# Tools and Techniques for Managing Virtual Machine Images

26 August 2008

Håvard Bjerke, Dimitar Shiyachki, Andreas Unterkircher, Irfan Habib

- ## Libfsimage
  - ### A library for creating VM images

- ## OS Farm
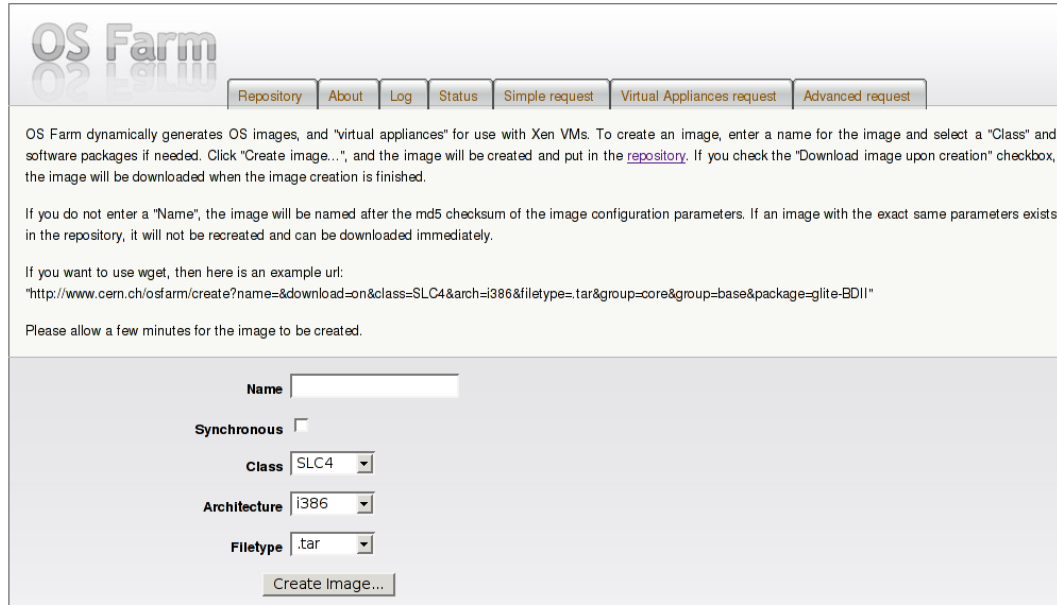  - ### Repository of tailored VM images

- ## Content Based Transfer
  - ### A technique for efficient transfer of VM images

- We need to test our Grid middleware on an increasing number of different OS flavours

- Need to increase the efficiency and reduce turnaround time of testing and certification

- Testing is most efficiently done in Virtual Machines

- A repository of unused VM images provides an easy way of setting up clean testing environments

# Libfsimage

- Driven by the need for testing software in a clean and reproducible environment

- Standalone tool and library for creating VM images

- Uses native package management tools in a chroot environment

- Can create many different Linux flavours

  - i386 and x86_64

  - Ubuntu, Debian, Fedora, Scientific Linux (CERN), CentOS

- Uses libfsimage for creating core images

- Adds a graphical user interface for configuring and adding software to images

- Provides a repository of stored images

- Optimizes the image creation process by caching shared data

- Web interface



- SOAP web service interface

- XML specification

# VM Image Repository



Images and their configuration are stored in repository for later retrieval

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

- Four different image stages: Core, Base, Image, and Virtual Appliance

- Image generation process is optimized by caching and sharing lower stages of an image

| Virtual appliance | |
|---|---|
| Base | Image |
| Core | |

- libfsimage is used to create *core*, or minimal, VM images
- *Core* consists of software critical to boot and access VM
  - Can be shared between VM image configurations
- *Base* consists of software critical for Virtual Appliances
  - Can be shared among Virtual Appliances

# Copy-on-Write staging

- Core and Base stages are kept in cache

- LVM snapshots (copy-on-write) are used for instantaneous staging

layer request

layer exists?

yes

no

Request inferior layer

snapshot

apply packages

return layer

- Any stage can be kept in cache

# Content-Based Transfer

- VM images are big
  - ~ 300 MB to several GB
  - Congests the network
  - Anything scheduled for a VM will have to wait for the image transfer to finish
- Observation from Content-based Addressing
  - Most images are relatively similar
  - Not always necessary to transfer the whole image; just transfer the delta

- Each file starts on a block boundary

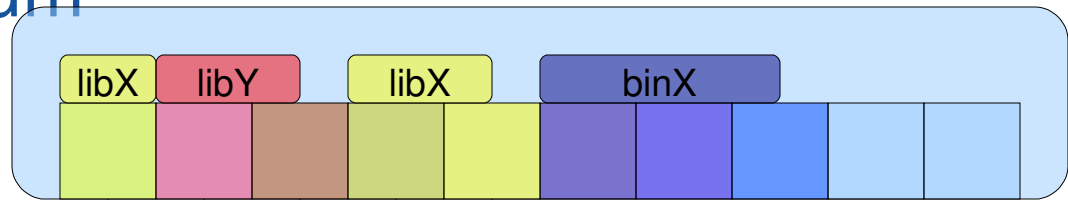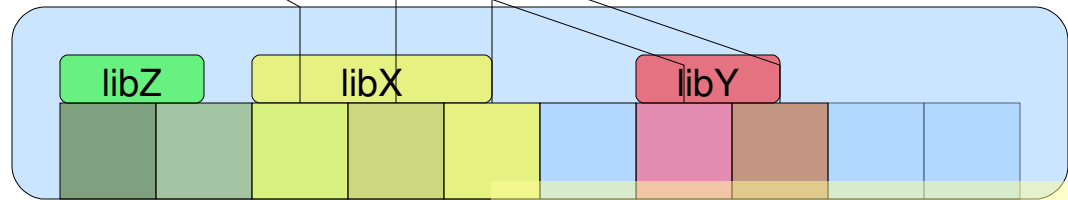- Identical blocks can be identified with a hash checksum

Image A

| libX | libY | | libX | | binX | | | | |

Image B

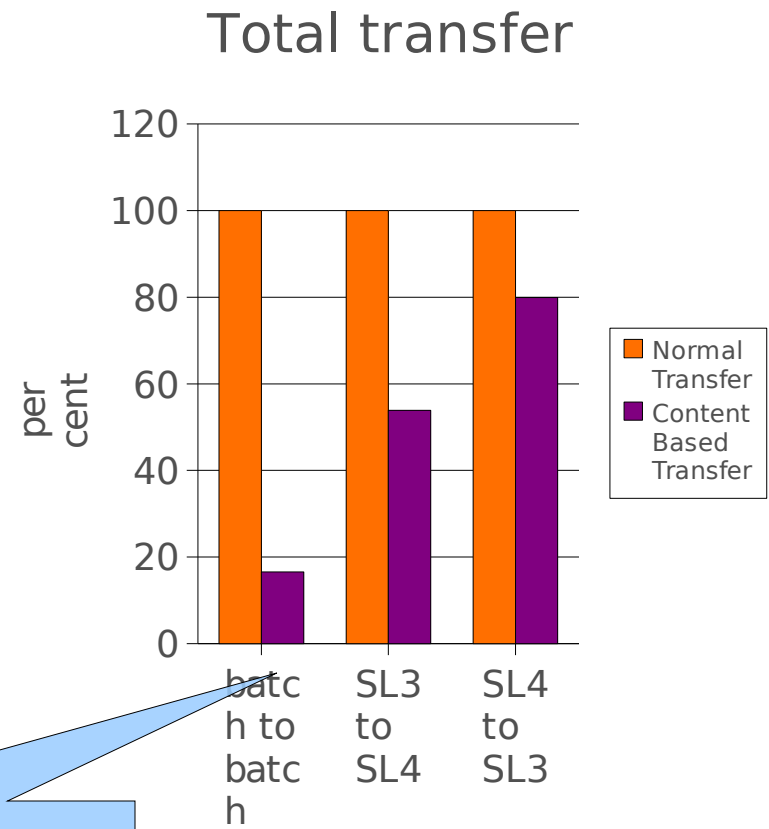| libZ | | libX | | | | | libY | | | |

Only 50 % of Image A needs to be transferred if Image B is already at the destination
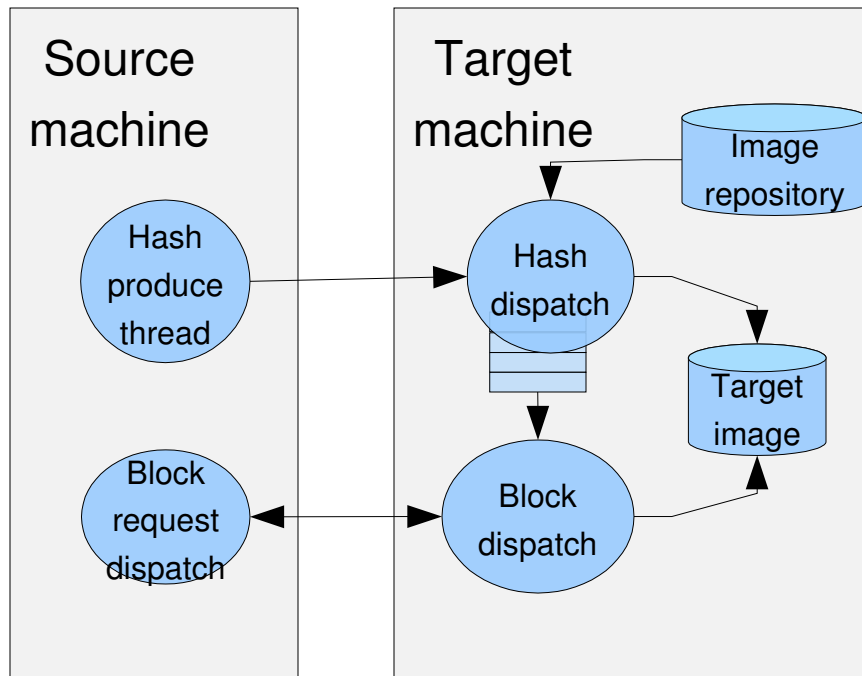
- Two typical batch machines (5.3 GB)
  - 84 % common blocks
- Scientific 3 (343 MB) and Scientific 4 (762 MB)
  - SL3 -> SL4
    - 48 % common blocks
  - SL4 -> SL3
    - 22 % common blocks

## Total transfer

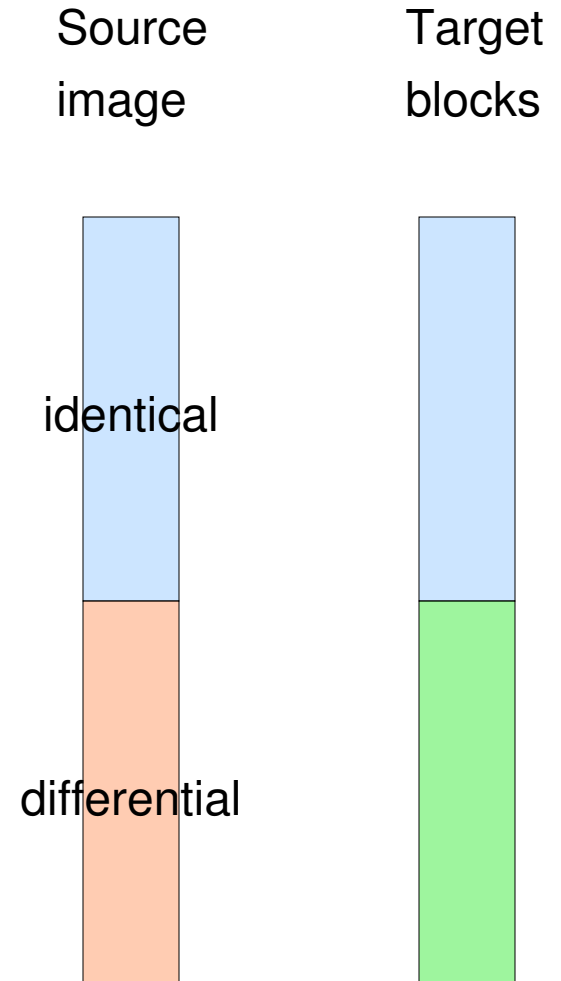Fraction of full image data needed to transfer, including hash table

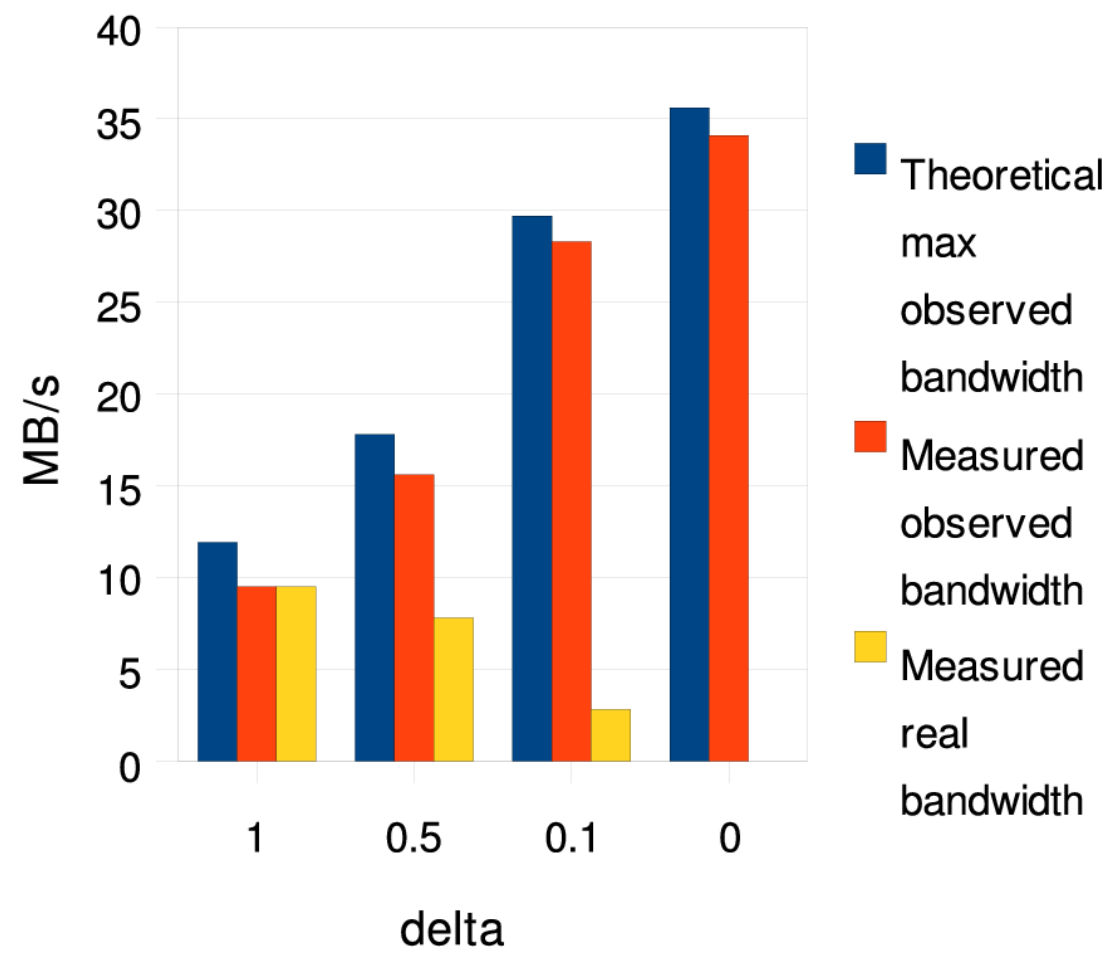- Generating hash tables for source file and target repository – linear cost

- Accessing hash tables
  - Java has a convenient constant-time hash table

- 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

# Implementation



- **Multi-threaded**

- **Hash calculation and data transfer pipelined**

- **Implemented in Java**

- A source image which to transfer consists of a *differential* and an *identical* part

- "Observed Bandwidth"
  - $|image|/time_{diff}$

- Theoretical maximum can be calculated

Source image

Target blocks

identical

differential

- Signed images
  - Signed by image author or image creation service
- Automated image testing
- Contextualization
  - One VM image should be able to run on any infrastructure, e.g. Amazon EC2, LCG, my laptop
- Open Virtual Machine Format
  - Open DMTF standard for VM images

- ## OS Farm

  - http://cern.ch/osfarm

- ## Content Based Transfer

  - http://hbjerke.web.cern.ch/hbjerke/cba/cba.xml