

# CernVM on the Cloud

Predrag Buncic, CERN/PH-SFT

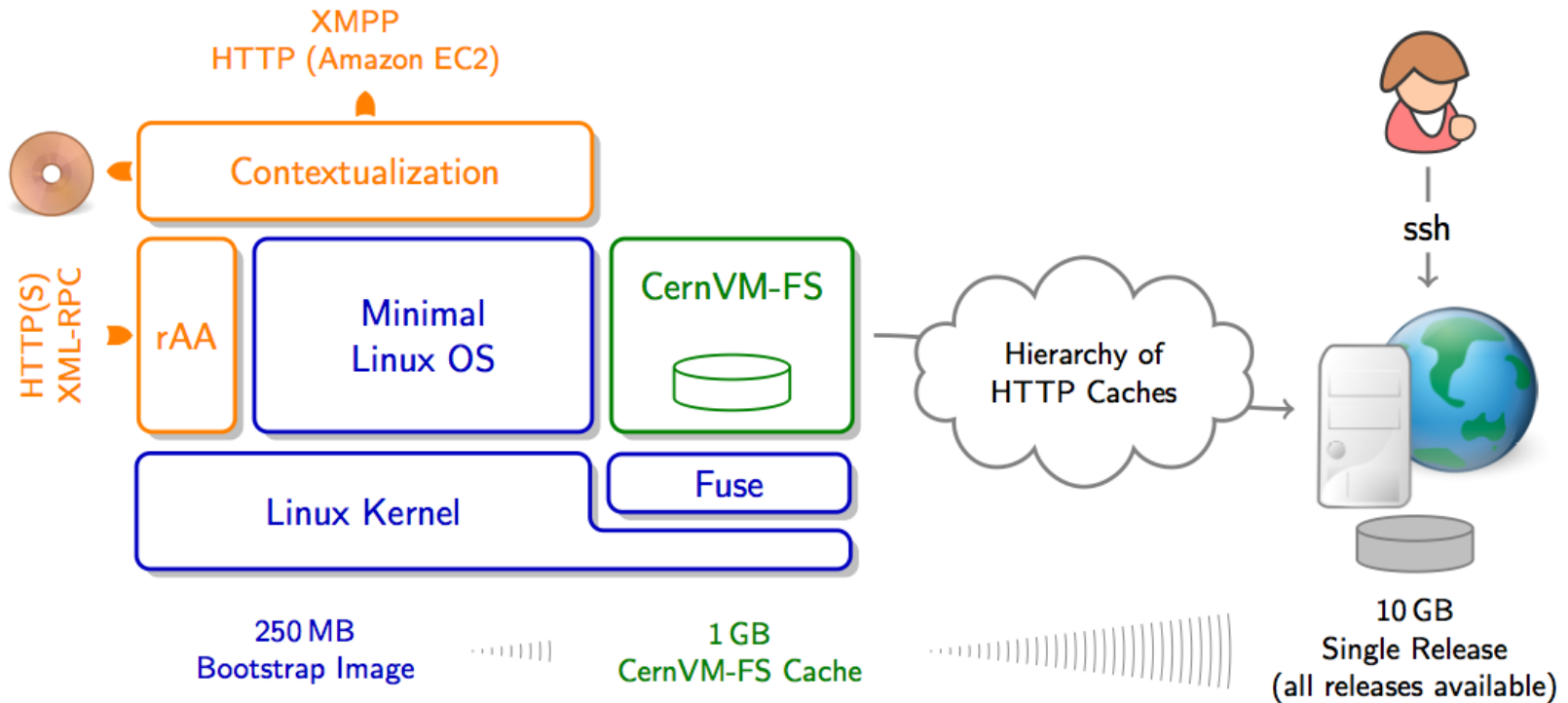


Is CernVM suitable for deployment on Grid and Cloud infrastructure?

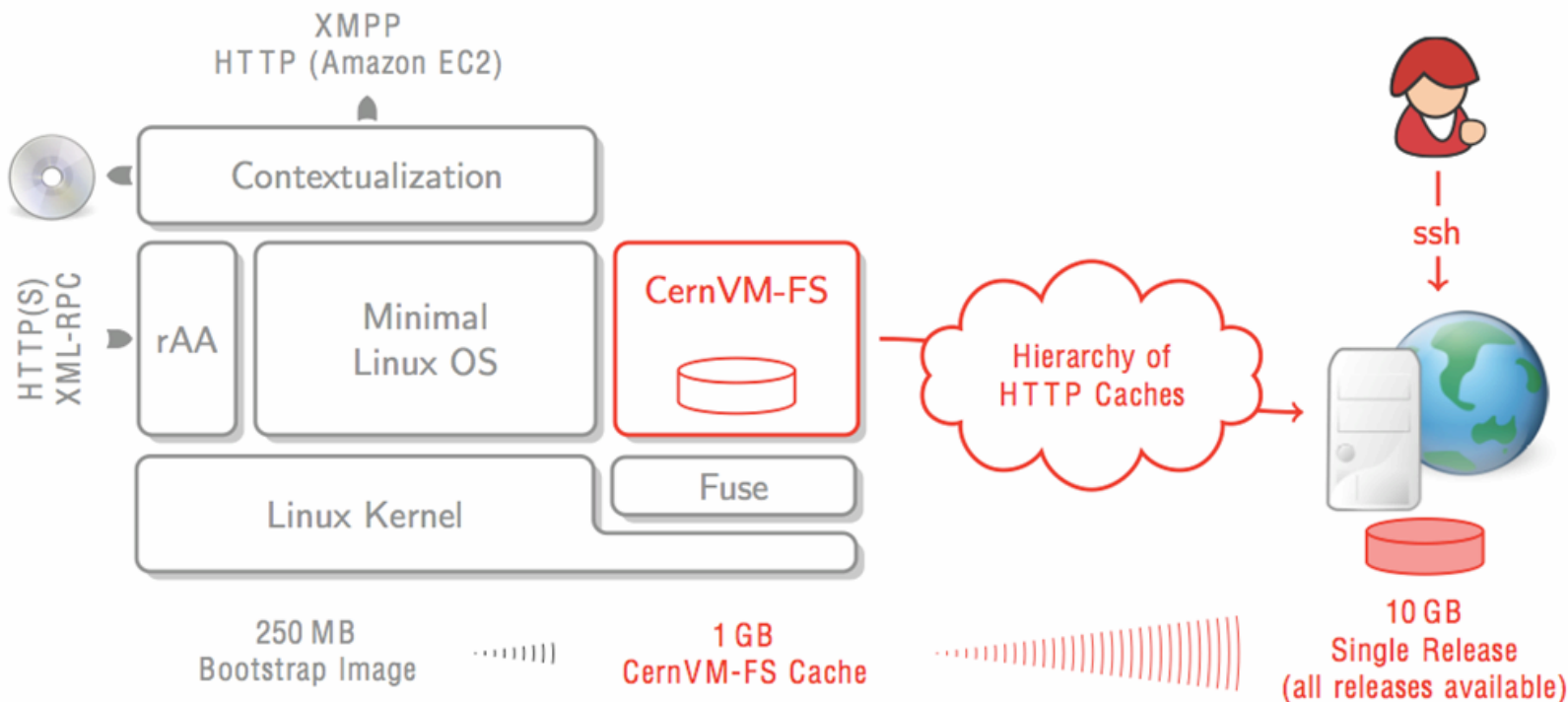
What are the benefits of going CernVM way comparing to more traditional<sup>1)</sup> approach to batch node virtualization?

1) Traditional approach:

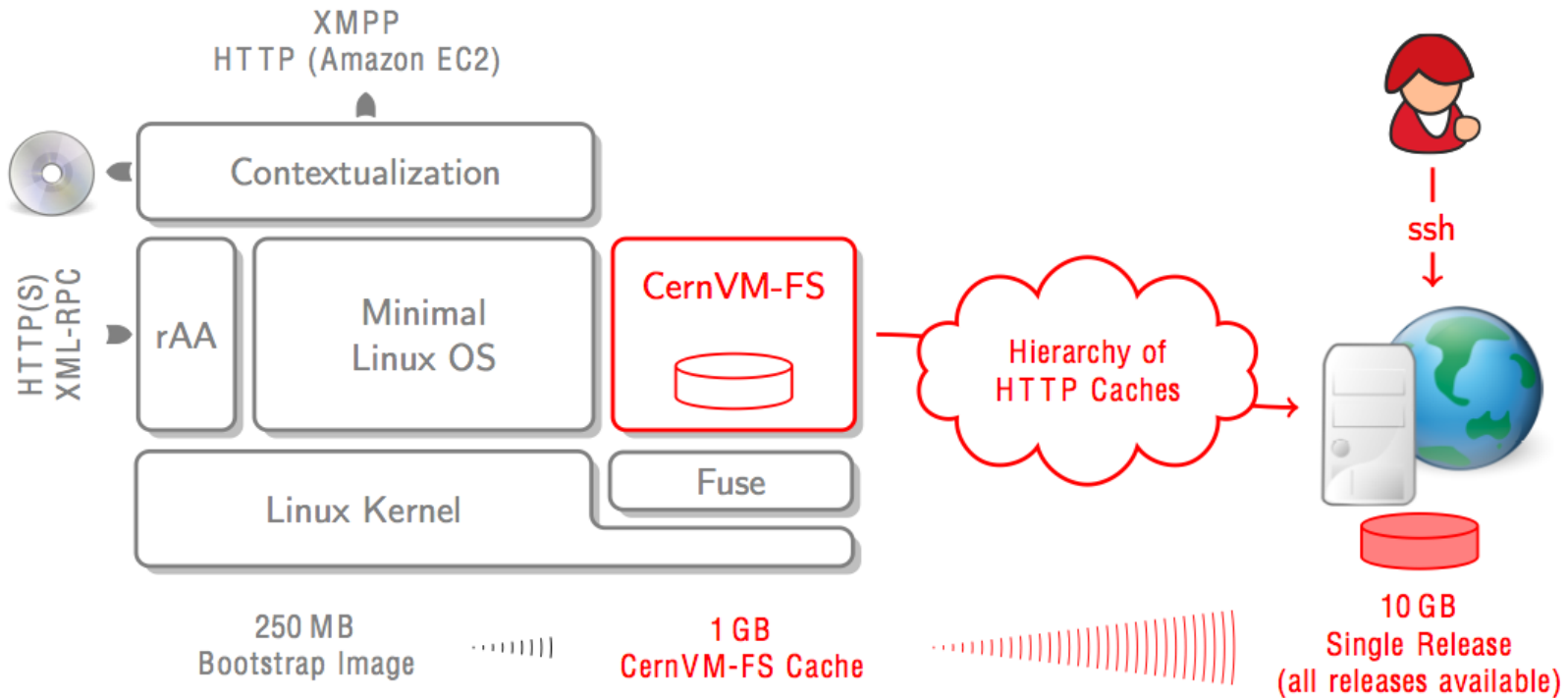
- Take “standard” batch node [2GB] and add experiment software [10GB] and generate VM image. Have experiment and security team certify the image, deploy it to all sites and worker nodes. Repeat this procedure 1-2 times per week and per experiment.



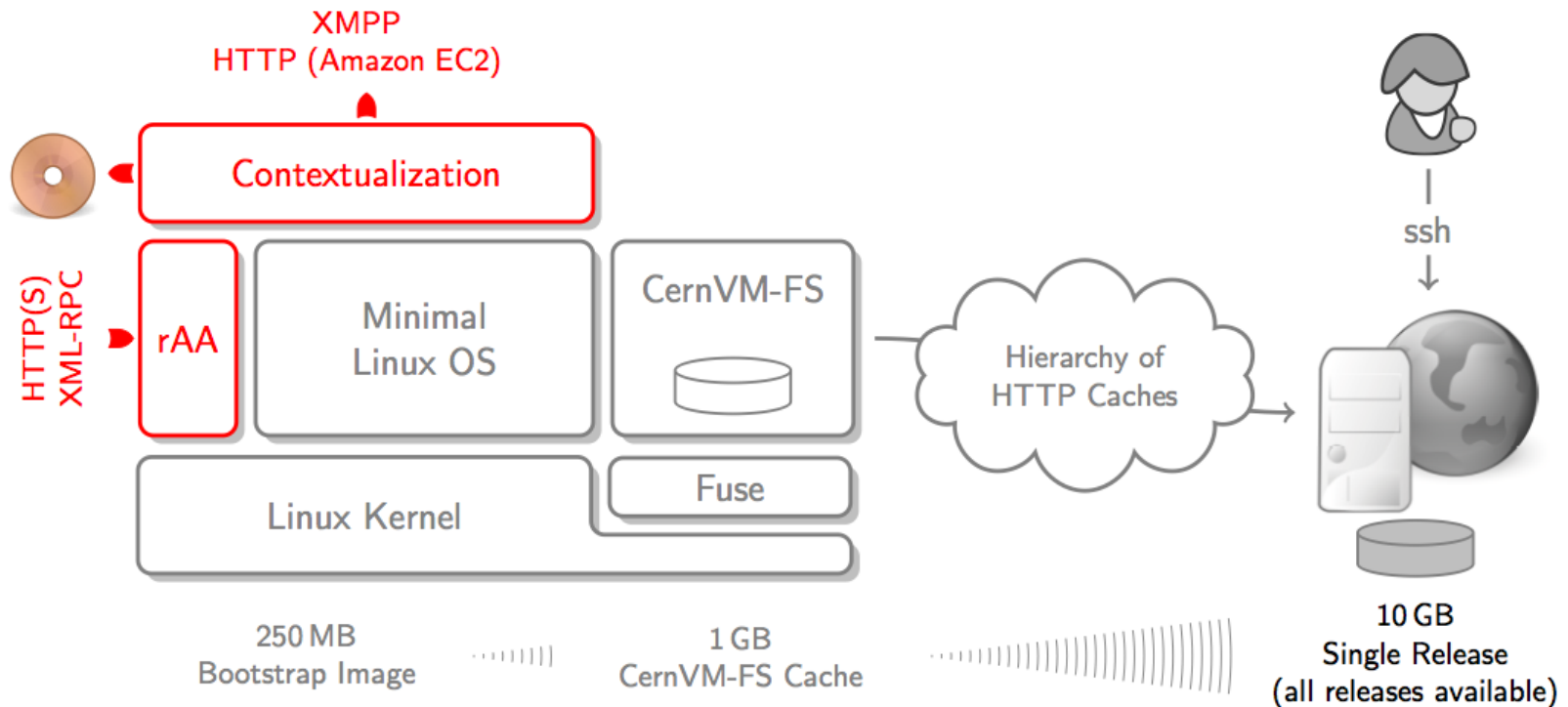
1. Minimal Linux OS (SL5)
2. CernVM-FS - HTTP network file system optimized for just in time delivery of experiment software
3. Flexible configuration and contextualization mechanism based on public Cloud API



- Just enough OS to run LHC applications
- Built using commercial tool (rBuilder by rPath)
  - Top-down approach - starting from application and automatically discovering dependencies
- Small images (250MB), easy to move around



- Experiment software is changing frequently and we want to avoid need to frequently update, certify and redistribute VM images with every release
- Only a small fraction of software release is really used
- Demonstrated scalability and reliability
- Now being deployed on across all Grid sites as the channel for software distributions



- There are several ways to contextualize CernVM
  - ✓ Web UI (for individual user)
  - ✓ CernVM Contextualization Agent
  - ✓ Hepix CDROM method
  - ☞ Amazon EC2 API user\_data method

## Syntax

```
ec2-run-instances ami_id [-n instance_count] [-g group [-g group ...]] [-k keypair]  
[-d user_data |-f user_data_file] [--addressing addressing_type] [--instance-type  
instance_type] [--availability-zone zone] [--kernel kernel_id] [--ramdisk ramdisk_id] [--  
block-device-mapping block_device_mapping] [--monitor] [--disable-api-termination] [--  
instance-initiated-shutdown-behavior behavior] [--placement-group placement-group] [--  
tenancy tenancy] [--subnet subnet] [--private-ip-address ip_address] [--client-token token]
```

***user\_data* specifies Base64-encoded MIME user data to be made available to the instance(s) in this reservation.**

Type: String  
Default: None



- Basic principles:
  - Owner of CernVM instance can contextualize and configure it to run arbitrary service as unprivileged user
  - Site can use HEPIX method to inject monitoring and accounting hooks w/o functionally modifying the image
- The contextualization is based on rPath *amiconfig* package extended with CernVM plugin
  - This tool will execute on boot time (before network services are available), parse user data and look for python style configuration blocks.
  - If match is found the corresponding plugin will process the options and execute configuration steps if needed.
- For more info on CernVM contextualization using EC2 API, see: <https://cernvm.cern.ch/project/trac/cernvm/wiki/EC2Contextualization>

```
[cernvm]
# list of ',' separated organizations/experiments (lowercase)
organisations = cms
# list of ',' separated repositories (lowercase)
repositories = cms,grid
# list of ',' separated user accounts to create <user:group:[password]>
users = cms:cms:
# CVMFS HTTP proxy
proxy = http://<host>:<port>;DIRECT
# install extra conary group
group_profile = group-cms
# script to be executed as given user: <user>:/path/to/script.sh
contextualization_command = cms:/path/to/script.sh
# list of ',' separated services to start
services = <list>
# extra environment variables to define
environment = CMS_SITECONFIG=EC2,CMS_ROOT=/opt/cms
```

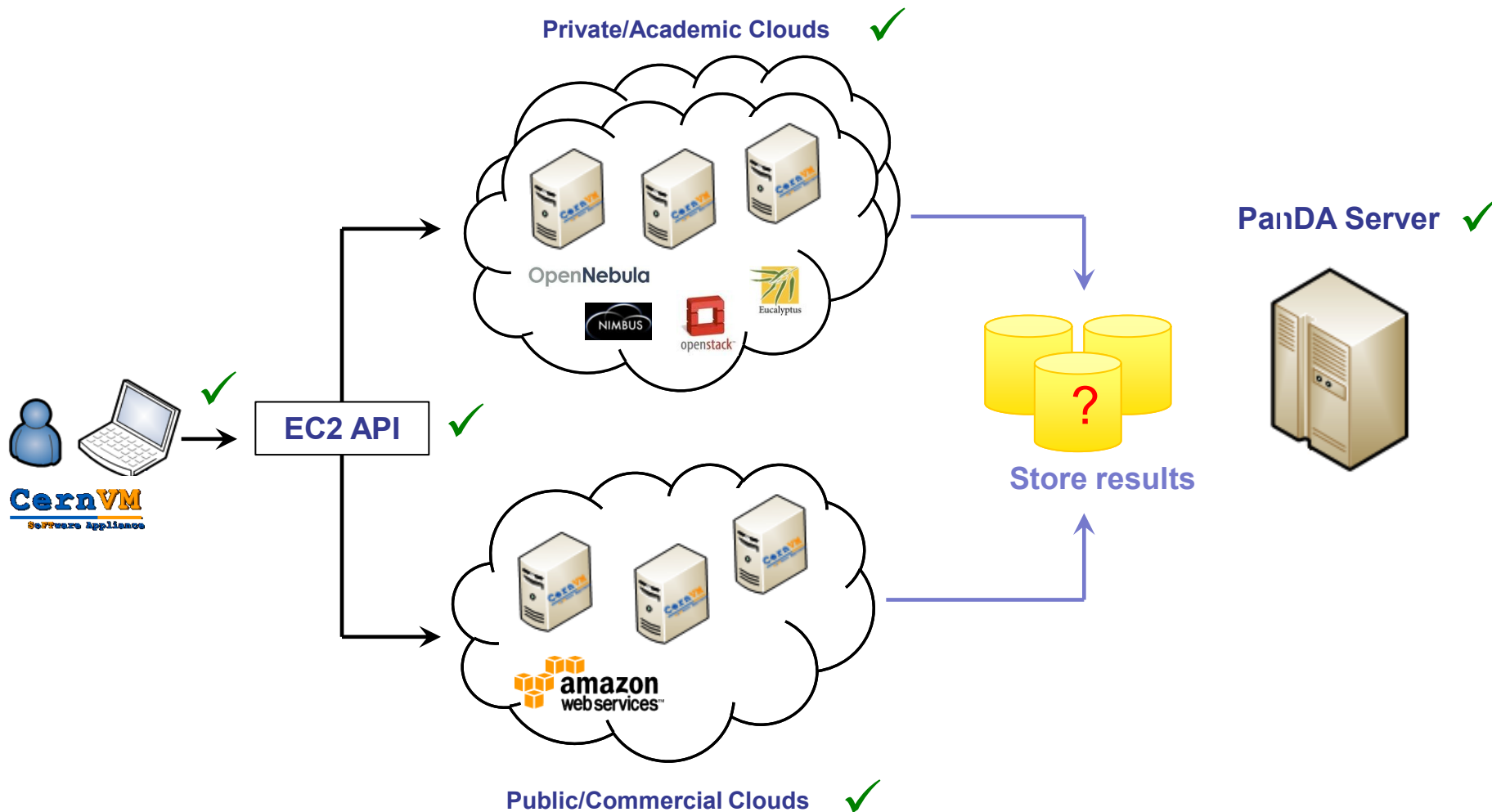
- **A tool to simplify interaction with multiple clouds supporting minimal EC2 API**
  - Aims to simplify instantiation of contextualized CernVM images using EC2 user\_data mechanism
- **Supports**
  - user\_data templates, bulk submission
  - delegation of proxy certificate in user\_data
- **Example:**

```
[pbuncic@localhost ~]$ cat panda-wn.tpl
[cernvm]
organisations = atlas
repositories  = atlas,grid,atlas-condb,sft
users = panda:atlas:
eos-server = eosatlas.cern.ch
contextualization_command =
panda: '/cvmfs/sft.cern.ch/lcg/external/experimental/panda-pilot/runPanda -m
/eos/atlas/vm/copilot -t 1'
[[IF X509_CERT]]x509-cert = [[X509_CERT]] [[ENDIF]]
```



Now we have CernVM OS, FS, Ctx, API...

What's next?

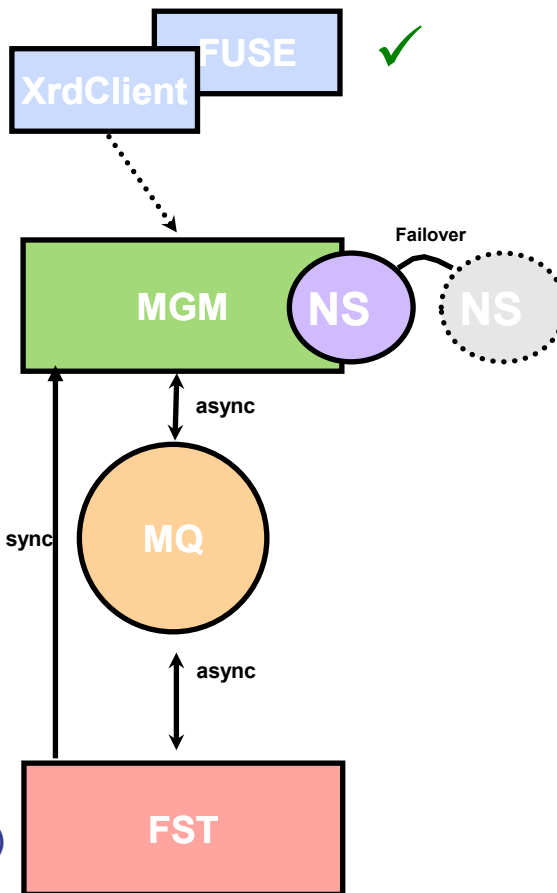


Clients  
XROOT client & FUSE  
KRB5 + X509 authenticated

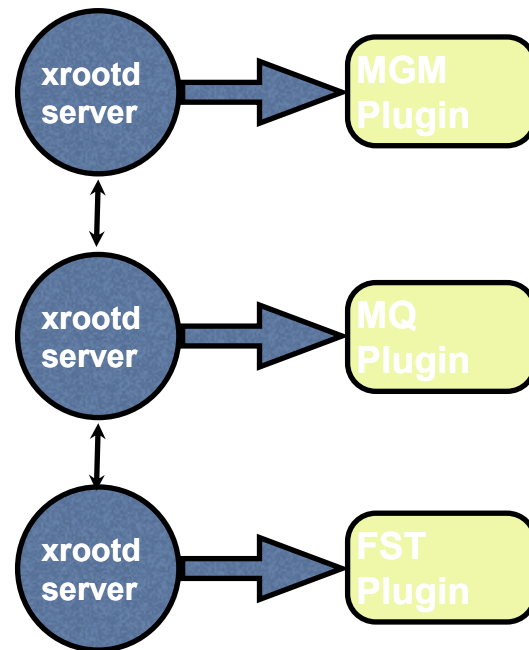
Management Server  
Pluggable Namespace, Quota, ACLs, HA  
Strong Authentication  
Capability Engine  
File Placement  
File Location

Message Queue  
Service State Messages  
File Transaction Reports

File Storage  
File & File Meta Data Store  
Capability Authorization  
File & Block Disk Error Detection (Scrubbing)  
Layout Plugins [ currently Raid-1(n) ]



Implemented as plugins in xrootd



Andreas Peters IT/DM

```
[root@vmbsq090700 ~]# df /eos
Filesystem      1K-blocks      Used Available Use% Mounted on
eos             3413376370848  23790724 3413352580124   1%
```

## Listing available images:

```
[pbuncic@localhost ~]$ cvm --region CERN -H ls -i
```

AMI	LOCATION	STATE	VISIBILITY	ARCH	TYPE
ami-00000008	hepix_sl155_x86_64_kvm	available	Public	i386	machine
ami-00000002	cernvm232_head_slc5_x86_64_kvm	available	Public	i386	machine
ami-00000004	cernvm231_slc5_x86_64_kvm	available	Public	i386	machine
ami-00000003	cernvm232_batch_slc5_x86_64_kvm	available	Public	i386	machine
ami-00000010	lxdev_slc6_quattor_slc6_x86_64_kvm	available	Public	i386	machine

## Support for multiple regions:

```
[pbuncic@localhost ~]$ cvm --region EC2 -H ls -i ami-5c3ec235
```

AMI	LOCATION	STATE	VISIBILITY	ARCH	TYPE
ami-5c3ec235	download.cernvm.cern.ch.s3.amazonaws.com/cernvm-2.3.1-x86_64_1899.img.manifest.xml	available	Public	x86_64	machine

## Proxy certificate

```
[pbuncic@localhost ~]$ lcg grid-proxy-init
Your identity: /DC=ch/DC=cern/OU=computers/CN=pilot/copilot.cern.ch
Creating proxy ..... Done
Your proxy is valid until: Thu May 19 06:43:50 2011
```

```
[pbuncic@localhost ~]$ lcg grid-proxy-info
subject   : /DC=ch/DC=cern/OU=computers/CN=pilot/copilot.cern.ch/CN=1766683191
issuer    : /DC=ch/DC=cern/OU=computers/CN=pilot/copilot.cern.ch
identity  : /DC=ch/DC=cern/OU=computers/CN=pilot/copilot.cern.ch
type      : Proxy draft (pre-RFC) compliant impersonation proxy
strength  : 512 bits
path      : /tmp/x509up_u500
timeleft  : 11:59:56
```

- **This proxy certificate is time limited and authorized only to**
  - Request jobs from dedicated PanDA queue
  - Write (but not delete) files in a given EOS directory



## Starting contextualized CernVM images on IxCloud:

```
[pbuncic@localhost ~]$ cvm --region CERN run ami-00000003 --proxy --template panda-wn:1  
r-47a5402e predrag default i-195 ami-00000003 128.142.192.62 128.142.192.62 pending  
default 0 m1.small 1970-01-01T01:00:00+01:00 default
```

```
[pbuncic@localhost ~]$ cvm --region CERN -H ls
```

ID	RID	OWNER	GROUP	DNS	STATE	KEY	TYPE
i-195	default	predrag	default	128.142.192.62	running	default	m1.small

```
[pbuncic@localhost ~]$ cvm --region CERN -H ls
```

ID	RID	OWNER	GROUP	DNS	STATE	KEY	TYPE
i-195	default	predrag	default	128.142.192.62	running	default	m1.small
i-196	default	predrag	default	128.142.192.63	running	default	m1.small
i-197	default	predrag	default	128.142.192.64	running	default	m1.small
i-198	default	predrag	default	128.142.192.65	running	default	m1.small
i-199	default	predrag	default	128.142.192.66	pending	default	m1.small
i-200	default	predrag	default	128.142.192.67	pending	default	m1.small
i-201	default	predrag	default	128.142.192.52	pending	default	m1.small
i-202	default	predrag	default	128.142.192.53	pending	default	m1.small
i-203	default	predrag	default	128.142.192.54	pending	default	m1.small
i-204	default	predrag	default	128.142.192.55	pending	default	m1.small
i-205	default	predrag	default	128.142.192.56	pending	default	m1.small
i-206	default	predrag	default	128.142.192.57	pending	default	m1.small

## Starting more contextualized CernVM images on EC2:

```
[pbuncic@localhost ~]$ cvm run ami-5c3ec235 -g default -t m1.large --kernel aki-9800e5f1 -
-key ami --proxy --template panda-wn:10
r-ad962dc1 392941794136 default i-f3b04a9d ami-5c3ec235 pending ami 0 m1.large
r-ad962dc1 392941794136 default i-f1b04a9f ami-5c3ec235 pending ami 1 m1.large
r-ad962dc1 392941794136 default i-cfb04aa1 ami-5c3ec235 pending ami 2 m1.large
r-ad962dc1 392941794136 default i-cdb04aa3 ami-5c3ec235 pending ami 3 m1.large
....
```

```
[pbuncic@localhost ~]$ cvm --region EC2 -H ls
```

ID	RID	OWNER	GROUP	DNS	STATE	KEY	TYPE
i-f3b04a9d	r-ad962dc1	392941794136	default	ec2-50-16-144-41	running	ami	m1.large
i-f1b04a9f	r-ad962dc1	392941794136	default	ec2-75-101-214-247	running	ami	m1.large
i-cfb04aa1	r-ad962dc1	392941794136	default	ec2-184-72-183-26	running	ami	m1.large
i-cdb04aa3	r-ad962dc1	392941794136	default	ec2-184-73-56-72	running	ami	m1.large
i-cbb04aa5	r-ad962dc1	392941794136	default	ec2-50-16-32-51	running	ami	m1.large
i-c9b04aa7	r-ad962dc1	392941794136	default	ec2-75-101-184-46	running	ami	m1.large
i-c7b04aa9	r-ad962dc1	392941794136	default	ec2-50-19-38-225	running	ami	m1.large
i-c5b04aab	r-ad962dc1	392941794136	default	ec2-50-16-105-241	running	ami	m1.large
i-c3b04aad	r-ad962dc1	392941794136	default	ec2-174-129-86-61	running	ami	m1.large
i-c1b04aaf	r-ad962dc1	392941794136	default	ec2-50-19-9-47	running	ami	m1.large

Jobs:								
<u>PandaID, Owner, Working group</u>	<u>Job</u>	<u>Status</u>	<u>Created</u>	<u>Time to start</u>	<u>Duration</u>	<u>Ended/Modified</u>	<u>Cloud/Site, Type</u>	<u>Priority</u>
<a href="#">1237433059</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:47	1 day, 15:31:21	0:00:34	05-18 13:18	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.35311e2f-8899-483e-9830-71095077949a</a>								
<a href="#">1237433058</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:47	1 day, 15:28:03	0:03:52	05-18 13:15	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.8bf86120-e688-494f-b20d-9a5d8bf830a0</a>								
<a href="#">1237433057</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:47	1 day, 15:26:30	0:05:26	05-18 13:13	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.20ef9979-d794-46d8-b34c-a953aad188fb</a>								
<a href="#">1237433056</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:47	1 day, 15:26:01	0:05:55	05-18 13:13	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.da608a8f-aa70-4226-8cb1-db7a200f174d</a>								
<a href="#">1237433055</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:47	1 day, 15:25:35	0:06:22	05-18 13:12	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.644be014-14f0-4b1d-a631-f3e41aa8c78d</a>								
<a href="#">1237433054</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:47	1 day, 15:25:10	0:06:47	05-18 13:12	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.e3ffff3d-7769-4d43-9ff9-7018eaae6c25</a>								
<a href="#">1237433045</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:46	1 day, 15:25:19	0:06:57	05-18 13:12	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.22b7d89d-5498-4434-a61c-a843d28fcd7b</a>								
<a href="#">1237433043</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:46	1 day, 15:24:50	0:07:34	05-18 13:11	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.d217d914-c77c-4c98-b339-709511e5e5cd</a>								
<a href="#">1237433042</a> <a href="#">pilot/copilot.cern.ch</a>	trans=csc_evgen_trf.py, pkg=AtlasProduction/15.6.5.5	running	2011-05-16 21:46	1 day, 15:23:49	0:08:36	05-18 13:10	<a href="#">US/CERN.CERNVM</a> , ptest	100
<b>Out:</b> <a href="#">panda.destDB.7572d85d-2873-4315-a4bb-c78616bc4eac</a>								

rBuilder CernVM - Published Rel... Elasticfox

Regions: CERN Credentials: CERN Account IDs

Instances Images KeyPairs Security Groups Elastic IPs Volumes and Snapshots Bundle Tasks Reserved Instances Virtual Private Clouds

Your Instances

Don't show Terminated Instances

Reservation ID	Ow...	Instance ID	AMI	ARI	State	Public DNS	Local Launc...	Availability Zone	Platform	Tag	VPC
default	predr...	i-195	...	eri-1FE...	running	128.142.192.62	0 ...	1970-01-01 0...	default		
default	predr...	i-196	...	eri-1FE...	running	128.142.192.63	0 ...	1970-01-01 0...	default		
default	predr...	i-197	...	eri-1FE...	running	128.142.192.64	0 ...	1970-01-01 0...	default		
default	predr...	i-198	...	eri-1FE...	running	128.142.192.65	0 ...	1970-01-01 0...	default		
default	predr...	i-199	...	eri-1FE...	running	128.142.192.66	0 ...	1970-01-01 0...	default		
default	predr...	i-200	...	eri-1FE...	running	128.142.192.67	0 ...	1970-01-01 0...	default		
default	predr...	i-201	...	eri-1FE...	running	128.142.192.52	0 ...	1970-01-01 0...	default		
default	predr...	i-202	...	eri-1FE...	running	128.142.192.53	0 ...	1970-01-01 0...	default		
default	predr...	i-203	...	eri-1FE...	running	128.142.192.54	0 ...	1970-01-01 0...	default		
default	predr...	i-204	...	eri-1FE...	running	128.142.192.55	0 ...	1970-01-01 0...	default		
default	predr...	i-205	...	eri-1FE...	running	128.142.192.56	0 ...	1970-01-01 0...	default		
default	predr...	i-206	...	eri-1FE...	running	128.142.192.57	0 ...	1970-01-01 0...	default		
default	predr...	i-207	...	eri-1FE...	running	128.142.192.57	0 ...	1970-01-01 0...	default		
default	predr...	i-208	...	eri-1FE...	running	128.142.192.57	0 ...	1970-01-01 0...	default		
default	predr...	i-209	...	eri-1FE...	running	128.142.192.57	0 ...	1970-01-01 0...	default		
default	predr...	i-210	...	eri-1FE...	running	128.142.192.57	0 ...	1970-01-01 0...	default		

rBuilder CernVM - Published Rel... Elasticfox

Regions: us-east-1 Credentials: 3929-4179-4136 Account IDs

Instances Images KeyPairs Security Groups Elastic IPs Volumes and Snapshots Bundle Tasks Reserved Instances Virtual Private Clouds

Your Instances

Don't show Terminated Instances

Reservation ID	Ow...	Instance ID	AMI	ARI	State	Public DNS	Local Launc...	Availability Zone	Platform	Tag	VPC
r-ad962dc1	3929...	i-f3b04a9d	ami-5c3ec235	...	running	ec2-50-16-144-41.comput...	0 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-f1b04a9f	ami-5c3ec235	...	running	ec2-75-101-214-247.com...	1 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-cfb04aa1	ami-5c3ec235	...	running	ec2-184-72-183-26.comp...	2 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-cdb04aa3	ami-5c3ec235	...	running	ec2-184-73-56-72.comput...	3 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-cbb04aa5	ami-5c3ec235	...	running	ec2-50-16-32-51.compute...	4 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-c9b04aa7	ami-5c3ec235	...	running	ec2-75-101-184-46.comp...	5 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-c7b04aa9	ami-5c3ec235	...	running	ec2-50-19-38-225.comput...	6 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-c5b04aab	ami-5c3ec235	...	running	ec2-50-16-105-241.comp...	7 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-c3b04aad	ami-5c3ec235	...	running	ec2-174-129-86-61.comp...	8 ...	2011-05-18 1...	us-east-1d		
r-ad962dc1	3929...	i-c1b04aaf	ami-5c3ec235	...	running	ec2-50-19-9-47.compute...	9 ...	2011-05-18 1...	us-east-1d		

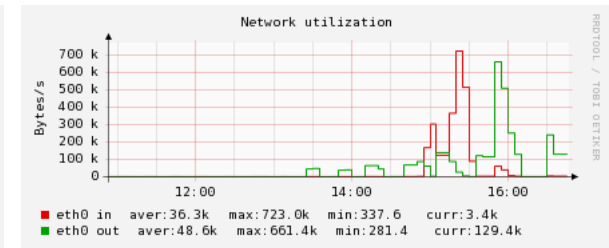
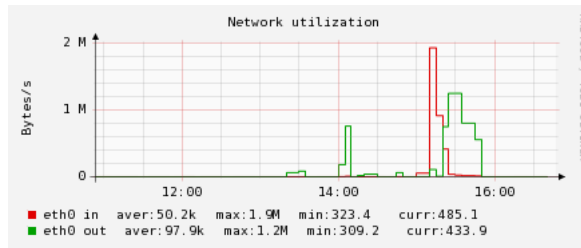
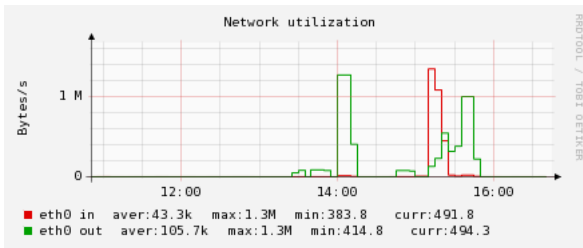
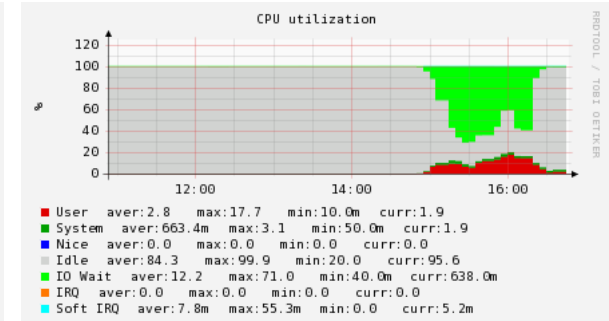
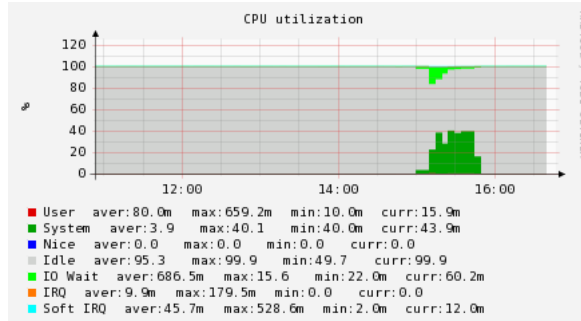
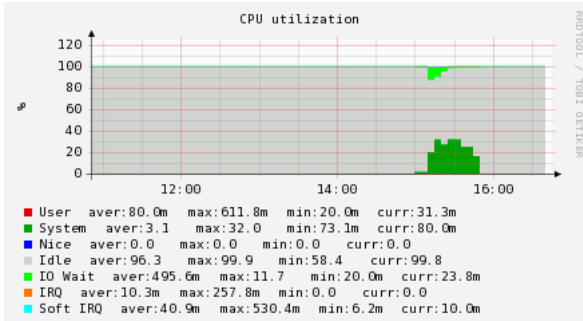
## No external connectivity to eosatlas.cern.ch (had to use local disk):

```
18 May 13:33:44| pilot.py      | Collecting WN info from: /tmp/panda-
pilot/3158/WORK/Panda_Pilot_3994_1305725624 (again)
18 May 13:33:44| pilot.py      | Job memory set by queuedata to:
18 May 13:33:44| pilot.py      | Local space limit: 5368709120 B
18 May 13:33:44| pilot.py      | !!FAILED!!1999!! Too little space left on local disk to run
job: 1967128576 B (need > 5368709120 B)
18 May 13:33:44| pilot.py      | Pilot was executed on host: ip-10-212-183-96
....
```

## EOD storage worked well on iCloud:

```
[root@vmbsq090700 ~]# df
Filesystem          1K-blocks      Used Available Use% Mounted on
/dev/hda1            13577676    4131152   8762240   33% /
none                 513176         0    513176    0% /dev/shm
eos                  3413376370848 23790724 3413352580124   1% /eos
```

```
[panda@vmbsq090700 copilot]$ ls -l *.root
024653d9-99f6-49e5-8dad-0b826f84fa3b_0.evgen.pool.root
02746a4b-ab5e-48a2-b5ff-4359f22b5da7_0.evgen.pool.root
0665543c-359a-4265-84ec-82ffdab2d8dd_0.evgen.pool.root
09719f1b-ae99-441f-9ee5-7fc265f099b4_0.evgen.pool.root
107acb02-7f9b-4b4d-bdaf-dc9b86c58dca_0.evgen.pool.root
```

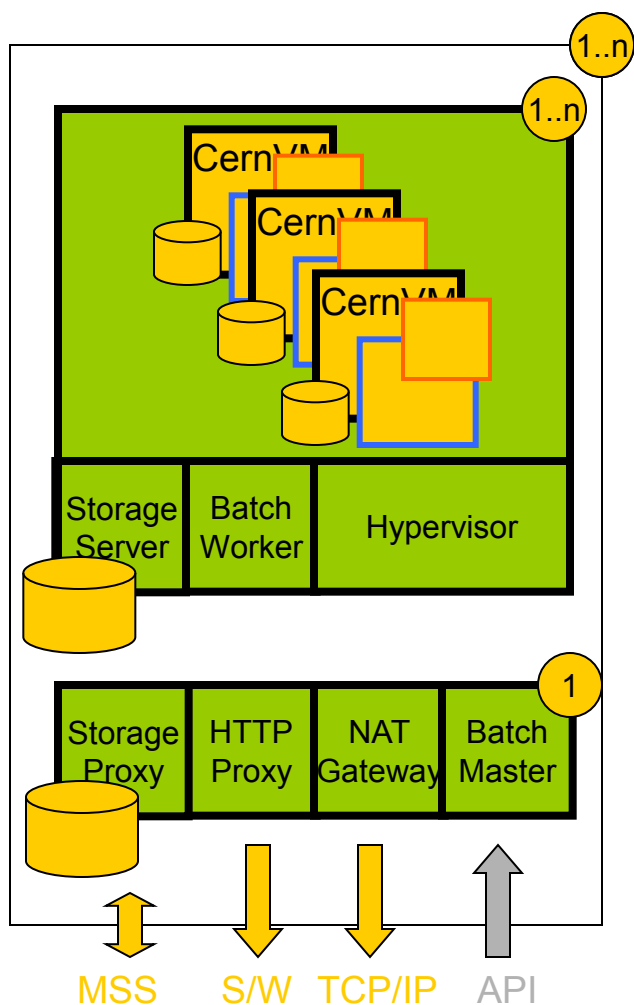


SLC5  
hypervisor

SLC6  
hypervisor

Work in progress, blue sky ideas...

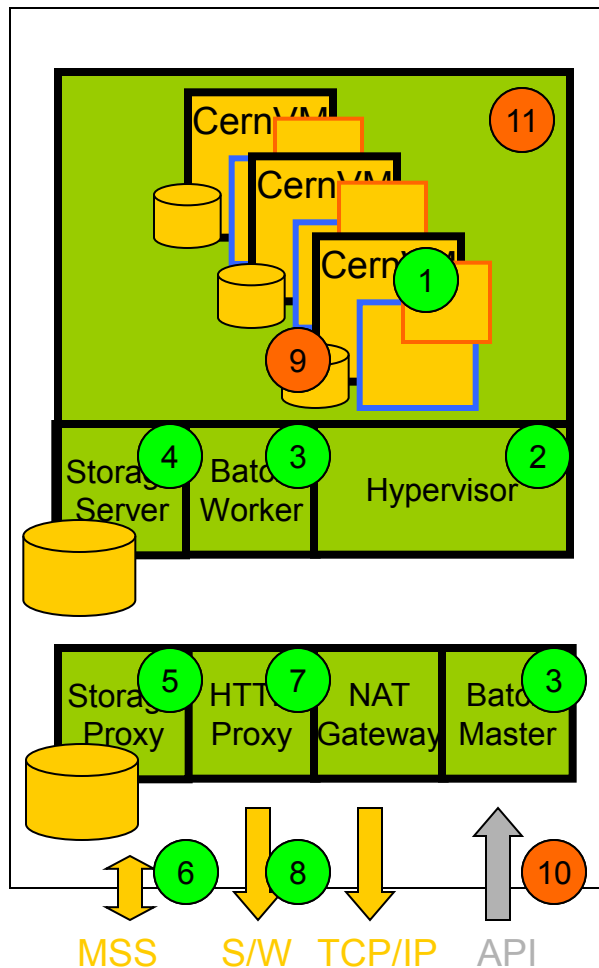
## Computer Center in a Box



- Common services hosted by front-end node
  - Batch master, NAT gateway, storage and HTTP proxy
- Each physical node
  - Contributes to common storage pool
  - Runs hypervisor, batch worker
  - Exports storage local storage to common pool
- Virtual Machines
  - Started by the suitable Cloud middleware
  - Only limited outgoing network connectivity via gateway node/HTTP proxy
  - Access to data files via POSIX (file system) layer
  - Software delivered to VMs from Web server repository
  - Built form recipes and components stored in strongly versioned repository
- Access to external mass storage via storage proxy
- End user API to submit jobs

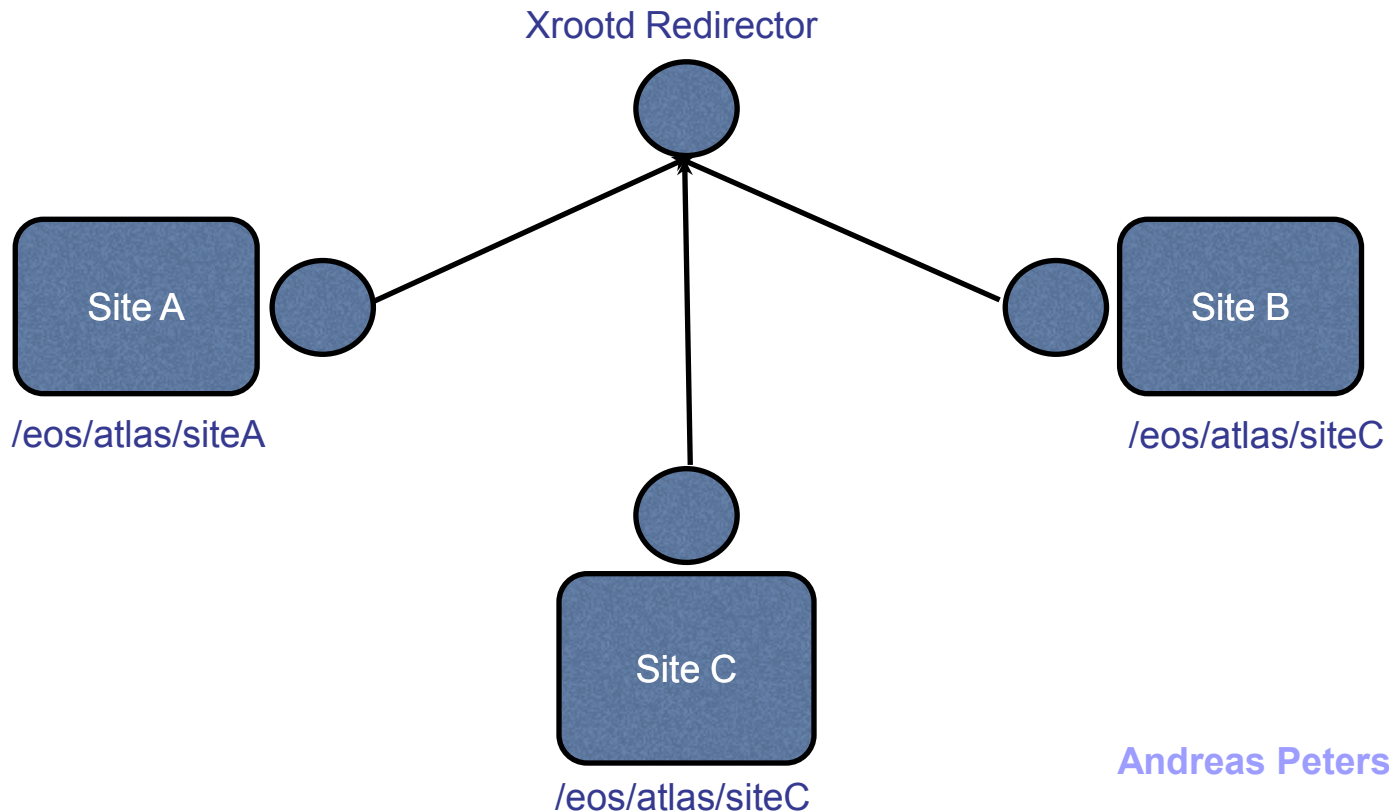


(long term data preservation use case)



- 1) CernVM for NA61
- 2) KVM - Linux native hypervisor
- 3) Condor - batch system that works well in dynamic environment and supports running jobs in VMs
- 4) Xrootd server - high performance file server, used by LHC experiments and distributed with ROOT
- 5) Xrootd redirector - each aggregates up to 64 servers, can cluster up to 64k servers
- 6) Xrootd supports access to MSS systems using SRM extension
- 7) Standard HTTP proxy (Squid)
- 8) CernVM-FS repository for software distribution
- 9) eosfs - POSIX File System for Xrootd
- 10) GANGA - user interface and front-end to various batch systems and Grid
- 11) Ganglia - monitoring

# Scaling EOS across the sites: Storage Cloud?



Andreas Peters IT/DM

👉 Storage Cloud plays fundamental role in Cloud concept, cannot be factored out as separate problem and solved later...

- DIRAC Pilot for LHCb...
  - and other experiments that use pilot job frameworks
- On demand Condor cluster
  - contextualization templates currently under testing
- PROOF on demand
- GridFactory scheduler
  - Ad-hoc cloud for small workgroups
- Volunteer computing using BOINC adapted to run jobs in CernVM
  - Ongoing activity with Theory Group at CERN
  - MonteCarlo validation project

- In combination with various **contextualization** options and **CernVM-FS**, just **one small image** can run frameworks of all LHC experiments and be **easily moved around** requiring far less frequent updates than traditional worker node
- **CernVM-FS** provides efficient, scalable, secure and maintenance free way to distribute software in CernVM and physical nodes alike
- **Flexible contextualization options** allow the same image to play different roles reducing the need for creation and certification of specialized images
- **Storage Cloud** plays fundamental role in Cloud concept, cannot be factored out as separate problem and solved later
- **IxCloud** is operational and usable, can run PanDA in CernVM with EOS as storage backend (thanks to Sebastian Goasguen and Ulrich Schwickerath)
  - Time to do some stress testing...