

# Managing Virtual Machine Images

2 November 2007

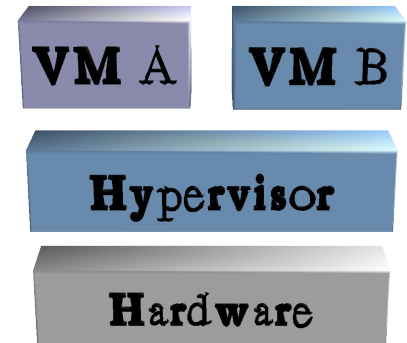
Håvard Bjerke





- HW virtualization introduction and motivation
- OS Farm
  - tool for creating and storing VM images
- Content Based Transfer
  - technique for efficient transfer of VM images

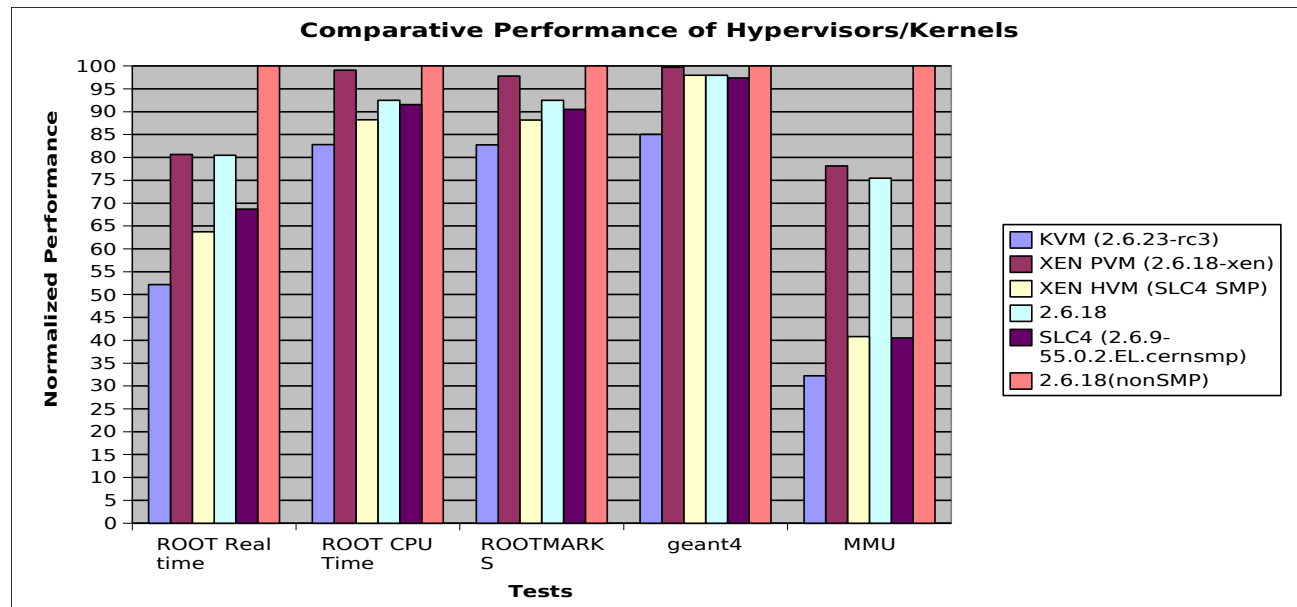
- Allows running several virtual machines (VMs) simultaneously on a single physical machine
- Classic consolidation scenario:
  - Run database and web server on the same machine
  - Run services in separate VMs, pinned to separate CPU cores
  - Save \$



- Benefits for Grid
  - Secure isolation
    - Small Trusted Computing Base in Xen
    - Isolate malicious software
  - Software flexibility
    - Better ability to satisfy requirements for execution environments
    - E.g. run both SLC3 and SLC4 on one physical node
  - Serialization, Live migration
    - Migrate essential services upon
      - hardware failure, or
      - maintenance

- HW acceleration
  - 1<sup>st</sup> generation Intel VTx CPU extensions
    - Allow full virtualization without binary rewriting or interpretation
    - A -1 or VMX privilege level
    - Already mainstream in Core architecture
  - 2<sup>nd</sup> generation Intel VTx CPU extensions
    - Add Extended Page Tables
    - Support guest VMs' page tables nested inside host's page tables
  - Intel VTd chipset extensions allow more efficient partitioning of I/O
    - Allocate device addresses to VMs

- Xen's virtual hardware has proved itself to be a good competitor to physical hardware
- Adds convenience while negligibly affecting performance



- Perhaps most significant (Grid related) developments

Past

- VM technology

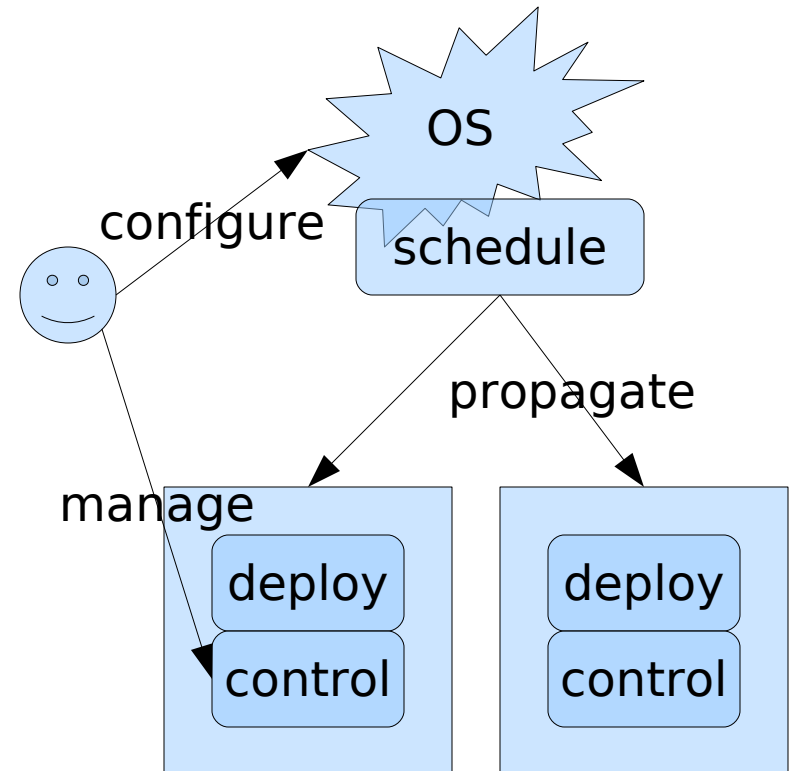
Present

- VM deployment
- Virtual appliances
- HW acceleration
- Testing environments

Future

- Authorization
- Security
- Contextualization
- Resource negotiation

- Offline
  - VM generation
  - VM scheduling
  - VM propagation
  - VM deployment
- Online
  - VM management
- Both
  - VM configuration





- Two tools already developed at CERN
  - SmartDomains
    - life-cycle management
  - vGrid
    - portal based
- Other models
  - Virtual Workspaces
    - VM scheduling and propagation
  - Batch system customization
    - LSF
    - Torque / MOAB scheduler

OS Farm



**CERN**  
**openlab**

- Various projects at CERN need quick deployment of clean machines
  - Software building and testing
- Need a rich set of Linux flavours
  - The default Linux distribution for CERN is Scientific Linux CERN 4 (SLC4)
    - SLC3 is being replaced by SLC4, but some software still depends on SLC3
    - SLC4 will be replaced by SLC5
  - Other flavours are needed for testing: Debian, Ubuntu, CentOS, Scientific Linux
- Windows, Solaris

- In a batch system, VM images can be used to satisfy execution environment requirements
  - A static set of requirements can be satisfied easily with a library of images
  - But requirements can have combinatorial complexity
  - Ultimately, users supplying their own images provides best coverage
  - Also better for software licenses

## OS Farm

[Repository](#)[About](#)[Log](#)[Status](#)[Simple request](#)[Virtual Appliances request](#)[Advanced request](#)

OS Farm dynamically generates OS images, and "virtual appliances" for use with Xen VMs. To create an image, enter a name for the image and select a "Class" and software packages if needed. Click "Create image...", and the image will be created and put in the [repository](#). If you check the "Download image upon creation" checkbox, the image will be downloaded when the image creation is finished.

If you do not enter a "Name", the image will be named after the md5 checksum of the image configuration parameters. If an image with the exact same parameters exists in the repository, it will not be recreated and can be downloaded immediately.

If you want to use wget, then here is an example url:

```
"http://www.cern.ch/osfarm/create?name=&download=on&class=SLC4&arch=i386&filetype=.tar&group=core&group=base&package=glite-BDII"
```

Please allow a few minutes for the image to be created.

Name

Synchronous

Class

Architecture

Filetype

- Base images
  - SLC 3 & 4 – standard at CERN
  - libfsimage – basis for several flavours
    - Debian and Red Hat based distributions
- Virtual appliances
  - gLite - Grid middleware
    - gLite-CE
    - gLite-WN
  - Quattor - fabric management
- 32 and 64 bit images
- tar or raw (\*.img) image format

## OS Farm

Repository
About
Log
Status
Simple request
Virtual Appliances request
Advanced request

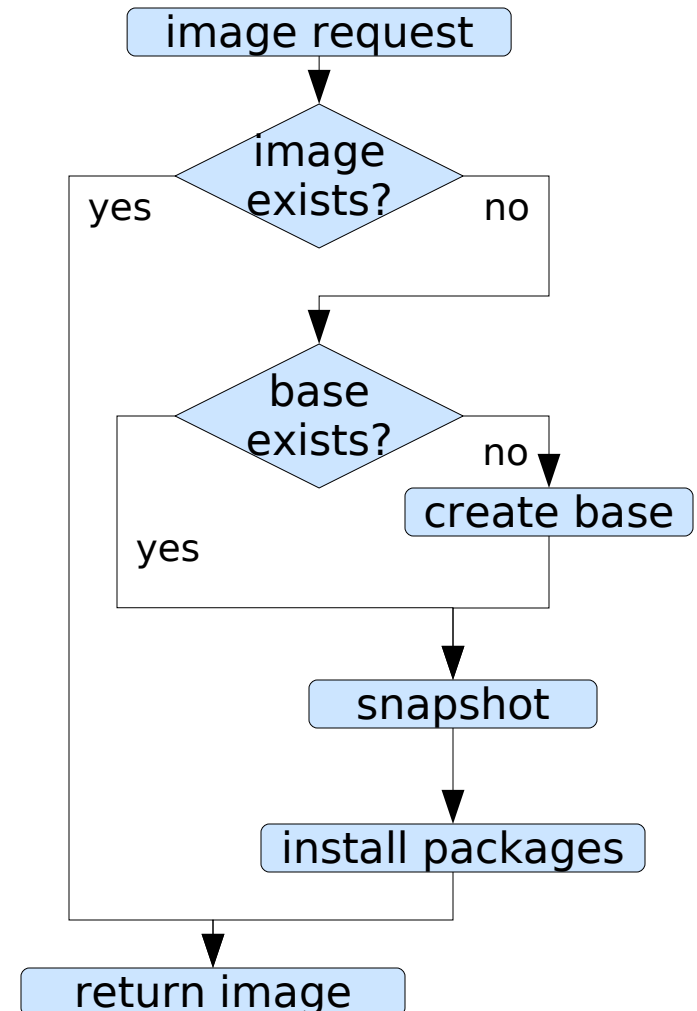
Location	Name	Class	Architecture	Filetype	Groups	Packages
<a href="#">download</a>		SLC4	i386	.img		<a href="#">delete</a>
<a href="#">download</a>	Test	SLC3	i386	.tar		<a href="#">delete</a>
<a href="#">download</a>	SLC3	SLC3	i386	.img		<a href="#">delete</a>
<a href="#">download</a>	sa301	SLC4	i386	.tar		<a href="#">delete</a>
<a href="#">download</a>	logo	SLC4	i386	.tar		<a href="#">delete</a>
<a href="#">download</a>	test	SLC4	i386	.tar		<a href="#">delete</a>
<a href="#">download</a>		glite-ce	i386	.tar		<a href="#">delete</a>
<a href="#">download</a>		SLC4	x86_64	.tar		<a href="#">delete</a>
<a href="#">download</a>	itmat	SLC4	x86_64	.tar		<a href="#">delete</a>
<a href="#">download</a>	image1	SLC4	i386	.img		<a href="#">delete</a>
<a href="#">download</a>		quattor-base	x86_64	.tar.gz		<a href="#">delete</a>
<a href="#">download</a>		SLC4	i386	.img	core base	<a href="#">delete</a>
<a href="#">download</a>		SLC4	i386	.tar	core base	glite-BDII <a href="#">delete</a>
<a href="#">download</a>		quattor-base	i386	.tar		<a href="#">delete</a>
<a href="#">download</a>		SLC3	i386	.tar.gz	core base	<a href="#">delete</a>
<a href="#">download</a>		SLC3	i386	.img	core base	<a href="#">delete</a>

Image configuration and image is stored in repository for later retrieval

Image configuration is stored in XML format

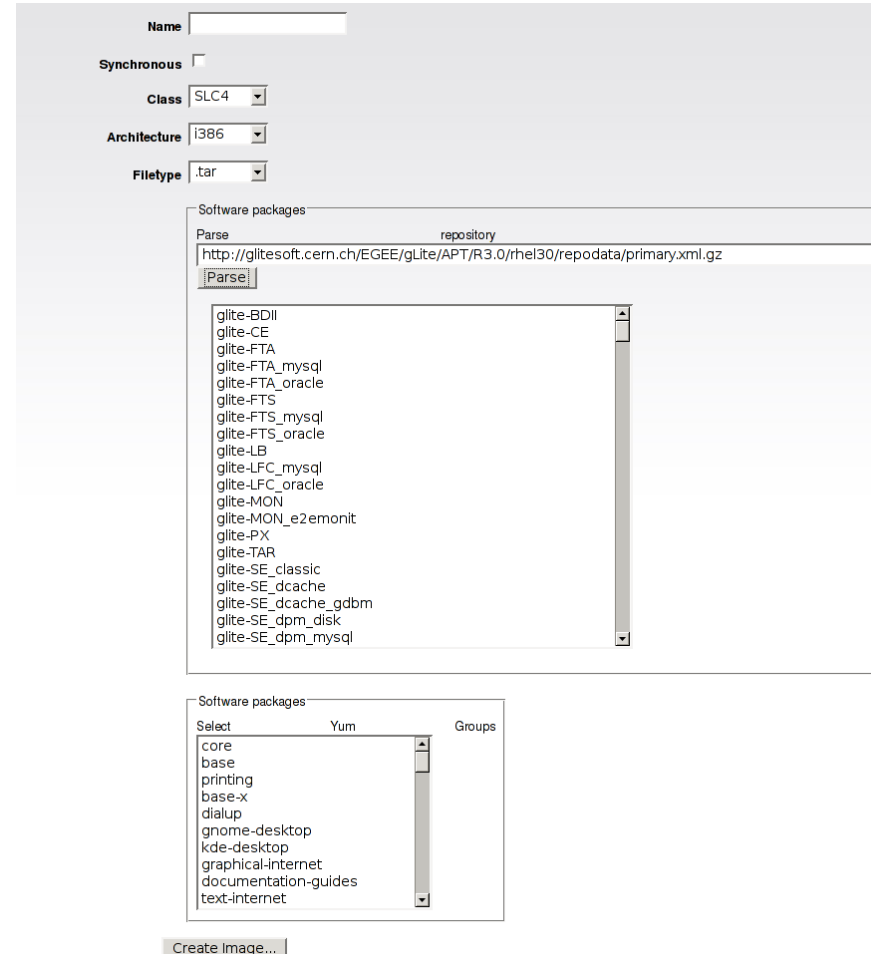
Each configuration is checksummed and compared to existing configurations  
 -> existing images are not recreated

- Images are created dynamically
- Base stages are kept in cache
- Uses LVM snapshots (copy-on-write) for instantaneous staging





- Most gLite and CERN fabric software reside in yum repositories
- Easy to add yum packages from yum repositories
- Supports most of CERNs fabric packages



The screenshot shows a web-based interface for configuring a yum repository. It includes the following fields and sections:

- Name:** An empty text input field.
- Synchronous:** A checkbox that is currently unchecked.
- Class:** A dropdown menu with "SLC4" selected.
- Architecture:** A dropdown menu with "i386" selected.
- Filetype:** A dropdown menu with ".tar" selected.
- Software packages:** A section containing:
  - A "Parse" button.
  - A text input field containing the URL: `http://glitesoft.cern.ch/EGEE/gLite/APT/R3.0/rhel30/repodata/primary.xml.gz`.
  - A "Parse" button below the URL field.
  - A scrollable list of package names: glite-BDII, glite-CE, glite-FTA, glite-FTA\_mysql, glite-FTA\_oracle, glite-FTS, glite-FTS\_mysql, glite-FTS\_oracle, glite-LB, glite-LFC\_mysql, glite-LFC\_oracle, glite-MON, glite-MON\_e2emonit, glite-PX, glite-TAR, glite-SE\_classic, glite-SE\_dcache, glite-SE\_dcache\_gdbm, glite-SE\_dpm\_disk, and glite-SE\_dpm\_mysql.
- Software packages:** A second section with a "Select" button and a scrollable list of yum groups: core, base, printing, base-x, dialup, gnome-desktop, kde-desktop, graphical-internet, documentation-guides, and text-internet.
- Groups:** A label positioned to the right of the group list.
- Create Image...:** A button located at the bottom of the interface.

- rPath – virtual appliance builder
  - Provides a massive selection of virtual appliances
  - Supports many different image formats
  - Roles separated into image authors and users
  - Proprietary build software
- OS Farm
  - Only one role – the user or a computer program
  - Dynamic image requests – image generated on the fly
  - Fully open source

- The Xen hypervisor isolates VMs from each other
  - As secure as or more secure than processes running in different address spaces
  - Xen has a small Trusted Computing Base
- If the user creates the VM image, she can have root privileges in the VM
  - Convenient for the user
  - Is it safe?
    - Promiscuous network access
  - Violations are traceable (encapsulation)

- Still, most administrators would be reluctant to allow users to deploy their own VM images on the grid
- Great effort goes into certifying SLC
- Running an uncertified OS on CERN Computer Centre hardware is risky
- Running an uncertified OS in an isolated VM is less risky

- **Solution: Secure images**
  - Controlled by admin (in control of propagation and deployment)
    - Administrator provides a static library of trusted images
    - Administrator provides a trusted image generation service
  - Trusted by certification, signed by
    - people
    - trusted image generation service

- There are several different OS image formats between vendors
  - VMWare (proprietary format)
  - Xen PVM
  - Xen HVM, KVM
- Images need host and site-specific configuration (network, etc.)
  - Images can be configured at deployment (parameterized) or pre-configured
- Goal: run my image on any Grid (LCG, Amazon EC2, my laptop)

- A DMTF (Distributed Management Task Force) Specification
  - Vendor neutral
  - Metadata specification
    - Virtual hardware description (CIM based)
    - CPU compatibility section
    - Feedback protocol
    - Service level requirements
  - Manifest
  - Certificate

- SOAP interface
  - Platform and language independent interface
- XML image descriptions
- Signed images (verify integrity)
- Content based image transfers



# Content Based Transfer

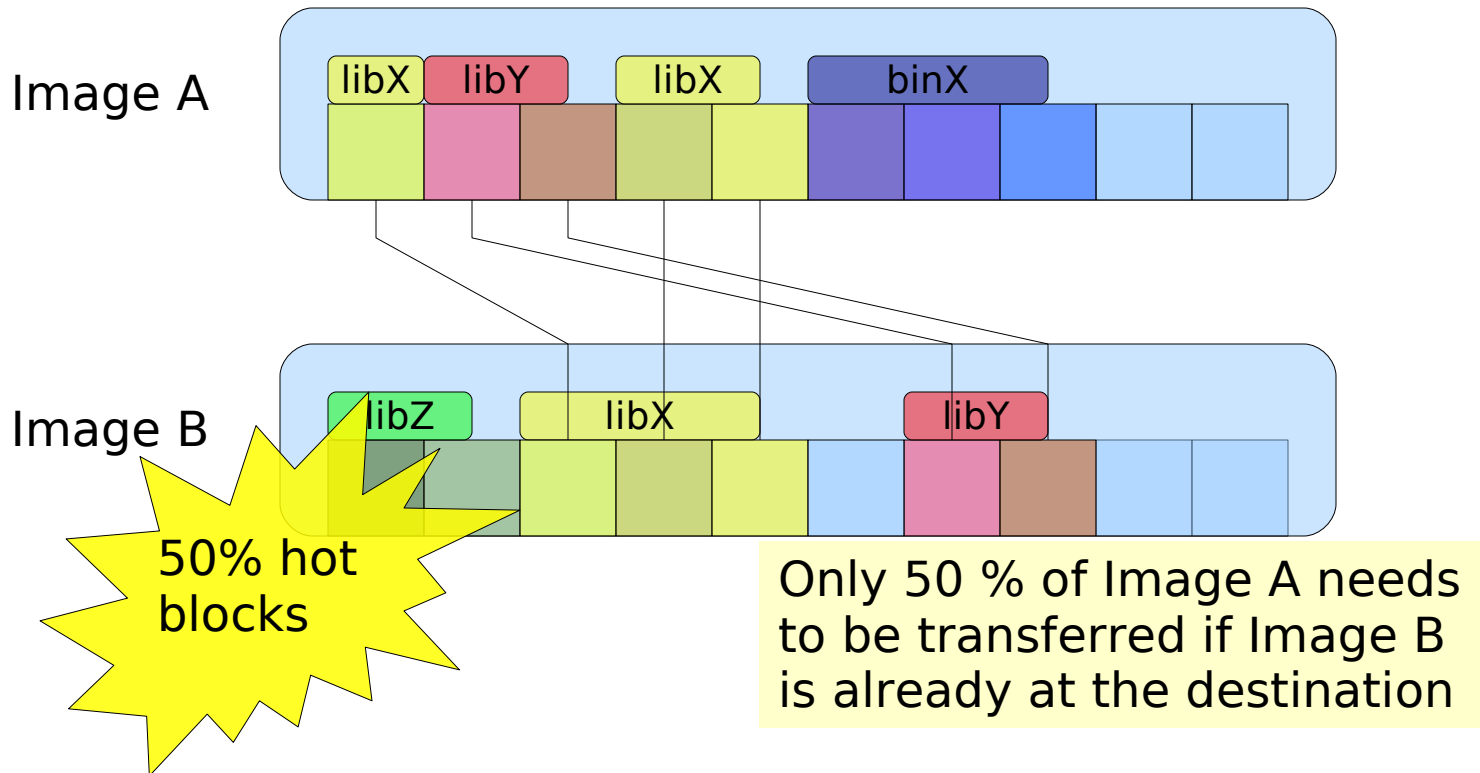


**CERN**  
**openlab**

- Jobs in the queue normally assume that files already exist at the execution node (e.g. through NFS) or that they will be staged when the job starts
- But VM images are big
  - ~ 300 MB to several GB
  - Jobs will have to wait for the image transfer to finish
  - Congests network

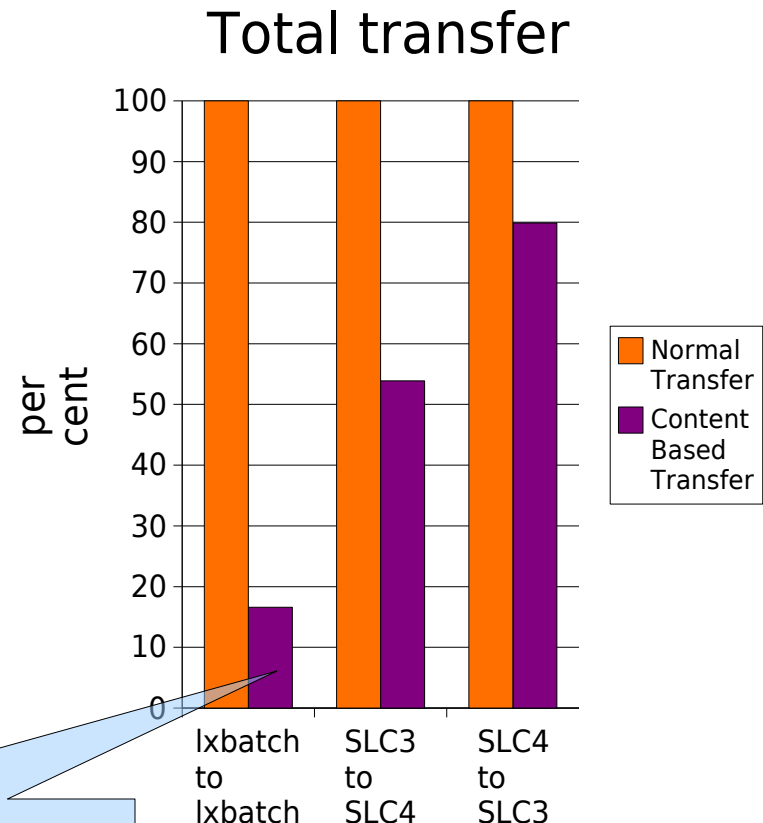
- Possible solutions?
  - Image cache
  - Transfer the image before the job is scheduled – image has to be on each possible target node
- Hard to satisfy each possible image configuration
- Observation from Content-based Addressing
  - Most images are relatively similar
  - No need to transfer the whole image; just transfer the delta

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



- Two typical batch machines (5.3 GB)
  - 84 % hot blocks
- SLC3 (343 MB) and SLC4 (762 MB)
  - SLC3 -> SLC4
    - 48 % hot blocks
  - SLC4 -> SLC3
    - 22 % hot blocks

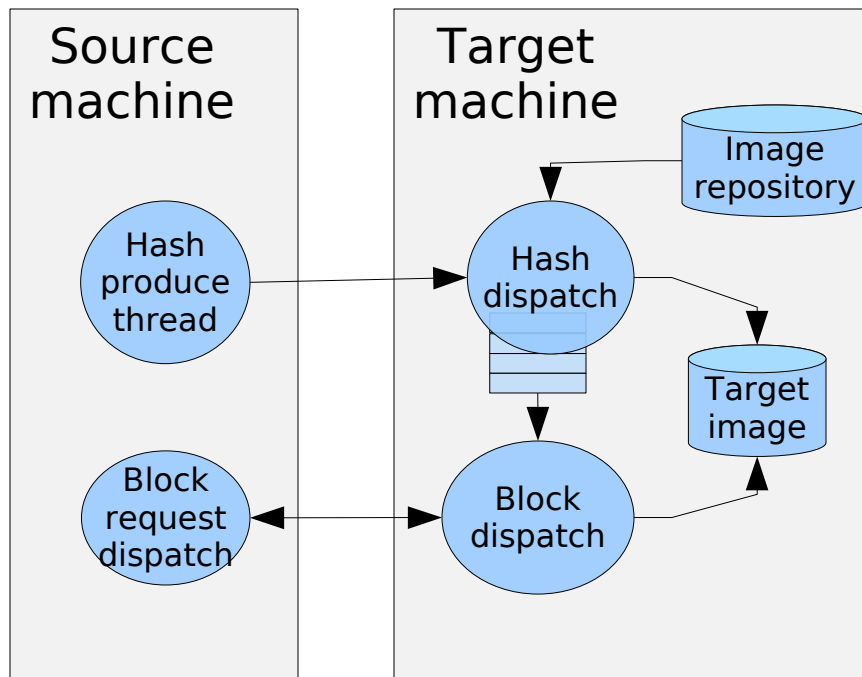
Fraction of full image data needed to transfer, including hash table



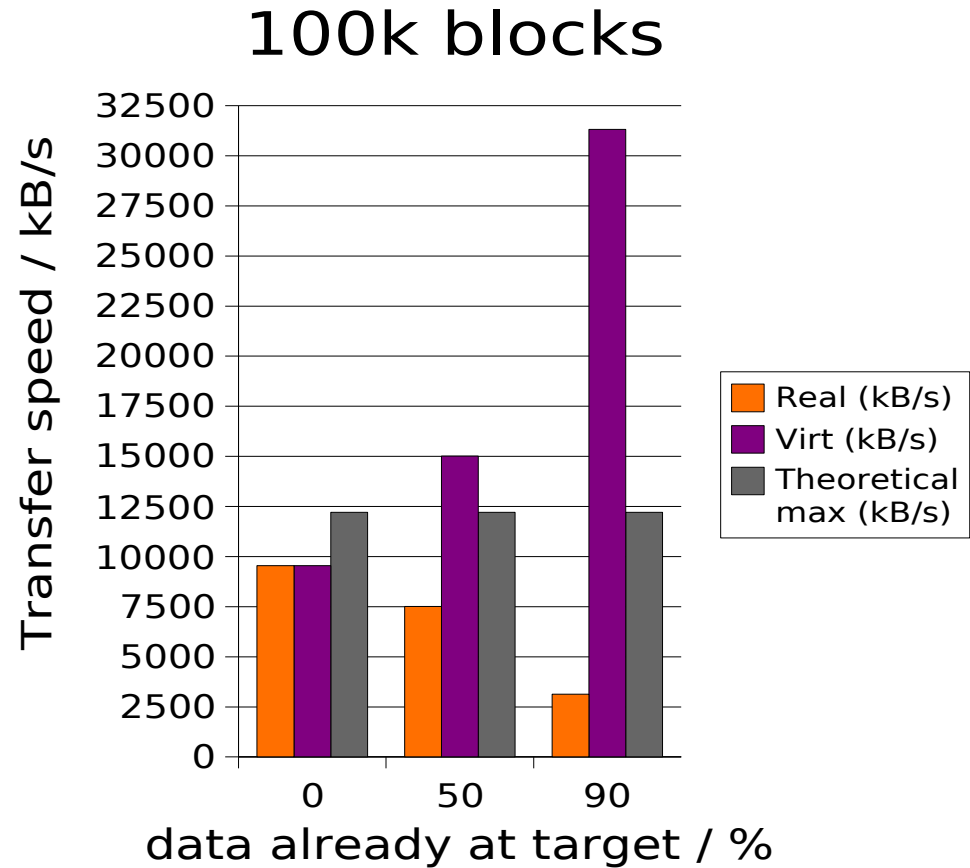
- Generating hash tables for source file and target repository
  - Linear
- Accessing hash tables
  - Java and Python have convenient constant-time hash tables
- Hash table data overhead
  - Depends on
    - hash function, e.g. SHA is 20 bytes
    - block size – usually 4096 bytes
  - 0.48 to 2.0 % of the image size

# Content Based Transfer

- Multithreaded
- Hash calculation and data transfer pipelined
- Implemented in Java (+ a Python prototype)

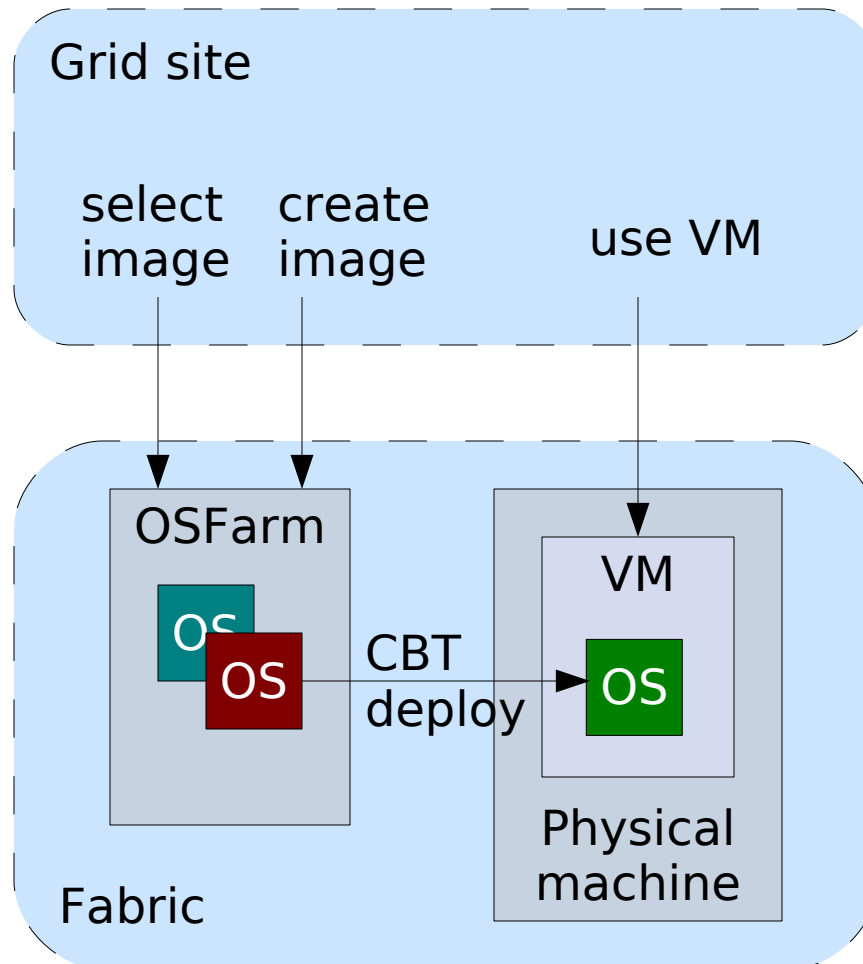


Virtual speed:  
full image size /  
time to transfer delta





# In control over Grid VM images



- OS Farm
  - <http://cern.ch/osfarm>
  
- Content Based Transfer
  - <http://hbjerke.web.cern.ch/hbjerke/cba/cba.xml>