

Vacuum Platform

A. McNab¹

¹*University of Manchester*

Abstract

This technical note describes components of the Vacuum Platform developed by GridPP 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 or other cloud 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 `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 [1] 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 using 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 `machinefeatures`
44 and `jobfeatures` metadata tags and using the `VacUserData` substitutions in the `user_data`
45 files supplied to VMs.

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

51 4 `$JOBOUTPUTS`

52 The `$JOBOUTPUTS` mechanism is an extension to Machine/Job Features by which the URL
53 of a location to which VMs can write status and log files is communicated to the VMs. This
54 value of the `$JOBOUTPUTS` URL should be given in the same way as the `$MACHINEFEATURES`
55 and `$JOBFEATURES` URLs, using a `VacUserData` substitution and an OpenStack/EC2
56 `joboutputs` metadata key.

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

62 4.1 Shutdown Messages

63 When VMs finish, they should write a `shutdown_message` file to
64 `$JOBOUTPUTS/shutdown_message` containing one line of text without a trailing
65 newline character. This text consists of a three digit shutdown message code in the range
66 100-999, a space, and then a human-readable description of the message code.

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

70 The shutdown codes are designed to be extensible by the insertion of intermediate
71 numbers for finer-grained reporting. This is similar to the three digit response codes of
72 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

73 5 Image URLs

74 Experiments should provide the HTTPS URL of the image file required to boot their VMs,
75 which VMLMs should use. VMLMs should support both standard CAs and International
76 Grid Trust Federation (IGTF) endorsed CAs when verifying the X.509 certificates used by
77 the relevant HTTPS webserver.

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

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

89 6 VacUserData templates

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

99 The VMLM must apply the following pattern based substitutions to the user_data tem-

100 plate supplied by the experiment. These patterns are all in the form `##user_data.XXX##`.

101 The following substitutions are performed automatically using data the VMLM holds
102 internally:

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

113 The VMLM must also provide a mechanism for the resource provider to specify strings
114 or files whose static values will be used in pattern substitutions required by the VM. These
115 patterns take the form `##user_data_option.XXX##` where XXX is an arbitrary string
116 consisting of letters, numbers, and underscores.

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

123 The VM must not assume that any other grid, HEP middleware, or scripts are running
124 as part of the VMLM and able to provide dynamic values for pattern substitutions. For
125 example, it must not require that resource providers provide proxies with VOMS attributes
126 to the VM. If this is needed, the VM should use the proxy provided to obtain the VOMS
127 credentials itself, using software managed by the experiment within the VM.

128 7 VacQuery

129 The VacQuery protocol specifies queries and status messages which can be sent over UDP
130 as short JSON documents.

131 The principal use of the VacQuery protocol is to allow Vac factories to gather informa-
132 tion from their neighbours about what VMs are running for what machinetypes. This is
133 done using the `machinetypes_query` and `machinetype_status` UDP messages. Factory and
134 machine message pairs are also supported which can be used for automated or manual
135 monitoring of Vac-based sites.

136 VacQuery queries sent to Vac daemons take the form of JSON documents in packets
137 directed to the unused UDP port 995.¹ Responses are sent to the UDP port from which the
138 query was sent. The protocol has been designed to keep JSON messages and IP headers
139 below the ethernet MTU of 1500 bytes to avoid fragmentation on local networks.

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

142 7.1 Factory messages

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

147 7.1.1 `factory_query`

148 The `factory_query` message is sent to a factory to request a `factory_status` message in
149 response.

150 **message_type** “factory_query”
151 **vac_version** Name and software version of Vac
152 **vacquery_version** Name and version of the VacQuery protocol
153 **space** Vac space name
154 **cookie** Freely chosen by the sender

155 7.1.2 `factory_status`

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

159 The format and units of the disk and memory values are aligned with the values
160 returned by the relevant system calls and the `/proc` interface.

161 **message_type** “factory_status”
162 **vac_version** Name and software version of Vac
163 **vacquery_version** Name and version of the VacQuery protocol
164 **cookie** Matching the value supplied by the recipient
165 **space** Vac space name
166 **factory** FQDN of the factory
167 **time_sent** Time in Unix seconds

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

site Name of the site registered in the GOCDB, or the Vac space name if the site is not registered
total_cpus Number of (logical) processors available to VMs
running_cpus Number of processors assigned to running VMs
running_machines Number of running VMs
total_machines Maximum possible number of VMs
total_hs06 Total HS06 available to VMs
boot_time The time when the factory booted up in Unix seconds
factory_heartbeat_time Time of the last heartbeat created by the VM factory agent in Unix seconds
responder_heartbeat_time Time of the last heartbeat created by the Vac-Query responder service in Unix seconds
mjf_heartbeat_time Time of the last heartbeat created by the HTTP Machine/Job Features service in Unix seconds
metadata_heartbeat_time Time of the last heartbeat created by the HTTP Metadata service in Unix seconds
vac_disk_avail_kb Free space available in Vac's workspace, in units of 1024 bytes
root_disk_avail_kb Free space available on the root partition, in units of 1024 bytes
vac_disk_avail_inodes Free inodes available in Vac's workspace
root_disk_avail_inodes Free inodes available on the root partition
load_average The 15 minute load average on the factory
kernel_version The kernel version of the factory
os_issue A string identifying the operating system (typically the first line of /etc/issue)
swap_used_kb Swap space in use on the factory, in units of 1024 bytes
swap_free_kb Free swap space, in units of 1024 bytes
mem_used_kb Physical memory in use on the factory, in units of 1024 bytes
mem_total_kb Free physical memory, in units of 1024 bytes

7.2 Machine messages

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

7.2.1 machines_query

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

205 **message_type** “machines_query”
 206 **vac_version** Name and software version of Vac
 207 **vacquery_version** Name and version of the VacQuery protocol
 208 **space** Vac space name
 209 **cookie** Freely chosen by the sender

210 **7.2.2 machine_status**

211 machine_status messages are returned in response to machines_query messages directed to
 212 a factory.

213 **message_type** “factory_status”
 214 **vac_version** Name and software version of Vac
 215 **vacquery_version** Name and version of the VacQuery protocol
 216 **cookie** Matching the value supplied by the recipient
 217 **space** Vac space name
 218 **factory** FQDN of the factory
 219 **num_machines** Number of machine_status messages to expect from this
 220 factory
 221 **time_sent** Time in Unix seconds
 222 **machine** Hostname of the VM slot
 223 **state** State of the current or most recent VM in this slot, as a string
 224 **uuid** Lowlevel UUID, as used by libvirtd
 225 **created_time** Unix time of the VM’s creation
 226 **started_time** Unix time the VM entered the running state
 227 **heartbeat_time** Unix time when the VM was last observed to be running
 228 (this is not the same as any heartbeat generated within the VM)
 229 **cpu_seconds** CPU seconds used by the VM
 230 **cpu_percentage** Recent CPU percentage use. May be over 100% for multit-
 231 processor VMs
 232 **hs06** Total HEPSPEC06 for the processors assigned to this VM
 233 **machinetype** Name of the machinetype
 234 **shutdown_message** Any shutdown message given by the last VM to run in
 235 this slot
 236 **shutdown_time** Unix time of the shutdown_message

237 **7.3 Machinetype messages**

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

241 7.3.1 machinetypes_query

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

244 **message_type** “machinetypes_query”
245 **vac_version** Name and software version of Vac
246 **vacquery_version** Name and version of the VacQuery protocol
247 **space** Vac space name
248 **cookie** Freely chosen by the sender

249 7.3.2 machinetype_status

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

253 **message_type** “factory_status”
254 **vac_version** Name and software version of Vac
255 **vacquery_version** Name and version of the VacQuery protocol
256 **cookie** Matching the value supplied by the recipient
257 **space** Vac space name
258 **factory** FQDN of the factory
259 **num_machinetypes** Number of machinetype_status messages to expect from
260 this factory
261 **time_sent** Time in Unix seconds
262 **machinetype** Name of the machinetype
263 **total_hs06** Total HEPSPEC06 of all the processors allocated to running VMs
264 for this machinetype on this factory
265 **total_machines** Number of running VMs for this machinetype on this factory
266 **num_before_fizzle** Number of running VMs which have not yet reached
267 fizzle_seconds
268 **shutdown_message** Shutdown message given by the most recently created
269 VM for this machinetype on this factory which has finished
270 **shutdown_time** Unix time of the shutdown_message
271 **shutdown_machine** Name of the VM slot associated with the shut-
272 down_message

273 7.4 VacMon services

274 VacMon services receive factory_status and machinetype_status messages from Vac daemons
275 on UDP port 8884. These may be used for Ganglia-style monitoring of individual sites

276 or groups of sites. As VacQuery messages are sent as JSON documents, they may be
277 conveniently recorded in data stores such as ElasticSearch.

278 8 APEL

279 VMLMs should support reporting of usage to the central APEL service with messages
280 of the type “APEL-individual-job-message”. These are the records used for conventional
281 grid sites, rather than those developed for cloud resources.

282 VLMs must include the following in the messages:

283 **FQAN** VOMS FQAN specified by the experiment when configuring the space
284 **SubmitHost** Must be of the form [space name] + "/" + [vmlm] + "-" +
285 [VLM host name], where “vmlm” is a lowercase name for the VMLM
286 software such as “vac”

287 **LocalJobId** VM UUID

288 **LocalUserId** VMLM hostname

289 **Queue** Name of the machinetype

290 **GlobalUserName** The space name converted to an X.500 DN with
291 DC components. For example, vac01.example.com would become
292 /DC=vac01/DC=example/DC=com

293 **InfrastructureDescription** Such as APEL-VAC or APEL-VCYCLE, with
294 APEL and then an uppercase name for the VMLM software.

295 **Processors** The number of virtual CPUs assigned to the VM.

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

299 In addition to grid-style APEL records generated by the VMLM, an underlying cloud
300 infrastructure may also be instrumented to submit cloud-style APEL usage records to the
301 central APEL service. This is especially likely at sites participating in the EGI Federated
302 Cloud. In this case, resource providers must ensure that double counting is avoided by
303 disabling reporting from the VMLM to the central APEL service.

304 9 GLUE2

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

308 **10 Summary**

309 Vacuum platform interfaces have been described ...

310 **References**

- 311 [1] M. Alef et al, HSF-TN-2016-02 “Machine/Job Features Specification” (HEP Software
312 Foundation)