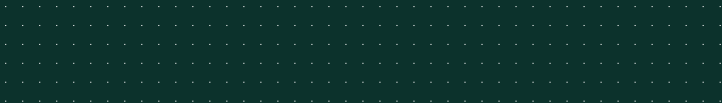


APRIL 25, OSTRAVA

Warewulf making cluster installations fast and reliable

Christian Goll <cgoll@suse.com>



Introduction

Warewulf is tool for managing Beowulf clusters

Beowulf

- old british poem

Beowulf cluster

- became popular in the 90.
- use of the shelf hardware
 - 486 & Linux **not** Cray & Unix
- Warewulf is a typo of werewolf
- Caos Linux & Warewulf were designed for Beowulf clusters
- Caos Linux evolved into CentOS
 - Gregory Kurtzner started developement of Warewulf v4



Introduction

HPC landscape

Top five Supercomputers

1	Frontier	EPYC 64C	AMD MI250X	Slingshot-11
2	Aurora	Xeon 9470	Intel GPU Max	Slingshot-11
3	Eagle	Xeon 8480	NVIDIA H100	NVIDIA Infiniband
4	Fugaku	A64FX 48C 2.2GHz	-	Tofu interconnect D
5	LUMI	EPYC 64C 2GHz	AMD MI250X	Slingshot-11

- only Fugaku uses non standard CPU
- others are Beowulf clusters with GPUs attached



Introduction

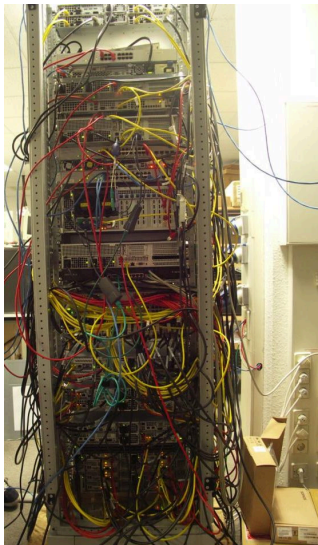
Beowulf cluster

Base components

- management node
- compute nodes
- management network
- network boot (PXE)

Optional components

- more compute nodes
- fast network interconnects
- central storage
- bmc/ipmi



Beowulf Cluster

- compute nodes are cattle
 - parallel/MPI require identical nodes
- hierarchical organization
- compute are not updated after boot process
- applications come from central storage
- applications are self compiled with tools like **spack** & **EasyBuild**



Warewulf description

Software stack

Warewulf components

Compute nodes

boot from network into tmpfs

Warewulfd delivers

- kernel & modules
- node image
- node configurations

wwctl cmd line tool

- manages node database
- manages node image



Copyright © SUSE 2024

external components

Dhcp server

- ISC dhcpd server
- dnsmasq

tftp

- tftp from
kernel.org
- dnsmasq

Optional

- nfs
- manage /etc/hosts ⁶

Warewulf configuration

database /etc/warewulf/nodes.conf

- plain yaml file
- easy backup
- can be version controlled
- external tools support
 - vim, ansible

Profiles

- stores identical values for collection of nodes
- values can be overridden on node basis

```
1 WW_INTERNAL: 45
2 nodeprofiles:
3   default:
4     comment: This profile is automatically
5     container name: leap
6     network devices:
7       default:
8         device: eth0
9 nodes:
10  n01:
11    profiles:
12    - default
13    network devices:
14      default:
15        hwaddr: 52:54:00:4e:cb:1d
16        ipaddr: 172.16.130.101
17  n02:
```



Warewulf configuration

command line database manipulation

add node

```
1 wwctl node add n01  
  -I 10.10.10.1
```

modify node

```
1 wwctl node set n01  
  --comment "Have fun"
```

list node

```
1 wwctl node list  
  n01 -a
```

	NODE	FIELD	PROFILE	VALUE
1	n01	Id	--	n01
2	n01	Comment	SUPERSEDED	Have fun
3	n01	ContainerName	default	leap
4	n01	Ipxe	--	(default)
5	n01	RuntimeOverlay	--	(generic)
6	n01	SystemOverlay	--	(wwinit)
7	n01	Root	--	(initramfs)
8	n01	Discoverable	--	false
9	n01	Init	--	(/sbin/init)
10	n01	Kernel.Args	--	(quiet crashkernel=)
11	n01	Profiles	--	default
12	n01	PrimaryNetDev	--	(default)
13	n01	NetDevs[default].Type	--	(ethernet)
14	n01	NetDevs[default].OnBoot	--	(true)
15	n01	NetDevs[default].Device	default	eth0
16	n01	NetDevs[default].Hwaddr	--	52:54:00:4e:cb:1d
17	n01	NetDevs[default].Ipaddr	--	172.16.100.101



Warewulf configuration

Templates& Overlays

Configuration templates

- based on go templates
- `{{.foo}}` replaced with variable `foo`
- exported go function can be called

Configuration overlays

- rendered templates packed into overlay
- overlay put on top of node image

Listing 1: `/etc/issue.ww` → `/etc/issue`

```
1 Warewulf Node:      {{.Id}}
2 Container:          {{.Container}}
3 {{ if .Kernel.Version }}Kernel:
4 Kernelargs:          {{.Kernel.Args}}
5
6 Network:
7 {{- range $devname, $netdev := .NetDevs}}
8     {{$devname}}:    {{$netdev.Device}}
9     {{$devname}}:    {{$netdev.IpCIDR}}
10 {{if $netdev.Ipaddr6 }}      {{$devname}}:
11 {{if $netdev.Hwaddr }}      {{$devname}}:
12 {{end}}
13
14 Greetings from {{ .Tags.Location }}
```



Warewulf configuration

Rendered template

User defined variables

- every profile/node can have user defined variables
- extra namespace for networks

Add location tag

```
1 wwctl node set n01 --tagadd="
  Location=Ostrava"
```

Render template

```
1 wwctl overlay show --render n01
  wwininit /etc/issue.ww
```

Listing 2: renderd /etc/issue

```
1 backupFile: true
2 writeFile: true
3 Filename: /etc/issue
4 Warewulf Node:      n01
5 Container:          leap
6 Kernelargs:         quiet crashkernel=no
7
8 Network:
9     default: eth0
10    default: 172.16.130.101/24
11    default: 52:54:00:4e:cb:1d
12
13 Greetings from Ostrava
```



Warewulf configuration

Overlays

warewulf defines two types of overlays

System overlay

- available on boot
- warewulf boot strap files
- static network configurations:
 - wicked
 - NetworkManager
 - other legacy scripts
- nfs mounts
- file system mounts

Runtime overlay

- updated regularly
- can be secured

User defined overlays

- users are encouraged to create own configuration templates
- can reside in system & runtime overlays

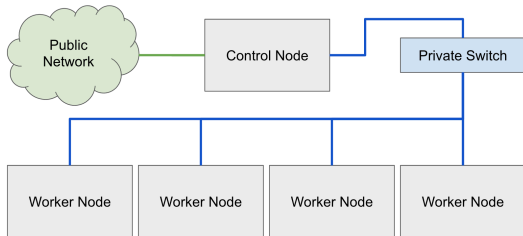


Warewulf configuration

Security

Assumptions

- private/cluster network is secure
- parallel filesystems & NFS
→ root on node means root everywhere



measurements

- node image & system overlay protected with Asset Tag
- system overlays **must** be downloaded from privileged port
- persistent malware installation is impossible as node images are ephemeral



Warewulf configuration

Node images

Elements

- complete OS images
- called containers in warewulf
- must be imported from:
 - chroot directory
 - OCI registry
- multiple different node images can be imported
- node images independent from host OS

registry.suse.com

- SUSE SLE 15SP5

registry.opensuse.org

- openSUSE Tumbleweed
- openSUSE Leap 15SP[3-5]

ghcr.io

- openSUSE Leap
- Rocky EL (8&9)
- Debian Bockworm



Warewulf configuration

Node image examples

Import SLE image from SUSE registry

```
1 ww4-host:# export WAREWULF_OCI_USERNAME=cgoll@suse.com
2 ww4-host:# export WAREWULF_OCI_PASSWORD=INTERNAL-USE-ONLY-xxxxxx
3 ww4-host:# wwctl container import docker://registry.suse.com/suse/hpc/
  warewulf4-x86_64/sle-hpc-node:latest sle-hpc
```

Import Leap image openSUSE registry

```
1 ww4-host:# wwctl container import docker://registry.opensuse.org/science/
  warewulf/leap-15.5/containers/kernel:latest leap
```

Execute shell in images

```
1 wwctl container shell sle-hpc
2 [sle-hpc] Warewulf>
```



Warewulf configuration

disk management

Needs following elements:

- disks
- partitions (on disk)
- filesystem (on partition)

Implementation

- use **ignition**
warewulf has its own service:
ignition-ww4-disk.service
- **not** called in `initrd`
- **before** `sysroot.mount`

Single partition

```
1 wwctl node set n01 --diskname /dev/  
  vda --diskwipe --partname scratch  
    --partcreate --fsname scratch --  
    fsformat btrfs --fspath /scratch  
    --fswipe
```

Add swap

```
1 wwctl node set n01 --diskname /dev/  
  vda --partname swap --partsize  
    =1024 --partnumber 1 --fsname  
    swap --fsformat swap --fspath  
    swap
```

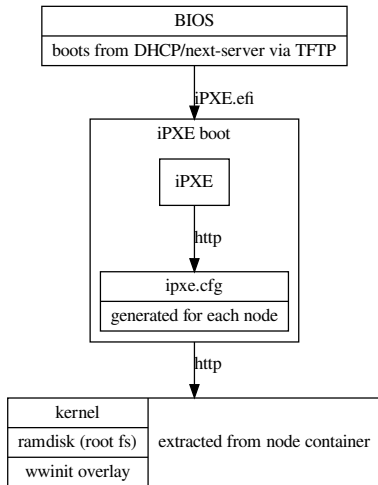


Warewulf boot

Boot process

boot with iPXE

- distribution iPXE binaries are used
- `tftp` transfers are small
- `kernel` is extracted from container/node image on the fly
- root fs is the container/node image configuration overlay added on top
- **no** secure boot

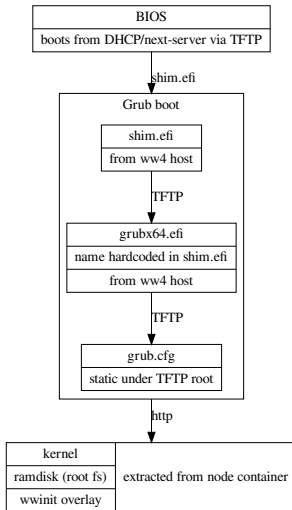


Warewulf boot

Secure boot

Boot with grub tftp

- grub & shim extracted from host
- secure boot only possible if host shim & grub can boot node kernel

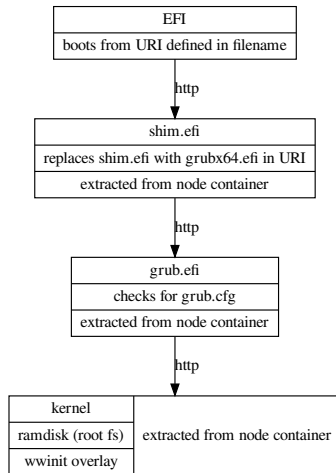


Warewulf secure boot

Cross product secure boot

Boot with grub http

- grub & shim extracted from node/container image
- secure boot with various distributions
- must be configured in BIOS



Warewulf

Development

Availability

- part of SLE since 15SP5
- SLE node/container available

upstream

github.com/warewulf/warewulf

Rocky Linux Foundation project

Stakeholders:

- SUSE
- CIQ
- Intel/openHPC



Thank your, for you attention!

