

# Vacuum Platform

A. McNab<sup>1</sup>

<sup>1</sup>*University of Manchester*

## Abstract

This technical note describes components of GridPP's Vacuum Platform for managing VMs, including the `$JOBOUTPUTS`, `VacQuery`, and `VacUserData` interfaces.

# 1 Introduction

This technical note describes components of GridPP’s Vacuum Platform for managing virtual machines (VMs) to run jobs for WLCG and other HEP experiments.

The `$JOBOUTPUTS`, `VacQuery`, and `VacUserData` interfaces are described, which have been developed for managing the VM lifecycle. These are used by two GridPP software systems, `Vac` and `Vcycle`, which can be described as VM lifecycle managers (VMLMs).

This note is written in terms of VMs, but the interfaces have been designed to be generalised to other forms of logical machine in the future, such as Docker containers and unikernels.

The term “resource provider” is used to refer to the entity which is able to take the decision about creating each VM. That is, the decision about whether resources will be provided or not. Typically, this is owner of an OpenStack tenancy managed by `Vcycle` or the manager of `Vac` VM factories. Beyond this may be many layers of resource provision in terms of legal owners of services and hardware and even funding agencies.

A location at which VMs can be created managed by one or more VMLMs which are cooperating to achieve a set of target shares is referred to as a “space”. This is equivalent to a Compute Element at a Grid site, and spaces must be given DNS names in DNS space available to the resource provider. However, it is not necessary to register the space name in the corresponding DNS zone. For `Vac`, a space is a set of VM factory machines which are communicating via `VacQuery` and may be said to be neighbours. For `Vcycle`, a space corresponds to an OpenStack tenancy or project, with a specified endpoint to contact and identity tokens to use.

## 2 Environment

Where possible, VMLMs should provide an approximation of OpenStack’s environment for VMs, which is derived from EC2. Any contextualization `user_data` file required and a metadata service should be provided via a “Magic IP” HTTP service at 169.254.169.254 from the point of view of the VM. Monolithic VM images which do not use a `user_data` file require a metadata service to be able to discover the URLs of the `$MACHINEFEATURES`, `$JOBFEATURES`, and `$JOBOUTPUTS` locations.

As not all IaaS cloud systems provide metadata services, VMs and VMLMs should also implement the standard `VacUserData` substitutions described in section 6. These include substitutions giving the URLs of the `$MACHINEFEATURES`, `$JOBFEATURES`, and `$JOBOUTPUTS` locations.

VMLMs must also support VMs which use Cloud Init contextualization.

## 3 Machine/Job Features

The Machine/Job Features (MJF) mechanism described in [?] allows resource providers to communicate information to batch jobs and virtual machines, including the number of

38 processors they are allocated and how long they may run for. The MJF terminology is  
39 derived from batch job environments, and job equates to virtual machine when applied to  
40 virtualized environments such as the Vacuum Platform.

41 Resource providers in the Vacuum Platform must make the MJF `$MACHINEFEATURES`  
42 and `$JOBFEATURES` locations available over HTTP(S) rather than through a shared filesys-  
43 tem, and should publish the URLs of these locations in OpenStack/EC2 metadata tags  
44 and using the `VacUserData` substitutions in the `user_data` files supplied to VMs.

45 The value of `$JOBFEATURES/job_id` should be set to the VM UUID by the VMLM as  
46 soon as it is known. For example, with Vac the UUID is chosen by VMLM and its value  
47 can be set when the first `$JOBFEATURES` key/values are created. However with Vcycle  
48 managing OpenStack, the VM UUID is only available after the VM has been created, and  
49 is then recorded by Vcycle in the `$JOBFEATURES` directory it provides.

## 50 4 `$JOBOUTPUTS`

51 The `$JOBOUTPUTS` mechanism is an extension to Machine/Job Features by which the URL  
52 of a location to which VMs can write status and log files is communicated to the VMs. This  
53 value of the `$JOBOUTPUTS` URL should be given in the same way as the `$MACHINEFEATURES`  
54 and `$JOBFEATURES` URLs.

55 Any log file which the VMs wish to make available to resource providers may be written  
56 to the `$JOBOUTPUTS` location, for later examination in case of problems. All of these files  
57 must have unique names, and they are all written to the same level (“directory”) of the  
58 URL space on the `$JOBOUTPUTS` HTTP(S) server. This mechanism is also used to provide  
59 the `shutdown_message` file described in the next section.

### 60 4.1 Shutdown Messages

61 When VMs finish, they should write a `shutdown_message` file to  
62 `$JOBOUTPUTS/shutdown_message` containing one line of text without a trailing  
63 newline character. This text consists of a three digit shutdown message code in the range  
64 100-999, a space, and then a human-readable description of the message code.

65 The message code (and not the human-readable description) will be used by the  
66 resource provider’s software to determine why the VM finished and whether to create  
67 more VMs of this type in the immediate future as slots become available.

68 The shutdown codes are designed to be extensible by the insertion of intermediate  
69 numbers for finer-grained reporting. This is similar to the three digit response codes of  
70 internet protocols such as SMTP and HTTP.

100	Shutdown as requested by the VM's host/hypervisor
200	Intended work completed ok
300	No more work available from task queue
400	Site/host/VM is currently banned/disabled from receiving more work
500	Problem detected with environment/VM provided by the site
600	Grid-wide problem with job agent or application within VM
700	Transient problem with job agent or application within VM

Table 1: Shutdown codes and messages

## 5 Image URLs

Experiments should provide the HTTPS URL of the image file required to boot their VMs, which VMLMs should use. VMLMs should support both standard CAs and IGTF endorsed-CAs when verifying the X.509 certificates used by the relevant HTTPS webserver.

To avoid overloading these webserver, VMLMs must cache images by Last-Modified time, and should use the HTTP If-Modified-Since mechanism when fetching images. If this header is used, then it is acceptable to check the URL for updates each time a VM is created.

Where the VMLM is unable to update the image used to boot the VMs itself, it should attempt to verify that the image being used is current and refuse to create new VMs with an old image. Typically this applies to IaaS cloud systems where users are unable to upload new images, or a manual upload step is required. VMLM authors should consider how resource providers will be made aware of this situation when it arises, but for scalability reasons, the VMLM should not rely on the experiment suffering from VMs starting and then failing due to an out of date VM image and then notifying resource providers.

## 6 VacUserData templates

In most cases, a generic image such as CernVM is used which then requires further contextualization as the VM starts using a user\_data file supplied by VMLM. The VMLM must be able to retrieve a template for the user\_data file from an HTTPS URL nominated by the experiment each time a VM is to be created. That is, without any caching. The VMLM must include an appropriate HTTP User-Agent header indicating the VMLM implementation and version when making this request to allow experiments to monitor which VMLMs are in use. The VMLM should support both standard CAs and IGTF endorsed-CAs when verifying the X.509 certificates used by the relevant HTTPS webserver.

The VMLM must apply the following pattern based substitutions to the user\_data template supplied by the experiment. These patterns are all in the form `##user_data.XXX##`.

The following substitutions are performed automatically using data the VMLM holds

98 internally:

```
99     ##user_data_space## Space name
100     ##user_data_machinetype## Name of the machinetype of this VM
101     ##user_data_machine_hostname## Hostname assigned to the VM by the
102         VMLM
103     ##user_data_manager_version## A string giving the VMLM version
104     ##user_data_manager_hostname## Hostname of the VMLM
105     ##user_data_manager_machinefeatures_url## $MACHINEFEATURES URL
106         (section 3)
107     ##user_data_manager_jobfeatures_url## $JOBFEATURES URL (section 3)
108     ##user_data_manager_joboutputs_url## $JOBOUTPUTS URL (section 4)
```

109 The VMLM must also provide a mechanism for the resource provider which uses the  
110 VMLM to specify strings or files which will be used in substitutions required by the VM.  
111 These patterns take the form `##user_data_option.XXX##` where XXX is an arbitrary  
112 string consisting of letters, numbers, and underscores.

113 If the VM requires an X.509 proxy, it must expect that the special pattern  
114 `##user_data_option_x509_proxy##` will be replaced by the PEM encoded X.509 cer-  
115 tificates and RSA private key which comprise the proxy. VMLM's should provide a  
116 mechanism for creating X.509 proxies dynamically for each VM instance from a host or  
117 robot certificate owned by the resource provider, with an X.509 proxy lifetime reflecting  
118 the maximum VM lifetime.

## 119 7 VacQuery

120 The principal use of the VacQuery protocol is to allow Vac factories to gather information  
121 from their neighbours about what VMs are running for what machinetypes. This is done using  
122 the `machinetypes_query` and `machinetype_status` UDP messages. Factory and machine  
123 message pairs are also supported which can be used for automated or manual monitoring  
124 of Vac-based sites.

125 VacQuery queries sent to Vac daemons take the form of JSON documents in packets  
126 directed to the unused UDP port 995<sup>1</sup>. The protocol has been designed to keep JSON  
127 message and IP header below the ethernet MTU of 1500 bytes to avoid fragmentation on  
128 local networks. Responses are sent to the UDP port from which the query was sent.

129 All dates/times in VacQuery messages are expressed as Unix seconds. That is, the  
130 number of seconds since 00:00:00 1st Jan 1970.

---

<sup>1</sup>In Roman numerals, V=5 and M=1000. 995 could be written as VM = 1000 - 5, although this violates conventions invented in modern times.

## 131 7.1 Factory messages

132 The factory messages `factory_query` and `factory_status` are intended for monitoring the  
133 state of the factories themselves, including generic Unix health metrics such as free disk  
134 and CPU load. As well as manual queries by administrators, these messages may also be  
135 used for automated Nagios-style monitoring and alarms.

### 136 7.1.1 `factory_query`

137 The `factory_query` message is sent to a factory to request a `factory_status` message in  
138 response.

139     **message\_type** “factory\_query”  
140     **vac\_version** Software version of Vac  
141     **vacquery\_version** Version of the VacQuery protocol  
142     **space** Vac space name  
143     **cookie** Freely chosen by the sender

### 144 7.1.2 `factory_status`

145 `factory_status` messages are returned in response to `factory_query` messages directed to  
146 a factory. They may also be generated spontaneously and sent to a VacMon service as  
147 described in section 7.4.

148     **message\_type** “factory\_status”  
149     **vac\_version** Software version of Vac  
150     **vacquery\_version** Version of the VacQuery protocol  
151     **cookie** Matching the value supplied by the recipient  
152     **space** Vac space name  
153     **factory** FQDN of the factory  
154     **time\_sent** Time in Unix seconds  
155     **total\_cpus** Number of (logical) processors available to VMs  
156     **running\_cpus** Number of processors assigned to running VMs  
157     **running\_machines** Number of running VMs  
158     **total\_machines** Maximum possible number of VMs  
159     **total\_hs06** Total HS06 available to VMs  
160     **vac\_disk\_avail\_kb** Free space available in Vac’s workspace, in units of 1024  
161         bytes  
162     **root\_disk\_avail\_kb** Free space available on the root partition, in units of  
163         1024 bytes  
164     **vac\_disk\_avail\_inodes** Free inodes available in Vac’s workspace

165 **root\_disk\_avail\_inodes** Free inodes available on the root partition  
 166 **load\_average** The 15 minute Unix load average on the factory  
 167 **kernel\_version** The running kernel version of the factory  
 168 **os\_issue** A string identifying the operating system (typically the first line of  
 169 /etc/issue)  
 170 **boot\_time** The time when the factory booted up in Unix seconds  
 171 **factory\_heartbeat\_time** Time of the last heartbeat created by the VM  
 172 factory agent in Unix seconds  
 173 **responder\_heartbeat\_time** Time of the last heartbeat created by the Vac-  
 174 Query responder service in Unix seconds  
 175 **mjf\_heartbeat\_time** Time of the last heartbeat created by the HTTP Ma-  
 176 chine/Job Features service in Unix seconds  
 177 **metadata\_heartbeat\_time** Time of the last heartbeat created by the HTTP  
 178 Metadata service in Unix seconds  
 179 **swap\_used\_kb** Swap space in use on the factory, in units of 1024 bytes  
 180 **swap\_free\_kb** Free swap space, in units of 1024 bytes  
 181 **mem\_used\_kb** Physical memory in use on the factory, in units of 1024 bytes  
 182 **mem\_total\_kb** Free physical memory, in units of 1024 bytes  
 183 **site** Name of the site registered in the GOCDB, or the Vac space name if the  
 184 site is not registered

## 185 7.2 Machine messages

186 The machines\_query (plural) and machine\_status (singular) messages can be used to create  
 187 views of the VMs running within a Vac space, similar to the views from the top command  
 188 of running processes on a single host.

### 189 7.2.1 machines\_query

190 The machines\_query message is sent to a factory to request a machine\_status message for  
 191 each of its VM slots.

192 **message\_type** “machines\_query”  
 193 **vac\_version** Software version of Vac  
 194 **vacquery\_version** Version of the VacQuery protocol  
 195 **space** Vac space name  
 196 **cookie** Freely chosen by the sender

### 197 7.2.2 machine\_status

198 machine\_status messages are returned in response to machines\_query messages directed to  
 199 a factory.

200     **message\_type** “factory\_status”  
201     **vac\_version** Software version of Vac  
202     **vacquery\_version** Version of the VacQuery protocol  
203     **cookie** Matching the value supplied by the recipient  
204     **space** Vac space name  
205     **factory** FQDN of the factory  
206     **num\_machines** Number of machine\_status messages to expect from this  
207         factory  
208     **time\_sent** Time in Unix seconds  
209     **machine** Hostname of the VM slot  
210     **state** State of the current or most recent VM in this slot, as a string  
211     **uuid** Lowlevel UUID, as used by libvirt  
212     **created\_time** Unix time of the VM’s creation  
213     **started\_time** Unix time the VM entered the running state  
214     **heartbeat\_time** Unix time when the VM was last observed to be running  
215         (this is not the same as any heartbeat generated within the VM)  
216     **cpu\_seconds** CPU seconds used by the VM  
217     **cpu\_percentage** Recent CPU percentage use. May be over 100% for multit-  
218         processor VMs  
219     **hs06** Total HEPSPEC06 for the number of processors assigned to this VM  
220     **machinetype** Name of the machinetype  
221     **shutdown\_message** Any shutdown message given by the last VM to run in  
222         this slot  
223     **shutdown\_time** Unix time of the shutdown\_message

## 224   **7.3   Machinetype messages**

225   The machinetypes\_query (plural) and machinetype\_status (singular) messages are used by  
226   factories to gather information from neighbours within the same Vac space about what  
227   they are running, and outcomes of recently started VMs which have finished.

### 228   **7.3.1   machinetypes\_query**

229   The machinetypes\_query message is sent to a factory to request a machinetype\_status  
230   message for each of the machinetypes it supports.

231     **message\_type** “machinetypes\_query”  
232     **vac\_version** Software version of Vac  
233     **vacquery\_version** Version of the VacQuery protocol  
234     **space** Vac space name  
235     **cookie** Freely chosen by the sender



### 236 7.3.2 machinetype\_status

237 machinetype\_status messages are returned in response to machinetypes\_query messages  
238 directed to a factory. They may also be generated spontaneously and sent to a VacMon  
239 service as described in section 7.4.

240 **message\_type** “factory\_status”  
241 **vac\_version** Software version of Vac  
242 **vacquery\_version** Version of the VacQuery protocol  
243 **cookie** Matching the value supplied by the recipient  
244 **space** Vac space name  
245 **factory** FQDN of the factory  
246 **num\_machinetypes** Number of machinetype\_status messages to expect from  
247 this factory  
248 **time\_sent** Time in Unix seconds  
249 **machinetype** Name of the machinetype  
250 **total\_hs06** Total HEPSPEC06 of all the processors allocated to running VMs  
251 for this machinetype on this factory  
252 **total\_machines** Number of running VMs for this machinetype on this factory  
253 **num\_before\_fizzle** Number of running VMs which have not yet reached  
254 fizzle\_seconds  
255 **shutdown\_message** Shutdown message given by the most recently created  
256 VM for this machinetype on this factory which has finished  
257 **shutdown\_time** Unix time of the shutdown\_message  
258 **shutdown\_machine** Name of the VM slot associated with the shut-  
259 down\_message

## 260 7.4 VacMon services

261 VacMon services receive factory\_status and machinetype\_status messages from Vac daemons  
262 on UDP port 8884. These may be used for Ganglia-style monitoring of individual sites or  
263 groups of sites. As VacQuery messages as JSON documents, they may be conveniently  
264 recorded in data stores such as Elasticsearch.

## 265 8 APEL

266 VMLMs should support reporting of usage to the central APEL service with messages of  
267 the type “APEL-individual-job-message”.

268 In this case, VLMs must include the following in the messages:

269 **FQAN** VOMS FQAN specified by the experiment when configuring the space

270       **SubmitHost** Must be of the form [space name] + "/" + [vmlm] + "-" +  
 271               [VMLM host name], where "vmlm" is a lowercase name for the VMLM  
 272               software such as "vac"  
 273       **LocalJobId** VM UUID  
 274       **LocalUserId** VMLM hostname  
 275       **Queue** Name of the machinetype  
 276       **GlobalUserName** The space name converted to an X.500 DN with  
 277               DC components. For example, vac01.example.com would become  
 278               /DC=vac01/DC=example/DC=com  
 279       **InfrastructureDescription** Such as APEL-VAC or APEL-VCYCLE, with  
 280               APEL and then an uppercase name for the VMLM software.  
 281       **Processors** The number of virtual CPUs assigned to the VM.

282       APEL Sync records must also be sent, and these can conveniently be generated by  
 283       each VMLM instance from an archive of the individual job messages, as the SubmitHost  
 284       is unique to the VMLM instance in both cases.

285       In addition to APEL records generated by the VMLM, an underlying cloud infrastrac-  
 286       ture may also be instrumented to submit APEL usage records to the central APEL service.  
 287       This is especially likely at sites participating in the EGI Federated Cloud. In this case,  
 288       resource providers must ensure that double counting is avoided by disabling reporting  
 289       from the VMLM to the central APEL service.

## 290   9   GLUE2

291       VMLMs should publish status information as JSON documents using the GLUE 2.0  
 292       schema, at an HTTPS URL advertised in the GOCDB service entry for the space, in the  
 293       Grid Information / URL field.

## 294 **10 Summary**

295 Vacuum platform interfaces have been described ...

## 296 **References**