





Enabling Grids for E-scienceE



## Tools and techniques for managing virtual machine images

*Andreas Unterkircher, Dimitar Shiyachki*  
*CERN Grid Deployment Group*  
*Havard Bjerke*

www.eu-eggee.org

EGEE and gLite are registered trademarks



Enabling Grids for E-scienceE

## Content

- **Motivation**
- **Tools developed**
  - VM image generation (libfsimage)
  - GUI for VM image generation (OSFarm)
  - Efficient transfer of VM images

Andreas Unterkircher

HEPIX Spring Meeting 2009 2




Enabling Grids for E-scienceE

## gLite certification process

- **GLite uses a continuous release process. Services are updated individually on top of a baseline release.**
- **Updates are added via a patch**
  - We use Savannah for bug and patch tracking
  - Has one or more bugs attached
  - Can also be used to introduce new features or services
  - Has all relevant information: OS, architecture, affected services, baseline release, configuration changes, rpm lists etc.
- **Patch is being certified**
  - Update and configure affected services
  - Run basic and regression tests
  - Verify if attached bugs are fixed and write regression tests.
  - Put patch into “certified” or “rejected”. For the former the patch is ready for release to pre production and later for production.

Andreas Unterkircher

HEPIX Spring Meeting 2009 3



Enabling Grids for E-scienceE

## Problems in patch certification

- **Certification of several patches at the same time can cause conflicts.**
- **A non functional patch may spoil the whole testbed**
- **Patch certification often fails already at an early stage (rpm installation, configuration)**
- **A failed patch can pollute a machine. A complete reinstallation is necessary.**
- **Many scenarios and interactions have to be considered: gLite 3.1/gLite 3.2, SL4/SL5/Debian, x86/x86\_64**
- **165 patches treated in 2008**

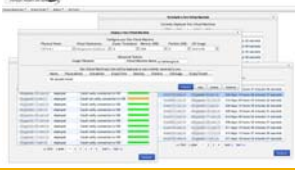
Andreas Unterkircher

HEPIX Spring Meeting 2009 4

**EGEE** Enabling Grids for E-science

## Virtualization for patch certification

- Patch certification can only be done efficiently by heavily using virtualization
- Already in summer 2006 we ran our own VM management system now known as VNode
- We collaborate with various groups in CERN IT
  - Openlab: OSFarm
  - FIO: Quattor profiles, SLC Xen support
  - Netops: hostnames for VMs



Andreas Unterkircher HEPIX Spring Meeting 2009 5

**EGEE** Enabling Grids for E-science

## Certification process

- We operate a certification grid testbed offering all services.
- Certifiers bring up virtual machines containing the updated versions of nodes affected by a patch and connect them to the testbed.

Andreas Unterkircher HEPIX Spring Meeting 2009 6

**EGEE** Enabling Grids for E-science

## Requirements for VM images

- We want to reproduce what is deployed at a site
  - Base OS installation + gLite node type
- Creating a VM images by scratching a physical machine is not practical
- A VM creation engine should run on an existing OS in a chroot environment
- VMs with different Linux flavors and architectures: SL(C)\*, Debian,... x86 and x86\_64

Andreas Unterkircher HEPIX Spring Meeting 2009 7

**EGEE** Enabling Grids for E-science

## VM image generation - libsfimage

- Python library for generating Linux file systems and populating them
  - Has command line interface
  - Produces images in a chroot environment
- Produces a tar.gz file that can be used as an image to boot with Xen
- Supported distributions: SL(C)3/4/5, Debian, Ubuntu, CentOS, Fedora on x86 and x86\_64
- Available in the xenvirt module in CERN's CVS

Andreas Unterkircher HEPIX Spring Meeting 2009 8

egEE Enabling Grids for E-science

## Example

```

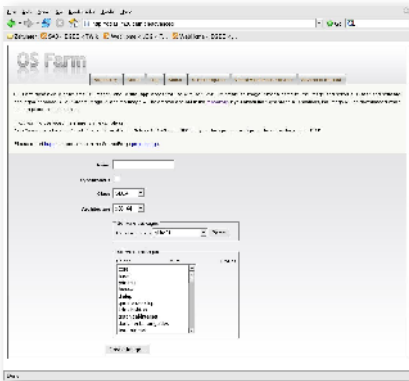
./genfs.py
-t SL4_x86_64
-d /tmp/SL4fs
-o /tmp/SL4_x86_64_img1.tar.gz
-g Base Core
-p python openldap-client
-w rootpwd

```

Andreas Unterkircher HEPIX Spring Meeting 2009 9

egEE Enabling Grids for E-science

## OSFarm



- Uses libsvmimg for creating core images
- GUI for configuring and adding software to images
- Optimizes the image creation process by caching shared data

Andreas Unterkircher HEPIX Spring Meeting 2009 10

egEE Enabling Grids for E-science

## OSFarm



Images and their configuration (XML) are stored in the repository for later retrieval

Each request is checksummed and compared to existing configurations

Andreas Unterkircher HEPIX Spring Meeting 2009 11

egEE Enabling Grids for E-science

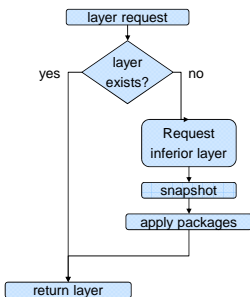
## Layered image generation

- **Image core:** minimal VM image
  - can be shared between VM image configurations
- **Image base:** contains software critical for virtual appliances
  - Can be shared between virtual appliances
- **Image generation process is optimized by caching and sharing lower stages of an image**

Andreas Unterkircher HEPIX Spring Meeting 2009 12

## Copy-on-write staging

- Base stages are kept in stage
- Uses LVM snapshots (copy-on-write) for instantaneous staging
- Tag of a cached image is the checksum of its configuration parameters

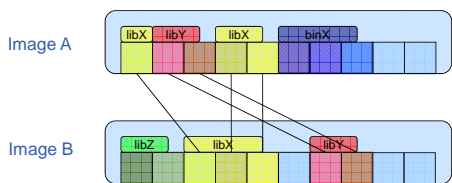


## Using images on the Grid

- Many different user communities use the Grid. It is difficult for the Worker Nodes to fulfill all the different requirements.
- Provide each job on the Grid a virtual machine with a dedicated OS setup
- With thousands of users image transfer becomes an issue
- **Observation**
  - Images are often similar
- **Content-Based transfer: don't transfer the whole image, just transfer the difference to what is already available**

## Content-Based transfer

- Each file starts on a block boundary
- Identical blocks can be identified with a hash checksum



Only 50% of image A needs to be transferred if image B is already at the destination

## Image comparisons

- **High commonality can be expected in practice**
  - Different Update versions of same base OS
  - Similar applications needed by users from same research area
- **SL 3 (343 MB) and SL 4 (762 MB)**
  - Transfer SLC4 to SLC3
    - 22 % common blocks
  - Transfer SLC3 to SLC4
    - 48 % common blocks

**eggee** Enabling Grids for E-science **Cost**

- **Generating hash tables for source file and target repository – linear cost**
- **Accessing hash tables**
- **Hash table data overhead**
  - Depends on
    - Hash function (e.g. SHA is 20 bytes)
    - Block size
  - Measured 0.5 – 2.0 % of the image size

Andreas Unterkircher HEPiX Spring Meeting 2009 17

**eggee** Enabling Grids for E-science **Implementation**

```

    graph LR
      subgraph Source_machine [Source machine]
        direction TB
        HPT[Hash produce thread]
        BRD[Block request dispatch]
      end
      HD[Hash dispatch]
      IR[(Image repository)]
      BD[Block dispatch]
      TI[(Target image)]
      HPT --> HD
      HD --> IR
      BRD --> BD
      BD --> TI
  
```

- Prototype implementation done By H. Bjerke at openlab
- Multi-threaded
- Implemented in Java

Andreas Unterkircher HEPiX Spring Meeting 2009 18

**eggee** Enabling Grids for E-science **Measurements**

- Disk read speed: 35.6 MB/s
- Network bandwidth: 11.9 MB/s
- Delta = difference between source and target

delta	Theoretical max observed bandwidth (MB/s)	Measured observed bandwidth (MB/s)	Measured real bandwidth (MB/s)
1	~12	~10	~10
0.5	~18	~15	~10
0.1	~30	~28	~10
0	~35	~33	~10

Andreas Unterkircher HEPiX Spring Meeting 2009 19

**eggee** Enabling Grids for E-science **Discussion**

Andreas Unterkircher DESY Virtualization Workshop 20