

HEPIX VWG Image transfer.

HEPIX virtualisation working group : 6 Month update

A short summary of the Documentation produced.

[HTML PDF A4 PDF Letter](#)

Owen Synge

HEPIX VWG Image transfer.

HEPIX fall 2011

HEPIX VWG Assumptions

- For new customers Cloud may be all we need.
 - But HEP is not a new customer.
- HEP Experiment software is currently partially trusted.
 - Sites allow NFS 3, sites are not ready for untrusted images.
 - NFS 4.1 with Kerberos security will be practical and possible.
 - X509 to Kerberos account mappers exist in the grid.
 - HEP experiments are not yet ready to abandon rshell style data access.
- Virtualising Worker node can be transparent for grid users.
 - We need to trust images more than with a cloud infrastructure.
- Grid/Batch Queues have high efficiency and high use.
 - We should demonstrate our ideas work with the grid.
- Accountancy is mature in the Grid.

HVWG Failures 1.

- We have not yet put enough effort into understanding the differences in cloud products.
 - Planning the mechanism of software evaluation.
 - Software.
 - WN on demand from INFN.
 - Stratus Lab marketplace.
 - Stratus Lab Open nebula.
 - OpenStack.
 - Big Grid.
 - Nimbus Infrastructure.
 - Evaluation
 - Logging.
 - Stability.
 - Scalability.
 - User management.
 - Security.
 - Accountancy.



Even planning comparison evaluation criteria is a lot of work.



HVWG Failures 2. We need to meet more often.

- We have not organized enough meetings without Tony.
- Tony proposes a face to face soon.
- The work days where a success.
 - Not enough sites can as yet host VM's in a production like way.
 - Expect this to change.
 - Stratus lab Market place seems to work.
 - Image list based transfer is working.
 - I think we need to get back to keeping minutes.
- A lot of cloud efforts are occurring.
 - We cant all even find all of them.

HVWG Failures 3. Clearing up misunderstandings

- Our work in HEPIX is intended to make it easy for scientists too.
 - Security is about enabling options (like signing an image list).
 - We welcome VO's to contribute VM for us to cache.
 - Providing guidelines how to get your image to resources.
 - Providing guidelines in making images site neutral.
- Some sites will not need to trust the VM Image.
 - User communities that access storage using secure protocols.
 - Sites with VLAN per user.
 - Giving users safe network isolation from each other.
- Some sites maybe working on alternatives to the grid accountancy.
 - Interesting work going on in big grid.
- If you have a different view how we carry on with this effort.

Four Areas of Focus

> Security Policies.

- Latest Draft

- <https://documents.egi.eu/public/ShowDocument?docid=771>
- “security-related policy requirements for the generation, distribution and operation of virtual machine images”
 - Policies are Valid even for Secured by VLAN systems.

> Image Creation.

- How to make an image XEN KVM neutral.

> Image Transfer.

- Most of this talk.

> Image Contextualization

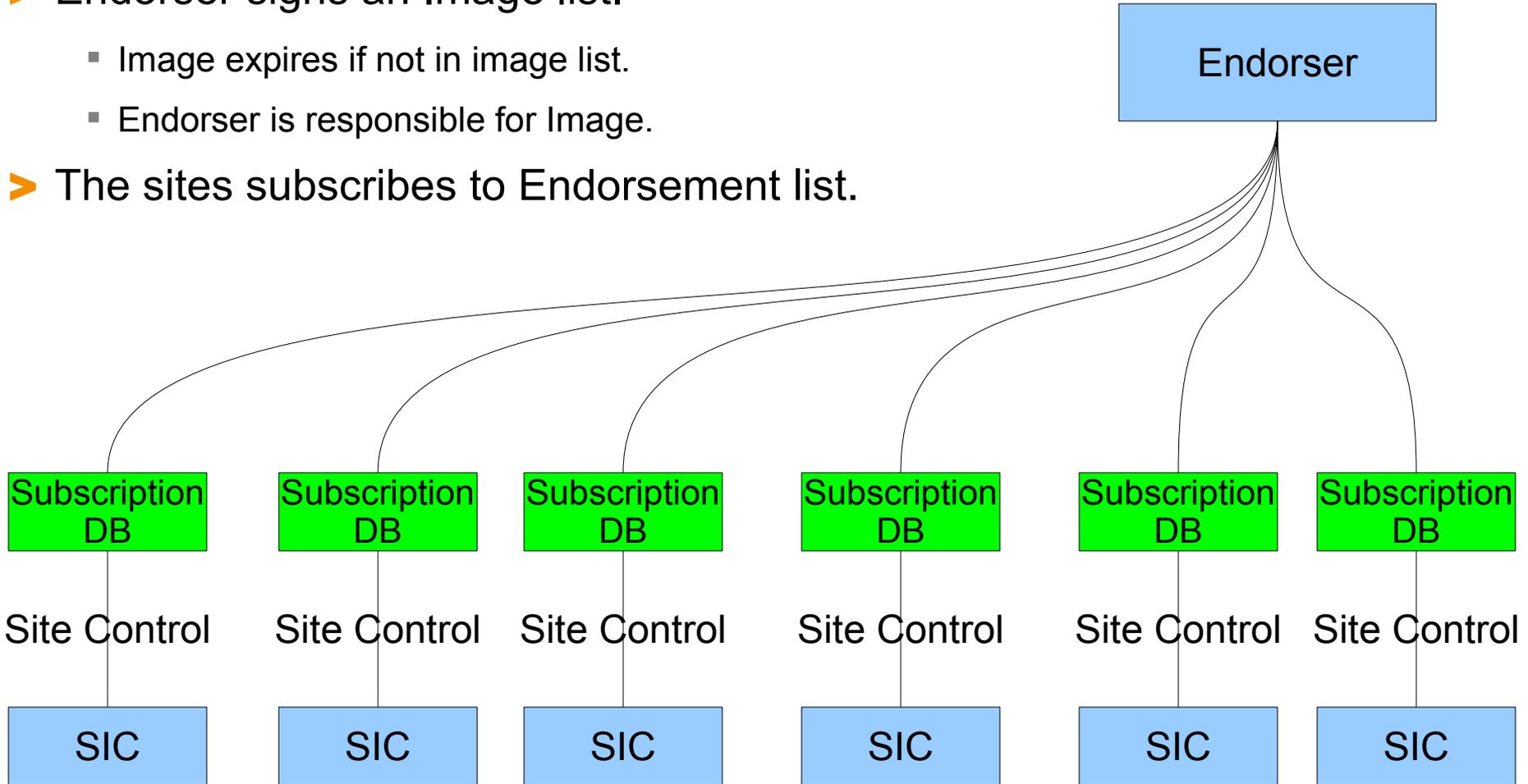
- Attach a ISO image (virtual CDROM) of site specific batch queue client software.
- Boot time scripts call ISO image.

HVWG model is Publish and Subscribe RSS.

> Endorser signs an Image list.

- Image expires if not in image list.
- Endorser is responsible for Image.

> The sites subscribes to Endorsement list.



Who Publishes/Subscribes?

> CernVM.

- 28 images for different VM formats.
- 12 months expire time.

> Stratus Market place.

> UVIC.

- 1 month lifetime.

> CERN

- Instantiated images from CernVM, UVIC.

> Edinburgh

- Instantiated images from via VMILS.
- Instantiated images from the Stratus Marketplace, CernVM,UVIC

> INFN

- Instantiated images via VMLIS on SL 5.



> Image to Meta data binding.

- Cryptographic hashes.
 - It is easy to compute the hash value for any given data.
 - It is infeasible to generate a message that has a given hash.
 - It is infeasible to modify a message without hash being changed.
 - It is infeasible to find two different messages with the same hash.
- Chose to use sha512 and file size to validate data.
 - Following Stratuslabs recommendation.
- Other hashes can be added.
 - If sha512 and size are later found to be too weak.
- URI to retrieve image.
 - Can be cached locally.
- Each image has a UUID
 - So we know which image is expired and which is upgraded.



Meta-data Security.

> Meta-data authenticity.

- X509 + signatures. (SMIME or XML signatures)
 - Gives non repudiation, and confidence in who endorsed.
 - Give tamper proof message.
 - Signature can be checked by all clients,
 - Allows checking of historic meta-data changes.
- Version number.
 - Prevents man in middle attacks.
 - Man In Middle attempts to return an old list blocked by this.
- UUID on Image and Image list
 - Allows messages to be identified.
 - So messages cannot effect each other.
 - So images can be expired and updated.



Meta-data subscription DB

> Mostly no admin interaction!

- Database is auto created.
- All subscriptions updated from a cron script.
- All data is derived from subscriptions to image lists.
 - So just need to store signed image lists which you should anyway.
- Migration is simply install a second in parallel.

> Simple RDBMS

- Almost No critical data to back up.
 - You might want to back up users and subscription list.

> Image cache as a client of the subscription data base.

- 3 Directories (should be on same file system).
 - Current, Expired, and Downloading directories.
- Current images are be validated.
- JSON catalog for each directory.



Whats been done in the last 6 months.

- HVWG Document generated nightly from DESY SVN.
 - HTML PDF A4 PDF Letter
 - <http://grid.desy.de/vm/hepix/vwg/doc/>
- Completed the Image cache code.
 - Initial version subscribed to all images in an image list.
- New feature: Images can be selected to be cached.
 - CERNVM has 28 different versions of same image.
 - Some of these where about 30G and not wanted.
- Test suite had to be developed.
 - Bugs fixes where getting to much work to replicate.
- Image transfer software is now stable and tested.
 - 2 data base schema evolutions since last meeting.



Whats been done in the last 6 months.

- Bug Fix: Subscription software now enforces UUID constraints.
 - Warn admin when unregistered images are added to a subscribed list.
 - Ignore image UUID's that clash with existing image UUID's.
- Enhancement: Logging has been upgraded
 - Now use standard python logging configuration.
 - Allowing admins to change logging.
- New Feature: Publisher User management enhanced in subscriber code.
 - Can now add delete user separately from subscriptions.
 - Can bind users to one or more subscription.
 - NIKEF feedback wanted changes to subscription model.
 - Not just based on URL, but on approved DN.
 - Now implemented.
- New Feature: Support Host certificates.
 - **Should I add this feature?**



What we learn through subscription and publishing.

- Code to publish image lists is almost trivial.
 - Minimal less than 70 lines of python (including the metadata and white space)
 - Command line code is considerably bigger as its easier to use.
- Subscribing to image lists is more complex than publishing.
 - Test suites have been needed.
- Test Image list updates are infrequent.
 - Our test image list providers are not updating regularly.
 - Test suites give confidence that updates are handled.
 - Hope this will change with more production VM hosting environments



Whats still to look at (so far)

- Network technology. (VLAN, subnets, monitoring at only network level)
- Contextualization
 - Can we just use the open nebular model?
 - Was our first plan.
 - Device drivers change with VHost setup.
 - First drive can be: sda,hda,vda
 - CDROM ISO mounting.
 - Need to come up with more conventions.
- Clouds and fair share.
 - UVIC cloud scheduler etc.
- Pilot factories and Virtualisation.
 - Can we just launch precooked images that hook to the pilot frameworks.
- Finalise our first iteration of documents.
- Stratus Lab's Image marketplace.
 - Need to look into this further. (many bug fixes since first looks)



Last words : Coming to the End of first Iteration.

- VMIC is replacing VMIC Strawman.
- HEPix Virtualisation Working Group is coming to end of first iteration.
 - We did what we were created for but its clear more work is needed.
 - We documented what we did.
 - <http://grid.desy.de/vm/hepix/vwg/doc/html/index.shtml>
 - <http://grid.desy.de/vm/hepix/vwg/doc/pdf/Book-a4.pdf>
 - <http://grid.desy.de/vm/hepix/vwg/doc/pdf/Book-letter.pdf>
- Glue supports Virtualised Execution Environments.
 - Should we start using this?
- I hope for changes to Cream CE to allow users to select images.
 - JDL Change/Addition?
 - We need support in Batch Queue integration.
- Image creation and Contextualization has only been demonstrated.
- Image List base VM Image caching is demonstrated.
 - Progress to putting it in EPEL.

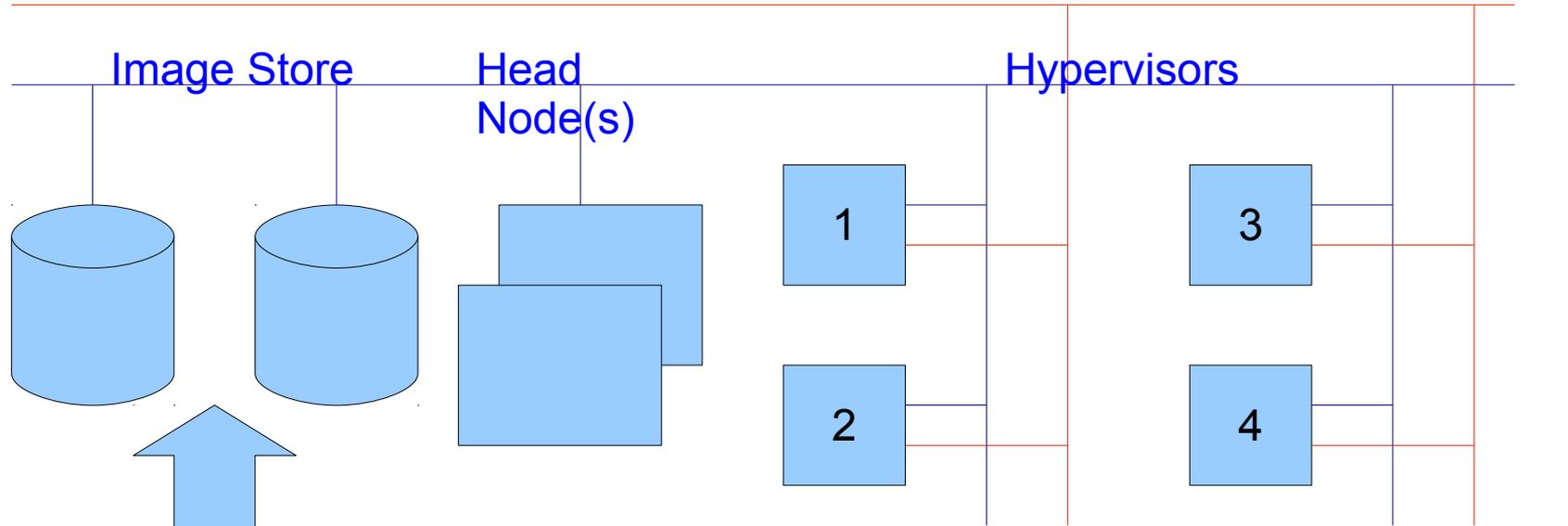


Summary (**Concepts matter not implementation**)

- Policy for dealing with Virtual Machine Images.
- Recipes for image creation as part of a Virtualised Grid Worker Node.
- Signed image lists define images that are endorsed.
 - First version of meta data is defined.
 - Non repudiation of image lists through signatures.
- Only Images on current Image list are endorsed.
 - This means images expire when not in current image list.
- Principle is generic to Clouds, Virtualised Worker Node.
- Implementation of Message Generation/Subscription exist.
 - Working on getting release into EPEL repository.
 - Further working starting in CERN/Academia Sinica
- We recommend the concept of **Signed image lists**.
 - And using Publish Subscribe model.



DESY 1 Cloud : Network View : Current + Planning



Blue (Current):

ONE: 1 HN + 5 Hs

Cloud Layer

256 private IPs

Planned:

As Cloud Admin Net

with

ONE: 2 Hns + 14 Hs

Red(Planned):

VLAN's (802.1Q +):

512 private Cloud IPs

512 public Cloud IPs

nnn Selected DESY

TODO:

High availability Head nodes

Storage is not yet fixed



DESY2 : DISH

- Desy ImageSHaring : Reaching pre production soon.
- Running on dCache WebDav Door.
- Using our Quattor to publish a worker node.
- Plan to publish fresh images weekly.

