

Vacuum Platform

A. McNab¹

¹*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

71 5 Image URLs

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

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

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85 then failing due to an out of date VM image and then notifying resource providers.

86 6 VacUserData templates

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

95 The VMLM must apply the following pattern based substitutions to the user_data tem-
96 plate supplied by the experiment. These patterns are all in the form **##user_data.XXX##**.

97 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 compromise 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 principle 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 995¹. 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 7.1 Factory messages

130 The factory messages `factory_query` and `factory_status` are intended for monitoring the
131 state of the factories themselves, including generic Unix health metrics such as free disk

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

132 and CPU load. As well as manual queries by administrators, these messages may also be
133 used for automated Nagios-style monitoring and alarms.

134 7.1.1 factory_query

135 The factory_query message is sent to a factory to request a factory_status message in
136 response.

137 **message_type** “factory_query”
138 **vac_version** Software version of Vac
139 **vacquery_version** Version of the VacQuery protocol
140 **space** Vac space name
141 **cookie** Freely chosen by the sender

142 7.1.2 factory_status

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

146 **message_type** “factory_status”
147 **vac_version** Software version of Vac
148 **vacquery_version** Version of the VacQuery protocol
149 **cookie** Matching the value supplied by the recipient
150 **space** Vac space name
151 **factory** FQDN of the factory
152 **time_sent** Time in Unix seconds (since 00:00:00 1st Jan 1970)
153 **total_cpus** Number of (logical) processors available to VMs
154 **running_cpus** Number of processors assigned to running VMs
155 **running_machines** Number of running VMs
156 **total_machines** Maximum possible number of VMs
157 **total_hs06** Total HS06 available to VMs
158 **vac_disk_avail_kb** Free space available in Vac’s workspace, in units of 1024
159 bytes
160 **root_disk_avail_kb** Free space available on the root partition, in units of
161 1024 bytes
162 **vac_disk_avail_inodes** Free inodes available in Vac’s workspace
163 **root_disk_avail_inodes** Free inodes available on the root partition
164 **load_average** The 15 minute Unix load average on the factory
165 **kernel_version** The running kernel version of the factory

166 **os_issue** A string identifying the operating system (typically the first line of
 167 /etc/issue)
 168 **boot_time** The time when the factory booted up in Unix seconds
 169 **factory_heartbeat_time** Time of the last heartbeat created by the VM
 170 factory agent in Unix seconds
 171 **responder_heartbeat_time** Time of the last heartbeat created by the Vac-
 172 Query responder service in Unix seconds
 173 **mjf_heartbeat_time** Time of the last heartbeat created by the HTTP Ma-
 174 chine/Job Features service in Unix seconds
 175 **metadata_heartbeat_time** Time of the last heartbeat created by the HTTP
 176 Metadata service in Unix seconds
 177 **swap_used_kb** Swap space in use on the factory, in units of 1024 bytes
 178 **swap_free_kb** Free swap space, in units of 1024 bytes
 179 **mem_used_kb** Physical memory in use on the factory, in units of 1024 bytes
 180 **mem_total_kb** Free physical memory, in units of 1024 bytes
 181 **site** Name of the site registered in the GOCDB, or the Vac space name if the
 182 site is not registered

183 **7.2 Machine messages**

184 The `machines_query` (plural) and `machine_status` (singular) messages can be used to create
 185 views of the VMs running within a Vac space, similar to the views from the `top` command
 186 of running processes on a single host.

187 **7.2.1 machines_query**

188 The `machines_query` message is sent to a factory to request a `machine_status` message for
 189 each of its VM slots.

190 **message_type** “machines_query”
 191 **vac_version** Software version of Vac
 192 **vacquery_version** Version of the VacQuery protocol
 193 **space** Vac space name
 194 **cookie** Freely chosen by the sender

195 **7.2.2 machine_status**

196 `machine_status` messages are returned in response to `machines_query` messages directed to
 197 a factory.

198 **message_type** “factory_status”
 199 **vac_version** Software version of Vac

vacquery_version Version of the VacQuery protocol
cookie Matching the value supplied by the recipient
space Vac space name
factory FQDN of the factory
num_machines Number of machine_status messages to expect from this factory
time_sent Time in Unix seconds (since 00:00:00 1st Jan 1970)
machine Hostname of the VM slot
state State of the current or most recent VM in this slot, as a string
uuid Lowlevel UUID, as used by libvirt
created_time Unix time of the VM's creation
started_time Unix time the VM entered the running state
heartbeat_time Unix time when the VM was last observed to be running (this is not the same as any heartbeat generated within the VM)
cpu_seconds CPU seconds used by the VM
cpu_percentage Recent CPU percentage use. May be over 100% for multi-processor VMs
hs06 Total HEPSPEC06 for the number of processors assigned to this VM
machinetype Name of the machinetype
shutdown_message Any shutdown message given by the last VM to run in this slot
shutdown_time Unix time of the shutdown_message

222 7.3 Machinetype messages

223 The machinetypes_query (plural) and machinetype_status (singular) messages are used by
 224 factories to gather information from neighbours within the same Vac space about what
 225 they are running, and outcomes of recently started VMs which have finished.

226 7.3.1 machinetypes_query

227 The machinetypes_query message is sent to a factory to request a machinetype_status
 228 message for each of the machinetypes it supports.

229 **message_type** "machinetypes_query"
 230 **vac_version** Software version of Vac
 231 **vacquery_version** Version of the VacQuery protocol
 232 **space** Vac space name
 233 **cookie** Freely chosen by the sender

234 7.3.2 machinetype_status

235 machinetype_status messages are returned in response to machinetypes_query messages
236 directed to a factory. They may also be generated spontaneously and sent to a VacMon
237 service as described in section 7.4.

238 **message_type** “factory_status”
239 **vac_version** Software version of Vac
240 **vacquery_version** Version of the VacQuery protocol
241 **cookie** Matching the value supplied by the recipient
242 **space** Vac space name
243 **factory** FQDN of the factory
244 **num_machinetypes** Number of machinetype_status messages to expect from
245 this factory
246 **time_sent** Time in Unix seconds (since 00:00:00 1st Jan 1970)
247 **machinetype** Name of the machinetype
248 **total_hs06** Total HEPSPEC06 of all the processors allocated to running VMs
249 for this machinetype on this factory
250 **total_machines** Number of running VMs for this machinetype on this factory
251 **num_before_fizzle** Number of running VMs which have not yet reached
252 fizzle_seconds
253 **shutdown_message** Shutdown message given by the most recently created
254 VM for this machinetype on this factory which has finished
255 **shutdown_time** Unix time of the shutdown_message
256 **shutdown_machine** Name of the VM slot associated with the shut-
257 down_message

258 7.4 VacMon services

259 VacMon services receive factory_status and machinetype_status messages from Vac daemons
260 on UDP port 8884. These may be used for Ganglia-style monitoring of individual sites or
261 groups of sites. As VacQuery messages as JSON documents, they may be conveniently
262 recorded in data stores such as ElasticSearch.

263 8 APEL

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

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

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

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

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

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

288 9 GLUE2

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

292 **10 Summary**

293 Vacuum platform interfaces have been described ...

294 **References**