

VM Storage

Drive I/O is one of the most time-consuming tasks that a service or program experiences, probably second only to network I/O. Proxmox presents an interesting challenge when it comes to managing bulk storage. Over the years, I have come across several ways of dealing with this issue.

Pass the whole honkin' drive through

Proxmox [allows you to passthrough a block device](#) into the VM. Issue is, this is a one to one relation; that is to say, only one of your VMs can access the disk(s) you choose. Of course, you could then share the disks via a network protocol and mount them on other VM guests, but there's a better solution below.

Also, from a data integrity standpoint, I have some (maybe unwarranted) reservations as to adding another layer of translation between your I/O requests and the drive controller. If the data were to be corrupted somehow during translation, you're in for a bad time even with software raid (your raid software is none-the-wiser to read/write errors []).

Pass through a PCIE HBA card with drives attached

This is similar to the previous method, however you are [passing a drive adapter card into the VM](#). This would avoid the I/O translation issue, as the VM gets direct memory access to the PCIE drive adapter, and therefore is able to speak to it directly. However, there is the same downside which is only a single VM has access to your drive(s).

Use a network based protocol

If we have the drives mounted on the host, and the host shares a network with a VM, how about we share "over the network" and mount the drive on the guest? This is the solution that I went with on many installations.

1. Create a private bridge network that is not attached to any NICs and assign the host an IP to use on the bridge
2. Create a network interface for each VM and connect it to the bridge
3. Assign an IP to the interface and mount the network share via the hosts' IP

If you use ZFS, an easy way of sharing a dataset is by setting the `sharenfs` property, which will automatically manage a zfs share for you. For example: `zfs set sharenfs=rw=192.168.254.0/24,no_root_squash pool/volume`

This method is pretty neat, as it allows multiple hosts to access one or more mount points from the host. When it comes to NFS, there are some configurations that need to be tuned, such as when it comes to user ID mapping (see [man 8 exports](#)) as well as the size of read and write requests.

Make sure that programs and services on the guest (other than the folder mounting software) do not have access to the Proxmox host's IP address over the internal network! This is especially important if you are using an unencrypted protocol such as NFS.

Unfortunately, when it comes to my own home server I have found this sort of setup is quite prohibitively demanding, requiring upwards of two threads worth of CPU shares when compared to the alternative (below). I find this is especially the case if a program is attempting to access multiple files simultaneously. However, with a powerful enough system I'd imagine dedicating a few cores to this won't be too bad.

Bind mount a folder on the host into a Linux Container (LXC)

Well gee, we're only running Linux on the guests anyway, how about we use LXCs and [just mount the drives directly](#)? Similar to how Podman/Docker containers can have volumes passed into them, we can pass through a folder into an LXC and have it show up... as a folder! This is the solution that I have settled on in the ol' homelab. Using `pct set <lxc_id> -mp0 /<host_mount_point>,mp=<guest_mount_point>` will configure your folder to show up in your LXC, no mess or fuss.

If the LXC is unprivileged, the folder and its contents may be inaccessible by the guest. This is due to the fact that the UID and GIDs of the contents of the drive are mapped to a number different than that of the host's UID and GID range (rootless Podman fanboys will be familiar with this). Take a peek at `/etc/subuid` and notice that the sub UIDs that the root user is able to [unshare](#) is shown as starting from the first number, and the number of UIDs that are available is shown as the first number. Same goes for sub GIDs, which are located in `/etc/subgid` **Ok, ok, TLDR; just recursively change the permission of your mount point to 10000:10000**

LXCs, among other things, are different from VMs in that instead of using CPU virtualization extensions they re-use the host's kernel and therefore can only run on Linux. Whether or not to use an LXC or a VM is kind of beyond what I'd like to write in this page though.

Revision #7

Created 2 July 2024 15:12:54 by GT

Updated 3 July 2024 00:27:46 by GT