

Managing Xen with SmartFrog

CERN openlab II quarterly review
31 January 2007
Preview of VHPC'07

Xavier Grehant





- On-demand execution environments must be:
 - Virtual
 - Distributed
 - Configurable
 - Composable

- For batch jobs (grid)
 - In contrast to Virtual Workspaces, Tycoon

- And software quality assurance (QA) tasks
 - In contrast to NMI builds & tests infrastructure

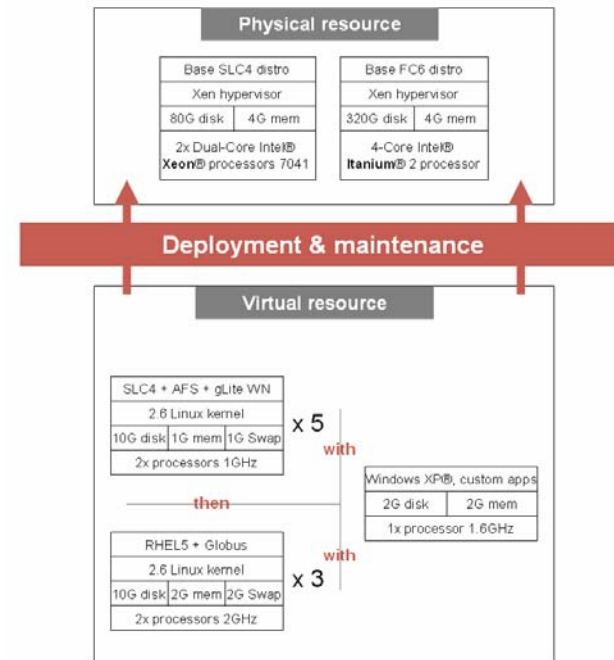
Why virtual resource management?

- Exploits benefits of virtual machines
 - Compatibility
 - Data isolation
 - Resource sharing and isolation
- Virtualization results in much more flexibility
- Xen enables automation
 - It does not provide the management system

- Xen
 - High performance
 - Advanced features
 - Popularity

- SmartFog
 - Description language
 - Configure and choreograph components
 - Tree of components with attributes
 - Daemons
 - Peer-to-peer network for deployment
 - Deployment engine
 - Interprets description
 - Dispatches work to daemons
 - Liveness, dependencies, references

- Xen VM deployment with SmartFrog
- Users submit a description to launch the pool of VMs
- SmartDomains automates deployment and management



- > sfStart localhost pool virtualPoolDesc.sf
- > sfTerminate localhost pool
- Simplicity on/off required for batch computing

- Other virtualization management systems:
 - Enterprise systems:
 - Platform VMO, Cassat Collage, OpenQRM, DynamicOE: Let admin define high-availability policies among apps
 - Open source systems:
 - Enomalism, Virtual Workspace + GPE: interface to Xen VMM

Usage: describing resource

```
volumeSize "1g";
swapSize "512m";
usingExistingVolumes false;
keepVolumes "statuquo";
saveImage true;
saveImageName "";
saveImagePath "/tmp";
saveImageExtension "tar";
makeFs "mkfs.ext3";
volumeBaseName "volume1";
tempMountPoint "/tmp/mnt";
kernel "/boot/vmlinuz-2.6-xen";
randisk "/boot/initrd-2.6-xen.img";
netmask "255.255.0.0";
memory 512;
vcpus 1;
domainLivenessDelay 2000;
domainLivenessFactor 3;
domainLivenessCheck true;
extra "fastboot nusb";
domainName "virtual-domain";
baseImage "slc3-WN.tgz";
ip "123.456.78.90";
gateway "98.765.43.21";
hostname "host";
```



Usage: describing resource

```
#include "org/smartfrog/components.sf"
#include "ch/cern/openlab/smartdomains/components.sf"
#include "org/smartfrog/services/shellscript/components.sf"

DefaultXenDomain extends XenDomain {
.   shell LAZY ATTRIB myShell;
.   kernel "/boot/vmlinuz-2.6-xen";
.   gateway " ";
.   netmask "255.255.0.0";
}

PhysicalHost extends Compound {
.   sfSyncTerminate true;
.   myShell extends BashShell;
}

sfConfig extends Compound {
.   sfSyncTerminate true;
.   computer1 extends PhysicalHost {
.       sfProcessHost " .cern.ch";
.       loop extends LoopbackStorageBackend {
.           shell LAZY ATTRIB myShell;
.           domainName "domainLoopback";
.           baseImage "/data/xen/slc3-smartfrog.img";
.       }
.       domain1 extends DefaultXenDomain {
.           domainName "domainLoopback";
.           ip " ";
.           hostname " -dom1";
.           storageBackend LAZY ATTRIB loop;
.       }
.       lvm extends LVMStorageBackend {
.           shell LAZY ATTRIB myShell;
.           domainName "domainLVM";
.           baseImage "/data/xen/slc3-smartfrog.img";
.       }
.       domain2 extends DefaultXenDomain {
.           domainName "domainLVM";
.           ip " ";
.           hostname " -dom2";
.           storageBackend LAZY ATTRIB lvm;
.       }
.   }
.   computer2 extends PhysicalHost {
.       sfProcessHost " .cern.ch";
.       ...
.   }
}
```


Web_GergoUI	hary_test	slc4TarTest	tarDistCreate	output																																					
rootProcess	sfDefault	oliver	joachim	Di-1	Di-2	Web_new_test	Di-5																																		
<ul style="list-style-type: none"> rootProcess <ul style="list-style-type: none"> *copy* sfDefault <ul style="list-style-type: none"> oliver <ul style="list-style-type: none"> *copy* *copy* joachim <ul style="list-style-type: none"> *copy* *copy* Di-1 <ul style="list-style-type: none"> *copy* *copy* Di-2 <ul style="list-style-type: none"> *copy* *copy* Web_new_test <ul style="list-style-type: none"> *copy* *copy* Di-5 <ul style="list-style-type: none"> *copy* *copy* Host-... .cern.ch Web_GergoUI <ul style="list-style-type: none"> *copy* *copy* hary_test <ul style="list-style-type: none"> *copy* *copy* slc4TarTest <ul style="list-style-type: none"> *copy* *copy* tarDistCreate <ul style="list-style-type: none"> *copy* *copy* VM <ul style="list-style-type: none"> myShell ctb-generic-1 VM <ul style="list-style-type: none"> myShell ctb-generic-3 VM-image VM 		<table border="1"> <thead> <tr> <th>Attribute</th> <th>Value</th> </tr> </thead> <tbody> <tr><td>sfCodeBase</td><td>"default"</td></tr> <tr><td>sfClass</td><td>"ch.cern.openlab.smartdo..."</td></tr> <tr><td>schema</td><td>LAZY ASSERT (shell APPLY (...</td></tr> <tr><td>shell</td><td>LAZY myShell</td></tr> <tr><td>volumeGroup</td><td>"vg1"</td></tr> <tr><td>volumeSize</td><td>"10g"</td></tr> <tr><td>savelImage</td><td>false</td></tr> <tr><td>savelImageExtension</td><td>"tar.gz"</td></tr> <tr><td>savelImagePath</td><td>"/tmp"</td></tr> <tr><td>baseImage</td><td>"/Nodes_status_Backup/tar..."</td></tr> <tr><td>domainName</td><td>"ctb-generic-3"</td></tr> <tr><td>sfLivenessDelay</td><td>240L</td></tr> <tr><td>sfLivenessFactor</td><td>2</td></tr> <tr><td>sfHost</td><td></td></tr> <tr><td>sfProcess</td><td>"rootProcess"</td></tr> <tr><td>undefined</td><td>"Creating swap volume."</td></tr> </tbody> </table>						Attribute	Value	sfCodeBase	"default"	sfClass	"ch.cern.openlab.smartdo..."	schema	LAZY ASSERT (shell APPLY (...	shell	LAZY myShell	volumeGroup	"vg1"	volumeSize	"10g"	savelImage	false	savelImageExtension	"tar.gz"	savelImagePath	"/tmp"	baseImage	"/Nodes_status_Backup/tar..."	domainName	"ctb-generic-3"	sfLivenessDelay	240L	sfLivenessFactor	2	sfHost		sfProcess	"rootProcess"	undefined	"Creating swap volume."
Attribute	Value																																								
sfCodeBase	"default"																																								
sfClass	"ch.cern.openlab.smartdo..."																																								
schema	LAZY ASSERT (shell APPLY (...																																								
shell	LAZY myShell																																								
volumeGroup	"vg1"																																								
volumeSize	"10g"																																								
savelImage	false																																								
savelImageExtension	"tar.gz"																																								
savelImagePath	"/tmp"																																								
baseImage	"/Nodes_status_Backup/tar..."																																								
domainName	"ctb-generic-3"																																								
sfLivenessDelay	240L																																								
sfLivenessFactor	2																																								
sfHost																																									
sfProcess	"rootProcess"																																								
undefined	"Creating swap volume."																																								

- Full configurability with base components attributes
 - Compared to:
 - Amazon EC2: same server, custom filesystem
 - Tycoon: same filesystem, custom resources
- Lifecycle management with components composition
 - Never seen before (acknowledged as issue in Xen roadmap)

- Specially suited for trusted community (P2P)
 - A computer bootstraps whole resource
 - Security system follows same scheme
- Predefine specialized components in description language
 - Extension mechanism, links
 - For specific usage, or simplicity of end-users descriptions
- Or provide a web interface
 - Hide descriptions, fill up missing fields
- Example: gLite testing

- Composite pattern:
 - Plug-in functionality
 - Scheduling, balancing, high-availability
 - Create higher-level structures
 - Virtual clusters
 - Modularity and reuse
- Peer-to-peer
 - Scope of an algorithm: the P2P network
 - As opposed to Tycoon where bidding scope is inside a physical host
 - No single point of failure

```
simpleScheduler extends Scheduler {
    hosts [|"host1", "host2", "host3"|];
}
VMs2Dispatch extends Schedulee {
    scheduler LAZY ATTRIB simpleScheduler;
    - extends VM {...}
}
```

- In the future, resource = VM

- SmartDomains uniqueness
 - Batch jobs tests: on / off
 - Distributed: workflows and lifecycle management
 - Peer-to-peer
 - Composition

- Applications:
 - Batch computing
 - QA tasks
 - Direct / specialized / enriched usage