vGPU配置安装手册

vWS 14.3+vSPhere 7.0.3+vCenter 7.0.3

修订记录

| Date | Version | Authors | Description |
|------------|---------|-----------|-------------|
| 2022.10.27 | V1.0 | Leon Wang | 概述 |
| 2022.11.28 | V1.5 | Leon Wang | 安装部署 |
| | | | |

1. NVIDIA vGPU概述

1.1 什么是NVIDIA vGPU

NVIDIA 虚拟 GPU (vGPU) 软件为众多工作负载(从图形丰富的虚拟工作站到数据科学和 AI)提供强大的 GPU 性能,使 IT 能够利用虚拟化的管理和安全优势以及现代工作负载所需的 NVIDIA GPU 的性能。 NVIDIA vGPU 软件安装在云或企业数据中心服务器的物理 GPU 上,会创建虚拟 GPU,这些 GPU 可以在多个虚拟机(可随时随地通过任意设备访问)之间共享。



1.2 NVIDIA vGPU软件产品分类

1.2.1 NVIDIA 虚拟计算服务器 (vCS)

仅支持CUDA, Guest OS仅支持Linux OS。加速基于 KVM 的基础架构上的虚拟化 AI 计算工作负载。若是基于VMWare VSPhere的基础架构,请查阅<u>NVIDIA AI Enterprise软件套件</u>。

AI、深度学习和数据科学工作流程需要出色的计算能力。借助新款 NVIDIA 数据中心 GPU(包括 NVIDIA A30 Tensor Core GPU), NVIDIA 虚拟计算服务器 (vCS) 可助力数据中心加速服务器虚拟化,以便可以在由 NVIDIA vGPU 技术驱动的虚拟机 (VM) 中运行计算密集程度极高的工作负载(例如人工智能、深度学习和数据科学)。

1.2.2 NVIDIA RTX 虚拟工作站 (vWS)

支持DiriectX, OpenGL, Vulkan等图形API, 同时支持CUDA/OpenCL计算API, 适用于使用图形应用程序的创意和技术专业人士的虚拟工作站。[NVIDIA RTX vWS 解决方案概述]

NVIDIA RTX 虚拟工作站 (vWS) 软件与我们世界领先的 GPU 相结合,可以为视觉计算提供强劲动力,从 数据中心或云向任何设备提供极其强大的虚拟工作站。数百万创意专业人员和技术专业人员可以随时随 地访问要求极高的应用程序,不仅可以获得堪比物理工作站的卓越性能,而且还可满足更高的安全性需 求。

1.2.3 NVIDIA 虚拟 PC (vPC)

支持DiriectX, OpenGL, Vulkan等图形API, 专注于桌面虚拟化,适用于使用办公程序和多媒体应用程序的知识工作者的虚拟桌面 (VDI)。[NVIDIA vPC和vApp 解决方案概述]

NVIDIA 虚拟 PC 软件和 NVIDIA GPU (包括 NVIDIA A16) 能够加速生产力应用,并实现令人难以置信的用户体验,因此当今员工可以随时随地无缝访问所需工具。

1.2.4 NVIDIA 虚拟应用程序 (vApp)

支持DiriectX, OpenG等图形API, 专注于应用程序虚拟化, 采用远程桌面会话主机 (RDSH) 解决方案的 应用程序流。[<u>NVIDIA vPC和vApp 解决方案概述</u>]

使用虚拟 GPU 和 NVIDIA 虚拟应用程序 (vApp) 软件提高工作效率,并利用远程桌面共享主机 (RDSH) 解 决方案(包括 Citrix 虚拟应用程序和 VMware Horizon vApp)加速应用程序流。NVIDIA vApp 允许用 户随时随地在任何设备上都能够以完整性能使用任何 Windows 应用程序。

1.2.5 vWS, vPC和vCS特性比较

| Configuration and Deployment | RTX vWS | vPC | vCS |
|--------------------------------------|--------------------|--------------------------------|-----------|
| Desktop Virtualization | ✓ | ✓ | |
| Server Virtualization | | | ✓ |
| Windows OS Support | ✓ | ✓ | |
| Linux OS Support | ~ | ~ | ✓ |
| NVIDIA Graphics Driver | ✓ | ✓ | |
| NVIDIA RTX Enterprise Driver | ✓ | | |
| NVIDIA Compute Driver | | | ✓ |
| Multi-vGPU/NVLink | ~ | | ✓ |
| ECC Reporting and Handling | ✓ | | ✓ |
| Page Retirement | ✓ | | ✓ |
| | | | |
| Display | RTX vWS | vPC | vCS |
| Maximum Hardware Rendered Display | Four 5K, Two 8K | Four QHD, Two 4K, One 5K | One 4K |
| Maximum Resolution | 7680×4302 | 5120x2880 | 4096x2160 |

| Advanced Professional Features | RTX vWS | vPC | vCS |
|-----------------------------------|--|------------|--|
| ISV Certifications | \checkmark | | |
| NVIDIA CUDA/OpenCL | \checkmark | | ✓ |
| Graphics Features and APIs | RTX vWS | vPC | vCS |
| | • | • | • |
| OpenGL Extensions (WebGL) | ~ | ~ | |
| Insitu Graphics/GL Support | | | ✓ |
| RTX Optimizations | ✓ | | |
| DirectX | ✓ | ✓ | |
| Vulkan Support | ✓ | | ✓ |
| Profiles | RTX vWS | vPC | vCS |
| Max Frame Buffer Supported | 48GB | 2GB | 80GB |
| Available Profiles | 0Q, 1Q, 2Q, 3Q, 4Q, 6Q, 8Q, 12Q, 16Q, 24Q, 32Q, 48Q | 0B, 1B, 2B | 4C, 5C, 6C, 8C, 10C, 12C, 16C, 20C, 24C,40C, 32C, 48C, 80C |

1.3 vGPU带给客户的收益

- 裸机性能 提供几乎与裸机环境无差别的性能。
- 管理和监控 利用常见的数据中心管理工具,例如实时迁移。
- 出色的资源利用率 使用部分或多 GPU 虚拟机 (VM) 实例调配 GPU 资源。
- 提高业务连续性 响应不断变化的业务需求和远程团队。

1.4 vGPU支持的硬件GPU

NVIDIA 虚拟 GPU (vGPU) 软件在 NVIDIA GPU 上运行。选择以下合适的 GPU 以满足您的需求。

https://www.nvidia.cn/data-center/graphics-cards-for-virtualization/

1.5 vGPU Profile

NVIDIA vGPU 软件允许您对 NVIDIA 数据中心 GPU 进行分割。 然后使用 vGPU 配置文件(Profile)将 这些虚拟 GPU 资源分配给虚拟机 (VM)。 选择正确的 vGPU 配置文件将提高您的 vGPU 环境的总拥有成 本、可扩展性、稳定性和性能。

vGPU分为不同的系列。 C-profile需要 NVIDIA vCS 许可证; Q-profile 需要 NVIDIA vWS 许可证,也可以与 vWS 许可证一起使用;B-profile需要 NVIDIA vPC许可证,也可以与 vWS 许可证一起使用;A-profile需要 NVIDIA vAPP许可证。下表提供了详细 vGPU 配置文件信息。

| Series | Optimal Workload |
|----------|--|
| Q-series | Virtual workstations for creative and technical professionals who require the performance and features of Quadro technology |
| C-series | Compute-intensive server workloads, such as artificial intelligence (AI), deep learning, or high-performance computing (HPC) ² , ³ |
| B-series | Virtual desktops for business professionals and knowledge workers |
| A-series | App streaming or session-based solutions for virtual applications users ⁶ |

物理GPU支持的vGPU Profile详细信息,请查阅<u>Virtual GPU Types Reference</u>。

1.6 快速选择适合您的vGPU产品



vGPU 14是当前最新版本,详细版本信息请查阅, https://docs.nvidia.com/grid/

| vGPU Software | vGPU Manager | Linux Driver | Windows Driver | Release Date |
|---------------|--------------|--------------|----------------|---------------|
| 14.3 | 510.108.03 | 510.108.03 | 513.91 | November 2022 |
| 14.2 | 510.85.03 | 510.85.02 | 513.46 | August 2022 |
| 14.1 | 510.73.06 | 510.73.08 | 512.78 | May 2022 |
| 14.0 | 510.47.03 | 510.47.03 | 511.65 | February 2022 |

Branch status: Production Branch supported until February 2023

详细的vGPU支持平台和软件兼容性列表,请查阅<u>Virtual GPU Software Supported Products</u> 对于NVAIE支持的平台和软件兼容性列表,请查阅<u>NVIDIA AI Enterprise Product Support Matrix</u>

1.7 如何购买vGPU Licnese

NVIDIA vGPU 软件产品可以购买带有支持更新和维护订阅 (SUMS) 的永久许可,也可以购买年度订阅。 永久许可证赋予用户无限期使用软件的权利,不会过期。所有具有永久许可的 NVIDIA vGPU 软件产品 必须与五年 SUMS 一起购买。一年期 SUMS 仅适用于续订。 年度订阅服务是一种更实惠的选择,可让 IT 部门更好地管理许可证数量的灵活性。在软件订阅许可期限内,具有年度订阅的 NVIDIA vGPU 软件产品与 SUMS 捆绑在一起。

| Entitlement | NVIDIA vGPU Production SUMS |
|------------------------------|---|
| Maintenance | Access to all maintenance releases, defect resolutions, and security patches for flexibility in upgrading as per the NVIDIA Virtual GPU Software Lifecycle Policy |
| Upgrades | Access to all new major version releases including feature enhancements and new hardware support |
| Long-term branch maintenance | Available for up to 3 years from general availability as per the NVIDIA Virtual GPU Software Lifecycle Policy |
| Direct support | Direct access to NVIDIA support engineering for timely resolution of customer-specific issues |
| Support availability | Customer support available during standard business hours Cases accepted 24 \times 7 |
| Knowledgebase access | \checkmark |
| Web support | \checkmark |
| E-mail support | \checkmark |
| Phone support | \checkmark |

更多详细信息,请查阅。 NVIDIA Virtual GPU Software Packaging, Pricing, and Licensing Guide

1.8 如何下载和申请vGPU POC

购买了 NVIDIA vGPU 软件的客户可以从 NVIDIA 许可门户下载驱动程序。 请检查您的 NVIDIA 授权证书,了解有关如何注册访问 NVIDIA 许可门户的信息,您可以在其中兑换产品激活密钥 (PAK) 和下载驱动程序。

如果您已经在 NVIDIA 许可门户中注册并兑换了 PAK,则可以通过在此处登录许可门户来访问您的vGPU 软件: //nvid.nvidia.com/dashboard/

如果您还没有购买 NVIDIA vGPU 软件并想注册免费试用 90 天,请点击"软件下载"按钮进行注册: //www.nvidia.com/object/vgpu-evaluation.html

1.9 Hands-on说明

本手册,将选取vWS产品,在VMWare vSPhere平台上,部署vGPU解决方案。提供完整的部署和配置流程,供NPN参考。

安装流程:

•安装环境准备

•vSphere 7.0.3安装

•vWS 软件获取

•vWS Host Driver安装

•vCenter安装

•VM安装和配置

•vWS Guest OS Dirver安装

•vWS Guest OS Driver 许可安装

2. 安装环境准备

•满足VMware 兼容性指南, 且支持NVIDIA GPU的硬件服务器

•具备vGPU硬件GPU支持列表的GPU

•如果是Ampere架构GPU,需要在服务器BIOS中开启如下设置:

•VT-D/IOMMU

•SR-IOV

•下载最新版本 VMWare vSphere 7.0 U3 <u>https://customerconnect.vmware.com/cn/downloads/info/</u> <u>slug/datacenter_cloud_infrastructure/vmware_vsphere/7_0</u>

•利用U盘或者DVD制作启动盘,引导裸金属服务器安装ESXi 7.0

3. VSphere安装

• 利用U盘或者DVD引导服务器,进入VMWare ESXi安装界面。

| | Load ing | ESXi | installer | | | |
|--|----------|------|---|------|------|----------|
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| < <u>ENTER</u> : Boot> Automatic boot in 3 seco | onds | | <shift+0:< th=""><th>Edit</th><th>boot</th><th>options></th></shift+0:<> | Edit | boot | options> |

• 安装完成后,重启服务器,进入ESXi登录界面。按F2进入ESXi,配置root用户密码和ESXi的IP地址,该IP地址用于后续用户通过Web界面登录这台安装了ESXi的Host。

| VMware ESXi 7.0.0 (VMKernel Release Duild 14096552) | |
|---|------------------------|
| Dell Inc. PowerEdge C6220 | |
| 2 x Intel(R) Xeon(R) CPU E5-2620 0 9 2.00GHz 64 Gi0 Menory | |
| To wanage this host, go to: https://0.0.0.0/ (Maiting for DMCP) https://(fe00::f2Mdia2ff:fe75:dblc]/ (STATIC) | |
| Marning: DHCP lookup failed. You way be unable to access this system un network configuration. | ntil you custonize its |
| #2) Controlline Sectors/View Loop | OTTO SHE SHOWS IN |



• 配置完成后,退出ESXi配置界面。然后通过另外一台PC的浏览器,登录这台ESXi Host,登录界面如下。登录后,可以通过点击"主机",查看Host配置详情。





• 在"主机"界面,选择管理—硬件—PCIE设备,检查当前GPU的SR-IOV配置为active状态,如否,点 击左上方"配置SR-IOV"来配置。



 注意,关于 ESXi 评估和许可模式:可以使用评估模式来浏览 ESXi 主机的全套功能。评估期60天。 https://docs.vmware.com/cn/VMware-vSphere/7.0/com.vmware.esxi.install.doc/GUID-17862 A54-C1D4-47A9-88AA-2A1A32602BC6.html

4. vWS软件获取

4.1 软件下载

• 在NVIDIA官网vGPU<u>申请</u>,或者购买授权后,在<u>授权许可网站</u>下载对应操作系统的最新版本驱动 包。

申请界面:



授权许可网站-软件下载



| ← → C (à ui.li | censing.nvidia.com/softw | /are | | | Q 🖻 🛊 | * 🗆 🛛 | (RK : |
|------------------------|--|--|---|--|--------------|---------------------------|------------|
| NVINFO Q Home | Workday < PID << | AE China Team T 📮 | Google 翻译 🛛 📥 Ni | PN 中文資料库 🌰 Technical Review 🥌 我的文件 - OneDri | Ta Bing Micr | osoft Tr | |
| CINIDIA, LICENSING | | | | NVOJA AMPLICATION HLB 💻 🖉 | | | ⊗ logout |
| ධ් DASHBOARD | | | | | | | |
| ENTITLEMENTS | Software Downlo | oads | | | | & ADDITIONA | LSOFTWARE |
| ■ UCENSE SERVERS > | | | | | | | _ |
| A NETWORK ENTITLEMENTS | | 🖉 FEATURED 🛛 👌 | ALLAVAILABLE | | | | |
| ID VIRTUAL GROUPS | | | | | | | へ <u>魚</u> |
| AL USER MANAGEMENT | | | | | | aparter (* 163433 | ·0· W |
| SOFTWARE DOWINLOADS | PLATFORM 🍸 🗘 | PLATFORM VERSION () | PRODUCT VERSION () | DESCRIPTION 🖓 🗘 | | DATE 0 | |
| E LEASES | VMware vSphere | 7.0 | 11.9 | Complete vGPU package for vSphere 7.0 including supported guest drivers | | Aug 2, 2022 | Download |
| B SERVICE INSTANCES | VMware vSphere | 7.0 | 13.4 | Complete vGPU 13.4 package for VMware vSphere 7.0 including supported guest drivers | | Aug 2, 2022 | Download |
| d' APORTS | VMware vSphere | 7.0 | 14.2 | Complete vGPU 14.2 package for VMware vSphere 7.0 including supported guest drivers | | Aug 2, 2022 | Download |
| ES EMAIL ALERTS | Libertha KAM | All Supported | 114 | Complete sci2(1) 13 g package for 1 burns 40M ALL including supported event drivers. | | Acce 2 2022 | Download |
| Q₂ SUPPORT | vGPU For information about the software NVIDIA vGPU documentation is a | are Mecycle for NV/DIA virsual GPU S available at: <u>tropp://docs.cv/dia.com</u> / | ofware vielt <u>topsziláso</u> zunoda.com stál | ngrafi sena index hani | | | |
| ((COLLAPSE | 10 v downloads per page | | | | ≪ < α-1 | e of 32 downloads) 1 of 4 | Fpages > ≫ |

4.2 软件内容

- 下载完成后,将获得vWS软件ZIP包,解压缩后,主要包含以下3个部分:
 - vWS用户手册 主要包含快速启动手册,用户手册,版本Release Notes等
 - 510.85.03-510.85.02-513.46-grid-gpumodeswitch-user-guide.pdf
 - 510.85.03-510.85.02-513.46-grid-licensing-user-guide.pdf
 - 510.85.03-510.85.02-513.46-grid-software-quick-start-guide.pdf
 - 510.85.03-510.85.02-513.46-grid-vgpu-release-notes-vmware-vsphere.pdf
 - 510.85.03-510.85.02-513.46-grid-vgpu-user-guide.pdf
 - 510.85.03-510.85.02-513.46-whats-new-vgpu.pdf
 - vWS Host Drivers 安装在ESXi Host上的驱动程序, 文件后缀为VIB。



○ vWS Guest Drivers - 安装在VM中Gust OS上的驱动程序,包含Windows和Linux版本。

| -GRID-vSphere-7.0-51 > Guest_Drivers v ひ の 捜索"Guest_Drivers" | | |
|---|--------|------------|
| □ 名称 | 类型 | 压缩大小 |
| 513.46_grid_win10_win11_server2019_server2022_64bit_international.exe | 应用程序 | 621,139 KB |
| nvidia-linux-grid-510_510.85.02_amd64.deb | DEB 文件 | 323,944 KB |
| nvidia-linux-grid-510-510.85.02-1.x86_64.rpm | RPM 文件 | 420,139 KB |
| NVIDIA-Linux-x86_64-510.85.02-grid.run | RUN 文件 | 328,034 KB |

5. vWS Host Driver安装

5.1 安装

•SSH登录ESXi Host

•将Host Driver文件解压缩,并上传VIB (The vGPU Manager vSphere Installation Bundles) 文件到 ESXi Host

•将进入ESXi 主机维护模式

•esxcli system maintenanceMode set --enable true

•安装VIB

•esxcli software vib install -d /vmfs/volumes/datastore/softwarecomponent.vib

•退出ESXi 主机维护模式

•esxcli system maintenanceMode set --enable false

•重启ESXi主机

Reboot

```
[root@esxi:~] esxcli software vib install -v directory/NVIDIA-vGPU-
VMware_ESXi_7.0_Host_Driver_510.85.03-10EM.700.0.0.8169922.vib
Installation Result
Message: Operation finished successfully.
Reboot Required: false
VIBs Installed: NVIDIA-vGPU-
VMware_ESXi_7.0_Host_Driver_510.85.03-10EM.700.0.0.8169922
VIBs Removed:
VIBs Removed:
VIBs Skipped:
```

5.2 验证

•验证Host Driver安装

•ESXi Host重启后, SSH登录, 使用nvidia-smi命令, 查看GPU信息。

```
[root@bogon:~] nvidia-smi
Mon Sep 19 16:15:45 2022
+-----
| NVIDIA-SMI 470.141.05 Driver Version: 470.141.05 CUDA Version: N/A
| GPU Name Persistence-M | Bus-Id Disp.A | Volatile Uncorr. EC
| Fan Temp Perf Pwr:Usage/Cap| Memory-Usage | GPU-Util Compute N
                      MIG N
           1
L
0 NVIDIA A40 On | 00000000:98:00.0 Off |
0% 38C P8 34W / 300W | 0MiB / 45634MiB | 0% Defau
                                     N/
                ---+-----
                             --+---
Processes:
GPU GI CI
          PID Type Process name
                                   GPU Memor
  ID ID
                                    Usage
No running processes found
[root@bogon:~]
```

6. vCenter安装

•<u>下载</u>vCenter ISO文件



•利用虚拟光驱加载ISO文件后,运行vcsa-ui-installer/win32/installer.exe。

安装界面启动后,选择"安装"



"设备部署目标",请选择希望安装VCenter的ESXi Host,填写Host IP和登录信息,VCenter将以VM的形式,部署在这个ESXi Host上。

| Vm 安果 - 第一阶段: 部署 vCenter Se | rver Appliance | | | |
|-----------------------------|-----------------------------|-------------------------|----------------|--|
| 1 周介 2 最佳用户许可协议 | 设备部署目标 182%和#10%25.300元年 | 包属中部署设备的 ESX0 王机威 vCerr | ter Server 宾德。 | |
| 3 设备起署目标 | ESXI主我名武 vCenter Server 名称 | 172.16.0.101 | Φ | |
| 4 设置设备透料机 | HTTPS INC | 443 | | |
| 5 选择部署大小 | R P& | root | • | |
| 6 选择数据存储 | 24 | | | |
| 7 配置网络设置 | L | | | |
| 8 #DIEAUG#1800 | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | RA 1-9 7-9 | |

注意: VCenter安装, ESXi Host至少需要500G存储空间

| vm 安装 - 第一阶段: 部署 设备 | | | | | | | | | | | | | |
|--|-----|-------------|------|---------|-----|----------|----|----------|---|----------|----|------|-----|
| 1 简介 | 选打 | 译数据有 | 民储 | | | | | | | | | | |
| 2 最终用户许可协议 3 设备部署目标 | . 5 | 装在可从目 | 标主机 | 访问的现有 | 数据有 | 储上 | | | | | | | |
| 4 设置设备虚拟机 | | 仅显示兼容 | 容的数据 | 存储 | | | | | | | | | |
| 5 选择部署大小 | | 部 | Τ | 类型 | Τ | 安量 | т | 可用 | т | 己豐务 | Ŧ | 精何置备 | т |
| 6 选择数据存储 | | latastore50 | 00G | VMFS-5 | | 457.75 (| GB | 94.03 GE | 3 | 363.72 0 | ЭB | 受支持 | |
| 7 配置网络设置 | | | | | | | | | | | | | 1 項 |
| 8 即将完成第1阶段 | ~ | 启用精简码 | 短模式 | i | | | | | | | | | |
| | 0\$ | 装在包含目 | 目标主机 | 的新 vSAN | 集群 | Ŀ | | | | | | | |

7. VM安装和配置

7.1 ESXi Host开启Shared Direct

•使用 vSphere Web Client 登录到 vCenter Server。

•在导航树中,选择您的 ESXi 主机并单击配置选项卡。

•从菜单中选择图形,然后单击主机图形选项卡。

•在主机图形选项卡上,单击编辑。

•重启ESXi主机

| Getling Started Summary Monitor Configure Permissions VMs Resource Pools Datastores Networks Update Manager | | | | | |
|--|---|--|---------------------------------|-------------------------------------|--|
| 4 Time Configuration Authentication Services Certificate Prover Management Advanced System Settings System Resource Reservation Security Profile System Swap | Host Graphics Caraphics Devices Graphics Devices NutDiATesia M60 NVIDIATesia M60 | Vender NVIDIA Corporation NVIDIA Corporation | Aslive Type Shared Shared | Cerfigured Type Shared Shared | Q Filter - Menory 7 98 08 7 99 08 |
| Host Profilo - Hardware Processors Memory Graphics | stem swap stem stem stem stem stem stem stem stem | | | | 2 itoms Deport - Cacopy - |





7.2 创建VM和安装Guest OS

详见<u>Create a Virtual Machine</u>

7.3 为VM挂载vGPU

•VM安装完成后,关闭VM电源。右键单击VM,选择 Edit Settings,再选择ADD NEW DEVICE,添加PCI Device

•选择vGPU profile,比如nvidia-a40-12q,关于vGPU 类型说明,请查看宣方文档。

•一张卡只能使用一个profile,比如本案例只有1张A40卡,第一个虚拟机启用了 a40-12q的vGPU profile,后面的虚拟机也只能使用 12q的profile。

Edit Settings | Wint0Pro-clone

Virtual Hardware VM Options

| | | ADD NEW DEVICE~ |
|-----------------------|------------------------|-----------------------------|
| > CPU | 6 ~ | Disks, Drives and Storage |
| > Memory | 20 | Hard Disk Culation Link |
| > Hard disk 1 | 128 GB v | RDM Disk |
| > SCSI controller 0 | LSI Logic SAS | Host USB Device |
| > Network adapter 1 | VM Network v | NVDIMM |
| > Network adapter 2 | VM Network ~ | Controllers |
| > CD/DVD drive 1 | Datastore ISO File 🗸 🗸 | SATA Controller |
| > USB xHCl controller | USB 3.1 | SCSI Controller |
| > Video card | Auto-detect settings v | Other Devices |
| > Security Devices | Not Configured | PCI Device |
| VMCI device | | Watchdog Timer |
| SATA controller 0 | AHCI | Serial Port |
| > Other | Accitional Hardware | Network Adapter |
| | | |
| | | CANCEL |

-1

Edit Settings | Win11

| | | | ADD NEW DEVIC |
|----------------------------------|---|---|---|
| CPU | 6 ~ | | G |
| Memory | 16 | ~ | GB ~ |
| Hard disk 1 | 200 | GB ~ | |
| SCSI controller 0 | LSI Logic SAS | | |
| Network adapter 1 | VM Network ~ | | Connect |
| CD/DVD drive 1 | Datastore ISO File | ~ | Connect |
| USB xHCl controller | USB 3.1 | | |
| New PCI device | NVIDIA GRID vGPU | nvidia_a40-1b | |
| | nvidia_440-10 nvidia_440-2b nvidia_440-1q nvidia_440-2q nvidia_440-3q | ual machine o CI/PCIe passti er guide for vi PCIe passthro | operations are nrough devices are irtual machine operation ough devices. |
| Video card | nvidia_a40-4q nvidia_a40-6q | s v | |
| Security Devices | nvidia_a40-8q | | |
| | | | |
| VMCI device | nvidia_a40-16q nvidia_a40-24q | | |
| VMCI device SATA controller 0 | nvidia_a40-16q nvidia_a40-24q nvidia_a40-48q nvidia_a40-1a | | |

8. vWS Guest OS Dirver安装

•加电VM,为VM安装Guest OS Driver

•将根据Guest OS的操作系统,选择之前Guest Drivers目录下的驱动程序,上载到VM中,安装驱动。安装成功后,可以通过nvidai-smi命令,查看vGPU的状态。

9. vWS Guest OS Driver 许可安装

9.1 NLS背景介绍

•在2021年8月伴随vGPU 13.0 的发布,NVIDIA推出了全新的软件License系统(NLS)来替代之前基于 Flexnet的License 服务。

•NLS提供两种License 服务实例类型来对vGPU客户进行软件授权。

•如果vGPU客户端可以访问互联网环境,vGPU可以访问NVIDIA 官方服务器上面的CLS服务直接进行授权,就可以选择使用CLS,使用CLS则无需在本地环境再搭建授权服务器,因此更加方便和快捷,但要注意公网访问的可靠性。



•另一种是在企业私网内部署本地DLS 服务器进行License授权。主要面向vGPU客户端全部在企业内网, 且不能访问互联网环境时应该采用的方式。



9.2 创建License Server

•登录用户企业账号管理门户NLP后,先创建License Server。

| 🚳 NVIDIA. LICENSING | NYERAPPLICATION HUB and intermigration with an intermitian and a second statements of the second statement of the second state |
|------------------------|--|
| 🛱 DASHBOARD | |
| ENTITLEMENTS | Create License Server ③ Helle? |
| 🗎 LICENSE SERVERS 🌒 🗸 | Crash a lorne server in NVIDA INFR-GER (in-definition/Cryspe) / Singe NVGM, NFR-GER (in-definition/Cryspe) |
| IIST SERVERS | |
| 📖 CREATE SERVER 🛛 🕘 | Create addrey terver () In the scene area in one assession as a stary internet even area of events, europe Create addrey rever |
| A NETWORK ENTITLEMENTS | Basic details →I Select features →I Preview server creation |
| TINTUAL GROUPS | (j) Enter a name, and description for this new license server |
| A USER MANAGEMENT | New |
| A SOFTWARE DOWNLOADS | Merim CLS-2022 0 |
| 置 EVENTS | Description |
| LEASES | Test Only 4 |
| SERVICE INSTANCES | |
| 🖉 API KEYS | |
| A SUPPORT | Express CLS Installation? The server will be installed on the default CLS service instance |
| | Next: Select features ->1 |

•添加您现有的License类型和分配的数量。

| ■ ENTITLEMENTS ■ LICENSE SERVERS ~ | Create License Server () Hele? Create a licing derive in NNDA INFR-SEX (IL-0011W0010276)(BW) / Group NNDIA INFR-SEX (IL-0011W0000276)(BW) |
|------------------------------------|---|
| LIST SERVERS | Create legacy server () If the loanse senter is to be installed on a legacy loansing system senter (zne-XLS), enable "Onster legacy server" |
| CREATE SERVER | |
| A NETWORK ENTITLEMENTS | Basic datails -> 🕗 Select features -> Previow server creation |
| D VIRTUAL GROUPS | ③ Select one or more entitlement features to add to the new license server |
| AL USER MANAGEMENT | ▼ workstation × |
| & SOFTWARE DOWNLOADS | |
| EVENTS | $\blacksquare NAME \ \forall \ \Diamond \qquad PRODUCT \ KEYID \ \forall \ \Diamond \qquad START \ DATE \ \forall \ \Diamond \qquad EXPIRATION \ \forall \ \Diamond \qquad AVAILABLE \ \forall \ \Diamond \qquad ADDED \ \Diamond \ \bullet \ \bullet$ |
| ELEASES | WUDA RTX Virtual Workstatien-5.0 1ug0misszt- frakútrgi- brywgen Active Nov 11, 2021 M Nov 11, 2022 31 20 |

•创建后的授权服务器,如下图。

| Merlin-CLS-2022 is ENABLED | ~ |
|---|---|
| Status: 🐻 ENABLED Type: NVIDIA Created: Sep 6, 2022 6:56 PM Modified: Sep 6, 2022 6:56 PM | |
| Bervice instance: 0011w000027i5ytgey-2022-09-06 10-56 CLS Install Status: III INSTALLED | |
| Description: Test Only | |
| Overview Server Features License Pools Fulfillment Conditions Leases | |
| ABOUT THIS SERVER | |
| License server Merlin-CLS-2022 is enabled and will serve leases, you can not make changes while the server is enabled in DisABLE SERVER | |

9.3 生成Token文件

•License Server创建成功后,到服务实例页面,生成用于vGPU客户的Token文件。

•生成License 配置Token以后,下载该tok文件。

```
•将Tok文件复制到vGPU客户端VM内。
```

| LICENSE SERVERS | View your service instances in NVIDIA INFR-GEN | (iic-0011w000027i5yiqay) | | |
|---|--|---|---------------------------------------|--|
| NETWORK ENTITLEMENTS | 😇 CLS | | | |
| VIRTUAL GROUPS | | | | |
| USER MANAGEMENT | | | | updated 🛞 8:32:07 PM <table-cell> 🍸 👱 🤅</table-cell> |
| SOFTWARE DOWNLOADS | NAME \bigtriangledown \diamondsuit | ENVIRONMENT T | status $\overline{\gamma}$ \Diamond | date created \bigtriangledown \diamondsuit |
| EVENTS | Merlin-CLSI-2022 (fa56ae40-e663-4600-e6b9-edaeed0b5297) | © CLS Default | Registered | Sep 6, 2022 6:56 PM |
| SERVICE INSTANCES | tost /st0/54/11.838/4.40+0.5+90.47408+0.91157) | ⇔ cls | Registered | Jul 19, 2022 10: Generate client config tok |
| API KEYS | in instance | | | Settings Delete |
| | | | | |
| ate a configuration token | Fulfillment class reference | es | | |
| te a configuration token cope references Search scope referen | Fulfillment class reference | | | |
| ate a configuration token cope references Search scope references SERVER NAME | Fulfillment class reference rces REFER | | | |
| Search scope references Search scope references Search scope references Merlin-CLS-2022 | Fulfillment class reference rces Y ↓ REFER 18898a | PS RENCE ♥ 0 72-b2a3-4526-b311-4df0c501049 | 0 | |
| a configuration token cope references P Search scope referen SERVER NAME Merlin-CLS-2022 | Fulfillment class reference sces S ♀ ♀ REFER 18898a ≪ < (1-1 of 1 so | PS RENCE ▽ 72-b2a3-4526-b31f-4df0c501049 ope references) 1 of 1 pages | a >>> | |
| te a configuration token cope references Search scope references SERVER NAME Merlin-CLS-2022 | Fulfillment class reference rces REFER 18898a < | RENCE \bigtriangledown \diamondsuit 72-b2a3-4526-b31f-4df0c501049 ope references) 1 of 1 pages $>$ AD CLIENT CONFIGURATION | a >>> TOKEN | |
| ate a configuration token cope references Search scope references Server NAME Merlin-CLS-2022 Client co | Transition token 09.06-2022-20 | PRENCE ♀ ↔ 72-b2a3-4526-b31f-4df0c501049 ope references) 1 of 1 pages > AD CLIENT CONFIGURATION ☆7 | a »» | |
| ate a configuration token cope references Search scope referen SERVER NAME Merlin-CLS-2022 Client_co blob-ittto | for client access to server resources Fulfillment class reference ces COS COS COS COS COS COS COS CO | PS RENCE ♡ 72-b2a3-4526-b31f-4df0c501049 ope references) 1 of 1 pages AD CLIENT CONFIGURATION © [™] -35-40.tok ia1-6PDb-4371 | a >>> TOKEN | |
| te a configuration token cope references Search scope references SERVER NAME Merlin-CLS-2022 dient_cope blob:http | The formation_token_09-06-2022-20 s://ui.licensing.nvidia.com/ef9ae5 | PS RENCE ♥ ♦ 72-b2a3-4526-b31f-4df0c501049 ope references) 1 of 1 pages > AD CLIENT CONFIGURATION €') +35-40.tok 5a1-6600b-4371 | a >> TOKEN | |

9.4 Guest OS授权

将Tok上传到VM后,请根据VM中Gust OS类型,完成vGPU授权。完成后,用户可以完整使用vGPU功能。

9.4.1 Linux

•cd /etc/nvidia, 复制gridd.conf.template为 gridd.conf.

•编辑gridd.conf,只需设置FeatureType的值为要请求的vGPU License类型编号,不要设置其中的ServerAddress。

•然后将下载到的Token文件复制到 /etc/nvidia/ClientConfigToken 目录中。

•systemctl restart nvidia-gridd重启 nvidia-gridd服务,应看到vGPU授权成功。

```
# Description: Set Feature to be enabled
# Data type: integer
# Possible values:
# 0 => for unlicensed state
# 1 => for NVIDIA vGPU (Optional, autodetected as per vGPU type)
# 2 => for NVIDIA RTX Virtual Workstation
# 4 => for NVIDIA Virtual Compute Server
# All other values reserved
FeatureType=2
```

| [root@linuxvm_000084000_80_1 ~]# systemctl restart nvidia-gridd.service [] [root@linuxvm_000084000_80_1 ~]# systemctl status nvidia-gridd.service [] |
|---|
| • nvidia-gridd.service - NVIDIA Grid Daemon |
| Loaded: loaded (/usr/lib/systemd/system/nvidia-gridd.service; enabled; vendor preset: disabled) |
| Active: active (running) since Mon 2021-11-15 12:40:40 CST; 8s ago |
| Process: 26582 ExecStopPost=/bin/rm -rf /var/run/nvidia-gridd (code=exited, status=0/SUCCESS) |
| Process: 26583 ExecStart=/usr/bin/nvidia-gridd (code=exited, status=0/SUCCESS) |
| Main PID: 26585 (nvidia-gridd) |
| Tasks: 4 (limit: 49632) |
| Memory: 1.5M |
| CGroup: /system.slice/nvidia-gridd.service |
| Losses /usr/bin/ordia_gridd |
| |
| Nov 15 12:4 ^图 形用户界面,文本 00084000 80 1 systemd[1]: Stopped NVIDIA Grid Daemon. |
| 000 15 12:4 00084000 80 1 systemd[1]: Starting NVIDIA Grid Daemon |
| Nov 15 12:4描述已自动生成 00084000 80 1 _nvidia-gridd[2585]; Started (26585) |
| Nov 15 12:40:40 linuxym 0000800000 1 system(1): Started NVIDIA Grid Daemon |
| Nov 15 12:40:40 linuxym_000004000 80 1 Jystema[1]. Started midstin of a parameter (ServerAddress) not |
| Nov 15 12:40:40 linuxym_000004000_00_1 nvidia-gridd[25555]; v60n1gridter national (50) |
| Nov 15 12:40:40 linuxym_000004000_00_1 hvidia-gridd[20505]. Ugopre service provider and pode_locked licensi |
| Nov 15 12:40:40 linuxym 000004000 00_1 nvidia-gridd[20505]. NIS initialized |
| Nov 15 12:40:40 linux/m_000004000_02_1 Nvidia-gridd[20505]. Accurring license (Info: ani cls licensing nvid |
| Nov 15 12.40.40 linux/m_00004000_02_1 Nvidia-gridd[2055], Liconse acquired successfully (Info; apic) |
| Nov is 12:40:46 indixin_000084000_80_1 Hvidia-gridd[20585]. Elense acquired successfully. (Inter apricis: |
| lines 1-21/21 (END) |

9.4.2 Windows

•将下载的.tok 文件复制到: C:\Program Files\NVIDIA Corporation\vGPU Licensing\ClientConfigToken 目录。

•重启NvDisplayContainer 服务。

•C:\Program Files\NVIDIA Corporation\NVSMI\nvidia-smi.exe -q 查看License状态。

| Services | | | | | | | × |
|--------------------|-----------------------------|---|--|---------|---------------------------|------------------------------|---|
| File Action View | Help | | | | | | |
| (+ +) 🖬 🗐 🖸 | à 🔒 🛛 🖬 🕨 🔲 💷 🕨 | | | | | | |
| 💁 Services (Local) | Services (Local) | | | | | | |
| | NVIDIA Display Container LS | Name | Description | Status | Startup Type | Log On As | ^ |
| | Stop the service | Network Setup Service Network Store Interface Service | The Network Setup Service This service delivers netwo | Running | Manual (Trig Automatic | Local Syste Local Service | |
| | Kestan the service | NVIDIA Display Container LS | Container service for NVID | Running | Automatic | Local Syste | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

| vGPU Software Licensed Product | |
|--------------------------------|---|
| Product Name | : NVIDIA RTX Virtual Workstation |
| License Status | : Licensed (Expiry: 2021-11-16 5:12:43 GMT) |
| TOMUDU | |

附录1:从 VM 创建模板

既然 VM 已经为 AI 训练和部署推理进行了vGPU的配置,IT 管理员的最终工作流程是创建一个 VM 模板,以便在未来快速部署 VM。 IT 管理员为 VM 创建模板,然后克隆模板以服务于多个 AI 从业者/工程师。以下步骤指导IT管理与创建 OVF 文件模板。

- 关闭虚拟机。
- 在 vCenter 中,右键单击新创建的虚拟机 -> 选择克隆 -> 选择"克隆到模板"。
- 添加名称文件夹 -> 选择计算资源 -> 添加存储 -> 选择您创建的来宾自定义规范 -> 单击完成。

附录2:参考文档

NVIDIA vGPU资料

•NVIDIA vGPU 文档首页 <u>https://docs.nvidia.com/grid/</u>

•vGPU支持Matrix <u>https://docs.nvidia.com/grid/latest/product-support-matrix/index.html</u>

•硬件GPU支持列表 <u>https://docs.nvidia.com/grid/gpus-supported-by-vgpu.html</u>

•vGPU 14 用户文档中心 https://docs.nvidia.com/grid/14.0/index.html

•vGPU中文(非官方) <u>http://vgpu.com.cn/index.html</u>

•NLS-2.0安装配置与PoC手册.pdf

链接: <u>https://pan.baidu.com/s/1aL0S3oyMZyDCafi</u>_IS3ssg

提取码: e592

VMWare VSPhere 7.0 U3

•vSphere7.0 提供了各种安装和设置选项。为确保成功部署 vSphere,应了解安装和设置选项以及任务的执行顺序。

•vSphere 的两个核心组件是 ESXi 和 vCenter Server。ESXi 是用于创建和运行虚拟机和虚拟设备的虚拟 化平台。vCenter Server 是一种服务,充当连接到网络的 ESXi 主机的中心管理员。使用 vCenter Server,您可以池化和管理多个主机的资源

•VMware 兼容性指南 <u>https://www.vmware.com/resources/compatibility/search.php</u>

•VMware ESXi 安装和设置 <u>https://docs.vmware.com/cn/VMware-vSphere/7.0/vsphere-esxi-703-inst</u> <u>allation-setup-guide.pdf</u>

•vCenter Server 和主机管理 <u>https://docs.vmware.com/cn/VMware-vSphere/7.0/vsphere-esxi-vcent</u> <u>er-server-703-host-management-guide.pdf</u>