

Grundlagen Docker

Christopher J. Ruwe

23. Januar 2020

1 Namespaces und Control Groups

Namespaces

Ein *Namespace*¹ ist ein Mechanismus, einer Prozessgruppe von $\{1 \dots p\} \mid p \in \mathbb{N}$ Prozessen eine Kapselung einer (globalen) Resource eines Systems (hier Betriebssystem) zur Verfügung zu stellen. Bei Betriebssystem-Namespaces wird dies vom Betriebssystem-*kernel*, konkret durch bestimmte *in-kernel* Datenstrukturen, vorgenommen.

Innerhalb dieser Kapselung haben die Prozesse einen aus Innen-Sicht exklusiven Zugriff auf die gekapselte Resource. Bei der gekapselten Resource kann es sich beispielsweise um eine *process ID* (PID) Hierarchie, eine separate *inter-process communication* (IPC) Domäne oder um eine separate Abbildung *user-* und *group-IDs* auf Berechtigungen handeln.

Damit erlaubt diese Kapselung eine quasi-Virtualisierung einer Resource durch den *kernel*, obwohl aus Außensicht die Resource von allen Prozessen (und damit Prozeßgruppen) des Systems geteilt wird.

DIESE PROZESS-ISOLIERUNG schließt so eine Lücke zwischen klassischen *time-sharing systems* und Virtualisierung durch einen Hypervisor.

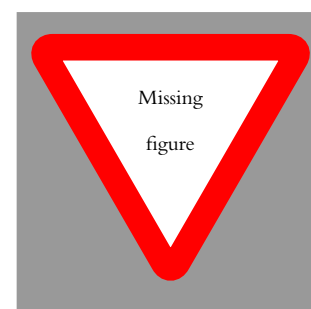
Bei *time-sharing* Systemen wie beispielsweise Unix oder Linux wird Prozessen durch den *kernel* Speicher und Rechenzeit zugeteilt. Die Prozesse sind wechselseitig „sichtbar“ und der Zugriff auf System-Ressourcen wird zwar durch den *kernel* gesteuert, ist grundsätzlich aber durch jeden Prozeß des Systems möglich.

Bei Virtualisierung durch einen Hypervisor werden die Ressourcen eines Computer-Systems verschiedenen Betriebssystem-Instanzen zugeteilt. Diese Instanzen haben alle eine separaten *kernel* und teilen wiederum ihren Prozessen die virtualisierten Ressourcen so wie auf einem nicht-virtualisierten System auch zu.

Prozess-Isolierung trennt durch *namespaces* stärker als klassische *time-sharing* Systeme. Alle isolierten Prozesse teilen sich jedoch den selben *kernel*, so daß weniger stark getrennt wird als durch Hypervisor-Virtualisierung. Die Trennung ist allerdings für sehr viele Anwendungsfälle hinreichend stark.

SEIT DEN SPÄTEN 90IGERN ist Prozess-Isolierung durch *namespaces* im produktiven Einsatz. Die bekanntesten Varianten sind FreeBSD Jails und Solaris Zones. Der Linux-*kernel* stellt beginnend mit den späten *kerneln* der 2.6er-Reihe eine Reihe von *namespaces* zur Verfügung.²

¹ *Linux Programmer's Manual: namespaces(7) - Overview of Linux Namespaces.* (2018). <http://man7.org/linux/man-pages/man7/namespaces.7.html>. Linux man-pages Project.



Formen der Prozess-Isolierung/-Virtualisierung.

time-sharing systems → Betriebssystem-Namespaces → Hypervisor I/II.

² `/proc/$$/ns`, wobei `$$` den aktuellen Prozess bezeichnet. Analog können *namespaces* für andere Prozesse unter `/proc` unter den jeweiligen PIDs eingesehen werden.

Die unter Linux bekannten Systeme zur Isolation von Prozessen LXC, LXK und auch docker nutzen eben diese Namespaces zur Kapselung von kompletten Betriebssystemen (ohne *kernel*, LXC, LXK) oder „lediglich“ von Applikationen (docker, rkt).

³ *Linux Programmer's Manual: pid_namespaces(7) - Overview of Linux PID Namespaces.* (2017). http://man7.org/linux/man-pages/man7/pid_namespaces.7.html. Linux man-pages Project.

Beware, here be dragons: Nicht jede Applikation eignet sich, *parent process* zu sein. Startet man beispielsweise einen Jenkins CI-Server mit PID 1 und *forkt* Jenkins aus der JVM Prozesse wie *make* oder *sh*, weil nicht-Java kompiliert wird und terminiert dieser Prozeß $\neq 0$, dann wird das geworfene Signal von der JVM oft nicht korrekt gefangen und die JVM terminiert. Hierzu benötigt man dann einen *init-Stub* wie *tiny* oder startet den Container mit *docker --init*, damit wenigstens Signale sauber gesammelt (*reaping*) werden.

⁴ *Linux Programmer's Manual: network_namespaces(7) - Overview of Linux Network Namespaces.* (2018). http://man7.org/linux/man-pages/man7/network_namespaces.7.html. Linux man-pages Project.

⁵ *Linux Programmer's Manual: svipc(7) - System V interprocess communication mechanisms.* (2016). <http://man7.org/linux/man-pages/man7/svipc.7.html>. Linux man-pages Project.; *Linux Programmer's Manual: mq_overview(7) - Overview of POSIX message queues.* (2017). http://man7.org/linux/man-pages/man7/mq_overview.7.html. Linux man-pages Project.

EIN *PID-namespace*³ isoliert *process* IDs ($1 \dots \$(\text{sysctl kernel.pid_max})$), so daß Prozesse in verschiedenen *PID-namespaces* aus Innen-Sicht, nicht aus Außensicht, kollisionsfrei die gleiche PID haben.

Der erste Prozess innerhalb eines solchen *namespaces* erhält die PID 1 analog zum klassischen *init*-Prozeß und wird zum *parent*-Prozeß aller nachfolgend aus diesem Prozeß *geforkten* Prozesse. Damit ist dieser Prozeß auch verantwortlich, gegebenenfalls Signale der *child*-Prozesse zu *reopen*. Wird dieser Prozeß terminiert, werden alle *child*-Prozesse in diesem *namespace* durch den *kernel* mit *SIGKILL* terminiert.

PID-namespaces können bis zu einer Tiefe von 32 in sich gekapselt (*nested*) werden. Die *process* ID von Prozessen ist für alle Prozesse innerhalb des *namespaces* und für alle Prozesse in „Vorgänger“-*namespaces* sichtbar. Deshalb können Operationen wie das senden eines Signals oder das „Einsteigen“ in einen *Namespace* „abwärts“, aber nicht „aufwärts“ erfolgen.

EIN *network namespace*⁴ isoliert Systemressourcen wie Netzwerk-*devices*, IPv4/6 *stacks*, *routing*-Tabellen, *firewall*-Regeln und den Konfigurationen des */proc/\$\$/net*-Verzeichnisses.

Ein physisches Netzwerk-*device* ist genau einem *network-namespace* zugeordnet. Virtuelle *ethernet-devices* werden immer als Paare in zwei verschiedenen *network-namespaces* erzeugt und können so als *tunnel* oder *bridge* dienen. Mit einem *veth* in einem virtuellen und einem physischen *namespace* beispielsweise kann über den Host ein *namespace* Zugriff auf angeschlossene Netzwerke des Hosts erhalten.

EIN *inter-process communication (IPC) namespace*⁵ isoliert SysV IPC Domänen, innerhalb derer Prozesse mit Signalen, *shared memory* und POSIX *message queues* kommunizieren können.

IPC *namespaces* erlauben damit, „Kommunikations-Verbünde“ von Prozessen zu erzeugen und, umgekehrt, Kommunikation zwischen Prozessen einzuschränken.

*mount namespaces*⁶ virtualisieren die Liste der *mount points*, die für Prozesse innerhalb dieses *namespaces* sichtbar sind. Ein *mount-namespace* gehört immer zu einem *user-namespace*.

Wird ein Prozess erzeugt, ist die initiale Liste der *mount points* die Liste des *parent* Prozesses. Nachfolgende *mount*- oder *umount*-Operationen innerhalb des *namespaces* verändern dann standardmäßig nur noch die Liste des Prozesses im jeweiligen *namespace*.

Es ist möglich, *mount-namespaces* mit verschiedenen Prozessgruppen zu teilen; in diesem Spezialfall können solche Operationen im *child-namespace* die *mounts* im *parent-namespace* verändern.

UNIX *time-sharing system namespaces* erzeugen separate *Identifizier* für *hostname* und *domainname*. Prozesse innerhalb eines UTS *namespaces* haben damit eigene Sichten auf die Datenstruktur *utsname*, die von *uname* zurückgegeben wird und gestatten „eigene“ System-Namen für derart isolierte Prozesse.

*User namespaces*⁷ erlauben separate Abbildungen von *user* und *group* IDs und damit Datei-Berechtigungen sowie *capabilities*⁸.

Damit können UID/GID eines Prozesses aus Innen- und Außensicht verschieden sein. Im häufigsten Fall kann so ein nicht-privilegiertes *user* (außen) Berechtigungen eines privilegierten users (*uid* == 0) innerhalb eines *namespaces* erhalten.

Wie schon *PID-namespaces* können *user namespaces* bis zu einer Tiefe von 32 in sich gekapselt (*nested*) werden.

*Control group namespaces*⁹ kapseln die Sicht auf die spezifischen *limits* und *quotas*, die dem Prozeß auf verschiedene Ressourcen-Typen gesetzt sind.¹⁰

Sie fügen der Prozess-Isolierung keine in dem Sinne „neue“ Funktion hinzu. Sie sind dennoch erforderlich, um zu verhindern, daß ein isolierter Prozeß über den Umweg seiner Sicht auf die Ressourcenkontrolle eine Sicht auf das kapselnde System erhält oder die *cgroups* seines *parent*-Prozesses steuern kann.

Control Groups

Ein *Control Group* (*cgroup*)¹¹ ist ein Mechanismus, für eine Prozessgruppe von $\{1 \dots p\} \mid p \in \mathbb{N}$ Prozessen den „Konsum“ von Systemressourcen wie *block device I/O*, *CPU time-slices* zu limitieren oder zu quotieren. So lassen sich Prozessgruppen Anteile dieser Ressourcen zuordnen. *cgroups* gehören immer zu einem bestimmten *kernel* Sub-System, im Wesentlichen *CPU*, *memory*, *block* und *network I/O*.

⁶ *Linux Programmer's Manual: mount_namespaces(7) - Overview of Linux Mount Namespaces*. (2018). http://man7.org/linux/man-pages/man7/mount_namespaces.7.html. Linux man-pages Project.

`/proc/$$/{mounts,mountinfo,mountstat}`

⁷ *Linux Programmer's Manual: user_namespaces(7) - Overview of Linux User Namespaces*. (2018). http://man7.org/linux/man-pages/man7/user_namespaces.7.html. Linux man-pages Project.

⁸ *Linux Programmer's Manual: capabilities(7) - Overview of Linux Capabilities*. (2018). <http://man7.org/linux/man-pages/man7/capabilities.7.html>. Linux man-pages Project.

⁹ *Linux Programmer's Manual: cgroup_namespaces(7) - Overview of Linux Cgroup Namespaces*. (2017). http://man7.org/linux/man-pages/man7/cgroup_namespaces.7.html. Linux man-pages Project.

¹⁰ `/proc/$$/{cgroup,mountinfo}`

¹¹ *Linux Programmer's Manual: cgroups(7) - Linux Control Groups*. (2018). <http://man7.org/linux/man-pages/man7/cgroups.7.html>. Linux man-pages Project.

MEHRERE CPU-CGROUP-CONTROLLER steuern den Konsum der Rechenleistung des Gesamtsystems:

Einer cgroup kann mit `CONFIG_CGROUP_SCHED` ein Minimal-Anteil an der CPU-Leistung eines Systems zugesichert werden, der dann, wenn das Rechner-System unter Last steht, für den Gebrauch der cgroup reserviert ist. Umgekehrt kann eine Obergrenze konfiguriert werden, die auch dann, wenn das System nicht belastet ist, nicht überschritten wird.

Ausgesteuert über `CONFIG_CGROUP_CPUACCT` kann die CPU-Nutzung der Prozesse einer cgroup abgerechnet werden. Unterschieden wird hier zwischen `user` (Zeit im *user-mode*) und `system` (Zeit im *kernel-mode*).

Mit `CONFIG_CPUSETS` können Prozesse „hart“ an spezifische Instanzen von Prozessoren, Prozessorkernen oder NUMA-Instanzen gebunden werden.

ÜBER DEN *cgroup memory-controller* kann der Konsum von *process* und *kernel-memory* sowie *swap* einer Prozessgruppe ausgesteuert werden. Der *hugetlb-controller* steuert die Nutzung von *huge pages* aus.

`net_prio` steuert Prioritäten (*per network-device*) aus. `net_cls` flaggt Pakete, die die cgroup verlassen, mit einer class-ID, die für firewall-Regeln oder für *traffic shaping* verwendet werden können.

`blkio` steuert den Bandbreiten-Verbrauch der übergebenen *block devices* durch entweder *time-slices* oder durch obere Grenzen der konsumierten Bandbreite.

ANDERE, SCHLECHTER ZU KLASSIFIZIERENDE *cgroup-controller* steuern die Erzeugung (`mknod`) oder den Zugriff auf *devices* (`devices`), die Fähigkeit zum `suspend/restore` (`freezer`), die Fähigkeit, *Monitoring-probes* auf Prozesse zu setzen (`perf_event`) und begrenzen die maximale Anzahl der Prozesse (`pids`)

Nutzung von Namespaces und Cgroups, Unterschiede zu anderen

Für das GNU/Linux-Betriebssystem existieren verschiedene Implementierungen von dort „Container“ genannten Mechanismen zur Isolation von Prozessen. Sie unterscheiden sich zum einen in der Konfigurierbarkeit der genutzten Isolations-Primitiven als auch im Anwendungsprofil.

`lxc/lxd`¹² container dienen der Provisionierung vollständiger Betriebssysteminstanzen (ohne kernel, `/sbin/init` hat PID 1) und erlauben Konfiguration der Parameter von namespace/cgroup-Isolation (bspw. *capabilities*).

Docker¹³ hingegen gestattet, die Isolierung aus verschiedenen *namespace*-Typen sehr flexibel zusammensetzen: Verschiedenen docker-container-Instanzen können bspw einen namespace wie den network-namespace oder den IPC-namespace teilen (so genutzt in Kubernetes) oder auch die Trennung vom Host-System ganz aufzuheben und beispielsweise den network-namespace des Hosts mitzunutzen (so genutzt für das `kubelet` einer Kubernetes *node*).

¹² *lxc(7) - Linux Containers*. (2019). <https://linuxcontainers.org/lxc/manpages/man7/lxc.7.html>. linux-containers.org; *LXD - System Container Manager*. (2019). <https://lxd.readthedocs.io/en/latest/>. linux-containers.org

¹³ Docker Inc. (2018). Docker Documentation. <https://docs.docker.com/>

AUS ANDEREN UNIXOIDEN SYSTEMEN sind unter Umständen Prozess-Isolations-Mechanismen wie FreeBSD Jails¹⁴ oder Solaris Zones¹⁵ bekannt. Für diese sind die zu nutzenden *namespace*-Primitiven fest. Es besteht keine Möglichkeit, dies zu konfigurieren und beispielsweise mehrere Zonen die selbe IPC-Kommunikations-Domäne gemeinsam zu nutzen und darüber Information auszutauschen.

Sowohl Jails als auch Zonen verfolgen einen anderen Einsatzzweck als *docker*-container. Dies stellt keinen Makel dar, ist aber für den Vergleich zu berücksichtigen. *docker* verfolgte von Anfang an den Zweck, verteilte Systeme auch aus eng anstatt nur lose gekoppelten Komponenten zu ermöglichen. Jails und Zonen sollten immer in sich geschlossene Gesamtsysteme bereitstellen und allenfalls über das Netzwerk zu verteilten Systemen zusammengefügt werden.

docker

docker zu betreiben bedeutet, Applikations-Container zu betreiben. Ein Applikations-Container ist ein durch die beschriebenen *namespace*-Mechanismen gekapselter Prozeß, wo das zugehörige Programm als eine *appliance*, als ein Image transportiert wird. Ein solches Image enthält alle Dateien (Programmtext und Bibliotheken), die zum Betrieb des Programmes benötigt werden. Oft enthalten Images allerdings erheblich mehr.

STRENG GENOMMEN müssen lediglich der Programmtext der Applikation, alle dynamisch gebundenen libraries und die Konfiguration vorliegen.

```

1 ldd /usr/bin/emacs25 | tr -d '\r' | sort
2 /lib64/ld-linux-x86-64.so.2 (0x00007ffa1de9a000)
3 libacl.so.1 => /lib/x86_64-linux-gnu/libacl.so.1 (0x00007ffa19476000)
4 libasound.so.2 => /usr/lib/x86_64-linux-gnu/libasound.so.2 (0x00007ffa1a430000)
5 libatk-1.0.so.0 => /usr/lib/x86_64-linux-gnu/libatk-1.0.so.0 (0x00007ffa15ea6000)
6 <...>
   /usr/bin/emacs25

```

Ein Image kann man sich dann als ein komprimiertes Archiv analog einer *.tar.gz* Datei vorstellen, das alle benötigten Dateien enthält. Dieses Image wird dann auf der Maschine, auf der der Container instanziiert werden soll, entpackt. Um die Menge der im Image transportierten Dateien zu reduzieren und damit die Übersichtlichkeit zu erhöhen, ist im „idealen“ *docker*-container nur ein einziges, statisch gebundenes *binary* enthalten.

In der Realität jedoch findet man in images oft vollständige Betriebssysteme (ohne *kernel*), die allenfalls in einer minimalen Variante vorliegen (*debian:jessie-slim*). Es ist unbestritten erheblich einfacher und weniger fehleranfällig, Programme und deren Abhängigkeiten mit dem Paketmanager einer Betriebssystem-Distribution zu installieren. Abhängigkeiten manuell aufzulösen und anderweitig in das Image zu kopieren stellt eine in vielen Fällen zu vermeidende Re-Implementierung dar.

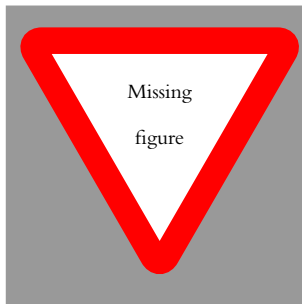
¹⁴ Kamp, P.-H. & Watson, R. N. M. (2000). Jails: Confining the Omnipotent Root. In *In Proc. 2nd Intl. SANE Conference; FreeBSD System Manager's Manual - jail(2)*. (2012). FreeBSD Project; *FreeBSD System Manager's Manual - jail(8)*. (2012). FreeBSD Project

¹⁵ Price, D. & Tucker, A. (2004). Solaris Zones: Operating System Support for Consolidating Commercial Workloads. In *Proceedings of 18th Large Installation System Administration Conference (LISA)* (S. 241–254); *Oracle Solaris 11.4 Reference Library - zones(7)*. (2016). https://docs.oracle.com/cd/E88353_01/html/E37853/zones-7.html. Oracle, Inc.

Listing 1: Auflösen der Bibliotheken, gegen die eine *binary* dynamisch gebunden wurde.

VOR INSTANZIIERUNG EINES CONTAINERS wird das Image auf das Zielsystem kopiert und entpackt. Bei Instanziierung wird hieraus eine Kopie erstellt und wird diese Kopie als Wurzelverzeichnis in den mount-namespace des Containers gelegt. Die konkrete Form der Kopie wird durch die Wahl des Storage für den docker-daemon, den *graph-volume driver* bestimmt.

Wählt man *vfs* als Storage, handelt es sich um eine tatsächliche Kopie. Dies ist eine schreib-intensive und deshalb langsame Methode, die für produktive Umgebungen nach Möglichkeit vermieden wird. Die Treiber *overlay* und *overlay2* ermöglichen es, das Wurzelverzeichnis des Containers aus einem oder mehreren *overlay-mounts* zusammenzusetzen. Dies ist erheblich weniger schreib-intensiv und daher um ein vielfaches performanter. Mit dem *devicemapper* kann man analog aus einem *logical volume manager* logische volumes konsumieren. Docker-Treiber für moderne Dateisysteme wie *ZFS* oder *btrfs* erlauben es, Images als separates Dateisystem anzulegen und im Bedarfsfall einen *snapshot* zu *clonen*.



Layered Images.

EINE SPEZIELLE EIGENSCHAFT VON IMAGES rechtfertigt diesen „Aufwand“. Images können gestapelt werden. Die Idee dabei ist, gemeinsame, „basische“ Anteile der paketierten appliance (Applikation + Libraries) als gemeinsame, identische Grundlage und damit das korrespondierende Image für mehr als eine appliance zu verwenden. Erst Differenzen im Aufbau von Appliances verlangen dann eigene, separate Images.

Dieser Aufbau läßt sich bei den *overlay*-Treibern schön beobachten:

Listing 2: Aufbau eines Docker-Images aus mehreren *Layers*.

```
docker inspect $SOMEDOCKER \
  | jq '.[].GraphDriver'
```

zeigt in Listing 3 die Schichtung eines Containers aus dem *nginx:latest* image, in Listing 4 die erheblich klarere, inhaltlich aber gleiche Ansicht bei Nutzung von *ZFS*. Storage wird hier als dataset in einer Art „verketteten Liste“ von *snapshots* und *clones* aufgebaut.

Die als *lower dir* angegebenen Verzeichnisse sind in absteigender Folge „übereinander geschichtet“ und unter dem Pfad *mergedDir* liegt das vollständig aufgebaute Wurzelverzeichnis der Appliance.

Listing 3: Aufbau der Storage-Schicht eines Docker-Containern bei overlay2. Auszug, kein legales JSON, dafür lesbar.

```

"GraphDriver": {
2  "Data": {
    "LowerDir": "/var/lib/docker/overlay2/e6daa<...>-init/diff
4          :/var/lib/docker/overlay2/efe96<...>/diff
          :/var/lib/docker/overlay2/01f84<...>/diff
6          :/var/lib/docker/overlay2/aca66<...>/diff
          :/var/lib/docker/overlay2/cc55a<...>/diff
8          :/var/lib/docker/overlay2/35a77<...>/diff
          :/var/lib/docker/overlay2/1f882<...>/diff
10         :/var/lib/docker/overlay2/6e29c<...>/diff
          :/var/lib/docker/overlay2/522d6<...>/diff
12         :/var/lib/docker/overlay2/d2d57<...>/diff
          :/var/lib/docker/overlay2/a3e53<...>/diff
14         :/var/lib/docker/overlay2/17304<...>/diff
          :/var/lib/docker/overlay2/67639<...>/diff
16         :/var/lib/docker/overlay2/dc8ad<...>/diff
          :/var/lib/docker/overlay2/7c9b8<...>/diff
18         :/var/lib/docker/overlay2/44c5f<...>/diff
          :/var/lib/docker/overlay2/bf38e<...>/diff
20         :/var/lib/docker/overlay2/e68db<...>/diff
          :/var/lib/docker/overlay2/ac010<...>/diff",
22  "MergedDir": "/var/lib/docker/overlay2/e6daa<...>/merged",
    "UpperDir": "/var/lib/docker/overlay2/e6daaf<...>/diff",
24  "WorkDir": "/var/lib/docker/overlay2/e6daa<...>/work"
  },
26  "Name": "overlay2"
}

```

Listing 4: Aufbau der Storage-Schicht eines Docker-Containern bei zfs.

```

{
2  "Data": {
    "Dataset": "rpool/var/lib/docker/647eb<...>",
4    "Mountpoint": "/var/lib/docker/zfs/graph/647eb8<...>"
  },
6  "Name": "zfs"
}

```

2 Übungen mit dem docker CLI

```

1 FROM debian:stretch
2
3 ENV DEBIAN_FRONTEND=noninteractive
4 RUN \
5     apt update \
6     && apt upgrade --yes \
7     && apt install \
8         --no-install-recommends \
9         --yes \
10        apt-transport-https \
11        ca-certificates \
12        curl \
13        gnupg \
14        less \
15        procps \
16    && apt clean \
17    && apt autoremove \
18    && rm -R /var/lib/apt/lists/* \
19    && rm -R /var/cache/apt

```

Listing 5: Aufbau eines Docker-Images aus einer Basis zzgl. der Installation von Tools.

¹⁶ In gut ausgebauten Build-Umgebungen wie einer CI/CD-Pipeline würde man vor jedem Build die Dockerfiles mit einem sogenannten Linter ¹⁷ prüfen. Damit werden Dockerfiles einer statischen Code-Analyse unterworfen, die bekannte Fehlerquellen (wie nicht-Löschen der Caches) moniert und daraus einen Fehler-Report erzeugt. Listing 6: Image Bau

¹⁷ github.com/hadolint. (2018). Haskell Dockerfile Linter. <https://github.com/hadolint>

¹⁸ Namespaces zu erzeugen erfordert root-Rechte.

Weitere Feinheiten des *builds* werden im Kurs über CI/CD Pipelines behandelt.

Listing 7: Image Bau FROM: scratch

```

1 FROM scratch
2
3 COPY <binary> /
4 COPY <config> /config.yaml
5
6 CMD ["/binary", "-c", "/config.yaml"]

```

Einige wenige praktischen Übungen sollten dies weiter vertiefen. Im ersten Schritt werden wir aus Dockerfile's Images bauen und damit Programme als Appliance paketieren und ausführen. Im zweiten Schritt werden wir kurz auf das Management des docker hosts und von Storage und Networking eingehen.

Images erzeugen

Ausgangspunkt eines builds ist üblicherweise eine Dockerfile. Diese Datei enthält in einer auf das Wesentliche beschränkten Syntax batch-ähnliche Aufrufe, wie die *layer*, aus denen letztendlich die Appliance aufgebaut wird, zu befüllen sind.

Dies Anweisungen in Listing 5 bauen auf einem Basis-Image auf (gegeben durch das statement FROM:) auf, setzen eine Umgebungsvariable (ENV), die den Paketmanager in den nicht-interaktiven Modus versetzt und installieren einige Pakete. Weil jedes Statement für sich ein separates *layer* erzeugt, und jedes *layer* für sich atomar abgeschlossen ist, löscht man nach Aufrufen des Paketmanagers die Caches. ¹⁶ Anderenfalls blieben die Caches bestehen und würden durch alle nachfolgenden Build-Schritte, die auf diesem Image aufbauen, mitgeschleppt.

Aus einer solchen Dockerfile wird das Image wie in Listing 6 erzeugt. ¹⁸

```

1 sudo docker build \
2     -f Dockerfile \
3     -t [<dockerregistry>/][<organization>/]<imagename>:[<version>|latest] \
4     <context_root_directory>

```

Die flags *-f* und *-t* erklären sich aus dem Beispiel. Es wird immer der sogenannte *build context*, oft *\$PWD* übergeben. Aus diesem *build context* bezieht der *dockerd* Dateien, die während des *builds* durch Primitiven wie *COPY* und *ADD* verwendet werden. Mit dem Kommando `docker tag -t <THETAG> <THEIMAGEHASH>` kann der Name eines Images im Nachhinein festgelegt werden.

In einer „leeren“ Installation wir der Befehl `docker images` (alte Form) oder `docker image ls` eine leere Ergebnismenge liefern. Nach erfolgreichem Build sollte das Image mit dem als Argument von *-t* angegebenen *tag* angezeigt werden.

Listing 7 zeigt die „reine Lehre“ (Leere?) für den Aufbau eines docker images. In einer wirklichen Micro-Service-Architektur befindet sich in

einem Docker-Container nur eine einzelnes *binary* (notwendigerweise dann statisch gebunden) und einige Konfigurationsdateien, pfiffigerweise vielleicht noch `/etc/passwd`, die Wurzelzertifikate und die *timezone* Definitionen.

Das Basis-Image `scratch` ist dabei ein „besonders“, nämlich leeres Image. So erzeugt man mit `docker` kleine und damit leicht transportable Images. Als ein Nebeneffekt erleichtert dies die Verwundbarkeits-Analyse ganz erheblich.

Alternativ zu der Methode `FROM: scratch` lässt sich mit dem Aufruf von `docker import` ein Image aus einem tar-ball erzeugen. Der Inhalt des tar-balls entspricht dann dem Inhalt des Wurzelverzeichnis. So könnte ein `docker-image` alternativ wie in Listing 8 erstellt werden.¹⁹ Weiterhin alternativ können Images auch mit `docker save` und `docker load` oder `docker export` in tar-balls importiert und exportiert werden. `docker commit` erzeugt neue

```

1 WHEN=$(date +%s)
2 FAMILY=ubuntu
3 DISTRO=bionic
4 TARGET=/${HOME}/images/${FAMILY}/${DISTRO}/${WHEN}
5 TAG="${FAMILY}-${DISTRO}-${WHEN}"
6
7 sudo debootstrap \
8 --merged-usr \
9 --variant=minbase \
10 --verbose ${TARGET} \
11 https://<aptproxy>:<port>/repository/apt-${DISTRO}
12
13 sudo tar \
14 -cPf $TAG.tar \
15 -C $TARGET \
16 ${PWD}
17
18 sudo docker import \
19 $TAG.tar \
20 $TAG

```

Images aus den Änderungen in einem Container. Dies wird im Aufbau von Images aus `Dockerfiles` durch den `dockerd` angewandt: Es wird aus der Basis ein Container erzeugt, im Container das Kommando aus der Dockerfile ausgeführt, `committed` und beendet. Dann folgt das nächste Kommando analog.²²

Container zu erzeugen, drin rumzurühren und dann manuell zu `commit` ist für mich persönlich ein Hinweis, daß *application state* nicht sauber isoliert ist oder die Installation einer Applikation interaktiv und daher schwer zu reproduzieren ist. Beides halte ich für *A Bad Thing™*.

Mit dem Kommando `docker history` lässt sich wie in Listing 9 gezeigt nachvollziehen, wie ein Image entstanden ist. Das Image für `docker-compose`²³ wurde offenbar von Linux Alpine²⁴ abgeleitet, in das danach das *binary*²⁵ kopiert wurde. Jede Zeile entspricht einem *layer* (Erinnerung: Listing 3) und man kann erkennen, daß jede Zeile einer `Dockerfile` zu einem separaten *layer* führt. Im Gegensatz dazu hat ein Image, das direkt aus einem tar-ball erzeugt wurde (`debootstrap` uä), nur ein *layer*.

¹⁹ Als Alternative zu eigenen Scripten könnte `hashicorp/packer`²⁰ oder ein dediziertes, prinzipiell gleich funktionierendes Build-Tool wie das `docker-maven-plugin`²⁰ verwendet werden.

²⁰ `hashicorp/packer`. (2019). Hashicorp Packer - Docker Builder. <https://www.packer.io/docs/builders/docker.html>

²¹ github.com/fabric8io/docker-maven-plugin. (2019). Maven plugin for running and creating Docker images. <https://github.com/fabric8io/docker-maven-plugin>

Listing 8: Bau eines Images aus einem tar-ball

²² Während eines *builds* kann man mit `watch sudo docker container ls` oder `watch sudo docker ps` beobachten, wie Container erzeugt und wieder beendet werden.

²³ `docker-compose` dient dem deklarativen Erzeugen von Kompositen aus mehreren Docker-Container und wird später (??) behandelt.

²⁴ Alpine Linux Development Team. (2019). Alpine Linux. Small. Simple. Secure. <https://alpinelinux.org>

²⁵ Python, paketierte mit `pyinstaller`²⁶, *dynamically linked, stripped*.

²⁶ Pyinstaller Development Team. (2019). Pyinstaller. <https://pyinstaller.org>

Listing 9: Untersuchung des Aufbaus eines Images

```

docker history docker/compose:1.21.2
2
IMAGE CREATED CREATED BY SIZE
4 86bf85b58bc8 9 months ago /bin/sh -c #(nop) ENTRYPOINT ["docker-compo..."] 0B
<missing> 9 months ago /bin/sh -c #(nop) COPY file:04bb859a7ea360f4... 10.9MB
6 <missing> 10 months ago /bin/sh -c apk update && apk add --no-cache ... 29.8MB
<missing> 10 months ago /bin/sh -c #(nop) ENV DOCKERBINS_SHA=1270dc... 0B
8 <missing> 10 months ago /bin/sh -c #(nop) ENV GLIBC=2.27-r0 0B
<missing> 12 months ago /bin/sh -c #(nop) CMD ["/bin/sh"] 0B
10 <missing> 12 months ago /bin/sh -c #(nop) ADD file:6edc55fb54ec9fc36... 3.96MB

```

Images transportieren

²⁷ Private Registries verlangen einen login vor der Nutzung, `docker login` und `docker logout`.

²⁸ Die bekanntesten sind der sogenannte Docker Hub ²⁹ und quay.io. ³⁰

²⁹ Docker Inc. (2019c). dockerhub. Build and Ship any Application Anywhere. <https://hub.docker.com>

³⁰ CoreOS, Inc. (2019). Quay [builds, analyzes, distributes] your container images. <https://quay.io>

³¹ Elastic Container Registry bei Amazon AWS, Azure Container Registry bei Microsoft und Google Container Registry bei Google. Einfallslos-langweilige Namen sind ein sicheres Kennzeichen des *mainstreams*.

³² Ausführliche Dokumentation: `docker pull(1)` - Pull an image or a repository from a registry. (2018). <https://docs.docker.com/engine/reference/commandline/pull/>. Docker Inc.

Funktional wäre es wohl unschädlich, aber üblicherweise möchte man nicht vor jeder Benutzung Images bauen. Stattdessen nutzt man sogenannte docker registries. Diese registries können privat ²⁷ oder öffentlich sein, verschiedene Firmen bieten diesbezüglich SaaS-Produkte an²⁸ und eigentlich alle Cloud-Provider haben Vergleichbares im Portfolio.³¹

Docker Images werden explizit (`docker pull nginx`) oder implizit (`docker run nginx`) gezogen. Mit der ersten Form werden die *layer* von der gewählten Docker Registry heruntergeladen und auf der lokalen Maschine gespeichert. Gleichzeitig wird, weil durch die fehlende Versionsangabe implizit `nginx:latest` verlangt wird, geprüft ob ein „jüngeres“ `:latest` verfügbar ist und dann, wenn, das Update durchgeführt. Mit der zweiten Form werden Images nur dann, wenn sie lokal nicht vorliegen, gezogen, ein implizites Update findet also nicht statt.³²

Wir erinnern uns an die Form des Namens eines Docker-Images aus dem Aufruf eines Builds (Listing 6, Seite8):

```
(-t): [<docker-registry>/][<organization>/]<imagename>[:<version>|latest]
```

`<docker-registry>` ist optional, entfällt dieser Teil des Arguments, wird angenommen, daß von der docker-registry hub.docker.io gezogen wird. Entfällt die `<organization>`, wird implizit eine spezielle, „*top-level*“ Organisation angenommen. Wird die Version nicht übergeben, wird automatisch `:latest` gezogen. Damit ist `docker pull nginx` (*at the time of writing*) eine Kurzform von `docker pull docker.io/nginx:1.15.8` Umgekehrt wird mit dem Aufruf von

```
docker push my.company.org/mydepartment/myimage:latest
docker push my.company.org/mydepartment/myimage:v1.2.3
```

ein Image in eine Docker Registry hochgeladen.

Images können mit `docker rmi <IMAGETAG>|<IMAGEHASH>` oder analog mit `docker image rm` gelöscht werden. Mit `docker search` können Images, allerdings nur von `docker.io` über die Kommandozeile gesucht werden. In aller Regel ist davon abzuraten, binaries ohne Herkunftsnachweis zu verwenden. Im Regelfall finden sich über die Webinterfaces der Registries Verweise auf das Quellcode-Repository³³, wo die Dockerfiles verwaltet und gepflegt werden. Dies erlaubt zum einen Kontrolle, zum anderen auch niederschwellige Möglichkeiten der Anpassung für den eigenen Bedarf.

³³ zB github.com/docker-library. (2019). Docker Official Images. <https://github.com/docker-library/official-images>

Container erzeugen und bedienen

Linux container und damit auch docker container nutzen die Isolations-Mechanismen des kernels (namespaces), um einen Prozess oder eine Prozessgruppe bezüglich der Nutzung spezifischer Ressourcen vom Rest des Systems zu trennen. Ein kein-Prozeß kann nicht in namespaces gestellt werden, weswegen zur Erzeugung eines Containers erst einmal ein Prozeß erzeugt werden muß.

Ein in einem container gekapselter Prozeß wird mit

```
docker run <imagename>[:<imageversion>] [program]
```

erzeugt. Dabei wird das Programm [program], wenn es übergeben wurde, gestartet, sonst das Programm unter CMD. Existiert ein ENTRYPOINT oder wurde docker mit

```
docker run \
  [--entrypoint <path>] \
  <imagename>[:<imageversion>] [program]
```

aufgerufen, wird stattdessen dieses Programm gestartet.

Der sogenannte ENTRYPOINT dient dem Start *vorbereitender* Programme, die „Hilfstätigkeiten“ wie das Setzen von *permissions* uÄ ausführen. Üblicherweise ruft der ENTRYPOINT danach das „eigentliche“ Programm auf, oft, nicht eine bash als parent-Prozeß stehen zu lassen, mit dem Aufruf von *exec*.

Möchte man zu Test-Zwecken bspw. ein aktuelles Ubuntu starten, könnte man mit

```
docker run \
  --interactive \
  --tty \
  ubuntu:18.04 /bin/bash
```

eine */bin/bash* in einem Container startet. Alle *binaries* würden dann in einem aus Ubuntu abgeleiteten Image bezogen. Hätte man lieber ein Image aus der RedHat-Familie, setze man statt *ubuntu:18.04* *centos:7.5* oder *fedora:27* ein.

Der Container könnte auch ohne *-tty* und *--interactive* gestartet werden.³⁴ Dies wäre allerdings für den Fall einer */bin/bash* wenig sinnvoll: Das Flag *--interactive* hält *STDIN* offen und *--tty* klebt *STDIN* des containers an ein (Pseudo)-Terminal. Ohne *--tty* kann man den Container starten und sogar extrem rudimentäre Eingaben (*ls anyone*) starten, Signale über das Terminal zu senden (*C^C*) ist aber nicht möglich (der container muß danach mit *docker rm -vf <containerId>* terminiert werden).

Ohne *--interactive* kann man den container ebenfalls starten, man sieht dann sogar den Eingabe-Prompt. Eingaben kann man wenigstens auf das Terminal leiten, allein, es nützt nichts: Sie gelangen im container nicht zur Ausführung. Erst die Kombination von *--interactive --tty* erlaubt eine „normale“ Interaktion mit der gestartet shell. Um einen nicht-interaktiv zu nutzenden Container wie *nginx* zu starten, benötigt man beides nicht.

Es haben sich hier mittlerweile verschiedene Geschmäcker herausgebildet, es kommt auch vor, daß einige mit ENTRYPOINT das „eigentliche“ Programm starten.

³⁴ Docker Inc. (2019b). Docker run reference. <https://docs.docker.com/engine/reference/run/>

Kommt der treiber `vfs` zum Einsatz, dann handelt es sich um Voll-Kopien und dann *wird* Platz konsumiert. Dies ist allerdings ein Spezialfall, aus Performanz-Gründen wird `vfs` produktiv nicht verwendet.

Ruft man `docker ps -a`, läßt man sich alle Container inklusive aller terminierten Container anzeigen. Container nach Abschluß stehen zu lassen hat Nebeneffekte, die man bei sorgsamem Betrieb besser vermeidet: Das Verzeichnis mit den gemergten binaries bleibt stehen und der Container kann wieder gestartet werden. Beides bricht mit dem Prinzip, daß container keinen persistenten Zustand haben sollen. Auch wenn bei Verwendung von aktuellen storage Mechanismen wie `overlay(2)`, `btrfs`, `zfs` oder `devicemapper` kein zusätzlicher Plattenplatz konsumiert wird, verdeckt dies dennoch das System.

Weil man das Löschen terminierter Container meistens dann doch vergisst, ruft man `docker` sinnvollerweise mit

```
docker run \
--name mybeautifulshellcontainer
--interactive \
--rm \
--tty \
ubuntu:18.04 /bin/bash
```

wodurch der Container nach Abschluss automatisch gelöscht wird. Containern Namen zu geben erhöht, ohne funktional erforderlich zu sein, die Übersicht ganz erheblich.

Möchte man nicht alle terminierten Container einzeln löschen, empfiehlt sich der Aufruf von `docker container prune`, der ohne weitere Unterscheidung alle terminierten Container abräumt.

BENÖTIGT MAN FÜR CONTAINER PERSISTENZ, so empfiehlt es sich, dies entweder in separaten `docker-volumes` vorzuhalten oder mit sogenannten *bind-mount*, `--volume=[[HOST-DIR:]CONTAINER-DIR[:OPTIONS]]` persistenten Speicher vom `docker-host` im `container` verfügbar zu halten. Beliebige viele `--volume` Parameter können übergeben werden, es ist jedoch darauf zu achten, nicht durch in der Reihenfolge später übergebene frühere zu überschreiben.

```
docker run \
--interactive \
--rm \
--tty \
--volume=<onhost>:<indocker>
ubuntu:18.04 /bin/bash
```

Mit den Optionen `:rw` (default) oder `:ro` kann angesteuert werden, ob das Volume *read-write* oder *read-only* gemounted werden soll. Weitere *tunables* sind in der Dokumentation³⁵ oder der `man`-page beschrieben.

³⁵ Docker Inc. (2019d). Manage data in Docker. <https://docs.docker.com/storage/volumes/>

OHNE NETZWERKANBINDUNG kann Software nicht gut genutzt werden. Einen webserver wie `nginx` von Netzwerk abgekoppelt zu nutzen ist wenig sinnvoll. Startet man bspw. `nginx` ohne gesonderte Parametrisierung der Konnektivität,

```
docker run \
--rm \
--volume=<somepath>/wwwdata:/var/www/html \
nginx
```

kann man bspw. mit `docker inspect` und `jq` dessen Anbindung untersuchen.

```
docker inspect <name|id> \
| jq '.'
docker inspect <name|id> \
| jq '.[].NetworkSettings.IPAddress'
```

Mit der so ermittelten IP kann man tatsächlich unter Port 80 zugreifen. Auf Hosts Dienste mit RFC1918³⁹ Adressen zu exponieren, die uU in nicht in den Unternehmens-Kontext integrierten Netzen stehen, wird allerdings wenig zielführend sein. Für die Öffentlichkeit bestimmte Dienste kann man so erst recht nicht betreiben.

Statt dessen exponiert man mit dem sogenannten *port-forwarding* Ports im Container auf dem Host. Wird `docker` mit dem Parameter `--publish` gerufen,

```
docker run \
--publish <HOSTPORT>:<CONTAINERPORT> \
--rm nginx
```

dann wird üblicherweise mit `iptables` wie in Listing 10 gezeigt der host port auf die IP des containers geNATed. Die Container-Configuration ist mit Listing 11 gezeigt.

```
sudo iptables-save \
2 | grep -E '8888|192.168.66.4|docker'

4 -A POSTROUTING -s 192.168.66.0/24 ! -o docker0 -j MASQUERADE
-A POSTROUTING -s 192.168.66.4/32 -d 192.168.66.4/32 -p tcp -m tcp --dport 80 -j
↪ MASQUERADE
6 -A DOCKER -i docker0 -j RETURN
-A DOCKER ! -i docker0 -p tcp -m tcp --dport 8888 -j DNAT --to-destination
↪ 192.168.66.4:80
8 -A FORWARD -o docker0 -m conntrack --ctstate RELATED,ESTABLISHED -j ACCEPT
-A FORWARD -o docker0 -j DOCKER
10 -A FORWARD -i docker0 ! -o docker0 -j ACCEPT
-A FORWARD -i docker0 -o docker0 -j ACCEPT
12 -A DOCKER -d 192.168.66.4/32 ! -i docker0 -o docker0 -p tcp -m tcp --dport 80 -j
↪ ACCEPT
-A DOCKER-ISOLATION-STAGE-1 -i docker0 ! -o docker0 -j DOCKER-ISOLATION-STAGE-2
14 -A DOCKER-ISOLATION-STAGE-2 -o docker0 -j DROP
```

Listing 10: Nutzung von `iptables` durch Docker.

Die Struktur des JSON aus Listing 11 zeigt es bereits: `NetworkSettings.Ports` referenziert ein Dictionary mehrerer Werte. `docker run` kann nämlich mehrere `--publish`-Parameter verarbeiten.⁴⁰

⁴⁰ Warum `.Ports.<no>/<proto>` ebenfalls ein array referenziert, sei dem geeigneten Leser anhand der Networking Dokumentation⁴¹ zur Übung empfohlen.
⁴¹ Docker Inc. (2019a). Container Networking. <https://docs.docker.com/config/containers/container-networking/>

```

docker inspect <name|id> \
| jq '.'
docker inspect <NGINXREVERSEPROXY> \
| jq '.[].NetworkSettings.Ports'

{
  "80/tcp": null,
  "8441/tcp": [
    {
      "HostIp": "0.0.0.0",
      "HostPort": "8441"
    }
  ],
  "8442/tcp": [
    {
      "HostIp": "0.0.0.0",
      "HostPort": "8442"
    }
  ],
  "8443/tcp": [
    {
      "HostIp": "0.0.0.0",
      "HostPort": "8443"
    }
  ]
}

```

Listing 11: Nutzung von iptables durch Docker.

Literatur

Literatur

- Alpine Linux Development Team. (2019). Alpine Linux. Small. Simple. Secure. <https://alpinelinux.org>
- CoreOS, Inc. (2019). Quay [builds, analyzes, distributes] your container images. <https://quay.io>
- Docker Inc. (2018). Docker Documentation. <https://docs.docker.com/>
- Docker Inc. (2019a). Container Networking. <https://docs.docker.com/config/containers/container-networking/>
- Docker Inc. (2019b). Docker run reference. <https://docs.docker.com/engine/reference/run/>
- Docker Inc. (2019c). dockerhub. Build and Ship any Application Anywhere. <https://hub.docker.com>
- Docker Inc. (2019d). Manage data in Docker. <https://docs.docker.com/storage/volumes/>
- docker pull(1) - Pull an image or a repository from a registry.* (2018). <https://docs.docker.com/engine/reference/commandline/pull/>. Docker Inc.
- FreeBSD System Manager's Manual - jail(2).* (2012). FreeBSD Project.
- FreeBSD System Manager's Manual - jail(8).* (2012). FreeBSD Project.
- github.com/docker-library. (2019). Docker Official Images. <https://github.com/docker-library/official-images>
- github.com/fabric8io/docker-maven-plugin. (2019). Maven plugin for running and creating Docker images. <https://github.com/fabric8io/docker-maven-plugin>
- github.com/hadolint. (2018). Haskell Dockerfile Linter. <https://github.com/hadolint>

- github.com/stedolan/jq. (2019). Command-line JSON processor. <https://github.com/stedolan/jq>
- Gregg, B. (2013). Virtualization Performance: Zones, KVM, Xen. <http://dtrace.org/blogs/brendan/2013/01/11/virtualization-performance-zones-kvm-xen/>.
- Gritton, J. (2010). Running the New FreeBSD Jails. http://2010.eurobsdcon.org/fileadmin/fe_user/jamie/qvo7qI.pdf. EuroBSDCon.
- hashicorp/packer. (2019). Hashicorp Packer - Docker Builder. <https://www.packer.io/docs/builders/docker.html>
- jq - Command-line JSON Processor(1)*. (2018). <http://manpages.ubuntu.com/manpages/bionic/man1/jq.1.html>. Canonical - Ubutu Manpage Repository.
- Kamp, P.-H. & Watson, R. N. M. (2000). Jails: Confining the Omnipotent Root. In *In Proc. 2nd Intl. SANE Conference*.
- Linux Programmer's Manual: capabilities(7) - Overview of Linux Capabilities*. (2018). <http://man7.org/linux/man-pages/man7/capabilities.7.html>. Linux man-pages Project.
- Linux Programmer's Manual: cgroup_namespaces(7) - Overview of Linux Cgroup Namespaces*. (2017). http://man7.org/linux/man-pages/man7/cgroup_namespaces.7.html. Linux man-pages Project.
- Linux Programmer's Manual: cgroups(7) - Linux Control Groups*. (2018). <http://man7.org/linux/man-pages/man7/cgroups.7.html>. Linux man-pages Project.
- Linux Programmer's Manual: clone(2) - Create a Child Process*. (2017). <http://man7.org/linux/man-pages/man2/clone.2.html>. Linux man-pages Project.
- Linux Programmer's Manual: mount_namespaces(7) - Overview of Linux Mount Namespaces*. (2018). http://man7.org/linux/man-pages/man7/mount_namespaces.7.html. Linux man-pages Project.
- Linux Programmer's Manual: mq_overview(7) - Overview of POSIX message queues*. (2017). http://man7.org/linux/man-pages/man7/mq_overview.7.html. Linux man-pages Project.
- Linux Programmer's Manual: namespaces(7) - Overview of Linux Namespaces*. (2018). <http://man7.org/linux/man-pages/man7/namespaces.7.html>. Linux man-pages Project.
- Linux Programmer's Manual: network_namespaces(7) - Overview of Linux Network Namespaces*. (2018). http://man7.org/linux/man-pages/man7/network_namespaces.7.html. Linux man-pages Project.
- Linux Programmer's Manual: pid_namespaces(7) - Overview of Linux PID Namespaces*. (2017). http://man7.org/linux/man-pages/man7/pid_namespaces.7.html. Linux man-pages Project.
- Linux Programmer's Manual: setns(2) - Reassociate Thread with a Namespace*. (2017). <http://man7.org/linux/man-pages/man2/setns.2.html>. Linux man-pages Project.
- Linux Programmer's Manual: svipc(7) - System V interprocess communication mechanisms*. (2016). <http://man7.org/linux/man-pages/man7/svipc.7.html>. Linux man-pages Project.

- Linux Programmer's Manual: unshare(2) - Disassociate Parts of the Process Execution Context.* (2018). <http://man7.org/linux/man-pages/man2/unshare.2.html>. Linux man-pages Project.
- Linux Programmer's Manual: user_namespaces(7) - Overview of Linux User Namespaces.* (2018). http://man7.org/linux/man-pages/man7/user_namespaces.7.html. Linux man-pages Project.
- Linux Programmer's Manual: veth(4) - Virtual Ethernet Device.* (2018). <http://man7.org/linux/man-pages/man4/veth.4.html>. Linux man-pages Project.
- Linux User Commands: nsenter(1) - Run Program with Namespaces of Other Processes.* (2013). <http://man7.org/linux/man-pages/man1/nsenter.1.html>. Linux man-pages Project.
- Linux User Commands: unshare(1) - Run Program with some Namespaces Unshared from Parent.* (2016). <http://man7.org/linux/man-pages/man1/unshare.1.html>. Linux man-pages Project.
- lxc(7) - Linux Containers.* (2019). <https://linuxcontainers.org/lxc/manpages/man7/lxc.7.html>. linux-containers.org.
- LXD - System Container Manager.* (2019). <https://lxd.readthedocs.io/en/latest/>. linux-containers.org.
- Moskowitz, R., Karrenberg, D., Rekhter, Y., Lear, E. & de Groot, G. J. (1996). Address Allocation for Private Internets. <https://rfc-editor.org/rfc/rfc1918.txt>
- Oracle Solaris 11.4 Reference Library - zones(7).* (2016). https://docs.oracle.com/cd/E88353_01/html/E37853/zones-7.html. Oracle, Inc.
- Price, D. & Tucker, A. (2004). Solaris Zones: Operating System Support for Consolidating Commercial Workloads. In *Proceedings of 18th Large Installation System Administration Conference (LISA)* (S. 241–254).
- Pyinstaller Development Team. (2019). Pyinstaller. <https://pyinstaller.org>
- Schenkeveld, P. (2010). Building Servers with NanoBSD, ZFS and Jails. In *Proceedings of AsiaBSDCon 2010*. <http://2010.asiabsdcon.org/papers/abc2010-P4A-paper.pdf>. Tokyo University of Science, Tokyo, Japan.
- Stephen Dolan. (2019). *jq Manual (development version)*. <https://stedolan.github.io/jq/manual/>.
- Teske, D. (2013). VIMAGE Jails on FreeBSD-8. <http://devinteske.com/vimage-jails-on-freebsd-8>.
- Tripathi, S., Droux, N., Belgaied, K. & Khare, S. (2009). Crossbow Virtual Wire: Network in a Box. In *Proceedings of the 23rd Conference on Large Installation System Administration (LISA)* (S. 4–4). LISA'09. Baltimore, MD: USENIX Association. <http://dl.acm.org/citation.cfm?id=1855698.1855702>
- Tripathi, S., Droux, N., Srinivasan, T. & Belgaied, K. (2009). Crossbow: From Hardware Virtualized NICs to Virtualized Networks. In *Proceedings of the 1st ACM Workshop on Virtualized Infrastructure Systems and Architectures* (S. 53–62). VISA '09. Barcelona, Spain: ACM. doi:10.1145/1592648.1592658. <http://doi.acm.org/10.1145/1592648.1592658>
- Zeeb, B. A. (2010). FreeBSD Jails - Notes Jotted on the Prison Wall. http://2010.eurobsdcon.org/fileadmin/fe_user/bz/ZDZZ8M.pdf. presented at EuroBSDCon 2010, Karlsruhe, Germany.

Zeeb, B. & Watson, R. (2010). The New VWorld – FreeBSD Jail Based Virtualization. http://www.bsdcn.org/2010/schedule/attachments/130_2010-bz-the-new-vworld.pdf. presented at BSDCan, Ottawa, Canada.