Scalable Cluster Management: LXC³

Erich Focht Andreas Jeutter NEC Deutschland *April 2023*

© NEC Corporation 2023

Motivation



The Problems

• Bare metal deployments \rightarrow platform specifics

- Stateful or stateless? Both!
- Node descriptions \rightarrow Data model!
- What shall be deployed? \rightarrow Node images
- Node images life cycle
- Nodes are not equal!
- Remote node power / platform control
- Management node(s) installation

Many steps, error prone. "Recipes" won't work.

More Comfort

- High Availability!
- Batch system and resource scheduler
- User management
- Performance Monitoring
- Health monitoring and Alerting
- License server

... Stuff

• Git server, registry, repository, application software store, ...

History



Perceus

- Perceus: enterprise spin-off from Warewulf
- Stateless / diskless compute nodes





Perceus

- Perceus: enterprise spin-off from Warewulf
- Stateless / diskless compute nodes





Diskless only!

LXC³ - Perceus

Added high-availability: active/passive with pacemaker/corosync

- Storage synchronization with DRBD
- Added stateful deployments (with disks)
- Replaced Perceus database with LXDaemon (hierarchical database)
- Added cluster services, resource scheduler, monitoring, alerting, power control, ...



LXC³ - Experience

- Installed at many customers
- Fulfilled expectations for small / medium clusters: O(500) nodes
- Covering most needs
- VNFS concept good, but ... versioning is expensive
- Automatization: helpful, but ... dependencies between components
 - Updating requires a lot of knowledge and care
- Concept not directly scalable to larger systems

Core requirements:

Scale to 2000 – 5000 nodes

High-availability extremely important

- Support heterogeneous setups with many node types
- Updating must be easy
- Base installation must be easy

LXC³ neo Architecture

Three or more master nodes
Cluster organized in groups
Services scaled → groups
Storage: distributed/parallel
Cluster database: distributed
Single point of control



LXC³ neo Implementation Choices

Services \rightarrow Microservices in Docker containers

- Select high performance alternatives
 - Eg. dnsmasq \rightarrow atftp + kea, perl tcp/ip server \rightarrow apache, kernel nfsd \rightarrow ganesha user space nfsd

\bullet High-availability \rightarrow Container orchestration

- Services move, but keep their IP addresses: MACVLAN
- Shared volumes with services data and state: parallel file system
- Database: distributed key-value store
- Keep: VNFS, Modules
 - but switch modules client/server from Perl to apache + curl



LXC³ neo Single Master Setup



LXC³ neo : yes, there are GUIs





E Potiener presay K	+						V O Private brow
+ - C & O D + H	calhest 10	Services					0 3 0 - 100x + 4 🕻 🖉 > D
portainer.io		Servi	ice list O				Q 🗇 A admin
		*	Services				C Basin for a service
			Name 1	Stack	Image 17	Scheduling Mode 11	Published Ports
E Destioned			apatra		%1,4path4237810.pos40078-2454	replaced $1 \neq 1$ EScale	dime.
@ Sacks					%c3,wttp.31713.gt6666a6-073	gibe s 7 s	
12 Services 49 Containers		0.	-		3x3,coredra 31703.g3938/14-18-3	replaced 1 / 1 Edam	Barries Barries
		.0.	-		tec3,8ea.31713.gf568f5e8-18.0	replated 1 / 1 Edam	
< Networks B Volumes		0.	portanae, agent	portanae	portamentagent tareat	global 3.7.3	
		0.	portanet, portanet	portainer	portane/portane/server	replaced 1.7.1 Elsave	12 commerce 12 commerce
E Seato		0.	per la		363,prix 31726.p07304240-2.0.16	replicated 1 / 1 Efforts	12 mm mm
		.0.9	soubyte-th		%c3,qualityte3+mi 31713.gm46mab=312	replated (ℓ) Effore	S lossion S lossion S S S lossion S lossion S rations in S rations in S
Settings		0.	No.		3x3,55593724.95734240-2.639	replaced $i, \ell \in \mathbf{I}$ (see	
							tanin (an taile 22 m
							*
							•











LXC³ neo Evaluation & Feedback

Scalability: provisioning is very fast

- Bring up 2000 diskless nodes in less than 100s
 - Most time spent in platform / BIOS
- If power circuits permit: just turn on all nodes at once
 - Previously: iterative process, with delays

Flexibility

- Admins are able to add services (as containers) into the master swarm
- Updates: easy for services / containers
- Add modules with host specific customizations
 - eg. ansible module

LXC³ neo Evaluation & Feedback

Complexity

- Steep learning curve if no experience with containers
- Network: macvlan fixed IP
 - Not supported by docker swarm
 - Requires workaround and additional (proxy-) containers

VNFS image handling

- Slow due to ... parallel file system + many small files
- VNFS "mount modify umount" administration hard to track

Parallel file system

- Overhead for parallel FS comes with a price...
- Issues when clocks slightly out of sync

Development, CI

Building: python packages, RPMs, container images
Automated testing on VMs and bare metal



- Networking without macvlan (ready)
 - Direct use of network interface not necessary: >25GbE interfaces are common
 - Swarm mesh network
 - CoreDNS aware of services location + load balancing
 - Eliminates need for proxy containers \rightarrow less containers running
 - Less IP addresses needed (for each master node only, not for services)
 - \rightarrow setup works on kubernetes, too (macvlans don't)

Networking without macvlan (ready)

before



Networking without macvlan (ready)

after

21



VNFS deploy from container image registry (ready)

before



\checkmark VNFS deploy from container image registry (ready)



Eliminate need for parallel filesystem
VNFSes are container images

• Changes trackable in git and reproducible

Networking simplified

Can run under kubernetes, too

Final Remarks

- Pure provisioning systems are "easy"
 - And come without high-availability, scaling, concept for additional admin services
 - In LXC³-neo we cared about the missing bits
- Convergence of cluster + cloud management methods
 - Cloud: moving target, evolving fast
 - Admins meanwhile got used to containers
 - Node images are also container images
- Trend: mixing HPC + AI workloads
 - AI: more ... interactive, running in containers
 - HPC: mostly bare metal, but ... why not run each job in containers, too?

Orchestrating a brighter world

