

原创

多选参数

于 2020-11-26 11:30:00 发布

189

收藏 6

版权

文章标签：

网络

docker

linux

java

操作系统

Namespace

Linux Namespace 是 Linux 提供的一种内核级别环境隔离的方法。这种隔离机制和 chroot 很类似，chroot 是把某个目录修改为根目录，从而无法访问外部的内容。Linux Namesapce 在此基础上，提供了对 UTS、IPC、Mount、PID、Network、User 等的隔离机制，如下所示。

分类	系统调用参数	相关内核版本
Mount Namespaces	CLONE_NEWNS	Linux 2.4.19
UTS Namespaces	CLONE_NEWUTS	Linux 2.6.19
IPC Namespaces	CLONE_NEWIPC	Linux 2.6.19
PID Namespaces	CLONE_NEWPID	Linux 2.6.19
Network Namespaces	CLONE_NEWNET	始于Linux 2.6.24 完成于 Linux 2.6.29
User Namespaces	CLONE_NEWUSER	始于 Linux 2.6.23 完成于 Linux 3.8)

★

Linux Namespace 官方文档：Namespaces in operation

”

namespace 有三个系统调用可以使用：

- clone() --- 实现线程的系统调用，用来创建一个新的进程，并可以通过设计上述参数达到隔离。
- unshare() --- 使某个进程脱离某个 namespace
- setns(int fd, int nstype) --- 把某进程加入到某个 namespace

下面使用这几个系统调用来演示 Namespace 的效果，更加详细地可以看 DOCKER基础技术：LINUX NAMESPACE（上）、DOCKER基础技术：LINUX NAMESPACE（下）。

UTS Namespace

UTS Namespace 主要是用来隔离主机名的，也就是每个容器都有自己的主机名。我们使用如下的代码来进行演示。注意：假如在容器内部没有设置主机名的话会使用主机的主机名的；假如在容器内部设置了主机名但是没有使用 CLONE_NEWUTS 的话那么改变的其实是主机的主机名。

```
1 #define _GNU_SOURCE
2 #include <sys/types.h>
3 #include <sys/wait.h>
4 #include <sys/mount.h>
5 #include <stdio.h>
6 #include <sched.h>
7 #include <signal.h>
8 #include <unistd.h>
9
10 #define STACK_SIZE (1024 * 1024)
```

目录

Namespace

- UTS Namespace
- PID Namespace
- IPC Namespace
- Mount Namespace
- User Namespace
- Network Namespace
- Namespace 情况查看

总结

巨人的肩膀

分类专栏

- 拿下数据结构与算法 2篇
- VUE 学习 5篇
- 杂文 3篇
- 可读代码编写炸鸡 9篇
- Redis 从入门到放弃 5篇
- Git 操作全系列 4篇
- 经典的 jQuery 2篇
- 爪哇 2篇

```

11 static char container_stack[STACK_SIZE];
12
13 char* const container_args[] = {
14     "/bin/bash",
15     NULL
16 };
17
18 int container_main(void* arg) {
19     printf("Container [%5d] - inside the container!\n", getpid());
20     sethostname("container_dawn", 15);
21     execv(container_args[0], container_args);
22     printf("Something's wrong!\n");
23     return 1;
24 }
25
26 int main() {
27     printf("Parent [%5d] - start a container!\n", getpid());
28     int container_id = clone(container_main, container_stack + STACK_SIZE,
29                             CLONE_NEWUTS | SIGCHLD, NULL);
30     waitpid(container_id, NULL, 0);
31     printf("Parent - container stopped!\n");
32     return 0;
33 }

```

```

root@dawn ~/docker_ws/mydocker
# ./mycontainer_uts
Parent [31904] - start a container!
Container [31905] - inside the container!
root@container_dawn:~/docker_ws/mydocker# hostname
container_dawn
root@container_dawn:~/docker_ws/mydocker# exit
exit
Parent - container stopped!
root@dawn ~/docker_ws/mydocker
# hostname
dawn

```

PID Namespace

每个容器都有自己的进程环境中，也就是相当于容器内进程的 PID 从 1 开始命名，此时主机上的 PID 其实也还是从 1 开始命名的，就相当于有两个进程环境：一个主机上的从 1 开始，另一个容器里的从 1 开始。

为啥 PID 从 1 开始就相当于进程环境的隔离了呢？因此在传统的 UNIX 系统中，PID 为 1 的进程是 init，地位特殊。它作为所有进程的父进程，有很多特权。另外，其还会检查所有进程的状态，我们知道如果某个进程脱离了父进程（父进程没有 wait 它），那么 init 就会负责回收资源并结束这个子进程。所以要想做到进程的隔离，首先需要创建出 PID 为 1 的进程。

但是，【kubernetes 里面的话】

```

1  int container_main(void* arg) {
2      printf("Container [%5d] - inside the container!\n", getpid());
3      sethostname("container_dawn", 15);
4      execv(container_args[0], container_args);
5      printf("Something's wrong!\n");
6      return 1;
7  }
8
9  int main() {
10     printf("Parent [%5d] - start a container!\n", getpid());
11     int container_id = clone(container_main, container_stack + STACK_SIZE,
12                             CLONE_NEWUTS | CLONE_NEWPID | SIGCHLD, NULL);

```

```
13     waitpid(container_id, NULL, 0);
14     printf("Parent - container stopped!\n");
15     return 0;
16 }
```

```
root@dawn ~/docker_ws/mydocker
# ./mycontainer_pid
Parent [32020] - start a container!
Container [ 1] - inside the container!
root@container_dawn:~/docker_ws/mydocker# echo $$
1
root@container_dawn:~/docker_ws/mydocker#
```

如果此时你在子进程的 shell 中输入 ps、top 等命令，我们还是可以看到所有进程。这是因为，ps、top 这些命令是去读 /proc 文件系统，由于此时文件系统并没有隔离，所以父进程和子进程通过命令看到的情况都是一样的。

IPC Namespace

常见的 IPC 有共享内存、信号量、消息队列等。当使用 IPC Namespace 把 IPC 隔离起来之后，只有同一个 Namespace 下的进程才能相互通信，因为主机的 IPC 和其他 Namespace 中的 IPC 都是看不到了。而这个的隔离主要是因为创建出来的 IPC 都会有一个唯一的 ID，那么主要对这个 ID 进行隔离就好了。

想要启动 IPC 隔离，只需要在调用 clone 的时候加上 CLONE_NEWIPC 参数就可以了。

```
1 int container_main(void* arg) {
2     printf("Container [%5d] - inside the container!\n", getpid());
3     sethostname("container_dawn", 15);
4     execv(container_args[0], container_args);
5     printf("Something's wrong!\n");
6     return 1;
7 }
8
9 int main() {
10     printf("Parent [%5d] - start a container!\n", getpid());
11     int container_id = clone(container_main, container_stack + STACK_SIZE,
12                             CLONE_NEWUTS | CLONE_NEWPID | CLONE_NEWIPC | SIGCHLD, NULL);
13     waitpid(container_id, NULL, 0);
14     printf("Parent - container stopped!\n");
15     return 0;
16 }
```

Mount Namespace

Mount Namespace 可以让容器有自己的 root 文件系统。需要注意的是，在通过 CLONE_NEWNS 创建 mount namespace 之后，父进程会把自己的文件结构复制给子进程中。所以当子进程中不重新 mount 的话，子进程和父进程的文件系统视图是一样的，假如想要改变容器进程的视图，一定需要重新 mount（这个是 mount namespace 和其他 namespace 不同的地方）。

另外，子进程中新的 namespace 中的所有 mount 操作都只影响自身的文件系统（注意这边是 mount 操作，而创建文件等操作都是会有所影响的），而不对外界产生任何影响，这样可以做到比较严格地隔离（当然这边是除 share mount 之外的）。

下面我们重新挂载子进程的 /proc 目录，从而可以使用 ps 来查看容器内部的情况。

```
1 int container_main(void* arg) {
2     printf("Container [%5d] - inside the container!\n", getpid());
3
4     sethostname("container_dawn", 15);
5
6     if (mount("proc", "/proc", "proc", 0, NULL) != 0) {
```

```

7         perror("proc");
8     }
9
10    execv(container_args[0], container_args);
11    printf("Something's wrong!\n");
12    return 1;
13 }
14
15 int main() {
16     printf("Parent [%5d] - start a container!\n", getpid());
17     int container_id = clone(container_main, container_stack + STACK_SIZE,
18                             CLONE_NEWUTS | CLONE_NEWPID | CLONE_NEWNS | SIGCHLD, NULL);
19     waitpid(container_id, NULL, 0);
20     printf("Parent - container stopped!\n");
21     return 0;
22 }

```

```

root@dawn ~/docker_ws/mydocker
# ./mycontainer_mount
Parent [32715] - start a container!
Container [ 1] - inside the container!
root@container_dawn:~/docker_ws/mydocker# ps -elf
F S UID          PID  PPID  C PRI  NI ADDR SZ WCHAN  STIME TTY          TIME CMD
4 S root           1      0  0  80   0 -  7100 wait   15:33 pts/25    00:00:00 /bin/bash
0 R root          11      1  0  80   0 - 11108 -      15:33 pts/25    00:00:00 ps -elf

```

★

这里会有个问题就是在退出子进程之后，当再次使用 `ps -elf` 的时候会报错，如下所示

```

root@dawn ~/docker_ws/mydocker
# ./mycontainer_mount
Parent [ 421] - start a container!
Container [ 1] - inside the container!
root@container_dawn:~/docker_ws/mydocker# ps -elf
F S UID          PID  PPID  C PRI  NI ADDR SZ WCHAN  STIME TTY          TIME CMD
4 S root           1      0  0  80   0 -  7100 wait   15:51 pts/25    00:00:00 /bin/bash
0 R root          11      1  0  80   0 - 11108 -      15:51 pts/25    00:00:00 ps -elf
root@container_dawn:~/docker_ws/mydocker# exit
exit
Parent - container stopped!
root@dawn ~/docker_ws/mydocker
# ps -elf
Error, do this: mount -t proc proc /proc
root@dawn ~/docker_ws/mydocker
# mount -t proc proc /proc

```

这是因为 `/proc` 是 share mount，对它的操作会影响所有的 mount namespace，可以看这里：

<http://unix.stackexchange.com/questions/281844/why-does-child-with-mount-namespace-affect-parent-mounts>

上面仅仅重新 mount 了 `/proc` 这个目录，其他的目录还是跟父进程一样视图的。一般来说，容器创建之后，容器进程需要看到的是一个独立的隔离环境，而不是继承宿主机的文件系统。接下来演示一个山寨镜像，来模仿 Docker 的 Mount Namespace。也就是给予进程实现一个较为完整的独立的 root 文件系统，让这个进程只能访问自己构成的文件系统的内容（想想我们平常使用 Docker 容器的样子）。

首先我们使用 `docker export` 将 `busybox` 镜像导出成一个 `rootfs` 目录，这个 `rootfs` 目录的情况如图所示，已经包含了 `/proc`、`/sys` 等特殊的目录。

```

root@dawn ~/docker_ws/mydocker
# tree rootfs -L 1
rootfs

```

```

— bin
— busybox.tar
— dev
— etc
— home
— mnt
— proc
— root
— run
— sys
— tmp
— usr
— var

12 directories, 1 file

```

之后我们在代码中将一些特殊目录重新挂载，并使用 `chroot()` 系统调用将进程的根目录改成上文的 `rootfs` 目录。

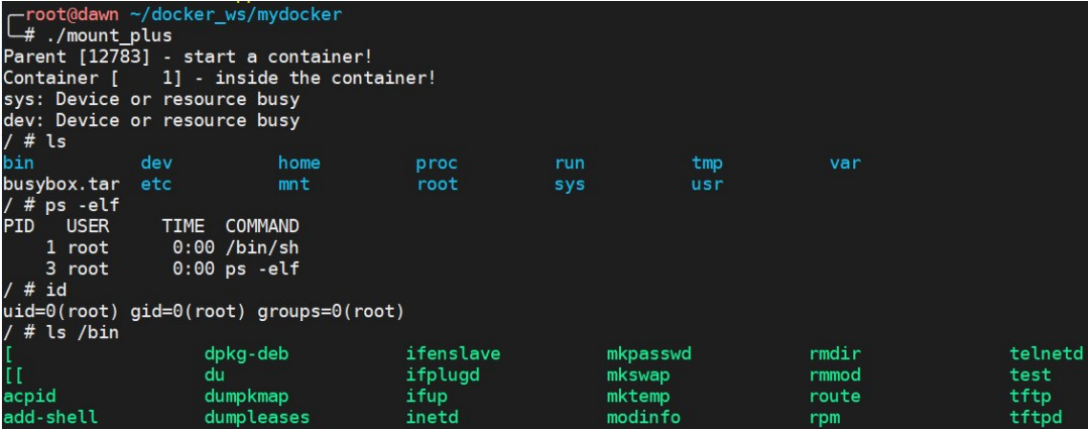
```

1 char* const container_args[] = {
2     "/bin/sh",
3     NULL
4 };
5
6 int container_main(void* arg) {
7     printf("Container [%5d] - inside the container!\n", getpid());
8     sethostname("container_dawn", 15);
9
10    if (mount("proc", "rootfs/proc", "proc", 0, NULL) != 0) {
11        perror("proc");
12    }
13    if (mount("sysfs", "rootfs/sys", "sysfs", 0, NULL) != 0) {
14        perror("sys");
15    }
16    if (mount("none", "rootfs/tmp", "tmpfs", 0, NULL) != 0) {
17        perror("tmp");
18    }
19    if (mount("udev", "rootfs/dev", "devtmpfs", 0, NULL) != 0) {
20        perror("dev");
21    }
22    if (mount("devpts", "rootfs/dev/pts", "devpts", 0, NULL) != 0) {
23        perror("dev/pts");
24    }
25    if (mount("shm", "rootfs/dev/shm", "tmpfs", 0, NULL) != 0) {
26        perror("dev/shm");
27    }
28    if (mount("tmpfs", "rootfs/run", "tmpfs", 0, NULL) != 0) {
29        perror("run");
30    }
31
32    if ( chdir("./rootfs") || chroot("./") != 0 ){
33        perror("chdir/chroot");
34    }
35
36    // 改变根目录之后，那么 /bin/bash 是从改变之后的根目录中搜索了
37    execv(container_args[0], container_args);
38    perror("exec");
39    printf("Something's wrong!\n");
40    return 1;
41 }
42
43 int main() {

```

```
44     printf("Parent [%5d] - start a container!\n", getpid());
45     int container_id = clone(container_main, container_stack + STACK_SIZE,
46                             CLONE_NEWUTS | CLONE_NEWPID | CLONE_NEWNS | SIGCHLD, NULL);
47     waitpid(container_id, NULL, 0);
48     printf("Parent - container stopped!\n");
49     return 0;
50 }
```

最后，查看实现效果如下图所示。



实际上，Mount Namespace 是基于 chroot 的不断改良才被发明出来的，chroot 可以算是 Linux 中第一个 Namespace。那么上面被挂载在容器根目录上、用来为容器镜像提供隔离后执行环境的文件系统，就是所谓的容器镜像，也被叫做 rootfs（根文件系统）。需要明确的是，rootfs 只是一个操作系统所包含的文件、配置和目录，并不包括操作系统内核。

User Namespace

容器内部看到的 UID 和 GID 和外部是不同的了，比如容器内部针对 dawn 这个用户显示的是 0，但是实际上这个用户在主机上应该是 1000。要实现这样的效果，需要把容器内部的 UID 和主机的 UID 进行映射，需要修改的文件是 /proc/<pid>/uid_map 和 /proc/<pid>/gid_map，这两个文件的格式是

ID-INSIDE-NS	ID-OUTSIDE-NS	LENGTH
--------------	---------------	--------

ID-INSIDE-NS：表示在容器内部显示的 UID 或 GID

ID-OUTSIDE-NS：表示容器外映射的真实的 UID 和 GID

LENGTH：表示映射的范围，一般为 1，表示一一对应

比如，下面就是将真实的 uid=1000 的映射为容器内的 uid =0：

```
1 | $ cat /proc/8353/uid_map
2 | 0      1000      1
```

再比如，下面则表示把 namesapce 内部从 0 开始的 uid 映射到外部从 0 开始的 uid，其最大范围是无符号 32 位整型（下面这条命令是在主机环境中输入的）。

```
1 | $ cat /proc/$$/uid_map
2 | 0      0 4294967295
```

默认情况，设置了 CLONE_NEWUSER 参数但是没有修改上述两个文件的话，容器中默认情况下显示为 65534，这是因为容器找不到真正的 UID，所以就设置了最大的 UID。如下面的代码所示：

```
1 #define _GNU_SOURCE
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <sys/types.h>
5 #include <sys/wait.h>
6 #include <sys/mount.h>
7 #include <sys/capability.h>
8 #include <stdio.h>
9 #include <sched.h>
10 #include <signal.h>
11 #include <unistd.h>
12
13 #define STACK_SIZE (1024 * 1024)
14
15 static char container_stack[STACK_SIZE];
16 char* const container_args[] = {
17     "/bin/bash",
18     NULL
19 };
20
21 int container_main(void* arg) {
22
23     printf("Container [%5d] - inside the container!\n", getpid());
24
25     printf("Container: eUID = %ld; eGID = %ld, UID=%ld, GID=%ld\n",
26         (long) geteuid(), (long) getegid(), (long) getuid(), (long) getgid());
27
28     printf("Container [%5d] - setup hostname!\n", getpid());
29
30     //set hostname
31     sethostname("container",10);
32
33     execv(container_args[0], container_args);
34     printf("Something's wrong!\n");
35     return 1;
36 }
37
38 int main() {
39     const int gid=getgid(), uid=getuid();
40
41     printf("Parent: eUID = %ld; eGID = %ld, UID=%ld, GID=%ld\n",
42         (long) geteuid(), (long) getegid(), (long) getuid(), (long) getgid());
43
44     printf("Parent [%5d] - start a container!\n", getpid());
45
46     int container_pid = clone(container_main, container_stack+STACK_SIZE,
47         CLONE_NEWUTS | CLONE_NEWPID | CLONE_NEWUSER | SIGCHLD, NULL);
48
49
50     printf("Parent [%5d] - Container [%5d]!\n", getpid(), container_pid);
51
52     printf("Parent [%5d] - user/group mapping done!\n", getpid());
53
54     waitpid(container_pid, NULL, 0);
55     printf("Parent - container stopped!\n");
56     return 0;
```

当我以 dawn 这个用户执行的该程序的时候，那么会显示如下图所示的效果。使用 root 用户的时候是同样的：

```
dawn@dawn:~$ ./user2
Parent: eUID = 1000; eGID = 1000, UID=1000, GID=1000
Parent [ 8062] - start a container!
Parent [ 8062] - Container [ 8063]!
Container [ 1] - inside the container!
Container: eUID = 65534; eGID = 65534, UID=65534, GID=65534
nobody@container:~$

root@dawn ~/docker_ws/mydocker
# ./user
Parent: eUID = 0; eGID = 0, UID=0, GID=0
Parent [ 8229] - start a container!
Parent [ 8229] - Container [ 8230]!
Container [ 1] - inside the container!
Container: eUID = 65534; eGID = 65534, UID=65534, GID=65534
nobody@container:~/docker_ws/mydocker$
```

接下去，我们要开始来实现映射的效果了，也就是让 dawn 这个用户在容器中显示为 0。代码是几乎完全拿耗子叔的博客上的，链接可见文末：

```
1 int pipefd[2];
2
3 void set_map(char* file, int inside_id, int outside_id, int len) {
4     FILE* mapfd = fopen(file, "w");
5     if (NULL == mapfd) {
6         perror("open file error");
7         return;
8     }
9     fprintf(mapfd, "%d %d %d", inside_id, outside_id, len);
10    fclose(mapfd);
11 }
12
13 void set_uid_map(pid_t pid, int inside_id, int outside_id, int len) {
14     char file[256];
15     sprintf(file, "/proc/%d/uid_map", pid);
16     set_map(file, inside_id, outside_id, len);
17 }
18
19 int container_main(void* arg) {
20
21     printf("Container [%5d] - inside the container!\n", getpid());
22
23     printf("Container: eUID = %ld; eGID = %ld, UID=%ld, GID=%ld\n",
24           (long) geteuid(), (long) getegid(), (long) getuid(), (long) getgid());
25
26     /* 等待父进程通知后再往下执行（进程间的同步） */
27     char ch;
28     close(pipefd[1]);
29     read(pipefd[0], &ch, 1);
30
31     printf("Container [%5d] - setup hostname!\n", getpid());
32     //set hostname
33     sethostname("container", 10);
34
35     //remount "/proc" to make sure the "top" and "ps" show container's information
36     mount("proc", "/proc", "proc", 0, NULL);
37 }
```



```

38     execv(container_args[0], container_args);
39     printf("Something's wrong!\n");
40     return 1;
41 }
42
43 int main() {
44     const int gid=getgid(), uid=getuid();
45
46     printf("Parent: eUID = %ld; eGID = %ld, UID=%ld, GID=%ld\n",
47           (long) geteuid(), (long) getegid(), (long) getuid(), (long) getgid());
48
49     pipe(pipefd);
50
51     printf("Parent [%5d] - start a container!\n", getpid());
52
53     int container_pid = clone(container_main, container_stack+STACK_SIZE,
54                              CLONE_NEWUTS | CLONE_NEWPID | CLONE_NEWNS | CLONE_NEWUSER | SIGCHLD, NULL);
55
56
57     printf("Parent [%5d] - Container [%5d]!\n", getpid(), container_pid);
58
59     //To map the uid/gid,
60     // we need edit the /proc/PID/uid_map (or /proc/PID/gid_map) in parent
61     set_uid_map(container_pid, 0, uid, 1);
62
63     printf("Parent [%5d] - user/group mapping done!\n", getpid());
64
65     /* 通知子进程 */
66     close(pipefd[1]);
67
68     waitpid(container_pid, NULL, 0);
69     printf("Parent - container stopped!\n");
70     return 0;
71 }

```

实现的最终效果如图所示，可以看到在容器内部将 dawn 这个用户 UID 显示为了 0 (root)，但其实这个容器中的 /bin/bash 进程还是以以一个普通用户，也就是 dawn 来运行的，只是显示出来的 UID 是 0，所以当查看 /root 目录的时候还是没有权限。

```

dawn@dawn:~$ ./user
Parent: eUID = 1000; eGID = 1000, UID=1000, GID=1000
Parent [ 8325] - start a container!
Parent [ 8325] - Container [ 8326]!
Parent [ 8325] - user/group mapping done!
Container [ 1] - inside the container!
Container: eUID = 0; eGID = 65534, UID=0, GID=65534
Container [ 1] - setup hostname!
root@container:~# id
uid=0(root) gid=65534(nogroup) groups=65534(nogroup)
root@container:~# ls /root/
ls: cannot open directory '/root/': Permission denied

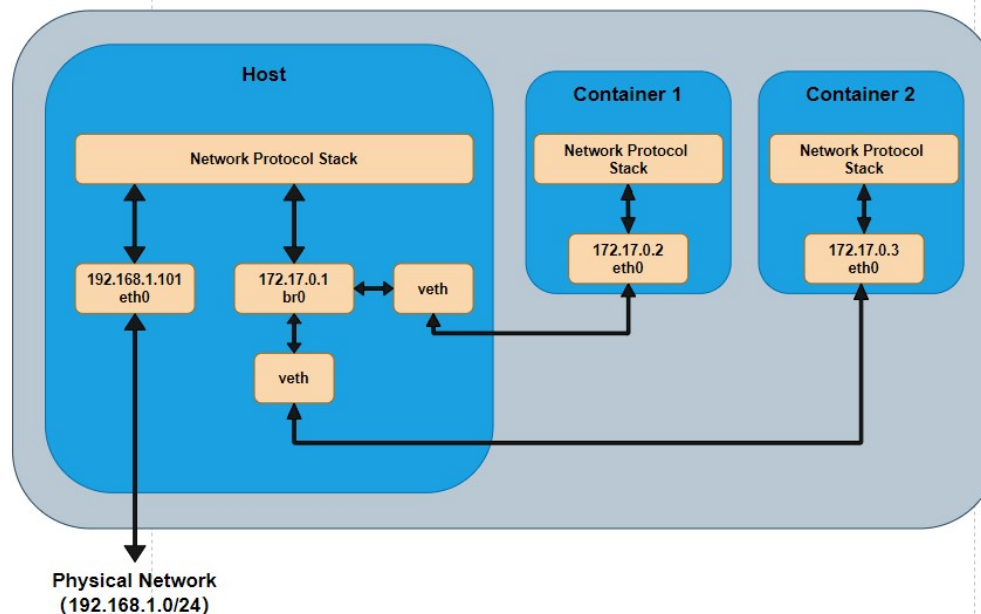
```

User Namespace 是以普通用户运行的，但是别的 Namespace 需要 root 权限，那么当使用多个 Namespace 该怎么办呢？我们可以先用一般用户创建 User Namespace，然后把这个一般用户映射成 root，那么在容器内用 root 来创建其他的 Namespace。

Network Namespace

隔离容器中的网络，每个容器都有自己的虚拟网络接口和 IP 地址。在 Linux 中，可以使用 ip 命令创建 Network Namespace (Docker 的源码中，它没有使用 ip 命令，而是自己实现了 ip 命令内的一些功能)。

下面就使用 ip 命令来讲解一下 Network Namespace 的构建，以 bridge 网络为例。bridge 网络的拓扑图一般如下图所示，其中 br0 是 Linux 网桥。



在使用 Docker 的时候，如果启动一个 Docker 容器，并使用 `ip link show` 查看当前宿主机上的网络情况，那么你会看到有一个 `docker0` 还有一个 `veth****` 的虚拟网卡，这个 `veth` 的虚拟网卡就是上图中 `veth`，而 `docker0` 就相当于上图中的 `br0`。

那么，我们可以使用下面这些命令即可创建跟 `docker` 类似的效果（参考自耗子叔的博客，链接见文末参考，结合上图加了一些文字）。

```
1 ## 1. 首先，我们先增加一个网桥 lxcbr0，模仿 docker0
2 brctl addbr lxcbr0
3 brctl stp lxcbr0 off
4 ifconfig lxcbr0 192.168.10.1/24 up #为网桥设置IP地址
5
6 ## 2. 接下来，我们要创建一个 network namespace，命名为 ns1
7
8 # 增加一个 namespace 命令为 ns1（使用 ip netns add 命令）
9 ip netns add ns1
10
11 # 激活 namespace 中的 loopback，即127.0.0.1（使用 ip netns exec ns1 相当于进入了 ns1 这个 namespace，那么 ip link set lo up
12 ip netns exec ns1 ip link set dev lo up
13
14 ## 3. 然后，我们需要增加一对虚拟网卡
15
16 # 增加一对虚拟网卡，注意其中的 veth 类型。这里有两个虚拟网卡：veth-ns1 和 lxcbr0.1，veth-ns1 网卡是要被安到容器中的，而 lxcbr0.1
17 ip link add veth-ns1 type veth peer name lxcbr0.1
18
19 # 把 veth-ns1 按到 namespace ns1 中，这样容器中就会有新的网卡了
20 ip link set veth-ns1 netns ns1
21
22 # 把容器里的 veth-ns1 改名为 eth0（容器外会冲突，容器内就不会了）
23 ip netns exec ns1 ip link set dev veth-ns1 name eth0
24
25 # 为容器中的网卡分配一个 IP 地址，并激活它
26 ip netns exec ns1 ifconfig eth0 192.168.10.11/24 up
27
28
29 # 上面我们把 veth-ns1 这个网卡按到了容器中，然后我们要把 lxcbr0.1 添加上网桥上
30 brctl addif lxcbr0 lxcbr0.1
```

```

31
32 # 为容器增加一个路由规则，让容器可以访问外面的网络
33 ip netns exec ns1 ip route add default via 192.168.10.1
34
35 ## 4. 为这个 namespace 设置 resolv.conf，这样，容器内就可以访问域名了
36 echo "nameserver 8.8.8.8" > conf/resolv.conf

```

上面基本上就相当于 docker 网络的原理，只不过：

Docker 不使用 ip 命令而是，自己实现了 ip 命令内的一些功能。

Docker 的 resolv.conf 没有使用这样的方式，而是将其写到指定的 resolv.conf 中，之后在启动容器的时候将其和 hostname、host 一起以只读的方式加载到容器的文件系统中。

docker 使用进程的 PID 来做 network namespace 的名称。

同理，我们还可以使用如下的方式为正在运行的 docker 容器增加一个新的网卡

```

1 ip link add peerA type veth peer name peerB
2 brctl addif docker0 peerA
3 ip link set peerA up
4 ip link set peerB netns ${container-pid}
5 ip netns exec ${container-pid} ip link set dev peerB name eth1
6 ip netns exec ${container-pid} ip link set eth1 up
7 ip netns exec ${container-pid} ip addr add ${ROUTEABLE_IP} dev eth1

```

Namespace 情况查看

Cgroup 的操作接口是文件系统，位于 `/sys/fs/cgroup` 中。假如想查看 namespace 的情况同样可以查看文件系统，namespace 主要查看 `/proc/<pid>/ns` 目录。

我们以上面的 [PID Namespace 程序](#PID Namespace) 为例，当这个程序运行起来之后，我们可以看到其 PID 为 11702。

```

root@dawn ~/docker_ws/mydocker
# ./mycontainer_pid
Parent [11702] - start a container!
Container [ 1 ] - inside the container!
root@container_dawn:~/docker_ws/mydocker#

```

之后，我们保持这个子进程运行，然后打开另一个 shell，查看这个程序创建的子进程的 PID，也就是容器中运行的进程在主机中的 PID。

最后，我们分别查看 `/proc/11702/ns` 和 `/proc/11703/ns` 这两个目录的情况，也就是查看这两个进程的 namespace 情况。可以看到其中 cgroup、ipc、mnt、net、user 都是同一个 ID，而 pid、uts 是不同的 ID。如果两个进程的 namespace 编号相同，那么表示这两个进程位于同一个 namespace 中，否则位于不同 namespace 中。

```

root@dawn ~/docker_ws/mydocker
# pstree -p 11702
mycontainer_pid(11702)─bash(11703)
root@dawn ~/docker_ws/mydocker
# ll /proc/11702/ns/
total 0
lrwxrwxrwx 1 root root 0 Oct 7 20:52 cgroup -> cgroup:[4026531835]
lrwxrwxrwx 1 root root 0 Oct 7 20:51 ipc -> ipc:[4026531839]
lrwxrwxrwx 1 root root 0 Oct 7 20:51 mnt -> mnt:[4026531840]
lrwxrwxrwx 1 root root 0 Oct 7 20:51 net -> net:[4026532009]
lrwxrwxrwx 1 root root 0 Oct 7 20:51 pid -> pid:[4026531836]
lrwxrwxrwx 1 root root 0 Oct 7 20:52 pid_for_children -> pid:[4026531836]
lrwxrwxrwx 1 root root 0 Oct 7 20:51 user -> user:[4026531837]
lrwxrwxrwx 1 root root 0 Oct 7 20:51 uts -> uts:[4026531838]
root@dawn ~/docker_ws/mydocker

```

```

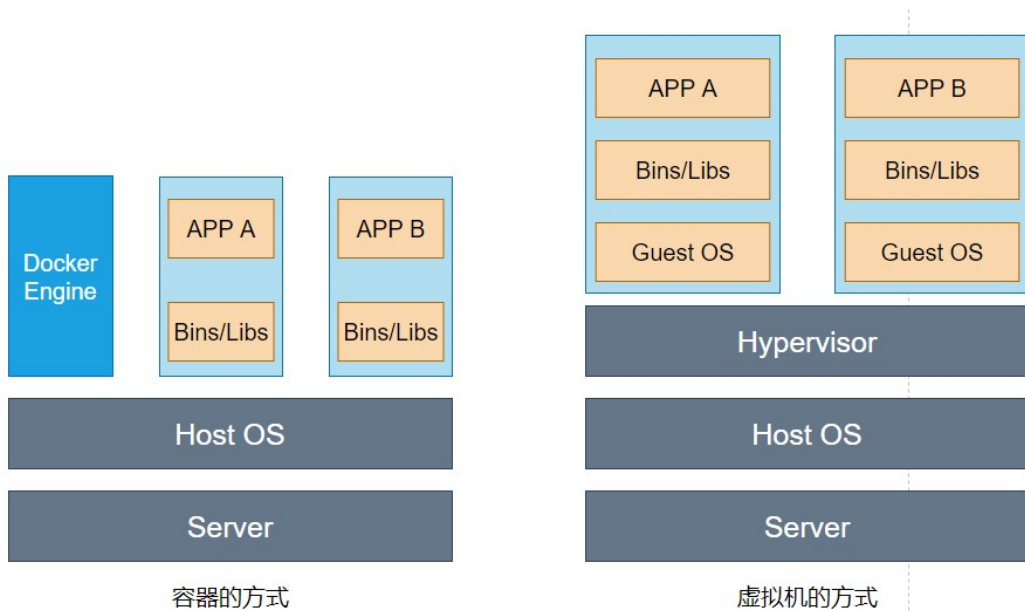
root@qawn:~/docker_ws/mydocker
└─# ll /proc/11703/ns/
total 0
lrwxrwxrwx 1 root root 0 0ct  7 20:52 cgroup -> cgroup:[4026531835]
lrwxrwxrwx 1 root root 0 0ct  7 20:51 ipc -> ipc:[4026531839]
lrwxrwxrwx 1 root root 0 0ct  7 20:51 mnt -> mnt:[4026531840]
lrwxrwxrwx 1 root root 0 0ct  7 20:51 net -> net:[4026532009]
lrwxrwxrwx 1 root root 0 0ct  7 20:51 pid -> pid:[4026532548]
lrwxrwxrwx 1 root root 0 0ct  7 20:52 pid_for_children -> pid:[4026532548]
lrwxrwxrwx 1 root root 0 0ct  7 20:51 user -> user:[4026531837]
lrwxrwxrwx 1 root root 0 0ct  7 20:51 uts -> uts:[4026532547]

```

如果可以查看 ns 的情况之外，这些文件一旦被打开，只要 fd 被占用着，即使 namespace 中所有进程都已经结束了，那么创建的 namespace 也会一直存在。比如可以使用 `mount --bind /proc/11703/ns/uts ~/uts`，让 11703 这个进程的 UTS Namespace 一直存在。

总结

Namespace 技术实际上修改了应用进程看待整个计算机“视图”，即它的“视图”已经被操作系统做了限制，只能“看到”某些指定的内容，这仅仅对应用进程产生了影响。**但是对宿主机来说，这些被隔离的进程，其实还是进程，跟宿主机上其他进程并无太大区别，都由宿主机统一管理。只不过这些被隔离的进程拥有额外设置过的 Namespace 参数。**那么 Docker 项目在这里扮演的，更多是旁路式的辅助和管理工作。如下左图所示



因此，相比虚拟机的方式，容器会更受欢迎。这是假如使用虚拟机的方式作为应用沙盒，那么必须要由 Hypervisor 来负责创建虚拟机，这个虚拟机是真实存在的，并且里面必须要运行一个完整的 Guest OS 才能执行用户的应用进程。这样就导致了采用虚拟机的方式之后，不可避免地带来额外的资源消耗和占用。根据实验，一个运行着 CentOS 的 KVM 虚拟机启动后，在不做优化的情况下，虚拟机就需要占用 100-200 MB 内存。此外，用户应用运行在虚拟机中，它对宿主机操作系统的调用就不可避免地要经过虚拟机软件的拦截和处理，这本身就是一层消耗，尤其对资源、网络和磁盘 IO 的损耗非常大。

而假如使用容器的方式，容器化之后应用本质还是宿主机上的一个进程，这也就意味着因为虚拟机化带来的性能损耗是不存在的；而另一方面使用 Namespace 作为隔离手段的容器并不需要单独的 Guest OS，这就使得容器额外的资源占用几乎可以忽略不计。

总得来说，“敏捷”和“高性能”是容器相对于虚拟机最大的优势，也就是容器能在 PaaS 这种更加细粒度的资源管理平台上大行其道的重要原因。

但是！基于 Linux Namespace 的隔离机制相比于虚拟化技术也有很多不足之处，其中最主要的问题就是隔离不彻底。

首先，容器只是运行在宿主机上的一种特殊进程，那么容器之间使用的还是同一个宿主机上的操作系统。尽管可以在容器里面通过 mount namesapce 单独挂载其他不同版本的操作系统文件，比如 centos、ubuntu，但是这并不能改变共享宿主机内核的事实。这就意味着你要在 windows 上运行 Linux 容器，或者在低版本的 Linux 宿主机上运行高版本的 Linux 容器都是行不通的。

而拥有虚拟机技术和独立 Guest OS 的虚拟机就要方便多了。

其次，在 Linux 内核中，有很多资源和对象都是不能被 namespace 化的，比如时间。假如你的容器中的程序使用 `settimeofday(2)` 系统调用修改了时间，整个宿主机的时间都会被随之修改。

相比虚拟机里面可以随意折腾的自由度，在容器里部署应用的时候，“什么能做，什么不能做” 是用户必须考虑的一个问题。之外，容器给应用暴露出来的攻击面是相当大的，应用“越狱”的难度也比虚拟机低很多。虽然，实践中可以使用 Seccomp 等技术对容器内部发起的所有系统调用进行过滤和甄别来进行安全加固，但这种方式因为多了一层对系统调用的过滤，也会对容器的性能产生影响。因此，在生产环境中没有人敢把运行在物理机上的 Linux 容器直接暴露到公网上。

另外，容器是一个“单进程”模型。容器的本质是一个进程，用户的应用进程实际上就是容器里 PID=1 的进程，而这个进程也是后续创建的所有进程的父进程。这也就意味着，在一个容器中，你没办法同时运行两个不同的应用，除非能事先找到一个公共的 PID=1 的程序来充当两者的父进程，比如使用 systemd 或者 supervisor。容器的设计更多是希望容器和应用同生命周期的，而不是容器还在运行，而里面的应用早已经挂了。

★

上面这段话个人的理解是：因为创建出子进程之后，子进程需要运行的，而此时父进程需要等待子进程运行结束，相当于只有子进程在运行。比如容器中的第一个进程往往就是业务需要的进程，也就是 entrypoint 指定的程序运行起来的进程。而创建出子进程会导致这个进程被暂停。即使将子进程改为后台执行，但是由于容器中 PID=1 的进程压根没有管理后台进程的能力，所以还是会有进程无法管理。

”

巨人的肩膀

极客时间——《深入剖析 Kubernetes》——张磊老师

DOCKER基础技术：LINUX NAMESPACE（上）

DOCKER基础技术：LINUX NAMESPACE（下）。



多选参数

关注

👍 1

💬 6

🌟 6

📄

|

🔗

Device namespace简介 - 基于Kernel namespace的设备虚拟化

世事难料，保持低调 4506

在移动设备上，虚拟化的需求正在逐渐增加。其一，移动设备配置配置越来越高，一些高端配置已和桌面设备接近；其二，用户对于移动设备使用场景的多样...

实现Docker容器的底层技术-cgroup 和namespace浅析

NathanniuBee的博客 1003

引言 docker启动的容器本质上是Host中的一个进程，cgroup和namespace是最重要的两项技术，cgroup主要实现资源的限额，而namespce则用来实现资...



请发表有价值的评论， 博客评论欢迎灌水，良好的社区氛围需大家一起维护。



评论

乎你： 服装赞助商---林品如，大佬也太秀了吧~~~ 1 年前 回复 ...



LaoYuanPython： 谢谢分享!点赞支持！希望点赞给博主带来好运气！祝博主新的一年平安健康、幸福快乐！ 本人正参与博客之星评选，1月24



日前每天都可投多票，敬请您支持！谢谢！ 投票链接：

1 https://bss.csdn.net/m/topic/blog_star2020/detail?username=laoyuanpython

或到老猿博文首页内的置顶博文跳转！ 1 年前 回复 ...

Java劝退师、： 感谢博主分享，讲的挺不错，希望后面有更多的文章~ 最近我也在学习写博客,有空来看看我呀，一起互相学习。期待你的关注与 支持 1 年前 回复 ...



多选参数

码龄6年



暂无认证

[illegible]

2021年 5篇

2020年 82篇

公安备案号11010502030143 京ICP备19004658号 京网文〔2020〕1039-165号 经营性网站备案信息 北京互联网违法和不良信息举报中心 网络110报警服务 中国互联网举报中心
家长监护 Chrome商店下载 ©1999-2022北京创新乐知网络技术有限公司 版权与免责声明 版权申诉 出版物许可证 营业执照