

Container-based CTA installation for storing midsize archive data

**Victor Kotlyar, Denis Syrko, Dmitry
Bolotin**

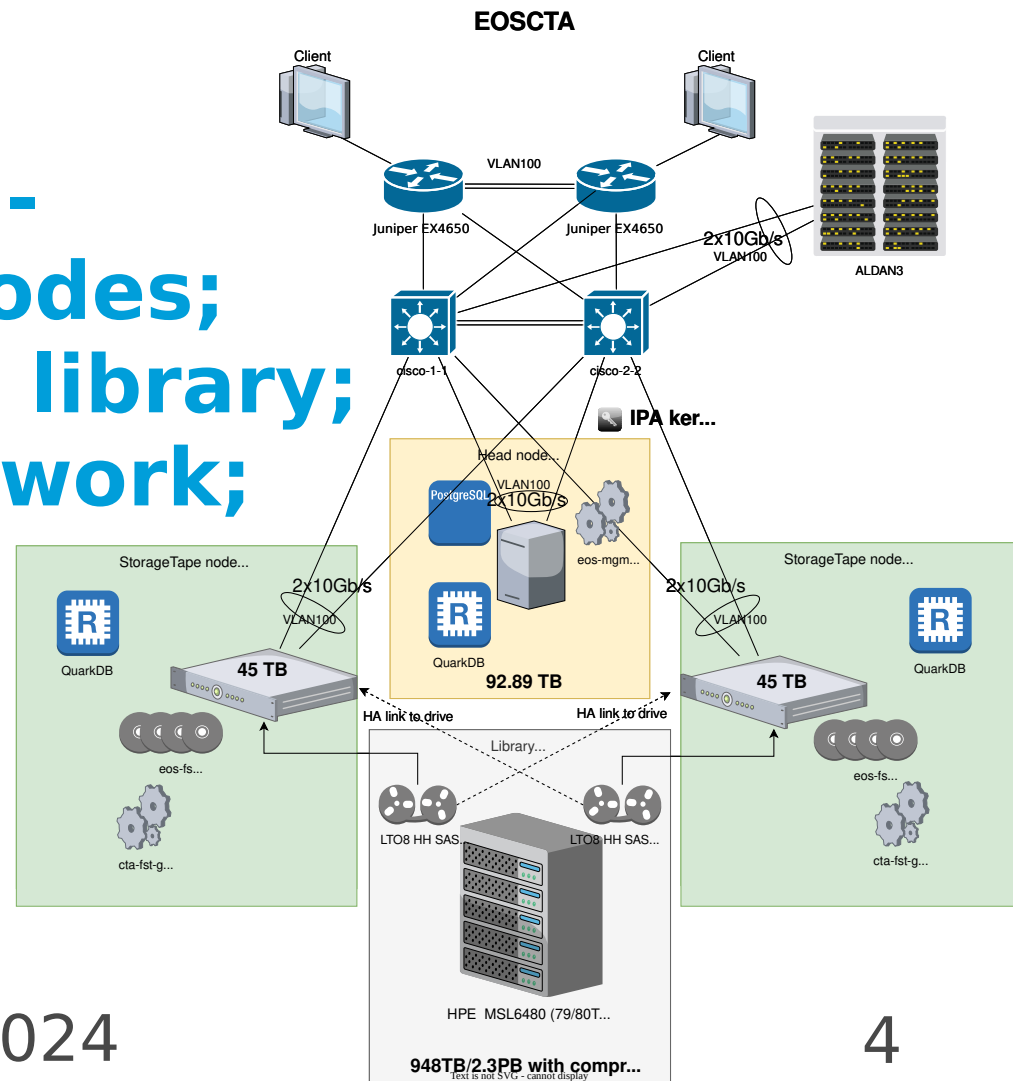
Outline

- Generic goal
- System architecture overview
- Hardware overview
- Software overview
- Installation tricks
- System usage
- Need to be done

System requirements:

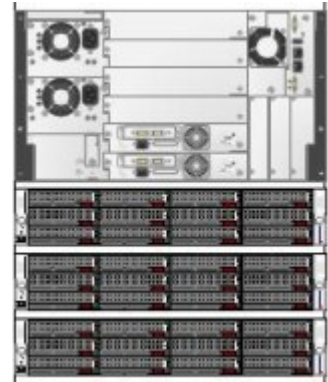
- **As small and effective as possible;**
- **For store archive data;**
- **For a backup service;**
- **Seamlessly integrated with an existing computing cluster for group based access pattern;**
- **Without any limits in size and features if needed in the future.**

- # System architecture:
- Three nodes cluster - Head + tape&disk nodes;
 - Modular expandable library;
 - Generic purpose network;
 - External IPA for AA;
 - Cluster and clients outside the system.



Hardware overview (12U):

- **Main storage -
HPE MSL6480 library with one module;
79 LTO8 tapes + 1 cleaning;
940TB (with compression in mind);
two LTO8 SAS3 HH drives;**
- **Disk storage -
6x12TB SATAIII per server 87.2TiB
disks ~ 10% of tapes**
- **All nodes are the same: 2x4216(64Th) 96GB**



Hardware overview from inside:

```
ansible@tape-1-3-1:~$ lsscsi -g
[4:0:0:0]    disk      ATA      CT480BX500SSD1   054   /dev/sda   /dev/sg0
[5:0:0:0]    disk      ATA      CT480BX500SSD1   054   /dev/sdb   /dev/sg1
[6:0:0:0]    disk      ATA      TOSHIBA MG07ACA1 0103  /dev/sdc   /dev/sg2
[6:0:1:0]    disk      ATA      TOSHIBA MG07ACA1 0103  /dev/sdd   /dev/sg3
[6:0:2:0]    disk      ATA      TOSHIBA MG07ACA1 0103  /dev/sde   /dev/sg4
[6:0:3:0]    disk      ATA      TOSHIBA MG07ACA1 0103  /dev/sdf   /dev/sg5
[6:0:4:0]    disk      ATA      TOSHIBA MG07ACA1 0103  /dev/sdg   /dev/sg6
[6:0:5:0]    disk      ATA      TOSHIBA MG07ACA1 0103  /dev/sdh   /dev/sg7
[6:0:6:0]    enclosu  SMC      SC826N4          100d  -          /dev/sg8
[15:0:0:0]   tape      HPE      Ultrium 8-SCSI   M571  /dev/st0   /dev/sg9
[15:0:1:0]   tape      HPE      Ultrium 8-SCSI   M571  /dev/st1   /dev/sg10
[15:0:1:1]   mediumx  HP       MSL6480          6.30  /dev/sch0  /dev/sg11
[15:0:2:0]   enclosu  BROADCOM VirtualSES       03    -          /dev/sg12
```

```
ansible@tape-1-3-1:~$ zpool list -vvv
NAME                SIZE  ALLOC  FREE  CKPOINT  EXPANDSZ  FRAG    CAP  DEDUP    HEALTH  ALTROOT
data                65.5T  59.3T  6.14T  -         -         1%   90%  1.00x    ONLINE  -
  raidz2            65.5T  59.3T  6.14T  -         -         1%  90.6%    -    ONLINE
    scsi-35000039b08e2be07  -      -      -      -         -         -    -      -    ONLINE
    scsi-35000039b08e2c970  -      -      -      -         -         -    -      -    ONLINE
    scsi-35000039b08e2c969  -      -      -      -         -         -    -      -    ONLINE
    scsi-35000039b08e2bc8b  -      -      -      -         -         -    -      -    ONLINE
    scsi-35000039b08e2be57  -      -      -      -         -         -    -      -    ONLINE
    scsi-35000039b08e2be24  -      -      -      -         -         -    -      -    ONLINE
```

Software overview:

- **Base operating system Ubuntu 20.04LTS standart setup;**
- **Separate docker container (CentOS7) for each service -**
 - head node: eos-mgm, eos-mq, cta-frontend, eos-qdb, PostgreSQL, nfsd;**
 - tape&disk node: eos-fst, cta-gcd, cta-taped, eos-qdb, cta-rmcd (on master)**
- **eos-fst on zfs z2 pool**

Software overview from inside:

```
ansible@tape-1-3-3:~$ sudo docker ps
```

CONTAINER ID	IMAGE	COMMAND	CREATED	STATUS	PORTS	NAMES
cb560b8fbcc6	eosmgm:latest	"/opt/run/bin/eos-mg..."	21 months ago	Up 2 months		eos-mgm-1-3-3
1f8360b6a353	eosmq:latest	"/opt/run/bin/eos-mq..."	21 months ago	Up 7 months	0.0.0.0:1097->1097/tcp, :::1097->1097/tcp	eos-mq-1-3-3
b95b2e9b5256	ctafontend:latest	"/opt/run/bin/cta-fr..."	21 months ago	Up 5 days	0.0.0.0:10955->10955/tcp, :::10955->10955/tcp	cta-frontend-1-3-3
31700c4c0caf	quarkdb:latest	"runuser --user xroo..."	21 months ago	Up 7 months	0.0.0.0:7777->7777/tcp, :::7777->7777/tcp	eos-qdb-1-3-3
1a4c0ad4e0cb	centos/postgresql-13-centos7	"container-entrypoin..."	21 months ago	Up 7 months	0.0.0.0:5432->5432/tcp, :::5432->5432/tcp	postgresql

```
ansible@tape-1-3-1:~$ sudo docker ps
```

CONTAINER ID	IMAGE	COMMAND	CREATED	STATUS	PORTS	NAMES
ed19748785c2	ctagcd:latest	"/opt/run/bin/cta-fs..."	17 months ago	Up 7 months		cta-gcd-1-3-1
92b420ef45e1	ctarmcd:latest	"/opt/run/bin/cta-rm..."	18 months ago	Up 7 months		cta-rmcd-1-3-1
247b23d892fc	fe0216285abd	"/opt/run/bin/cta-ta..."	20 months ago	Up 7 months		cta-taped-1-3-1
3bea991048c2	eosfst:latest	"/opt/run/bin/eos-fs..."	21 months ago	Up 7 months		eos-fst-1-3-1
601eddeb7f59	quarkdb:latest	"runuser --user xroo..."	21 months ago	Up 7 months	0.0.0.0:7777->7777/tcp, :::7777->7777/tcp	eos-qdb-1-3-1

```
ansible@tape-1-3-2:~$ sudo docker ps
```

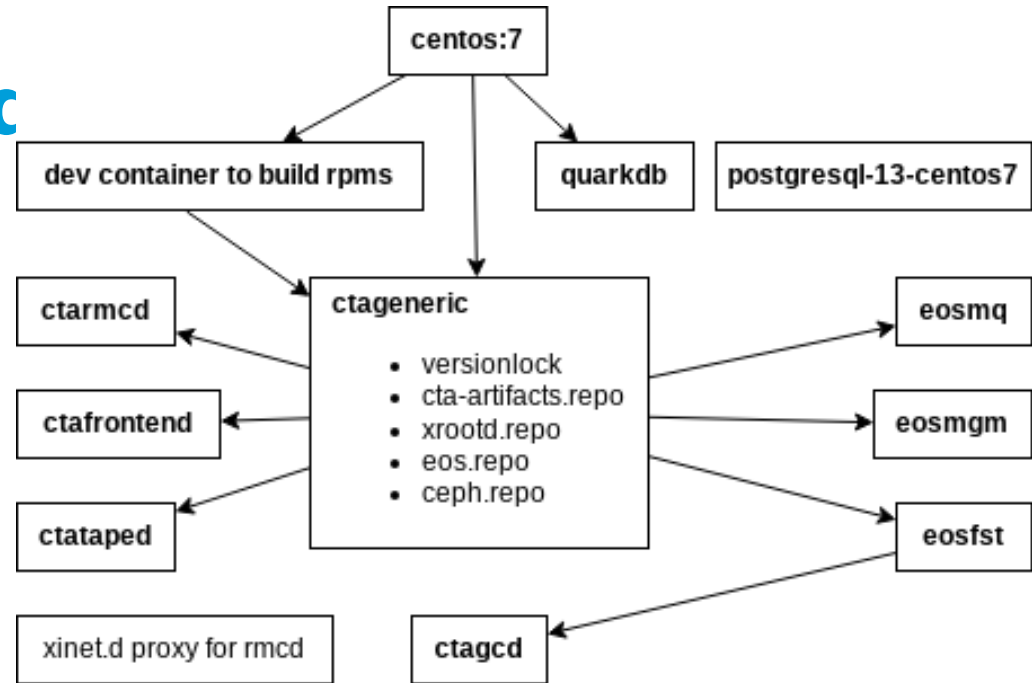
CONTAINER ID	IMAGE	COMMAND	CREATED	STATUS	PORTS	NAMES
2b0f8b7b5208	ctagcd:latest	"/opt/run/bin/cta-fs..."	17 months ago	Up 7 months		cta-gcd-1-3-2
f047470da4c1	ctataped:latest	"/opt/run/bin/cta-ta..."	20 months ago	Up 2 days		cta-taped-1-3-2
4865672bda87	eosfst:latest	"/opt/run/bin/eos-fs..."	20 months ago	Up 7 months		eos-fst-1-3-2
1056cb893330	quarkdb:latest	"runuser --user xroo..."	21 months ago	Up 7 months	0.0.0.0:7777->7777/tcp, :::7777->7777/tcp	eos-qdb-1-3-2

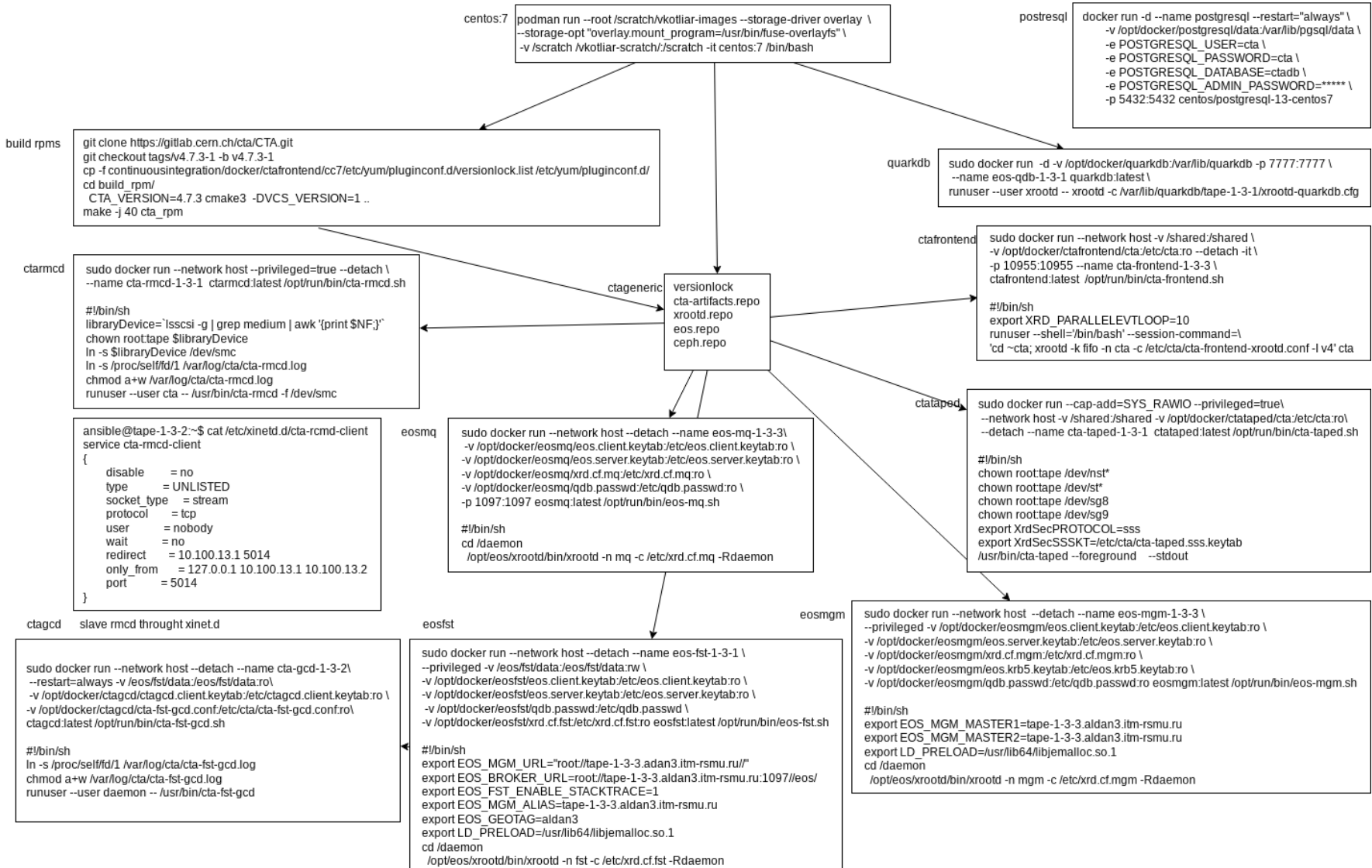
Installations tricks:

- **All information to build containers has been taken from sources or from the official docs;**
- **Based on images creation - runs on base image; make modifications inside container; save image as a service base image; use a service base image to run service.**

Installations tricks: container images hierarchy

- Centos:7 was used
- CTA rpms were built;
- Ctageneric with version lock as a prototype for all others;
- Each container is one service system





Installations wish list (maybe some already addressed):

- Would be better to be a container friendly (no file to stdout redirect);**
- Allow to use only config files without env variables;**
- Restrict as much as possible privileges and avoid their exscalations;**
- One rmcd daemon per tape library;**
- Allow to use tape drive with disks (one node transfer optimisations)**

System usage from cluster:

- For backup and archiving we use direct eos and xrootd commands (tried webdav to put files but there is no use case in our setup);**
- Store files as long as possible on disks;**
- Authorisation group based so in EOS there is a role for each tapeops group;**
- To prepare archives we use 4G tar files;**
- CEPH - ZFS - LT08 all of the use similar compression LZ4-LZ4-LZS very convinient**

System usage from inside:

```
# eos-mgm-1-3-3
eos mkdir /eos/cta/aldan3
eos chown vkotliar:aldan3_tapeops /eos/cta/aldan3
eos chmod 770 /eos/cta/aldan3
eos attr set sys.acl=u:0:rwx+dp,g:aldan3_tapeops:rwx+dp,z:'!u,u:0:+u /eos/cta/aldan3
eos vid ls
  eos vid set membership vkotliar -gids 1516200179,1516200054
  eos vid set membership dbolotin -gids 1516200179,1516200054
  eos vid set membership taped -gids 1516200179,1516200054,1516200081
eos vid ls
```

```
# name: Create archive files
cat /root/fstab | \
grep home | \
grep ceph | \
grep "\/$userNameForBackup " | \
awk '{print $2;}' | \
xargs -iDIR /bin/bash -c 'mkdir $tempDirForBackup/DIR ; cd $tempDirForBackup/DIR ; DNAME=DIR; NAME="${DNAME////_}";
$TAR_PATH -c --ignore-failed-read --multi-volume --tape-length=4G --file=$NAME.tar.{0..10000} DIR'
```

```
# name: Copy to EOS
rc=3
export EOS_MGM_URL="root://tape-1-3-3.aldan3.itm-rsmu.ru"
for HDIR in `ls .` ;
do kinit -k -t /home/vkotliar/krb5.keytab vkotliar ;
set +e
  eos --role vkotliar aldan3_tapeops mkdir $eosDirForBackup;
set -e
  eos --role vkotliar aldan3_tapeops cp -r $HDIR $eosDirForBackup;
done
```

Needs to be done:

- Upgrade to the “latest greatest”
AL9
EOS5
CTA5

Thank you!

Questions?