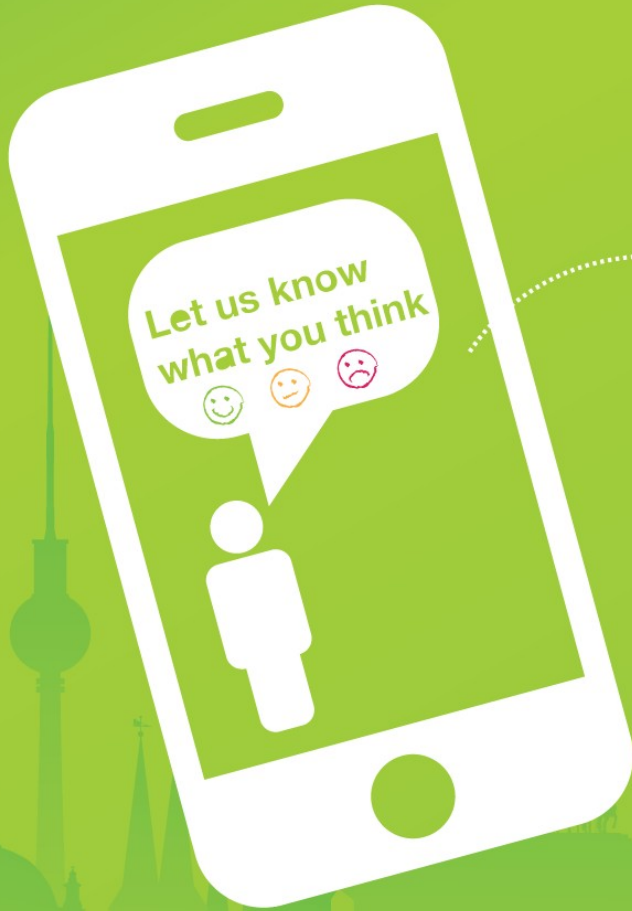


# goto;

conference



**Click 'engage'  
to rate sessions  
and ask questions**



Follow us on Twitter @GOTOber

[www.gotober.com](http://www.gotober.com)

# Nils Magnus, LinuxTag Association

## Docker Security:

Who can we trust, what should we verify?



**where .com meets .org**  
Open Source events since 1996

The background of the slide is a photograph of several stacked shipping containers in a yard. Most containers are reddish-brown, but one in the middle-left is white. A semi-transparent grey rectangle is overlaid on the center of the image, containing the title text in white. The text is split into two lines: 'Ubiquitous Containers:' on the top line and 'Excitement vs. production use' on the bottom line. The font is a clean, sans-serif typeface.

# Ubiquitous Containers: Excitement vs. production use



**where .com meets .org**  
Open Source events since 1996





# **53% of decision-makers have concerns about security**

Source: Forrester/Red Hat, January 2015

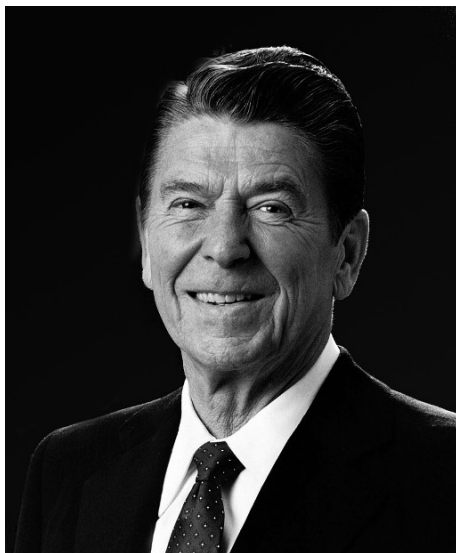
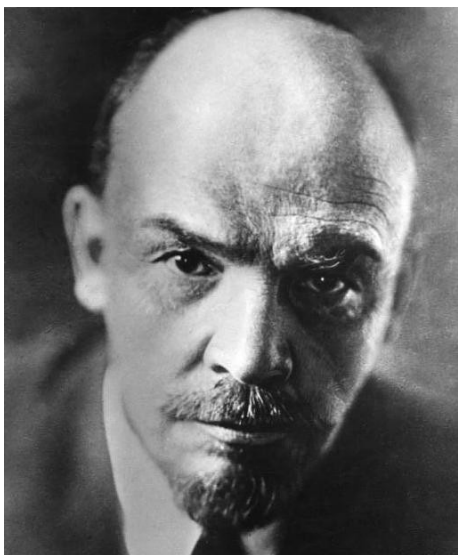
## **Do container contain?**



**where .com meets .org**  
Open Source events since 1996

# Доверяй, но проверяй

(„doveryai, no proveryai“)



*„trust,  
but verify“*



where .com meets .org  
Open Source events since 1996

# „Is Docker secure?“

- The term „secure“ depends on your security objectives
- To use Docker is most often more secure compared to not using Docker
- What do you plan to protect your software from?

→ Most of the time: **Isolation** (to the host and other containers)

*„Protect against mistake, not abuse!“*

*– from the Docker, Inc. documentation*

*„Containers do not contain!“*

*– Dan Walsh, Red Hat*



**where .com meets .org**  
Open Source events since 1996

# Isolation on Several Levels

	Physical Host	Virtual Maschine	Container
Shared Ressources	Share a network	Share host hardware	Share linux kernel
Attack scenario	Attack hosts and ports	Attack hypervisor	Attack kernel isolation via syscall (namespaces, cgroups, ...)
Security Controls	Portfilter, firewalls, segmentation of networks	Good hypervisor	Inside containermanager, SE-Linux, capabilities
Complexity of Mitigation	Easy, best Practices	Complex, but possible to manage from a centralalized point	Tricky in detail due to large attack surface



**where .com meets .org**  
Open Source events since 1996

# Four Threat Types

- Tech Attacks
- Architecture
- Sources
- User



**where .com meets .org**  
Open Source events since 1996



An aerial, high-angle photograph of a dense urban landscape, likely a city like Amsterdam, showing a complex grid of buildings and a winding river. The image is in grayscale and serves as the background for the slide.

# Threat 1: Complex attack surface



**where .com meets .org**  
Open Source events since 1996

# Capabilities

- Withdrawing of capabilities prevents restoration of status quo ante (even with root permissions):

```
# getpcaps $$
Capabilities for `22424': =
cap_chown,cap_dac_override,cap_dac_read_search,cap_fowner,cap_fsetid,cap_kill,cap_setgid,cap_setuid,cap_setpcap,cap_linux_immutable,cap_net_bind_service,cap_net_broadcast,cap_net_admin,cap_net_raw,cap_ipc_lock,cap_ipc_owner,cap_sys_module,cap_sys_rawio,cap_sys_chroot,cap_sys_ptrace,cap_sys_pacct,cap_sys_admin,cap_sys_boot,cap_sys_nice,cap_sys_resource,cap_sys_time,cap_sys_tty_config,cap_mknod,cap_lease,cap_audit_write,cap_audit_control,cap_setfcap,cap_mac_override,cap_mac_admin,cap_syslog,cap_wake_alarm,cap_block_suspend+ep

# docker run -it ubuntu /bin/bash
root@39ed301e0731:/# getpcaps $$
Capabilities for `1': =
cap_chown,cap_dac_override,cap_fowner,cap_fsetid,cap_kill,cap_setgid,cap_setuid,cap_setpcap,cap_net_bind_service,cap_net_raw,cap_sys_chroot,cap_mknod,cap_audit_write,cap_setfcap+eip
```



where .com meets .org  
Open Source events since 1996

# Issues with Capabilities

- Large number of usecases
- Only less than 40 capabilities are defined
- Semantics not very well defined
- Example: „CAP\_SYS\_PACCT“ has 30+ use cases, ranging from „random device management“ to „turning DMA on/off in xd driver“



**where .com meets .org**  
Open Source events since 1996

# Syscalls

- System call API to the kernel is rapidly growing:
- `Linux/arch/arm/include/asm/unistd.h`:  
`22 #define __NR_syscalls (392)`
- Each call is a potential attack vector into the kernel and thus to the host
- Within the kernel a single, tiny bug exploits the whole isolation



**where .com meets .org**  
Open Source events since 1996

# Issues with Syscalls

- Sebastian Kramer from the Openwall project released in June 2014 a proof-of-concept („Shocker“), enabling him to escape Docker 0.11 (predecessor of version 1.0)
- Docker creates a new filesystem context and **bindmounts** new „/“.
- Container and host share within the kernel the **the same *struct fs*** in order to maintain bindmounts.
- Do you know **syscall *open\_by\_handle\_at()***? To use it, you need **CAP\_DAC\_OVERRIDE**, which Docker had at that time.
- The resulting resources allowed you to **traverse the inodes of the host**. That enabled you to read */etc/shadow*, for example.



where .com meets .org  
Open Source events since 1996



# Namespaces

- Virtualize/isolate important system resources like
  - PIDs,
  - network interfaces,
  - UIDs,
  - hostnames and more
- Old way: Access global variable within the kernel for a resource
- New way (namespace enabled): Ask a nsproxy for the resource as which is inside your current namespace
- All access paths inside the kernel need to be scrutinized



**where .com meets .org**  
Open Source events since 1996

# Namespace example

- „hostname“ command → syscall uname(2) → kernel space

- kernel/sys.c:

```
1141 SYSCALL_DEFINE1(newuname, struct new_utsname __user *, name)
1146     if (copy_to_user(name, utsname(), sizeof *name))
```

- include/linux/utsname.h:

```
72 static inline struct new_utsname *utsname(void)
74     return &current->nsproxy->uts_ns->name;
```

- „cat /proc/version“ command → open() → kernel space → procfs

- proc/version.c:

```
8 static int version_proc_show(struct seq_file *m, void *v)
10     seq_printf(m, linux_proc_banner,
11     utsname()-sysname, [...]
```



**where .com meets .org**  
Open Source events since 1996

# Bind- and other mounts

```
root@5a5ec53ca213:/# mount
none on / type aufs (rw,relatime,si=39574450792819a9,dio,dirperm1)
proc on /proc type proc (rw,nosuid,nodev,noexec,relatime)
tmpfs on /dev type tmpfs (rw,nosuid,mode=755)
devpts on /dev/pts type devpts (rw,nosuid,noexec,relatime,gid=5,mode=620,ptmxmode=666)
sysfs on /sys type sysfs (ro,nosuid,nodev,noexec,relatime)
tmpfs on /sys/fs/cgroup type tmpfs (ro,nosuid,nodev,noexec,relatime,mode=755)
cgroup on /sys/fs/cgroup/cpuset type cgroup (ro,nosuid,nodev,noexec,relatime,cpuset)
cgroup on /sys/fs/cgroup/cpu type cgroup (ro,nosuid,nodev,noexec,relatime,cpu)
cgroup on /sys/fs/cgroup/cpuacct type cgroup (ro,nosuid,nodev,noexec,relatime,cpuacct)
cgroup on /sys/fs/cgroup/memory type cgroup (ro,nosuid,nodev,noexec,relatime,memory)
cgroup on /sys/fs/cgroup/devices type cgroup (ro,nosuid,nodev,noexec,relatime,devices)
cgroup on /sys/fs/cgroup/freezer type cgroup (ro,nosuid,nodev,noexec,relatime,freezer)
cgroup on /sys/fs/cgroup/net_cls type cgroup (ro,nosuid,nodev,noexec,relatime,net_cls)
cgroup on /sys/fs/cgroup/blkio type cgroup (ro,nosuid,nodev,noexec,relatime,blkio)
cgroup on /sys/fs/cgroup/perf_event type cgroup (ro,nosuid,nodev,noexec,relatime,perf_event)
cgroup on /sys/fs/cgroup/net_prio type cgroup (ro,nosuid,nodev,noexec,relatime,net_prio)
cgroup on /sys/fs/cgroup/hugetlb type cgroup (ro,nosuid,nodev,noexec,relatime,hugetlb)
systemd on /sys/fs/cgroup/systemd type cgroup (ro,nosuid,nodev,noexec,relatime,name=systemd)
/dev/disk/by-uuid/79bd203c-aea7-4564-b00d-7ac555e31168 on /etc/resolv.conf type ext4 (rw,relatime,errors=remount-ro,data=ordered)
/dev/disk/by-uuid/79bd203c-aea7-4564-b00d-7ac555e31168 on /etc/hostname type ext4 (rw,relatime,errors=remount-ro,data=ordered)
/dev/disk/by-uuid/79bd203c-aea7-4564-b00d-7ac555e31168 on /etc/hosts type ext4 (rw,relatime,errors=remount-ro,data=ordered)
shm on /dev/shm type tmpfs (rw,nosuid,nodev,noexec,relatime,size=65536k)
mqueue on /dev/mqueue type mqueue (rw,nosuid,nodev,noexec,relatime)
devpts on /dev/console type devpts (rw,nosuid,noexec,relatime,gid=5,mode=620,ptmxmode=000)
proc on /proc/asound type proc (ro,nosuid,nodev,noexec,relatime)
proc on /proc/bus type proc (ro,nosuid,nodev,noexec,relatime)
proc on /proc/fs type proc (ro,nosuid,nodev,noexec,relatime)
proc on /proc/irq type proc (ro,nosuid,nodev,noexec,relatime)
proc on /proc/sys type proc (ro,nosuid,nodev,noexec,relatime)
proc on /proc/sysrq-trigger type proc (ro,nosuid,nodev,noexec,relatime)
tmpfs on /proc/kcore type tmpfs (rw,nosuid,mode=755)
tmpfs on /proc/timer_stats type tmpfs (rw,nosuid,mode=755)
```



where .com meets .org  
Open Source events since 1996

# Cgroups

- Have only limited effect to isolation needs
- Restrict consumption of several ressources
  - RAM
  - CPU
  - I/O
  - Network bandwidth
- Can be useful if DoS scenarios are feasible
- Evaluate if ressource allocation interferes with overall system architecture



**where .com meets .org**  
Open Source events since 1996



# Threat 2: Insufficient Architecture



**where .com meets .org**  
Open Source events since 1996



# Daemon

- Docker manages all container operations by means of a permanent daemon with a REST-API.
- Uses Unix domain socket per default, runs as root user.
- Option -H binds daemon to a TCP port. Necessary for orchestration. Access control is important (both authorization and authentication).
- If not protected by SSL/TLS, exposing this port is dangerous.
- New project Rocket (rkt) addresses these issues, but is still in an early stage work in progress.



# Application Architecture

- Single containers are fun, but effectively neat toys
- Mature applications have resilience objectives:
  - Scale out (being able to deal with any number of containers)
  - High availability (failover if single containers die)
  - Load balancing (distribute workloads to idle containers)
  - Statelessness (so you can replace and upgrade any part of the application)
  - Separation of duties
- Application design, orchestration, and integration into a CI/CD pipeline are serious tasks by themselves



**where .com meets .org**  
Open Source events since 1996

# Security Architecture

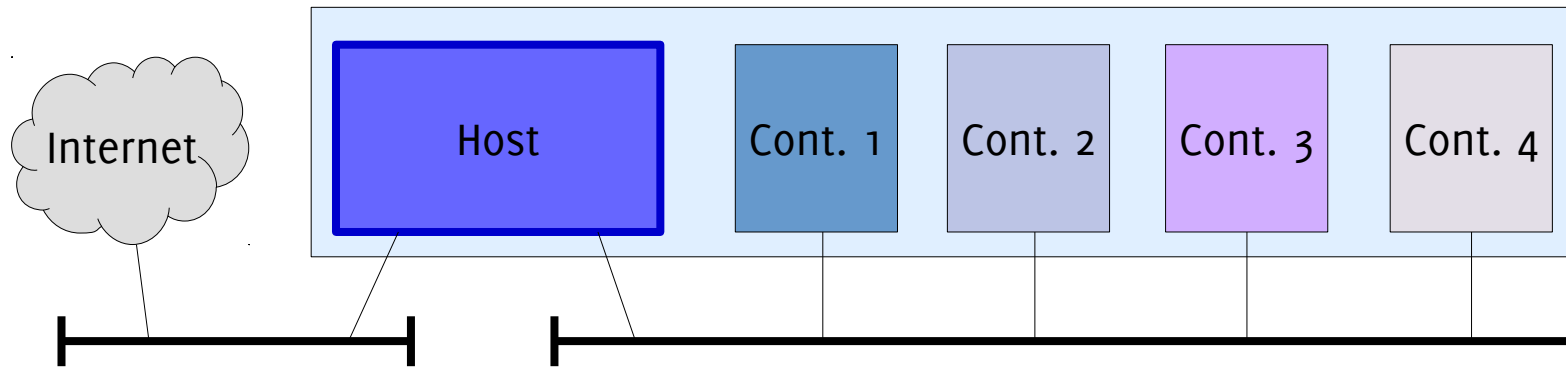
- Defense in Depth:  
Never rely on a single measure to protect your application and data
  - Container implement isolation
  - Firewalls provide additional network access
  - Encryptions protects data
- Single applications in a single container
- No need to SSH into containers in production (pet vs. cattle)



**where .com meets .org**  
Open Source events since 1996

# Network Security

- Per default, all containers share a common bridge.
- All container are thus part of the same segment/subnet: there is no special separation on the network layer between single containers.
- A host firewall is not sufficient to prevent from neighbor attacks:



- Communication of cluster management (etcd, k8s etc.) needs to be authentic and confidential (no default)



**where .com meets .org**  
Open Source events since 1996



# Threat 3: Unauthentic sources



**where .com meets .org**  
Open Source events since 1996



# Images

- Images are a convenient feature of Docker
- **Malware shipped in containers**
  - Necessity to check content
  - Experimental feature: signed images (since 1.3). Rkt has this built in
- **Malware shipped in packages**
  - Patchmanagement remains important
  - Update path for active containers (pet vs. Cattle)
- Run your own Repository? Don't store credentials!



**where .com meets .org**  
Open Source events since 1996

# CI pipeline

- **Build your own images (really easy with Dockerfiles)**
  - Version control for Dockerfiles
  - Integrate build into CI pipeline
  - Run your own image registry
  - Never have credentials inside your containers
- Validate sources for containers and packages



**where .com meets .org**  
Open Source events since 1996



# Threat 4: Turning off security features



**where .com meets .org**  
Open Source events since 1996

# What could possibly go wrong?

```
# sudo docker run --privileged=true -it ubuntu
```

- Containers are not „small desktops“
- Don't try to enable every feature that exists in a legacy distro
- Bind mounts can be nasty → use a data storage or object store
- There's no need to access raw hardware features inside containers



**where .com meets .org**  
Open Source events since 1996



# Four essential facts and tips about container security:

(1) A lot of security measures for isolation are built-in, more are to come.



(2) Containers protect against mistake, not abuse!


(3) Don't put all your eggs in one basket!

(4) Never rely on sources you don't trust!



where .com meets .org  
Open Source events since 1996





Questions about  
container security?  
– Send me an email:  
[magnus@linuxtag.org](mailto:magnus@linuxtag.org)

**Credits:**

Thanks for great footage under  
free and open licenses from Wikipedia, MorgueFile,  
Federal archive of Germany, Executive Office of the  
President of the United States, and Toronto Public  
Library.



**where .com meets .org**  
Open Source events since 1996



*Please*

# Remember to rate session

*Thank you!*

