

Docker (Swarm) and NFS volumes

Posted by: Anthony Dahanne in DevOps November 22nd, 2017

Investigating how to use shared volumes with Docker (Swarm), I decided to take a look at NFS volume on premises way to share folders.

With Docker, you have 3 different syntaxes to mount NFS volumes :

- simple container (via `docker volume create + docker run`)
- single service (via `docker service create`)
- complete stack (`docker deploy -f stack.yml`)

I actually had some trouble mounting NFS volumes, especially with images that COPY files into declar

So, without further ado...

Mounting a NFS share to a Docker container

Let's start creating a folder under our nfs share (/nfs)

```
1 #(host) mkdir /nfs/test
2 #(host) ls -al
3 drwxr-xr-x 2 nfsnobody nfsnobody 4096 Nov 18 18:37 test
4 #(host) touch /nfs/test/hello-test
```

Now we have a file named « hello-test » under /nfs/test/

Let's create a new docker volume linked to this nfs share, and then let's spin up a new container that

```
01 #(host) docker volume create --driver local --opt type=nfs --opt
   o=addr=nfs.my.corporate.network,rw --opt device=:/nfs/test test
02 test
03 #(host) docker run --rm -it -v test:/data alpine
04 #(container) touch /data/hello-from-container
05 #(container) ls -al /data
06 -rw-r--r-- 1 nobody nobody 0 Nov 19 02:44 hello-from-container
07 -rw-r--r-- 1 nobody nobody 0 Nov 19 02:43 hello-test
08 #(host) ls -al /nfs/test/
09 -rw-r--r-- 1 nfsnobody nfsnobody 0 Nov 18 18:44 hello-from-containe
10 -rw-r--r-- 1 nfsnobody nfsnobody 0 Nov 18 18:43 hello-test
```

So far, so good ! (did you notice the host root user mapping to nfsnobody and the container root user this a bit later)

Things unfortunately get weirder when you mount to a defined VOLUME location.

For example, with this Dockerfile :

```
1 FROM alpine
2 VOLUME /data
```

```
1 #(host) docker build -t volume:defined .
2 #(host) docker run --rm -it -v test:/data alpine
```

```
1 | # (host) docker build -t volume:defined .
2 | # (host) docker run --rm -it -v test:/data alpine
```

Everything works great.

Now if you try to do the same thing with an image that copied something to this volume :

```
1 | FROM alpine
2 | VOLUME /data/
3 | COPY empty.file /data/empty.file
```

```
1 | # (host) docker build -t volume:defined-with-copy .
```

and try to map /data to the nfs volume, this time you'll get :

```
1 | 2
2 | # (host) docker run --rm -it -v test:/data volume:defined-with-copy
3 | docker: Error response from daemon: chown /var/lib/docker/volumes/t
   | permitted.
```

To get out of this issue, you need the nocopy option :

```
1 | # (host) docker run --rm -it -v test:/data:nocopy volume:defined-wit
2 | # (container) ls -al /data
3 | -rw-r--r-- 1 nobody nobody 0 Nov 19 02:44 hello-from-container
4 | -rw-r--r-- 1 nobody nobody 0 Nov 19 02:43 hello-test
```

Notice how the original empty.file from the Docker file was « nocopy »ed

So we got away with our issue, well, except something that was provided by the image (empty.file), is

Let's clean up first the previous experimentation

```
1 | docker volume rm test
2 | test
```

and move on to services.

Mounting a NFS share to a Docker service

This use case is more interesting than the previous one, since you only need to mount once to the service, and any container created by this service, on any host part of the swarm. (here the service will just ls /data in its entrypoint overriding from a service creation – the parsed command is ... not obvious...)

```
1 | # docker service create --mount 'type=volume,src=test,volume-driver=local' test
```

```
1 | # docker service logs -f test
2 | test.1.4wteqs8zsd99@worker01 | hello-from-container
3 | test.1.4wteqs8zsd99@worker01 | hello-test
```

So far so good, but now with the infamous VOLUME with copy :

```
1 | docker service create --mount 'type=volume,src=test,volume-driver=local' test
```

The service won't create the container ! To see that, use service ps :

```
1 | docker service create --mount 'type=volume,src=test,volume-driver=local'
```

The service won't create the container ! To see that, use service ps :

```
1 | # docker service ps test --no-trunc
2 | ID NAME IMAGE NODE DESIRED STATE CURRENT STATE ERROR PORTS
3 | i3 test.1 volume:defined-with-copy worker01 Ready Rejected 2 seconds
   | /volumes/test/_data: operation not permitted"
```

So adding a volume-nocopy=true should solve the issue :

```
1 | # docker service create --mount 'type=volume,src=test,volume-driver=local,volume-nocopy=true'
```

Since, per the documentation :

« By default, if you attach an empty volume to a container, and files or directories already existed at the destination (dst

), the Engine copies those files and directories into the volume, allowing the host to access them. Set the volume-nocopy

to disables copying files from the container's filesystem to the volume and mount the empty volume. And

and indeed, the service is now started :

```
1 | # docker service logs -f test
2 | test.1.on5ahp585oqy@worker01 | hello-from-container
3 | test.1.on5ahp585oqy@worker01 | hello-test
```

now, to have this service interact with others, let's create a stack

Mounting a NFS share in a Docker Stack

To create a stack, we'll create a docker-compose.yml file this time :

```
01 | version: '3.4'
02 |
03 | networks:
04 |   terracotta-net:
05 |     driver: overlay
06 |
07 | volumes:
08 |   test:
09 |     driver: local
10 |     driver_opts:
11 |       type: nfs
12 |       o: addr=nfs.my.corporate.network,rw
13 |       device: ":/nfs/test"
14 | services:
15 |   producer:
16 |     image : volume:defined-with-copy
17 |     command: |
18 |       /bin/sh -c "
19 |         while true;
20 |         do touch /data/hello-from-producer;
```

```

16     image : volume:defined-with-copy
17     command: |
18         /bin/sh -c "
19             while true;
20                 do touch /data/hello-from-producer;
21                 echo '/data/hello-from-producer written';
22                 sleep 5;
23                 rm /data/hello-from-producer;
24                 echo '/data/hello-from-producer deleted';
25                 sleep 5;
26             done
27         "
28     volumes:
29         - type: volume
30           source: test
31           target: /data
32           volume:
33             nocopy: true
34     consumer:
35     image : volume:defined-with-copy
36     command : |
37         /bin/sh -c "
38             while true; do ls -al /data/; sleep 1; done
39         "
40     volumes:
41         - type: volume
42           source: test
43           target: /data
44           volume:
45             nocopy: true

```

So a producer will write a file to the nfs share, and a consumer will ls -al the folder where the file is su and see how well that goes :

```

01 # docker stack deploy test --compose-file docker-compose.yml
02 Creating network test_default
03 Creating service test_producer
04 Creating service test_consumer
05 # docker service logs -f test_producer
06 test_producer.1@worker01 | /data/hello-from-producer written
07 test_producer.1@worker01 | /data/hello-from-producer deleted
08 [...]
09 # docker service logs -f test_consumer
10 test_consumer.1@worker02 | drwxr-xr-x 2 nobody nobody 4096 Nov 20 (
11 test_consumer.1@worker02 | drwxr-xr-x 19 root root 218 Nov 20 04:39
12 test_consumer.1@worker02 | total 4
13 test_consumer.1@worker02 | drwxr-xr-x 2 nobody nobody 4096 Nov 20 (
14 test_consumer.1@worker02 | drwxr-xr-x 19 root root 218 Nov 20 04:40
15 test_consumer.1@worker02 | -rw-r--r-- 1 nobody nobody 0 Nov 20 04:40
16 test_consumer.1@worker02 | total 4
17 [...]

```

So that worked pretty well : you'll notice the nocopy syntax that is this time :

```

1 volumes:
2   - type: volume
3     source: test
4     target: /data
5     volume:
6       nocopy: true

```

If you make the **mistake** of putting it like this :

```

1 volumes:

```

If you make the **mistake** of putting it like this :

```
1 volumes:
2   - type: volume
3     source: test
4     target: /data:nocopy
```

you'll probably end up with :

```
1 worker01 Shutdown Failed 19 seconds ago "starting container failed:
```

since there is no folder named /nfs/test:nocopy !

Was that so terrible ?

Well, the different syntaxes are a bit confusing, and the error messages sometimes are...

Oh, I almost forgot, about the users ownerships : if you run an image that creates a new user and use write the file to the NFS share : well, according to your NFS share setup, it could be that the user create a whole different user than nfsnobody !

In that case make sure all your images create this user the same way on Docker nodes that are configured nicely with the NFS share !