



# Clustered and distributed storage

with  
commodity hardware  
and open source software

**Phil Cryer**

BHL Developer, Systems Analyst  
BHL Europe Technical Board Meeting  
25-27 August 2010, NHM London

# BHL data, on our cluster



BHL's first cluster in Woods Hole

- **Hardware** - commodity servers
  - (6) six 4U sized cabinets
  - (24) twenty-four 1.5TB hard drives in each cabinet



# BHL data, on our cluster

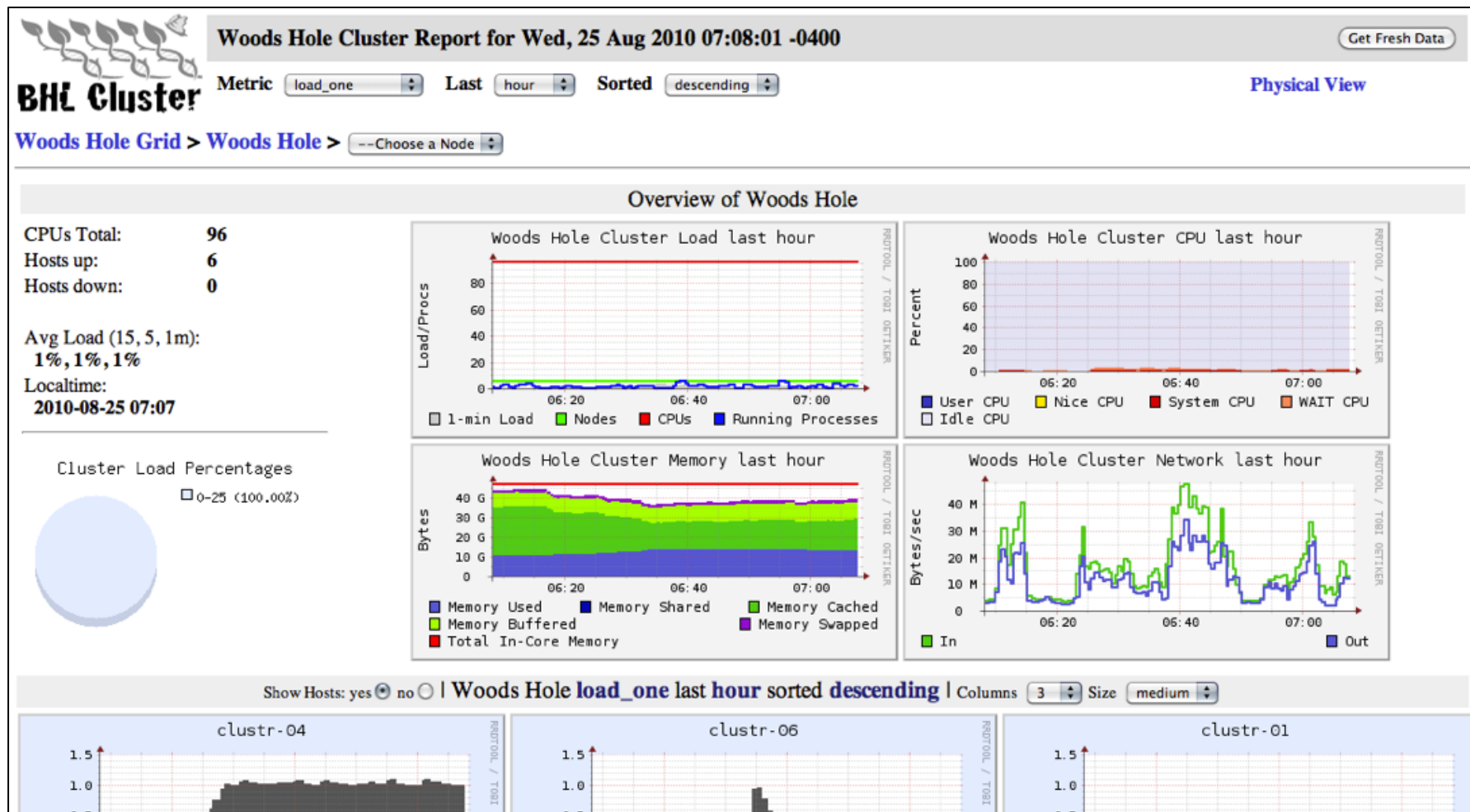


BHL's first cluster in Woods Hole

- **Hardware** - commodity servers
  - (6) six 4U sized cabinets
  - (24) twenty-four 1.5TB hard drives in each cabinet
- **Software** - open source software
  - operating system is Debian GNU/Linux (squeeze)
  - filesystem - ext4
    - supports filesystems up to 1 EB (1000 PB) and max file size of 16 TB
  - clustered file system - GlusterFS (3.0.4)
    - all drives run in a networked/RAID1 setup
    - all files are replicated and redundantly copied across the cluster
    - New: Acquia is using GlusterFS for their Drupal SaaS implementation
  - monitoring - Monit, Ganglia for alerts and reporting



# BHL data, on our cluster



<http://whbhl01.ubio.org/ganglia>



# BHL data, on our cluster



BHL's first cluster in Woods Hole

- **Hardware** - commodity servers
  - (6) six 4U sized cabinets
  - (24) twenty-four 1.5TB hard drives in each cabinet
- **Software** - open source software
  - operating system is Debian GNU/Linux (squeeze)
  - filesystem - ext4
    - supports filesystems up to 1 EB (1000 PB) and max file size of 16 TB
  - clustered file system - GlusterFS (3.0.4)
    - all drives run in a networked/RAID1 setup
    - all files are replicated and redundantly copied across the cluster
    - New: Acquia is using GlusterFS for their Drupal SaaS implementation
  - monitoring - Monit, Ganglia for alerts and reporting
- **Capacity** - cluster has 97TB of replicated/distributed storage
  - currently using 66TB of data for 78492 books
  - a full record for a book can be 24MB - 3GB



# Files from a record



```
# ls -lh /mnt/glusterfs/www/a/actasocietatisssc26suom
```

**total 649M**

```
-rwxr-xr-x 1 www-data www-data 19M 2009-07-10 01:55 actasocietatisssc26suom_abbyy.gz
-rwxr-xr-x 1 www-data www-data 28M 2009-07-10 06:53 actasocietatisssc26suom_bw.pdf
-rwxr-xr-x 1 www-data www-data 1.3K 2009-06-12 10:21 actasocietatisssc26suom_dc.xml
-rwxr-xr-x 1 www-data www-data 18M 2009-07-10 03:05 actasocietatisssc26suom.djvu
-rwxr-xr-x 1 www-data www-data 1.3M 2009-07-10 06:54 actasocietatisssc26suom_djvu.txt
-rwxr-xr-x 1 www-data www-data 14M 2009-07-10 02:08 actasocietatisssc26suom_djvu.xml
-rwxr-xr-x 1 www-data www-data 4.4K 2009-12-14 04:42 actasocietatisssc26suom_files.xml
-rwxr-xr-x 1 www-data www-data 20M 2009-07-09 18:57 actasocietatisssc26suom_flippy.zip
-rwxr-xr-x 1 www-data www-data 285K 2009-07-09 18:52 actasocietatisssc26suom.gif
-rwxr-xr-x 1 www-data www-data 193M 2009-07-09 18:51 actasocietatisssc26suom_jp2.zip
-rwxr-xr-x 1 www-data www-data 5.7K 2009-06-12 10:21 actasocietatisssc26suom_marc.xml
-rwxr-xr-x 1 www-data www-data 2.0K 2009-06-12 10:21 actasocietatisssc26suom_meta.mrc
-rwxr-xr-x 1 www-data www-data 416 2009-06-12 10:21 actasocietatisssc26suom_metasource.xml
-rwxr-xr-x 1 www-data www-data 2.2K 2009-12-01 12:20 actasocietatisssc26suom_meta.xml
-rwxr-xr-x 1 www-data www-data 279K 2009-12-14 04:42 actasocietatisssc26suom_names.xml
-rwxr-xr-x 1 www-data www-data 324M 2009-07-09 13:28 actasocietatisssc26suom_orig_jp2.tar
-rwxr-xr-x 1 www-data www-data 34M 2009-07-10 04:35 actasocietatisssc26suom.pdf
-rwxr-xr-x 1 www-data www-data 365K 2009-07-09 13:28 actasocietatisssc26suom_scandata.xml
```

# Initial file population



Populating a cluster with our data at the Internet Archive

- Looked at many options
  - ship a pre-populated server (Sun Thumper with 48TB capacity)
  - shipping individual external hard-drives
  - download the files on our own



# Initial file population



	191Mb	381Mb	572Mb	763Mb	954Mb			
clustr-02		<= ia311215.us.archive.org		85.7Mb	54.2Mb	41.7Mb		
clustr-02		<= ia311030.us.archive.org		24.6Mb	19.1Mb	13.4Mb		
clustr-02		<= ia310807.us.archive.org		22.3Mb	19.0Mb	15.4Mb		
clustr-02		<= ia360606.us.archive.org		18.1Mb	16.7Mb	16.5Mb		
clustr-02		<= ia350603.us.archive.org		17.9Mb	13.9Mb	13.6Mb		
clustr-02		<= ia331404.us.archive.org		14.3Mb	11.0Mb	12.8Mb		
clustr-02		<= ia331419.us.archive.org		9.52Mb	6.89Mb	7.21Mb		
clustr-02		<= ia301506.us.archive.org		9.44Mb	4.67Mb	3.11Mb		
clustr-02		<= ia301531.us.archive.org		2.67Mb	2.43Mb	2.66Mb		
clustr-02		<= vpnreserved2.mbl.edu		208b	208b	208b		
128.128.175.255		<= geldoc319.mbl.edu		0b	366b	203b		
clustr-02		<= dns1.mbl.edu		704b	246b	352b		
128.128.163.255		<= DHCP160103.mbl.edu		936b	187b	47b		
128.128.175.255		<= bpccolorphoto.mbl.edu		0b	183b	92b		
128.128.171.255		<= 128.128.171.146		0b	183b	46b		
255.255.255.255		<= 192.168.56.1		0b	54b	14b		
128.128.171.255		<= 128.128.168.27		0b	54b	14b		
255.255.255.255		<= 10.0.1.1		0b	0b	691b		
255.255.255.255		<= 10.0.2.1		0b	0b	461b		
128.128.175.255		<= 175wins134.mbl.edu		0b	0b	216b		
128.128.175.255		<= squalus.mbl.edu		0b	0b	163b		
255.255.255.255		<= *		0b	0b	138b		
128.128.163.255		<= 128.128.162.165		0b	0b	109b		
128.128.171.255		<= 128.128.170.216		0b	0b	52b		
128.128.171.255		<= 128.128.168.227		0b	0b	52b		
TX:	cumm:	37.0MB	peak:	3.39Mb	rates:	3.39Mb	2.43Mb	2.12Mb
RX:		2.20GB		205Mb		205Mb	148Mb	126Mb
TOTAL:		2.24GB		208Mb		208Mb	150Mb	128Mb



# Initial file population

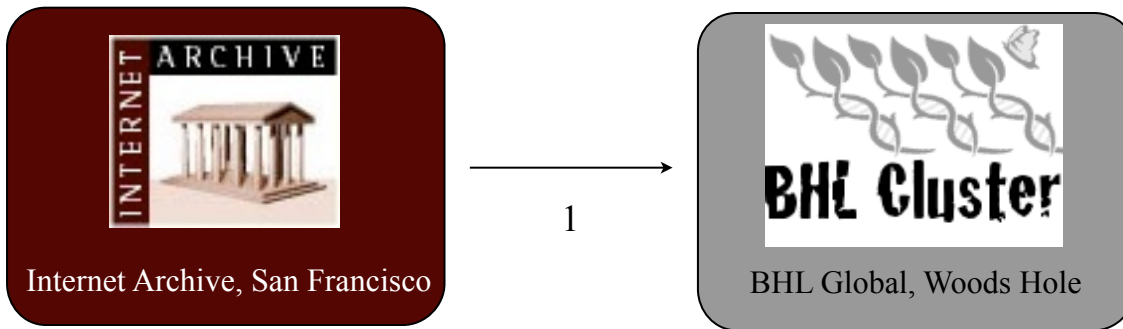


Populating a cluster with our data at the Internet Archive

- Looked at many options
  - ship a pre-populated server (Sun Thumper with 48TB capacity)
  - shipping individual external hard-drives
  - download the files on our own
- Path of least resistance, we wrote a script and used the Internet2 connection at the Marine Biology Laboratory (Woods Hole) to download directly to the first cluster
  - knew it would take forever to download (but it took longer)
  - needed space to download files (cluster buildout)
  - networking issues in Woods Hole (overloaded local router)
  - file verification (checksums that don't...)
- **Lessons learned** - would we do it again? Probably not.
- Current propagation method
  - initial distribution - mailing external drives (1, 5)
  - syncing of the changes for future content (smaller bites)



# Code: grabbyd



Automated process to continuously download the latest BHL data

- Uses **subversion** to get an updated list of new BHL content as IA identifiers  
<http://code.google.com/p/bhl-bits/source/browse/#svn/trunk/iaidentifiers>
- An enhanced version of the original download script to transfer the data
  - **grabbyd** - a script that parses the latest iaidentifiers list, determines the IDs of the new data and downloads the data to the cluster
  - Will provide detailed reporting with status pages and/or another method (webapp, email, RSS, XML, etc)

Code available (open sourced, BSD licensed):

[1] <http://code.google.com/p/bhl-bits/source/browse/trunk/utilities/grabby/grabbyd>

# Code: grabbyd + reporting



**Current transfer rate 10043.09 KB/s**

**Clustered data store size 66 of 97 Terabytes**

**RUN1** download completed at 15:13:26 NUM: 3136/9950 ID: anatomischeranze47anat Number of files: 18 Size of files: 463M total

**RUN2** download completed at 15:25:47 NUM: 570/10000 ID: blainesoutlineso00blai Number of files: 18 Size of files: 959M total

**RUN3** download completed at 14:38:35 NUM: 866/10000 ID: controlseriesbul1556mass Number of files: 18 Size of files: 1.4G total

cluster status as of Tue Aug 24 15:30:01 EDT 2010 - updated every 15 minutes

<http://cluster.biodiversitylibrary.org/>

# Replication|Replication



Why do we need replication?

- First BHL stored everything at the Internet Archive in San Francisco
  - no backup or safety net
  - limited in what we could do with, and serve, our data
- Now with our first BHL cluster, we gain
  - redundancy - will be able to serve from the cluster and fall back to IA if needed
  - analytics - the files are 'local' to parse through, discover new relationships
  - serving options - geo-location, eventually will be able to serve from closest server



# Replication|Replication



Why do we need replication?

- First BHL stored everything at the Internet Archive in San Francisco
  - no backup or safety net
  - limited in what we could do with, and serve, our data
- Now with our first BHL cluster, we gain
  - redundancy - will be able to serve from the cluster and fall back to IA if needed
  - analytics - the files are 'local' to parse through, discover new relationships
  - serving options - geo-location, eventually will be able to serve from closest server
- Next - share the data with everyone
  - Europe
  - Australia
  - China
  - etc...
- Provide safe harbor
  - lots of copies...





# Code: bhl-sync



## Open source **Dropbox** model

- uses and implements many open source projects
  - **inotify** - a subsystem within the Linux kernel that extends the filesystem to notice changes to the filesystem and report them to applications (in the kernel since 2.6.13 (2005))
  - **lsyncd** - an open source project that provides a wrapper into inotify
  - **OpenSSH** - secure file transfer
  - **rsync** - long term, proven syncing subsystem

# Code: bhl-sync



## Open source **Dropbox** model

- uses and implements many open source projects
  - **inotify** - a subsystem within the Linux kernel that extends the filesystem to notice changes to the filesystem and report them to applications (in the kernel since 2.6.13 (2005))
  - **lsyncd** - an open source project that provides a wrapper into inotify
  - **OpenSSH** - secure file transfer
  - **rsync** - long term, proven syncing subsystem

## What does bhl-sync do?

- runs lsyncd as a daemon that notices kernel events and kicks off rync over OpenSSH to mirror data to designated remote servers
- the only requirement on the remote system is a secure login for a normal user (using a key based OpenSSH) keeping the process neutral and not requiring any other specific technologies (OS, applications, filesystem) on the remote system (cross-platform)
- want to mirror BHL? it's now possible (you just need a lot of storage)

Code available (open sourced, BSD licensed):

<http://code.google.com/p/bhl-bits/source/browse/trunk/utilities/bhl-sync.sh>

# Code: bhl-sync + status



## bhl-sync

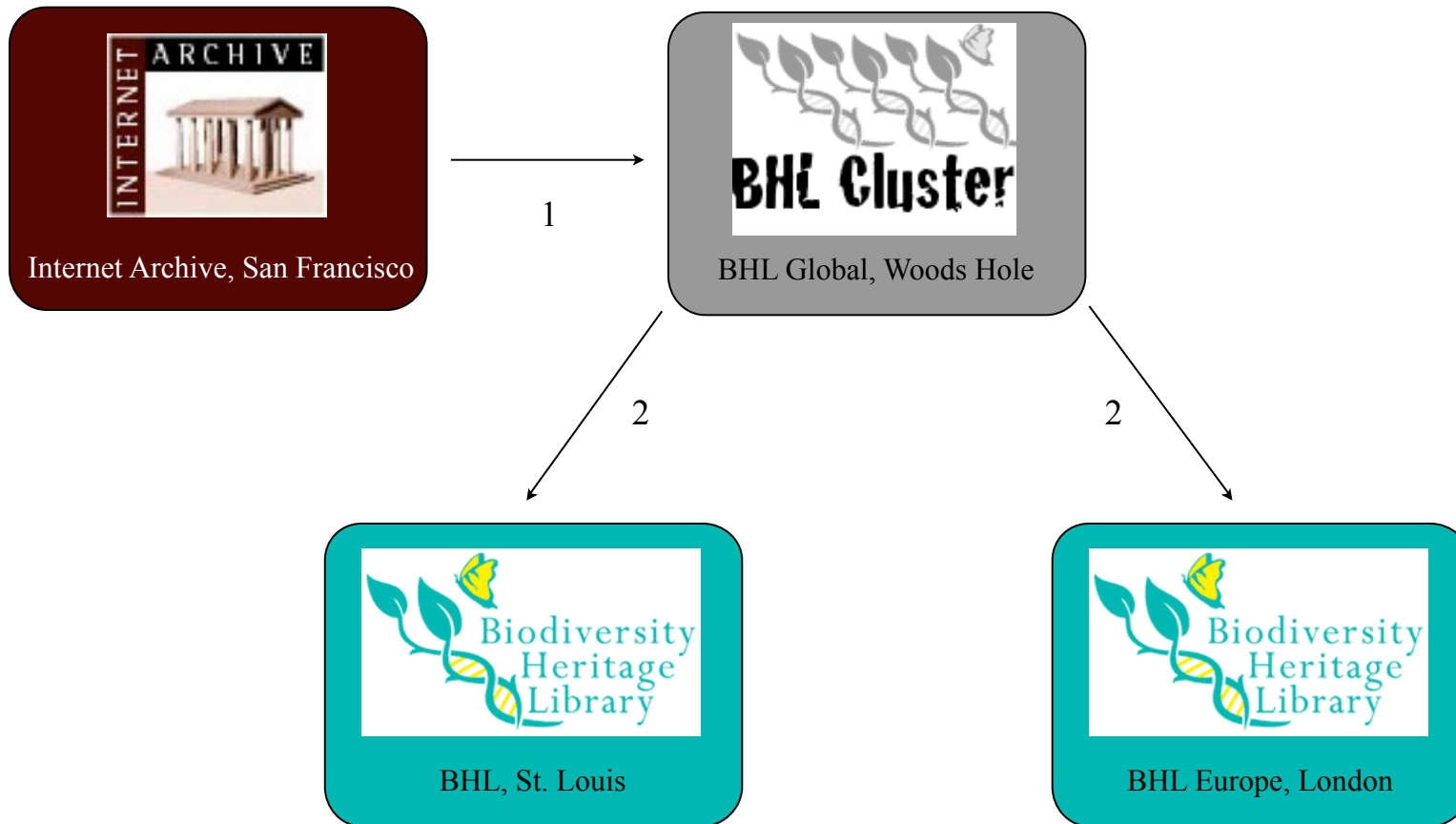
---

St. Louis ← Woods Hole → London  
**Synced**                      **24476 files**                      **Synced**

updated: Fri Aug 20 13:02:18 EDT 2010 (updated hourly)

<http://bit.ly/09-bhl-sync>

# BHL content distribution

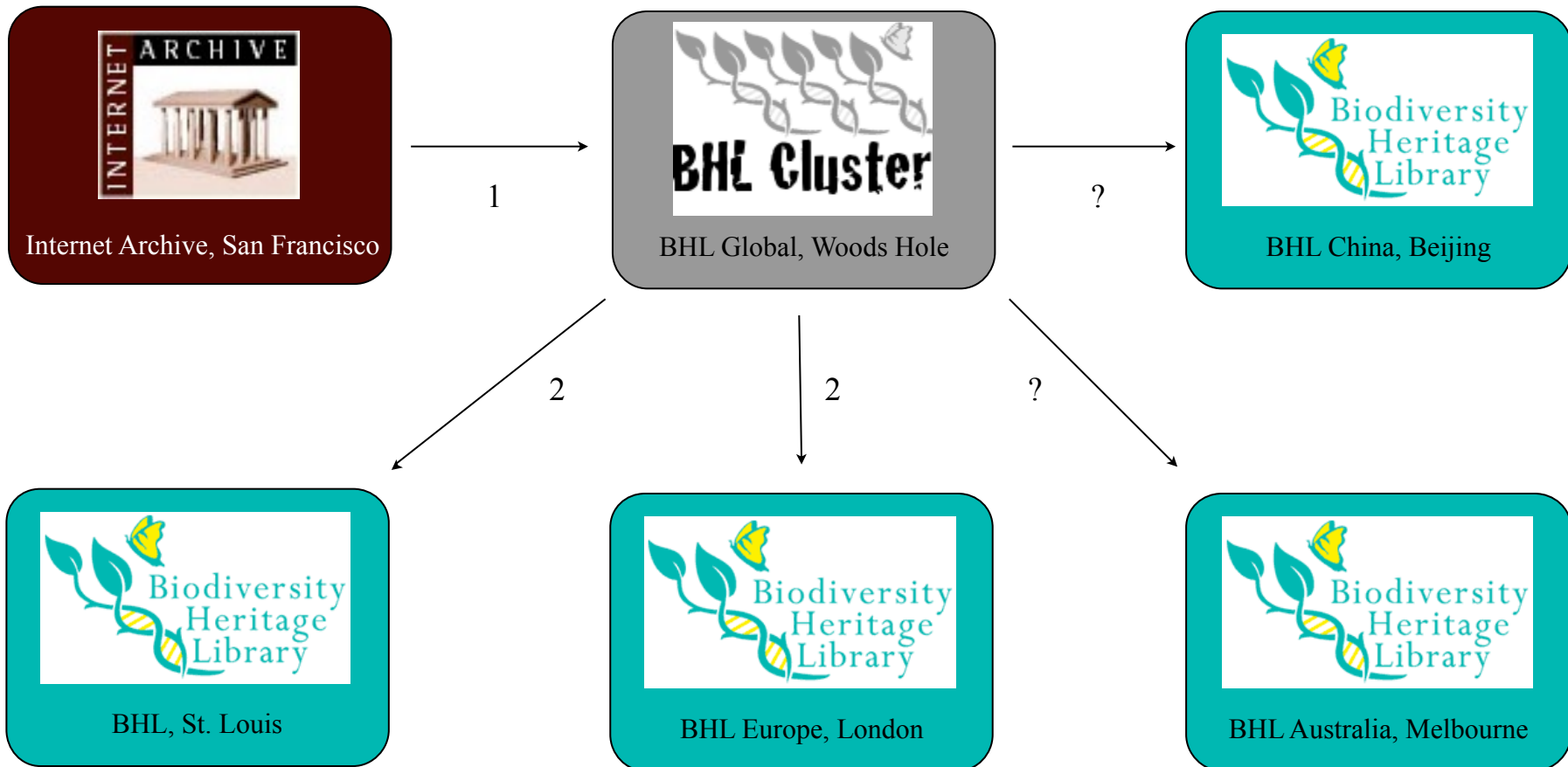


Code available (open sourced, BSD licensed):

[1] <http://code.google.com/p/bhl-bits/source/browse/trunk/utilities/grabby/grabbyd>

[2] <http://code.google.com/p/bhl-bits/source/browse/trunk/utilities/bhl-sync.sh>

# BHL content distribution



Code available (open sourced, BSD licensed):

[1] <http://code.google.com/p/bhl-bits/source/browse/trunk/utilities/grabby/grabbyd>

[2] <http://code.google.com/p/bhl-bits/source/browse/trunk/utilities/bhl-sync.sh>



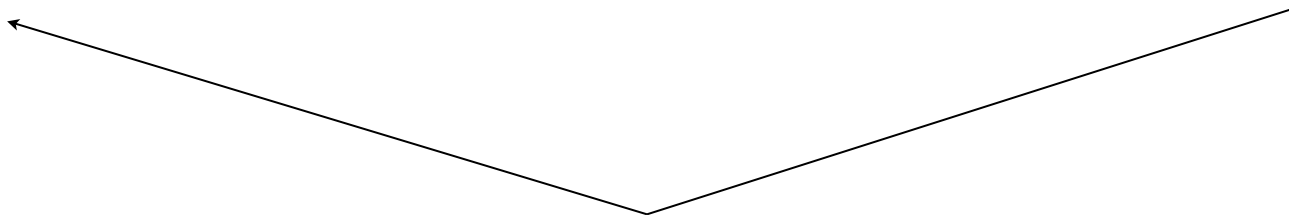
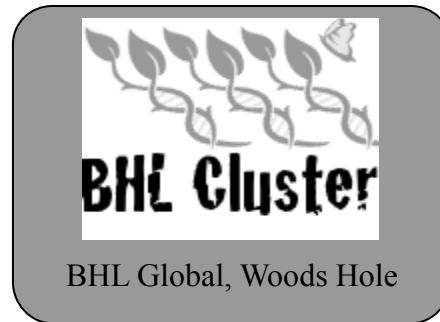
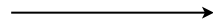
# Other replication challenges



- Deleting content - "going dark"
  - this can be data that is removed from search indexes, but still retrievable via URI
  - or deleted data not available (requires a separate sync process)
- New content coming in from other sources
  - Localization of content - maybe it all can't be shared?
  - National nodes consideration

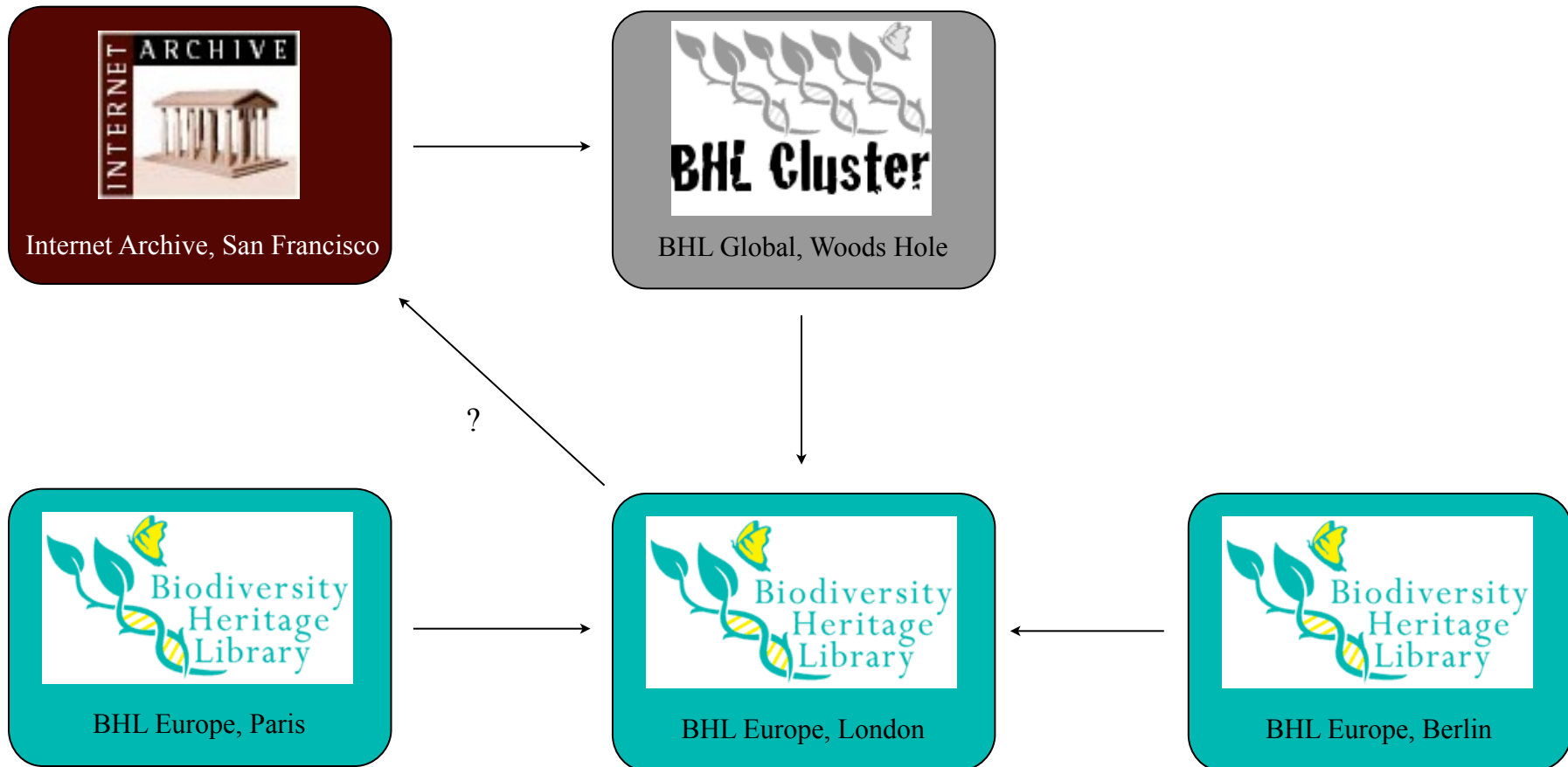


# BHL content + local data



Content sourced from China, scanned by  
Internet Archive, replicated into BHL Global

# BHL content + regional data



Content sourced from BHL Europe partners may, or may not, be passed back to Internet Archive and BHL Global

# Fedora-commons integration



Integrated digital repository-centered platform

- Enables storage, access and management of virtually any kind of digital content
- can be a base for software developers to build tools and front ends on for sharing, reuse and displaying data online
- Is free, community supported, open source software



# Fedora-commons integration



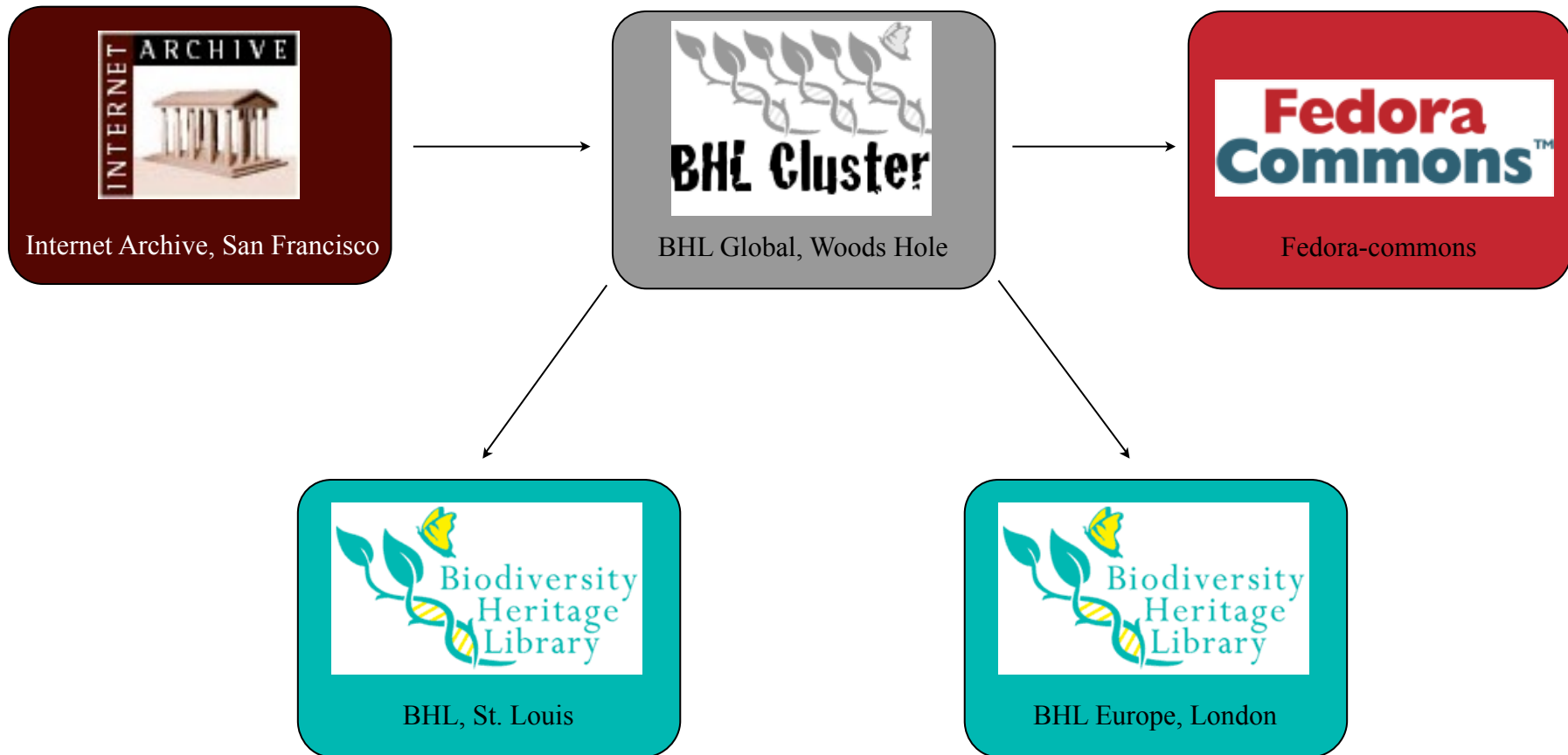
Integrated digital repository-centered platform

- Enables storage, access and management of virtually any kind of digital content
- can be a base for software developers to build tools and front ends on for sharing, reuse and displaying data online
- Is free, community supported, open source software
- Creates and maintains a persistent, stable, digital archive
  - provides backup, redundancy and disaster recovery
  - complements (doesn't replace or put any demands upon) existing architecture by incorporating open standards
  - stores data in a neutral manner, allowing for an independent disaster recovery option
  - shares data via OAI, REST based interface

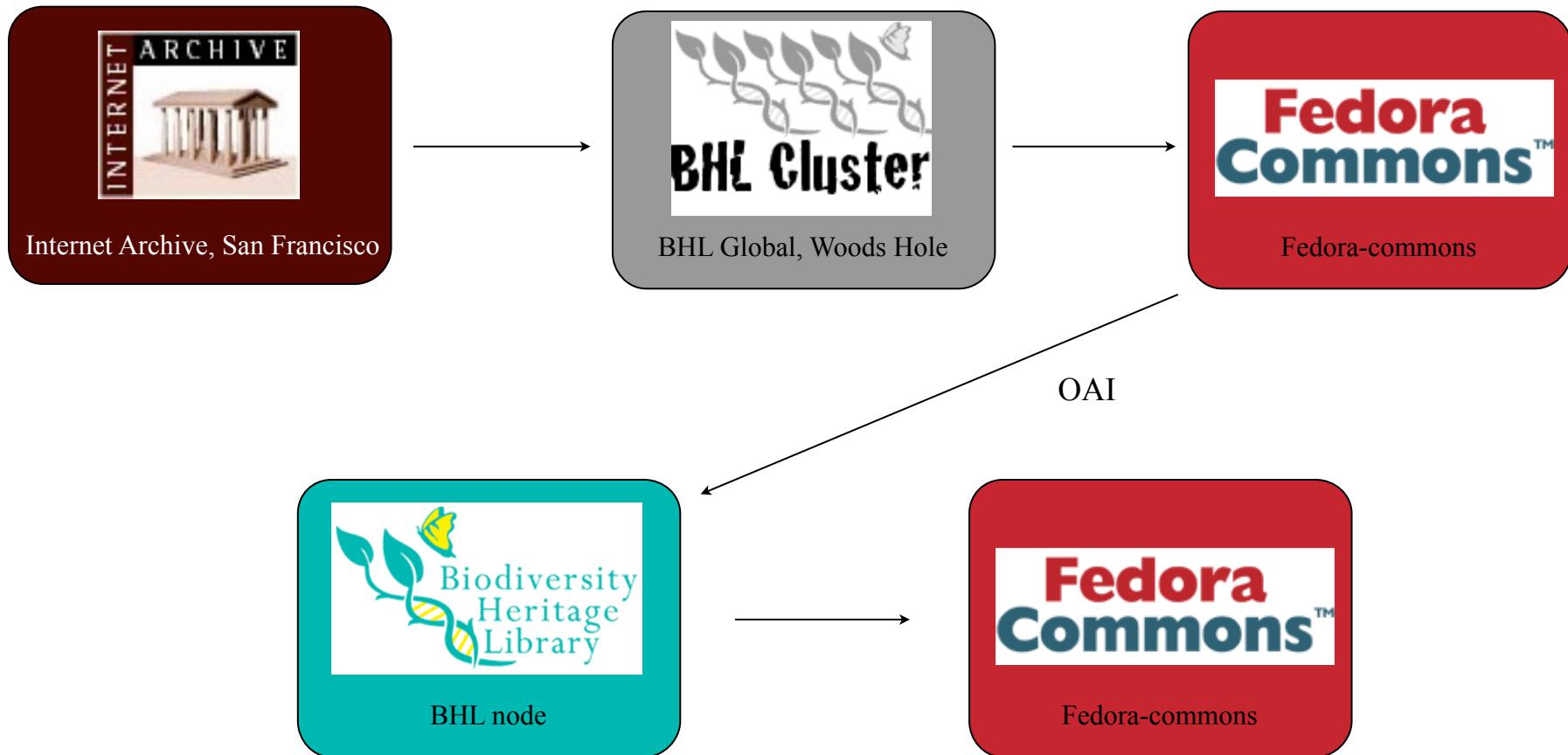




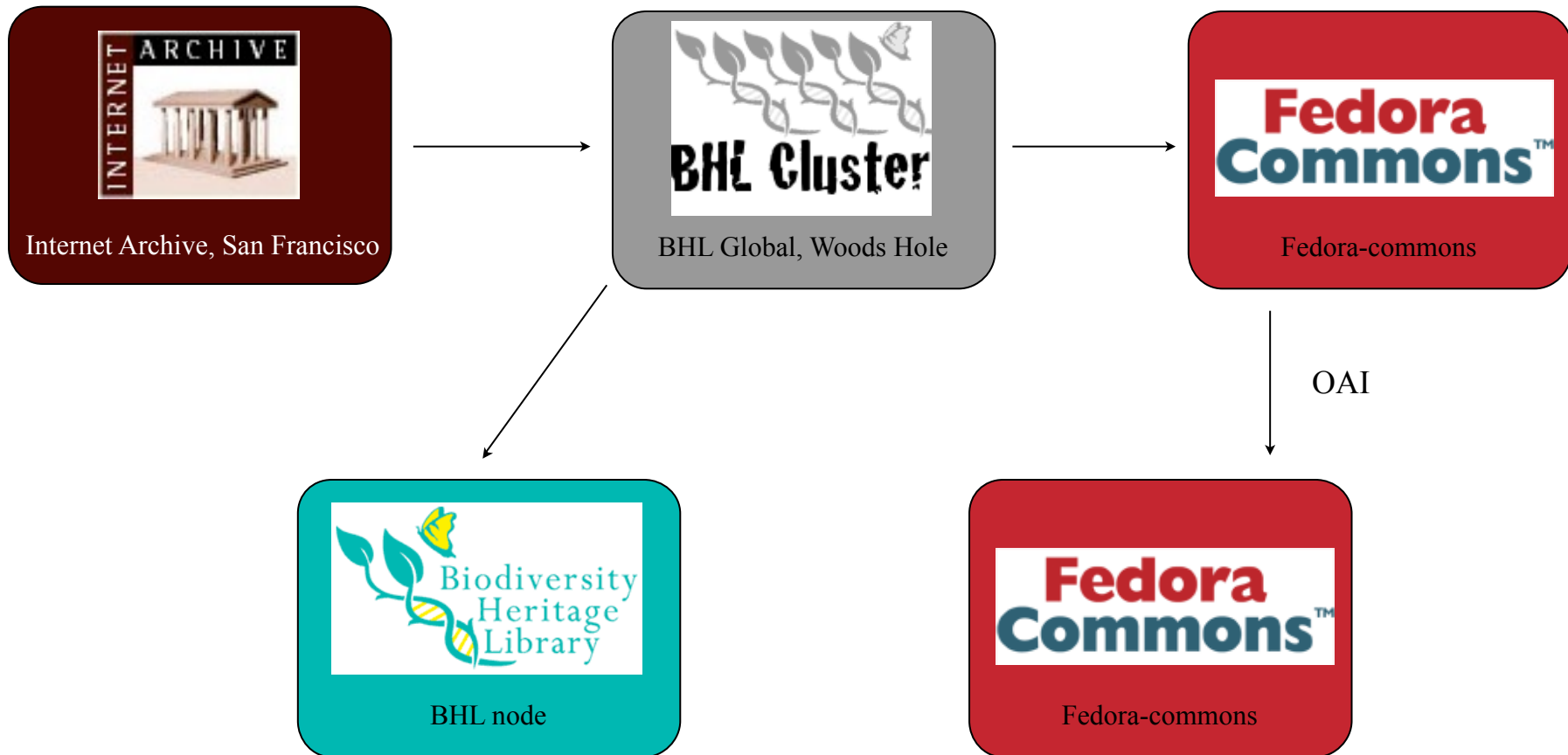
# BHL content distribution



# BHL content distribution



# BHL content distribution



# Thanks + questions



Thanks to Adrian Smales,  
Chris Sleep (NMH), Chris Freeland,  
Tom Garnett (BHL) and Cathy Norton,  
Anthony Goddard, Woods Hole  
networking admins (MBL) for their  
work and support of this project.



email [phil.cryer@mobot.org](mailto:phil.cryer@mobot.org)  
skype [phil.cryer](#)  
twitter [@fak3r](#)

slides available on [slideshare](#)