

Synnefo - Feature # 505

Status:	Closed	Priority:	Medium
Author:	Giorgos Gousios	Category:	logic
Created:	05/10/2011	Assignee:	Giorgos Gousios
Updated:	06/03/2011	Due date:	
Subject:	Μηχανισμός επανυπολογισμού κατάστασης VM		
Description	<p>Κατά τη λειτουργία του συστήματος μπορεί να υπάρξουν στιγμές κατά τις οποίες η κατάσταση ενός VM στη βάση του συστήματος είναι μη συγχρονισμένη με την κατάσταση του ίδιου VM στο Ganeti. Ο μηχανισμός αυτός θα πρέπει να ελέγχει περιοδικά (?) ή όταν αντιμετωπιστεί κάποιο πρόβλημα και να συγχρονίζει τις 2 καταστάσεις.</p>		

History

#1 - 05/26/2011 12:10 pm - Giorgos Gousios

Κάτι που νομίζω θέλει συζήτηση και ίσως αλλαγή τώρα που έχουμε αλλάξει τον τρόπο που ενημερώνουμε το status του VM, είναι τα πεδία που χρειάζονται στη βάση και οι καταστάσεις που υποστηρίζουμε.

Το OpenStack 1.1 (Apr 25, 2011), υποστηρίζει τα παρακάτω.

ACTIVE, BUILD, REBUILD, SUSPENDED, QUEUE_RESIZE, PREP_RESIZE, RESIZE, VERIFY_RESIZE, PASSWORD, RESCUE, REBOOT, HARD_REBOOT, DELETE_IP, and UNKNOWN.

Στο Ganeti, το οποίο πλέον θεωρούμε ως το authoritative source για την κατάστασή μας, δεν έχω βρει κάπου ομαδοποιημένες τις καταστάσεις που μπορεί να υπάρξει ένα VM. Γενικά φαίνεται να υποστηρίζει τουλάχιστο τις παρακάτω τερματικές καταστάσεις, όπως φαίνεται από lib/query.py:

ERROR_nodeoffline, ERROR_nodedown, ERROR_wrongnode, running, ERROR_up, ERROR_down, ADMIN_down

Από όσο καταλαβαίνω από τις συζητήσεις μας και τον κώδικα που έχω δει, το ganeti στην ουσία δεν κρατάει κάπου το state του μηχανήματος και δεν έχει κάποιο state machine για το τι επιτρέπεται. Το κάνει infer από το αποτέλεσμα της τελευταίας δουλειάς που έχει τρέξει για το VM.

Το Σύννεφο προς το παρόν στην ουσία κάνει mirror τη συμπεριφορά του ganeti (6 πεδία που κρατάνε το job status για το τελευταίο job που ξεκίνησε ο χρήστης, όπου job στην ουσία μια ενέργεια που τροποποιεί την κατάσταση του vm στο ganeti) και 2 flags (suspended και deleted) που κρατάνε κάποιες καταστάσεις που ο συνδιασμός των 6 πεδίων δεν είναι ικανός να αναπαραστήσει (ή κάτι άλλο που δεν καταλαβαίνω). Επιπλέον, ο μηχανισμός που κάνει infer την κατάσταση από τα σύνολο 8 πεδία που έχουμε στη βάση μπορεί να δώσει στο API προς απάντηση είναι οι εξής καταστάσεις: UNKNOWN, REBOOT, BUILD, ERROR, STOPPED, ACTIVE, DELETED. Επιπλέον, το UI έχει κάποιες ενδιάμεσες μεταβατικές καταστάσεις (BUILDING, REBOOTING κτλ) που δεν αντικατοπτρίζονται πουθενά.

Νομίζω ότι το όλο σύστημα μπορεί να γίνει πολύ πιο απλό, αν κάνουμε τις 2 παρακάτω παραδοχές:

- Οι καταστάσεις που υποστηρίζει το σύστημα είναι μόνο αυτές που υποστηρίζει το API
- Το οποιοδήποτε backend είναι υπεύθυνο να μας ενημερώνει για την οποιαδήποτε καινούργια κατάσταση (πάλι από το set των καταστάσεων που υποστηρίζουμε μόνο).

Έτσι, στην ουσία, το Σύννεφο απλά καταγράφει την τελευταία γνωστή κατάσταση και ενημερώνει τον όποιο API client για αυτή. Είναι δουλειά του backend (ganeti, vmware, ?) να ενημερώνει το frontend σχετικά με την τρέχουσα κατάσταση γιατί εκεί υπάρχει γνώση των όποιων εσωτερικών καταστάσεων και των μεταβάσεών τους.

Καταλαβαίνω ότι το μέχρι τώρα design είχε βασιστεί στην ιδέα ότι το backend είναι αποκλειστικά ganeti, αλλά νομίζω ότι το σύστημα θα ήταν πολύ πιο απλό όπως το περιγράψω.

#2 - 05/26/2011 12:14 pm - Panagiotis Louridas

Giorgos Gousios wrote:

Κάτι που νομίζω θέλει συζήτηση και ίσως αλλαγή τώρα που έχουμε αλλάξει τον τρόπο που ενημερώνουμε το status του VM, είναι τα πεδία που χρειάζονται στη βάση
Νομίζω ότι το όλο σύστημα μπορεί να γίνει πολύ πιο απλό, αν κάνουμε τις 2 παρακάτω παραδοχές:
- Οι καταστάσεις που υποστηρίζει το σύστημα είναι μόνο αυτές που υποστηρίζει το API

Εννοείς το OpenStack API ή το RSAPI;

- Το οποιοδήποτε backend είναι υπεύθυνο να μας ενημερώνει για την οποιαδήποτε καινούργια κατάσταση (πάλι από το set των καταστάσεων που υποστηρίζουμε μόνο).

Καμμία διαφωνία εδώ, αφού μόνο αυτό ξέρει την πραγματική κατάσταση.

#3 - 05/26/2011 01:26 pm - Giorgos Gousios

- Οι καταστάσεις που υποστηρίζει το σύστημα είναι μόνο αυτές που υποστηρίζει το API

Εννοείς το OpenStack API ή το RSAPI;

Όποιο υποστηρίζουμε, το OpenStack αν δεν κάνω λάθος.

Να προσθέσω στα παραπάνω ότι στο σύστημα ως έχει, ο υπολογισμός της κατάστασης θα πρέπει να μετακινηθεί σε κάποιο κεντρικό σημείο στην πλευρά του backend. Επίσης, πρέπει να σχεδιάσουμε τα ελάχιστα μηνύματα που ενδιαφέρουν το σύστημά μας και να βάλουμε τα hooks στο ganeti να μεταφράζουν τα job status του σε αυτά που ακούμε εμείς. Προς το παρόν, αυτό είναι πολύ δουλειά και αλλάζει πολύ το σύστημα για να πάει στην v0.5. Προτείνω να αφήσουμε το ticket ανοιχτό και να το κοιτάξουμε σοβαρά μετά το release.

#4 - 05/26/2011 07:56 pm - Giorgos Gousios

- % Done changed from 0 to 60

Δυστυχώς, δεν γίνεται να προχωρήσουμε εύκολα με αυτό το feature καθώς το ganeti επιστρέφει λάθος JSON

- μονά εισαγωγικά, τα διπλά είναι δικιά μου προσθήκη για να μπορέσω να το περάσω από pretty printer
- μαργαριτάρια όπως το "dry_run":None, "static":False

φαντάζομαι ότι με κάποιο μη έγκυρο τρόπο έχουν κάνει κάποιο εσωτερικό data structure serialize. Αν υπάρχει κάποιος εύκολος τρόπος για να κάνω de-serialize, παρακαλώ ενημερώστε με.

```
1 {
2   "status":"success",
3   "ops":[
4     {
5       "dry_run":None,
6       "instances":[
7         "gousiosg-1"
8       ],
9       "priority":0,
```

```

10  "debug_level":0,
11  "OP_ID":"OP_INSTANCE_QUERY_DATA",
12  "static":False
13  }
14 ],
15 "end_ts":[
16  1306427618,
17  616642
18 ],
19 "start_ts":[
20  1306427618,
21  413629
22 ],
23 "summary":[-
24  "INSTANCE_QUERY_DATA"
25 ],
26 "received_ts":[
27  1306427618,
28  388367
29 ],
30 "opresult":[-
31  {
32    "gousiosg-1":{
33      "config_state":"up",
34      "network_port":12359,
35      "serial_no":2,
36      "be_instance":{
37        "auto_balance":True,
38        "vcpus":1,
39        "memory":1024
40      },
41      "be_actual":{
42        "auto_balance":True,
43        "vcpus":1,
44        "memory":1024
45      },
46      "name":"gousiosg-1",
47      "pnode":"store67",
48      "hypervisor":"kvm",
49      "disks":[-
50        {
51          "logical_id":[-
52            "ganeti",
53            "37d717a4-7f86-4dbd-9834-31630e1e8948.disk0"
54          ],
55          "sstatus":None,
56          "dev_type":"lvm",
57          "pstatus":[-
58            "/dev/ganeti/37d717a4-7f86-4dbd-9834-31630e1e8948.disk0",
59            254,
60            17,
61            None,
62            None,

```

```

63     False,
64     1
65 ],
66 "mode":"rw",
67 "physical_id":["
68     "ganes",
69     "37d717a4-7f86-4dbd-9834-31630e1e8948.disk0"
70 ],
71 "children":["
72
73 ],
74 "iv_name":"disk/0",
75 "size":2000
76 }
77 ],
78 "uuid":"10182022-2110-4240-b0be-9eb4ef5da114",
79 "hv_actual":{
80     "nic_type":"paravirtual",
81     "use_chroot":False,
82     "vnc_x509_path": "",
83     "vnc_bind_address":"127.0.0.1",
84     "usb_mouse": "",
85     "migration_downtime":30,
86     "security_model":"none",
87     "cdrom_image_path": "",
88     "boot_order":"disk",
89     "vhost_net":False,
90     "disk_cache":"default",
91     "kernel_path": "",
92     "acpi":True,
93     "vnc_x509_verify":False,
94     "vnc_tls":False,
95     "use_localtime":False,
96     "security_domain": "",
97     "serial_console":False,
98     "kvm_flag": "",
99     "vnc_password_file": "",
100    "disk_type":"paravirtual",
101    "kernel_args":"ro",
102    "root_path":"/dev/vda1",
103    "initrd_path": "",
104    "mem_path": ""
105 },
106 "hv_instance":{
107
108 },
109 "nics":[
110 [
111     None,
112     "aa:00:00:fc:99:26",
113     "bridged",
114     "br0"
115 ]

```

```
116     ],
117     "snodes":[
118
119     ],
120     "disk_template":"plain",
121     "mtime":1306417204.811044,
122     "run_state":"up",
123     "os_instance":{
124
125     },
126     "os":"debootstrap+default",
127     "os_actual":{
128
129     },
130     "ctime":1306417137.888135
131   }
132 }
133 ],
134 "opstatus":[
135   "success"
136 ],
137 "oplog":[
138   [
139
140   ]
141 ],
142 "id":"18160"
143}
144
```

#5 - 05/27/2011 10:40 am - Giorgos Gousios

Βοηθάει βέβαια το να μην αντιμετωπίζεις τα return types από Python calls σαν return types απο HTTP calls :-). Credits to vkoukis for the heads up.

#6 - 06/03/2011 04:42 pm - Vangelis Koukis

- Status changed from New to Closed

Έχει υλοποιηθεί μηχανισμός reconciliation.

Ενεργοποιείται με λήψη απαντήσεων σε αιτήσεις OP_INSTANCE_QUERY_DATA, που στέλνονται μαζικά από mgmt command του Django, περιοδικά.

Ο διαχειριστής μπορεί κατά βούληση (gnt-instance info) να κάνει το ίδιο, αν χρειαστεί.