# **Configuration Explained**

| An | A-Z quide | to P | acemaker's | Confid | guration | <b>Options</b> |
|----|-----------|------|------------|--------|----------|----------------|
|    |           |      |            |        |          |                |

, Written by the Pacemaker project contributors

#### Configuration Explained: An A-Z guide to Pacemaker's Configuration Options

by

#### Abstract

The purpose of this document is to definitively explain the concepts used to configure Pacemaker. To achieve this, it will focus exclusively on the XML syntax used to configure Pacemaker's Cluster Information Base (CIB).

Copyright © 2009-2020 The Pacemaker project contributors.

The text of and illustrations in this document are licensed under version 4.0 or later of the Creative Commons Attribution-ShareAlike International Public License ("CC-BY-SA")<sup>1</sup>.

In accordance with CC-BY-SA, if you distribute this document or an adaptation of it, you must provide the URL for the original version.

In addition to the requirements of this license, the following activities are looked upon favorably:

- 1. If you are distributing Open Publication works on hardcopy or CD-ROM, you provide email notification to the authors of your intent to redistribute at least thirty days before your manuscript or media freeze, to give the authors time to provide updated documents. This notification should describe modifications, if any, made to the document.
- 2. All substantive modifications (including deletions) be either clearly marked up in the document or else described in an attachment to the document.
- 3. Finally, while it is not mandatory under this license, it is considered good form to offer a free copy of any hardcopy or CD-ROM expression of the author(s) work.

<sup>&</sup>lt;sup>1</sup> An explanation of CC-BY-SA is available at https://creativecommons.org/licenses/by-sa/4.0/

# **Table of Contents**

| Preface  | . xi |
|--|------|
| Document Conventions                                 | . xi |
| Typographic Conventions                              | . xi |
| Pull-quote Conventions                               | xii  |
| Notes and Warnings                                   | xiv  |
| We Need Feedback!                                    |      |
| 1. Read-Me-First                                     | 1    |
| The Scope of this Document                           | 1    |
| What Is Pacemaker?                                   |      |
| Cluster Architecture                                 |      |
| Pacemaker Architecture                               |      |
| Node Redundancy Designs                              |      |
| 2. Cluster-Wide Configuration                        |      |
| Configuration Layout                                 |      |
| CIB Properties                                       |      |
| Cluster Options                                      |      |
| 3. Cluster Nodes                                     |      |
| Defining a Cluster Node                              |      |
| Where Pacemaker Gets the Node Name                   |      |
|  |      |
| Node Attributes                                      |      |
| Setting and querying node attributes                 |      |
| Special node attributes                              |      |
| 4. Cluster Resources                                 |      |
| What is a Cluster Resource?                          |      |
| Resource Classes                                     |      |
| Open Cluster Framework                               |      |
| Linux Standard Base                                  |      |
| Systemd  |      |
| Upstart  |      |
| System Services                                      |      |
| STONITH  |      |
| Nagios Plugins                                       |      |
| Resource Properties                                  |      |
| Resource Options                                     | . 19 |
| Resource Meta-Attributes                             | . 19 |
| Setting Global Defaults for Resource Meta-Attributes | . 22 |
| Resource Instance Attributes                         | . 22 |
| Resource Operations                                  | . 23 |
| Operation Properties                                 | . 24 |
| Monitoring Resources for Failure                     |      |
| Monitoring Resources When Administration is Disabled |      |
| Setting Global Defaults for Operations               |      |
| When Implicit Operations Take a Long Time            |      |
| Multiple Monitor Operations                          |      |
| Disabling a Monitor Operation                        |      |
| 5. Resource Constraints                              |      |
| Scores   |      |
| Infinity Math  |      |
| Deciding Which Nodes a Resource Can Run On           |      |
| Location Properties                                  |      |
| Asymmetrical "Opt-In" Clusters                       |      |
|  |      |

| Symmetrical "Opt-Out" Clusters                            |    |
|---|----|
| What if Two Nodes Have the Same Score                     |    |
| Specifying the Order in which Resources Should Start/Stop | 32 |
| Ordering Properties                                       |    |
| Optional and mandatory ordering                           | 34 |
| Placing Resources Relative to other Resources             |    |
| Colocation Properties                                     |    |
| Mandatory Placement                                       |    |
| Advisory Placement  |    |
| Colocation by Node Attribute                              |    |
| Resource Sets   |    |
| Ordering Sets of Resources                                |    |
| Ordered Set   |    |
| Ordering Multiple Sets                                    |    |
|   |    |
| Resource Set OR Logic                                     |    |
| Colocating Sets of Resources                              |    |
| 6. Fencing  |    |
| What Is Fencing?  |    |
| Why Is Fencing Necessary?                                 |    |
| Fence Devices   |    |
| Fence Agents  |    |
| When a Fence Device Can Be Used                           |    |
| Limitations of Fencing Resources                          |    |
| Special Options for Fencing Resources                     |    |
| Unfencing   |    |
| Fence Devices Dependent on Other Resources                | 48 |
| Configuring Fencing                                       | 49 |
| Example Fencing Configuration                             | 50 |
| Fencing Topologies  | 52 |
| Example Dual-Layer, Dual-Device Fencing Topologies        | 53 |
| Remapping Reboots   | 58 |
| 7. Alerts   | 59 |
| Alert Agents  |    |
| Alert Recipients  |    |
| Alert Meta-Attributes                                     |    |
| Alert Instance Attributes                                 |    |
| Alert Filters   |    |
| Using the Sample Alert Agents                             |    |
| Writing an Alert Agent                                    |    |
| 8. Rules  |    |
| Rule Properties   |    |
| Node Attribute Expressions                                |    |
| Date/Time Expressions                                     |    |
| *   |    |
| Date Specifications                                       |    |
| Durations   |    |
| Example Time-Based Expressions                            |    |
| Resource Expressions                                      |    |
| Example Resource-Based Expressions                        |    |
| Operation Expressions                                     |    |
| Example Operation-Based Expressions                       |    |
| Using Rules to Determine Resource Location                |    |
| Location Rules Based on Other Node Properties             |    |
| Using score-attribute Instead of score                    |    |
| Using Rules to Define Options                             |    |

| Using Rules to Control Resource Options                    |      |
|--|------|
| Using Rules to Control Resource Defaults                   |      |
| Using Rules to Control Cluster Options                     |      |
| 9. Advanced Configuration                                  |      |
| Specifying When Recurring Actions are Performed            | . 78 |
| Handling Resource Failure                                  | . 78 |
| Failure Counts   | . 78 |
| Failure Response   | . 79 |
| Moving Resources   | . 80 |
| Moving Resources Manually                                  | . 80 |
| Moving Resources Due to Connectivity Changes               | . 82 |
| Migrating Resources  |      |
| Tracking Node Health                                       |      |
| Node Health Attributes                                     |      |
| Node Health Strategy                                       |      |
| Measuring Node Health                                      |      |
| Reloading Services After a Definition Change               |      |
| 10. Advanced Resource Types                                |      |
| Groups - A Syntactic Shortcut                              |      |
| Group Properties   |      |
| Group Options  |      |
| Group Instance Attributes                                  |      |
| Group Contents   |      |
| Group Constraints  |      |
| Group Stickiness   |      |
| Clones - Resources That Can Have Multiple Active Instances |      |
| Anonymous versus Unique Clones                             |      |
| Promotable clones  |      |
| Clone Properties   |      |
| •  |      |
| Clone Options  |      |
|  |      |
| Clone Instance Attributes                                  |      |
|  |      |
| Clone Stickiness   |      |
| Clone Resource Agent Requirements                          |      |
| Monitoring Promotable Clone Resources                      |      |
| Determining Which Instance is Promoted                     |      |
| Bundles - Isolated Environments                            |      |
| Bundle Prerequisites                                       |      |
| Bundle Properties  |      |
| Bundle Container Properties                                |      |
| Bundle Network Properties                                  |      |
| Bundle Storage Properties                                  |      |
| Bundle Primitive   |      |
| Bundle Node Attributes                                     |      |
| Bundle Meta-Attributes                                     | 109  |
| Limitations of Bundles                                     |      |
| 11. Reusing Parts of the Configuration                     |      |
| Reusing Resource Definitions                               |      |
| Configuring Resources with Templates                       |      |
| Using Templates in Constraints                             |      |
| Using Templates in Resource Sets                           |      |
| Reusing Rules, Options and Sets of Operations              |      |
| Tagging Configuration Elements                             | 114  |

## Configuration Explained

| Configuring Tags   | 115  |
|--|------|
| Using Tags in Constraints and Resource Sets                              | 115  |
| 12. Utilization and Placement Strategy                                   | 116  |
| Utilization attributes   | 116  |
| Placement Strategy   | 117  |
| Allocation Details   | 118  |
| Which node is preferred to get consumed first when allocating resources? | 118  |
| Which node has more free capacity?                                       | 118  |
| Which resource is preferred to be assigned first?                        |      |
| Limitations and Workarounds  | 119  |
| 13. Access Control Lists (ACLs)  | 120  |
| ACL Prerequisites  | 120  |
| ACL Configuration  | 120  |
| ACL Roles  | 120  |
| ACL Targets and Groups   | 121  |
| ACL Examples   | 122  |
| 14. Status — Here be dragons   | 125  |
| Node Status  | 125  |
| Transient Node Attributes  | 126  |
| Operation History  | 126  |
| Simple Operation History Example   | 128  |
| Complex Operation History Example  | 129  |
| 15. Multi-Site Clusters and Tickets                                      | 131  |
| Challenges for Multi-Site Clusters                                       | 131  |
| Conceptual Overview  | 131  |
| Ticket   | 132  |
| Dead Man Dependency  | 132  |
| Cluster Ticket Registry  | 132  |
| Configuration Replication  | 133  |
| Configuring Ticket Dependencies  | 133  |
| Managing Multi-Site Clusters   | 134  |
| Granting and Revoking Tickets Manually                                   | 134  |
| Granting and Revoking Tickets via a Cluster Ticket Registry              | 134  |
| General Management of Tickets  | 135  |
| For more information   | 136  |
| A. Sample Configurations   | 137  |
| Empty  | 137  |
| Simple   | 137  |
| Advanced Configuration   |      |
| B. Revision History  | 140  |
| Index  | 1/12 |

# **List of Figures**

| 1.1. Example Cluster Stack  | 2 |
|---|---|
| 1.2. Internal Components  |   |
| 1.3. Active/Passive Redundancy  |   |
| 1.4. Shared Failover  |   |
| 1.5. N to N Redundancy  |   |
| 5.1. Visual representation of the four resources' start order for the above constraints   |   |
| 5.2. Visual representation of the start order for two ordered sets of unordered resources |   |
| 5.3. Visual representation of the start order for the three sets defined above            |   |
| 5.4. Visual representation of the above example (resources are placed from left to right) |   |

# **List of Tables**

| 2.1. CIB Properties  |      |
|--|------|
| 2.2. Cluster Options   |      |
| 3.1. Node attributes with special significance                                 |      |
| 4.1. Properties of a Primitive Resource  |      |
| 4.2. Meta-attributes of a Primitive Resource                                   |      |
| 4.3. Properties of an Operation  |      |
| 5.1. Attributes of a rsc_location Element                                      |      |
| 5.2. Attributes of a rsc_order Element   |      |
| 5.3. Attributes of a rsc_colocation Constraint                                 |      |
| 5.4. Attributes of a resource_set Element                                      |      |
| 6.1. Additional Properties of Fencing Resources                                | . 44 |
| 6.2. Properties of Fencing Levels  | . 53 |
| 7.1. Meta-Attributes of an Alert   | . 60 |
| 7.2. Environment variables passed to alert agents                              | . 63 |
| 8.1. Attributes of a rule Element  | . 65 |
| 8.2. Attributes of an expression Element                                       | 66   |
| 8.3. Built-in Node Attributes  |      |
| 8.4. Attributes of a date_expression Element                                   | . 68 |
| 8.5. Attributes of a date_spec Element   | . 69 |
| 8.6. Attributes of a duration Element  | . 70 |
| 8.7. Attributes of an rsc_expression Element                                   | . 72 |
| 8.8. Attributes of an op_expression Element                                    | 72   |
| 9.1. Common Options for a ping Resource  | . 82 |
| 9.2. Allowed Values for Node Health Attributes                                 |      |
| 9.3. Node Health Strategies  | 86   |
| 10.1. Properties of a Group Resource   |      |
| 10.2. Properties of a Clone Resource   |      |
| 10.3. Clone-specific configuration options                                     |      |
| 10.4. Additional colocation constraint options for promotable clone resources  |      |
| 10.5. Additional colocation set options relevant to promotable clone resources |      |
| 10.6. Additional ordered set options relevant to promotable clone resources    |      |
| 10.7. Role implications of OCF return codes                                    |      |
| 10.8. Environment variables supplied with Clone notify actions                 |      |
| 10.9. Extra environment variables supplied for promotable clones               |      |
| 10.10. XML Attributes of a bundle Element                                      |      |
| 10.11. XML attributes of a docker, podman, or rkt Element                      |      |
| 10.12. XML attributes of a network Element                                     |      |
| 10.13. Attributes of a port-mapping Element                                    |      |
| 10.14. Attributes of a storage-mapping Element                                 |      |
| 13.1. Properties of an ACL Role  |      |
| 13.2. Properties of an ACL Permission  |      |
| 13.3. Properties of an ACL Target  |      |
| 13.4. Properties of an ACL Group   |      |
| 13.5. Properties of an ACL Role Reference                                      |      |
| 14.1. Authoritative Sources for State Information                              |      |
| 14.2. Node Status Fields   |      |
| 14.3. Contents of an lrm rsc op job  |      |

# **List of Examples**

| 2.1. An empty configuration   | . 5 |
|---|-----|
| 3.1. Example Corosync cluster node entry  | 12  |
| 3.2. Result of using crm_attribute to specify which kernel pcmk-1 is running                        | 13  |
| 4.1. A system resource definition   |     |
| 4.2. An OCF resource definition   | 19  |
| 4.3. An LSB resource with cluster options   | 21  |
| 4.4. An example OCF resource with instance attributes   | 22  |
| 4.5. Displaying the metadata for the Dummy resource agent template                                  | 22  |
| 4.6. An OCF resource with a non-default start timeout   |     |
| 4.7. An OCF resource with a recurring health check  | 26  |
| 4.8. An OCF resource with custom timeouts for its implicit actions                                  |     |
| 4.9. An OCF resource with two recurring health checks, performing different levels of checks        |     |
| specified via OCF_CHECK_LEVEL.  | 28  |
| 4.10. Example of an OCF resource with a disabled health check                                       | 28  |
| 5.1. Opt-in location constraints for two resources  | 31  |
| 5.2. Opt-out location constraints for two resources   |     |
| 5.3. Constraints where a resource prefers two nodes equally   |     |
| 5.4. Optional and mandatory ordering constraints  |     |
| 5.5. Mandatory colocation constraint for two resources  |     |
| 5.6. Mandatory anti-colocation constraint for two resources   |     |
| 5.7. Advisory colocation constraint for two resources   |     |
| 5.8. A set of 3 resources   |     |
| 5.9. A chain of ordered resources   | 37  |
| 5.10. A chain of ordered resources expressed as a set   | 37  |
| 5.11. Ordered sets of unordered resources   | 38  |
| 5.12. Advanced use of set ordering - Three ordered sets, two of which are internally unordered      | 38  |
| 5.13. Resource Set "OR" logic: Three ordered sets, where the first set is internally unordered with |     |
| "OR" logic  | 39  |
| 5.14. Colocation chain as individual constraints, where A is placed first, then B, then C, then D   | 39  |
| 5.15. Equivalent colocation chain expressed using resource_set                                      | 40  |
| 5.16. Using colocated sets to specify a shared dependency   | 40  |
| 5.17. Colocation in which the members of the middle set have no interdependencies, and the last     |     |
| set listed applies only to instances in the master role   | 41  |
| 6.1. Obtaining a list of Fence Agent Parameters   | 50  |
| 6.2. An IPMI-based Fencing Resource   | 52  |
| 6.3. Fencing topology with different devices for different nodes                                    | 53  |
| 7.1. Simple alert configuration   | 59  |
| 7.2. Alert configuration with recipient   | 59  |
| 7.3. Alert configuration with meta-attributes   | 60  |
| 7.4. Alert configuration with instance attributes   | 61  |
| 7.5. Alert configuration to receive only node events and fencing events                             | 61  |
| 7.6. Alert configuration to be called when certain node attributes change                           | 62  |
| 7.7. Sending cluster events as SNMP traps   | 62  |
| 7.8. Sending cluster events as e-mails  | 62  |
| 8.1. True if now is any time in the year 2005   |     |
| 8.2. Equivalent expression  |     |
| 8.3. 9am-5pm Monday-Friday  |     |
| 8.4. 9am-6pm Monday through Friday or anytime Saturday  |     |
| 8.5. 9am-5pm or 9pm-12am Monday through Friday  |     |
| 8.6. Mondays in March 2005  |     |
| ·   | 71  |

| 8.9. Provider doesn't apply to non-OCF resources  8.10. True for all monitor actions  8.11. True for all monitor actions with a 10 second interval |          |
|--|----------|
|  | 73       |
| 8.11. True for all monitor actions with a 10 second interval   |          |
| **   | 73       |
| 8.12. Prevent resource "webserver" from running on node3   | 73       |
| 8.13. Prevent resource "webserver" from running on node3 using rule  | 73       |
| 8.14. A sample nodes section for use with score-attribute  |          |
| 8.15. Defining different resource options based on the node name   |          |
| 8.16. Change resource-stickiness during working hours  |          |
| 8.17. Set all IPaddr2 resources to stopped   |          |
| 8.18. Set all monitor action timeouts to 7 seconds   |          |
| 8.19. Set the monitor action timeout on all IPaddr2 resources with a given monitor interval to 8   |          |
| seconds  | 76       |
| 8.20. Schedule a maintenance window for 9 to 11 p.m. CDT Sept. 20, 2019  |          |
| 9.1. Specifying a Base for Recurring Action Intervals  |          |
| 9.2. An example ping cluster resource that checks node connectivity once every minute  |          |
| 9.3. Don't run a resource on unconnected nodes   |          |
| 9.4. Run only on nodes connected to three or more ping targets.  |          |
| 9.5. Prefer the node with the most connected ping nodes  |          |
| 9.6. How the cluster translates the above location constraint  |          |
| 9.7. A more complex example of choosing a location based on connectivity   |          |
| 9.8. The DRBD agent's logic for supporting reload  |          |
| 9.9. The DRBD Agent Advertising Support for the reload Operation   |          |
| 9.10. Parameter that can be changed using reload   |          |
| 10.1. A group of two primitive resources   | 80<br>80 |
| 10.2. How the cluster sees a group resource  |          |
| 10.3. Some constraints involving groups  |          |
| 10.4. A clone that runs a web server on all nodes  |          |
| 10.4. A clone that runs a web server on an nodes   |          |
|  |          |
| 10.6. Constraints involving promotable clone resources   |          |
| 10.7. Colocate C and D with A's and B's master instances   |          |
| 10.8. Start C and D after first promoting A and B  |          |
| 10.9. Notification variables   |          |
| 10.10. Monitoring both states of a promotable clone resource   |          |
| 10.11. Explicitly preferring node1 to be promoted to master  |          |
| 10.12. A bundle for a containerized web server   |          |
| 11.1. Resource template for a migratable Xen virtual machine   |          |
| 11.2. Xen primitive resource using a resource template   |          |
| 11.3. Equivalent Xen primitive resource not using a resource template  |          |
| 11.4. Xen resource overriding template values  |          |
| 11.5. Referencing rules from other constraints   |          |
| 11.6. Referencing attributes, options, and operations from other resources   |          |
| 11.7. Tag referencing three resources  |          |
| 11.8. Constraint using a tag   |          |
| 11.9. Equivalent constraints without tags  |          |
| 12.1. Specifying CPU and RAM capacities of two nodes   |          |
| 12.2. Specifying CPU and RAM consumed by several resources   | 17       |
| 14.1. A bare-bones status entry for a healthy node <b>cl-virt-1</b>  |          |
| 14.2. A set of transient node attributes for node <b>cl-virt-1</b>   | 26       |
| 14.3. A record of the apcstonith resource  |          |
| 14.4. A monitor operation (determines current state of the apcstonith resource)  | 28       |
| 14.5. Resource history of a pingd clone with multiple jobs   |          |
| 15.1. Constraint that fences node if ticketA is revoked  |          |
| 15.2. Constraint that demotes rsc1 if ticketA is revoked   | 33       |

### Configuration Explained

| 15.3. Ticket constraint for multiple resources                                  | 133 |
|---|-----|
| A.1. An Empty Configuration   | 137 |
| A.2. A simple configuration with two nodes, some cluster options and a resource | 137 |
| A.3. An advanced configuration with groups, clones and STONITH                  | 138 |

# **Preface**

### **Table of Contents**

| Document Conventions    | хi  |
|-------------------------|-----|
| Typographic Conventions | хi  |
| Pull-quote Conventions  | xii |
| Notes and Warnings      | xiv |
| We Need Feedback!       | xiv |

## **Document Conventions**

This manual uses several conventions to highlight certain words and phrases and draw attention to specific pieces of information.

# **Typographic Conventions**

Four typographic conventions are used to call attention to specific words and phrases. These conventions, and the circumstances they apply to, are as follows.

Mono-spaced Bold

Used to highlight system input, including shell commands, file names and paths. Also used to highlight keys and key combinations. For example:

To see the contents of the file my\_next\_bestselling\_novel in your current working directory, enter the **cat my\_next\_bestselling\_novel** command at the shell prompt and press **Enter** to execute the command.

The above includes a file name, a shell command and a key, all presented in mono-spaced bold and all distinguishable thanks to context.

Key combinations can be distinguished from an individual key by the plus sign that connects each part of a key combination. For example:

Press **Enter** to execute the command.

Press Ctrl+Alt+F2 to switch to a virtual terminal.

The first example highlights a particular key to press. The second example highlights a key combination: a set of three keys pressed simultaneously.

If source code is discussed, class names, methods, functions, variable names and returned values mentioned within a paragraph will be presented as above, in mono-spaced bold. For example:

File-related classes include filesystem for file systems, file for files, and dir for directories. Each class has its own associated set of permissions.

#### Proportional Bold

This denotes words or phrases encountered on a system, including application names; dialog-box text; labeled buttons; check-box and radio-button labels; menu titles and submenu titles. For example:

Choose System  $\rightarrow$  Preferences  $\rightarrow$  Mouse from the main menu bar to launch Mouse Preferences. In the Buttons tab, select the Left-handed mouse check box and click Close to switch the primary mouse button from the left to the right (making the mouse suitable for use in the left hand).

To insert a special character into a gedit file, choose Applications  $\rightarrow$  Accessories  $\rightarrow$  Character Map from the main menu bar. Next, choose Search  $\rightarrow$  Find... from the Character Map menu bar, type the name of the character in the Search field and click Next. The character you sought will be highlighted in the Character Table. Double-click this highlighted character to place it in the Text to copy field and then click the Copy button. Now switch back to your document and choose Edit  $\rightarrow$  Paste from the gedit menu bar.

The above text includes application names; system-wide menu names and items; application-specific menu names; and buttons and text found within a GUI interface, all presented in proportional bold and all distinguishable by context.

```
Mono-spaced Bold Italic or Proportional Bold Italic
```

Whether mono-spaced bold or proportional bold, the addition of italics indicates replaceable or variable text. Italics denotes text you do not input literally or displayed text that changes depending on circumstance. For example:

To connect to a remote machine using ssh, type **ssh username@domain.name** at a shell prompt. If the remote machine is example.com and your username on that machine is john, type **ssh john@example.com**.

The **mount -o remount file-system** command remounts the named file system. For example, to remount the /home file system, the command is **mount -o remount /home**.

To see the version of a currently installed package, use the **rpm -q package** command. It will return a result as follows: **package-version-release**.

Note the words in bold italics above: username, domain.name, file-system, package, version and release. Each word is a placeholder, either for text you enter when issuing a command or for text displayed by the system.

Aside from standard usage for presenting the title of a work, italics denotes the first use of a new and important term. For example:

Publican is a *DocBook* publishing system.

## **Pull-quote Conventions**

Terminal output and source code listings are set off visually from the surrounding text.

Output sent to a terminal is set in mono-spaced roman and presented thus:

```
books Desktop documentation drafts mss photos stuff svn books_tests Desktop1 downloads images notes scripts svgs
```

Source-code listings are also set in mono-spaced roman but add syntax highlighting as follows:

```
package org.jboss.book.jca.ex1;
import javax.naming.InitialContext;
```

## **Notes and Warnings**

Finally, we use three visual styles to draw attention to information that might otherwise be overlooked.

#### Note

Notes are tips, shortcuts or alternative approaches to the task at hand. Ignoring a note should have no negative consequences, but you might miss out on a trick that makes your life easier.

### **Important**

Important boxes detail things that are easily missed: configuration changes that only apply to the current session, or services that need restarting before an update will apply. Ignoring a box labeled "Important" will not cause data loss but may cause irritation and frustration.

### Warning

Warnings should not be ignored. Ignoring warnings will most likely cause data loss.

# We Need Feedback!

If you find a typographical error in this manual, or if you have thought of a way to make this manual better, we would love to hear from you! Please submit a report in Bugzilla<sup>1</sup> against the product Pacemaker.

When submitting a bug report, be sure to mention the manual's identifier: Pacemaker\_Explained

If you have a suggestion for improving the documentation, try to be as specific as possible when describing it. If you have found an error, please include the section number and some of the surrounding text so we can find it easily.

<sup>1</sup> http://bugs.clusterlabs.org

# Chapter 1. Read-Me-First

### **Table of Contents**

| The Scope of this Document | . 1 |
|----------------------------|-----|
| What Is Pacemaker?         | . 1 |
| Cluster Architecture       |     |
| Pacemaker Architecture     | . 2 |
| Node Redundancy Designs    |     |

# 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 at all in this document <sup>1</sup>.

Users may be interested in other parts of the Pacemaker documentation set [https://www.clusterlabs.org/pacemaker/doc/], 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.

## 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>2</sup> It is maintained by the ClusterLabs [https://www.ClusterLabs.org/] 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 anticolocation
- Support for advanced service types, such as *clones* (services that need to be active on multiple nodes), *stateful resources* (clones that can run in one of two modes), and containerized services

<sup>&</sup>lt;sup>1</sup> I hope, however, that the concepts explained here make the functionality of these tools more easily understood.

<sup>&</sup>lt;sup>2</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

### **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 [http://www.corosync.org/] 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.

Figure 1.1. Example Cluster Stack

## **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

#### Figure 1.2. Internal Components

The Pacemaker master process (pacemakerd) spawns all the other daemons, and respawns them if they unexpectedly exit.

The *Cluster Information Base* (CIB) is an XML [https://en.wikipedia.org/wiki/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-attrd) 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.

#### 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 |

| Old name          | New name          |
|-------------------|-------------------|
| lrmd              | pacemaker-execd   |
| stonithd          | pacemaker-fenced  |
| pacemaker_remoted | pacemaker-remoted |

# **Node Redundancy Designs**

Pacemaker supports practically any node redundancy configuration [https://en.wikipedia.org/wiki/High-availability\_cluster#Node\_configurations] 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 [https://en.wikipedia.org/wiki/Distributed\_Replicated\_Block\_Device:] 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.

#### Figure 1.3. Active/Passive Redundancy

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.

#### Figure 1.4. Shared Failover

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.

Figure 1.5. N to N Redundancy

# **Chapter 2. Cluster-Wide Configuration**

### **Table of Contents**

| Configuration Layout | . 5 |
|----------------------|-----|
| CIB Properties       | . 6 |
| Cluster Options      | 6   |

# **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:

#### **Example 2.1. An empty configuration**

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

- cib: The entire CIB is enclosed with a cib tag. Certain fundamental settings are defined as attributes of this tag.
  - 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.

# **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 2.1. CIB Properties** 

| Field            | 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 (admin_epoch, epoch, num_updates) tuple to replace the configuration on all the nodes — which makes setting them, and setting them correctly, very important. admin_epoch is never modified by the cluster; you can use this to make the configurations on any inactive nodes obsolete. <i>Never set this value to zero</i> . 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 none, the cluster will not verify that updates conform to the DTD (nor reject ones that don't). This option can be useful when operating a mixed-version cluster during an upgrade.  |  |
| 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 no-quorum-policy 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.   |  |

# **Cluster Options**

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

They are grouped into sets within the crm\_config section, and, in advanced configurations, there may be more than one set. (This will be described later in the section on Chapter 8, 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 2.2. 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 rules (as the #cluster-name node attribute). It may also be used by higher-level tools when displaying cluster information, and by certain resource agents (for example, the ocf:heartbeat:GFS2 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-<br>infrastructure |         | The messaging stack on which Pacemaker is currently running. Determined automatically by the cluster. Used for informational and diagnostic purposes.   |
| no-quorum-policy           | stop    | What to do when the cluster does not have quorum. Allowed values:  • ignore: continue all resource management  • freeze: continue resource management, but don't recover resources from nodes not in the affected partition  • stop: stop all resources in the affected cluster partition  • demote: demote promotable resources and stop all other resources in the affected cluster partition  • suicide: fence all nodes in the affected cluster partition |
| 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.   |
| migration-limit            | -1      | The number of live migration actions that the cluster is allowed to execute in parallel on a node. A value of -1 means unlimited.   |
| symmetric-cluster          | TRUE    | Can all resources run on any node by default?   |
| stop-all-<br>resources     | FALSE   | Should the cluster stop all resources?  |

| Option                       | Default | Description  |
|------------------------------|---------|--|
| stop-orphan-<br>resources    | TRUE    | Should deleted resources be stopped? This value takes precedence over is-managed (i.e. even unmanaged resources will be stopped if deleted from the configuration when this value is TRUE).  |
| stop-orphan-<br>actions      | TRUE    | Should deleted actions be cancelled?   |
| start-failure-is-<br>fatal   | TRUE    | Should a failure to start a resource on a particular node prevent further start attempts on that node? If FALSE, the cluster will decide whether the same node is still eligible based on the resource's current failure count and migration-threshold (see the section called "Handling Resource Failure").   |
| enable-startup-<br>probes    | TRUE    | Should the cluster check for active resources during startup?  |
| maintenance-mode             | FALSE   | Should the cluster refrain from monitoring, starting and stopping resources?   |
| stonith-enabled              | TRUE    | Should failed nodes and nodes with resources that can't be stopped be shot? If you value your data, set up a STONITH device and enable this.   |
|                              |         | If true, or unset, the cluster will refuse to start resources unless one or more STONITH resources have been configured. If false, unresponsive nodes are immediately assumed to be running no resources, and resource takeover to online nodes starts without any further protection (which means <i>data loss</i> if the unresponsive node still accesses shared storage, for example). See also the requires meta-attribute in the section called "Resource Options".   |
| stonith-action               | reboot  | Action to send to STONITH device. Allowed values are reboot and off. The value poweroff is also allowed, but is only used for legacy devices.  |
| stonith-timeout              | 60s     | How long to wait for STONITH actions (reboot, on, off) to complete   |
| stonith-max-<br>attempts     | 10      | How many times fencing can fail for a target before the cluster will no longer immediately re-attempt it.  |
| stonith-watchdog-<br>timeout | 0       | If nonzero, along with have-watchdog=true automatically set by the cluster, when fencing is required, watchdog-based self-fencing will be performed via SBD without requiring a fencing resource explicitly configured. If stonith-watchdog-timeout is set to a positive value, unseen nodes are assumed to self-fence within this much time. WARNING: It must be ensured that this value is larger than the SBD_WATCHDOG_TIMEOUT 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's strongly recommended that SBD_WATCHDOG_TIMEOUT is set to the same value on all nodes. If stonith-watchdog-timeout is set to |

| Option                 | Default | Description  |
|------------------------|---------|--|
|                        |         | a negative value, and SBD_WATCHDOG_TIMEOUT is set, twice that value will be used. WARNING: In this case, it's essential (currently not verified by pacemaker) that SBD_WATCHDOG_TIMEOUT is set to the same value on all nodes.   |
| concurrent-<br>fencing | FALSE   | Is the cluster allowed to initiate multiple fence actions concurrently?  |
| 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 stop to attempt to immediately stop pacemaker and stay stopped, or panic to attempt to immediately reboot the local node, falling back to stop on failure. The default is likely to be changed to panic in a future release. (since 2.0.3)   |
| priority-fencing-delay | 0       | Apply specified delay for the fencings that are 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, which is especially meaningful under split-brain of 2-node cluster. A promoted resource instance takes the base priority + 1 on calculation if the base priority is not 0. Any static/random delays that are introduced by pcmk_delay_base/max configured for the corresponding fencing resources will be added to this delay. This delay should be significantly greater than, safely twice, the maximum pcmk_delay_base/max. By default, priority fencing delay is disabled. (since 2.0.4) |
| 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.   |

| Option                       | Default | Description   |
|------------------------------|---------|---|
| pe-error-series-<br>max      | -1      | The number of PE inputs resulting in ERRORs to save. Used when reporting problems. A value of -1 means unlimited (report all).  |
| pe-warn-series-<br>max       | -1      | The number of PE inputs resulting in WARNINGs to save. Used when reporting problems. A value of -1 means unlimited (report all).  |
| pe-input-series-<br>max      | -1      | The number of "normal" PE inputs to save. Used when reporting problems. A value of -1 means unlimited (report all).   |
| placement-<br>strategy       | default | How the cluster should allocate resources to nodes (see Chapter 12, <i>Utilization and Placement Strategy</i> [116]). Allowed values are default, utilization, balanced, and minimal.   |
| node-health-<br>strategy     | none    | How the cluster should react to node health attributes (see the section called "Tracking Node Health"). Allowed values are none, migrate-on-red, only-green, progressive, and custom.   |
| enable-acl                   | FALSE   | Whether access control lists (ACLs) (see Chapter 13, ACLs) can be used to authorize modifications to the CIB.   |
| node-health-base             | 0       | The base health score assigned to a node. Only used when node-health-strategy is progressive.   |
| node-health-green            | 0       | The score to use for a node health attribute whose value is green. Only used when node-health-strategy is progressive or custom.  |
| node-health-<br>yellow       | 0       | The score to use for a node health attribute whose value is yellow. Only used when node-health-strategy is progressive or custom.   |
| node-health-red              | 0       | The score to use for a node health attribute whose value is red. Only used when node-health-strategy is progressive or custom.  |
| cluster-recheck-<br>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. However, it will also recheck the cluster after this amount of inactivity. This has two goals: rules with date_spec are only guaranteed to be checked this often, and it also serves as a fail-safe for certain classes 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 shutdown-lock-limit, if set). Stonith resources and Pacemaker Remote connections |

| Option                            | Default | Description  |
|-----------------------------------|---------|--|
|                                   |         | are never locked. Clone and bundle instances and the master role of promotable clones are currently never locked, though support could be added in a future release. Locks may be manually cleared using therefresh option of crm_resource (both the resource and node must be specified; this works with remote nodes if their connection resource's target-role is set to Stopped, but not if Pacemaker Remote is stopped on the remote node without disabling the connection resource). (since 2.0.4) |
| shutdown-lock-<br>limit           | 0       | If shutdown-lock 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 target-role is set to Stopped.) (since 2.0.4)  |
| remove-after-stop                 | FALSE   | Advanced Use Only: Should the cluster remove resources from the LRM after they are stopped? Values other than the default are, at best, poorly tested and potentially dangerous.   |
| startup-fencing                   | TRUE    | Advanced Use Only: Should the cluster shoot unseen nodes? Not using the default is very unsafe!  |
| election-timeout                  | 2min    | Advanced Use Only: If you need to adjust this value, it probably indicates the presence of a bug.  |
| shutdown-<br>escalation           | 20min   | Advanced Use Only: If you need to adjust this value, it probably indicates the presence of a bug.  |
| join-integration-<br>timeout      | 3min    | Advanced Use Only: If you need to adjust this value, it probably indicates the presence of a bug.  |
| join-<br>finalization-<br>timeout | 30min   | Advanced Use Only: If you need to adjust this value, it probably indicates the presence of a bug.  |
| transition-delay                  | Os      | Advanced Use Only: Delay cluster recovery for the configured interval to allow for additional/related events to occur. 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.   |

# **Chapter 3. Cluster Nodes**

### **Table of Contents**

| Defining a Cluster Node              | - 12 |
|--------------------------------------|------|
| Where Pacemaker Gets the Node Name   |      |
| Node Attributes                      | 12   |
| Setting and querying node attributes | 13   |
| Special node attributes              |      |

# **Defining a Cluster Node**

Each node in the cluster will have an entry in the nodes section containing its UUID, uname, and type.

#### Example 3.1. 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 communications and membership data.

## Where Pacemaker Gets the Node Name

Traditionally, Pacemaker required nodes to be referred to by the value returned by uname -n. This can be problematic for services that require the uname -n to be a specific value (e.g. for a licence file).

This requirement has been relaxed for clusters using Corosync 2.0 or later. The name Pacemaker uses is:

- 1. The value stored in corosync.conf under ring0\_addr in the nodelist, if it does not contain an IP address; otherwise
- 2. The value stored in corosync.conf under name in the nodelist; otherwise
- 3. The value of uname -n

Pacemaker provides the crm\_node -n command which displays the name used by a running cluster.

If a Corosync **nodelist** is used, crm\_node --name-for-id *number* is also available to display the name used by the node with the corosync **nodeid** of *number*, for example: crm\_node --name-for-id 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.

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 of what XML configuration would be generated if an administrator ran this command:

# Example 3.2. Result of using crm\_attribute to specify which kernel pcmk-1 is running

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
```

By specifying --type nodes the admin tells the cluster that this attribute is persistent across reboots. There are also transient attributes which are kept in the status section and are "forgotten" whenever the node leaves the cluster. Administrators can use this section by specifying --type status.

## 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 the section called "Node Attribute Expressions" [66].

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 3.1. Node attributes with special significance

| Name           | Description   |  |
|----------------|---|--|
| fail-count-*   | Attributes whose names start with fail-count- 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 crm_failcount or crm_resourcecleanup commands rather than directly. |  |
| last-failure-* | Attributes whose names start with last-failure- are managed by the cluster to track when particular resource operations have most recently failed on this node. These should be cleared via the   |  |

| Name                           | Description  |  |
|--------------------------------|--|--|
|                                | crm_failcount or crm_resourcecleanup commands rather than directly.  |  |
| maintenance                    | Similar to the maintenance-mode cluster option, 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.                        |  |
| probe_complete                 | This is managed by the cluster to detect when nodes need to be reprobed, and should never be used directly.  |  |
| resource-<br>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 resource-discovery 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 #site-name node attribute used in rules. (If not set, the value of the cluster-name cluster option will be used as #site-name instead.)  |  |
| standby                        | If true, the node is in standby mode. This is typically set and queried via the crm_standby 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 #digests- are managed by the cluster to detect when unfencing 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.  |  |

### Warning

Restarting pacemaker on a node that is in single-node maintenance mode will likely lead to undesirable effects. If maintenance 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.

# **Chapter 4. Cluster Resources**

### **Table of Contents**

| What is a Cluster Resource?                          |
|--|
| Resource Classes                                     |
| Open Cluster Framework                               |
| Linux Standard Base                                  |
| Systemd  |
| Upstart  |
| System Services                                      |
| STONITH  |
| Nagios Plugins                                       |
| Resource Properties                                  |
| Resource Options                                     |
| Resource Meta-Attributes                             |
| Setting Global Defaults for Resource Meta-Attributes |
| Resource Instance Attributes                         |
| Resource Operations                                  |
| Operation Properties                                 |
| Monitoring Resources for Failure                     |
| Monitoring Resources When Administration is Disabled |
| Setting Global Defaults for Operations               |
| When Implicit Operations Take a Long Time            |
| Multiple Monitor Operations                          |
| Disabling a Monitor Operation                        |

## 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 section. More complex forms, such as groups and clones, are described in later sections.

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.

## **Resource Classes**

Pacemaker supports several classes of agents:

OCF

- LSB
- Upstart
- Systemd
- · 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.

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 [http://www.linux-ha.org/wiki/OCF\_Resource\_Agents] and the *Resource Agents* section 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

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

The resource-agents source code includes the ocf-tester script, which can be useful in this regard.

<sup>&</sup>lt;sup>3</sup> See 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.

Resource Agents section 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**

Some newer distributions have replaced the old "SysV" [http://en.wikipedia.org/wiki/Init#SysV-style] style of initialization daemons and scripts with an alternative called Systemd [http://www.freedesktop.org/wiki/Software/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 newer distributions have replaced the old "SysV" [http://en.wikipedia.org/wiki/Init#SysV-style] style of initialization daemons (and scripts) with an alternative called Upstart [http://upstart.ubuntu.com/].

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.

### **Important**

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

## **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.

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

This is particularly useful when the cluster contains a mix of systemd, upstart, and 1sb.

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

- 1. an LSB init script
- 2. a Systemd unit file
- 3. an Upstart job

### **STONITH**

The STONITH class is used exclusively for fencing-related resources. This is discussed later in Chapter 6, Fencing.

# **Nagios Plugins**

Nagios Plugins <sup>5</sup> allow us to monitor services on remote hosts.

Pacemaker is able to do remote monitoring with the plugins if they are present.

A common use case is to configure them as resources belonging to a resource container (usually a virtual machine), and the container will be restarted if any of them has failed. Another use is to configure them as ordinary resources to be used for monitoring hosts or services via the network.

The supported parameters are same as the long options of the plugin.

# **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 4.1. Properties of a Primitive Resource

| Field    | Description   |  |
|----------|---|--|
| id       | Your name for the resource  |  |
| class    | The standard the resource agent conforms to. Allowed values: lsb, nagios, ocf, service, stonith, systemd, upstart   |  |
| type     | The name of the Resource Agent you wish to use. E.g. IPaddr or Filesystem   |  |
| 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 heartbeat 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:

<sup>&</sup>lt;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.

### Example 4.1. 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!

#### **Example 4.2. An OCF resource definition**

# **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 4.2. 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.   |
| target-role | Started | <ul> <li>What state should the cluster attempt to keep this resource in? Allowed values:</li> <li>Stopped: Force the resource to be stopped</li> <li>Started: Allow the resource to be started (and in the case of promotable clone resources, promoted to master if appropriate)</li> <li>Slave: Allow the resource to be started, but only in Slave mode if the resource is promotable</li> <li>Master: Equivalent to Started</li> </ul> |
| is-managed  | TRUE    | Is the cluster allowed to start and stop the resource? Allowed values: true, false   |
| maintenance | FALSE   | Similar to the maintenance-mode cluster option, but for a single resource. If true, the resource will not be started, stopped, or monitored on any node. This differs from is-managed in that monitors will not be run. Allowed values: true, false  |

| Field                   | Default  | Description   |
|-------------------------|--|---|
| resource-<br>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.   |
| requires                | quorum for resources with a class of stonith, otherwise unfencing if unfencing is active in the cluster, otherwise fencing if stonithenabled is true, otherwise quorum | Conditions under which the resource can be started Allowed values:  • nothing: can always be started  • quorum: The cluster can only start this resource if a majority of the configured nodes are active  • fencing: The cluster can only start this resource if a majority of the configured nodes are active and any failed or unknown nodes have been fenced  • unfencing: The cluster can only start this resource if a majority of the configured nodes are active and any failed or unknown nodes have been fenced and only on nodes that have been unfenced |
| migration-<br>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 constrast, the cluster treats INFINITY (the default) as a very large but finite number. This option has an effect only if the failed operation specifies on-fail as restart (the default), and additionally for failed start operations, if the cluster property start-failure-is-fatal is false.   |
| failure-<br>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.   |
| multiple-<br>active     | stop_start   | What should the cluster do if it ever finds the resource active on more than one node? Allowed values:  • block: mark the resource as unmanaged  • stop_only: stop all active instances and leave them that way  • stop_start: stop all active instances and start the resource in one location only  |

| Field                              | Default  | Description   |
|------------------------------------|--|---|
| 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 the section called "Migrating Resources")   |
| container-<br>attribute-<br>target |  | Specific to bundle resources; see the section called "Bundle Node Attributes"   |
| 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. WARNING: This value cannot overlap with any resource or node IDs. |
| remote-port                        | 3121   | If remote-node 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 remote-<br>node                                 | If remote-node 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.   |
| remote-<br>connect-<br>timeout     | 60s  | If remote-node 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
# crm_resource -m -r Email -p multiple-active -v block
```

the resulting resource definition might be:

#### Example 4.3. An LSB resource with cluster options

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 rsc\_defaults section with crm\_attribute. For example,

```
# crm_attribute --type rsc_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:

#### Example 4.4. An example OCF resource with instance attributes

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.

#### Example 4.5. 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="1.0">
<version>1.0</version>

<longdesc>
This is a Dummy Resource Agent. It does absolutely nothing except keep track of whether its running or not.
Its purpose in life is for testing and to serve as a template for RA writers.

NB: Please pay attention to the timeouts specified in the actions
```

section below. They should be meaningful for the kind of resource

```
the agent manages. They should be the minimum advised timeouts,
but they shouldn't/cannot cover _all_ possible resource
instances. So, try to be neither overly generous nor too stingy,
but moderate. The minimum timeouts should never be below 10 seconds.
</longdesc>
<shortdesc>Example stateless resource agent/shortdesc>
<parameters>
<parameter name="state" unique="1">
<longdesc>
Location to store the resource state in.
</longdesc>
<shortdesc>State file</shortdesc>
<content type="string" default="/var/run/Dummy-default.state" />
<parameter name="fake" unique="0">
<longdesc>
Fake attribute that can be changed to cause a reload
</longdesc>
<shortdesc>Fake attribute that can be changed to cause a reload</shortdesc>
<content type="string" default="dummy" />
</parameter>
<parameter name="op_sleep" unique="1">
Number of seconds to sleep during operations. This can be used to test how
the cluster reacts to operation timeouts.
</longdesc>
<shortdesc>Operation sleep duration in seconds.
<content type="string" default="0" />
</parameter>
</parameters>
<actions>
<action name="start"
                          timeout="20" />
<action name="stop"
                           timeout="20" />
<action name="monitor"
                           timeout="20" interval="10" depth="0"/>
<action name="reload"
                           timeout="20" />
<action name="migrate to"
                           timeout="20" />
<action name="migrate_from" timeout="20" />
<action name="validate-all" timeout="20" />
</actions>
</resource-agent>
```

# **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).

#### Example 4.6. An OCF resource with a non-default start timeout

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 4.3. 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 monitor, start, and stop.  |  |  |
| 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 monitor.   |  |  |
| timeout  |   | How long to wait before declaring the action has failed  |  |  |
| on-fail  | Varies by action:  • stop: fence if stonith-enabled is true or block otherwise  • demote: on-fail of the monitor action with role set to Master, if present, enabled, and configured to a value | <ul> <li>The action to take if this action ever fails. Allowed values:</li> <li>ignore: Pretend the resource did not fail.</li> <li>block: Don't perform any further operations on the resource.</li> <li>stop: Stop the resource and do not start it elsewhere.</li> <li>demote: Demote the resource, without a full restart. This is valid only for promote actions, and for monitor actions with both a nonzero interval and role set to Master; for any</li> </ul> |  |  |

| Field              | Default   | Description   |
|--------------------|---|---|
|                    | other than demote, or restart otherwise  • all other actions: restart | other action, a configuration error will be logged, and the default behavior will be used.  • restart: Stop the resource and start it again (possibly on a different node).   |
|                    |   | • fence: STONITH the node on which the resource failed.   |
|                    |   | standby: Move <i>all</i> resources away from the node on which the resource failed.   |
| enabled            | TRUE  | If false, ignore this operation definition. This is typically used to pause a particular recurring monitor operation; for instance, it can complement the respective resource being unmanaged (is-managed=false), as this alone will not block any configured monitoring. Disabling the operation does not suppress all actions of the given type. Allowed values: true, false. |
| record-<br>pending | TRUE  | If true, 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 the section called "Setting Global Defaults for Operations"). Allowed values: true, false.   |
| role               |   | Run the operation only on node(s) that the cluster thinks should be in the specified role. This only makes sense for recurring monitor operations.  Allowed (case-sensitive) values: Stopped, Started, and in the case of promotable clone resources, Slave and Master.   |

#### 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 master scores have not changed, it will be promoted again.

There is no direct equivalent of migration—threshold for the master 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 master role from a node with any failed promote or master monitor:

```
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 master 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.

#### Example 4.7. An OCF resource with a recurring health check

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:

• 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.

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

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.

#### Example 4.8. An OCF resource with custom timeouts for its implicit actions

# **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.

# Example 4.9. An OCF resource with two recurring health checks, performing different levels of checks specified via OCF\_CHECK\_LEVEL.

```
<pri><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 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 4.10. Example of an OCF resource with a disabled health check

# **Chapter 5. Resource Constraints**

#### **Table of Contents**

| Scores  | 29 |
|---|----|
| Infinity Math   | 29 |
| Deciding Which Nodes a Resource Can Run On                | 3( |
| Location Properties                                       | 3( |
| Asymmetrical "Opt-In" Clusters                            | 31 |
| Symmetrical "Opt-Out" Clusters                            | 32 |
| What if Two Nodes Have the Same Score                     | 32 |
| Specifying the Order in which Resources Should Start/Stop | 32 |
| Ordering Properties                                       | 33 |
| Optional and mandatory ordering                           | 34 |
| Placing Resources Relative to other Resources             | 34 |
| Colocation Properties                                     | 34 |
| Mandatory Placement                                       | 35 |
| Advisory Placement  |    |
| Colocation by Node Attribute                              | 36 |
| Resource Sets   | 36 |
| Ordering Sets of Resources                                | 37 |
| Ordered Set   | 37 |
| Ordering Multiple Sets                                    | 37 |
| Resource Set OR Logic                                     |    |
| Colocating Sets of Resources                              | 39 |

# **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).

# **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 5.1. Attributes of a rsc\_location Element

| Attribute   | Default | Description   |  |
|-------------|---------|---|--|
| id          |         | A unique name for the constraint (required)   |  |
| rsc         |         | The name of the resource to which this constraint applies.  A location constraint must either have a rsc, have a rsc- pattern, or contain at least one resource set.  |  |
| rsc-pattern |         | A pattern matching the names of resources to which this constraint applies. The syntax is the same as POSIX [http://pubs.opengroup.org/onlinepubs/9699919799/basedefs/V1_chap09.html#tag_09_04] extended regular expressions, with the addition of an initial ! indicating that resources not matching the pattern are selected. If the regular expression contains submatches, and the constraint is governed by a rule, the submatches can be referenced as %0 through %9 in the rule's score-attribute or a rule expression's attribute. A location constraint must either have a rsc, have a rsc-pattern, or contain at least one resource set. |  |
| node        |         | The name of the node to which this constraint applies. A location constraint must either have a node and score, or contain at least one rule.   |  |
| score       |         | Positive values indicate a preference for running the affected resource(s) on node — the higher the value, the stronger the preference. Negative values indicate the resource(s) should avoid this node (a value of -INFINITY changes "should" to "must"). A location constraint must either have a node and score, or contain at least one rule.   |  |

| Attribute              | Default | Description   |
|------------------------|---------|---|
| resource-<br>discovery | always  | 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.  • always: Always perform resource discovery for the specified resource on this node.  • never: Never perform resource discovery for the specified resource on this node. This option should generally be used with a -INFINITY score, although that is not strictly required.  • exclusive: Perform resource discovery for the specified resource only on this node (and other nodes similarly marked as exclusive). Multiple location constraints using exclusive discovery for the same resource across different nodes creates a subset of nodes resource-discovery is exclusive to. If a resource is marked for exclusive discovery on one or more nodes, that resource is only allowed to be placed within that subset of nodes. |

#### 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.

#### **Example 5.1. Opt-in location constraints for two resources**

<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.

#### **Example 5.2. Opt-out location constraints for two resources**

#### 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.

#### Example 5.3. 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.

# 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 the section called "Placing Resources Relative to other Resources"), or alternatively, a group (see the section called "Groups - A Syntactic Shortcut").

# **Ordering Properties**

Table 5.2. Attributes of a rsc\_order Element

| Field            | Default  | Description  |  |
|------------------|--|--|--|
| id               |  | A unique name for the constraint   |  |
| first            |  | Name of the resource that the then resource depends on   |  |
| then             |  | Name of the dependent resource   |  |
| first-<br>action | start  | The action that the first resource must complete before then-action can be initiated for the then resource. Allowed values: start, stop, promote, demote.  |  |
| then-<br>action  | value of<br>first-<br>action                     | The action that the then resource can execute only after the first-action on the first resource has completed. Allowed values: start, stop, promote, demote.   |  |
| kind             | Mandatory  | How to enforce the constraint. Allowed values:   |  |
|                  |  | • Mandatory: then-action will never be initiated for the then resource unless and until first-action successfully completes for the first resource.  |  |
|                  |  | • Optional: 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 thenaction is allowed on the then resource regardless of the state of the first resource, but if both actions happen to be scheduled at the same time, they will be ordered. |  |
|                  |  | • Serialize: Ensure that the specified actions are never performed concurrently for the specified resources. Firstaction and then-action 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.   |  |
| symmetrica       | TRUE for Mandatory and Optional kinds. FALSE for | 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). Serialize orders cannot be symmetrical.   |  |

| Field | Default   | Description |
|-------|-----------|-------------|
|       | Serialize |             |
|       | kind.     |             |

Promote and demote apply to the master role of promotable 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:

#### Example 5.4. 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.

# 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.

Another thing to keep in mind is that, assuming A is colocated with B, the cluster will take into account A's preferences when deciding which node to choose for B.

For a detailed look at exactly how this occurs, see Colocation Explained [http://clusterlabs.org/doc/Colocation\_Explained.pdf].

#### **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 the section called "Specifying the Order in which Resources Should Start/Stop") *and* a colocation constraint, or alternatively, a group (see the section called "Groups - A Syntactic Shortcut").

### **Colocation Properties**

<sup>&</sup>lt;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 5.3. Attributes of a rsc colocation Constraint

| Field              | Default | Description  |  |
|--------------------|---------|--|--|
| id                 |         | A unique name for the constraint (required).   |  |
| rsc                |         | The name of a resource that should be located relative to with-rsc (required).   |  |
| 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 rsc (required).  |  |
| node-<br>attribute | #uname  | The node attribute that must be the same on the node running rsc and the node running with-rsc for the constraint to be satisfied. (For details, see the section called "Colocation by Node Attribute".) |  |
| score              |         | 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".              |  |

# **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:

#### **Example 5.5. 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".

#### **Example 5.6. 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.

#### **Example 5.7. 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.

### **Resource Sets**

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

#### Example 5.8. A set of 3 resources

Resource sets are valid inside rsc\_location, rsc\_order (see the section called "Ordering Sets of Resources"), rsc\_colocation (see the section called "Colocating Sets of Resources"), and rsc\_ticket (see the section called "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 5.4. Attributes of a resource set Element

| Field       | Default | Description   |  |
|-------------|---------|---|--|
| id          |         | A unique name for the set   |  |
| sequential  | true    | Whether the members of the set must be acted on in order.  Meaningful within rsc_order and rsc_colocation.  |  |
| require-all | 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 rsc_order. |  |
| role        |         | Limit the effect of the constraint to the specified role.  Meaningful within rsc_location, rsc_colocation and rsc_ticket.   |  |

| Field  | Default | Description   |  |
|--------|---------|---|--|
| action |         | Limit the effect of the constraint to the specified action.  Meaningful within rsc_order. |  |
| score  |         | Advanced use only. Use a specific score for this set within the constraint.               |  |

# **Ordering Sets of Resources**

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

#### Example 5.9. 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>
```

Figure 5.1. Visual representation of the four resources' start order for the above constraints

#### **Ordered Set**

To simplify this situation, resource sets (see the section called "Resource Sets") can be used within ordering constraints:

#### Example 5.10. A chain of ordered resources expressed as a set

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.

#### Example 5.11. Ordered sets of unordered resources

Figure 5.2. 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.

# Example 5.12. 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>
```

Figure 5.3. 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.

# Example 5.13. 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.

# **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 the section called "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.

# Example 5.14. Colocation chain as individual constraints, where A is placed first, then B, then C, then D

```
<constraints>
```

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

#### Example 5.15. Equivalent colocation chain expressed using resource set

#### 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.

#### Example 5.16. Using colocated sets to specify a shared dependency

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

#### 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.

# Example 5.17. Colocation in which the members of the middle set have no interdependencies, and the last set listed applies only to instances in the master role

```
<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="Master">
        <resource_ref id="A"/>
        <resource_ref id="B"/>
      </resource set>
    </rsc_colocation>
</constraints>
```

Figure 5.4. 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.

# **Chapter 6. Fencing**

#### **Table of Contents**

| 42 |
|----|
| 42 |
| 43 |
| 43 |
| 43 |
| 43 |
| 44 |
| 48 |
| 48 |
| 49 |
| 5( |
| 52 |
| 53 |
| 58 |
|    |

# 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).

# 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 also is used when a resource cannot otherwise be stopped. If a failed resource fails to stop, it cannot be recovered elsewhere. Fencing the resource's node is the only way to ensure the resource is recoverable.

Users may also configure the on-fail property of any resource operation to fencing, in which case the cluster will fence the resource's node if the operation fails.

### **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).

# **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.

# 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 targetrole to Stopped, then no node can use that device. If mandatory location constraints prevent 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. The best practice is to make the constraint optional (i.e. a finite negative score rather than -INFINITY), so that the node can fence itself if no other nodes can.

# **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 may only be evaluated 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 significant user demand.

# **Special Options 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 6.1. Additional Properties of Fencing Resources** 

| Field           | Type   | Default  | Description  |
|-----------------|--------|--|--|
| stonith-timeout | NA     | NA   | Older versions used this to override the default period to wait for a STONITH (reboot, on, off) action to complete for this device. It has been replaced by the pcmk_reboot_timeout and pcmk_off_timeout properties.   |
| provides        | string |  | Any special capability provided by the fence device. Currently, only one such capability is meaningful: unfencing (see the section called "Unfencing").  |
| pcmk_host_map   | string |  | A mapping of host names to ports numbers for devices that do not support host names. Example: node1:1;node2:2,3 tells the cluster to use port 1 for node1 and ports 2 and 3 for node2. If pcmk_host_check is explicitly set to static-list, either this or pcmk_host_list must be set. |
| pcmk_host_list  | string |  | A list of machines controlled by this device. If pcmk_host_check is explicitly set to static-list, either this or pcmk_host_map must be set.   |
| pcmk_host_check | string | A value<br>appropriate to<br>other configuration<br>options and device<br>capabilities (see<br>note below) | How to determine which machines are controlled by the device. Allowed values:  • dynamic-list: query the device via the "list" command  • static-list: check the pcmk_host_list or pcmk_host_map attribute   |

| Field              | Type    | Default   | Description  |
|--------------------|---------|---|--|
|                    |         |   | status: query the device via the     "status" command      none: assume every device can     fence every machine   |
| pcmk_delay_max     | time    | 0s  | Enable a random delay of up to the time specified before executing fencing actions. This is sometimes used in two-node clusters to ensure that the nodes don't fence each other at the same time. The overall delay introduced by pacemaker is derived from this random delay value adding a static delay so that the sum is kept below the maximum delay.   |
| pcmk_delay_base    | time    | 0s  | Enable a static delay before executing fencing actions. This can be used e.g. 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. |
| pcmk_action_limit  | integer | 1   | The maximum number of actions that can be performed in parallel on this device, if the cluster option concurrent-fencing is true1 is unlimited.  |
| pcmk_host_argument | string  | port otherwise<br>plug if supported<br>according to the<br>metadata of the<br>fence agent | Advanced use only. Which parameter should be supplied to the fence agent to identify the node to be fenced. Some devices support neither the standard plug nor the deprecated port parameter, or may provide additional ones. Use this to specify an alternate, device-specific parameter. A value of none tells the cluster not to supply any additional parameters.  |

| Field              | Type    | Default | Description  |
|--------------------|---------|---------|--|
| pcmk_reboot_action | string  | reboot  | Advanced use only. 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.   |
| pcmk_reboot_timeou | time    | 60s     | Advanced use only. Specify an alternate timeout to use for reboot actions instead of the value of stonithtimeout. Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.   |
| pcmk_reboot_retrie | integer | 2       | Advanced use only. The maximum number of times to retry the reboot 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. |
| pcmk_off_action    | string  | off     | Advanced use only. 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.  |
| pcmk_off_timeout   | time    | 60s     | Advanced use only. Specify an alternate timeout to use for off actions instead of the value of stonith-timeout. Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.   |
| pcmk_off_retries   | integer | 2       | Advanced use only. The maximum number of times to retry the off 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.    |

| Field               | Type                 | Default | Description   |
|---------------------|----------------------|---------|---|
| pcmk_list_action    | string               | list    | Advanced use only. 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.   |
| pcmk_list_timeout   | time                 | 60s     | Advanced use only. Specify an alternate timeout to use for list actions instead of the value of stonith-timeout. Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.   |
| pcmk_list_retries   | integer              | 2       | Advanced use only. The maximum number of times to retry the list 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.    |
| pcmk_monitor_action | entring              | monitor | Advanced use only. 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.   |
| pcmk_monitor_timed  | wime                 | 60s     | Advanced use only. Specify an alternate timeout to use for monitor actions instead of the value of stonith—timeout. Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.  |
| pcmk_monitor_retri  | <del>खंड</del> ieger | 2       | Advanced use only. The maximum number of times to retry the monitor 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. |

| Field              | Type     | Default | Description  |
|--------------------|----------|---------|--|
| pcmk_status_action | string   | status  | Advanced use only. 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.   |
| pcmk_status_timeou | time     | 60s     | Advanced use only. Specify an alternate timeout to use for status actions instead of the value of stonith—timeout. Some devices need much more or less time to complete than normal. Use this to specify an alternate, device-specific timeout.  |
| pcmk_status_retrie | sinteger | 2       | Advanced use only. The maximum number of times to retry the status 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. |

#### Note

The default value for pcmk\_host\_check is static-list if either pcmk\_host\_list or pcmk\_host\_map is configured. If neither of those are configured, the default is dynamic-list if the fence device supports the list action, or status if the fence device supports the status action but not the list action. If none of those conditions apply, the default is none.

# **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.

# **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 dependedon 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.

# **Configuring Fencing**

1. Find the correct driver:

```
# stonith admin --list-installed
```

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
```

- 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 man pacemaker-fenced 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 man pacemaker-fenced for details.
- 6. 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 man pacemaker-fenced for details.
- 7. Upload it into the CIB using cibadmin:

```
# cibadmin -C -o resources --xml-file stonith.xml
```

8. Set stonith-enabled to true:

```
# crm_attribute -t crm_config -n stonith-enabled -v true
```

9. Once the stonith resource is running, you can test it by executing the following (although you might want to stop the cluster on that machine first):

```
# stonith_admin --reboot nodename
```

# **Example Fencing Configuration**

Assume we have a chassis containing four nodes and an IPMI device active on 192.0.2.1. We would choose the fence\_ipmilan driver, and obtain the following list of parameters:

#### **Example 6.1. Obtaining a list of Fence Agent Parameters**

```
# stonith_admin --metadata -a fence_ipmilan
<resource-agent name="fence_ipmilan" shortdesc="Fence agent for IPMI over LAN">
  <symlink name="fence_ilo3" shortdesc="Fence agent for HP iLO3"/>
  <symlink name="fence_ilo4" shortdesc="Fence agent for HP iLO4"/>
  <symlink name="fence_idrac" shortdesc="Fence agent for Dell iDRAC"/>
  <symlink name="fence_imm" shortdesc="Fence agent for IBM Integrated Management M
  <longdesc>
  </longdesc>
  <vendor-url>
  </vendor-url>
  <parameters>
    <parameter name="auth" unique="0" required="0">
      <getopt mixed="-A"/>
      <content type="string"/>
      <shortdesc>
      </shortdesc>
    </parameter>
    <parameter name="ipaddr" unique="0" required="1">
      <getopt mixed="-a"/>
      <content type="string"/>
      <shortdesc>
      </shortdesc>
    </parameter>
    <parameter name="passwd" unique="0" required="0">
      <getopt mixed="-p"/>
      <content type="string"/>
      <shortdesc>
      </shortdesc>
    </parameter>
    <parameter name="passwd_script" unique="0" required="0">
      <getopt mixed="-S"/>
      <content type="string"/>
      <shortdesc>
      </shortdesc>
    </parameter>
    <parameter name="lanplus" unique="0" required="0">
      <getopt mixed="-P"/>
      <content type="boolean"/>
```

```
<shortdesc>
  </shortdesc>
</parameter>
<parameter name="login" unique="0" required="0">
  <getopt mixed="-l"/>
  <content type="string"/>
 <shortdesc>
  </shortdesc>
</parameter>
<parameter name="action" unique="0" required="0">
  <getopt mixed="-o"/>
  <content type="string" default="reboot"/>
  <shortdesc>
  </shortdesc>
</parameter>
<parameter name="timeout" unique="0" required="0">
  <getopt mixed="-t"/>
 <content type="string"/>
 <shortdesc>
  </shortdesc>
</parameter>
<parameter name="cipher" unique="0" required="0">
  <getopt mixed="-C"/>
  <content type="string"/>
 <shortdesc>
  </shortdesc>
</parameter>
<parameter name="method" unique="0" required="0">
  <getopt mixed="-M"/>
 <content type="string" default="onoff"/>
  <shortdesc>
  </shortdesc>
</parameter>
<parameter name="power_wait" unique="0" required="0">
  <getopt mixed="-T"/>
 <content type="string" default="2"/>
  <shortdesc>
  </shortdesc>
</parameter>
<parameter name="delay" unique="0" required="0">
 <getopt mixed="-f"/>
  <content type="string"/>
  <shortdesc>
  </shortdesc>
</parameter>
<parameter name="privlvl" unique="0" required="0">
  <getopt mixed="-L"/>
 <content type="string"/>
  <shortdesc>
  </shortdesc>
</parameter>
<parameter name="verbose" unique="0" required="0">
  <getopt mixed="-v"/>
  <content type="boolean"/>
```

```
<shortdesc>
      </shortdesc>
    </parameter>
  </parameters>
  <actions>
    <action name="on"/>
    <action name="off"/>
    <action name="reboot"/>
    <action name="status"/>
    <action name="diag"/>
    <action name="list"/>
    <action name="monitor"/>
    <action name="metadata"/>
    <action name="stop" timeout="20s"/>
    <action name="start" timeout="20s"/>
  </actions>
</resource-agent>
```

Based on that, we would create a fencing resource fragment that might look like this:

#### **Example 6.2. An IPMI-based Fencing Resource**

# **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 6.2. Properties of Fencing Levels** 

| Field                | Description  |
|----------------------|--|
| id                   | A unique name for the level  |
| target               | The name of a single node to which this level applies  |
| target-pattern       | An extended regular expression (as defined in POSIX [http://pubs.opengroup.org/onlinepubs/9699919799/basedefs/V1_chap09.html#tag_09_04]) matching the names of nodes to which this level applies |
| target-<br>attribute | The name of a node attribute that is set (to target-value) for nodes to which this level applies   |
| target-value         | The node attribute value (of target-attribute) 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  |

#### Example 6.3. Fencing topology with different devices for different nodes

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

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

- 3 nodes (2 active prod-mysql nodes, 1 prod\_mysql-rep in standby for quorum purposes)
- the active nodes have an IPMI-controlled power board reached at 192.0.2.1 and 192.0.2.2
- the active nodes also have two independent PSUs (Power Supply Units) connected to two independent PDUs (Power Distribution Units) reached at 198.51.100.1 (port 10 and port 11) and 203.0.113.1 (port 10 and port 11)
- the first fencing method uses the fence\_ipmi agent
- the second fencing method uses the fence\_apc\_snmp agent targetting 2 fencing devices (one per PSU, either port 10 or 11)
- fencing is only implemented for the active nodes and has location constraints
- · fencing topology is set to try IPMI fencing first then default to a "sure-kill" dual PDU fencing

In a normal failure scenario, STONITH will first select fence\_ipmi to try to kill the faulty node. Using a fencing topology, if that first method fails, STONITH will then move on to selecting fence\_apc\_snmp twice:

- · once for the first PDU
- · again for the second PDU

</instance attributes>

</primitive>

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

**First fencing method: single IPMI device.** Each cluster node has it own dedicated IPMI channel that can be called 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" v</pre>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-action" name="action" v</pre>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-login" name="login" val</pre>
    <nvpair id="fence prod-mysql1 ipmi-instance attributes-passwd" name="passwd" v</pre>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-verbose" name="verbose"</pre>
    <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-pcmk_host_list" name="p</pre>
    <nvpair id="fence prod-mysql1 ipmi-instance attributes-lanplus" name="lanplus"</pre>
  </instance attributes>
</primitive>
<primitive class="stonith" id="fence_prod-mysq12_ipmi" type="fence_ipmilan">
  <instance_attributes id="fence_prod-mysql2_ipmi-instance_attributes">
    <nvpair id="fence_prod-mysq12_ipmi-instance_attributes-ipaddr" name="ipaddr" v</pre>
    <nvpair id="fence prod-mysql2 ipmi-instance attributes-action" name="action" v</pre>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-login" name="login" val</pre>
    <nvpair id="fence prod-mysql2 ipmi-instance attributes-passwd" name="passwd" v</pre>
    <nvpair id="fence_prod-mysq12_ipmi-instance_attributes-verbose" name="verbose"</pre>
    <nvpair id="fence_prod-mysql2_ipmi-instance_attributes-pcmk_host_list" name="p</pre>
    <nvpair id="fence prod-mysql2 ipmi-instance attributes-lanplus" name="lanplus"</pre>
```

**Second fencing method: dual PDU devices.** Each cluster node also has two distinct power channels controlled by two distinct PDUs. That means a total of 4 fencing devices configured as follows:

```
• Node 1, PDU 1, PSU 1 @ port 10
```

- Node 1, PDU 2, PSU 2 @ port 10
- Node 2, PDU 1, PSU 1 @ port 11
- Node 2, PDU 2, PSU 2 @ port 11

The matching fencing agents are configured as follows:

```
<primitive class="stonith" id="fence_prod-mysql1_apc1" type="fence_apc_snmp">
  <instance_attributes id="fence_prod-mysql1_apcl-instance_attributes">
    <nvpair id="fence_prod-mysql1_apc1-instance_attributes-ipaddr" name="ipaddr" v</pre>
    <nvpair id="fence_prod-mysql1_apc1-instance_attributes-action" name="action" v</pre>
    <nvpair id="fence_prod-mysql1_apc1-instance_attributes-port" name="port" value</pre>
    <nvpair id="fence_prod-mysql1_apc1-instance_attributes-login" name="login" val</pre>
    <nvpair id="fence_prod-mysqll_apcl-instance_attributes-passwd" name="passwd" v</pre>
    <nvpair id="fence_prod-mysql1_apc1-instance_attributes-pcmk_host_list" name="p</pre>
  </instance_attributes>
</primitive>
<primitive class="stonith" id="fence_prod-mysql1_apc2" type="fence_apc_snmp">
  <instance_attributes id="fence_prod-mysql1_apc2-instance_attributes">
    <nvpair id="fence_prod-mysql1_apc2-instance_attributes-ipaddr" name="ipaddr" v</pre>
    <nvpair id="fence_prod-mysql1_apc2-instance_attributes-action" name="action" v</pre>
    <nvpair id="fence_prod-mysql1_apc2-instance_attributes-port" name="port" value</pre>
    <nvpair id="fence_prod-mysql1_apc2-instance_attributes-login" name="login" val</pre>
    <nvpair id="fence_prod-mysql1_apc2-instance_attributes-passwd" name="passwd" v</pre>
    <nvpair id="fence_prod-mysql1_apc2-instance_attributes-pcmk_host_list" name="p</pre>
  </instance_attributes>
</primitive>
<primitive class="stonith" id="fence_prod-mysql2_apc1" type="fence_apc_snmp">
  <instance_attributes id="fence_prod-mysql2_apc1-instance_attributes">
    <nvpair id="fence_prod-mysql2_apc1-instance_attributes-ipaddr" name="ipaddr" v</pre>
    <nvpair id="fence_prod-mysql2_apcl-instance_attributes-action" name="action" v</pre>
    <nvpair id="fence_prod-mysql2_apcl-instance_attributes-port" name="port" value</pre>
    <nvpair id="fence_prod-mysql2_apc1-instance_attributes-login" name="login" val</pre>
    <nvpair id="fence_prod-mysql2_apcl-instance_attributes-passwd" name="passwd" v</pre>
    <nvpair id="fence_prod-mysql2_apc1-instance_attributes-pcmk_host_list" name="p</pre>
  </instance_attributes>
</primitive>
<pri><primitive class="stonith" id="fence_prod-mysql2_apc2" type="fence_apc_snmp">
  <instance_attributes id="fence_prod-mysql2_apc2-instance_attributes">
    <nvpair id="fence_prod-mysq12_apc2-instance_attributes-ipaddr" name="ipaddr" v</pre>
    <nvpair id="fence_prod-mysql2_apc2-instance_attributes-action" name="action" v</pre>
    <nvpair id="fence_prod-mysq12_apc2-instance_attributes-port" name="port" value</pre>
    <nvpair id="fence_prod-mysq12_apc2-instance_attributes-login" name="login" val</pre>
    <nvpair id="fence_prod-mysq12_apc2-instance_attributes-passwd" name="passwd" v</pre>
    <nvpair id="fence_prod-mysql2_apc2-instance_attributes-pcmk_host_list" name="p</pre>
  </instance_attributes>
</primitive>
```

**Location Constraints.** To prevent STONITH from trying to run a fencing agent on the same node it is supposed to fence, constraints are placed on all the fencing primitives:

<constraints>

**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 devices="fence_prod-mysql1_ipmi" id="fencing-2" index="1" target=
    <fencing-level devices="fence_prod-mysql1_apc1,fence_prod-mysql1_apc2" id="fenci
    <fencing-level devices="fence_prod-mysql2_ipmi" id="fencing-0" index="1" target=
    <fencing-level devices="fence_prod-mysql2_apc1,fence_prod-mysql2_apc2" id="fenci
</fencing-topology>
```

Please note, in fencing-topology, the lowest index value determines the priority of the first fencing method.

**Final configuration.** Put together, the configuration looks like this:

```
<cib admin_epoch="0" crm_feature_set="3.0.7" epoch="292" have-quorum="1" num_updat</pre>
  <configuration>
    <crm config>
      <cluster_property_set id="cib-bootstrap-options">
        <nvpair id="cib-bootstrap-options-stonith-enabled" name="stonith-enabled"</pre>
        <nvpair id="cib-bootstrap-options-stonith-action" name="stonith-action" va</pre>
        <nvpair id="cib-bootstrap-options-expected-quorum-votes" name="expected-qu</pre>
      </cluster_property_set>
    </crm_config>
    <nodes>
      <node id="prod-mysql1" uname="prod-mysql1">
      <node id="prod-mysql2" uname="prod-mysql2"/>
      <node id="prod-mysql-rep1" uname="prod-mysql-rep1"/>
        <instance_attributes id="prod-mysql-rep1">
          <nvpair id="prod-mysql-rep1-standby" name="standby" value="on"/>
        </instance_attributes>
      </node>
    </nodes>
    <resources>
      <pri><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="ipa</pre>
          <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-action" name="act</pre>
          <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-login" name="logi</pre>
          <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-passwd" name="pas</pre>
          <nvpair id="fence_prod-mysq11_ipmi-instance_attributes-verbose" name="ve</pre>
          <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-pcmk_host_list" n</pre>
          <nvpair id="fence_prod-mysql1_ipmi-instance_attributes-lanplus" name="la</pre>
        </instance_attributes>
      </primitive>
```

```
<primitive class="stonith" id="fence_prod-mysq12_ipmi" type="fence_ipmilan">
     <instance attributes id="fence prod-mysql2 ipmi-instance attributes">
       <nvpair id="fence_prod-mysq12_ipmi-instance_attributes-ipaddr" name="ipa</pre>
       <nvpair id="fence_prod-mysq12_ipmi-instance_attributes-action" name="act</pre>
       <nvpair id="fence_prod-mysq12_ipmi-instance_attributes-login" name="logi</pre>
       <nvpair id="fence_prod-mysq12_ipmi-instance_attributes-passwd" name="pas</pre>
       <nvpair id="fence_prod-mysq12_ipmi-instance_attributes-verbose" name="ve</pre>
       <nvpair id="fence prod-mysq12 ipmi-instance attributes-pcmk host list" n</pre>
       <nvpair id="fence_prod-mysq12_ipmi-instance_attributes-lanplus" name="la</pre>
     </instance_attributes>
   </primitive>
   <pri><primitive class="stonith" id="fence_prod-mysql1_apc1" type="fence_apc_snmp"</pr>
     <instance attributes id="fence prod-mysql1 apc1-instance attributes">
       <nvpair id="fence_prod-mysql1_apc1-instance_attributes-ipaddr" name="ipa</pre>
       <nvpair id="fence_prod-mysql1_apc1-instance_attributes-action" name="act</pre>
       <nvpair id="fence_prod-mysql1_apc1-instance_attributes-port" name="port"</pre>
       <nvpair id="fence_prod-mysql1_apc1-instance_attributes-login" name="logi</pre>
       <nvpair id="fence_prod-mysql1_apc1-instance_attributes-passwd" name="pas</pre>
       <nvpair id="fence_prod-mysql1_apc1-instance_attributes-pcmk_host_list" n</pre>
     </instance attributes>
   </primitive>
   <pri><primitive class="stonith" id="fence_prod-mysql1_apc2" type="fence_apc_snmp"</pr>
     <instance_attributes id="fence_prod-mysql1_apc2-instance_attributes">
       <nvpair id="fence prod-mysgl1 apc2-instance attributes-ipaddr" name="ipa</pre>
       <nvpair id="fence_prod-mysq11_apc2-instance_attributes-action" name="act</pre>
       <nvpair id="fence_prod-mysq11_apc2-instance_attributes-port" name="port"</pre>
       <nvpair id="fence_prod-mysql1_apc2-instance_attributes-login" name="logi</pre>
       <nvpair id="fence_prod-mysq11_apc2-instance_attributes-passwd" name="pas</pre>
       <nvpair id="fence_prod-mysq11_apc2-instance_attributes-pcmk_host_list" n</pre>
     </instance attributes>
   </primitive>
   <pri><primitive class="stonith" id="fence_prod-mysql2_apc1" type="fence_apc_snmp"</pr>
     <instance_attributes id="fence_prod-mysql2_apc1-instance_attributes">
       <nvpair id="fence_prod-mysq12_apc1-instance_attributes-ipaddr" name="ipa</pre>
       <nvpair id="fence_prod-mysq12_apc1-instance_attributes-action" name="act</pre>
       <nvpair id="fence_prod-mysq12_apc1-instance_attributes-port" name="port"</pre>
       <nvpair id="fence_prod-mysq12_apc1-instance_attributes-login" name="logi</pre>
       <nvpair id="fence_prod-mysq12_apc1-instance_attributes-passwd" name="pas</pre>
       <nvpair id="fence_prod-mysq12_apc1-instance_attributes-pcmk_host_list" n</pre>
     </instance_attributes>
   </primitive>
   <pri><primitive class="stonith" id="fence_prod-mysq12_apc2" type="fence_apc_snmp"</pr>
     <instance_attributes id="fence_prod-mysql2_apc2-instance_attributes">
       <nvpair id="fence_prod-mysq12_apc2-instance_attributes-ipaddr" name="ipa</pre>
       <nvpair id="fence_prod-mysq12_apc2-instance_attributes-action" name="act</pre>
       <nvpair id="fence_prod-mysq12_apc2-instance_attributes-port" name="port"</pre>
       <nvpair id="fence_prod-mysq12_apc2-instance_attributes-login" name="logi</pre>
       <nvpair id="fence prod-mysgl2 apc2-instance attributes-passwd" name="pas</pre>
       <nvpair id="fence_prod-mysq12_apc2-instance_attributes-pcmk_host_list" n</pre>
     </instance attributes>
   </primitive>
</resources>
 <constraints>
   <rsc_location id="l_fence_prod-mysql1_ipmi" node="prod-mysql1" rsc="fence_pr
```

# **Remapping Reboots**

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

- 1. If the chosen fencing device does not support the reboot command, the cluster will ask it to perform off instead.
- 2. 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).

## Chapter 7. Alerts

### **Table of Contents**

| Alert Agents                  | 59 |
|-------------------------------|----|
| Alert Recipients              | 59 |
| Alert Meta-Attributes         | 60 |
| Alert Instance Attributes     | 61 |
| Alert Filters                 | 61 |
| Using the Sample Alert Agents | 62 |
| Writing an Alert Agent        |    |

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

### **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.).

### **Example 7.1. Simple alert 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.

### **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.

#### **Example 7.2. Alert configuration with recipient**

```
</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.

### **Alert Meta-Attributes**

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

Table 7.1. Meta-Attributes of an Alert

| Meta-Attribute       | Default       | Description   |
|----------------------|---------------|---|
| timestamp-<br>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 date(1) 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.

### **Example 7.3. 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"</pre>
                    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"</pre>
                        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"</pre>
                         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.

### **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.

#### **Example 7.4. Alert configuration with instance attributes**

### **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.

### Example 7.5. Alert configuration to receive only node events and fencing events

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.

### Example 7.6. Alert configuration to be called when certain node attributes change

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.

## **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.

#### Example 7.7. Sending cluster events as SNMP traps

#### Example 7.8. Sending cluster events as e-mails

## **Writing an Alert Agent**

Table 7.2. Environment variables passed to alert agents

| Environment Variable    | Description   |
|-------------------------|---|
| CRM_alert_kind          | The type of alert (node, fencing, resource, or attribute)   |
| CRM_alert_version       | The version of Pacemaker sending the alert  |
| CRM_alert_recipient     | The configured recipient  |
| 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_timestamp     | A timestamp created prior to executing the agent, in the format specified by the timestamp-format 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_ep  | <b>Thh</b> same time as CRM_alert_timestamp, expressed as the integer number of seconds since January 1, 1970. This (along with CRM_alert_timestamp_usec) 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_us  | The same time as CRM_alert_timestamp, expressed as the integer number of microseconds since CRM_alert_timestamp_epoch.  |
| CRM_alert_node          | Name of affected node   |
| CRM_alert_desc          | Detail about event. For node alerts, this is the node's current state (member or lost). For fencing alerts, this is a summary of the requested fencing operation, including origin, target, and fencing operation error code, if any. For resource alerts, this is a readable string equivalent of CRM_alert_status.                                    |
| CRM_alert_nodeid        | ID of node whose status changed (provided with node alerts only)  |
| CRM_alert_task          | The requested fencing or resource operation (provided with fencing and resource alerts only)  |
| CRM_alert_rc            | The numerical return code of the fencing or resource operation (provided with fencing and resource alerts only)   |
| CRM_alert_rsc           | The name of the affected resource (resource alerts only)  |

| Environment Variable   | Description   |
|------------------------|---|
| CRM_alert_interval     | The interval of the resource operation (resource alerts only)   |
| CRM_alert_target_rc    | The expected numerical return code of the operation (resource alerts only)  |
| CRM_alert_status       | A numerical code used by Pacemaker to represent the operation result (resource alerts only)   |
| CRM_alert_exec_time    | The (wall-clock) time, in milliseconds, that it took to execute the action. If the action timed out, CRM_alert_status will be 2, CRM_alert_desc will be "Timed Out", and this value will be the action timeout. May not be supported on all platforms. (resource alerts only) |
| CRM_alert_attribute_na | in the name of the node attribute that changed (attribute alerts only)  |
| CRM_alert_attribute_va | <b>The</b> new value of the node attribute that changed (attribute 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

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.

## **Chapter 8. Rules**

### **Table of Contents**

| Rule Properties                               |
|---|
| Node Attribute Expressions                    |
| Date/Time Expressions                         |
| Date Specifications                           |
| Durations                                     |
| Example Time-Based Expressions                |
| Resource Expressions                          |
| Example Resource-Based Expressions            |
| Operation Expressions                         |
| Example Operation-Based Expressions           |
| Using Rules to Determine Resource Location    |
| Location Rules Based on Other Node Properties |
| Using score-attribute Instead of score        |
| Using Rules to Define Options                 |
| Using Rules to Control Resource Options       |
| Using Rules to Control Resource Defaults      |
| Using Rules to Control Cluster Options        |

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 metaattributes, 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 section will make this clearer.

## **Rule Properties**

Table 8.1. Attributes of a rule Element

| Attribute | Default | Description                           |
|-----------|---------|---------------------------------------|
| id        |         | A unique name for the rule (required) |

| Attribute       | Default | Description  |
|-----------------|---------|--|
| role            | Started | The rule is in effect only when the resource is in the specified role. Allowed values are Started, Slave, and Master. A rule with role="Master" 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 score and score-attribute 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 score and scoreattribute may be used.   |
| boolean-op      | and     | If this rule contains more than one condition, a value of and specifies that the rule evaluates to true only if all conditions are true, and a value of or 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.

## **Node Attribute Expressions**

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

Table 8.2. Attributes of an expression Element

| Field     | Default  | Description   |
|-----------|--|---|
| id        |  | A unique name for the expression (required)   |
| attribute |  | The node attribute to test (required)   |
| type      | The default type for lt, gt, lte, and gte operations is number if either value contains a decimal point character, or integer otherwise. The | How the node attributes should be compared. Allowed values are string, integer, number, and version. integer truncates floating-point values if necessary before performing an integer comparison. number performs a floating-point comparison. |

| Field        | Default     | Description  |
|--------------|-------------|--|
|              | default     |  |
|              | type for    |  |
|              | all other   |  |
|              | operations  |  |
|              | is          |  |
|              | string.     |  |
|              | If a        |  |
|              | numeric     |  |
|              | parse fails |  |
|              | for either  |  |
|              | value, then |  |
|              | the values  |  |
|              | compared    |  |
|              | as type     |  |
|              | string.     |  |
| operation    |             | The comparison to perform (required). Allowed values:  |
|              |             | • lt: True if the node attribute value is less than the comparison value                                     |
|              |             | • gt: True if the node attribute value is greater than the comparison value                                  |
|              |             | • lte: True if the node attribute value is less than or equal to the comparison value                        |
|              |             | • gte: True if the node attribute value is greater than or equal to the comparison value                     |
|              |             | • eq: True if the node attribute value is equal to the comparison value                                      |
|              |             | • ne: True if the node attribute value is not equal to the comparison value                                  |
|              |             | • defined: True if the node has the named attribute  |
|              |             | <ul> <li>not_defined: True if the node does not have the<br/>named attribute</li> </ul>                      |
| value        |             | User-supplied value for comparison (required for operations other than defined and not_defined)              |
| value-source | literal     | How the value is derived. Allowed values:  |
|              |             | • literal: value is a literal string to compare against  |
|              |             | param: value is the name of a resource parameter to compare against (only valid in location constraints)     |
|              |             | meta: value is the name of a resource meta-attribute to compare against (only valid in location constraints) |

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 8.3. Built-in Node Attributes** 

| Name              | Value   |
|-------------------|---|
| #uname            | Node name   |
| #id               | Node ID   |
| #kind             | Node type. Possible values are cluster, remote, and container.  Kind is remote for Pacemaker Remote nodes created with the ocf:pacemaker:remote resource, and container for Pacemaker Remote guest nodes and bundle nodes |
| #is_dc            | "true" if this node is a Designated Controller (DC), "false" otherwise  |
| #cluster-<br>name | The value of the cluster-name cluster property, if set  |
| #site-name        | The value of the site-name node attribute, if set, otherwise identical to #cluster-name   |
| #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.   |

## **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 8.4. Attributes of a date\_expression Element

| Field     | Description  |
|-----------|--|
| id        | A unique name for the expression (required)  |
| start     | A date/time conforming to the ISO8601 [http://en.wikipedia.org/wiki/ISO_8601] specification. May be used when operation is in_range (in which case at least one of start or end must be specified) or gt (in which case start is required).  |
| end       | A date/time conforming to the ISO8601 [http://en.wikipedia.org/wiki/ISO_8601] specification. May be used when operation is in_range (in which case at least one of start or end must be specified) or lt (in which case end is required).  |
| operation | Compares the current date/time with the start and/or end date, depending on the context. Allowed values:  • gt: True if the current date/time is after start  • lt: True if the current date/time is before end  • in_range: True if the current date/time is after start (if specified) and before either end (if specified) or start plus the value of the duration element (if one is contained in the date_expression) |

| Field | Description   |
|-------|---|
|       | • date_spec: True if the current date/time matches the specification given in the contained date_spec element (described below) |

#### **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 8.5. Attributes of a date\_spec Element

| Field     | Description   |  |
|-----------|---|--|
| id        | A unique name for the object (required)   |  |
| hours     | Allowed values: 0-23 (where 0 is midnight and 23 is 11 p.m.)  |  |
| monthdays | Allowed values: 1-31 (depending on month and year)  |  |
| weekdays  | Allowed values: 1-7 (where 1 is Monday and 7 is Sunday)   |  |
| yeardays  | Allowed values: 1-366 (depending on the year)   |  |
| months    | Allowed values: 1-12  |  |
| weeks     | Allowed values: 1-53 (depending on weekyear)  |  |
| years     | Year according to the Gregorian calendar  |  |
| weekyears | 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 years="2005" or weekyears="2004" |  |
| moon      | 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 8.6. Attributes of a duration Element

| Field   | Description   |  |
|---------|---|--|
| id      | A unique name for this duration 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:

#### Example 8.1. True if now is any time in the year 2005

#### Example 8.2. Equivalent expression

#### Example 8.3. 9am-5pm Monday-Friday

```
</rule>
```

Please note that the 16 matches up to 16:59:59, as the numeric value (hour) still matches!

### Example 8.4. 9am-6pm Monday through Friday or anytime Saturday

#### Example 8.5. 9am-5pm or 9pm-12am Monday through Friday

#### **Example 8.6. Mondays in March 2005**

#### 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="2005-03-31T23:59:59" to avoid confusion.

#### Example 8.7. A full moon on Friday the 13th

</rule>

## **Resource Expressions**

An rsc\_expression 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 8.7. Attributes of an rsc\_expression Element

| Field    | Description  |  |
|----------|--|--|
| id       | A unique name for the expression (required)  |  |
| class    | The standard name to be matched against resource agents  |  |
| provider | If given, the vendor to be matched against resource agents. This only makes sense for agents using the OCF spec. |  |
| 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:

#### Example 8.8. True for all ocf:heartbeat:IPaddr2 resources

#### Example 8.9. Provider doesn't apply to non-OCF resources

## **Operation Expressions**

An op\_expression is a rule condition based on an action of some resource agent. This rule is only valid within an op\_defaults context.

Table 8.8. Attributes of an op\_expression Element

| Field    | Description  |  |
|----------|--|--|
| id       | A unique name for the expression (required)  |  |
| name     | The action name to match against. This can be any action supported by the resource agent; common values include monitor, start, and stop (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:

### **Example 8.10. True for all monitor actions**

#### Example 8.11. True for all monitor actions with a 10 second interval

### **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:

### Example 8.12. Prevent resource "webserver" from running on node3

The constraint can be more verbosely written using a rule:

#### Example 8.13. Prevent resource "webserver" from running on node3 using rule

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:

#### Example 8.14. A sample nodes section for use with score-attribute

<nodes>

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

### **Using Rules to Define Options**

Rules may be used to control a variety of options:

- Cluster options (cluster\_property\_set elements)
- Node attributes (as instance\_attributes or utilization elements inside a node element)
- Resource options (as 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 inside an op or op\_defaults element)

### **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.

#### Example 8.15. Defining different resource options based on the node name

```
id="mySpecialRsc" class="ocf" type="Special" provider="me">
```

```
<instance_attributes id="special-node1" score="3">
    <rule id="nodel-special-case" score="INFINITY" >
     <expression id="nodel-special-case-expr" attribute="#uname"</pre>
       operation="eq" value="node1"/>
    </rule>
    <nvpair id="nodel-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"</pre>
       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, score defaults to zero, and objects with an equal score are processed in 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:

- 1. special-nodel has the highest score (3) and so is evaluated first; its rule evaluates to true, so interface is set to ethl.
- 2. special-node2 is evaluated next with score 2, but its rule evaluates to false, so it is ignored.
- 3. 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.

#### Example 8.16. Change resource-stickiness during working hours

```
</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.

#### Example 8.17. Set all IPaddr2 resources to stopped

#### Example 8.18. Set all monitor action timeouts to 7 seconds

## Example 8.19. Set the monitor action timeout on all IPaddr2 resources with a given monitor interval to 8 seconds

### **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 difference is that because they are cluster options, one cannot (or should not, because they won't work) use attribute-based expressions. 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.

#### Example 8.20. Schedule a maintenance window for 9 to 11 p.m. CDT Sept. 20, 2019

### **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.

## **Chapter 9. Advanced Configuration**

### **Table of Contents**

| Specifying When Recurring Actions are Performed | 78 |
|---|----|
| Handling Resource Failure                       | 78 |
| Failure Counts                                  | 78 |
| Failure Response                                | 79 |
| Moving Resources                                | 80 |
| Moving Resources Manually                       | 80 |
| Moving Resources Due to Connectivity Changes    | 82 |
| Migrating Resources                             | 84 |
| Tracking Node Health                            | 85 |
| Node Health Attributes                          | 85 |
| Node Health Strategy                            | 85 |
| Measuring Node Health                           |    |
| Reloading Services After a Definition Change    |    |

# **Specifying When Recurring Actions are Performed**

By default, recurring actions are scheduled relative to when the resource started. So if your resource was last started at 14:32 and you have a backup set to be performed every 24 hours, then the backup will always run in the middle of the business day — hardly desirable.

To specify a date and time that the operation should be relative to, set the operation's intervalorigin. The cluster uses this point to calculate the correct start-delay such that the operation will occur at origin + (interval \* N).

So, if the 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 with a start delay of 11 hours and 28 minutes. If the resource is moved to another node before 2am, then the operation is cancelled.

The value specified for interval and interval-origin can be any date/time conforming to the ISO8601 standard [http://en.wikipedia.org/wiki/ISO\_8601]. By way of example, to specify an operation that would run on the first Monday of 2009 and every Monday after that, you would add:

#### **Example 9.1. Specifying a Base for Recurring Action Intervals**

<op id="my-weekly-action" name="custom-action" interval="P7D" interval-origin="200</pre>

### 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 nodel -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=N for a resource, it will be banned from the original node after N failures.

#### 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.

<sup>&</sup>lt;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.

### **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 falure timeout.

For example, a setting of migration-threshold=2 and failure-timeout=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 concept: 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 STONITH is enabled, then the cluster will fence the node in order to be able to start the resource elsewhere. If STONITH is not enabled, 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 the failure timeout.

### **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, one 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 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, 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 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 the section called "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. E.g.

```
Email prefers sles-1 with a score of -INFINITY

Email prefers sles-2 with a score of 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 a test whether list of machines (specified by DNS hostname or IPv4/IPv6 address) are reachable and uses the results to maintain a node attribute called pingd by default. <sup>2</sup>

#### Note

Older versions of Pacemaker used a different agent **ocf:pacemaker:pingd** which is now deprecated in favor of **ping**. If your version of Pacemaker does not contain the **ping** resource agent, download the latest version from https://github.com/ClusterLabs/pacemaker/tree/master/extra/resources/ping

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 9.1. Common Options for a ping Resource

| Field      | 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 host names, IPv4 and IPv6 addresses.  |

## Example 9.2. An example ping cluster resource that checks node connectivity once every minute

<sup>&</sup>lt;sup>2</sup> The attribute name is customizable, in order to allow multiple ping groups to be defined.

```
</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 Chapter 8, 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 (e.g. the service network's default gateway), to prevent the cluster from running a resource on any unconnected node.

#### Example 9.3. Don't run a resource on unconnected nodes

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.

#### Example 9.4. Run only on nodes connected to three or more ping targets.

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

#### Example 9.5. Prefer the node with the 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:

#### Example 9.6. 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).

#### Example 9.7. A more complex example of choosing a location based on connectivity

### **Migrating Resources**

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

However, some types of resources, such as Xen virtual guests, 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, and 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 is it 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 must use an OCF style agent.
- 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 metadata.
- 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.

### **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 9.2. Allowed Values for Node Health Attributes

| Value   | Intended significance  |  |
|---------|--|--|
| red     | This indicator is unhealthy  |  |
| yellow  | This indicator is becoming unhealthy   |  |
| green   | This indicator is healthy  |  |
| integer | 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 9.3. Node Health Strategies** 

| Value          | Effect   |  |
|----------------|--|--|
| none           | Do not track node health attributes at all.  |  |
| migrate-on-red | Assign the value of -INFINITY to red, and 0 to yellow and green. This will cause all resources to move off the node if any attribute is red  |  |
| only-green     | Assign the value of -INFINITY to red and yellow, and 0 to green This will cause all resources to move off the node if any attribute is red or yellow.  |  |
| progressive    | Assign the value of the node-health-red cluster option to red, the value of node-health-yellow to yellow, and the value of node-health-green to green. Each node is additionally assigne a score of node-health-base (this allows resources to start even some attributes are yellow). This strategy gives the administrator fin control over how important each value is. |  |
| custom         | Track node health attributes using the same values as progressive for red, yellow, and green, but do not take them into account. The administrator is expected to implement a policy by defining rules (see Chapter 8, Rules) referencing node health attributes.  |  |

### **Measuring Node Health**

</primitive>

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 ways are resource agents or separate daemons.

Pacemaker provides examples that can be used directly or as a basis for custom code. The ocf:pacemaker:HealthCPU and ocf:pacemaker:HealthSMART resource agents set node health attributes based on CPU and disk parameters. The ipmiservicelogd daemon sets node health attributes based on IPMI values (the ocf:pacemaker:SystemHealth resource agent can be used to manage the daemon as a cluster resource).

In order to take advantage of this feature - firstly add the resource to your cluster, preferably as a cloned resource to constantly measure health on all nodes:

```
</clone>
```

This way attrd\_updater will set proper status for each node running this resource. Any attribute matching "#health-[a-zA-z]+" will force cluster to migrate all resources from unhealthy node and place it on other nodes according to all constraints defined in your cluster.

When the node is no longer faulty you can force the cluster to restart the cloned resource on faulty node and make it available to take resources, in this case since we are using HealthIOWait provider:

```
# attrd_updater -n "#health-iowait" -U "green" --node="<nodename>" -d "60s"
```

## Reloading Services After a Definition Change

The cluster automatically detects changes to the definition of services it manages. The normal response is to stop the service (using the old definition) and start it again (with the new definition). This works well, but some services 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:

1. Accept the reload operation and perform any required actions. *The actions here depend completely on your application!* 

### Example 9.8. The DRBD agent's logic for supporting reload

```
case $1 in
    start)
        drbd_start
        ;;
    stop)
        drbd_stop
        ;;
    reload)
        drbd_reload
    monitor)
        drbd monitor
        ;;
        drbd_usage
        exit $OCF_ERR_UNIMPLEMENTED
esac
exit $?
```

2. Advertise the reload operation in the actions section of its metadata

#### Example 9.9. The DRBD Agent Advertising Support for the reload Operation

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

```
Master/Slave OCF Resource Agent for DRBD
  </longdesc>
  <actions>
   <action name="start"
                          timeout="240" />
   <action name="reload" timeout="240" />
    <action name="promote" timeout="90" />
    <action name="demote" timeout="90" />
    <action name="notify" timeout="90" />
   <action name="stop"
                         timeout="100" />
    <action name="meta-data"
                               timeout="5" />
    <action name="validate-all" timeout="30" />
  </actions>
</resource-agent>
```

3. Advertise one or more parameters that can take effect using reload.

Any parameter with the unique set to 0 is eligible to be used in this way.

### Example 9.10. Parameter that can be changed using reload

Once these requirements are satisfied, the cluster will automatically know to reload the resource (instead of restarting) when a non-unique field changes.

#### Note

Metadata will not be re-read unless the resource needs to be started. This may mean that the resource will be restarted the first time, even though you changed a parameter with unique=0.

#### Note

If both a unique and non-unique field are changed simultaneously, the resource will still be restarted.

## **Chapter 10. Advanced Resource Types**

### **Table of Contents**

| Groups - A Syntactic Shortcut                              | 89  |
|--|-----|
| Group Properties   | 9(  |
| Group Options  | 9(  |
| Group Instance Attributes                                  | 9(  |
| Group Contents   | 91  |
| Group Constraints  | 91  |
| Group Stickiness   | 91  |
| Clones - Resources That Can Have Multiple Active Instances | 91  |
| Anonymous versus Unique Clones                             |     |
| Promotable clones  | 91  |
| Clone Properties   | 92  |
| Clone Options  | 92  |
| Clone Contents   | 93  |
| Clone Instance Attributes                                  | 93  |
| Clone Constraints  | 93  |
| Clone Stickiness   | 96  |
| Clone Resource Agent Requirements                          | 96  |
| Monitoring Promotable Clone Resources                      | .02 |
| Determining Which Instance is Promoted                     | 03  |
| Bundles - Isolated Environments                            | .03 |
| Bundle Prerequisites                                       | .04 |
| Bundle Properties  | .04 |
| Bundle Container Properties                                | .04 |
| Bundle Network Properties                                  | 05  |
| Bundle Storage Properties                                  | .06 |
| Bundle Primitive   | 07  |
| Bundle Node Attributes                                     | .08 |
| Bundle Meta-Attributes                                     | .09 |
| Limitations of Bundles 1                                   | 09  |

## **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.

#### **Example 10.1. A group of two primitive resources**

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:

#### Example 10.2. How the cluster sees a group resource

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 10.1. 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 the section called "Resource Options" for information about those properties.

### **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 the section called "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.

#### **Example 10.3. Some constraints involving groups**

### **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.

# **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>

### **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**

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

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 *master* and *slave*, <sup>2</sup> but is agnostic to what the service calls them or what they do.

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

### **Clone Properties**

Table 10.2. 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 10.3. Clone-specific configuration options

| Field               | Default                    | Description  |
|---------------------|----------------------------|--|
| globally-<br>unique | false                      | If true, each clone instance performs a distinct function  |
| clone-max           | number of nodes in cluster | The maximum number of clone instances that can be started across the entire cluster  |
| clone-node-<br>max  | 1                          | If globally-unique is true, 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.   |
| notify              | false                      | Call the resource agent's notify action for all active instances, before and after starting or stopping any clone instance. The resource agent must support this action. Allowed values: false, true   |
| ordered             | false                      | If true, clone instances must be started sequentially instead of in parallel Allowed values: false, true   |
| interleave          | false                      | When this clone is ordered relative to another clone, if this option is false (the default), the ordering is relative to <i>all</i> instances of the other clone, whereas if this option is true, the ordering is relative only to instances on the same node. Allowed values: false, true |

<sup>&</sup>lt;sup>2</sup> These are historical terms that will eventually be replaced, but the extensive use of them and the need for backward compatibility makes it a long process. You may see examples using a master tag instead of a clone tag with the promotable meta-attribute set to true; the master tag is supported, but deprecated, and will be removed in a future version. You may also see such services referred to as *multi-state* or *stateful*; these means the same thing as *promotable*.

| Field                 | Default | Description  |
|-----------------------|---------|--|
| promotable            | false   | If true, clone instances can perform a special role that Pacemaker will manage via the resource agent's promote and demote actions. The resource agent must support these actions. Allowed values: false, true |
| promoted-max          | 1       | If promotable is true, the number of instances that can be promoted at one time across the entire cluster  |
| promoted-<br>node-max | 1       | If promotable is true and globally-unique is false, the number of clone instances can be promoted at one time on a single node   |

For backward compatibility, master-max and master-node-max are accepted as aliases for promoted-max and promoted-node-max, but are deprecated since 2.0.0, and support for them will be removed in a future version.

### **Clone Contents**

Clones must contain exactly one primitive or group resource.

#### Example 10.4. A clone that runs a web server on all nodes

### 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 Attributes**

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.

#### **Example 10.5. Some constraints involving clones**

```
<constraints>
    <rsc_location id="clone-prefers-nodel" rsc="apache-clone" node="nodel" score="
    <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")</pre>
```

```
</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 master role, and colocation constraints may contain rsc-role and/or with-rsc-role fields.

Table 10.4. Additional colocation constraint options for promotable clone resources

| Field             | Default | Description   |
|-------------------|---------|---|
| rsc-role          | Started | An additional attribute of colocation constraints that specifies the role that rsc must be in. Allowed values: Started, Master, Slave.      |
| with-rsc-<br>role | Started | An additional attribute of colocation constraints that specifies the role that with-rsc must be in. Allowed values: Started, Master, Slave. |

#### **Example 10.6. Constraints involving promotable clone resources**

In the example above, myApp will wait until one of the database copies has been started and promoted to master 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 (master or slave). In the example above, the cluster will choose a location based on where database is running as a master, and if there are multiple master 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**

Table 10.5. Additional colocation set options relevant to promotable clone resources

| Field | Default | Description   |  |
|-------|---------|---|--|
| role  | Started | The role that <i>all members</i> of the set must be in. Allowed values: Started, Master, Slave. |  |

In the following example B's master must be located on the same node as A's master. Additionally resources C and D must be located on the same node as A's and B's masters.

#### Example 10.7. Colocate C and D with A's and B's master instances

## **Using Promotable Clone Resources in Ordered Sets**

Table 10.6. Additional ordered set options relevant to promotable clone resources

| Field  | Default | Description  |
|--------|---------|--|
| action |         | An additional attribute of ordering constraint sets that specifies the action that applies to <i>all members</i> of the set. Allowed values: start, stop, promote, demote. |

#### Example 10.8. Start C and D after first promoting A and B

In the above example, B cannot be promoted to a master role 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 mode called slave. From there the cluster will decide which instances to promote to master.

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 10.7. Role implications of OCF return codes

| Monitor Return Code | Description      |
|---------------------|------------------|
| OCF_NOT_RUNNING     | Stopped          |
| OCF_SUCCESS         | Running (Slave)  |
| OCF_RUNNING_MASTER  | Running (Master) |
| OCF_FAILED_MASTER   | Failed (Master)  |
| Other               | Failed (Slave)   |

#### **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 10.8. Environment variables supplied with Clone notify actions

| Variable                                     | Description                              |
|--|--|
| OCF_RESKEY_CRM_meta_notify_type              | Allowed values: pre, post                |
| OCF_RESKEY_CRM_meta_notify_operation         | Allowed values: start, stop              |
| 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:

#### **Example 10.9. 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
- $\bullet \ \ Resources \ that \ were \ stopped: \verb§OCF_RESKEY_CRM_meta_notify_stop_resource \\$

#### **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 10.9. Extra environment variables supplied for promotable clones

| Variable                                    | Description   |
|---|---|
| OCF_RESKEY_CRM_meta_notify_master_resource  | Resources that are running in Master mode           |
| OCF_RESKEY_CRM_meta_notify_slave_resource   | Resources that are running in Slave mode            |
| 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_master_uname     | Nodes on which resources are running in Master mode |
| OCF_RESKEY_CRM_meta_notify_slave_uname      | Nodes on which resources are running in Slave mode  |

#### **Interpretation of Promotable Notification Variables**

#### **Pre-notification (demote):**

- Active resources: \$OCF\_RESKEY\_CRM\_meta\_notify\_active\_resource
- Master resources: \$OCF\_RESKEY\_CRM\_meta\_notify\_master\_resource
- Slave resources: \$OCF\_RESKEY\_CRM\_meta\_notify\_slave\_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
- $\bullet \ \ Resources \ to \ be \ stopped: \verb§OCF_RESKEY_CRM_meta_notify_stop_resource \\$

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

- Active resources: \$OCF\_RESKEY\_CRM\_meta\_notify\_active\_resource
- Master resources:
  - \$OCF\_RESKEY\_CRM\_meta\_notify\_master\_resource
  - minus \$OCF\_RESKEY\_CRM\_meta\_notify\_demote\_resource
- Slave resources: \$OCF\_RESKEY\_CRM\_meta\_notify\_slave\_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
- Master resources:
  - \$OCF\_RESKEY\_CRM\_meta\_notify\_master\_resource
  - minus \$OCF\_RESKEY\_CRM\_meta\_notify\_demote\_resource
- Slave resources:
  - \$OCF\_RESKEY\_CRM\_meta\_notify\_slave\_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
- $\bullet \ \ Resources \ that \ were \ demoted: \verb§OCF_RESKEY_CRM_meta_notify_demote_resource$
- 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

- Master resources:
  - \$OCF\_RESKEY\_CRM\_meta\_notify\_master\_resource
  - minus \$OCF\_RESKEY\_CRM\_meta\_notify\_demote\_resource
- Slave resources:
  - \$OCF\_RESKEY\_CRM\_meta\_notify\_slave\_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
- Master resources:
  - \$OCF\_RESKEY\_CRM\_meta\_notify\_master\_resource
  - minus \$OCF\_RESKEY\_CRM\_meta\_notify\_demote\_resource
  - plus \$OCF\_RESKEY\_CRM\_meta\_notify\_promote\_resource
- Slave resources:
  - \$OCF\_RESKEY\_CRM\_meta\_notify\_slave\_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 slave role, and an additional one with role="master" will cover the master role.

#### Example 10.10. Monitoring both states of a promotable clone resource

#### **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\_master utility, 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 as masters.

#### Example 10.11. Explicitly preferring node1 to be promoted to master

## **Bundles - Isolated Environments**

Pacemaker supports a special syntax for launching a container [https://en.wikipedia.org/wiki/Operating-system-level\_virtualization] with any infrastructure it requires: the *bundle*.

Pacemaker bundles support Docker [https://www.docker.com/], podman [https://podman.io/], and rkt [https://coreos.com/rkt/] container technologies. <sup>3</sup>

#### Example 10.12. 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"</pre>
            host-netmask="24"
            host-interface="eth0">
      <port-mapping id="httpd-port" port="80"/>
   </network>
   <storage>
      <storage-mapping id="httpd-syslog"</pre>
                        source-dir="/dev/log"
                        target-dir="/dev/log"
                        options="rw"/>
      <storage-mapping id="httpd-root"</pre>
                        source-dir="/srv/html"
                        target-dir="/var/www/html"
                        options="rw,Z"/>
      <storage-mapping id="httpd-logs"</pre>
```

<sup>&</sup>lt;sup>3</sup> Docker is a trademark of Docker, Inc. No endorsement by or association with Docker, Inc. is implied.

## **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,orocf: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 10.10. XML Attributes of a bundle Element

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

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

## **Bundle Container Properties**

Table 10.11. XML attributes of a docker, podman, or rkt Element

| Attribute             | Default   | Description   |
|-----------------------|---|---|
| image                 |   | Container image tag (required)  |
| replicas              | Value of promoted-max if that is positive, else 1 | A positive integer specifying the number of container instances to launch   |
| replicas-per-<br>host | 1   | A positive integer specifying the number of container instances allowed to run on a single node   |
| promoted-max          | 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 master role |

| Attribute   | Default   | Description  |
|-------------|---|--|
| network     |   | If specified, this will be passed to the docker run, podman run, or rkt run command as the network setting for the container.  |
| run-command | /usr/sbin/pacemaker-<br>remoted if bundle contains a<br>primitive, otherwise none | This command will be run inside the container when launching it ("PID 1"). If the bundle contains a primitive, this command <i>must</i> start pacemakerremoted (but could, for example, be a script that does other stuff, too). |
| options     |   | Extra command-line options to pass to the docker run, podman run, or rkt run 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 10.12. XML attributes of a network Element

| Attribute      | Default | Description   |
|----------------|---------|---|
| add-host       | TRUE    | If TRUE, and ip-range-start is used, Pacemaker will automatically ensure that /etc/hosts inside the containers has entries for each replica name and its assigned IP.   |
| ip-range-start |         | If specified, Pacemaker will create an implicit ocf:heartbeat:IPaddr2 resource for each container instance, starting with this IP address, using up to replicas 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 ip-range-start is specified, the IP addresses are created with this CIDR netmask (as a number of bits).  |
| host-interface |         | If ip-range-start is specified, the IP addresses are created on this host interface (by default, it will be determined from the IP address).  |

| Attribute    | Default | Description  |
|--------------|---------|--|
| control-port | 3121    | If the bundle contains a primitive, 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 ip-range-start (in which case replicas-per-host must be 1), or when the bundle may run on a Pacemaker Remote node that is already listening on the default port. Any PCMK_remote_port 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 10.13. 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 ip-range-start is specified) will be forwarded to the container network. Exactly one of port or range must be specified in a port-mapping.  |
| internal-port | value of port | If port and this are specified, connections to port 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 <code>first_port-last_port</code> ) on the host network (on the container's assigned IP address, if <code>ip-range-start</code> is specified) will be forwarded to the same ports in the container network. Exactly one of port or range must be specified in a <code>port-mapping</code> . |

#### 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 10.14. 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 sourcedir and source-dir-root must be specified in a storage-mapping.  |
| 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 replica name. Exactly one of source-dir and source-dir-root must be specified in a storage-mapping. |
| 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 masters 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. 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 controlport 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 the section called "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.

# Chapter 11. Reusing Parts of the Configuration

#### **Table of Contents**

| Reusing Resource Definitions                  | 110 |
|---|-----|
| Configuring Resources with Templates          | 110 |
| Using Templates in Constraints                | 112 |
| Using Templates in Resource Sets              | 112 |
| Reusing Rules, Options and Sets of Operations | 113 |
| Tagging Configuration Elements                | 114 |
| Configuring Tags                              | 115 |
| Using Tags in Constraints and Resource Sets   | 115 |

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.

## **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.

#### **Example 11.1. Resource template for a migratable Xen virtual machine**

```
<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.

#### Example 11.2. Xen primitive resource using a resource template

In the example above, the new primitive vml will inherit everything from vm-template. For example, the equivalent of the above two examples would be:

#### Example 11.3. 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" va</pre>
  </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="vml-instance_attributes">
    <nvpair id="vml-instance_attributes-name" name="name" value="vml"/>
    <nvpair id="vml-instance_attributes-xmfile" name="xmfile" value="/etc/xen/shar</pre>
  </instance attributes>
</primitive>
```

If you want to overwrite some attributes or operations, add them to the particular primitive's definition.

#### Example 11.4. Xen resource overriding template values

<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 the section called "Specifying the Order in which Resources Should Start/Stop")
- colocation constraints (see the section called "Placing Resources Relative to other Resources")
- rsc\_ticket constraints (for multi-site clusters as described in the section called "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 section:

<rsc\_colocation id="vm-template-colo-base-rsc" rsc="vm-template" rsc-role="Started"</pre>

#### 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_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>
```

## 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

Then instead of duplicating the rule for all your other resources, you can instead specify:

#### **Example 11.5. Referencing rules from other constraints**

#### **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:

#### Example 11.6. 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="myOtherlRsc" 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
```

## **Tagging Configuration Elements**

operations.

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.

#### **Example 11.7. Tag referencing three resources**

```
<tags>
  <tag id="all-vms">
      <obj_ref id="vm1"/>
      <obj_ref id="vm2"/>
      <obj_ref id="vm3"/>
      </tag>
</tag>
```

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.

#### Example 11.8. 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:

#### **Example 11.9. 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="vm2" 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.

# Chapter 12. Utilization and Placement Strategy

#### **Table of Contents**

| Utilization attributes   | . 116 |
|--|-------|
| Placement Strategy   | 117   |
| Allocation Details   | . 118 |
| Which node is preferred to get consumed first when allocating resources? | . 118 |
| Which node has more free capacity?                                       | . 118 |
| Which resource is preferred to be assigned first?                        | . 118 |
| Limitations and Workarounds  |       |

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.

## **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.

#### Example 12.1. Specifying CPU and RAM capacities of two nodes

#### Example 12.2. 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.

## **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 <i>only</i> 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 <i>and</i> 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 <i>only</i> 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.

## **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 the section called "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.

## **Limitations and Workarounds**

The type of problem Pacemaker is dealing with here is known as the knapsack problem [http://en.wikipedia.org/wiki/Knapsack\_problem] and falls into the NP-complete [http://en.wikipedia.org/wiki/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 section:

- 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.

# **Chapter 13. Access Control Lists** (ACLs)

#### **Table of Contents**

| ACL Prerequisites      | 120 |
|------------------------|-----|
| ACL Configuration      | 120 |
| ACL Roles              | 120 |
| ACL Targets and Groups | 121 |
| ACL Examples           | 122 |

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.

## **ACL Prerequisites**

In order to use ACLs:

- The Pacemaker software must have been compiled with ACL support. If the output of the command pacemakerd --features contains acls, your installation supports ACLs.
- Desired users must have user accounts in the haclient group on all nodes in the cluster.
- If your CIB was created before Pacemaker 1.1.12, it may need to be updated to the current schema using cibadmin --upgrade in order to use the syntax documented here.
- The enable-acl cluster option must be set to true.

## **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.

## **ACL Roles**

An ACL role is a collection of permissions allowing or denying access to particular portions of the CIB.

Table 13.1. Properties of an ACL Role

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

An acl\_role element may contain any number of acl\_permission elements.

Table 13.2. Properties of an ACL Permission

| Attribute   | Description   |
|-------------|---|
| id          | A unique name for the permission (required)   |
| description | Arbitrary text (not used by Pacemaker)  |
| kind        | The access being granted. Allowed values are read, write, and deny. A value of write 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 object-type, xpath, and reference must be specified for a permission.)   |
| attribute   | If specified, the permission applies only to object-type elements that have this attribute set (to any value). If not specified, the permission applies to all object-type elements. May only be used with object-type.   |
| reference   | The ID of an XML element in the CIB to which the permission applies. (Exactly one of object-type, xpath, and reference must be specified for a permission.)   |
| xpath       | An XPath [https://www.w3.org/TR/xpath-10/] 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 object-type, xpath, and reference 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.

## **ACL Targets and Groups**

ACL targets correspond to user accounts on the system.

Table 13.3. Properties of an ACL Target

| Attribute | Description                                 |
|-----------|---|
| id        | The name of a user on the system (required) |

ACL groups may be specified, but are not currently used by Pacemaker. This is expected to change in a future version.

Table 13.4. Properties of an ACL Group

| Attribute | Description                                  |
|-----------|--|
| id        | The name of a group on the system (required) |

Each acl\_target and acl\_group element may contain any number of role elements.

Table 13.5. Properties of an ACL Role Reference

| Attribute | Description   |
|-----------|---|
|           | The id of an acl_role 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 other user accounts, when enable-acl is true, permission to all parts of the CIB is denied by default (permissions must be explicitly granted).

## **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"</pre>
           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"</pre>
           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 transie
    xpath="//instance_attributes/nvpair[@name='standby']"/>
<acl_permission id="minimal-maintenance" kind="read"
    description="allow reading maintenance node attribute (permanent or tra
    xpath="//nvpair[@name='maintenance']"/>
<acl permission id="minimal-target-role" kind="read"</pre>
    description="allow reading resource target roles"
    xpath="//resources//meta_attributes/nvpair[@name='target-role']"/>
<acl permission id="minimal-is-managed" kind="read"</pre>
    description="allow reading resource managed status"
    xpath="//resources//meta_attributes/nvpair[@name='is-managed']"/>
<acl_permission id="minimal-deny-instance-attributes" kind="deny"</pre>
    xpath="//instance_attributes"/>
<acl_permission id="minimal-deny-meta-attributes" kind="deny"</pre>
    xpath="//meta_attributes"/>
<acl_permission id="minimal-deny-operations" kind="deny"
    xpath="//operations"/>
<acl_permission id="minimal-deny-utilization" kind="deny"</pre>
    xpath="//utilization"/>
<acl_permission id="minimal-nodes" kind="read"
    description="allow reading node names/IDs (attributes are denied separa
    xpath="/cib/configuration/nodes"/>
<acl_permission id="minimal-resources" kind="read"
    description="allow reading resource names/agents (parameters are denied
    xpath="/cib/configuration/resources"/>
<acl permission id="minimal-deny-constraints" kind="deny"</pre>
    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"</pre>
    xpath="/cib/configuration/op_defaults"/>
<acl permission id="minimal-deny-rsc defaults" kind="deny"</pre>
    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.

## Chapter 14. Status — Here be dragons

#### **Table of Contents**

| Node Status                       | 123 |
|-----------------------------------|-----|
| Transient Node Attributes         | 120 |
| Operation History                 | 120 |
| Simple Operation History Example  | 128 |
| Complex Operation History Example | 129 |

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.

## **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.

#### Example 14.1. A bare-bones status entry for a healthy node cl-virt-1

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 14.1. Authoritative Sources for State Information** 

| CIB Object           | Authoritative Source |
|----------------------|----------------------|
| node_state           | pacemaker-controld   |
| transient_attributes | pacemaker-attrd      |
| lrm                  | 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 14.2. Node Status Fields** 

| Field | Description   |
|-------|---|
| id    | Unique identifier for the node. Corosync-based clusters use a numeric |
|       | counter.  |

| Field                | Description  |
|----------------------|--|
| uname                | The node's name as known by the cluster  |
| in_ccm               | Is the node a member at the cluster communication layer? Allowed values: true, false.          |
| crmd                 | Is the node a member at the pacemaker layer? Allowed values: online, offline.                  |
| crm-debug-<br>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: down, pending, member, banned. |
| expected             | Expected value for join.   |

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.

## **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.

#### Example 14.2. A set of transient node attributes for node cl-virt-1

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. 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).

## **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.

<sup>&</sup>lt;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.

#### Example 14.3. A record of the apostonith resource

<lrm\_resource id="apcstonith" type="apcmastersnmp" 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 14.3. Contents of an lrm\_rsc\_op job

| Field                       | Description   |
|-----------------------------|---|
| id                          |   |
|                             | Identifier for the job constructed from the resource's id, operation and interval.  |
| call-id                     |   |
|                             | The job's ticket number. Used as a sort key to determine the order in which the jobs were executed.   |
| operation                   |   |
|                             | The action the resource agent was invoked with.   |
| interval                    |   |
|                             | The frequency, in milliseconds, at which the operation will be repeated. A one-off job is indicated by 0.   |
| op-status                   |   |
|                             | The job's status. Generally this will be either 0 (done) or -1 (pending). Rarely used in favor of rc-code.  |
| rc-code                     |   |
|                             | The job's result. Refer to the <i>Resource Agents</i> section of <i>Pacemaker Administration</i> for details on what the values here mean and how they are interpreted. |
| last-run                    |   |
|                             | Machine-local date/time, in seconds since epoch, at which the job was executed. For diagnostic purposes.  |
| last-rc-change              |   |
|                             | Machine-local date/time, in seconds since epoch, at which the job first returned the current value of rc-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 LRMd. For diagnostic purposes.   |
| queue-time  crm_feature_set | · -   |

| Field            | Description  |
|------------------|--|
|                  | The version which this job description conforms to. Used when processing op-digest.  |
| 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 transition-magic (below).   |
| transition-magic |  |
|                  | A concatenation of the job's op-status, rc-code and transition-key. 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-origin |  |
|                  | The origin of the current values. For diagnostic purposes.   |

## **Simple Operation History Example**

## Example 14.4. A monitor operation (determines current state of the apcstonith resource)

```
<lrm_resource id="apcstonith" type="apcmastersnmp" class="stonith">
    <lrm_rsc_op id="apcstonith_monitor_0" operation="monitor" call-id="2"
        rc-code="7" op-status="0" interval="0"
        crm-debug-origin="do_update_resource" crm_feature_set="3.0.1"
        op-digest="2e3da9274d3550dc6526fb24bfcbcba0"
        transition-key="22:2:7:2668bbeb-06d5-40f9-936d-24cb7f87006a"
        transition-magic="0:7;22:2:7:2668bbeb-06d5-40f9-936d-24cb7f87006a"
        last-run="1239008085" last-rc-change="1239008085" exec-time="10" queue-time="0"
</lrm_resource>
```

In the above example, the job is a non-recurring monitor operation often referred to as a "probe" for the apostonith 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**

#### Example 14.5. Resource history of a pingd clone with multiple jobs

```
<lrm_resource id="pingd:0" type="pingd" class="ocf" provider="pacemaker">
  <lrm_rsc_op id="pingd:0_monitor_30000" operation="monitor" call-id="34"</pre>
    rc-code="0" op-status="0" interval="30000"
    crm-debug-origin="do_update_resource" crm_feature_set="3.0.1"
    transition-key="10:11:0:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-run="1239009741" last-rc-change="1239009741" exec-time="10" queue-time="0
  <lrm_rsc_op id="pingd:0_stop_0" operation="stop"</pre>
    crm-debug-origin="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-run="1239009741" last-rc-change="1239009741" exec-time="10" queue-time="0
  <lrm_rsc_op id="pingd:0_start_0" operation="start" call-id="33"</pre>
    rc-code="0" op-status="0" interval="0"
    crm-debug-origin="do_update_resource" crm_feature_set="3.0.1"
    transition-key="31:11:0:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-run="1239009741" last-rc-change="1239009741" exec-time="10" queue-time="0
  <lrm_rsc_op id="pingd:0_monitor_0" operation="monitor" call-id="3"</pre>
    rc-code="0" op-status="0" interval="0"
    crm-debug-origin="do_update_resource" crm_feature_set="3.0.1"
    transition-key="23:2:7:2668bbeb-06d5-40f9-936d-24cb7f87006a"
    last-run="1239008085" last-rc-change="1239008085" exec-time="20" queue-time="0
  </lrm_resource>
```

When more than one job record exists, it is important to first sort them by call-id before interpreting them.

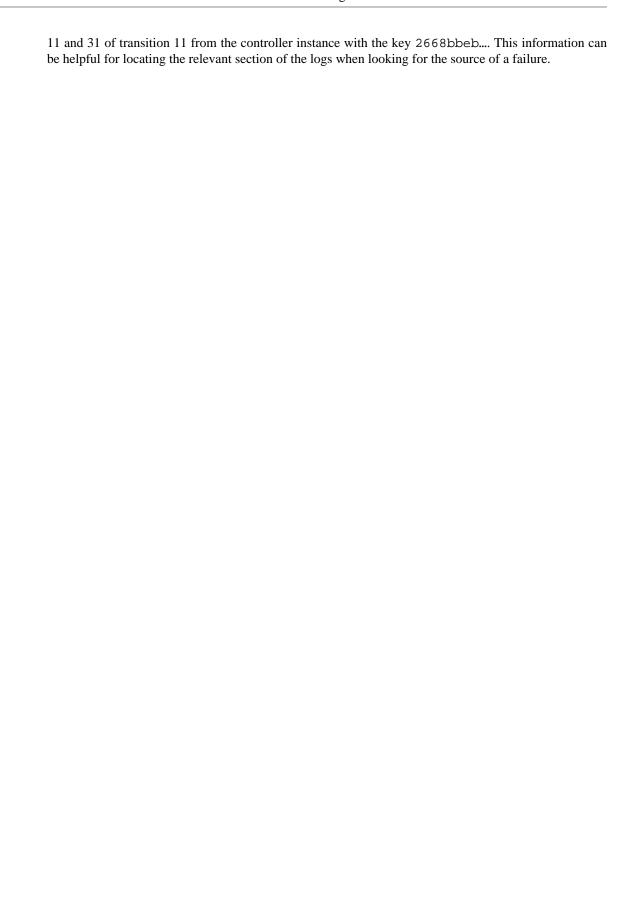
Once sorted, the above example can be summarized as:

- 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



# Chapter 15. Multi-Site Clusters and Tickets

#### **Table of Contents**

| Challenges for Multi-Site Clusters                          | 131 |
|---|-----|
| Conceptual Overview   | 131 |
| Ticket  | 132 |
| Dead Man Dependency   | 132 |
| Cluster Ticket Registry                                     | 132 |
| Configuration Replication                                   | 133 |
| Configuring Ticket Dependencies                             | 133 |
| Managing Multi-Site Clusters                                | 134 |
| Granting and Revoking Tickets Manually                      | 134 |
| Granting and Revoking Tickets via a Cluster Ticket Registry | 134 |
| General Management of Tickets                               | 135 |
| For more information  | 136 |

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).

## **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.

## **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 multisite 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.

## **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 master mode to slave mode.

#### Example 15.1. Constraint that fences node if ticketA is revoked

```
<rsc_ticket id="rsc1-req-ticketA" rsc="rsc1" ticket="ticketA" loss-policy="fence"/</pre>
```

The example above creates a constraint with the ID rscl-req-ticketA. It defines that the resource rscl depends on ticketA and that the node running the resource should be fenced if ticketA is revoked.

If resource rsc1 were a promotable resource (i.e. it could run in master or slave mode), you might want to configure that only master mode depends on ticketA. With the following configuration, rsc1 will be demoted to slave mode if ticketA is revoked:

#### Example 15.2. Constraint that demotes rsc1 if ticketA is revoked

```
<rsc_ticket id="rsc1-req-ticketA" rsc="rsc1" rsc-role="Master" ticket="ticketA" lo</pre>
```

You can create multiple rsc\_ticket constraints to let multiple resources depend on the same ticket. However, rsc\_ticket also supports resource sets (see the section called "Resource Sets"), so one can easily list all the resources in one rsc\_ticket constraint instead.

#### **Example 15.3. Ticket constraint for multiple resources**

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.

## **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 [https://github.com/ClusterLabs/booth] here as an example of software that can be used with pacemaker as a Cluster Ticket Registry. Booth implements the Raft [http://en.wikipedia.org/wiki/Raft\_%28computer\_science%29] 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
```

## For more information

- SUSE's Geo Clustering quick start [https://www.suse.com/documentation/sle-ha-geo-12/art\_ha\_geo\_quick/data/art\_ha\_geo\_quick.html]
- Booth [https://github.com/ClusterLabs/booth]

## **Appendix A. Sample Configurations**

#### **Table of Contents**

| Empty                  | 137 |
|------------------------|-----|
| Simple                 | 137 |
| Advanced Configuration | 138 |

## **Empty**

#### **Example A.1. An Empty Configuration**

## **Simple**

## Example A.2. 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="</pre>
  <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</pre>
    </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.

## **Advanced Configuration**

#### Example A.3. An advanced configuration with groups, clones and STONITH

```
<cib crm_feature_set="3.0.7" validate-with="pacemaker-1.2" admin_epoch="1" epoch="</pre>
  <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>
      <pri><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">
            <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"</pre>
        score="INFINITY"/>
      <rsc colocation id="group-with-ip" rsc="myGroup" with-rsc="myAddr"</pre>
        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>
  <status/>
</cib>
```

## **Appendix B. Revision History**

**Revision History** 

Revision 1-0 19 Oct 2009 AndrewBeekhof<andrew@beekhof.net>

Import from Pages.app

AndrewBeekhof<andrew@beekhof.net> Revision 2-0 26 Oct 2009

Cleanup and reformatting of docbook xml complete

Revision 3-0 Tue Nov 12 2009 AndrewBeekhof<andrew@beekhof.net>

Split book into chapters and pass validation

Re-organize book for use with Publican [https://fedorahosted.org/publican/]

Revision 3-1 Tue Nov 12 2009 TanjaRothSUSE<taroth@suse.com>,

LarsMarowsky-

BreeSUSE<1mb@suse.com>, YanGaoSUSE<ygao@suse.com>,

ThomasSchraitleSUSE<toms@suse.com>,

DejanMuhamedagicSUSE<dmuhamedagic@suse

Utilization chapter

Resource Templates chapter Multi-Site Clusters chapter

Revision 3-2 Fri Nov 4 2011 PhilippMarekLINBit<philipp.marek@linbit

Extensive style, formatting, and indexing updates

Revision 4-0 Mon Oct 8 2012 AndrewBeekhof<andrew@beekhof.net>

Converted to asciidoc [http://www.methods.co.nz/asciidoc] (which is converted to docbook for use with

Publican [https://fedorahosted.org/publican/])

Revision 5-0 Mon Feb 23 2015 KenGaillot<kgaillot@redhat.com>

Update for clarity, stylistic consistency and current command-line syntax

Revision 6-0 Tue Dec 8 2015 KenGaillot<kgaillot@redhat.com>

Update for Pacemaker 1.1.14

Revision 7-0 Tue May 3 2016 KenGaillot<kgaillot@redhat.com>

Update for Pacemaker 1.1.15

Revision 7-1 Fri Oct 28 2016 KenGaillot<kgaillot@redhat.com>

Overhaul upgrade documentation, and document node health strategies

Tue Oct 25 2016 Revision 8-0 KenGaillot<kgaillot@redhat.com>

Update for Pacemaker 1.1.16

Tue Jul 11 2017 Revision 9-0 KenGaillot<kgaillot@redhat.com>

Update for Pacemaker 1.1.17

Revision 10-0 Fri Oct 6 2017 KenGaillot<kgaillot@redhat.com>

Update for Pacemaker 1.1.18

Revision 11-0 Fri Jan 12 2018 KenGaillot<kgaillot@redhat.com>

Update for Pacemaker 2.0.0

Revision 12-0 Mon Jan 28 2019 KenGaillot<kgaillot@redhat.com>,

> ReidWahl<nwahl@redhat.com>, JanPokorný < jpokorny@redhat.com>

Update for Pacemaker 2.0.1, remove "Further Reading" and "FAQ" sections, and add minor clarifications

and reformatting

Revision 12-1 Mon May 13 2019 KenGaillot<kgaillot@redhat.com>,

MaciejSobkowiak<msobkowiak@egnyte.com>

Document podman support, cluster-name cluster option, and new HealthIOWait agent, with other minor

clarifications and corrections

Revision 12-2 Wed Jun 19 2019 KenGaillot<kgaillot@redhat.com>

Add chapter for ACLs

Revision 13-0 Tue Oct 15 2019 KenGaillot<kgaillot@redhat.com>

#### **Revision History**

Overhaul fencing, rules, and constraints chapters, elaborate on various options, and update for Pacemaker 2.0.3

Revision 13-1 Wed Apr 1 2020

YanGaoSUSE<ygao@suse.com>

Document new feature priority-fencing-delay

| Index   | Status  |
|---|---|
| IIIGCA  | call-id, 127  |
|   | crm-debug-origin, 128<br>crm_feature_set, 127         |
| Symbols   | exec-time, 127  |
| #digests-, 14                                     | id, 127   |
| Node attribute, 14                                | interval, 127   |
| #node-unfenced, 14                                | last-rc-change, 127                                   |
| Node attribute, 14                                | last-run, 127   |
|   | op-digest, 128  |
| A   | op-status, 127  |
| access control list, 120                          | operation, 127  |
| acl_group   | queue-time, 127                                       |
| id, 122   | rc-code, 127  |
| acl_permission                                    | transition-key, 128                                   |
| attribute, 121                                    | transition-magic, 128                                 |
| description, 121                                  | action, 95  |
| id, 121   | Ordering Constraints, 95                              |
| kind, 121   | action attribute, 37                                  |
| object-type, 121                                  | resource_set element, 37                              |
| reference, 121                                    | Action Property, 24, 24, 24, 24, 25, 25, 25, 25       |
| xpath, 121  | Action Status, 127, 127, 127, 127, 127, 127, 127, 127 |
| acl_role  | 127, 127, 127, 128, 128, 128, 128                     |
| description, 120                                  | active_resource, 97                                   |
| id, 120   | Notification Environment Variable, 97                 |
| acl_target  | active_uname, 97                                      |
| id, 121   | Notification Environment Variable, 97                 |
| role  | add-host attribute, 105                               |
| id, 122   | network element, 105                                  |
| ACL, 120  | admin_epoch, 6  |
| acl_group, 122                                    | Cluster Option, 6                                     |
| id, 122   | Alert   |
| acl_permission, 121, 121, 121, 121, 121, 121, 121 | Option  |
| attribute, 121                                    | timeout, 60   |
| description, 121                                  | timestamp-format, 60                                  |
| id, 121   | Alerts, 59  |
| kind, 121   | Asymmetrical Clusters, 31                             |
| object-type, 121                                  | attribute, 12, 121                                    |
| reference, 121                                    | #digests-, 14   |
| xpath, 121  | #node-unfenced, 14                                    |
| acl_role, 120, 120                                | acl_permission, 121                                   |
| description, 120                                  | fail-count-, 13                                       |
| id, 120   | last-failure-, 14<br>maintenance, 14                  |
| acl_target, 121<br>id, 121                        | probe_complete, 14                                    |
| Action, 23  | resource-discovery-enabled, 14                        |
| Property  | shutdown, 14  |
| enabled, 25, 25                                   | site-name, 14   |
| id, 24  | standby, 14   |
| interval, 24                                      | terminate, 14   |
| name, 24  | attribute attribute, 66                               |
| on-fail, 25                                       | expression element, 66                                |
| role, 25  | attribute_name, 64                                    |
| timeout, 24                                       | attribute_value, 64                                   |
|   |   |

| В  | batch-limit, 7                 |
|--|--------------------------------|
| batch-limit, 7                                   | cluster-delay, 9               |
| Cluster Option, 7                                | cluster-ipc-limit, 9           |
| boolean-op attribute, 66                         | cluster-recheck-interval, 10   |
| rule element, 66                                 | concurrent-fencing, 9          |
| Bundle, 103, 103, 103, 103                       | Configuration Version, 6       |
| Container, 104                                   | dc-deadtime, 9                 |
| Meta-attributes, 109                             | election-timeout, 11           |
| Networking, 105                                  | enable-acl, 10                 |
| Node Attributes, 108                             | enable-startup-probes, 8       |
| Prerequisites, 104                               | epoch, 6                       |
| Primitive, 107                                   | fence-reaction, 9              |
| Storage, 106                                     | join-finalization-timeout, 11  |
| bundle element, 104, 104, 104                    | join-integration-timeout, 11   |
| description attribute, 104                       | maintenance-mode, 8            |
| id attribute, 104                                | migration-limit, 7             |
|  | no-quorum-policy, 7            |
| C  | node-health-base, 10           |
| _  | node-health-green, 10          |
| call-id, 127                                     | node-health-red, 10            |
| Action Status, 127                               | node-health-strategy, 10       |
| cib-last-written, 6                              | node-health-yellow, 10         |
| Cluster Property, 6                              | num_updates, 6                 |
| class, 15, 18                                    | pe-error-series-max, 10        |
| Resource, 18                                     | pe-input-series-max, 10        |
| class attribute, 72                              | pe-warn-series-max, 10         |
| rsc_expression element, 72                       | placement-strategy, 10         |
| Clone  | priority-fencing-delay, 9      |
| Option   | remove-after-stop, 11          |
| clone-max, 92                                    | shutdown-escalation, 11        |
| clone-min, 92                                    | shutdown-lock, 11              |
| clone-node-max, 92                               | shutdown-lock-limit, 11        |
| globally-unique, 92                              | start-failure-is-fatal, 8      |
| interleave, 92                                   | startup-fencing, 11            |
| notify, 92                                       | stonith-action, 8              |
| ordered, 92                                      | stonith-enabled, 8             |
| promotable, 93                                   | stonith-max-attempts, 8        |
| promoted-max, 93                                 | stonith-timeout, 8             |
| promoted-node-max, 93                            | stonith-watchdog-timeout, 8    |
| Property   | stop-all-resources, 7          |
| id, 92   | stop-orphan-actions, 8         |
| Clone Option, 92, 92, 92, 92, 92, 92, 93, 93, 93 | stop-orphan-resources, 8       |
| Clone Property, 92                               | symmetric-cluster, 7           |
| Clone Resources, 91                              | transition-delay, 11           |
| clone-max, 92                                    | validate-with, 6               |
| Clone Option, 92                                 | Property                       |
| clone-min, 92                                    | cib-last-written, 6            |
| Clone Option, 92                                 | cluster-infrastructure, 7      |
| clone-node-max, 92                               | cluster-name, 7                |
| Clone Option, 92                                 | dc-uuid, 6                     |
| Clones, 91, 96                                   | dc-version, 7                  |
| Cluster, 6                                       | have-quorum, 6                 |
| Option   | Setting Options with Rules, 77 |
| admin epoch, 6                                   | <u> </u>                       |

| Cluster Option, 6, 6, 6, 6, 7, 7, 7, 7, 7, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, | attribute_value, 64                                       |
|---|---|
| 8, 8, 8, 9, 9, 9, 9, 9, 10, 10, 10, 10, 10, 10, 10, 10, 10,                       | desc, 63  |
| 10, 10, 11, 11, 11, 11, 11, 11, 11, 11  | exec_time, 64   |
| Cluster Property, 6, 6, 6, 7, 7, 7  | interval, 64  |
| cluster-delay, 9  | kind, 63  |
| Cluster Option, 9   | node, 63  |
| cluster-infrastructure, 7   | nodeid, 63  |
| Cluster Property, 7   | rc, 63  |
| cluster-ipc-limit, 9  | recipient, 63   |
| Cluster Option, 9   | rsc, 63   |
| cluster-name, 7   | status, 64  |
| Cluster Property, 7   | target_rc, 64   |
| cluster-recheck-interval, 10  | task, 63  |
| Cluster Option, 10  | timestamp, 63   |
| Colocation Constraint, 34   | timestamp_epoch, 63                                       |
| rsc_colocation element, 34  | timestamp_usec, 63  |
| concurrent-fencing, 9   | version, 63   |
| Cluster Option, 9   | CRM_alert_node_   |
| Configuration, 42, 42   | sequence, 63  |
| Configuration, 42, 42 Configuration Version, 6                                    | crm_feature_set, 127                                      |
| Cluster, 6  | Action Status, 127  |
| Constraint, 29  | custom, 86  |
| •   | custom, 80  |
| Colocation Constraint, 34   | <b>D</b>  |
| rsc_colocation element, 34  | D   |
| Location Constraint, 30   | dampen, 82  |
| Resource Discovery, 31  | Ping Resource Option, 82                                  |
| rsc_location element, 30  | Date Specification, 69                                    |
| Ordering Constraint, 32   | Date/Time Expression, 68                                  |
| rsc_order element, 33   | Date Specification, 69                                    |
| Resource Set, 36  | Duration, 70  |
| Rule, 65  | date_expression element, 68, 68, 68, 68, 69               |
| Constraints   | end attribute, 68   |
| Ordering  | id attribute, 68  |
| action, 95  | operation attribute, 69                                   |
| role, 95  | start attribute, 68                                       |
| rsc-role, 94  | date_spec element, 69, 69, 69, 69, 69, 69, 69, 69, 69, 69 |
| with-rsc-role, 94   | 69  |
| Container, 104  | hours attribute, 69                                       |
| Docker  | id attribute, 69  |
| Bundle, 103   | monthdays attribute, 69                                   |
| podman  | months attribute, 69                                      |
| Bundle, 103   | moon attribute, 69  |
| rkt   | weekdays attribute, 69                                    |
| Bundle, 103   | weeks attribute, 69                                       |
| control-port attribute, 106   | weekyears attribute, 69                                   |
| network element, 106  | yeardays attribute, 69                                    |
| Controlling Cluster Options, 77   | years attribute, 69                                       |
| crm-debug-origin, 126, 128  | dc-deadtime, 9  |
| Action Status, 128  | Cluster Option, 9   |
| Node Status, 126  | dc-uuid, 6  |
| crmd, 126   | Cluster Property, 6                                       |
| Node Status, 126  | dc-version, 7   |
| CRM_alert_  | Cluster Property, 7                                       |
| attribute name, 64  | Clusici Hoporty, /  |

| demote_resource, 99                                    | node, 63                                       |
|--|--|
| Notification Environment Variable, 99                  | nodeid, 63                                     |
| demote_uname, 99                                       | rc, 63   |
| Notification Environment Variable, 99                  | recipient, 63                                  |
| desc, 63   | rsc, 63  |
| description, 120, 121                                  | status, 64                                     |
| acl_permission, 121                                    | target_rc, 64                                  |
| acl_role, 120  | task, 63                                       |
| description attribute, 104                             | timestamp, 63                                  |
| bundle element, 104                                    | timestamp_epoch, 63                            |
| Determine by Rules, 73                                 | timestamp_usec, 63                             |
| Determine Resource Location, 73                        | version, 63                                    |
| devices, 53  | CRM_alert_node_                                |
| fencing-level, 53                                      | sequence, 63                                   |
| Docker   | OCF_RESKEY_CRM_meta_notify_                    |
| Bundle, 103  | active_resource, 97                            |
| docker element, 104, 104, 104, 104, 104, 105, 105, 105 | active_uname, 97                               |
| image attribute, 104                                   | demote_resource, 99                            |
| network attribute, 105                                 | demote_uname, 99                               |
| options attribute, 105                                 | inactive_resource, 97                          |
| promoted-max attribute, 104                            | master_resource, 99                            |
| replicas attribute, 104                                | master_uname, 99                               |
| replicas-per-host attribute, 104                       | operation, 97                                  |
| run-command attribute, 105                             | promote_resource, 99                           |
| Duration, 70   | promote_uname, 99                              |
| duration element, 70, 70, 70, 70, 70, 70, 70, 70       | slave_resource, 99                             |
| hours attribute, 70                                    | slave_uname, 99                                |
| id attribute, 70                                       | start_resource, 97                             |
| minutes attribute, 70                                  | start_uname, 97                                |
| months attribute, 70                                   | stop_resource, 97                              |
| seconds attribute, 70                                  | stop_uname, 97                                 |
| weeks attribute, 70                                    | type, 97                                       |
| years attribute, 70                                    | epoch, 6                                       |
|  | Cluster Option, 6                              |
| E  | exec-time, 127                                 |
|  | Action Status, 127                             |
| election-timeout, 11                                   | exec_time, 64                                  |
| Cluster Option, 11 enable-acl, 10                      | expected, 126                                  |
| Cluster Option, 10                                     | Node Status, 126                               |
| enable-startup-probes, 8                               | expression element, 66, 66, 66, 66, 67, 67, 67 |
| * *  | attribute attribute, 66                        |
| Cluster Option, 8 enabled, 25, 25                      | id attribute, 66                               |
| Action Property, 25, 25                                | operation attribute, 67                        |
| end attribute, 68                                      | type attribute, 66                             |
|  | value attribute, 67                            |
| date_expression element, 68<br>Environment Variable    | value-source attribute, 67                     |
|  |  |
| CRM_alert_   | F  |
| attribute_name, 64                                     | <del>-</del>                                   |
| attribute_value, 64                                    | fail-count-, 13                                |
| desc, 63   | Node attribute, 13                             |
| exec_time, 64  | failure-timeout, 20                            |
| interval, 64   | Resource Option, 20                            |
| kind, 63   | feedback                                       |

| contact information for this manual, xiv                     | green, 85  |
|--|--|
| fence-reaction, 9  | Group Property   |
| Cluster Option, 9  | id, 90   |
| Fencing, 44, 44, 44, 44, 45, 45, 45, 45, 45, 46, 46, 46, 46, | Group Resource Property, 90                                    |
| 46, 46, 46, 46, 47, 47, 47, 47, 47, 47, 48, 48, 48           | Group Resources, 89  |
| Configuration, 42  | Groups, 89, 91   |
| fencing-level  |  |
| devices, 53  | Н  |
| id, 53   | have-quorum, 6   |
| index, 53  | Cluster Property, 6  |
| target, 53   | host-interface attribute, 105                                  |
| target-attribute, 53, 53                                     | network element, 105   |
| target-pattern, 53   | host-netmask attribute, 105                                    |
| Property   | network element, 105   |
| pcmk_action_limit, 45  | host_list, 82  |
| pcmk_delay_base, 45  | Ping Resource Option, 82                                       |
| pcmk_delay_max, 45   | hours attribute, 69, 70  |
| pcmk_host_argument, 45                                       | date_spec element, 69  |
| pcmk_host_check, 45  | duration element, 70   |
| pcmk_host_list, 44   | duration cionicity 70  |
| pcmk_host_map, 44  | 1  |
| pcmk_list_action, 47   | 1  |
| pcmk_list_retries, 47  | id, 18, 24, 53, 90, 92, 120, 121, 121, 122, 122, 125, 127      |
| pcmk_list_timeout, 47  | acl_group, 122   |
| pcmk_monitor_action, 47                                      | acl_permission, 121  |
| pcmk_monitor_retries, 47                                     | acl_role, 120  |
| pcmk_monitor_timeout, 47                                     | acl_target, 121  |
| pcmk_off_action, 46  | Action Property, 24  |
| pcmk_off_retries, 46   | Action Status, 127   |
| pcmk_off_timeout, 46   | Clone Property, 92   |
| pcmk_reboot_action, 46                                       | fencing-level, 53  |
| pcmk_reboot_retries, 46                                      | Group Resource Property, 90                                    |
| pcmk_reboot_timeout, 46                                      | Node Status, 125   |
| pcmk_status_action, 48                                       | Resource, 18   |
| pcmk_status_retries, 48                                      | role, 122  |
| pcmk_status_timeout, 48                                      | id attribute, 30, 33, 35, 36, 65, 66, 68, 69, 70, 72, 72, 104, |
| provides, 44   | 106, 107   |
| stonith-timeout, 44, 46, 46                                  | bundle element, 104  |
| fencing-level, 53, 53, 53, 53, 53, 53                        | date_expression element, 68                                    |
| devices, 53  | date_spec element, 69  |
| id, 53   | duration element, 70   |
| index, 53  | expression element, 66   |
| target, 53   | op_expression element, 72                                      |
| target-attribute, 53, 53                                     | port-mapping element, 106                                      |
| target-pattern, 53   | resource_set element, 36                                       |
| first attribute, 33  | rsc_colocation element, 35                                     |
| rsc_order element, 33  | rsc_expression element, 72                                     |
| first-action attribute, 33                                   | rsc_location element, 30                                       |
| rsc_order element, 33  | rsc_order element, 33  |
|  | rule element, 65   |
| G  | storage-mapping element, 107                                   |
| globally-unique, 92  | image attribute, 104, 104, 104                                 |
| Clone Option, 92   | docker element, 104  |
| Cione Option, 72   | podman element, 104  |

| rkt element, 104                         | M   |
|--|---|
| inactive_resource, 97                    | maintenance, 14, 19                           |
| Notification Environment Variable, 97    | Node attribute, 14                            |
| index, 53                                |   |
| fencing-level, 53                        | Resource Option, 19                           |
| interleave, 92                           | maintenance-mode, 8                           |
| Clone Option, 92                         | Cluster Option, 8                             |
| internal-port attribute, 106             | master_resource, 99                           |
| port-mapping element, 106                | Notification Environment Variable, 99         |
| interval, 24, 64, 127                    | master_uname, 99                              |
| Action Property, 24                      | Notification Environment Variable, 99         |
| Action Status, 127                       | Meta-attributes, 109                          |
| interval attribute, 72                   | migrate-on-red, 86                            |
| op_expression element, 72                | migration-limit, 7                            |
| in_ccm, 126                              | Cluster Option, 7                             |
|  | migration-threshold, 20                       |
| Node Status, 126                         | Resource Option, 20                           |
| ip-range-start attribute, 105            | minutes attribute, 70                         |
| network element, 105                     | duration element, 70                          |
| is-managed, 19                           | monthdays attribute, 69                       |
| Resource Option, 19                      | date_spec element, 69                         |
| _  | months attribute, 69, 70                      |
| J  | date_spec element, 69                         |
| join, 126                                | duration element, 70                          |
| Node Status, 126                         | moon attribute, 69                            |
| join-finalization-timeout, 11            | date_spec element, 69                         |
| Cluster Option, 11                       | Moving, 80                                    |
| join-integration-timeout, 11             | Resources, 80                                 |
| Cluster Option, 11                       | multiple-active, 20                           |
| Claster Option, 11                       | Resource Option, 20                           |
| <b>V</b>                                 | multiplier, 82                                |
| K  | Ping Resource Option, 82                      |
| kind, 63, 121                            | Ting Resource Option, 62                      |
| acl_permission, 121                      | NI .  |
| kind attribute, 33                       | N   |
| rsc_order element, 33                    | Nagios Plugins, 18                            |
|  | Resources, 18                                 |
| L  | name, 24                                      |
| last-failure-, 14                        | Action Property, 24                           |
| Node attribute, 14                       | name attribute, 72                            |
|  | op_expression element, 72                     |
| last-rc-change, 127                      | network attribute, 105, 105, 105              |
| Action Status, 127                       | docker element, 105                           |
| last-run, 127                            | podman element, 105                           |
| Action Status, 127                       | rkt element, 105                              |
| Linux Standard Base                      | network element, 105, 105, 105, 105, 105, 106 |
| Resources, 16                            | add-host attribute, 105                       |
| Location                                 | control-port attribute, 106                   |
| Determine by Rules, 73                   | host-interface attribute, 105                 |
| Location Constraint, 30                  | host-netmask attribute, 105                   |
| Resource Discovery, 31                   | ip-range-start attribute, 105                 |
| rsc_location element, 30                 | Networking, 105                               |
| Location Relative to Other Resources, 34 | <u>~</u>                                      |
| LSB, 16                                  | no-quorum-policy, 7                           |
| Resources, 16                            | Cluster Option, 7                             |
|  | Node  |

| attribute, 12  | notify, 92                            |
|--|---------------------------------------|
| #digests-, 14  | Clone Option, 92                      |
| #node-unfenced, 14   | num_updates, 6                        |
| fail-count-, 13  | Cluster Option, 6                     |
| last-failure-, 14  |                                       |
| maintenance, 14  | 0                                     |
| probe_complete, 14   | •                                     |
| resource-discovery-enabled, 14                                 | object-type, 121                      |
| shutdown, 14   | acl_permission, 121                   |
| site-name, 14  | OCF, 16                               |
| standby, 14  | Resources, 16                         |
| terminate, 14  | OCF_FAILED_MASTER, 97                 |
| Score, 29  | OCF_NOT_RUNNING, 97                   |
| Status, 125  | OCF_RESKEY_CRM_meta_notify_           |
| crm-debug-origin, 126  | active_resource, 97                   |
| crmd, 126  | active_uname, 97                      |
| expected, 126  | demote_resource, 99                   |
| •  | demote_uname, 99                      |
| id, 125  | inactive_resource, 97                 |
| in_ccm, 126  | master_resource, 99                   |
| join, 126  | master_uname, 99                      |
| uname, 126   | operation, 97                         |
| node, 63   | promote_resource, 99                  |
| Node attribute, 13, 14, 14, 14, 14, 14, 14, 14, 14, 14, 14     | promote_uname, 99                     |
| node attribute, 30   | slave_resource, 99                    |
| rsc_location element, 30                                       | slave_uname, 99                       |
| Node Attribute Expression, 66                                  | start_resource, 97                    |
| Node Attributes, 108   | start_uname, 97                       |
| Node health  | stop_resource, 97                     |
| custom, 86   | stop_uname, 97                        |
| green, 85  | type, 97                              |
| migrate-on-red, 86   | OCF_RUNNING_MASTER, 97                |
| none, 86   | OCF_SUCCESS, 97                       |
| only-green, 86   | on-fail, 25                           |
| progressive, 86  | Action Property, 25                   |
| red, 85  | only-green, 86                        |
| score, 85  | op-digest, 128                        |
| yellow, 85   | Action Status, 128                    |
| Node Status, 125, 126, 126, 126, 126, 126, 126                 | op-status, 127                        |
| node-attribute attribute, 35                                   | Action Status, 127                    |
| rsc_colocation element, 35                                     | Open Cluster Framework                |
| node-health-base, 10   | =                                     |
| Cluster Option, 10   | Resources, 16                         |
| node-health-green, 10  | operation, 97, 127                    |
| Cluster Option, 10   | Action Status, 127                    |
| node-health-red, 10  | Notification Environment Variable, 97 |
| Cluster Option, 10   | operation attribute, 67, 69           |
| node-health-strategy, 10                                       | date_expression element, 69           |
| Cluster Option, 10   | expression element, 67                |
| node-health-yellow, 10   | Operation History, 126                |
| Cluster Option, 10   | Opt-In Clusters, 31                   |
| nodeid, 63   | Opt-Out Clusters, 32                  |
| none, 86   | Option                                |
| Notification Environment Variable, 97, 97, 97, 97, 97, 97, 97, | admin_epoch, 6                        |
| 97. 97. 97. 99. 99. 99. 99. 99. 99. 99                         | batch-limit, 7                        |
| ,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,                        |                                       |

| clone-max, 92                 | stonith-action, 8                     |
|-------------------------------|---------------------------------------|
| clone-min, 92                 | stonith-enabled, 8                    |
| clone-node-max, 92            | stonith-max-attempts, 8               |
| cluster-delay, 9              | stonith-timeout, 8                    |
| cluster-ipc-limit, 9          | stonith-watchdog-timeout, 8           |
| cluster-recheck-interval, 10  | stop-all-resources, 7                 |
| concurrent-fencing, 9         | stop-orphan-actions, 8                |
| Configuration Version, 6      | stop-orphan-resources, 8              |
| dampen, 82                    | symmetric-cluster, 7                  |
| dc-deadtime, 9                | target-role, 19                       |
| election-timeout, 11          | timeout, 60                           |
| enable-acl, 10                | timestamp-format, 60                  |
| enable-startup-probes, 8      | transition-delay, 11                  |
| epoch, 6                      | validate-with, 6                      |
| failure-timeout, 20           | options attribute, 105, 105, 105, 107 |
| fence-reaction, 9             | docker element, 105                   |
| globally-unique, 92           | podman element, 105                   |
| host_list, 82                 | rkt element, 105                      |
| interleave, 92                | storage-mapping element, 107          |
| is-managed, 19                | op_expression element, 72, 72, 72     |
| join-finalization-timeout, 11 | id attribute, 72                      |
| join-integration-timeout, 11  | interval attribute, 72                |
| maintenance, 19               | name attribute, 72                    |
| maintenance-mode, 8           | ordered, 92                           |
| migration-limit, 7            | Clone Option, 92                      |
| migration-threshold, 20       | Ordering                              |
| multiple-active, 20           | action, 95                            |
| multiplier, 82                | role, 95                              |
| no-quorum-policy, 7           | rsc-role, 94                          |
| node-health-base, 10          | with-rsc-role, 94                     |
| node-health-green, 10         | Ordering Constraint, 32               |
| node-health-red, 10           | rsc_order element, 33                 |
| node-health-strategy, 10      | Ordering Constraints, 94, 94, 95, 95  |
| node-health-yellow, 10        |                                       |
| notify, 92                    | P                                     |
| num_updates, 6                | -                                     |
| ordered, 92                   | pcmk_action_limit, 45                 |
| pe-error-series-max, 10       | Fencing, 45                           |
| pe-input-series-max, 10       | pcmk_delay_base, 45                   |
| pe-warn-series-max, 10        | Fencing, 45<br>pcmk_delay_max, 45     |
| placement-strategy, 10        | Fencing, 45                           |
| priority, 19                  | <u> </u>                              |
| priority-fencing-delay, 9     | pcmk_host_argument, 45 Fencing, 45    |
| promotable, 93                | <u> </u>                              |
| promoted-max, 93              | pcmk_host_check, 45                   |
| promoted-node-max, 93         | Fencing, 45                           |
| remove-after-stop, 11         | pcmk_host_list, 44 Fencing, 44        |
| requires, 20                  | _                                     |
| resource-stickiness, 20       | pcmk_host_map, 44                     |
| shutdown-escalation, 11       | Fencing, 44                           |
| shutdown-lock, 11             | pcmk_list_action, 47                  |
| shutdown-lock-limit, 11       | Fencing, 47                           |
| start-failure-is-fatal, 8     | pcmk_list_retries, 47                 |
| startup-fencing, 11           | Fencing, 47                           |

| pcmk_list_timeout, 47                                  | id attribute, 106                     |
|--|---------------------------------------|
| Fencing, 47  | internal-port attribute, 106          |
| pcmk_monitor_action, 47                                | port attribute, 106                   |
| Fencing, 47  | range attribute, 106                  |
| pcmk_monitor_retries, 47                               | Prerequisites, 104                    |
| Fencing, 47  | Primitive, 107                        |
| pcmk_monitor_timeout, 47                               | priority, 19                          |
| Fencing, 47  | Resource Option, 19                   |
| pcmk_off_action, 46                                    | priority-fencing-delay, 9             |
| Fencing, 46  | Cluster Option, 9                     |
| pcmk_off_retries, 46                                   | probe_complete, 14                    |
| Fencing, 46  | Node attribute, 14                    |
| pcmk_off_timeout, 46                                   | progressive, 86                       |
| Fencing, 46  | Promotable, 91                        |
| pcmk_reboot_action, 46                                 | promotable, 93                        |
| Fencing, 46  | Clone Option, 93                      |
| pcmk_reboot_retries, 46                                | Promotable Clone Resources, 91        |
| Fencing, 46  | promoted-max, 93                      |
| <u> </u>   | ±                                     |
| pcmk_reboot_timeout, 46                                | Clone Option, 93                      |
| Fencing, 46  | promoted-max attribute, 104, 104, 104 |
| pcmk_status_action, 48                                 | docker element, 104                   |
| Fencing, 48  | podman element, 104                   |
| pcmk_status_retries, 48                                | rkt element, 104                      |
| Fencing, 48  | promoted-node-max, 93                 |
| pcmk_status_timeout, 48                                | Clone Option, 93                      |
| Fencing, 48  | promote_resource, 99                  |
| pe-error-series-max, 10                                | Notification Environment Variable, 99 |
| Cluster Option, 10                                     | promote_uname, 99                     |
| pe-input-series-max, 10                                | Notification Environment Variable, 99 |
| Cluster Option, 10                                     | Property                              |
| pe-warn-series-max, 10                                 | cib-last-written, 6                   |
| Cluster Option, 10                                     | class, 18                             |
| Ping Resource  | cluster-infrastructure, 7             |
| Option   | cluster-name, 7                       |
| dampen, 82   | dc-uuid, 6                            |
| host_list, 82  | dc-version, 7                         |
| multiplier, 82   | enabled, 25, 25                       |
| Ping Resource Option, 82, 82, 82                       | have-quorum, 6                        |
| placement-strategy, 10                                 | id, 18, 24, 92                        |
| Cluster Option, 10                                     | interval, 24                          |
| podman   | name, 24                              |
| Bundle, 103  | on-fail, 25                           |
| podman element, 104, 104, 104, 104, 104, 105, 105, 105 | pcmk_action_limit, 45                 |
| image attribute, 104                                   | pcmk_delay_base, 45                   |
| network attribute, 105                                 | pcmk_delay_max, 45                    |
| options attribute, 105                                 | pcmk_host_argument, 45                |
| promoted-max attribute, 104                            | pcmk_host_check, 45                   |
| replicas attribute, 104                                | pcmk_host_list, 44                    |
| replicas-per-host attribute, 104                       | pcmk_host_map, 44                     |
| run-command attribute, 105                             | pcmk_list_action, 47                  |
| port attribute, 106                                    | pcmk_list_retries, 47                 |
| port-mapping element, 106                              | pcmk_list_timeout, 47                 |
| port-mapping, 106                                      | pcmk_monitor_action, 47               |
| port-mapping element, 106, 106, 106, 106               | pcmk_monitor_retries, 47              |
| 1 11 0 · · · · · · · · · · · · · · · · ·               |                                       |

| pcmk_monitor_timeout, 47   | Alerts, 59  |
|--|---|
| pcmk_off_action, 46  | Bundle, 103   |
| pcmk_off_retries, 46   | Container, 104  |
| pcmk_off_timeout, 46   | Meta-attributes, 109  |
| pcmk_reboot_action, 46   | Networking, 105   |
| pcmk_reboot_retries, 46  | Node Attributes, 108  |
| pcmk_reboot_timeout, 46  | Prerequisites, 104  |
| pcmk_status_action, 48   | Primitive, 107  |
| pcmk_status_retries, 48  | Storage, 106  |
| pcmk_status_timeout, 48  | class, 15   |
| provider, 18   | Clones, 91  |
| provides, 44   | Constraint, 29  |
| role, 25   | Group Property  |
| stonith-timeout, 44, 46, 46  | id, 90  |
| timeout, 24  | Groups, 89  |
| type, 18   | Location  |
| provider, 18   | Determine by Rules, 73  |
| Resource, 18   | Location Relative to Other Resources, 34  |
| provider attribute, 72   | LSB, 16   |
| rsc_expression element, 72   | Moving, 80  |
| provides, 44   | Nagios Plugins, 18  |
| Fencing, 44  | OCF, 16   |
|  | Option  |
| Q  | failure-timeout, 20   |
| •  | is-managed, 19  |
| queue-time, 127  | maintenance, 19   |
| Action Status, 127   | migration-threshold, 20   |
| <b>D</b>   | multiple-active, 20   |
| R  | priority, 19  |
| range attribute, 106   | requires, 20  |
| port-mapping element, 106  | resource-stickiness, 20   |
| rc, 63   | target-role, 19   |
| rc-code, 127   | Promotable, 91  |
| 10 0000, 127   | 1 10111014010, 91   |
| Action Status, 127   | Property  |
|  |   |
| Action Status, 127   | Property  |
| Action Status, 127 recipient, 63   | Property class, 18  |
| Action Status, 127 recipient, 63 red, 85   | Property class, 18 id, 18 provider, 18  |
| Action Status, 127 recipient, 63 red, 85 reference, 121  | Property<br>class, 18<br>id, 18   |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121  | Property class, 18 id, 18 provider, 18 type, 18   |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104   | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36  |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11   | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29  |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104   | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32  |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104   | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18  |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104   | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17  |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104 rkt element, 104  | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17 Systemd, 17  |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104 rkt element, 104 replicas-per-host attribute, 104, 104, 104   | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17 Systemd, 17 Upstart, 17  |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104 rkt element, 104 replicas-per-host attribute, 104, 104, 104 docker element, 104   | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17 Systemd, 17 Upstart, 17 Resource Discovery, 31   |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104 rkt element, 104 replicas-per-host attribute, 104, 104, 104 docker element, 104 replicas-per-host attribute, 104, 104, 104 replicas-per-host attribute, 104 require-all attribute, 36   | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17 Systemd, 17 Upstart, 17 Resource Discovery, 31 Resource Option, 19, 19, 19, 19, 20, 20, 20, 20, 20   |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104 rkt element, 104 replicas-per-host attribute, 104, 104, 104 docker element, 104 replicas-per-host attribute, 104, 104, 104 replicas-per-host attribute, 104 require-all attribute, 36 resource_set element, 36                                  | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17 Systemd, 17 Upstart, 17 Resource Discovery, 31 Resource Option, 19, 19, 19, 20, 20, 20, 20, 20 Resource Set, 36, 36  |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104 rkt element, 104 replicas-per-host attribute, 104, 104, 104 docker element, 104 replicas-per-host attribute, 104, 104, 104 replicas-per-host attribute, 104 require-all attribute, 36 resource_set element, 36 requires, 20                     | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17 Systemd, 17 Upstart, 17 Resource Discovery, 31 Resource Option, 19, 19, 19, 20, 20, 20, 20, 20 Resource Set, 36, 36 resource-discovery attribute, 31   |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104 rkt element, 104 replicas-per-host attribute, 104, 104, 104 docker element, 104 replicas-per-host attribute, 104, 104, 104 replicas-per-host attribute, 104 require-all attribute, 36 resource_set element, 36 requires, 20 Resource Option, 20 | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17 Systemd, 17 Upstart, 17 Resource Discovery, 31 Resource Option, 19, 19, 19, 20, 20, 20, 20, 20 Resource-discovery attribute, 31 rsc_location element, 31   |
| Action Status, 127 recipient, 63 red, 85 reference, 121 acl_permission, 121 remove-after-stop, 11 Cluster Option, 11 replicas attribute, 104, 104, 104 docker element, 104 podman element, 104 rkt element, 104 replicas-per-host attribute, 104, 104, 104 docker element, 104 replicas-per-host attribute, 104, 104, 104 replicas-per-host attribute, 104 require-all attribute, 36 resource_set element, 36 requires, 20                     | Property class, 18 id, 18 provider, 18 type, 18 Resource Set, 36 Score, 29 Start Order, 32 STONITH, 18 System Services, 17 Systemd, 17 Upstart, 17 Resource Discovery, 31 Resource Option, 19, 19, 19, 19, 20, 20, 20, 20, 20 Resource Set, 36, 36 resource-discovery attribute, 31 rsc_location element, 31 resource-discovery-enabled, 14 |

| Groups, 91   | resource-discovery attribute, 31                      |
|--|---|
| Resource Option, 20                                  | rsc attribute, 30                                     |
| Resources, 16, 16, 16, 16, 17, 17, 17, 18, 18, 80    | rsc-pattern attribute, 30                             |
| resource_set element, 36, 36, 36, 36, 37, 37         | score attribute, 30                                   |
| action attribute, 37                                 | rsc_order element, 33, 33, 33, 33, 33, 33, 33, 33, 33 |
| id attribute, 36                                     | first attribute, 33                                   |
| require-all attribute, 36                            | first-action attribute, 33                            |
| role attribute, 36                                   | id attribute, 33                                      |
| score attribute, 37                                  | kind attribute, 33                                    |
| sequential attribute, 36                             | symmetrical attribute, 33                             |
| Return Code  | then attribute, 33                                    |
| OCF_FAILED_MASTER, 97                                | then-action attribute, 33                             |
| OCF_NOT_RUNNING, 97                                  | Rule, 65  |
| OCF_RUNNING_MASTER, 97                               | Controlling Cluster Options, 77                       |
| OCF_SUCCESS, 97                                      | Date/Time Expression, 68                              |
| rkt  | Date Specification, 69                                |
| Bundle, 103  | Duration, 70  |
| rkt element, 104, 104, 104, 104, 104, 105, 105, 105  | Determine Resource Location, 73                       |
| image attribute, 104                                 | Node Attribute Expression, 66                         |
| network attribute, 105                               | rule element, 65, 65, 66, 66, 66, 66                  |
| options attribute, 105                               | boolean-op attribute, 66                              |
| promoted-max attribute, 104                          | id attribute, 65                                      |
| replicas attribute, 104                              | role attribute, 66                                    |
| replicas-per-host attribute, 104                     | score attribute, 66                                   |
| run-command attribute, 105                           | score-attribute attribute, 66                         |
| role, 25, 95, 122                                    | run-command attribute, 105, 105, 105                  |
| Action Property, 25                                  | docker element, 105                                   |
| id, 122  | podman element, 105                                   |
| Ordering Constraints, 95                             | rkt element, 105                                      |
| role attribute, 36, 66                               |   |
| resource_set element, 36                             | S   |
| rule element, 66                                     | Score, 29, 29   |
| rsc, 63  | score, 85   |
| rsc attribute, 30, 35                                | score attribute, 30, 35, 37, 66                       |
| rsc_colocation element, 35                           | resource_set element, 37                              |
| rsc_location element, 30                             | rsc_colocation element, 35                            |
| rsc-pattern attribute, 30                            | rsc_location element, 30                              |
| rsc_location element, 30                             | rule element, 66                                      |
| rsc-role, 94   | score-attribute attribute, 66                         |
| Ordering Constraints, 94                             | rule element, 66                                      |
| rsc_colocation element, 34, 34, 35, 35, 35, 35, 35   | seconds attribute, 70                                 |
| id attribute, 35                                     | duration element, 70                                  |
| node-attribute attribute, 35                         | sequence, 63  |
| rsc attribute, 35                                    | sequential attribute, 36                              |
| score attribute, 35                                  | resource_set element, 36                              |
| with-rsc attribute, 35                               | Setting Options with Rules, 77                        |
| rsc_expression element, 72, 72, 72, 72               | shutdown, 14  |
| class attribute, 72                                  | Node attribute, 14                                    |
| id attribute, 72                                     | shutdown-escalation, 11                               |
| provider attribute, 72                               | Cluster Option, 11                                    |
| type attribute, 72                                   | shutdown-lock, 11                                     |
| rsc_location element, 30, 30, 30, 30, 30, 30, 30, 31 | Cluster Option, 11                                    |
| id attribute, 30                                     | shutdown-lock-limit, 11                               |
| node attribute 30                                    | ,   |

| Cluster Option, 11                    | stonith-max-attempts, 8                               |
|---------------------------------------|---|
| site-name, 14                         | Cluster Option, 8                                     |
| Node attribute, 14                    | stonith-timeout, 8, 44, 46, 46                        |
| slave_resource, 99                    | Cluster Option, 8                                     |
| Notification Environment Variable, 99 | Fencing, 44, 46, 46                                   |
| slave_uname, 99                       | stonith-watchdog-timeout, 8                           |
| Notification Environment Variable, 99 | Cluster Option, 8                                     |
| source-dir attribute, 107             | stop-all-resources, 7                                 |
| storage-mapping element, 107          | Cluster Option, 7                                     |
| source-dir-root attribute, 107        | stop-orphan-actions, 8                                |
| storage-mapping element, 107          | Cluster Option, 8                                     |
| standby, 14                           | stop-orphan-resources, 8                              |
| Node attribute, 14                    | Cluster Option, 8                                     |
| start attribute, 68                   | stop_resource, 97                                     |
| date_expression element, 68           | Notification Environment Variable, 97                 |
| Start Order, 32                       | stop_uname, 97  |
| start-failure-is-fatal, 8             | Notification Environment Variable, 97                 |
| Cluster Option, 8                     | Storage, 106  |
| startup-fencing, 11                   | storage element, 106                                  |
| Cluster Option, 11                    | storage-mapping element, 106, 107, 107, 107, 107, 107 |
| start_resource, 97                    | id attribute, 107                                     |
| Notification Environment Variable, 97 | options attribute, 107                                |
| start_uname, 97                       | source-dir attribute, 107                             |
|                                       | source-dir-root attribute, 107                        |
| Notification Environment Variable, 97 |   |
| status, 64                            | target-dir attribute, 107                             |
| Status, 125                           | symmetric-cluster, 7                                  |
| call-id, 127                          | Cluster Option, 7                                     |
| crm-debug-origin, 126, 128            | symmetrical attribute, 33                             |
| crmd, 126                             | rsc_order element, 33                                 |
| crm_feature_set, 127                  | Symmetrical Clusters, 32                              |
| exec-time, 127                        | System Service  |
| expected, 126                         | Resources, 17   |
| id, 125, 127                          | System Services, 17                                   |
| interval, 127                         | Systemd, 17   |
| in_ccm, 126                           | Resources, 17   |
| join, 126                             |   |
| last-rc-change, 127                   | T   |
| last-run, 127                         | target, 53  |
| op-digest, 128                        | fencing-level, 53                                     |
| op-status, 127                        | target-attribute, 53, 53                              |
| operation, 127                        | fencing-level, 53, 53                                 |
| queue-time, 127                       | target-dir attribute, 107                             |
| rc-code, 127                          | storage-mapping element, 107                          |
| transition-key, 128                   | target-pattern, 53                                    |
| transition-magic, 128                 | fencing-level, 53                                     |
| uname, 126                            | target-role, 19                                       |
| Status of a Node, 125                 | Resource Option, 19                                   |
| STONITH, 18                           | target_rc, 64   |
| Configuration, 42                     | 6 - 1   |
| Resources, 18                         | task, 63  |
| stonith-action, 8                     | terminate, 14   |
| Cluster Option, 8                     | Node attribute, 14                                    |
| stonith-enabled, 8                    | then attribute, 33                                    |
| Cluster Option, 8                     | rsc_order element, 33                                 |

| then-action attribute, 33             | attribute attribute          |
|---------------------------------------|------------------------------|
| rsc_order element, 33                 | expression element, 66       |
| timeout, 24, 60                       | boolean-op attribute         |
| Action Property, 24                   | rule element, 66             |
| timestamp, 63                         | class attribute              |
| timestamp-format, 60                  | rsc_expression element, 72   |
| timestamp_epoch, 63                   | description attribute        |
| timestamp_usec, 63                    | bundle element, 104          |
| transition-delay, 11                  | end attribute                |
| Cluster Option, 11                    | date_expression element, 68  |
| transition-key, 128                   | first attribute              |
| Action Status, 128                    | rsc_order element, 33        |
| transition-magic, 128                 | first-action attribute       |
| Action Status, 128                    | rsc_order element, 33        |
| type, 18, 97                          | hours attribute              |
| Notification Environment Variable, 97 | date_spec element, 69        |
| Resource, 18                          | duration element, 70         |
| type attribute, 66, 72                | id attribute                 |
| expression element, 66                | bundle element, 104          |
| rsc_expression element, 72            | date_expression element, 68  |
|                                       | date_spec element, 69        |
| U                                     | duration element, 70         |
| uname, 126                            | expression element, 66       |
| Node Status, 126                      | op_expression element, 72    |
| Upstart, 17                           | port-mapping element, 106    |
| Resources, 17                         | resource_set element, 36     |
| ,                                     | rsc_colocation element, 35   |
| V                                     | rsc_expression element, 72   |
| validate-with, 6                      | rsc_location element, 30     |
| Cluster Option, 6                     | rsc_order element, 33        |
| value attribute, 67                   | rule element, 65             |
| expression element, 67                | storage-mapping element, 107 |
| value-source attribute, 67            | image attribute              |
| expression element, 67                | docker element, 104          |
| -                                     | podman element, 104          |
| version, 63                           | rkt element, 104             |
| W                                     | internal-port attribute      |
|                                       | port-mapping element, 106    |
| weekdays attribute, 69                | interval attribute           |
| date_spec element, 69                 | op_expression element, 72    |
| weeks attribute, 69, 70               | kind attribute               |
| date_spec element, 69                 | rsc_order element, 33        |
| duration element, 70                  | minutes attribute            |
| weekyears attribute, 69               | duration element, 70         |
| date_spec element, 69                 | monthdays attribute          |
| with-rsc attribute, 35                | date_spec element, 69        |
| rsc_colocation element, 35            | months attribute             |
| with-rsc-role, 94                     | date_spec element, 69        |
| Ordering Constraints, 94              | duration element, 70         |
| V                                     | moon attribute               |
| X                                     | date_spec element, 69        |
| XML attribute                         | name attribute               |
| action attribute                      | op_expression element, 72    |
| resource set element, 37              | network attribute            |

| docker element, 105                                 | rkt element, 105             |
|---|------------------------------|
| podman element, 105                                 | score attribute              |
| rkt element, 105                                    | resource_set element, 37     |
| network element                                     | rsc_colocation element, 35   |
| add-host attribute, 105                             | rsc_location element, 30     |
| control-port attribute, 106                         | rule element, 66             |
| host-interface attribute, 105                       | score-attribute attribute    |
| host-netmask attribute, 105                         | rule element, 66             |
| ip-range-start attribute, 105                       | seconds attribute            |
| node attribute                                      | duration element, 70         |
| rsc_location element, 30                            | sequential attribute         |
| node-attribute                                      | resource_set element, 36     |
| rsc_colocation element, 35                          | source-dir attribute         |
| operation attribute                                 | storage-mapping element, 107 |
| date_expression element, 69                         | source-dir-root attribute    |
| expression element, 67                              | storage-mapping element, 107 |
| options attribute                                   | start attribute              |
| docker element, 105                                 | date_expression element, 68  |
| podman element, 105                                 | symmetrical attribute        |
| rkt element, 105                                    | rsc_order element, 33        |
| storage-mapping element, 107                        | target-dir attribute         |
| port attribute                                      | storage-mapping element, 107 |
| -   | then attribute               |
| port-mapping element, 106<br>promoted-max attribute | rsc_order element, 33        |
| docker element, 104                                 | then-action attribute        |
|   |                              |
| podman element, 104                                 | rsc_order element, 33        |
| rkt element, 104                                    | type attribute               |
| provider attribute                                  | expression element, 66       |
| rsc_expression element, 72                          | rsc_expression element, 72   |
| range attribute                                     | value attribute              |
| port-mapping element, 106                           | expression element, 67       |
| replicas attribute                                  | value-source attribute       |
| docker element, 104                                 | expression element, 67       |
| podman element, 104                                 | weekdays attribute           |
| rkt element, 104                                    | date_spec element, 69        |
| replicas-per-host attribute                         | weeks attribute              |
| docker element, 104                                 | date_spec element, 69        |
| podman element, 104                                 | duration element, 70         |
| rkt element, 104                                    | weekyears attribute          |
| require-all attribute                               | date_spec element, 69        |
| resource_set element, 36                            | with-rsc attribute           |
| resource-discovery attribute                        | rsc_colocation element, 35   |
| rsc_location element, 31                            | yeardays attribute           |
| role attribute                                      | date_spec element, 69        |
| resource_set element, 36                            | years attribute              |
| rule element, 66                                    | date_spec element, 69        |
| rsc attribute                                       | duration element, 70         |
| rsc_colocation element, 35                          | XML element                  |
| rsc_location element, 30                            | add-host attribute           |
| rsc-pattern attribute                               | network element, 105         |
| rsc_location element, 30                            | bundle element, 104          |
| run-command attribute                               | description attribute, 104   |
| docker element, 105                                 | id attribute, 104            |
| podman element, 105                                 | control-port attribute       |

| network element, 106             | options attribute, 105           |
|----------------------------------|----------------------------------|
| date_expression element, 68      | promoted-max attribute, 104      |
| end attribute, 68                | replicas attribute, 104          |
| id attribute, 68                 | replicas-per-host attribute, 104 |
| operation attribute, 69          | run-command attribute, 105       |
| start attribute, 68              | port-mapping, 106                |
| date_spec element, 69            | port-mapping element             |
| hours attribute, 69              | id attribute, 106                |
| id attribute, 69                 | internal-port attribute, 106     |
| monthdays attribute, 69          | port attribute, 106              |
| months attribute, 69             | range attribute, 106             |
| moon attribute, 69               | resource_set element             |
| weekdays attribute, 69           | action attribute, 37             |
| weeks attribute, 69              | id attribute, 36                 |
| weekyears attribute, 69          | require-all attribute, 36        |
| yeardays attribute, 69           | role attribute, 36               |
| years attribute, 69              | score attribute, 37              |
| docker element, 104              | sequential attribute, 36         |
| image attribute, 104             | rkt element, 104                 |
| network attribute, 105           | image attribute, 104             |
|                                  | =                                |
| options attribute, 105           | network attribute, 105           |
| promoted-max attribute, 104      | options attribute, 105           |
| replicas attribute, 104          | promoted-max attribute, 104      |
| replicas-per-host attribute, 104 | replicas attribute, 104          |
| run-command attribute, 105       | replicas-per-host attribute, 104 |
| duration element, 70             | run-command attribute, 105       |
| hours attribute, 70              | rsc_colocation element, 34       |
| id attribute, 70                 | id attribute, 35                 |
| minutes attribute, 70            | node-attribute attribute, 35     |
| months attribute, 70             | rsc attribute, 35                |
| seconds attribute, 70            | score attribute, 35              |
| weeks attribute, 70              | with-rsc attribute, 35           |
| years attribute, 70              | rsc_expression element           |
| expression element, 66           | class attribute, 72              |
| attribute attribute, 66          | id attribute, 72                 |
| id attribute, 66                 | provider attribute, 72           |
| operation attribute, 67          | type attribute, 72               |
| type attribute, 66               | rsc_location element, 30         |
| value attribute, 67              | id attribute, 30                 |
| value-source attribute, 67       | node attribute, 30               |
| host-interface attribute         | resource-discovery attribute, 31 |
| network element, 105             | rsc attribute, 30                |
| host-netmask attribute           | rsc-pattern attribute, 30        |
| network element, 105             | score attribute, 30              |
| ip-range-start attribute         | rsc_order element, 33            |
| network element, 105             | first attribute, 33              |
| network element, 105             | first-action attribute, 33       |
| op_expression element            | id attribute, 33                 |
| id attribute, 72                 | kind attribute, 33               |
| interval attribute, 72           | symmetrical attribute, 33        |
| name attribute, 72               | then attribute, 33               |
| podman element, 104              | then-action attribute, 33        |
| image attribute, 104             | rule element, 65                 |
| network attribute, 105           | boolean-op attribute. 66         |
|                                  |                                  |

```
id attribute, 65
role attribute, 66
score attribute, 66
score-attribute attribute, 66
storage element, 106
storage-mapping element, 106
id attribute, 107
options attribute, 107
source-dir attribute, 107
source-dir-root attribute, 107
target-dir attribute, 107
xpath, 121
acl_permission, 121
```

#### Y

yeardays attribute, 69 date\_spec element, 69 years attribute, 69, 70 date\_spec element, 69 duration element, 70 yellow, 85