

Resource Management Daemon

Authors

Joseph Gasparakis

Dakshina Ilangovan

1 Introduction

Increasing the number of CPU cores and storage capacity in data centers has greatly increased the density of software components per server. This high density causes software components to compete for finite resources on the platform, possibly producing a negative impact on the performance, service assurance, and the end user's general quality of experience. The Resource Management Daemon (RMD) is a software component that arbitrates between resources and software components.

RMD runs on the server and maps resources (or pools of them) to software components according to a certain policy. For each workload, RMD exposes a RESTful API to receive a policy from external management software (such as Orchestrator, container engine, customer management plane, and so forth) and enforces that policy on the platform.

RMD handles arbitration locally on the compute node (rather than sending telemetry to the upper layers of the stack to make the decision there). The reason is because some use cases, such as latency sensitive applications, require immediate mitigation of resource contention. Certainly one could see how a system that required a remote entity to analyze status and then determine an appropriate action, possibly making the decision there and re-configuring the platform, could create performance issues in cases where low latency and fast response are essential.

This document describes RMD architecture and installation, provides a usage scenario, discusses plugins, and provides an example of using RMD for the Intel® Resource Director Technology (Intel® RDT).

RMD enables a flexible and scalable implementation for resource management. It provides a unified centralized decision making mechanism that can be used in multiple environments.

- RMD can be implemented using any management/policy methodology
- RMD can manage any resource. More resources can be managed over time, per need, using the same RMD engine.
- RMD can be used for implementing resource management/policy on any type of compute node resource.

This document is part of the Network Transformation Experience Kit, which is available at: <https://networkbuilders.intel.com/>

Table of Contents

1	Introduction.....	1
1.1	Terminology.....	3
1.2	Reference Documents	3
2	Architecture.....	3
2.1	RMD Implementation Example	5
3	Integration with Orchestration Layers.....	6
3.1	RMD integration with OpenStack*.....	6
4	Installation Guide.....	6
4.1	Prerequisites	6
4.2	Install as Executable.....	7
4.3	Install from Source.....	7
5	Usage Example	7
5.1	CAT – Node Policy	7
5.2	Querying Cache Inventory Information.....	8
5.3	Workload Policies to Request LLC Allocation	8
5.3.1	Changing and Adding Policies.....	8
5.3.2	Pre-defined Policy.....	8
5.3.3	Custom Policy.....	9
5.3.4	Hospitality Score of a Workload.....	9
6	Plugins	9
7	Summary	10

Figures

Figure 1.	RMD Internal Architecture.....	4
Figure 2.	RMD System Level Architecture	4
Figure 3.	Logical Grouping of Cache Ways into CLOS.....	5
Figure 4.	RMD Partitioning of Cache Ways into Cache Pools using Intel® RDT CLOS Constructs.....	5

Tables

Table 1.	Terminology.....	3
Table 2.	Reference Documents	3
Table 3.	Hospitality Score Calculation.....	9

1.1 Terminology

Table 1. Terminology

Abbreviation	Description
CAT	Cache Allocation Technology
CLOS	Class of Service
NFV	Network Functions Virtualization
PAM	Pluggable Authentication Modules
QoS	Quality of Service
RDT	Intel® Resource Director Technology (Intel® RDT)
RMD	Resource Management Daemon
RPC	Remote Procedure Calls
VNF	Virtual Network Function
vRouter	Virtual Router
vSwitch	Virtual Switch

1.2 Reference Documents

Table 2. Reference Documents

Reference	Source
Intel® Resource Director Technology	https://www.intel.com/content/www/us/en/architecture-and-technology/resource-director-technology.html
Resource Management Daemon repository	https://github.com/intel/rmd
Introduction to Cache Allocation Technology in the Intel® Xeon® Processor E5 v4 Family	https://software.intel.com/en-us/articles/introduction-to-cache-allocation-technology
Introduction to Memory Bandwidth Monitoring in the Intel® Xeon® Processor E5 v4 Family	https://software.intel.com/en-us/articles/introduction-to-memory-bandwidth-monitoring

2 Architecture

RMD is a daemon that runs on the server, one instance per compute node. It exposes a RESTful API for other software layers to make requests and push policies that perform mapping between workloads and resources. These layers need to authenticate themselves through Pluggable Authentication Modules (PAM). There is also HTTPS support accessing the REST API.

The policy can be static or dynamic:

- Static policy is a one off allocation of resources. RMD allocates these resources to the associated workload and then ensures that the workload is still running. If the workload stops running, RMD releases the resources so that they are available for consumption by the next workload that is scheduled on the server.
- When a dynamic policy is received, RMD constantly monitors and re-allocates resources as described in the policy. One can describe an algorithm in that policy on how RMD should do this. As a very simple example, for a multi-way cache, the policy might request for the workload to receive two cache ways from the processor's last level cache (LLC), but to allocate an extra cache way if the throughput on a specific network interface increases above a certain threshold, or to release the extra cache way if the throughput decreases below the threshold.

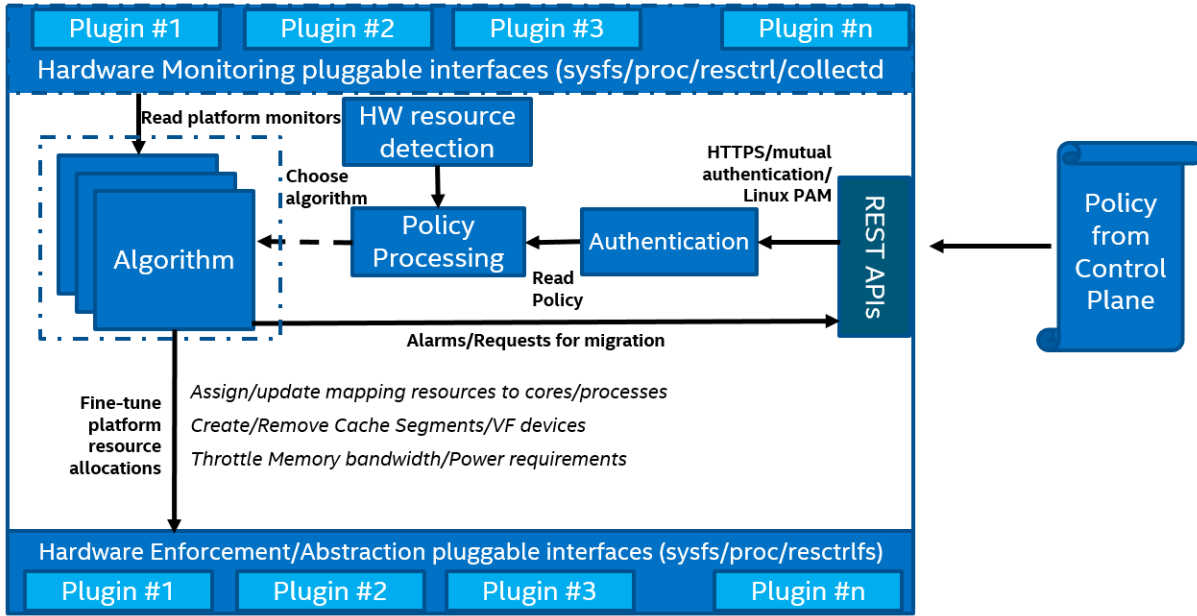


Figure 1. RMD Internal Architecture

Figure 1 above shows the internal architecture of RMD. When a policy is received by an authenticated software entity, it is processed and validated against a hardware resource detection mechanism (if the resources exist, are available, and RMD can control them.). Provided everything is valid, RMD spawns a thread that works in a loop, monitoring resources constantly, and re-allocating them as needed and if the policy is of dynamic nature.

Monitoring and enforcement (re-allocation) of resources is done through pluggable interfaces. For RMD to monitor and control a new resource, you must implement the right plugin and RMD will be able to support this resource onward. For more details, see Section 6.

Figure 2 below shows that structurally, RMD is broken into user and root engines. The user engine provides the RESTful API, runs the monitoring threads, and runs the re-enforcing threads. The engine with elevated privileges (Super user process) performs the actual monitoring and allocations and uses Remote Procedure Calls (RPC) to communicate the information to and from the user engine. Pluggable Authentication Modules (PAM) authentication is also handled by the root process.

RMD uses nosql databases for persisting policy information, shown below in Figure 2 as Bolt and Mongo databases.) If RMD terminates abnormally, it reads the previously stored database information and restores the resource allocations. It also uses configuration files in order to organize resources at start up.

Note: Figure 2 shows that Monitor threads and the Redis datastore are WIP.

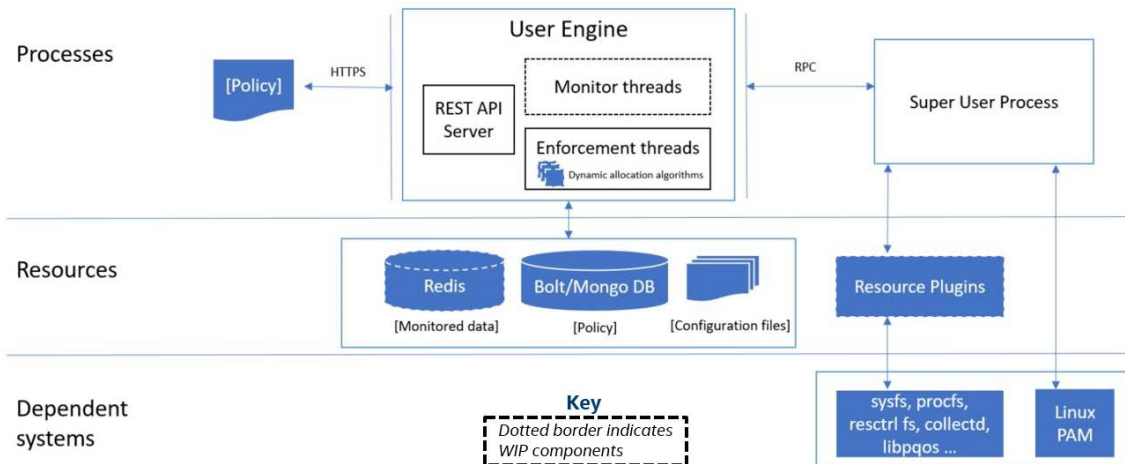


Figure 2. RMD System Level Architecture

2.1 RMD Implementation Example

RMD currently supports management of allocations to Last Level Cache (LLC) using Intel® RDT-CAT (Intel® Resource Director Technology – Cache Allocation Technology). For more information, see <https://www.intel.com/content/www/us/en/architecture-and-technology/resource-director-technology.html>

RMD controls the Last Level Cache on a host using Cache Allocation Technology (CAT) from Intel® Resource Director Technology (Intel® RDT). CAT introduces the concept of describing multiple LLC available on a system. On current generation systems there is one L3 cache per socket but CAT allows for multiple separate L3 caches, representing each LLC using a Cache ID.

Each LLC is measured in cache ways, for example a four-way cache has four cache ways. The size of a cache way is the total size of the LLC divided by the number of cache ways available. The cache ways are logically partitioned using Class of Service (CLOS). This segment of LLC is assigned for sole use by one or more processes or cores. The number of available CLOS is platform dependent.

Currently the allocation of cache ways to processes or cores with CLOS can be controlled using MSR registers or Linux `resctrl` file system starting kernel version 4.10.

Figure 3 illustrates a logical grouping of cache ways into CLOS that determine how cache ways are shared.

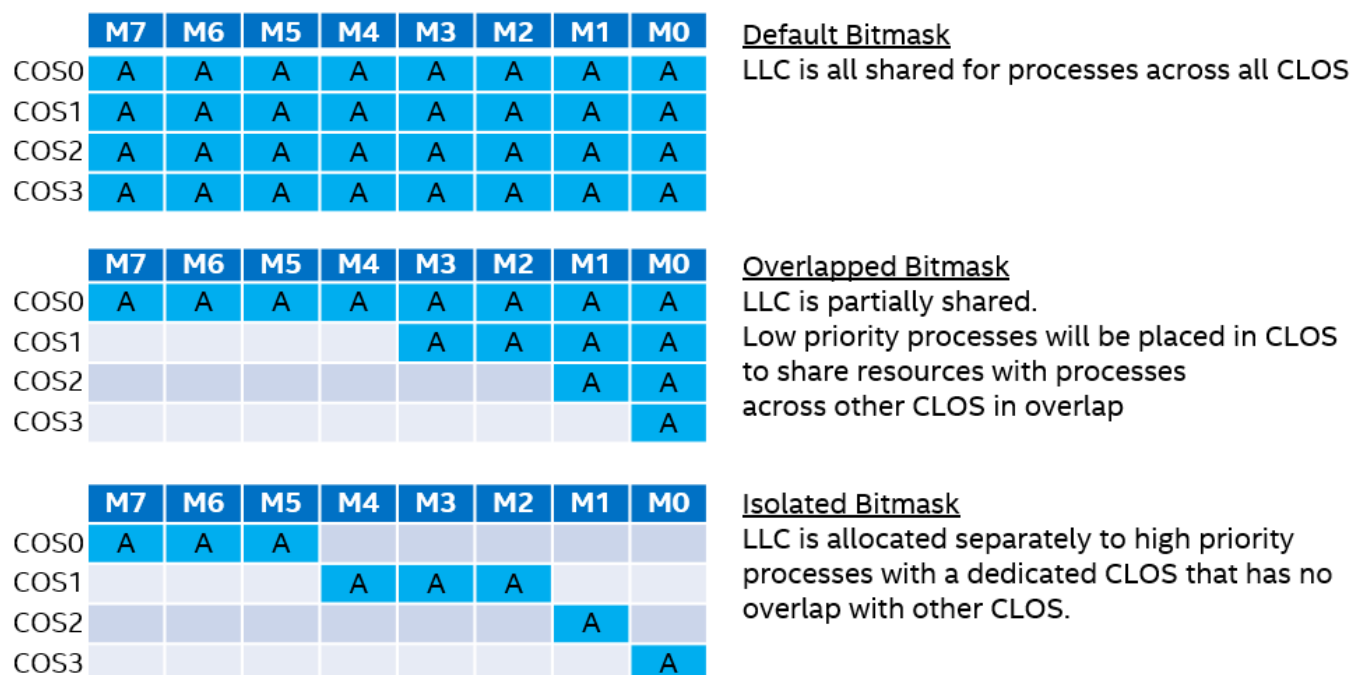


Figure 3. Logical Grouping of Cache Ways into CLOS

RMD provides hooks into the `resctrl` file system to further associate QoS meta-data to Intel® RDT CLOS that can be provided. With features to create CLOS for processes such that cache ways can be allocated in dedicated, shared or overlapped manner, RMD can allocate cache ways for processes into three QoS based cache pools. Refer to Figure 4.

All Cache Ways			
OS Cache Pool	Shared Infrastructure Cache Pool		
	Guaranteed Cache Pool Max=Min>0	Best Effort Cache Pool Max>Min>0	Shared Cache Pool Max=Min=0

Figure 4. RMD Partitioning of Cache Ways into Cache Pools using Intel® RDT CLOS Constructs

The **Guaranteed** cache pool allocates dedicated cache ways for all processes requesting from it. This request is provided to RMD as `Max_Cache=Min_Cache`, `Max_Cache>0`, and `Min_Cache>0`.

The **Best Effort** cache pool allocates cache ways within the window of `Min_Cache` and `Max_Cache` described where `Max_Cache>0` and `Min_Cache>0`. RMD attempts to allocate `Max_Cache` upon request and can oversubscribe processes on host by changing allocations based on usage. This cache pool is leveraged with dynamic re-allocations.

User Guide | Resource Management Daemon

The **Shared** cache pool allocates all processes to a fixed size CLOS (set of cache ways). All processes in this pool share their cache ways. Low priority/short-lived processes can be subscribed to this pool. The number of processes in this pool is configurable in RMD to avoid any contention.

Two other statically configured cache pools are provided:

- The **OS** cache pool is dedicated for operating system processes.
- The **Shared Infrastructure** cache pool is configured in overlap with all available cache pools except the OS cache pool. This enables data sharing between infrastructure processes subscribed to other processes that can be scheduled on other cache pools. A typical example would be for a network functions virtualization (NFV) use case with virtual switch (vSwitch) or virtual router (vRouter) processes that require sharing of data in LLC with other virtual network functions (VNFs).

3 Integration with Orchestration Layers

From an orchestration perspective, interaction with a physical host is limited to the initial allocation of resources and the characterization of nodes to improve the placement of workloads. Runtime handling of workload requirements is not part of the core functionality.

RMD acting as an external component to orchestration environments like OpenStack* or Kubernetes* offers the functionality to manage different node capabilities and provide configurability to specify the policies and changes required by the sys-admin, the tenant or the workload deployment requirements. It provides a unified API for the resource orchestrator to enable high frequency and autonomic control to the node within a set of guardrails defined by a node policy to improve the workload functional and non-functional characteristics.

3.1 RMD integration with OpenStack*

Nova is the compute component in OpenStack*. It is responsible to provide awareness to OpenStack on the compute node capabilities and also provisions and manages the lifecycle of workloads (VNFs) scheduled on it.

RMD as a component on each compute node external to OpenStack services executes policies from two actors, namely the Sys-admin and the VNF Manager (VNF-M). The compute node needs to be provisioned with a set of RMD resource plugins to support the Sys-admin and the VNF-M requests.

The Sys-admin can provide the “Node Policy” to enable a subset of host resources and supply node level configuration for those resources. The VNF-M supplies VNF requirements for the resources that can be provisioned or controlled by RMD during the lifecycle of the VNF.

The interaction from the Sys-admin happens through the transfer of “Node Policy” to the compute node which will be utilized by RMD. The RMD component is in charge of verifying the applicability of the policy and to communicate back to Nova the characteristics enabled in the node.

The VNF-M is able to request the reinforcement of its policy on a node by using the Compute API provided by Nova. The requirement should go to the Nova Compute agent, and then to the RMD component, which should be able to apply the specific changes to the hardware using RMD resource plugins. The VNF demanding specific characteristics will be delivered to the right compute node that provides the matching specs.

4 Installation Guide

RMD is currently supported as a system daemon running on generic Linux platforms/x86 platforms.

The current release supports the static allocation and QoS management of a LLC hardware resource using Intel® RDT – CAT technology.

4.1 Prerequisites

The prerequisites for deploying RMD are as follows:

- RMD is supported on select x86 CPU SKUs that enable Intel® RDT – CAT.
 - RMD is supported on select Linux kernel versions (4.10 or greater) that enable `resctrl fs`. For details, refer to https://www.kernel.org/doc/Documentation/x86/intel_rdt_ui.txt
1. To enable `resctrl fs` with Intel® RDT – CAT support add the following to the GRUB command line:

```
rdt=13cat
```
 2. Verify this on the deployment node using the command line:

```
cat /proc/filesystems | grep resctrl
```
 3. Mount the `resctrl fs` using the command line:

```
mount -t resctrl resctrl /sys/fs/resctrl
```
 4. Install the Linux PAM, Berkeley DB and OpenSSL packages.
 5. To install the Debian* or Ubuntu* OS, enter on the command line:

```
sudo apt-get install openssl libpam0g-dev db-util
```

User Guide | Resource Management Daemon

6. To install the Red Hat* Linux* OS, enter on the command line:

```
sudo dnf install openssl pam-devel db4-utils
```

4.2 Install as Executable

RMD follows versioned history and with each stable release, an executable will be released on Github in the location:

<https://github.com/intel/rmd/releases>

The executable file must be executed with super user privileges and can be added to the system path.

4.3 Install from Source

The RMD source code can be cloned from the Github repository: <https://github.com/intel/rmd.git>

To build RMD from source, the system should support the Golang environment. For instructions, refer to <https://golang.org/doc/install> and ensure `rmd` repo is cloned to `GOPATH src` directory.

The following commands are a shorthand to get started with building and running RMD:

1. Download RMD source code dependencies and build:

```
make
```
2. Install RMD binary to build/ directory:

```
scripts/install.sh
```
3. Run `rmd` from build directory in debug mode:

```
rmd -d
```
4. Test RMD REST interfaces exposed over HTTP in debug mode and query LLC information:

```
wget http://localhost:8081/v1/cache/l3
```

The above steps will help users and developers to quickly get started with building and running RMD using default configuration files. For detailed documentation, refer to the links below.

- Configuration guide: <https://github.com/intel/rmd/blob/master/docs/ConfigurationGuide.md>
- User guide: <https://github.com/intel/rmd/blob/master/docs/UserGuide.md>
- Developer guide: <https://github.com/intel/rmd/blob/master/docs/DeveloperQuickStart.md>

5 Usage Example

5.1 CAT – Node Policy

In order to associate cache ways to workloads efficiently using RMD, the host should be prepared by applying “Node Policy” using the RMD configuration file.

The RMD configuration file supports fields to specify cache allocation groups and quantities for reserved processes (such as operating system) and infrastructure processes in general.

The default configuration file is installed at `/usr/local/etc/rmd/rmd.toml`.

```
[OSGroup] # mandatory
cacheways = 1
cpuset = "0-1"

[InfraGroup] # optional
cacheways = 19
cpuset = "2-3"
# array or comma-separated values
tasks = ["ovs*"] # just support Wildcards

[CachePool]
shrink = false # whether allow to shrink cache ways in best effort pool
max_allowed_shared = 10 # max allowed workload in shared pool, default is 10
guarantee = 10
besteffort = 7
shared = 2
```

The configuration file directives are explained below:

- **OSGroup**: cache ways reserved for use by the operating system usage.
- **InfraGroup**: cache ways dedicated for use by identified infrastructure process. These cache ways are allocated to be shared with every other workload scheduled on the system.
- **CachePool**: Cache allocation requests are grouped into one of three groups:
 - **guarantee**: allocate dedicated cache ways to a workload (`max_cache == min_cache > 0`)
 - **besteffort**: dedicated minimum number of cache ways for workload which can grow to maximum number of cache ways in an overlapped manner (`max_cache > min_cache > 0`)

User Guide | Resource Management Daemon

- **shared**: all workloads in this pool share cache ways (`max_cache == min_cache = 0`)

Flags are:

- **shrink**: flag to indicate whether to shrink cache ways already allocated to workloads in **besteffort** pool upon arrival of a new workload.
- **max_allowed_shared**: threshold to define the max allowed workloads in shared cache pool

Refer to [Figure 4](#) and [Section 2.1](#) for additional information.

On a host that supports up to 20 cache ways, the configuration file will create the following cache bit mask layout.

```
OSGroup: 0000 0000 0000 0000 0001
InfraGroup: 1111 1111 1111 1111 1110
```

The layout of the remaining cache ways available for cache pools represented by the cache bit mask is:

```
guarantee: 0000 0000 0111 1111 1110
besteffort: 0011 1111 1000 0000 0000
shared:    1100 0000 0000 0000 0000
```

5.2 Querying Cache Inventory Information

The LLC information on the node deploying RMD can be queried using the following REST API calls:

```
$ curl -i http://127.0.0.1:8081/v1/cache/
$ curl -i http://127.0.0.1:8081/v1/cache/l3
$ curl -i http://127.0.0.1:8081/v1/cache/l3/0
```

5.3 Workload Policies to Request LLC Allocation

5.3.1 Changing and Adding Policies

The sys-admin can change and add new policies by editing the following line in the yaml file which is pointed to in the configuration file.

```
policy_path = "etc/rmd/policy.toml"
```

5.3.2 Pre-defined Policy

The sys-admin can provide predefined LLC allocation policies based on the target platform characteristics. The predefined allocation policy represents cache pools of Guaranteed, Best-effort and Shared using Gold, Silver and Bronze levels respectively.

The following example is a predefined policy which can be referenced in workload policy to RMD to request LLC.

Note: The following `catpolicy` file includes product codenames for processors which have been released by Intel.

For desktop and mobile, Broadwell is branded as 5th Generation Intel® Core™ processors. For server class processors, Broadwell is branded as Intel® Xeon® E3 v4, Intel® Xeon® E5 v4, and Intel® Xeon® E7 v4 processors.

For desktop and mobile, Skylake is branded as 6th Generation Intel® Core™ i3, Intel® Core™ i5, and Intel® Core™ i7 processors. For workstations, Skylake is branded as Intel® Xeon® E3 v5 processors.

```
catpolicy:
  broadwell:
    - gold:
      - MaxCache: 4
      - MinCache: 4
    - silver-bf:
      - MaxCache: 2
      - MinCache: 1
    - bronze-shared:
      - MaxCache: 0
      - MinCache: 0
  skylake:
    - gold:
      - MaxCache: 3
      - MinCache: 3
    - silver-bf:
      - MaxCache: 2
      - MinCache: 1
    - bronze-shared:
      - MaxCache: 0
      - MinCache: 0
```


Example Workload Request Command

The following command creates a workload request to RMD with gold policy for a workload described with process id 78377:

```
$ curl -H "Content-Type: application/json" --request POST --data \
  '{"task_ids":["78377"], "policy": "gold"}' \
  http://127.0.0.1:8081/v1/workloads
```

5.3.3 Custom Policy

To specify a custom policy to request LLC allocation the payload body should contain the following. A workload can be described as a set of task_ids or set of core_ids

- **task_ids**: set of process ids describing the workload
- **core_ids**: set of core ids the workload process ids run on
- **max_cache**: maximum number of cache ways the workload can run on
- **min_cache**: minimum number of cache ways the workload must run on

Example Workload Request Command

The following command creates a workload request to RMD with minimum and maximum cache values for a workload described with process id 78377:

```
$ curl -H "Content-Type: application/json" --request POST --data \
  '{"task_ids":["78377"], "max_cache": 4, "min_cache": 4}' \
  http://127.0.0.1:8081/v1/workloads
```

5.3.4 Hospitality Score of a Workload

The hospitality score API provides a score in the range of 0 to 100 to indicate how hospitable the host is to provide the desired LLC allocation values for the workload.

The score is calculated as follows:

Table 3. Hospitality Score Calculation

Request	Hospitality Score	LLC Pool
max_cache == min_cache > 0	[0 100]	Guarantee
max_cache == min_cache == 0	[0 100]	Shared
max_cache > min_cache > 0	[0, 100]	Best-effort

Example: The following code will display the hospitality score for a workload that has desired LLC allocation values of min=2 and max=2. Since two LLC ways can be provided on the host, the score is returned as 100.

```
$ curl -H "Content-Type: application/json" --request POST --data \
  '{"max_cache": 2, "min_cache": 2}' \
  http://127.0.0.1:8888/v1/hospitality
{
  "score": {
    "13": {
      "0": 100,
      "1": 100
    }
  }
}
```

6 Plugins

RMD is a centralized arbitrator to perform resource management on the platform. For this it needs access to system resources for enforcement as well as monitoring. For ease of use for platform and vendor developers, RMD defines standard interfaces for each resource it can control. Platform and vendor developers can write plugins that can conform to these provided standard interfaces. We can expect multiple versions for a single plugin for a resource as well as multiple plugins for a resource. Plugins can be combined to handle both write and read but must be registered separately as write and read plugins with RMD for better management.

User Guide | Resource Management Daemon

Using the plugin approach:

- RMD administrators can control the list of resources supported by RMD for enforcement and monitoring through configuration. It is the responsibility of administrators to ensure the plugins configured are compatible with RMD and the platform. For example, compatibility might be required for Golang, the kernel, hardware versions like CPU model, and so forth.
- Platform and vendor developers should provide template configuration and documentation to explain the features provided by the plugin as well as any dependencies required.
- For a resource, the configuration should allow only one plugin each for write and read operations.

Different resources on the system will have different needs and the configurability and control will be highly coupled to the resource controller implementation. RMD plugins reduce the dependency on platform or vendor specific implementation details. It uses its configuration to load the correct plugin based on platform or vendor resource version, if provided.

Usually the RMD plugins are pieces of code that read or write into standard kernel facilities such as the pseudo filesystems (`proc/sysfs/resctrlfs` and others) but they can be implemented by accessing resources in different ways: through communicating with other daemons on the platform, calling libraries, and writing directly to Model Specific Registers (MSRs).

To support a newly introduced resource, you need to do a bit of research on how this resource is exposed. Starting from the kernel pseudo filesystem is usually a good approach, but there might be other ways that a resource can be accessed, such as a command line tool (such as `sysctl` or `ethtool`) or via another daemon or library.

It's important to experiment to get an understanding of how this resource affects workloads. As an example, refer to this investigation for Cache Allocation Technology [here](#). Then you can proceed with writing a plugin that monitors and enforces the resource in a meaningful way.

7 Summary

Resource Management Daemon (RMD) enables a flexible and scalable implementation for resource management. It provides a unified centralized decision making mechanism that can be used in multiple environments.

This document described RMD architecture and installation, provided a usage scenario, discussed plugins, and provided an example of using RMD for the Intel® Resource Director Technology (Intel® RDT).



You may not use or facilitate the use of this document in connection with any infringement or other legal analysis concerning Intel products described herein. You agree to grant Intel a non-exclusive, royalty-free license to any patent claim thereafter drafted which includes subject matter disclosed herein.

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.

Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.

This document contains information on products, services and/or processes in development. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

The products and services described may contain defects or errors known as errata which may cause deviations from published specifications. Current characterized errata are available on request. No product or component can be absolutely secure.

Copies of documents which have an order number and are referenced in this document may be obtained by calling 1-800-548-4725 or by visiting www.intel.com/design/literature.htm.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Learn more at <http://www.intel.com/> or from the OEM or retailer.

Intel, the Intel logo, Intel Core, and Xeon are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.

* Other names and brands may be claimed as the property of others.