How do we make sure that resources are only started once?
- How do we make sure that quorum can be reached between the different sites and a split-brain scenario avoided?
- How do we manage failover between sites?
- How do we deal with high latency in case of resources that need to be stopped?

In the following sections, learn how to meet these challenges.

### 2.15.2 Conceptual Overview

Multi-site clusters can be considered as “overlay” clusters where each cluster site corresponds to a cluster node in a traditional cluster. The overlay cluster can be managed by a CTR in order to guarantee that any cluster resource will be active on no more than one cluster site. This is achieved by using *tickets* that are treated as failover domain between cluster sites, in case a site should be down.

The following sections explain the individual components and mechanisms that were introduced for multi-site clusters in more detail.

#### Ticket

Tickets are, essentially, cluster-wide attributes. A ticket grants the right to run certain resources on a specific cluster site. Resources can be bound to a certain ticket by `rsc_ticket` constraints. Only if the ticket is available at a site can the respective resources be started there. Vice versa, if the ticket is revoked, the resources depending on that ticket must be stopped.

The ticket thus is similar to a *site quorum*, i.e. the permission to manage/own resources associated with that site. (One can also think of the current `have-quorum` flag as a special, cluster-wide ticket that is granted in case of node majority.)

Tickets can be granted and revoked either manually by administrators (which could be the default for classic enterprise clusters), or via the automated CTR mechanism described below.

A ticket can only be owned by one site at a time. Initially, none of the sites has a ticket. Each ticket must be granted once by the cluster administrator.

The presence or absence of tickets for a site is stored in the CIB as a cluster status. With regards to a certain ticket, there are only two states for a site: `true` (the site has the ticket) or `false` (the site does not have the ticket). The absence of a certain ticket (during the initial state of the multi-site cluster) is the same as the value `false`.

#### Dead Man Dependency

A site can only activate resources safely if it can be sure that the other site has deactivated them. However after a ticket is revoked, it can take a long time until all resources depending on that ticket are stopped “cleanly”, especially in case of cascaded resources. To cut that process short, the concept of a *Dead Man Dependency* was introduced.

If a dead man dependency is in force, if a ticket is revoked from a site, the nodes that are hosting dependent resources are fenced. This considerably speeds up the recovery process of the cluster and makes sure that resources can be migrated more quickly.

This can be configured by specifying a `loss-policy="fence"` in `rsc_ticket` constraints.

### Cluster Ticket Registry

A CTR is a coordinated group of network daemons that automatically handles granting, revoking, and timing out tickets (instead of the administrator revoking the ticket somewhere, waiting for everything to stop, and then granting it on the desired site).

Pacemaker does not implement its own CTR, but interoperates with external software designed for that purpose (similar to how resource and fencing agents are not directly part of pacemaker).

Participating clusters run the CTR daemons, which connect to each other, exchange information about their connectivity, and vote on which sites gets which tickets.

A ticket is granted to a site only once the CTR is sure that the ticket has been relinquished by the previous owner, implemented via a timer in most scenarios. If a site loses connection to its peers, its tickets time out and recovery occurs. After the connection timeout plus the recovery timeout has passed, the other sites are allowed to re-acquire the ticket and start the resources again.

This can also be thought of as a “quorum server”, except that it is not a single quorum ticket, but several.

### Configuration Replication

As usual, the CIB is synchronized within each cluster, but it is *not* synchronized across cluster sites of a multi-site cluster. You have to configure the resources that will be highly available across the multi-site cluster for every site accordingly.

## 2.15.3 Configuring Ticket Dependencies

The `rsc_ticket` constraint lets you specify the resources depending on a certain ticket. Together with the constraint, you can set a **loss-policy** that defines what should happen to the respective resources if the ticket is revoked.

The attribute **loss-policy** can have the following values:

- **fence**: Fence the nodes that are running the relevant resources.
- **stop**: Stop the relevant resources.
- **freeze**: Do nothing to the relevant resources.
- **demote**: Demote relevant resources that are running in the promoted role.

#### Constraint that fences node if ticketA is revoked

```
<rsc_ticket id="rsc1-req-ticketA" rsc="rsc1" ticket="ticketA" loss-policy="fence"/>
```

The example above creates a constraint with the ID `rsc1-req-ticketA`. It defines that the resource `rsc1` depends on `ticketA` and that the node running the resource should be fenced if `ticketA` is revoked.

If resource `rsc1` were a promotable resource, you might want to configure that only being in the promoted role depends on `ticketA`. With the following configuration, `rsc1` will be demoted if `ticketA` is revoked:

**Constraint that demotes rsc1 if ticketA is revoked**

```
<rsc_ticket id="rsc1-req-ticketA" rsc="rsc1" rsc-role="Promoted" ticket="ticketA" loss-policy=
  ↪ "demote"/>
```

You can create multiple `rsc_ticket` constraints to let multiple resources depend on the same ticket. However, `rsc_ticket` also supports resource sets (see *Resource Sets*), so one can easily list all the resources in one `rsc_ticket` constraint instead.

**Ticket constraint for multiple resources**

```
<rsc_ticket id="resources-dep-ticketA" ticket="ticketA" loss-policy="fence">
  <resource_set id="resources-dep-ticketA-0" role="Started">
    <resource_ref id="rsc1"/>
    <resource_ref id="group1"/>
    <resource_ref id="clone1"/>
  </resource_set>
  <resource_set id="resources-dep-ticketA-1" role="Promoted">
    <resource_ref id="ms1"/>
  </resource_set>
</rsc_ticket>
```

In the example above, there are two resource sets, so we can list resources with different roles in a single `rsc_ticket` constraint. There's no dependency between the two resource sets, and there's no dependency among the resources within a resource set. Each of the resources just depends on `ticketA`.

Referencing resource templates in `rsc_ticket` constraints, and even referencing them within resource sets, is also supported.

If you want other resources to depend on further tickets, create as many constraints as necessary with `rsc_ticket`.

## 2.15.4 Managing Multi-Site Clusters

### Granting and Revoking Tickets Manually

You can grant tickets to sites or revoke them from sites manually. If you want to re-distribute a ticket, you should wait for the dependent resources to stop cleanly at the previous site before you grant the ticket to the new site.

Use the `crm_ticket` command line tool to grant and revoke tickets.

To grant a ticket to this site:

```
# crm_ticket --ticket ticketA --grant
```

To revoke a ticket from this site:

```
# crm_ticket --ticket ticketA --revoke
```

**Important:** If you are managing tickets manually, use the `crm_ticket` command with great care, because

it cannot check whether the same ticket is already granted elsewhere.

---

### Granting and Revoking Tickets via a Cluster Ticket Registry

We will use [Booth](#) here as an example of software that can be used with pacemaker as a Cluster Ticket Registry. Booth implements the [Raft](#) algorithm to guarantee the distributed consensus among different cluster sites, and manages the ticket distribution (and thus the failover process between sites).

Each of the participating clusters and *arbitrators* runs the Booth daemon **boothd**.

An *arbitrator* is the multi-site equivalent of a quorum-only node in a local cluster. If you have a setup with an even number of sites, you need an additional instance to reach consensus about decisions such as failover of resources across sites. In this case, add one or more arbitrators running at additional sites. Arbitrators are single machines that run a booth instance in a special mode. An arbitrator is especially important for a two-site scenario, otherwise there is no way for one site to distinguish between a network failure between it and the other site, and a failure of the other site.

The most common multi-site scenario is probably a multi-site cluster with two sites and a single arbitrator on a third site. However, technically, there are no limitations with regards to the number of sites and the number of arbitrators involved.

**Boothd** at each site connects to its peers running at the other sites and exchanges connectivity details. Once a ticket is granted to a site, the booth mechanism will manage the ticket automatically: If the site which holds the ticket is out of service, the booth daemons will vote which of the other sites will get the ticket. To protect against brief connection failures, sites that lose the vote (either explicitly or implicitly by being disconnected from the voting body) need to relinquish the ticket after a time-out. Thus, it is made sure that a ticket will only be re-distributed after it has been relinquished by the previous site. The resources that depend on that ticket will fail over to the new site holding the ticket. The nodes that have run the resources before will be treated according to the **loss-policy** you set within the **rsc\_ticket** constraint.

Before the booth can manage a certain ticket within the multi-site cluster, you initially need to grant it to a site manually via the **booth** command-line tool. After you have initially granted a ticket to a site, **boothd** will take over and manage the ticket automatically.

---

**Important:** The **booth** command-line tool can be used to grant, list, or revoke tickets and can be run on any machine where **boothd** is running. If you are managing tickets via Booth, use only **booth** for manual intervention, not **crm\_ticket**. That ensures the same ticket will only be owned by one cluster site at a time.

---

### Booth Requirements

- All clusters that will be part of the multi-site cluster must be based on Pacemaker.
- Booth must be installed on all cluster nodes and on all arbitrators that will be part of the multi-site cluster.
- Nodes belonging to the same cluster site should be synchronized via NTP. However, time synchronization is not required between the individual cluster sites.

### General Management of Tickets

Display the information of tickets:



```
# crm_ticket --info
```

Or you can monitor them with:

```
# crm_mon --tickets
```

Display the `rsc_ticket` constraints that apply to a ticket:

```
# crm_ticket --ticket ticketA --constraints
```

When you want to do maintenance or manual switch-over of a ticket, revoking the ticket would trigger the loss policies. If `loss-policy="fence"`, the dependent resources could not be gracefully stopped/demoted, and other unrelated resources could even be affected.

The proper way is making the ticket *standby* first with:

```
# crm_ticket --ticket ticketA --standby
```

Then the dependent resources will be stopped or demoted gracefully without triggering the loss policies.

If you have finished the maintenance and want to activate the ticket again, you can run:

```
# crm_ticket --ticket ticketA --activate
```

### 2.15.5 For more information

- [SUSE's Geo Clustering quick start](#)
- [Booth](#)

## 2.16 Sample Configurations

### 2.16.1 Empty

#### An Empty Configuration

```
<cib crm_feature_set="3.0.7" validate-with="pacemaker-1.2" admin_epoch="1" epoch="0" num_updates=
  ↪"0">
  <configuration>
    <crm_config/>
    <nodes/>
    <resources/>
    <constraints/>
  </configuration>
  <status/>
</cib>
```

### 2.16.2 Simple

## A simple configuration with two nodes, some cluster options and a resource

```

<cib crm_feature_set="3.0.7" validate-with="pacemaker-1.2" admin_epoch="1" epoch="0" num_updates=
  ↪ "0">
  <configuration>
    <crm_config>
      <cluster_property_set id="cib-bootstrap-options">
        <nvpair id="option-1" name="symmetric-cluster" value="true"/>
        <nvpair id="option-2" name="no-quorum-policy" value="stop"/>
        <nvpair id="option-3" name="stonith-enabled" value="0"/>
      </cluster_property_set>
    </crm_config>
    <nodes>
      <node id="xxx" uname="c001n01" type="normal"/>
      <node id="yyy" uname="c001n02" type="normal"/>
    </nodes>
    <resources>
      <primitive id="myAddr" class="ocf" provider="heartbeat" type="IPAddr">
        <operations>
          <op id="myAddr-monitor" name="monitor" interval="300s"/>
        </operations>
        <instance_attributes id="myAddr-params">
          <nvpair id="myAddr-ip" name="ip" value="192.0.2.10"/>
        </instance_attributes>
      </primitive>
    </resources>
    <constraints>
      <rsc_location id="myAddr-prefer" rsc="myAddr" node="c001n01" score="INFINITY"/>
    </constraints>
    <rsc_defaults>
      <meta_attributes id="rsc_defaults-options">
        <nvpair id="rsc-default-1" name="resource-stickiness" value="100"/>
        <nvpair id="rsc-default-2" name="migration-threshold" value="10"/>
      </meta_attributes>
    </rsc_defaults>
    <op_defaults>
      <meta_attributes id="op_defaults-options">
        <nvpair id="op-default-1" name="timeout" value="30s"/>
      </meta_attributes>
    </op_defaults>
  </configuration>
  <status/>
</cib>

```

In the above example, we have one resource (an IP address) that we check every five minutes and will run on host c001n01 until either the resource fails 10 times or the host shuts down.

### 2.16.3 Advanced Configuration

## An advanced configuration with groups, clones and STONITH

```

<cib crm_feature_set="3.0.7" validate-with="pacemaker-1.2" admin_epoch="1" epoch="0" num_updates=
↳ "0">
  <configuration>
    <crm_config>
      <cluster_property_set id="cib-bootstrap-options">
        <nvpair id="option-1" name="symmetric-cluster" value="true"/>
        <nvpair id="option-2" name="no-quorum-policy" value="stop"/>
        <nvpair id="option-3" name="stonith-enabled" value="true"/>
      </cluster_property_set>
    </crm_config>
    <nodes>
      <node id="xxx" uname="c001n01" type="normal"/>
      <node id="yyy" uname="c001n02" type="normal"/>
      <node id="zzz" uname="c001n03" type="normal"/>
    </nodes>
    <resources>
      <primitive id="myAddr" class="ocf" provider="heartbeat" type="IPaddr">
        <operations>
          <op id="myAddr-monitor" name="monitor" interval="300s"/>
        </operations>
        <instance_attributes id="myAddr-attrs">
          <nvpair id="myAddr-attr-1" name="ip" value="192.0.2.10"/>
        </instance_attributes>
      </primitive>
      <group id="myGroup">
        <primitive id="database" class="lsb" type="oracle">
          <operations>
            <op id="database-monitor" name="monitor" interval="300s"/>
          </operations>
        </primitive>
        <primitive id="webserver" class="lsb" type="apache">
          <operations>
            <op id="webserver-monitor" name="monitor" interval="300s"/>
          </operations>
        </primitive>
      </group>
      <clone id="STONITH">
        <meta_attributes id="stonith-options">
          <nvpair id="stonith-option-1" name="globally-unique" value="false"/>
        </meta_attributes>
        <primitive id="stonithclone" class="stonith" type="external/ssh">
          <operations>
            <op id="stonith-op-mon" name="monitor" interval="5s"/>
          </operations>
          <instance_attributes id="stonith-attrs">
            <nvpair id="stonith-attr-1" name="hostlist" value="c001n01,c001n02"/>
          </instance_attributes>
        </primitive>
      </clone>
    </resources>
    <constraints>
      <rsc_location id="myAddr-prefer" rsc="myAddr" node="c001n01"
        score="INFINITY"/>
      <rsc_colocation id="group-with-ip" rsc="myGroup" with-rsc="myAddr"
        score="INFINITY"/>
    </constraints>
    <op_defaults>
      <meta_attributes id="op_defaults-options">
        <nvpair id="op-default-1" name="timeout" value="30s"/>
      </meta_attributes>
    </op_defaults>
    <rsc_defaults>
      <meta_attributes id="rsc_defaults-options">
        <nvpair id="rsc-default-1" name="resource-stickiness" value="100"/>
        <nvpair id="rsc-default-2" name="migration-threshold" value="10"/>
      </meta_attributes>
    </rsc_defaults>
  </configuration>
</cib>

```



## INDEX

- `genindex`
- `search`



## Symbols

#digests

node attribute, 21

#node-unfenced

node attribute, 21

## A

Access Control List (ACL), 125

acl\_group, 127

acl\_permission, 126

acl\_role, 126

acl\_target, 127

acls, 126

role, 127

acl\_group

id (attribute), 127

XML element, 127

acl\_permission

attribute (attribute), 126

description (attribute), 126

id (attribute), 126

kind (attribute), 126

object-type (attribute), 126

reference (attribute), 126

XML element, 126

xpath (attribute), 126

acl\_role

description (attribute), 126

id (attribute), 126

XML element, 126

acl\_target

id (attribute), 127

XML element, 127

acls

XML element, 126

action

property, enabled, 35

property, id, 34

property, interval, 34

property, name, 34

property, on-fail, 35

property, record-pending, 35

property, role, 35

property, timeout, 34

resource\_set attribute, 46

status, call-id, 132

status, crm-debug-origin, 133

status, crm\_feature\_set, 133

status, exec-time, 133

status, id, 132

status, interval, 132

status, last-rc-change, 133

status, op-digest, 133

status, op-status, 132

status, operation, 132

status, queue-time, 133

status, rc-code, 132

status, transition-key, 133

status, transition-magic, 133

add-host

network attribute, 112

admin\_epoch

cib, 13

agent

alert, 71

alert, 71

agent, 71

environment variables, 75

filters, 73

instance attributes, 72

meta-attribute, timeout, 72

meta-attribute, timestamp-format, 72

meta-attributes, 72

recipient, 71

sample agents, 74

XML element, 71

alerts

XML element, 71

Asymmetrical Clusters, 41

attribute

acl\_permission attribute, 126

action (resource\_set), 46

add-host (network), 112

attribute (acl\_permission), 126

- control-port (network), 112
- description (acl\_permission), 126
- description (acl\_role), 126
- description (bundle), 110
- expression, 77
- first (rsc\_order), 42
- first-action (rsc\_order), 42
- host-interface (network), 112
- host-netmask (network), 112
- id (acl\_group), 127
- id (acl\_permission), 126
- id (acl\_role), 126
- id (acl\_target), 127
- id (bundle), 110
- id (port-mapping), 113
- id (resource\_set), 46
- id (role), 127
- id (rsc\_colocation), 43
- id (rsc\_location), 40
- id (rsc\_order), 42
- id (storage-mapping), 113
- image (docker), 111
- image (podman), 111
- image (rkt), 111
- influence (rsc\_colocation), 44
- internal-port (port-mapping), 113
- ip-range-start (network), 112
- kind (acl\_permission), 126
- kind (rsc\_order), 42
- network (docker), 111
- network (podman), 111
- network (rkt), 111
- node (rsc\_location), 40
- node-attribute (rsc\_colocation), 44
- object-type (acl\_permission), 126
- options (docker), 111
- options (podman), 111
- options (rkt), 111
- options (storage-mapping), 113
- port (port-mapping), 113
- promoted-max (docker), 111
- promoted-max (podman), 111
- promoted-max (rkt), 111
- range (port-mapping), 113
- reference (acl\_permission), 126
- replicas (docker), 111
- replicas (podman), 111
- replicas (rkt), 111
- replicas-per-host (docker), 111
- replicas-per-host (podman), 111
- replicas-per-host (rkt), 111
- require-all (resource\_set), 46
- resource-discovery (rsc\_location), 40
- role (resource\_set), 46

- rsc (rsc\_colocation), 44
- rsc (rsc\_location), 40
- rsc-pattern (rsc\_location), 40
- run-command (docker), 111
- run-command (podman), 111
- run-command (rkt), 111
- score (resource\_set), 46
- score (rsc\_colocation), 44
- score (rsc\_location), 40
- sequential (resource\_set), 46
- source-dir (storage-mapping), 113
- source-dir-root (storage-mapping), 113
- symmetrical (rsc\_order), 43
- target-dir (storage-mapping), 113
- then (rsc\_order), 42
- then-action (rsc\_order), 42
- with-rsc (rsc\_colocation), 44
- XML element, 73
- xpath (acl\_permission), 126

## B

- batch-limit
  - cluster option, 15
- boolean-op
  - rule, 77
- bundle
  - attribute, description, 110
  - attribute, id, 110
  - networking, 111
  - XML element, 110

## C

- call-id
  - action status, 132
- cib
  - admin\_epoch, 13
  - cib-last-written, 14
  - dc-uuid, 14
  - epoch, 13
  - have-quorum, 14
  - num\_updates, 13
  - validate-with, 13
  - XML element, 12
- cib-last-written
  - cib, 14
- class
  - resource, 27
  - rsc\_expression, 83
- clone, 98
  - constraint, 101
  - environment variables, 104, 106
  - option, clone-max, 99
  - option, clone-min, 99
  - option, clone-node-max, 99



- option, globally-unique, 99
- option, interleave, 100
- option, notify, 100
- option, ordered, 100
- option, promotable, 100
- option, promoted-max, 100
- option, promoted-node-max, 100
- options, 99
- ordering constraint, rsc-role, 44
- ordering constraint, with-rsc-role, 44
- property, id, 99
- resource-stickiness, 103
- XML element, 99
- clone-max
  - clone option, 99
- clone-min
  - clone option, 99
- clone-node-max
  - clone option, 99
- cluster option
  - batch-limit, 15
  - cluster-delay, 17
  - cluster-infrastructure, 14
  - cluster-ipc-limit, 17
  - cluster-name, 14
  - cluster-recheck-interval, 18
  - concurrent-fencing, 16
  - dc-deadtime, 17
  - dc-version, 14
  - election-timeout, 18
  - enable-acl, 17
  - enable-startup-probes, 15
  - fence-reaction, 16
  - join-finalization-timeout, 19
  - join-integration-timeout, 19
  - maintenance-mode, 15
  - migration-limit, 15
  - no-quorum-policy, 15
  - node-health-base, 17
  - node-health-green, 17
  - node-health-red, 18
  - node-health-strategy, 17, 22
  - node-health-yellow, 18
  - pe-error-series-max, 17
  - pe-input-series-max, 17
  - pe-warn-series-max, 17
  - placement-strategy, 17
  - priority-fencing-delay, 17
  - remove-after-stop, 18
  - rule, 85, 88
  - shutdown-escalation, 19
  - shutdown-lock, 18
  - shutdown-lock-limit, 18
  - start-failure-is-fatal, 15
  - startup-fencing, 18
  - stonith-action, 16
  - stonith-enabled, 16
  - stonith-max-attempts, 16
  - stonith-timeout, 16
  - stonith-watchdog-timeout, 16
  - stop-all-resources, 15
  - stop-orphan-actions, 15
  - stop-orphan-resources, 15
  - symmetric-cluster, 15
  - transition-delay, 19
- cluster-delay
  - cluster option, 17
- cluster-infrastructure
  - cluster option, 14
- cluster-ipc-limit
  - cluster option, 17
- cluster-name
  - cluster option, 14
- cluster-recheck-interval
  - cluster option, 18
- colocation, 43
- concurrent-fencing
  - cluster option, 16
- configuration
  - XML element, 12
- constraint, 39
  - colocation, 43
  - location, 39
  - ordering, 42
  - resource set, 46
  - rsc\_colocation, 43
  - rsc\_location, 39
  - rsc\_order, 42
- control-port
  - network attribute, 112
- critical
  - resource option, 28
- crm-debug-origin
  - action status, 133
  - node status, 131
- CRM\_alert\_attribute\_name, 75
- CRM\_alert\_attribute\_value, 75
- CRM\_alert\_desc, 75
- CRM\_alert\_exec\_time, 75
- CRM\_alert\_interval, 75
- CRM\_alert\_kind, 75
- CRM\_alert\_node, 75
- CRM\_alert\_nodeid, 75
- CRM\_alert\_rc, 75
- CRM\_alert\_recipient, 75
- CRM\_alert\_rsc, 75
- CRM\_alert\_sequence, 75
- CRM\_alert\_status, 75

- CRM\_alert\_target\_rc, 75
- CRM\_alert\_task, 75
- CRM\_alert\_timestamp, 75
- CRM\_alert\_timestamp\_epoch, 75
- CRM\_alert\_timestamp\_usec, 75
- CRM\_alert\_version, 75
- crm\_feature\_set
  - action status, 133
- crmd
  - node status, 131
- custom
  - node-health-strategy value, 22

## D

- dampen
  - ocf:pacemaker:ping resource parameter, 93
- date specification, 79
- date\_expression
  - end, 79
  - id, 79
  - operation, 79
  - start, 79
  - XML element, 78
- date\_spec
  - hours, 80
  - id, 80
  - monthdays, 80
  - months, 80
  - moon, 80
  - weekdays, 80
  - weeks, 80
  - weekyears, 80
  - XML element, 79
  - yeardays, 80
  - years, 80
- dc-deadtime
  - cluster option, 17
- dc-uuid
  - cib, 14
- dc-version
  - cluster option, 14
- description
  - acl\_permission attribute, 126
  - acl\_role attribute, 126
  - bundle attribute, 110
- devices
  - fencing-level, 68
- docker
  - attribute, image, 111
  - attribute, network, 111
  - attribute, options, 111
  - attribute, promoted-max, 111
  - attribute, replicas, 111
  - attribute, replicas-per-host, 111

- attribute, run-command, 111
- XML element, 111
- duration, 80
  - hours, 81
  - id, 81
  - minutes, 81
  - months, 81
  - seconds, 81
  - weeks, 81
  - XML element, 80
  - years, 81

## E

- election-timeout
  - cluster option, 18
- enable-acl
  - cluster option, 17
- enable-startup-probes
  - cluster option, 15
- enabled
  - action property, 35
- end
  - date\_expression, 79
- environment variable
  - alert agents, 75
  - CRM\_alert\_attribute\_name, 75
  - CRM\_alert\_attribute\_value, 75
  - CRM\_alert\_desc, 75
  - CRM\_alert\_exec\_time, 75
  - CRM\_alert\_interval, 75
  - CRM\_alert\_kind, 75
  - CRM\_alert\_node, 75
  - CRM\_alert\_nodeid, 75
  - CRM\_alert\_rc, 75
  - CRM\_alert\_recipient, 75
  - CRM\_alert\_rsc, 75
  - CRM\_alert\_sequence, 75
  - CRM\_alert\_status, 75
  - CRM\_alert\_target\_rc, 75
  - CRM\_alert\_task, 75
  - CRM\_alert\_timestamp, 75
  - CRM\_alert\_timestamp\_epoch, 75
  - CRM\_alert\_timestamp\_usec, 75
  - CRM\_alert\_version, 75
  - OCF\_RESKEY\_CRM\_meta\_notify\_active\_resource, 104
  - OCF\_RESKEY\_CRM\_meta\_notify\_active\_uname, 104
  - OCF\_RESKEY\_CRM\_meta\_notify\_demote\_resource, 106
  - OCF\_RESKEY\_CRM\_meta\_notify\_demote\_uname, 106
  - OCF\_RESKEY\_CRM\_meta\_notify\_inactive\_resource, 104

- OCF\_RESKEY\_CRM\_meta\_notify\_operation, why necessary, 53
  - 104 fencing-level, 67
  - OCF\_RESKEY\_CRM\_meta\_notify\_promote\_resource, 68
  - 106 id, 68
  - OCF\_RESKEY\_CRM\_meta\_notify\_promote\_uname, 68
  - 106 index, 68
  - OCF\_RESKEY\_CRM\_meta\_notify\_promoted\_resource, 68
  - 106 target, 68
  - OCF\_RESKEY\_CRM\_meta\_notify\_promoted\_uname, 68
  - 106 target-attribute, 68
  - OCF\_RESKEY\_CRM\_meta\_notify\_promoted\_uname, 68
  - 106 target-pattern, 68
  - OCF\_RESKEY\_CRM\_meta\_notify\_promoted\_uname, 68
  - 106 target-value, 68
  - OCF\_RESKEY\_CRM\_meta\_notify\_start\_resource, 67
  - 104 fencing-topology, 67
  - OCF\_RESKEY\_CRM\_meta\_notify\_start\_uname, 42
  - 104 first, rsc\_order attribute, 42
  - OCF\_RESKEY\_CRM\_meta\_notify\_start\_uname, 42
  - 104 first-action, rsc\_order attribute, 42
  - OCF\_RESKEY\_CRM\_meta\_notify\_stop\_resource, 42
  - 104 G
  - OCF\_RESKEY\_CRM\_meta\_notify\_stop\_uname, 99
  - 104 globally-unique, clone option, 99
  - OCF\_RESKEY\_CRM\_meta\_notify\_type, 99
  - 104 green, node health attribute value, 22
  - OCF\_RESKEY\_CRM\_meta\_notify\_unpromoted\_resource, 98
  - 106 group, option, ordered, 98
  - OCF\_RESKEY\_CRM\_meta\_notify\_unpromoted\_uname, 97
  - 106 id, 97
  - OCF\_RESKEY\_CRM\_meta\_notify\_unpromoted\_uname, 98
  - 106 property, id, 97
  - epoch, resource-stickiness, 98
  - cib, 13 XML element, 97
  - exec-time
  - action status, 133
  - expected
  - node status, 131
  - expression
  - attribute, 77
  - id, 77
  - operation, 78
  - type, 77
  - value, 78
  - value-source, 78
  - XML element, 77
- ## F
- fail-count
  - node attribute, 21
  - failure-timeout
  - resource option, 29
  - fence-reaction
  - cluster option, 16
  - fencing, 53
  - agent, 54
  - alert, 71
  - configuration, 61
  - device, 54
  - special instance attributes, 55
  - topology, 67
  - unfencing, 59
- ## H
- have-quorum
  - cib, 14
  - host-interface
  - network attribute, 112
  - host-netmask
  - network attribute, 112
  - host\_list
  - ocf:pacemaker:ping resource parameter, 93
  - hours
  - date\_spec, 80
  - duration, 81
- ## I
- id
  - acl\_group attribute, 127
  - acl\_permission attribute, 126
  - acl\_role attribute, 126
  - acl\_target attribute, 127
  - action property, 34
  - action status, 132
  - bundle attribute, 110
  - clone property, 99
  - date\_expression, 79
  - date\_spec, 80
  - duration, 81
  - expression, 77
  - fencing-level, 68

- group property, 97
- op\_expression, 83
- port-mapping attribute, 113
- resource, 27
- resource\_set attribute, 46
- role attribute, 127
- rsc\_colocation attribute, 43
- rsc\_expression, 83
- rsc\_location attribute, 40
- rsc\_order attribute, 42
- rule, 77
- storage-mapping attribute, 113

image

- docker attribute, 111
- podman attribute, 111
- rkt attribute, 111

in\_ccm

- node status, 131

index

- fencing-level, 68

influence

- rsc\_colocation attribute, 44

instance attribute

- alert instance attributes, 72
- rule, 85

interleave

- clone option, 100

internal-port

- port-mapping attribute, 113

interval

- action property, 34
- action status, 132
- interval-origin, 89
- op\_expression, 83

interval-origin

- operation attribute, 89

ip-range-start

- network attribute, 112

is-managed

- resource option, 28

## J

join

- node status, 131

join-finalization-timeout

- cluster option, 19

join-integration-timeout

- cluster option, 19

## K

kind

- acl\_permission attribute, 126
- rsc\_order attribute, 42

## L

last-failure

- node attribute, 21

last-rc-change

- action status, 133

Linux Standard Base

- resources, 25

location constraint, 39

- rule, 84

LSB

- resources, 25

## M

maintenance

- node attribute, 21
- resource option, 28

maintenance-mode

- cluster option, 15

meta-attribute

- alert meta-attributes, 72
- rule, 85
- timeout (alert), 72
- timestamp-format (alert), 72

migrate-on-red

- node-health-strategy value, 22

migration-limit

- cluster option, 15

migration-threshold

- resource meta-attribute, 90
- resource option, 29

minutes

- duration, 81

monthdays

- date\_spec, 80

months

- date\_spec, 80
- duration, 81

moon

- date\_spec, 80

multiple-active

- resource option, 30

multiplier

- ocf:pacemaker:ping resource parameter, 93

## N

Nagios Plugins

- resources, 26

name

- action property, 34
- op\_expression, 83

network

- attribute
  - control-port, 112
  - host-interface, 112

- host-netmask, 112
- attribute, add-host, 112
- attribute, ip-range-start, 112
- docker attribute, 111
- podman attribute, 111
- rkt attribute, 111
- XML element, 111
- no-quorum-policy
  - cluster option, 15
- node
  - alert, 71
  - attribute, 19
  - health, 21
  - rsc\_location attribute, 40
  - score, 39
  - standby mode, 91
  - status, 130
  - status, crm-debug-origin, 131
  - status, crmd, 131
  - status, expected, 131
  - status, in\_ccm, 131
  - status, join, 131
  - status, uname, 131
- node attribute, 19
  - #digests, 21
  - #node-unfenced, 21
  - fail-count, 21
  - health, 22
  - health (green), 22
  - health (red), 22
  - health (score), 22
  - health (yellow), 22
  - last-failure, 21
  - maintenance, 21
  - probe\_complete, 21
  - resource-discovery-enabled, 21
  - rule, 85
  - rule expression, 77
  - shutdown, 21
  - site-name, 21
  - standby, 21
  - terminate, 21
- node-attribute
  - rsc\_colocation attribute, 44
- node-health-base
  - cluster option, 17
- node-health-green
  - cluster option, 17
- node-health-red
  - cluster option, 18
- node-health-strategy
  - cluster option, 17, 22
  - custom, 22
  - migrate-on-red, 22

- none, 22
- only-green, 22
- progressive, 22
- node-health-yellow
  - cluster option, 18
- none
  - node-health-strategy value, 22
- notify
  - clone option, 100
  - environment variables, 104
- num\_updates
  - cib, 13

## O

- object-type
  - acl\_permission attribute, 126
- OCF
  - resources, 24
- OCF return code
  - OCF\_FAILED\_PROMOTED, 104
  - OCF\_NOT\_RUNNING, 104
  - OCF\_RUNNING\_PROMOTED, 104
  - OCF\_SUCCESS, 104
- ocf:pacemaker:ping resource, 92
  - dampen parameter, 93
  - host\_list parameter, 93
  - multiplier parameter, 93
- OCF\_FAILED\_PROMOTED, 104
- OCF\_NOT\_RUNNING, 104
- OCF\_RESKEY\_CRM\_meta\_notify\_active\_resource, 104
- OCF\_RESKEY\_CRM\_meta\_notify\_active\_uname, 104
- OCF\_RESKEY\_CRM\_meta\_notify\_demote\_resource, 106
- OCF\_RESKEY\_CRM\_meta\_notify\_demote\_uname, 106
- OCF\_RESKEY\_CRM\_meta\_notify\_inactive\_resource, 104
- OCF\_RESKEY\_CRM\_meta\_notify\_operation, 104
- OCF\_RESKEY\_CRM\_meta\_notify\_promote\_resource, 106
- OCF\_RESKEY\_CRM\_meta\_notify\_promote\_uname, 106
- OCF\_RESKEY\_CRM\_meta\_notify\_promoted\_resource, 106
- OCF\_RESKEY\_CRM\_meta\_notify\_promoted\_uname, 106
- OCF\_RESKEY\_CRM\_meta\_notify\_start\_resource, 104
- OCF\_RESKEY\_CRM\_meta\_notify\_start\_uname, 104

- OCF\_RESKEY\_CRM\_meta\_notify\_stop\_resource,  
104
  - OCF\_RESKEY\_CRM\_meta\_notify\_stop\_uname,  
104
  - OCF\_RESKEY\_CRM\_meta\_notify\_type, 104
  - OCF\_RESKEY\_CRM\_meta\_notify\_unpromoted\_resource,  
106
  - OCF\_RESKEY\_CRM\_meta\_notify\_unpromoted\_uname,  
106
  - OCF\_RUNNING\_PROMOTED, 104
  - OCF\_SUCCESS, 104
  - on-fail
    - action property, 35
  - only-green
    - node-health-strategy value, 22
  - op-digest
    - action status, 133
  - op-status
    - action status, 132
  - op\_expression
    - id, 83
    - interval, 83
    - name, 83
    - XML element, 83
  - Open Cluster Framework
    - resources, 24
  - operation
    - action status, 132
    - date\_expression, 79
    - expression, 78
    - failure count, 89
    - failure recovery, 89
    - interval-origin, 89
    - rule expression, 83
    - start-delay, 89
  - operation defaults
    - rule, 85
  - Operation History, 132
  - Opt-In Clusters, 41
  - Opt-Out Clusters, 41
  - option
    - clone-max (clone), 99
    - clone-min (clone), 99
    - clone-node-max (clone), 99
    - globally-unique (clone), 99
    - interleave (clone), 100
    - notify (clone), 100
    - ordered (clone), 100
    - ordered (group), 98
    - promotable (clone), 100
    - promoted-max (clone), 100
    - promoted-node-max (clone), 100
  - options
    - clone, 99
    - docker attribute, 111
    - podman attribute, 111
    - rkt attribute, 111
    - storage-mapping attribute, 113
  - ordered
    - clone option, 100
    - group option, 98
  - ordering constraint
    - rsc-role (clone), 44
    - with-rsc-role (clone), 44
- ## P
- pcmk\_action\_limit, 56
  - pcmk\_delay\_base, 56
  - pcmk\_delay\_max, 56
  - pcmk\_host\_argument, 57
  - pcmk\_host\_check, 56
  - pcmk\_host\_list, 56
  - pcmk\_host\_map, 55
  - pcmk\_list\_action, 58
  - pcmk\_list\_retries, 58
  - pcmk\_list\_timeout, 58
  - pcmk\_monitor\_action, 58
  - pcmk\_monitor\_retries, 58
  - pcmk\_monitor\_timeout, 58
  - pcmk\_off\_action, 57
  - pcmk\_off\_retries, 57
  - pcmk\_off\_timeout, 57
  - pcmk\_reboot\_action, 57
  - pcmk\_reboot\_retries, 57
  - pcmk\_reboot\_timeout, 57
  - pcmk\_status\_action, 58
  - pcmk\_status\_retries, 59
  - pcmk\_status\_timeout, 59
  - pe-error-series-max
    - cluster option, 17
  - pe-input-series-max
    - cluster option, 17
  - pe-warn-series-max
    - cluster option, 17
  - ping resource, 92
  - placement-strategy
    - cluster option, 17
  - podman
    - attribute, image, 111
    - attribute, network, 111
    - attribute, options, 111
    - attribute, promoted-max, 111
    - attribute, replicas, 111
    - attribute, replicas-per-host, 111
    - attribute, run-command, 111
    - XML element, 111
  - port
    - port-mapping attribute, 113

- port-mapping
  - attribute, id, 113
  - attribute, internal-port, 113
  - attribute, port, 113
  - attribute, range, 113
  - XML element, 112
- priority
  - resource option, 28
- priority-fencing-delay
  - cluster option, 17
- probe\_complete
  - node attribute, 21
- progressive
  - node-health-strategy value, 22
- promotable
  - clone option, 100
  - environment variables, 106
- promotable clone, 99
  - constraint, 101
- promoted-max
  - clone option, 100
  - docker attribute, 111
  - podman attribute, 111
  - rkt attribute, 111
- promoted-node-max
  - clone option, 100
- property
  - id (clone), 99
  - id (group), 97
- provider
  - resource, 27
  - rsc\_expression, 83
- provides, 55
- Q**
- queue-time
  - action status, 133
- R**
- range
  - port-mapping attribute, 113
- rc-code
  - action status, 132
- recipient
  - XML element, 71
- record-pending
  - action property, 35
- red
  - node health attribute value, 22
- reference
  - acl\_permission attribute, 126
- reload, 96
- reload-agent, 96
- remove-after-stop
  - cluster option, 18
- replicas
  - docker attribute, 111
  - podman attribute, 111
  - rkt attribute, 111
- replicas-per-host
  - docker attribute, 111
  - podman attribute, 111
  - rkt attribute, 111
- require-all
  - resource\_set attribute, 46
- requires
  - resource option, 29
- Resource
  - Nagios Plugins, 26
  - STONITH, 26
  - System Services, 26
  - Systemd, 25
  - Upstart, 25
- resource, 24
  - action, 34
  - alert, 71
  - bundle, 111, 114, 115
  - class, 24
  - clone, 98
  - constraint, 39
  - failure count, 89
  - failure recovery, 89
  - location relative to other resources, 43
  - LSB, 25
  - migration-threshold, 90
  - move, 90
  - OCF, 24
  - operation, 34
  - option, critical, 28
  - option, failure-timeout, 29
  - option, is-managed, 28
  - option, maintenance, 28
  - option, migration-threshold, 29
  - option, multiple-active, 30
  - option, priority, 28
  - option, requires, 29
  - option, resource-stickiness, 28
  - option, target-role, 28
  - promotable, 99
  - property, class, 27
  - property, id, 27
  - property, provider, 27
  - property, type, 27
  - resource set, 46
  - rule expression, 82
  - score, 39
  - start order, 42
- resource defaults

- rule, 85
- resource-discovery
  - rsc\_location attribute, 40
- resource-discovery-enabled
  - node attribute, 21
- resource-stickiness
  - clone, 103
  - group, 98
  - resource option, 28
- resource\_set
  - attribute, action, 46
  - attribute, id, 46
  - attribute, require-all, 46
  - attribute, role, 46
  - attribute, score, 46
  - attribute, sequential, 46
  - XML element, 46
- return code, 104
- rkt
  - attribute, image, 111
  - attribute, network, 111
  - attribute, options, 111
  - attribute, promoted-max, 111
  - attribute, replicas, 111
  - attribute, replicas-per-host, 111
  - attribute, run-command, 111
  - XML element, 111
- role
  - action property, 35
  - id (attribute), 127
  - resource\_set attribute, 46
  - rule, 77
  - XML element, 127
- rsc
  - rsc\_colocation attribute, 44
  - rsc\_location attribute, 40
- rsc-pattern
  - rsc\_location attribute, 40
- rsc-role
  - clone ordering constraint, 44
- rsc\_colocation
  - attribute, id, 43
  - attribute, influence, 44
  - attribute, node-attribute, 44
  - attribute, rsc, 44
  - attribute, score, 44
  - attribute, with-rsc, 44
  - XML element, 43
- rsc\_expression
  - class, 83
  - id, 83
  - provider, 83
  - type, 83
  - XML element, 82

- rsc\_location
  - attribute, id, 40
  - attribute, node, 40
  - attribute, resource-discovery, 40
  - attribute, rsc, 40
  - attribute, rsc-pattern, 40
  - attribute, score, 40
  - XML element, 39
- rsc\_order
  - attribute, first, 42
  - attribute, first-action, 42
  - attribute, id, 42
  - attribute, kind, 42
  - attribute, symmetrical, 43
  - attribute, then, 42
  - attribute, then-action, 42
  - constraint, 42
  - XML element, 42
- rule, 76
  - boolean-op, 77
  - cluster option, 85, 88
  - date/time expression, 78
  - id, 77
  - instance attribute, 85
  - location constraint, 84
  - meta-attribute, 85
  - node attribute, 85
  - node attribute expression, 77
  - operation defaults, 85
  - operation expression, 83
  - resource defaults, 85
  - resource expression, 82
  - role, 77
  - score, 77
  - score-attribute, 77
  - XML element, 76
- run-command
  - docker attribute, 111
  - podman attribute, 111
  - rkt attribute, 111

## S

- score
  - node health attribute value, 22
  - resource\_set attribute, 46
  - rsc\_colocation attribute, 44
  - rsc\_location attribute, 40
  - rule, 77
- score-attribute
  - rule, 77
- seconds
  - duration, 81
- select
  - XML element, 73



- select\_attributes
  - XML element, 73
- select\_fencing
  - XML element, 73
- select\_nodes
  - XML element, 73
- select\_resources
  - XML element, 73
- sequential
  - resource\_set attribute, 46
- shutdown
  - node attribute, 21
- shutdown-escalation
  - cluster option, 19
- shutdown-lock
  - cluster option, 18
- shutdown-lock-limit
  - cluster option, 18
- site-name
  - node attribute, 21
- source-dir
  - storage-mapping attribute, 113
- source-dir-root
  - storage-mapping attribute, 113
- standby
  - node attribute, 21
- standby mode, 91
- start
  - date\_expression, 79
- start-delay
  - operation attribute, 89
- start-failure-is-fatal
  - cluster option, 15
- startup-fencing
  - cluster option, 18
- status, 130
- STONITH, 53
  - resources, 26
- stonith-action
  - cluster option, 16
- stonith-enabled
  - cluster option, 16
- stonith-max-attempts
  - cluster option, 16
- stonith-timeout, 55
  - cluster option, 16
- stonith-watchdog-timeout
  - cluster option, 16
- stop-all-resources
  - cluster option, 15
- stop-orphan-actions
  - cluster option, 15
- stop-orphan-resources
  - cluster option, 15

- storage-mapping
  - attribute, id, 113
  - attribute, options, 113
  - attribute, source-dir, 113
  - attribute, source-dir-root, 113
  - attribute, target-dir, 113
- symmetric-cluster
  - cluster option, 15
- symmetrical
  - rsc\_order attribute, 43
- Symmetrical Clusters, 41
- System Service
  - resources, 26
- Systemd
  - resources, 25

## T

- target
  - fencing-level, 68
- target-attribute
  - fencing-level, 68
- target-dir
  - storage-mapping attribute, 113
- target-pattern
  - fencing-level, 68
- target-role
  - resource option, 28
- target-value
  - fencing-level, 68
- terminate
  - node attribute, 21
- then
  - rsc\_order attribute, 42
- timeout
  - action property, 34
  - alert meta-attribute, 72
- timestamp-format
  - alert meta-attribute, 72
- transition-delay
  - cluster option, 19
- transition-key
  - action status, 133
- transition-magic
  - action status, 133
- type
  - expression, 77
  - resource, 27
  - rsc\_expression, 83

## U

- uname
  - node status, 131
- unfencing, 59
- Upstart

resources, 25

## V

validate-with

cib, 13

value

expression, 78

value-source

expression, 78

## W

weekdays

date\_spec, 80

weeks

date\_spec, 80

duration, 81

weekyears

date\_spec, 80

with-rsc

rsc\_colocation attribute, 44

with-rsc-role

clone ordering constraint, 44

## X

XML element

acl\_group, 127

acl\_permission, 126

acl\_role, 126

acl\_target, 127

acls, 126

alert, 71

alerts, 71

attribute, 73

bundle, 110

cib, 12

clone, 99

configuration, 12

date\_expression, 78

date\_spec, 79

docker, 111

duration, 80

expression, 77

group, 97

network, 111

op\_expression, 83

podman, 111

port-mapping, 112

recipient, 71

resource\_set, 46

rkt, 111

role, 127

rsc\_colocation, 43

rsc\_expression, 82

rsc\_location, 39

rsc\_order, 42

rule, 76

select, 73

select\_attributes, 73

select\_fencing, 73

select\_nodes, 73

select\_resources, 73

XML element, status, 130

xpath

acl\_permission attribute, 126

## Y

yeardays

date\_spec, 80

years

date\_spec, 80

duration, 81

yellow

node health attribute value, 22