

# PCI passthrough via OVMF

---

The Open Virtual Machine Firmware ([OVMF \(https://github.com/tianocore/tianocore.github.io/wiki/OVMF\)](https://github.com/tianocore/tianocore.github.io/wiki/OVMF)) is a project to enable UEFI support for virtual machines. Starting with Linux 3.9 and recent versions of [QEMU](#), it is now possible to passthrough a graphics card, offering the VM native graphics performance which is useful for graphic-intensive tasks.

**Related articles**

[Intel GVT-g](#)

Provided you have a desktop computer with a spare GPU you can dedicate to the host (be it an integrated GPU or an old OEM card, the brands do not even need to match) and that your hardware supports it (see [#Prerequisites](#)), it is possible to have a VM of any OS with its own dedicated GPU and near-native performance. For more information on techniques see the background [presentation \(pdf\) \(https://www.linux-kvm.org/images/b/b3/01x09b-VFIOandYou-small.pdf\)](https://www.linux-kvm.org/images/b/b3/01x09b-VFIOandYou-small.pdf).

## Contents

---

### Prerequisites

#### Setting up IOMMU

- Enabling IOMMU

- Ensuring that the groups are valid

- Gotchas

  - Plugging your guest GPU in an unisolated CPU-based PCIe slot

#### Isolating the GPU

- Binding vfio-pci via device ID

- Loading vfio-pci early

- Verifying that the configuration worked

#### Setting up an OVMF-based guest VM

- Configuring libvirt

- Setting up the guest OS

- Attaching the PCI devices

- Passing keyboard/mouse via Evdev

- Gotchas

  - Using a non-EFI image on an OVMF-based VM

#### Performance tuning

- CPU pinning

  - CPU topology

  - XML examples

    - 4c/1t CPU w/o Hyperthreading Example

[4c/2t Intel/AMD CPU example \(after ComboPI AGESA bios update\)](#)

[4c/2t AMD CPU example \(Before ComboPi AGESA bios update\)](#)

[Huge memory pages](#)

[Transparent huge pages](#)

[Static huge pages](#)

[Dynamic huge pages](#)

[CPU frequency governor](#)

[Isolating pinned CPUs](#)

[With isolcpus kernel parameter](#)

[Dynamically isolating CPUs](#)

[Improving performance on AMD CPUs](#)

[Virtio disk](#)

[Drivers](#)

[Considerations](#)

[IO threads](#)

[Examples with libvirt](#)

[Virtio network](#)

[Further tuning](#)

[Special procedures](#)

[Using identical guest and host GPUs](#)

[Script variants](#)

[Passthrough all GPUs but the boot GPU](#)

[Passthrough selected GPU](#)

[Script installation](#)

[Passing the boot GPU to the guest](#)

[Using Looking Glass to stream guest screen to the host](#)

[Adding IVSHMEM Device to VM](#)

[Installing the IVSHMEM Host to Windows guest](#)

[Getting a client](#)

[Swap peripherals to and from the Host](#)

[Bypassing the IOMMU groups \(ACS override patch\)](#)

[Plain QEMU without libvirt](#)

[Passing through other devices](#)

[USB controller](#)

[Passing VM audio to host via PulseAudio](#)

[QEMU 4.0/4.2+ audio changes](#)

[Passing VM audio to host via Scream](#)

[Using Scream with a bridged network](#)

[Adding the IVSHMEM device to use Scream with IVSHMEM](#)

Configuring the Windows guest for IVSHMEM

Configuring the host

Physical disk/partition

Gotchas

Passing through a device that does not support resetting

Complete setups and examples

Troubleshooting

QEMU 4.0: Unable to load graphics drivers/BSOD/Graphics stutter after driver install using Q35

QEMU 5.0: host-passthrough with newer kernel than 5.4 when using Zen2 processors: Windows 10 BSOD loop: KERNEL SECURITY CHECK FAILURE

"Error 43: Driver failed to load" on Nvidia GPUs passed to Windows VMs

"Error 43: Driver failed to load" with mobile (Optimus/max-q) nvidia GPUs

"BAR 3: cannot reserve [mem]" error in dmesg after starting VM

UEFI (OVMF) compatibility in VBIOS

Slowed down audio pumped through HDMI on the video card

No HDMI audio output on host when intel\_iommu is enabled

X does not start after enabling vfio\_pci

Chromium ignores integrated graphics for acceleration

VM only uses one core

Passthrough seems to work but no output is displayed

Host lockup after VM shutdown

Host lockup if guest is left running during sleep

Cannot boot after upgrading ovmf

QEMU via cli pulseaudio stuttering/delay

Bluescreen at boot since Windows 10 1803

AMD Ryzen / BIOS updates (AGESA) yields "Error: internal error: Unknown PCI header type '127'"

Navi 10 AMD GPU not resetting properly yielding "Error: internal error: Unknown PCI header type '127'" (Separate issue from the one above)

Host crashes when hotplugging Nvidia card with USB

Host unable to boot and stuck in black screen after enabling vfio

AER errors when passing through PCIe USB hub

See also

## Prerequisites

---

A VGA Passthrough relies on a number of technologies that are not ubiquitous as of today and might not be available on your hardware. You will not be able to do this on your machine unless the following requirements are met :

- Your CPU must support hardware virtualization (for kvm) and IOMMU (for the passthrough itself)
  - **List of compatible Intel CPUs (Intel VT-x and Intel VT-d)** ([https://ark.intel.com/SeArch/FeatureFilter?productType=873&0\\_VTD=True](https://ark.intel.com/SeArch/FeatureFilter?productType=873&0_VTD=True))
  - All AMD CPUs from the Bulldozer generation and up (including Zen) should be compatible.
    - CPUs from the K10 generation (2007) do not have an IOMMU, so you **need** to have a motherboard with a **890FX** (<https://support.amd.com/TechDocs/43403.pdf#page=18>) or **990FX** (<https://support.amd.com/TechDocs/48691.pdf#page=21>) chipset to make it work, as those have their own IOMMU.
- Your motherboard must also support IOMMU
  - Both the chipset and the BIOS must support it. It is not always easy to tell at a glance whether or not this is the case, but there is a fairly comprehensive list on the matter on the [Xen wiki](https://wiki.xen.org/wiki/VTd_HowTo) ([https://wiki.xen.org/wiki/VTd\\_HowTo](https://wiki.xen.org/wiki/VTd_HowTo)) as well as [Wikipedia:List of IOMMU-supporting hardware](#).
- Your guest GPU ROM must support UEFI.
  - If you can find **any ROM in this list** (<https://www.techpowerup.com/vgabios/>) that applies to your specific GPU and is said to support UEFI, you are generally in the clear. All GPUs from 2012 and later should support this, as Microsoft made UEFI a requirement for devices to be marketed as compatible with Windows 8.

You will probably want to have a spare monitor or one with multiple input ports connected to different GPUs (the passthrough GPU will not display anything if there is no screen plugged in and using a VNC or Spice connection will not help your performance), as well as a mouse and a keyboard you can pass to your VM. If anything goes wrong, you will at least have a way to control your host machine this way.

## Setting up IOMMU

### Note:

- IOMMU is a generic name for Intel VT-d and AMD-Vi.
- VT-d stands for *Intel Virtualization Technology for Directed I/O* and should not be confused with VT-x *Intel Virtualization Technology*. VT-x allows one hardware platform to function as multiple “virtual” platforms while VT-d improves security and reliability of the systems and also improves performance of I/O devices in virtualized environments.

Using IOMMU opens to features like PCI passthrough and memory protection from faulty or malicious devices, see [Wikipedia:Input-output memory management unit#Advantages](#) and [Memory Management \(computer programming\): Could you explain IOMMU in plain English?](#) (<https://www.quora.com/Memory-Management-computer-programming/Could-you-explain-IOMMU-in-plain-English>).

## Enabling IOMMU

Ensure that AMD-Vi/Intel VT-d is supported by the CPU and enabled in the BIOS settings. Both normally show up

alongside other CPU features (meaning they could be in an overclocking-related menu) either with their actual names ("VT-d" or "AMD-Vi") or in more ambiguous terms such as "Virtualization technology", which may or may not be explained in the manual.

Enable IOMMU support by setting the correct [kernel parameter](#) depending on the type of CPU in use:

- For Intel CPUs (VT-d) set `intel_iommu=on`
- For AMD CPUs (AMD-Vi) set `amd_iommu=on`

You should also append the `iommu=pt` parameter. This will prevent Linux from touching devices which cannot be passed through.

After rebooting, check `dmesg` to confirm that IOMMU has been correctly enabled:

```
dmesg | grep -i -e DMAR -e IOMMU

[ 0.000000] ACPI: DMAR 0x00000000BDCB1CB0 0000B8 (v01 INTEL BDW 00000001 INTL 00000001)
[ 0.000000] Intel-IOMMU: enabled
[ 0.028879] dmar: IOMMU 0: reg_base_addr fed90000 ver 1:0 cap c0000020660462 ecap f0101a
[ 0.028883] dmar: IOMMU 1: reg_base_addr fed91000 ver 1:0 cap d2008c20660462 ecap f010da
[ 0.028950] IOAPIC id 8 under DRHD base 0xfed91000 IOMMU 1
[ 0.536212] DMAR: No ATSR found
[ 0.536229] IOMMU 0 0xfed90000: using Queued invalidation
[ 0.536230] IOMMU 1 0xfed91000: using Queued invalidation
[ 0.536231] IOMMU: Setting RMRR:
[ 0.536241] IOMMU: Setting identity map for device 0000:00:02.0 [0xbf000000 - 0xcf1ffffff]
[ 0.537490] IOMMU: Setting identity map for device 0000:00:14.0 [0xbdea8000 - 0xbdeb6fff]
[ 0.537512] IOMMU: Setting identity map for device 0000:00:1a.0 [0xbdea8000 - 0xbdeb6fff]
[ 0.537530] IOMMU: Setting identity map for device 0000:00:1d.0 [0xbdea8000 - 0xbdeb6fff]
[ 0.537543] IOMMU: Prepare 0-16MiB unity mapping for LPC
[ 0.537549] IOMMU: Setting identity map for device 0000:00:1f.0 [0x0 - 0xffffffff]
[ 2.182790] [drm] DMAR active, disabling use of stolen memory
```

## Ensuring that the groups are valid

The following script should allow you to see how your various PCI devices are mapped to IOMMU groups. If it does not return anything, you either have not enabled IOMMU support properly or your hardware does not support it.

```
#!/bin/bash
shopt -s nullglob
for g in /sys/kernel/iommu_groups/*; do
    echo "IOMMU Group ${g##*/}:"
    for d in $g/devices/*; do
        echo -e "\t${lspci -nns ${d##*/}}"
    done;
done;
```

Example output:

```
IOMMU Group 1:
    00:01.0 PCI bridge: Intel Corporation Xeon E3-1200 v2/3rd Gen Core processor PCI Express
Root Port [8086:0151] (rev 09)
```

```
IOMMU Group 2:
    00:14.0 USB controller: Intel Corporation 7 Series/C210 Series Chipset Family USB xHCI Host Controller [8086:0e31] (rev 04)
IOMMU Group 4:
    00:1a.0 USB controller: Intel Corporation 7 Series/C210 Series Chipset Family USB Enhanced Host Controller #2 [8086:0e2d] (rev 04)
IOMMU Group 10:
    00:1d.0 USB controller: Intel Corporation 7 Series/C210 Series Chipset Family USB Enhanced Host Controller #1 [8086:0e26] (rev 04)
IOMMU Group 13:
    06:00.0 VGA compatible controller: NVIDIA Corporation GM204 [GeForce GTX 970] [10de:13c2] (rev a1)
    06:00.1 Audio device: NVIDIA Corporation GM204 High Definition Audio Controller [10de:0fb] (rev a1)
```

An IOMMU group is the smallest set of physical devices that can be passed to a virtual machine. For instance, in the example above, both the GPU in 06:00.0 and its audio controller in 6:00.1 belong to IOMMU group 13 and can only be passed together. The frontal USB controller, however, has its own group (group 2) which is separate from both the USB expansion controller (group 10) and the rear USB controller (group 4), meaning that [any of them could be passed to a VM without affecting the others](#).

## Gotchas

### Plugging your guest GPU in an unisolated CPU-based PCIe slot

Not all PCI-E slots are the same. Most motherboards have PCIe slots provided by both the CPU and the PCH. Depending on your CPU, it is possible that your processor-based PCIe slot does not support isolation properly, in which case the PCI slot itself will appear to be grouped with the device that is connected to it.

```
IOMMU Group 1:
    00:01.0 PCI bridge: Intel Corporation Xeon E3-1200 v2/3rd Gen Core processor PCI Express Root Port (rev 09)
    01:00.0 VGA compatible controller: NVIDIA Corporation GM107 [GeForce GTX 750] (rev a2)
    01:00.1 Audio device: NVIDIA Corporation Device 0fbc (rev a1)
```

This is fine so long as only your guest GPU is included in here, such as above. Depending on what is plugged in to your other PCIe slots and whether they are allocated to your CPU or your PCH, you may find yourself with additional devices within the same group, which would force you to pass those as well. If you are ok with passing everything that is in there to your VM, you are free to continue. Otherwise, you will either need to try and plug your GPU in your other PCIe slots (if you have any) and see if those provide isolation from the rest or to install the ACS override patch, which comes with its own drawbacks. See [#Bypassing the IOMMU groups \(ACS override patch\)](#) for more information.

**Note:** If they are grouped with other devices in this manner, pci root ports and bridges should neither be bound to vfio at boot, nor be added to the VM.

## Isolating the GPU

In order to assign a device to a virtual machine, this device and all those sharing the same IOMMU group must have their

driver replaced by a stub driver or a VFIO driver in order to prevent the host machine from interacting with them. In the case of most devices, this can be done on the fly right before the VM starts.

However, due to their size and complexity, GPU drivers do not tend to support dynamic rebinding very well, so you cannot just have some GPU you use on the host be transparently passed to a VM without having both drivers conflict with each other. Because of this, it is generally advised to bind those placeholder drivers manually before starting the VM, in order to stop other drivers from attempting to claim it.

The following section details how to configure a GPU so those placeholder drivers are bound early during the boot process, which makes said device inactive until a VM claims it or the driver is switched back. This is the preferred method, considering it has less caveats than switching drivers once the system is fully online.

**Warning:** Once you reboot after this procedure, whatever GPU you have configured will no longer be usable on the host until you reverse the manipulation. Make sure the GPU you intend to use on the host is properly configured before doing this - your motherboard should be set to display using the host GPU.

Starting with Linux 4.1, the kernel includes `vfio-pci`. This is a VFIO driver, meaning it fulfills the same role as `pci-stub` did, but it can also control devices to an extent, such as by switching them into their D3 state when they are not in use.

## Binding `vfio-pci` via device ID

`Vfio-pci` normally targets PCI devices by ID, meaning you only need to specify the IDs of the devices you intend to passthrough. For the following IOMMU group, you would want to bind `vfio-pci` with `10de:13c2` and `10de:0fbb`, which will be used as example values for the rest of this section.

```
IOMMU Group 13:
    06:00.0 VGA compatible controller: NVIDIA Corporation GM204 [GeForce GTX 970] [10de:13c2] (rev a1)
    06:00.1 Audio device: NVIDIA Corporation GM204 High Definition Audio Controller [10de:0fbb] (rev a1)}
```

### Note:

- You cannot specify which device to isolate using vendor-device ID pairs if the host GPU and the guest GPU share the same pair (i.e. : if both are the same model). If this is your case, read [#Using identical guest and host GPUs](#) instead.
- If, as noted in [#Plugging your guest GPU in an unisolated CPU-based PCIe slot](#), your `pci` root port is part of your IOMMU group, you **must not** pass its ID to `vfio-pci`, as it needs to remain attached to the host to function properly. Any other device within that group, however, should be left for `vfio-pci` to bind with.

Two methods exist for providing the device IDs. Specifying them via [kernel parameters](#) has the advantage of being able to easily edit, remove, or undo any breaking changes via your boot loader:

```
vfio-pci.ids=10de:13c2,10de:0fbb
```

Alternatively, the IDs may be added to a modprobe conf file. Since these conf files are embedded in the initramfs image, any changes require regenerating a new image each time:

```
/etc/modprobe.d/vfio.conf
-----
options vfio-pci ids=10de:13c2,10de:0fbb
```

## Loading vfio-pci early

Since Arch's [linux](https://www.archlinux.org/packages/?name=linux) (<https://www.archlinux.org/packages/?name=linux>) has vfio-pci built as a module, we need to force it to load early before the graphics drivers have a chance to bind to the card. To ensure that, add `vfio_pci`, `vfio`, `vfio_iommu_type1`, and `vfio_virqfd` to [mkinitcpio](#):

```
/etc/mkinitcpio.conf
-----
MODULES=(... vfio_pci vfio vfio_iommu_type1 vfio_virqfd ...)
```

**Note:** If you also have another driver loaded this way for [early modesetting](#) (such as `nouveau`, `radeon`, `amdgpu`, `i915`, etc.), all of the aforementioned VFIO modules must precede it.

Also, ensure that the modconf hook is included in the HOOKS list of `mkinitcpio.conf`:

```
/etc/mkinitcpio.conf
-----
HOOKS=(... modconf ...)
```

Since new modules have been added to the initramfs configuration, you must [regenerate the initramfs](#).

## Verifying that the configuration worked

Reboot and verify that vfio-pci has loaded properly and that it is now bound to the right devices.

```
$ dmesg | grep -i vfio
[ 0.329224] VFIO - User Level meta-driver version: 0.3
[ 0.341372] vfio_pci: add [10de:13c2[ffff:ffff]] class 0x0000000/00000000
[ 0.354704] vfio_pci: add [10de:0fbb[ffff:ffff]] class 0x0000000/00000000
[ 2.061326] vfio-pci 0000:06:00.0: enabling device (0100 -> 0103)
```

It is not necessary for all devices (or even expected device) from `vfio.conf` to be in dmesg output. Sometimes a device does not appear in output at boot but actually is able to be visible and operatable in guest VM.

```
$ lspci -nnk -d 10de:13c2
```

```
06:00.0 VGA compatible controller: NVIDIA Corporation GM204 [GeForce GTX 970] [10de:13c2] (rev a1)
Kernel driver in use: vfio-pci
Kernel modules: nouveau nvidia
```

```
$ lspci -nnk -d 10de:0fbb
```

```
06:00.1 Audio device: NVIDIA Corporation GM204 High Definition Audio Controller [10de:0fbb] (rev a1)
Kernel driver in use: vfio-pci
Kernel modules: snd_hda_intel
```

## Setting up an OVMF-based guest VM

OVMF is an open-source UEFI firmware for QEMU virtual machines. While it is possible to use SeaBIOS to get similar results to an actual PCI passthrough, the setup process is different and it is generally preferable to use the EFI method if your hardware supports it.

### Configuring libvirt

[Libvirt](#) is a wrapper for a number of virtualization utilities that greatly simplifies the configuration and deployment process of virtual machines. In the case of KVM and QEMU, the frontend it provides allows us to avoid dealing with the permissions for QEMU and make it easier to add and remove various devices on a live VM. Its status as a wrapper, however, means that it might not always support all of the latest qemu features, which could end up requiring the use of a wrapper script to provide some extra arguments to QEMU.

Install [qemu](https://www.archlinux.org/packages/?name=qemu) (<https://www.archlinux.org/packages/?name=qemu>), [libvirt](https://www.archlinux.org/packages/?name=libvirt) (<https://www.archlinux.org/packages/?name=libvirt>), [edk2-ovmf](https://www.archlinux.org/packages/?name=edk2-ovmf) (<https://www.archlinux.org/packages/?name=edk2-ovmf>), and [virt-manager](https://www.archlinux.org/packages/?name=virt-manager) (<https://www.archlinux.org/packages/?name=virt-manager>).

You can now [enable](#) and [start](#) `libvirtd.service` and its logging component `virtlogd.socket`.

You may also need to [activate the default libvirt network](https://wiki.libvirt.org/page/Networking#NAT_forwarding_.28aka_.22virtual_networks.22.29) ([https://wiki.libvirt.org/page/Networking#NAT\\_forwarding\\_.28aka\\_.22virtual\\_networks.22.29](https://wiki.libvirt.org/page/Networking#NAT_forwarding_.28aka_.22virtual_networks.22.29)).

### Setting up the guest OS

The process of setting up a VM using `virt-manager` is mostly self-explanatory, as most of the process comes with fairly comprehensive on-screen instructions.

However, you should pay special attention to the following steps:

- When the VM creation wizard asks you to name your VM (final step before clicking "Finish"), check the "Customize before install" checkbox.
- In the "Overview" section, **set your firmware to "UEFI"** (<https://i.imgur.com/73r2ctM.png>). If the option is grayed out, make sure that:
  - Your hypervisor is running as a system session and not a user session. This can be verified **by clicking, then hovering** (<https://i.ibb.co/N1XZCdp/Deepin-Screenshot-select-area-20190125113216.png>) over the session in virt-manager. If you are accidentally running it as a user session, you must open a new connection by clicking "File" > "Add Connection..", then select the option from the drop-down menu station "QEMU/KVM" and not "QEMU/KVM user session".
  - In the "CPUs" section, change your CPU model to "host-passthrough". If it is not in the list, you will have to type it by hand. This will ensure that your CPU is detected properly, since it causes libvirt to expose your CPU capabilities exactly as they are instead of only those it recognizes (which is the preferred default behavior to make CPU behavior easier to reproduce). Without it, some applications may complain about your CPU being of an unknown model.
  - If you want to minimize IO overhead, it's easier to setup **#Virtio disk** before installing

The rest of the installation process will take place as normal using a standard QXL video adapter running in a window. At this point, there is no need to install additional drivers for the rest of the virtual devices, since most of them will be removed later on. Once the guest OS is done installing, simply turn off the virtual machine. It is possible you will be dropped into the UEFI menu instead of starting the installation upon powering your VM for the first time. Sometimes the correct ISO file was not automatically detected and you will need to manually specify the drive to boot. By typing exit and navigating to "boot manager" you will enter a menu that allows you to choose between devices.

## Attaching the PCI devices

With the installation done, it is now possible to edit the hardware details in libvirt and remove virtual integration devices, such as the spice channel and virtual display, the QXL video adapter, the emulated mouse and keyboard and the USB tablet device. Since that leaves you with no input devices, you may want to bind a few USB host devices to your VM as well, but remember to **leave at least one mouse and/or keyboard assigned to your host** in case something goes wrong with the guest. At this point, it also becomes possible to attach the PCI device that was isolated earlier; simply click on "Add Hardware" and select the PCI Host Devices you want to passthrough. If everything went well, the screen plugged into your GPU should show the OVMF splash screen and your VM should start up normally. From there, you can setup the drivers for the rest of your VM.

## Passing keyboard/mouse via Evdev

If you do not have a spare mouse or keyboard to dedicate to your guest, and you do not want to suffer from the video overhead of Spice, you can setup evdev to share them between your Linux host and your virtual machine.

**Note:** Press both left and right **Ctrl** keys at the same time to swap control between the host and the guest.

First, modify the libvirt configuration. If you encounter failed to get domain [vmname] error, add `-c qemu:///system` arguments to `virsh`.

```
$ virsh edit [vmname]
```

```
<domain type='kvm'>
```

to

```
$ virsh edit [vmname]
```

```
<domain type='kvm' xmlns:qemu='http://libvirt.org/schemas/domain/qemu/1.0'>
```

Next, find your keyboard and mouse devices in `/dev/input/by-id/`. You may find multiple devices associated to your mouse or keyboard, so try `cat /dev/input/by-id/device_id` and either hit some keys on the keyboard or wiggle your mouse to see if input comes through, if so you have got the right device. Now add those devices to your configuration right before the closing `</domain>` tag:

```
$ virsh edit [vmname]
```

```
...
<qemu:commandline>
<qemu:arg value='-object' />
<qemu:arg value='input-linux,id=mouse1,evdev=/dev/input/by-id/MOUSE_NAME' />
<qemu:arg value='-object' />
<qemu:arg value='input-linux,id=kbd1,evdev=/dev/input/by-id/KEYBOARD_NAME,grab_all=on,repeat=on' />
</qemu:commandline>
...
```

Replace `MOUSE_NAME` and `KEYBOARD_NAME` with your device id. You will also need to include these devices in your qemu config, and setting the user and group to one that has access to your input devices:

```
/etc/libvirt/qemu.conf
```

```
...
user = "<your_user>"
group = "kvm"
...
cgroup_device_acl = [
"/dev/kvm",
"/dev/input/by-id/KEYBOARD_NAME",
"/dev/input/event0", # Destination of the symbolic link above
"/dev/input/by-id/MOUSE_NAME",
"/dev/input/event1", # Destination of the symbolic link above
"/dev/null", "/dev/full", "/dev/zero",
"/dev/random", "/dev/urandom",
"/dev/ptmx", "/dev/kvm", "/dev/kqemu",
"/dev/rtc", "/dev/hpet", "/dev/sev"
]
...
```

Then ensure that the user you provided has access to the `kvm` and `input` [user groups](#). [Restart](#) `libvirtd.service`. Now you can startup the guest OS and test swapping control of your mouse and keyboard between the host and guest by pressing both the left and right control keys at the same time.

You may also consider switching from PS/2 to Virtio inputs in your configurations:

```
$ virsh edit [vmname]

...
<input type='mouse' bus='virtio'>
  <address type='pci' domain='0x0000' bus='0x00' slot='0x0e' function='0x0' />
</input>
<input type='keyboard' bus='virtio'>
  <address type='pci' domain='0x0000' bus='0x00' slot='0x0f' function='0x0' />
</input>
<input type='mouse' bus='ps2' />
<input type='keyboard' bus='ps2' />
...
```

Next startup the guest OS and install the virtIO drivers for those devices.

## Gotchas

### Using a non-EFI image on an OVMF-based VM

The OVMF firmware does not support booting off non-EFI mediums. If the installation process drops you in a UEFI shell right after booting, you may have an invalid EFI boot media. Try using an alternate Linux/Windows image to determine if you have an invalid media.

## Performance tuning

---

Most use cases for PCI passthroughs relate to performance-intensive domains such as video games and GPU-accelerated tasks. While a PCI passthrough on its own is a step towards reaching native performance, there are still a few adjustments on the host and guest to get the most out of your VM.

### CPU pinning

The default behavior for KVM guests is to run operations coming from the guest as a number of threads representing virtual processors. Those threads are managed by the Linux scheduler like any other thread and are dispatched to any available CPU cores based on niceness and priority queues. As such, the local CPU cache benefits (L1/L2) are lost each time the host scheduler reschedules the virtual CPU thread on a different physical CPU. This can noticeably harm performance on the guest. CPU pinning aims to resolve this by limiting which physical CPUs the virtual CPUs are allowed to run on. The ideal setup is a one to one mapping such that the virtual CPU cores match physical CPU cores while taking hyperthreading/SMT into account.

**Note:** For certain users enabling CPU pinning may introduce stuttering and short hangs, especially with the MuQSS scheduler (present in linux-ck and linux-zen kernels). You might want to try disabling pinning first if you experience similar issues, which effectively trades maximum performance for responsiveness at all times.

## CPU topology

Most modern CPUs support hardware multitasking, also known as hyper-threading on Intel CPUs or SMT on AMD CPUs. Hyper-threading/SMT is simply a very efficient way of running two threads on one CPU core at any given time. You will want to take into consideration that the CPU pinning you choose will greatly depend on what you do with your host while your VM is running.

To find the topology for your CPU run `lscpu -e`:

**Note:** Pay special attention to the 4th column "**CORE**" as this shows the association of the Physical/Logical CPU cores

`lscpu -e` on a 6c/12t Ryzen 5 1600:

CPU	NODE	SOCKET	CORE	L1d:L1i:L2:L3	ONLINE	MAXMHZ	MINMHZ
0	0	0	0	0:0:0:0	yes	3800.0000	1550.0000
1	0	0	0	0:0:0:0	yes	3800.0000	1550.0000
2	0	0	1	1:1:1:0	yes	3800.0000	1550.0000
3	0	0	1	1:1:1:0	yes	3800.0000	1550.0000
4	0	0	2	2:2:2:0	yes	3800.0000	1550.0000
5	0	0	2	2:2:2:0	yes	3800.0000	1550.0000
6	0	0	3	3:3:3:1	yes	3800.0000	1550.0000
7	0	0	3	3:3:3:1	yes	3800.0000	1550.0000
8	0	0	4	4:4:4:1	yes	3800.0000	1550.0000
9	0	0	4	4:4:4:1	yes	3800.0000	1550.0000
10	0	0	5	5:5:5:1	yes	3800.0000	1550.0000
11	0	0	5	5:5:5:1	yes	3800.0000	1550.0000

**Note:** Ryzen 3000 ComboPi AGESA changes topology to match Intel example, even on prior generation CPUs. Above valid only on older AGESA.

`lscpu -e` on a 6c/12t Intel 8700k:

CPU	NODE	SOCKET	CORE	L1d:L1i:L2:L3	ONLINE	MAXMHZ	MINMHZ
0	0	0	0	0:0:0:0	yes	4600.0000	800.0000
1	0	0	1	1:1:1:0	yes	4600.0000	800.0000
2	0	0	2	2:2:2:0	yes	4600.0000	800.0000
3	0	0	3	3:3:3:0	yes	4600.0000	800.0000
4	0	0	4	4:4:4:0	yes	4600.0000	800.0000
5	0	0	5	5:5:5:0	yes	4600.0000	800.0000
6	0	0	0	0:0:0:0	yes	4600.0000	800.0000
7	0	0	1	1:1:1:0	yes	4600.0000	800.0000
8	0	0	2	2:2:2:0	yes	4600.0000	800.0000
9	0	0	3	3:3:3:0	yes	4600.0000	800.0000
10	0	0	4	4:4:4:0	yes	4600.0000	800.0000
11	0	0	5	5:5:5:0	yes	4600.0000	800.0000

As we see above, with AMD **Core 0** is sequential with **CPU 0 & 1**, whereas Intel places **Core 0** on **CPU 0 & 6**.

If you do not need all cores for the guest, it would then be preferable to leave at the very least one core for the host. Choosing which cores one to use for the host or guest should be based on the specific hardware characteristics of your CPU, however **Core 0** is a good choice for the host in most cases. If any cores are reserved for the host, it is recommended to pin the emulator and iotreads, if used, to the host cores rather than the VCPUs. This may improve performance and reduce latency for the guest since those threads will not pollute the cache or contend for scheduling with the guest VCPU threads. If all cores are passed to the guest, there is no need or benefit to pinning the emulator or iotreads.

## XML examples

**Note:** Do not use the **iothread** lines from the XML examples shown below if you have not added an **iothread** to your disk controller. **iothread**'s only work on **virtio-scsi** or **virtio-blk** devices.

### 4c/1t CPU w/o Hyperthreading Example

```
$ virsh edit [vmname]
...
<vcpu placement='static'>4</vcpu>
<cputune>
  <vcpupin vcpu='0' cpuset='0' />
  <vcpupin vcpu='1' cpuset='1' />
  <vcpupin vcpu='2' cpuset='2' />
  <vcpupin vcpu='3' cpuset='3' />
</cputune>
...
```

### 4c/2t Intel/AMD CPU example (after ComboPI AGESA bios update)

```
$ virsh edit [vmname]
...
<vcpu placement='static'>8</vcpu>
<iothreads>1</iothreads>
<cputune>
  <vcpupin vcpu='0' cpuset='2' />
  <vcpupin vcpu='1' cpuset='8' />
  <vcpupin vcpu='2' cpuset='3' />
  <vcpupin vcpu='3' cpuset='9' />
  <vcpupin vcpu='4' cpuset='4' />
  <vcpupin vcpu='5' cpuset='10' />
  <vcpupin vcpu='6' cpuset='5' />
  <vcpupin vcpu='7' cpuset='11' />
  <emulatorpin cpuset='0,6' />
  <iothreadpin iothread='1' cpuset='0,6' />
</cputune>
...
<topology sockets='1' cores='4' threads='2' />
...
```

## 4c/2t AMD CPU example (Before ComboPi AGESA bios update)

```
$ virsh edit [vmname]

...
<vcpu placement='static'>8</vcpu>
<iothreads>1</iothreads>
<cputune>
  <vcupin vcpu='0' cpuset='2' />
  <vcupin vcpu='1' cpuset='3' />
  <vcupin vcpu='2' cpuset='4' />
  <vcupin vcpu='3' cpuset='5' />
  <vcupin vcpu='4' cpuset='6' />
  <vcupin vcpu='5' cpuset='7' />
  <vcupin vcpu='6' cpuset='8' />
  <vcupin vcpu='7' cpuset='9' />
  <emulatorpin cpuset='0-1' />
  <iothreadpin iothread='1' cpuset='0-1' />
</cputune>
...
  <topology sockets='1' cores='4' threads='2' />
...
```

**Note:** If further CPU isolation is needed, consider using the **isolepus** kernel command-line parameter on the unused physical/logical cores.

If you do not intend to be doing any computation-heavy work on the host (or even anything at all) at the same time as you would on the VM, you may want to pin your VM threads across all of your cores, so that the VM can fully take advantage of the spare CPU time the host has available. Be aware that pinning all physical and logical cores of your CPU could induce latency in the guest VM.

## Huge memory pages

When dealing with applications that require large amounts of memory, memory latency can become a problem since the more memory pages are being used, the more likely it is that this application will attempt to access information across multiple memory "pages", which is the base unit for memory allocation. Resolving the actual address of the memory page takes multiple steps, and so CPUs normally cache information on recently used memory pages to make subsequent uses on the same pages faster. Applications using large amounts of memory run into a problem where, for instance, a virtual machine uses 4 GiB of memory divided into 4 KiB pages (which is the default size for normal pages) for a total of 1.04 million pages, meaning that such cache misses can become extremely frequent and greatly increase memory latency. Huge pages exist to mitigate this issue by giving larger individual pages to those applications, increasing the odds that multiple operations will target the same page in succession.

## Transparent huge pages

QEMU will use 2MiB sized transparent huge pages automatically without any explicit configuration in QEMU or Libvirt, subject to some important caveats. When using VFIO the pages are locked in at boot time and transparent huge pages are allocated up front when the VM first boots. If the kernel memory is highly fragmented, or the VM is using a majority of the remaining free memory, it is likely that the kernel will not have enough 2MiB pages to fully satisfy the allocation. In

such a case, it silently fails by using a mix of 2MiB and 4KiB pages. Since the pages are locked in VFIO mode, the kernel will not be able to convert those 4KiB pages to huge after the VM starts either. The number of available 2MiB huge pages available to THP is the same as via the [#Dynamic huge pages](#) mechanism described in the following sections.

To check how much memory THP is using globally:

```
$ grep AnonHugePages /proc/meminfo
-----
AnonHugePages: 8091648 kB
```

To check a specific QEMU instance. QEMU's PID must be substituted in the grep command:

```
$ grep -P 'AnonHugePages:\s+(?!0)\d+' /proc/[PID]/smaps
-----
AnonHugePages: 8087552 kB
```

In this example, the VM was allocated 8388608KiB of memory, but only 8087552KiB was available via THP. The remaining 301056KiB are allocated as 4KiB pages. Aside from manually checking, there is no indication when partial allocations occur. As such, THP's effectiveness is very much dependent on the host system's memory fragmentation at the time of VM startup. If this trade off is unacceptable or strict guarantees are required, [#Static huge pages](#) is recommended.

Arch kernels have THP compiled in and enabled by default with `/sys/kernel/mm/transparent_hugepage/enabled` set to `madvise` mode.

## Static huge pages

While transparent huge pages should work in the vast majority of cases, they can also be allocated statically during boot. This should only be needed to make use 1 GiB hugepages on machines that support it, since transparent huge pages normally only go up to 2 MiB.

**Warning:** Static huge pages lock down the allocated amount of memory, making it unavailable for applications that are not configured to use them. Allocating 4 GiBs worth of huge pages on a machine with 8 GiB of memory will only leave you with 4 GiB of available memory on the host **even when the VM is not running**.

To allocate huge pages at boot, one must simply specify the desired amount on their kernel command line with `hugepages=x`. For instance, reserving 1024 pages with `hugepages=1024` and the default size of 2048 KiB per huge page creates 2 GiB worth of memory for the virtual machine to use.

If supported by CPU page size could be set manually. 1 GiB huge page support could be verified by `grep pdpe1gb /proc/cpuinfo`. Setting 1 GiB huge page size via kernel parameters: `default_hugepagesz=1G hugepagesz=1G hugepages=X`.

Also, since static huge pages can only be used by applications that specifically request it, you must add this section in your

libvirt domain configuration to allow kvm to benefit from them :

```
$ virsh edit [vmname]
-----
...
<memoryBacking>
  <hugepages/>
</memoryBacking>
...
```

## Dynamic huge pages

Hugepages could be allocated manually via `vm.nr_overcommit_hugepages` [sysctl](#) parameter.

```
/etc/sysctl.d/10-kvm.conf
-----
vm.nr_hugepages = 0
vm.nr_overcommit_hugepages = num
```

Where *num* - is the number of huge pages, which default size if 2 MiB. Pages will be automatically allocated, and freed after VM stops.

More manual way:

```
# echo num > /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages
# echo num > /sys/kernel/mm/hugepages/hugepages-1048576kB/nr_hugepages
```

For 2 MiB and 1 GiB page size respectively. And they should be manually freed in the same way.

It is hardly recommended to drop caches, compact memory and wait couple of seconds before starting VM, as there could be not enough free contiguous memory for required huge pages blocks. Especially after some uptime of the host system.

```
# echo 3 > /proc/sys/vm/drop_caches
# echo 1 > /proc/sys/vm/compact_memory
```

Theoretically, 1 GiB pages works as 2 MiB. But practically - no guaranteed way was found to get contiguous 1 GiB memory blocks. Each consequent request of 1 GiB blocks lead to lesser and lesser dynamically allocated count.

## CPU frequency governor

Depending on the way your [CPU governor](#) is configured, the VM threads may not hit the CPU load thresholds for the frequency to ramp up. Indeed, KVM cannot actually change the CPU frequency on its own, which can be a problem if it does not scale up with vCPU usage as it would result in underwhelming performance. An easy way to see if it behaves correctly is to check if the frequency reported by `watch lscpu` goes up when running a CPU-intensive task on the guest. If you are indeed experiencing stutter and the frequency does not go up to reach its reported maximum, it may be

due to [cpu scaling being controlled by the host OS \(https://lime-technology.com/forum/index.php?topic=46664.msg447678#msg447678\)](https://lime-technology.com/forum/index.php?topic=46664.msg447678#msg447678). In this case, try setting all cores to maximum frequency to see if this improves performance. Note that if you are using a modern intel chip with the default pstate driver, cpupower commands will be ineffective, so monitor `/proc/cpuinfo` to make sure your cpu is actually at max frequency.

## Isolating pinned CPUs

CPU pinning by itself will not prevent other host processes from running on the pinned CPUs. Properly isolating the pinned CPUs can reduce latency in the guest VM.

### With `isolcpus` kernel parameter

In this example, let us assume you are using CPUs 4-7. Use the [kernel parameters](#) `isolcpus nohz_full rcu_nocbs` to completely isolate the CPUs from the kernel. For example:

```
isolcpus=4-7 nohz_full=4-7 rcu_nocbs=4-7
```

Then, run `qemu-system-x86_64` with `taskset` and `chrt`:

```
# chrt -r 1 taskset -c 4-7 qemu-system-x86_64 ...
```

The `chrt` command will ensure that the task scheduler will round-robin distribute work (otherwise it will all stay on the first cpu). For `taskset`, the CPU numbers can be comma- and/or dash-separated, like "0,1,2,3" or "0-4" or "1,7-8,10" etc.

See [this Removeddit mirror of a Reddit thread \(https://www.removeddit.com/r/VFIO/comments/6vgtpx/high\\_dpc\\_latency\\_and\\_audio\\_stuttering\\_on\\_windows/dm0sfto/\)](https://www.removeddit.com/r/VFIO/comments/6vgtpx/high_dpc_latency_and_audio_stuttering_on_windows/dm0sfto/) for more info. ([The original thread \(https://www.reddit.com/r/VFIO/comments/6vgtpx/high\\_dpc\\_latency\\_and\\_audio\\_stuttering\\_on\\_windows/dm0sfto/\)](https://www.reddit.com/r/VFIO/comments/6vgtpx/high_dpc_latency_and_audio_stuttering_on_windows/dm0sfto/) is worthless because of deleted comments.)

### Dynamically isolating CPUs

The `isolcpus` kernel parameter will permanently reserve CPU cores, even when the guest is not running. A more flexible alternative is to use the `cset` tool from [cpuset-git \(https://aur.archlinux.org/packages/cpuset-git/\)](https://aur.archlinux.org/packages/cpuset-git/)<sup>AUR</sup> to dynamically isolate CPUs when starting the guest.

See this [vfio-users post \(https://www.redhat.com/archives/vfio-users/2016-September/msg00072.html\)](https://www.redhat.com/archives/vfio-users/2016-September/msg00072.html) for more info, as well as this [blog post \(https://rokups.github.io/#!pages/gaming-vm-performance.md\)](https://rokups.github.io/#!pages/gaming-vm-performance.md) and this [script \(https://github.com/PassthroughPOST/VFIO-Tools/blob/master/libvirt\\_hooks/hooks/cset.sh\)](https://github.com/PassthroughPOST/VFIO-Tools/blob/master/libvirt_hooks/hooks/cset.sh) for working examples.

## Improving performance on AMD CPUs

Previously, Nested Page Tables (NPT) had to be disabled on AMD systems running KVM to improve GPU performance because of a **very old bug** (<https://sourceforge.net/p/kvm/bugs/230/>), but the trade off was decreased CPU performance, including stuttering.

There is a **kernel patch** (<https://patchwork.kernel.org/patch/10027525/>) that resolves this issue, which was accepted into kernel 4.14-stable and 4.9-stable. If you are running the official **linux** (<https://www.archlinux.org/packages/?name=linux>) or **linux-lts** (<https://www.archlinux.org/packages/?name=linux-lts>) kernel the patch has already been applied (make sure you are on the latest). If you are running another kernel you might need to manually patch yourself.

**Note:** Several Ryzen users (see [this Reddit thread \(https://www.reddit.com/r/VFIO/comments/78i3jx/possible\\_fix\\_for\\_the\\_npt\\_issue\\_discussed\\_on\\_iommu/\)](https://www.reddit.com/r/VFIO/comments/78i3jx/possible_fix_for_the_npt_issue_discussed_on_iommu/)) have tested the patch, and can confirm that it works, bringing GPU passthrough performance up to near native quality.

Starting with QEMU 3.1 the TOPOEXT cpuid flag is disabled by default. In order to use hyperthreading(SMT) on AMD CPUs you need to manually enable it:

```
<cpu mode='host-passthrough' check='none'>
<topology sockets='1' cores='4' threads='2'/>
<feature policy='require' name='topoext'/>
</cpu>
```

commit: <https://git.qemu.org/?p=qemu.git;a=commit;h=7210a02c58572b2686a3a8d610c6628f87864aed>

## Virtio disk

The default disk types are SATA or IDE emulation out of the box. These controllers offer maximum compatibility but are not suited for efficient virtualization. Two accelerated models exist: `virtio-scsi` for SCSI emulation and passthrough, or `virtio-blk` for a more basic block device emulation.

## Drivers

- Linux guests should support these out of the box on any modern kernel
- macOS has rudimentary `virtio-blk` support starting in Mojave via `AppleVirtIO.kext` but performance is poor and TRIM does not appear to work. SATA is still suggested
- Windows needs the **Windows virtio drivers** (<https://docs.fedoraproject.org/en-US/quick-docs/creating-windows-virtual-machines-using-virtio-drivers/>). `virtio-scsi` uses the `vioscsi` driver. `virtio-blk` uses the `viostor` driver
- Windows can be installed directly onto these disks by selecting 'load driver' on the installer disk selection menu. The windows iso and virtio driver iso should both be attached as regular SATA/IDE cdroms during the installation process
- To switch boot disks to virtio on an existing Windows installation:
  - `virtio-blk`: Add a temporary disk with bus `virtio`, boot windows & load the driver for the disk, then shutdown and switch the boot disk disk bus to `virtio`

- `virtio-scsi` : Add a scsi controller with model `virtio` , boot windows & load the driver for the controller, then shutdown and switch the boot disk bus to `scsi` (not `virtio`)

## Considerations

- `virtio-scsi` TRIM support is mature, all versions should support it. Traditionally, `virtio-scsi` has been the preferred approach for this reason
- `virtio-blk` TRIM support is new, this requires `qemu 4.0+`, guest linux kernel `5.0+`, guest windows drivers `0.1.173+`
- Thin provisioning works by enabling TRIM on a sparse image file: `discard='unmap'` . Unused blocks will be freed and the disk usage will drop (works on both raw and `qcow2`). Actual on-disk size of a sparse image file may be checked with `du /path/to/disk.img`
- Thin provisioning can also work with block storage such as `zfs zvols` or `thin lvm`
- Virt queue count will influence the number of threads inside the guest kernel used for IO processing, suggest using `queues='4'` or more
- Native mode ( `io='native'` ) uses a single threaded model based on linux AIO, is a bit more CPU efficient but may have lower peak performance and does not allow host side caching to be used
- Threaded mode ( `io='threads'` ) will spawn dozens of threads on demand as the disk is used. This is less efficient but may perform better if there are enough host cores available to run them, and allows for host side caching to be used
- Modern versions of `libvirt` will group the dynamic worker threads created when using threaded mode in with the `iothread=1` cgroup for pinning purposes. Very old versions of `libvirt` left these in the emulator cgroup

## IO threads

An IO thread is a dedicated thread for processing disk events, rather than using the main `qemu` emulator loop. This should not be confused with the worker threads spawned on demand with `io='threads'` .

- You can only use one `iothread` per disk controller. The thread must be assigned to a specific controller with `iothread='X'` in the `<driver>` tag. Furthermore, extra & unassigned `iothreads` will not be used and do nothing
- In the case of `virtio-scsi` , there is one controller for multiple scsi disks. The `iothread` is assigned on the controller: `<controller><driver iothread='X'>`
- In the case of `virtio-blk` , each disk has its own controller. The `iothread` is assigned in the driver tag under the disk itself: `<disk><driver iothread='X'>`
- Since emulated disks incur a significant amount of CPU overhead, that can lead to `vcpu` stuttering under high disk load (especially high random IOPS). In this case it helps to pin the IO to different core(s) than your `vcpus` with `<iothreadpin>`

## Examples with libvirt

`virtio-scsi` + `iothread` + worker threads + host side writeback caching + full disk block device backend:

```
<domain>
```

```
<devices>
  <disk type='block' device='disk'>
    <driver name='qemu' type='raw' cache='writeback' io='threads' discard='unmap' />
    <source dev='/dev/disk/by-id/ata-Samsung_SSD_840_EVO_1TB_S1D9NSAF206396F' />
    <target dev='sda' bus='scsi' />
  </disk>
  <controller type='scsi' index='0' model='virtio-scsi'>
    <driver iothread='1' queues='8' />
  </controller>
```

virtio-blk + iothread + native aio + no host caching + raw sparse image backend:

```
<domain>
  <devices>
    <disk type='file' device='disk'>
      <driver name='qemu' type='raw' cache='none' io='native' discard='unmap' iothread='1' queues
='8' />
      <source dev='/var/lib/libvirt/images/pool/win10.img' />
      <target dev='vda' bus='virtio' />
    </disk>
```

Creating the iothreads:

```
<domain>
  <iothreads>1</iothreads>
```

Pinning iothreads:

```
<domain>
  <cputune>
    <iothreadpin iothread='1' cpuset='0-1,6-7' />
```

## Virtio network

The default NIC models rtl8139 or e1000 can be a bottleneck for gigabit+ speeds and have a significant amount of CPU overhead compared to `virtio-net`.

- Select `virtio` as the model for the NIC with libvirt or use the `virtio-net-pci` device in bare qemu
- Windows needs the NetKVM driver from [Windows virtio drivers \(https://docs.fedoraproject.org/en-US/quick-docs/creating-windows-virtual-machines-using-virtio-drivers/\)](https://docs.fedoraproject.org/en-US/quick-docs/creating-windows-virtual-machines-using-virtio-drivers/)
- Virtio uses `vhost-net` by default for in-kernel packet processing without exiting to userspace
- Multiqueue can be enabled for a further speedup with multiple connections but typically will not boost single stream speeds. For libvirt add `<driver queues='8' />` under the interface tag
- Zero copy transmit may also be enabled on `macvtap` by setting the module parameter `vhost_net.experimental_zcopytx=1` but this may actually have worse performance, see [commit \(https://github.com/torvalds/linux/commit/098eadce3c622c07b328d0a43dda379b38cf7c5e\)](https://github.com/torvalds/linux/commit/098eadce3c622c07b328d0a43dda379b38cf7c5e)

Libvirt example with a bridge:

```
<interface type='bridge'>
  <mac address="52:54:00:6d:6e:2e"/>
  <source bridge='br0' />
  <model type='virtio' />
  <driver queues='8' />
</interface>
```

## Further tuning

More specialized VM tuning tips are available at [Red Hat's Virtualization Tuning and Optimization Guide \(https://access.redhat.com/documentation/en-US/Red\\_Hat\\_Enterprise\\_Linux/7/html-single/Virtualization\\_Tuning\\_and\\_Optimization\\_Guide/index.html\)](https://access.redhat.com/documentation/en-US/Red_Hat_Enterprise_Linux/7/html-single/Virtualization_Tuning_and_Optimization_Guide/index.html).

## Special procedures

---

Certain setups require specific configuration tweaks in order to work properly. If you are having problems getting your host or your VM to work properly, see if your system matches one of the cases below and try adjusting your configuration accordingly.

## Using identical guest and host GPUs

Due to how vfio-pci uses your vendor and device id pair to identify which device they need to bind to at boot, if you have two GPUs sharing such an ID pair you will not be able to get your passthrough driver to bind with just one of them. This sort of setup makes it necessary to use a script, so that whichever driver you are using is instead assigned by pci bus address using the `driver_override` mechanism.

### Script variants

#### Passthrough all GPUs but the boot GPU

Here, we will make a script to bind vfio-pci to all GPUs but the boot gpu. Create the script `/usr/local/bin/vfio-pci-override.sh`:

```
#!/bin/sh

for i in /sys/bus/pci/devices/*/boot_vga; do
  if [ $(cat "$i") -eq 0 ]; then
    GPU="{i%/boot_vga}"
    AUDIO="$(echo "$GPU" | sed -e "s/0$/1/")"
    echo "vfio-pci" > "$GPU/driver_override"
    if [ -d "$AUDIO" ]; then
      echo "vfio-pci" > "$AUDIO/driver_override"
    fi
  fi
done

modprobe -i vfio-pci
```

## Passthrough selected GPU

In this case we manually specify the GPU to bind.

```
#!/bin/sh

DEVS="0000:03:00.0 0000:03:00.1"

if [ ! -z "$(ls -A /sys/class/iommu)" ]; then
    for DEV in $DEVS; do
        echo "vfio-pci" > /sys/bus/pci/devices/$DEV/driver_override
    done
fi

modprobe -i vfio-pci
```

## Script installation

Edit `/etc/mkinitcpio.conf` :

1. Add `modconf` to the **HOOKS** array and `/usr/local/bin/vfio-pci-override.sh` to the **FILES** array.

Edit `/etc/modprobe.d/vfio.conf` :

1. Add the following line: `install vfio-pci /usr/local/bin/vfio-pci-override.sh`
2. **Regenerate the `initramfs`** and reboot.

## Passing the boot GPU to the guest

The GPU marked as `boot_vga` is a special case when it comes to doing PCI passthroughs, since the BIOS needs to use it in order to display things like boot messages or the BIOS configuration menu. To do that, it makes [a copy of the VGA boot ROM which can then be freely modified \(https://www.redhat.com/archives/vfio-users/2016-May/msg00224.html\)](https://www.redhat.com/archives/vfio-users/2016-May/msg00224.html). This modified copy is the version the system gets to see, which the passthrough driver may reject as invalid. As such, it is generally recommended to change the boot GPU in the BIOS configuration so the host GPU is used instead or, if that is not possible, to swap the host and guest cards in the machine itself.

## Using Looking Glass to stream guest screen to the host

It is possible to make VM share the monitor, and optionally a keyboard and a mouse with a help of [Looking Glass \(http://looking-glass.hostfission.com/\)](http://looking-glass.hostfission.com/).

## Adding IVSHMEM Device to VM

Looking glass works by creating a shared memory buffer between a host and a guest. This is a lot faster than streaming frames via localhost, but requires additional setup.

With your VM turned off open the machine configuration

```
$ virsh edit [vmname]
...
<devices>
  ...
  <shmem name='looking-glass'>
    <model type='ivshmem-plain' />
    <size unit='M'>32</size>
  </shmem>
</devices>
...
```

You should replace 32 with your own calculated value based on what resolution you are going to pass through. It can be calculated like this:

```
width x height x 4 x 2 = total bytes
total bytes / 1024 / 1024 = total mebibytes + 2
```

For example, in case of 1920x1080

```
1920 x 1080 x 4 x 2 = 16,588,800 bytes
16,588,800 / 1024 / 1024 = 15.82 MiB + 2 = 17.82
```

The result must be **rounded up** to the nearest power of two, and since 17.82 is bigger than 16 we should choose 32.

Next create a configuration file to create the shared memory file on boot

```
/etc/tmpfiles.d/10-looking-glass.conf
-----
f      /dev/shm/looking-glass 0660  user  kvm  -
```

Replace user with your username.

Ask systemd-tmpfiles to create the shared memory file now without waiting to next boot

```
# systemd-tmpfiles --create /etc/tmpfiles.d/10-looking-glass.conf
```

## Installing the IVSHMEM Host to Windows guest

Currently Windows would not notify users about a new IVSHMEM device, it would silently install a dummy driver. To actually enable the device you have to go into device manager and update the driver for the device under the "System Devices" node for "PCI standard RAM Controller". Download the signed driver [from Red Hat \(https://fedorapeople.org/groups/virt/virtio-win/direct-downloads/upstream-virtio/\)](https://fedorapeople.org/groups/virt/virtio-win/direct-downloads/upstream-virtio/).

Once the driver is installed you must download a matching [looking-glass-host \(https://github.com/gnif/LookingGlass/releases\)](https://github.com/gnif/LookingGlass/releases) that matches the client you installed from AUR and start it on your guest. In order to run it you would also need to install Microsoft Visual C++ Redistributable from [Microsoft \(https://www.visualstudio.com/downloads/\)](https://www.visualstudio.com/downloads/) It is also

possible to make it start automatically on VM boot by editing the HKEY\_LOCAL\_MACHINE\SOFTWARE\Microsoft\Windows\CurrentVersion\Run registry and adding a path to the downloaded executable.

## Getting a client

Looking glass client can be installed from AUR using [looking-glass](https://aur.archlinux.org/packages/looking-glass/) (<https://aur.archlinux.org/packages/looking-glass/>)<sup>AUR</sup> or [looking-glass-git](https://aur.archlinux.org/packages/looking-glass-git/) (<https://aur.archlinux.org/packages/looking-glass-git/>)<sup>AUR</sup> packages.

You can start it once the VM is set up and running

```
$ looking-glass-client
```

If you do not want to use Spice to control the guest mouse and keyboard you can disable the Spice server.

```
$ looking-glass-client -s
```

Additionally you may want to start Looking Glass Client as full screen, otherwise the image may be scaled down resulting in poor image fidelity.

```
$ looking-glass-client -F
```

Launch with the `--help` option for further information.

## Swap peripherals to and from the Host

Looking Glass includes a Spice client in order to control mouse movement on the Windows guest. However this may have too much latency for certain applications, such as gaming. An alternative method is passing through specific USB devices for minimal latency. This allows for switching the devices between host and guest.

First create a .xml file for the device(s) you wish to pass-through, which libvirt will use to identify the device.

```
~/VFI0input/input_1.xml  
-----  
<hostdev mode='subsystem' type='usb' managed='no'>  
  <source>  
    <vendor id='0x[Before Colon]'/>  
    <product id='0x[After Colon]'/>  
  </source>  
</hostdev>
```

Replace [Before/After Colon] with the contents of the 'lsusb' command, specific to the device you want to pass-through.

For instance my mouse is `Bus 005 Device 002: ID 1532:0037 Razer USA, Ltd` so I would replace

vendor id with 1532, and product id with 1037.

Repeat this process for any additional USB devices you want to pass-through. If your mouse / keyboard has multiple entries in `lsusb`, perhaps if it is wireless, then create additional xml files for each.

**Note:** Do not forget to change the path & name of the script(s) above and below to match your user and specific system.

Next a bash script file is needed to tell libvirt what to attach/detach the USB devices to the guest.

```
~/VFI0input/input_attach.sh
-----
#!/bin/bash
virsh attach-device [VM-Name] [USBdevice]
```

Replace [VM-Name] with the name of your virtual machine, which can be seen under virt-manager. Additionally replace [USBdevice] with the **full** path to the .xml file for the device you wish to pass-through. Add additional lines for more than 1 device. For example here is my script:

```
~/VFI0input/input_attach.sh
-----
#!/bin/bash
virsh attach-device win10 /home/$USER/.VFI0input/input_mouse.xml
virsh attach-device win10 /home/$USER/.VFI0input/input_keyboard.xml
```

Next duplicate the script file and replace `attach-device` with `detach-device`. Ensure both scripts are executable with `chmod +x $script.sh`

This 2 script files can now be executed to attach or detach your USB devices from the host to the guest VM. It is important to note that they may need to be executed as root. To run the script from the Windows VM, one possibility is using [PuTTY](#) to [SSH](#) into the host, and execute the script. On Windows PuTTY comes with `plink.exe` which can execute singular commands over SSH before then logging out, instead of opening a SSH terminal, all in the background.

```
detach_devices.bat
-----
"C:\Program Files\PuTTY\plink.exe" root@$HOST_IP -pw $ROOTPASSWORD /home/$USER/.VFI0input/input_d
etach.sh
```

Replace `$HOST_IP` with the Host [IP Address](#) and `$ROOTPASSWORD` with the root password.

**Warning:** This method is insecure if somebody has access to your VM, since they could open the file and read your password. It is advisable to use [SSH keys](#) instead!

You may also want to execute the script files using key binds. On Windows one option is [Autohotkey](#) (<https://autohotke>

[y.com/](#)), and on the Host [Xbindkeys](#). Because of the need to run the scripts as root, you may also need to use [Polkit](#) or [Sudo](#) which can both be used to authenticate specific executables as able to run as root without needing a password.

## Bypassing the IOMMU groups (ACS override patch)

If you find your PCI devices grouped among others that you do not wish to pass through, you may be able to separate them using Alex Williamson's ACS override patch. Make sure you understand [the potential risk \(https://vfio.blogspot.com/2014/08/iommu-groups-inside-and-out.html\)](https://vfio.blogspot.com/2014/08/iommu-groups-inside-and-out.html) of doing so.

You will need a kernel with the patch applied. The easiest method to acquiring this is through the [linux-vfio \(http://s://aur.archlinux.org/packages/linux-vfio/\)](http://s://aur.archlinux.org/packages/linux-vfio/)<sup>AUR</sup> package.

In addition, the ACS override patch needs to be enabled with kernel command line options. The patch file adds the following documentation:

```
pcie_acs_override =
    [PCIE] Override missing PCIe ACS support for:
    downstream
        All downstream ports - full ACS capabilities
    multifunction
        All multifunction devices - multifunction ACS subset
    id:nnn:nnn
        Specific device - full ACS capabilities
        Specified as vid:did (vendor/device ID) in hex
```

The option `pcie_acs_override=downstream,multifunction` should break up as many devices as possible.

After installation and configuration, reconfigure your [bootloader kernel parameters](#) to load the new kernel with the `pcie_acs_override=` option enabled.

## Plain QEMU without libvirt

Instead of setting up a virtual machine with the help of libvirt, plain QEMU commands with custom parameters can be used for running the VM intended to be used with PCI passthrough. This is desirable for some use cases like scripted setups, where the flexibility for usage with other scripts is needed.

To achieve this after [#Setting up IOMMU](#) and [#Isolating the GPU](#), follow the [QEMU](#) article to setup the virtualized environment, [enable KVM](#) on it and use the flag `-device vfio-pci,host=07:00.0` replacing the identifier (07:00.0) with your actual device's ID that you used for the GPU isolation earlier.

For utilizing the OVMF firmware, make sure the [edk2-ovmf \(https://www.archlinux.org/packages/?name=edk2-ovmf\)](https://www.archlinux.org/packages/?name=edk2-ovmf) package is installed, copy the UEFI variables from `/usr/share/edk2-ovmf/x64/OVMF_VARS.fd` to temporary location like `/tmp/MY_VARS.fd` and finally specify the OVMF paths by appending the following parameters to the QEMU command (order matters):

- drive if=pflash,format=raw,readonly,file=/usr/share/edk2-ovmf/x64
- /OVMF\_CODE.fd
- for the actual OVMF firmware binary, note the readonly option
- -drive if=pflash,format=raw,file=/tmp/MY\_VARS.fd for the variables

**Note:** QEMU's default SeaBIOS can be used instead of OVMF, but it is not recommended as it can cause issues with passthrough setups.

It is recommended to study the QEMU article for ways to enhance the performance by using the [virtio drivers](#) and other further configurations for the setup.

You also might have to use the `-cpu host,kvm=off` parameter to forward the host's CPU model info to the VM and fool the virtualization detection used by Nvidia's and possibly other manufacturers' device drivers trying to block the full hardware usage inside a virtualized system.

## Passing through other devices

---

### USB controller

If your motherboard has multiple USB controllers mapped to multiple groups, it is possible to pass those instead of USB devices. Passing an actual controller over an individual USB device provides the following advantages :

- If a device disconnects or changes ID over the course of an given operation (such as a phone undergoing an update), the VM will not suddenly stop seeing it.
- Any USB port managed by this controller is directly handled by the VM and can have its devices unplugged, replugged and changed without having to notify the hypervisor.
- Libvirt will not complain if one of the USB devices you usually pass to the guest is missing when starting the VM.

Unlike with GPUs, drivers for most USB controllers do not require any specific configuration to work on a VM and control can normally be passed back and forth between the host and guest systems with no side effects.

**Warning:** Make sure your USB controller supports resetting: [#Passing through a device that does not support resetting](#)

You can find out which PCI devices correspond to which controller and how various ports and devices are assigned to each one of them using this command :

```
$ for usb_ctrl in /sys/bus/pci/devices/*/usb*; do pci_path=${usb_ctrl%/*}; iommu_group=$(readlink $pci_path/iommu_group); echo "Bus $(cat $usb_ctrl/busnum) --> ${pci_path###*/} (IOMMU group ${iommu_group###*/})"; lsusb -s ${usb_ctrl##*/usb};; echo; done
```

```
Bus 1 --> 0000:00:1a.0 (IOMMU group 4)
Bus 001 Device 004: ID 04f2:b217 Chicony Electronics Co., Ltd Lenovo Integrated Camera (0.3MP)
Bus 001 Device 007: ID 0a5c:21e6 Broadcom Corp. BCM20702 Bluetooth 4.0 [ThinkPad]
Bus 001 Device 008: ID 0781:5530 SanDisk Corp. Cruzer
Bus 001 Device 002: ID 8087:0024 Intel Corp. Integrated Rate Matching Hub
```

```

Bus 001 Device 001: ID 1d6b:0002 Linux Foundation 2.0 root hub
Bus 2 --> 0000:00:1d.0 (IOMMU group 9)
Bus 002 Device 006: ID 0451:e012 Texas Instruments, Inc. TI-Nspire Calculator
Bus 002 Device 002: ID 8087:0024 Intel Corp. Integrated Rate Matching Hub
Bus 002 Device 001: ID 1d6b:0002 Linux Foundation 2.0 root hub

```

This laptop has 3 USB ports managed by 2 USB controllers, each with their own IOMMU group. In this example, Bus 001 manages a single USB port (with a SanDisk USB pendrive plugged into it so it appears on the list), but also a number of internal devices, such as the internal webcam and the bluetooth card. Bus 002, on the other hand, does not appear to manage anything except for the calculator that is plugged into it. The third port is empty, which is why it does not show up on the list, but is actually managed by Bus 002.

Once you have identified which controller manages which ports by plugging various devices into them and decided which one you want to passthrough, simply add it to the list of PCI host devices controlled by the VM in your guest configuration. No other configuration should be needed.

**Note:** If your USB controller does not support resetting, is not in an isolated group, or is otherwise unable to be passed through then it may still be possible to accomplish similar results through [udev](#) rules. See [\[1\] \(https://github.com/olavmrk/usb-libvirt-hotplug\)](#) which allows any device connected to specified USB ports to be automatically attached to a virtual machine.

## Passing VM audio to host via PulseAudio

It is possible to route the virtual machine's audio to the host as an application using libvirt. This has the advantage of multiple audio streams being routable to one host output, and working with audio output devices that do not support passthrough. [PulseAudio](#) is required for this to work.

First, remove the comment from the `#user = ""` line. Then add your username in the quotations. This tells QEMU which user's pulseaudio stream to route through.

```

/etc/libvirt/qemu.conf
-----
user = "example"

```

Next, modify the libvirt configuration

```

$ virsh edit [vmname]
-----
<domain type='kvm'>

```

to

```

$ virsh edit [vmname]

```

```
<domain type='kvm' xmlns:qemu='http://libvirt.org/schemas/domain/qemu/1.0'>
```

Then set the QEMU PulseAudio variables at the bottom of the libvirt xml file

```
$ virsh edit [vmname]
```

```
</devices>
</domain>
```

to

```
$ virsh edit [vmname]
```

```
</devices>
  <qemu:commandline>
    <qemu:arg value="-audiodev"/>
    <qemu:arg value="pa,id=snd0,server=/run/user/1000/pulse/native"/>
  </qemu:commandline>
</domain>
```

You can also use the /tmp directory if you have multiple users accessing PulseAudio

```
$ virsh edit [vmname]
```

```
</devices>
  <qemu:commandline>
    <qemu:arg value="-audiodev"/>
    <qemu:arg value="pa,id=snd0,server=unix:/tmp/pulse-socket"/>
  </qemu:commandline>
</domain>
```

You may want to change 1000 under the user directory to your current user uid (which can be found by running the `id` command. Remember to save the file and exit it without ending the process before continuing, otherwise the changes will not register. If you get the message `Domain [vmname] XML configuration edited.` after exiting, it means that your changes have been applied. Virtual Machine audio will now be routed through the host as an application. A application such as [pavucontrol](https://www.archlinux.org/packages/?name=pavucontrol) (<https://www.archlinux.org/packages/?name=pavucontrol>) can be used to control the output device.

### QEMU 4.0/4.2+ audio changes

As of QEMU 4.0 and made official with [QEMU 4.2](https://www.kraxel.org/blog/2020/01/qemu-sound-audiodev/) (<https://www.kraxel.org/blog/2020/01/qemu-sound-audiodev/>) the `-audiodev` parameter is now to be used with PulseAudio sound passthrough, any previous environmental variables you have defined in your XML **need** to be removed e.g. `qemu:env` to properly use this.

You will also need to change the chipset accordingly with each QEMU update to how your VM is set up, i.e. `pc-q35-4.2` or `pc-i440fx-4.2` to advantage of the changes in QEMU:

```
$ virsh edit [vmname]

<os>
  <type arch='x86_64' machine='pc-q35-4.2'>hvm</type>
  ...
</os>
...
<os>
  <type arch='x86_64' machine='pc-i440fx-4.2'>hvm</type>
  ...
</os>
```

The latest audiodev changes in QEMU also suggest using the ICH9 instead of ICH6 or AC97 for audio emulation to take advantage of this newer code, it is recommended to use the parameters as shown below because Libvirt currently does not define the `id=` parameter correctly:

```
$ virsh edit [vmname]

<qemu:arg value="-device"/>
<qemu:arg value="ich9-intel-hda,bus=pcie.0,addr=0x1b"/>
<qemu:arg value="-device"/>
<qemu:arg value="hda-micro,audiodev=hda"/>
<qemu:arg value="-audiodev"/>
<qemu:arg value="pa,id=hda,server=unix:/tmp/pulse-socket"/>
```

QEMU 4.2 also has 5.1 channel audio support via `usb-audio` emulation, this can be enabled as shown below:

```
$ virsh edit [vmname]

<qemu:arg value="-device"/>
<qemu:arg value="usb-audio,audiodev=usb,multi=on"/>
<qemu:arg value="-audiodev"/>
<qemu:arg value="pa,id=usb,server=unix:/tmp/pulse-socket,out.mixing-engine=off"/>
```

### Note:

- You can have multiple audio backends, by simply specifying `-audiodev` multiple times in your XML and by assigning them different ids. This can be useful for a use case of having two identical backends. With PulseAudio each backend is a separate stream and can be routed to different output devices on the host (using a pulse mixer app like [pavucontrol](https://www.archlinux.org/packages/?name=pavucontrol) (<https://www.archlinux.org/packages/?name=pavucontrol>) or [pulsemixer](http://www.archlinux.org/packages/?name=pulsemixer) (<http://www.archlinux.org/packages/?name=pulsemixer>)).
- USB 3 emulation is needed in Libvirt/QEMU to enable the `usb-audio` parameter
- It is recommended to enable MSI interrupts with a tool such as [this \(https://github.com/CHEF-KOCH/MSI-utility\)](https://github.com/CHEF-KOCH/MSI-utility) on the ICH9 audio device to mitigate any crackling, stuttering, speedup, or no audio at all after VM restart
- The `-audiodev=hda` parameter can be changed to your liking, but as shown above, the names must match each other or the VM will not start
- If audio is crackling/stuttering/speedup etc. is still present you may want to adjust parameters such as `buffer-length` and `timer-period`, more information on these

parameters and more can be found in the [qemu\(1\) \(https://jlk.fjfi.cvut.cz/arch/ma npages/man/qemu.1\)](https://jlk.fjfi.cvut.cz/arch/ma npages/man/qemu.1) manual

- Some audio chipsets such as [Realtek alc1220 \(https://bugzilla.kernel.org/show\\_bug.cgi?id=195303\)](https://bugzilla.kernel.org/show_bug.cgi?id=195303) may also have issues out of the box so do consider this when using any audio emulation with QEMU
- Improper pinning or heavy host usage without using [isolcpus](#) can also influence sound bugs, especially while gaming in a VM

## Passing VM audio to host via Scream

It is possible to pass VM audio through a bridged network such as the one provided by Libvirt or by adding a IVSHMEM device to the host by using a application called [Scream \(https://github.com/duncanthrax/scream\)](https://github.com/duncanthrax/scream). This section will only cover using PulseAudio as a receiver on the host. See the project page for more details and instructions on other methods.

### Using Scream with a bridged network

#### Note:

- This is the *preferred* way to use this, although results may vary per user
- It is recommend to use the [#Virtio network](#) adapter while using Scream, other virtual adapters provided by QEMU such as **e1000e** may lead to poor performance

To use scream via your network you will want to find your bridge name via `ip -a`, in most cases it will be called **br0** or **virbr0**. Below is a example of the command needed to start the Scream application:

```
$ scream-pulse -i virbr0 &
```

**Warning:** This will not work with a **macvtap bridge** as that does not allow host to guest communication, also make sure you have the proper firewall ports open for it to communicate with the VM

### Adding the IVSHMEM device to use Scream with IVSHMEM

With the VM turned off, edit the machine configuration

```
$ virsh edit [vmname]

...
<devices>
  ...
  <shmem name='scream-ivshmem'>
    <model type='ivshmem-plain' />
    <size unit='M'>2</size>
  </shmem>
</devices>
```

```
...
```

In the above configuration, the size of the IVSHMEM device is 2MB (the recommended amount). Change this as needed.

Now refer to [#Adding IVSHMEM Device to VM](#) to configure the host to create the shared memory file on boot, replacing `looking-glass` with `scream-ivshmem`.

## Configuring the Windows guest for IVSHMEM

The correct driver must be installed for the IVSHMEM device on the guest. See [#Installing the IVSHMEM Host to Windows guest](#). Ignore the part about `looking-glass-host`.

Install the [Scream \(https://github.com/duncanthrax/scream/releases\)](https://github.com/duncanthrax/scream/releases) virtual audio driver on the guest. If you have secure boot enabled for your VM, you may need to disable it.

Using the registry editor, set the DWORD `HKLM\SYSTEM\CurrentControlSet\Services\Scream\Options\UseIVSHMEM` to the size of the IVSHMEM device in MB. Note that scream identifies its IVSHMEM device using its size, so make sure there is only one device of that size.

## Configuring the host

Install [scream \(https://aur.archlinux.org/packages/scream/\)](https://aur.archlinux.org/packages/scream/)<sup>AUR</sup>.

Create the systemd user service file to control the reciever

```
~/.config/systemd/user/scream-ivshmem-pulse.service

[Unit]
Description=Scream IVSHMEM pulse reciever
After=pulseaudio.service
Wants=pulseaudio.service

[Service]
Type=simple
ExecStartPre=/usr/bin/truncate -s 0 /dev/shm/scream-ivshmem
ExecStartPre=/usr/bin/dd if=/dev/zero of=/dev/shm/scream-ivshmem bs=1M count=2
ExecStart=/usr/bin/scream -m /dev/shm/scream-ivshmem

[Install]
WantedBy=default.target
```

Edit `count=2` with the size of the IVSHMEM device in MB.

Now start the service with

```
$ systemctl start --user scream-ivshmem-pulse
```

To have it automatically start on next login, enable the service

```
$ systemctl enable --user scream-ivshmem-pulse
```

## Physical disk/partition

Raw and qcow2 especially can have noticeable overhead for heavy IO. A whole disk or a partition may be used directly to bypass the filesystem and improve I/O performance. If you wish to dual boot the guest OS natively you would need to pass the entire disk without any partitioning. It is suggested to use `/dev/disk/by-` paths to refer to the disk since `/dev/sdX` entries can change between boots. To find out which disk/partition is associated with the one you would like to pass:

```
$ ls -l /dev/disk/by-id
-----
ata-ST1000LM002-9VQ14L_Z0501SZ9 -> ../../sdd
```

See [#Virtio disk](#) on how to add these with libvirt XML. You can also add the disk with Virt-Manager's **Add Hardware** menu and then type the disk you want in the **Select or create custom storage** box, e.g. `/dev/disk/by-id/ata-ST1000LM002-9VQ14L_Z0501SZ9`

## Gotchas

### Passing through a device that does not support resetting

When the VM shuts down, all devices used by the guest are deinitialized by its OS in preparation for shutdown. In this state, those devices are no longer functional and must then be power-cycled before they can resume normal operation. Linux can handle this power-cycling on its own, but when a device has no known reset methods, it remains in this disabled state and becomes unavailable. Since Libvirt and Qemu both expect all host PCI devices to be ready to reattach to the host before completely stopping the VM, when encountering a device that will not reset, they will hang in a "Shutting down" state where they will not be able to be restarted until the host system has been rebooted. It is therefore recommended to only pass through PCI devices which the kernel is able to reset, as evidenced by the presence of a `reset` file in the PCI device sysfs node, such as `/sys/bus/pci/devices/0000:00:1a.0/reset`.

The following bash command shows which devices can and cannot be reset.

```
for iommu_group in $(find /sys/kernel/iommu_groups/ -maxdepth 1 -mindepth 1 -type d);do echo "IOMMU group $(basename "$iommu_group")"; for device in $(\ls -l "$iommu_group"/devices/); do if [[ -e "$iommu_group"/devices/"$device"/reset ]]; then echo -n "[RESET]"; fi; echo -n '$\t';lspci -nns "$device"; done; done
IOMMU group 0
    00:00.0 Host bridge [0600]: Intel Corporation Xeon E3-1200 v2/Ivy Bridge DRAM Controller [8086:0158] (rev 09)
IOMMU group 1
    00:01.0 PCI bridge [0604]: Intel Corporation Xeon E3-1200 v2/3rd Gen Core processor PCI E
```

```
xpress Root Port [8086:0151] (rev 09)
  01:00.0 VGA compatible controller [0300]: NVIDIA Corporation GK208 [GeForce GT 720] [10de:1288] (rev a1)
  01:00.1 Audio device [0403]: NVIDIA Corporation GK208 HDMI/DP Audio Controller [10de:0e0f] (rev a1)
IOMMU group 2
  00:14.0 USB controller [0c03]: Intel Corporation 7 Series/C210 Series Chipset Family USB xHCI Host Controller [8086:1e31] (rev 04)
IOMMU group 4
[RESET] 00:1a.0 USB controller [0c03]: Intel Corporation 7 Series/C210 Series Chipset Family USB Enhanced Host Controller #2 [8086:1e2d] (rev 04)
IOMMU group 5
[RESET] 00:1b.0 Audio device [0403]: Intel Corporation 7 Series/C210 Series Chipset Family High Definition Audio Controller [8086:1e20] (rev 04)
IOMMU group 10
[RESET] 00:1d.0 USB controller [0c03]: Intel Corporation 7 Series/C210 Series Chipset Family USB Enhanced Host Controller #1 [8086:1e26] (rev 04)
IOMMU group 13
  06:00.0 VGA compatible controller [0300]: NVIDIA Corporation GM204 [GeForce GTX 970] [10de:13c2] (rev a1)
  06:00.1 Audio device [0403]: NVIDIA Corporation GM204 High Definition Audio Controller [10de:0fbb] (rev a1)
```

This signals that the xHCI USB controller in 00:14.0 cannot be reset and will therefore stop the VM from shutting down properly, while the integrated sound card in 00:1b.0 and the other two controllers in 00:1a.0 and 00:1d.0 do not share this problem and can be passed without issue.

## Complete setups and examples

---

For many reasons users may seek to see [complete passthrough setup examples](#).

These examples offer a supplement to existing hardware compatibility lists. Additionally, if you have trouble configuring a certain mechanism in your setup, you might find these examples very valuable. Users there have described their setups in detail, and some have provided examples of their configuration files as well.

We encourage those who successfully build their system from this resource to help improve it by contributing their builds. Due to the many different hardware manufacturers involved, the seemingly significant lack of sufficient documentation, as well as other issues due to the nature of this process, community contributions are necessary.

## Troubleshooting

---

If your issue is not mentioned below, you may want to browse [QEMU#Troubleshooting](#).

### QEMU 4.0: Unable to load graphics drivers/BSOD/Graphics stutter after driver install using Q35

Starting with QEMU 4.0, the Q35 machine type changes the default `kernel_irqchip` from `off` to `split` which breaks some guest devices, such as nVidia graphics (the driver fails to load / black screen / code 43 / graphics stutters, usually when mouse moving). Switch to full KVM mode instead by adding `<ioapic driver='kvm'/>` under libvirt's `<features>` tag in your VM configuration or by adding `kernel_irqchip=on` in the `-machine`

QEMU arg.

## QEMU 5.0: host-passthrough with newer kernel than 5.4 when using Zen2 processors: Windows 10 BSOD loop: KERNEL SECURITY CHECK FAILURE

Starting with QEMU 5.0 and added features to zen2 on newer kernels than 5.4 will cause a BSOD loop of: **KERNEL SECURITY CHECK FAILURE**. The cause is missing *amd-stibp* flag that QEMU does not support in version 5.0. Recommended fix is to use the cpu model "host-model" because manually editing host-passthrough to work is unstable and should be used just for benchmarking or testing purposes. To get host-passthrough to work you need to have the following inside your VM xml configuration file:

```
<domain type='kvm' xmlns:qemu='http://libvirt.org/schemas/domain/qemu/1.0'>
...
<qemu:commandline>
  <qemu:arg value='-cpu' />
  <qemu:arg value='host,topoext=on,hv-time,hv-relaxed,hv-vapic,hv-spinlocks=0x1fff,kvm=off,host-cache-info=on,l3-cache=off,-amd-stibp,hv-vendor-id=driverError43Bypass' />
</qemu:commandline>
```

You can get cpu arg values for your zen2 cpu with the command `grep "\-cpu" /var/log/libvirt/qemu/win10.log` assuming your virtual machine is named win10.

## "Error 43: Driver failed to load" on Nvidia GPUs passed to Windows VMs

### Note:

- This may also fix SYSTEM\_THREAD\_EXCEPTION\_NOT\_HANDLED boot crashes related to Nvidia drivers.
- This may also fix problems under linux guests.

Since version 337.88, Nvidia drivers on Windows check if an hypervisor is running and fail if it detects one, which results in an Error 43 in the Windows device manager. Starting with QEMU 2.5.0 and libvirt 1.3.3, the `vendor_id` for the hypervisor can be spoofed, which is enough to fool the Nvidia drivers into loading anyway. All one must do is add `hv_vendor_id=whatever` to the hypervisor parameters in their QEMU command line, or by adding the following line to their libvirt domain configuration. The `vendor_id` can be [any string value up to 12 characters \(https://libvirt.org/formatdomain.html#elementsFeatures\)](https://libvirt.org/formatdomain.html#elementsFeatures).

```
$ virsh edit [vmname]
```

```
...
<features>
  <hyperv>
```

```

    ...
    <vendor_id state='on' value='whatever' />
    ...
</hyperv>
...
<kvm>
  <hidden state='on' />
</kvm>
</features>
...

```

Users with older versions of QEMU and/or libvirt will instead have to disable a few hypervisor extensions, which can degrade performance substantially. If this is what you want to do, do the following replacement in your libvirt domain config file.

```

$ virsh edit [vmname]
-----
...
<features>
  <hyperv>
    <relaxed state='on' />
    <vapic state='on' />
    <spinlocks state='on' retries='8191' />
  </hyperv>
  ...
</features>
...
<clock offset='localtime'>
  <timer name='hypervclock' present='yes' />
</clock>
...

```

```

...
<clock offset='localtime'>
  <timer name='hypervclock' present='no' />
</clock>
...
<features>
  <kvm>
    <hidden state='on' />
  </kvm>
  ...
  <hyperv>
    <relaxed state='off' />
    <vapic state='off' />
    <spinlocks state='off' />
  </hyperv>
  ...
</features>
...

```

## "Error 43: Driver failed to load" with mobile (Optimus/max-q) nvidia GPUs

This error occurs because the Nvidia driver wants to check the status of the power supply. If no battery is present, the driver does not work. Whether Libvirt or Quemu, by default none of them provide the possibility to simulate a battery.

You can however create and add a custom acpi table file to the virtual machine which will do the work.

First you have to create the custom acpi table file by pasting the following base64 string [here \(https://base64.guru/converter/decode/file\)](https://base64.guru/converter/decode/file) and save the result file as SSDT1.dat:

```
U1NEVKAAAAAB9EJPQ0hTAEJYUENTU0RUAQAAAELOVEwYEBkgoA8AFVwuX1NCX1BDSTAGABMBi5f
U0JfUENJMFuCTwVCQVQwCF9ISUQMqDAMCghfVU\EABQJX1NUQQckCh8UK19CSUYApBIjDQELcBcL
cBcBC9A5C1gCCyWBCjwKPA0ADQANTElPTgANABQsX0JTVACKEgoEAAALcBcL0DK=
```

Next you must add the processed file to the virtual machine:

```
<qemu:commandline>
  <qemu:arg value="-acpitable"/>
  <qemu:arg value="file=/path/to/your/SSDT1.dat"/>
</qemu:commandline>
```

Source ([https://www.reddit.com/r/VFIO/comments/ebo2uk/nvidia\\_geforce\\_rtx\\_2060\\_mobile\\_success\\_qemu\\_ovmf/](https://www.reddit.com/r/VFIO/comments/ebo2uk/nvidia_geforce_rtx_2060_mobile_success_qemu_ovmf/))

## "BAR 3: cannot reserve [mem]" error in dmesg after starting VM

With respect to [this article \(https://www.linuxquestions.org/questions/linux-kernel-70/kernel-fails-to-assign-memory-to-pcie-device-4175487043/\)](https://www.linuxquestions.org/questions/linux-kernel-70/kernel-fails-to-assign-memory-to-pcie-device-4175487043/):

If you still have code 43 check dmesg for memory reservation errors after starting up VM, if you have similar it could be the case:

```
vfio-pci 0000:09:00.0: BAR 3: cannot reserve [mem 0xf0000000-0xf1ffffff 64bit pref]
```

Find out a PCI Bridge your graphic card is connected to. This will give actual hierarchy of devices:

```
$ lspci -t
```

Before starting VM run following lines replacing IDs with actual from previous output.

```
# echo 1 > /sys/bus/pci/devices/0000\:00\:03.1/remove
# echo 1 > /sys/bus/pci/rescan
```

**Note:** Probably setting [kernel parameter](https://proxmox.com/wiki/Pci_passthrough#BAR_3:_can.27t_reserve_.5Bmem.5D_error) `video=efifb:off` is required as well. Source ([https://proxmox.com/wiki/Pci\\_passthrough#BAR\\_3:\\_can.27t\\_reserve\\_.5Bmem.5D\\_error](https://proxmox.com/wiki/Pci_passthrough#BAR_3:_can.27t_reserve_.5Bmem.5D_error))

In addition try adding kernel parameter `pci=realloc` which also [helps with hotplugging issues \(https://github.com/Dunedan/mbp-2016-linux/issues/60#issuecomment-396311301\)](https://github.com/Dunedan/mbp-2016-linux/issues/60#issuecomment-396311301).

## UEFI (OVMF) compatibility in VBIOS

With respect to [this article \(https://pve.proxmox.com/wiki/Pci\\_passthrough#How\\_to\\_known\\_if\\_card\\_is\\_UEFI\\_280vmf.29\\_compatible\)](https://pve.proxmox.com/wiki/Pci_passthrough#How_to_known_if_card_is_UEFI_280vmf.29_compatible):

Error 43 can be caused by the GPU's VBIOS without UEFI support. To check whenever your VBIOS supports it, you will have to use `rom-parser` :

```
$ git clone https://github.com/awilliam/rom-parser
$ cd rom-parser && make
```

Dump the GPU VBIOS:

```
# echo 1 > /sys/bus/pci/devices/0000:01:00.0/rom
# cat /sys/bus/pci/devices/0000:01:00.0/rom > /tmp/image.rom
# echo 0 > /sys/bus/pci/devices/0000:01:00.0/rom
```

And test it for compatibility:

```
$ ./rom-parser /tmp/image.rom
-----
Valid ROM signature found @600h, PCIR offset 190h
  PCIR: type 0 (x86 PC-AT), vendor: 10de, device: 1184, class: 030000
  PCIR: revision 0, vendor revision: 1
Valid ROM signature found @fa00h, PCIR offset 1ch
  PCIR: type 3 (EFI), vendor: 10de, device: 1184, class: 030000
  PCIR: revision 3, vendor revision: 0
    EFI: Signature Valid, Subsystem: Boot, Machine: X64
Last image
```

To be UEFI compatible, you need a "type 3 (EFI)" in the result. If it is not there, try updating your GPU VBIOS. GPU manufacturers often share VBIOS upgrades on their support pages. A large database of known compatible and working VBIOSes (along with their UEFI compatibility status!) is available on [TechPowerUp \(https://www.techpowerup.com/vgabios/\)](https://www.techpowerup.com/vgabios/).

Updated VBIOS can be used in the VM without flashing. To load it in QEMU:

```
-device vfio-pci,host=07:00.0,.....,romfile=/path/to/your/gpu/bios.bin \
```

And in libvirt:

```
<hostdev>
  ...
  <rom file='/path/to/your/gpu/bios.bin' />
  ...
</hostdev>
```

One should compare VBIOS versions between host and guest systems using [nvflash \(https://www.techpowerup.com/do\\_nvflash\)](https://www.techpowerup.com/do_nvflash)

[wnload/nvidia-nvflash/](#) (Linux versions under *Show more versions*) or [GPU-Z \(https://www.techpowerup.com/download/techpowerup-gpu-z/\)](#) (in Windows guest). To check the currently loaded VBIOS:

```
$ ./nvflash --version
...
Version           : 80.04.XX.00.97
...
UEFI Support      : No
UEFI Version      : N/A
UEFI Variant Id   : N/A ( Unknown )
UEFI Signer(s)    : Unsigned
...
```

And to check a given VBIOS file:

```
$ ./nvflash --version NV299MH.rom
...
Version           : 80.04.XX.00.95
...
UEFI Support      : Yes
UEFI Version      : 0x10022 (Jul  2 2013 @ 16377903 )
UEFI Variant Id   : 0x0000000000000004 ( GK1xx )
UEFI Signer(s)    : Microsoft Corporation UEFI CA 2011
...
```

If the external ROM did not work as it should in the guest, you will have to flash the newer VBIOS image to the GPU. In some cases it is possible to create your own VBIOS image with UEFI support using [GOPUpd \(https://www.win-raid.com/t892f16-AMD-and-Nvidia-GOP-update-No-requests-DIY.html\)](https://www.win-raid.com/t892f16-AMD-and-Nvidia-GOP-update-No-requests-DIY.html) tool, however this is risky and may result in GPU brick.

**Warning:** Failure during flashing may "brick" your GPU - recovery may be possible, but rarely easy and often requires additional hardware. **DO NOT** flash VBIOS images for other GPU models (different boards may use different VBIOSes, clocks, fan configuration). If it breaks, you get to keep all the pieces.

In order to avoid the irreparable damage to your graphics adapter it is necessary to unload the NVIDIA kernel driver first:

```
# modprobe -r nvidia_modeset nvidia
```

Flashing the VBIOS can be done with:

```
# ./nvflash romfile.bin
```

**Warning: DO NOT** interrupt the flashing process, even if it looks like it is stuck. Flashing should take about a minute on most GPUs, but may take longer.

## Slowed down audio pumped through HDMI on the video card

For some users VM's audio slows down/starts stuttering/becomes demonic after a while when it is pumped through HDMI on the video card. This usually also slows down graphics. A possible solution consists of enabling MSI (Message Signaled-Based Interrupts) instead of the default (Line-Based Interrupts).

In order to check whether MSI is supported or enabled, run the following command as root:

```
# lspci -vs $device | grep 'MSI:'
```

where ``$device`` is the card's address (e.g. ``01:00.0``).

The output should be similar to:

```
Capabilities: [60] MSI: Enable- Count=1/1 Maskable- 64bit+
```

A `-` after `Enable` means MSI is supported, but not used by the VM, while a `+` says that the VM is using it.

The procedure to enable it is quite complex, instructions and an overview of the setting can be found [here \(https://forums.guru3d.com/showthread.php?t=378044\)](https://forums.guru3d.com/showthread.php?t=378044).

On a linux guest you can use `modinfo` to see if there is option to enable MSI (for example: `"modinfo snd_hda_intel |grep msi"`). If there is, one can enable it by adding the relevant option to a custom `omodprobe` file - in `"/etc/modprobe.d/snd-hda-intel.conf"` inserting `"options snd-hda-intel enable_msi=1"`

Other hints can be found on the [lime-technology's wiki \(https://lime-technology.com/wiki/index.php/UnRAID\\_6/VM\\_Guest\\_Support#Enable\\_MSI\\_for\\_Interrupts\\_to\\_Fix\\_HDMI\\_Audio\\_Support\)](https://lime-technology.com/wiki/index.php/UnRAID_6/VM_Guest_Support#Enable_MSI_for_Interrupts_to_Fix_HDMI_Audio_Support), or on this article on [VFIO tips and tricks \(https://vfio.blogspot.it/2014/09/vfio-interrupts-and-how-to-coax-windows.html\)](https://vfio.blogspot.it/2014/09/vfio-interrupts-and-how-to-coax-windows.html).

A UI tool called [MSI Utility \(FOSS Version 2\) \(https://github.com/CHEF-KOCH/MSI-utility\)](https://github.com/CHEF-KOCH/MSI-utility) works with Windows 10 64-bit and simplifies the process.

In order to fix the issues enabling MSI on the `0` function of a nVidia card (`01:00.0 VGA compatible controller: NVIDIA Corporation GM206 [GeForce GTX 960] (rev a1) (prog-if 00 [VGA controller])`) was not enough; it will also be required to enable it on the other function (`01:00.1 Audio device: NVIDIA Corporation Device 0fba (rev a1)`) to fix the issue.

## No HDMI audio output on host when intel\_iommu is enabled

If after enabling `intel_iommu` the HDMI output device of Intel GPU becomes unusable on the host then setting the option `igfx_off` (i.e. `intel_iommu=on,igfx_off`) might bring the audio back, please read [intel-iommu.html \(https://www.kernel.org/doc/html/latest/x86/intel-iommu.html#graphics-problems\)](https://www.kernel.org/doc/html/latest/x86/intel-iommu.html#graphics-problems) for details about setting `igfx_off`.

## X does not start after enabling vfio\_pci

This is related to the host GPU being detected as a secondary GPU, which causes X to fail/crash when it tries to load a driver for the guest GPU. To circumvent this, a Xorg configuration file specifying the BusID for the host GPU is required. The correct BusID can be acquired from lspci or the Xorg log. [Source → \(https://www.redhat.com/archives/vfio-users/2016-August/msg00025.html\)](https://www.redhat.com/archives/vfio-users/2016-August/msg00025.html)

```
/etc/X11/xorg.conf.d/10-intel.conf
```

```
Section "Device"
    Identifier "Intel GPU"
    Driver "modesetting"
    BusID "PCI:0:2:0"
EndSection
```

## Chromium ignores integrated graphics for acceleration

Chromium and friends will try to detect as many GPUs as they can in the system and pick which one is preferred (usually discrete NVIDIA/AMD graphics). It tries to pick a GPU by looking at PCI devices, not OpenGL renderers available in the system - the result is that Chromium may ignore the integrated GPU available for rendering and try to use the dedicated GPU bound to the `vfio-pci` driver, and unusable on the host system, regardless of whenever a guest VM is running or not. This results in software rendering being used (leading to higher CPU load, which may also result in choppy video playback, scrolling and general un-smoothness).

This can be fixed by [explicitly telling Chromium which GPU you want to use](#).

## VM only uses one core

For some users, even if IOMMU is enabled and the core count is set to more than 1, the VM still only uses one CPU core and thread. To solve this enable "Manually set CPU topology" in `virt-manager` and set it to the desirable amount of CPU sockets, cores and threads. Keep in mind that "Threads" refers to the thread count per CPU, not the total count.

## Passthrough seems to work but no output is displayed

Make sure if you are using `virt-manager` that UEFI firmware is selected for your virtual machine. Also, make sure you have passed the correct device to the VM.

## Host lockup after VM shutdown

This issue seems to primarily affect users running a Windows 10 guest and usually after the VM has been run for a prolonged period of time: the host will experience multiple CPU core lockups (see [\[2\] \(https://bbs.archlinux.org/viewtopic.php?id=206050&p=2\)](https://bbs.archlinux.org/viewtopic.php?id=206050&p=2)). To fix this try enabling Message Signal Interrupts on the GPU passed through to the guest. A good guide for how to do this can be found in [\[3\] \(https://forums.guru3d.com/threads/windows-line-based-vs-message-sigaled-based-interrupts.378044/\)](https://forums.guru3d.com/threads/windows-line-based-vs-message-sigaled-based-interrupts.378044/). You can also download this application for windows here [\[4\] \(https://github.com/\)](https://github.com/)

[TechtonicSoftware/MSIInterruptEnabler](#)) that should make the process easier.

## Host lockup if guest is left running during sleep

VFIO-enabled virtual machines tend to become unstable if left running through a sleep/wakeup cycle and have been known to cause the host machine to lockup when an attempt is then made to shut them down. In order to avoid this, one can simply prevent the host from going into sleep while the guest is running using the following libvirt hook script and systemd unit. The hook file needs executable permissions to work.

```
/etc/libvirt/hooks/qemu
```

```
#!/bin/bash

OBJECT="$1"
OPERATION="$2"
SUBOPERATION="$3"
EXTRA_ARG="$4"

case "$OPERATION" in
    "prepare")
        systemctl start libvirt-nosleep@"$OBJECT"
        ;;
    "release")
        systemctl stop libvirt-nosleep@"$OBJECT"
        ;;
esac
```

```
/etc/systemd/system/libvirt-nosleep@.service
```

```
[Unit]
Description=Preventing sleep while libvirt domain "%i" is running

[Service]
Type=simple
ExecStart=/usr/bin/systemd-inhibit --what=sleep --why="Libvirt domain \"%i\" is running" --who=%U
--mode=block sleep infinity
```

## Cannot boot after upgrading ovmf

If you cannot boot after upgrading from [ovmf](https://www.archlinux.org/packages/?name=ovmf) (<https://www.archlinux.org/packages/?name=ovmf>) [\[broken link\]](#): replaced by [edk2-ovmf](https://www.archlinux.org/packages/?name=edk2-ovmf) (<https://www.archlinux.org/packages/?name=edk2-ovmf>) version 1:r23112.018432f0ce-1 then you need to remove the old `*VARS.fd` file in `/var/lib/libvirt/qemu/nvram/` :

```
# mv /var/lib/libvirt/qemu/nvram/vmname_VARS.fd /var/lib/libvirt/qemu/nvram/vmname_VARS.fd.old
```

See [FS#57825](https://bugs.archlinux.org/task/57825) (<https://bugs.archlinux.org/task/57825>) for further details.

## QEMU via cli pulseaudio stuttering/delay

Using following flags for the audio device and chipset might help if you are running into the stuttering/delay audio issues when running QEMU via cli:

```
qemu-system-x86_64 \  
-machine pc-i440fx-3.0 \  
-device hda-micro \  
-soundhw hda \  
-...
```

As noted in [QEMU 3.0 audio changes](#)<sup>[[broken link](#): invalid section]</sup> the specified chipset will include a series of audio patches.

Setting `QEMU_AUDIO_TIMER_PERIOD` to values higher than 100 might also help (did not test value lower than 100).

## Bluescreen at boot since Windows 10 1803

Since Windows 10 1803 there is a problem when you are using "host-passthrough" as cpu model. The machine cannot boot and is either boot looping or you get a bluescreen. You can workaround this by:

```
# echo 1 > /sys/module/kvm/parameters/ignore_msrs
```

To make it permanently you can create a modprobe file `kvm.conf` :

```
options kvm ignore_msrs=1
```

## AMD Ryzen / BIOS updates (AGESA) yields "Error: internal error: Unknown PCI header type '127'"

AMD users have been experiencing breakage of their KVM setups after updating the BIOS on their motherboard. There is a kernel [patch \(https://clbin.com/VCiYJ\)](https://clbin.com/VCiYJ), (see [Kernel/Arch Build System](#) for instruction on compiling kernels with custom patches) that can resolve the issue as of now (7/28/19), but this is not the first time AMD has made an error of this very nature, so take this into account if you are considering updating your BIOS in the future as a VFIO user.

## Navi 10 AMD GPU not resetting properly yielding "Error: internal error: Unknown PCI header type '127'" (Separate issue from the one above)

Users with a Navi architecture GPU are able to start up their KVM setup but after shutting down the VM they are not able to start it up are faced with the unknown PCI header error this forces the user to restart the host to be able to resolve the issue. This is due to the NAVI reset bug which is explained here along with a kernel patch [\[5\] \(https://forum.level1techs.com/t/navi-reset-kernel-patch/147547\)](https://forum.level1techs.com/t/navi-reset-kernel-patch/147547) (The solution to this bug requires a kernel patch as detailed and can use [Kernel/Arch Build System](#) for instruction)

## Host crashes when hotplugging Nvidia card with USB

If attempting to hotplug an Nvidia card with a USB port, you may have to blacklist the `i2c_nvidia_gpu` driver. Do this by adding the line `blacklist i2c_nvidia_gpu` to `/etc/modprobe.d/blacklist.conf`.

## Host unable to boot and stuck in black screen after enabling vfio

If debug kernel messages during boot are enabled by adding with the `debug ignore_loglevel` [kernel parameters](#), you may see boot stuck with the last message similar to

```
vfio-pci 0000:01:00.0: vgaarb: changed VGA decodes: olddecodes=io+mem,decodes=io+mem:owns=none
```

This can be resolved by disconnecting the passthroughed GPU with your monitor. You may reconnect the passthroughed GPU to a monitor after host is booted.

## AER errors when passing through PCIe USB hub

In some cases passing through a PCIe USB hub, such as one connected to the guest GPU, might fail with AER errors similar to the following:

```
kernel: pcieport 0000:00:01.1: AER: Uncorrected (Non-Fatal) error received: 0000:00:01.1
kernel: pcieport 0000:00:01.1: AER: PCIe Bus Error: severity=Uncorrected (Non-Fatal), type=Transaction Layer, (Requester ID)
kernel: pcieport 0000:00:01.1: AER: device [8086:1905] error status/mask=00100000/00000000
kernel: pcieport 0000:00:01.1: AER: [20] UnsupReq (First)
kernel: pcieport 0000:00:01.1: AER: TLP Header: 00000000 00000000 00000000 00000000
kernel: pcieport 0000:00:01.1: AER: device recovery successful
```

In that case you might want to try booting the kernel with the parameter `pci=noaer` to have it ignore these errors altogether.

## See also

- [Discussion on Arch Linux forums \(https://bbs.archlinux.org/viewtopic.php?id=162768\)](https://bbs.archlinux.org/viewtopic.php?id=162768) | [Archived link \(https://archive.is/kZYMt\)](https://archive.is/kZYMt)
- [User contributed hardware compatibility list \(https://docs.google.com/spreadsheet/ccc?key=0Aryg5nO-kBebdFozaW9tUWdVd2VHM0lvck95TUIpMIE\)](https://docs.google.com/spreadsheet/ccc?key=0Aryg5nO-kBebdFozaW9tUWdVd2VHM0lvck95TUIpMIE)
- [Example script \(https://pastebin.com/rcnUZCv7\)](https://pastebin.com/rcnUZCv7) from <https://www.youtube.com/watch?v=37D2bRsthfl>
- [Complete tutorial for PCI passthrough \(https://vfio.blogspot.com/\)](https://vfio.blogspot.com/)
- [VFIO users mailing list \(https://www.redhat.com/archives/vfio-users/\)](https://www.redhat.com/archives/vfio-users/)
- [#vfio-users on freenode \(ircs://chat.freenode.net/vfio-users\)](https://ircs://chat.freenode.net/vfio-users)
- [YouTube: Level1Linux - GPU Passthrough for Virtualization with Ryzen \(https://www.youtube.com/watch?v=aLeWg11ZBn0\)](https://www.youtube.com/watch?v=aLeWg11ZBn0)

- [/r/VFIO: A subreddit focused on vfio \(https://www.reddit.com/r/VFIO\)](https://www.reddit.com/r/VFIO)
  - [GVT-d: passthrough of an entire integrated GPU \(https://github.com/intel/gvt-linux/wiki/GVTd\\_Setup\\_Guide\)](https://github.com/intel/gvt-linux/wiki/GVTd_Setup_Guide)
- 

Retrieved from "[https://wiki.archlinux.org/index.php?title=PCI\\_passthrough\\_via\\_OVMF&oldid=625389](https://wiki.archlinux.org/index.php?title=PCI_passthrough_via_OVMF&oldid=625389)"

---

**This page was last edited on 15 July 2020, at 21:41.**

Content is available under [GNU Free Documentation License 1.3 or later](#) unless otherwise noted.