PRESENTED BY MASAFUMI OHTA @masafumiohta

GPU ON OPENSTACK - GPUインターナルクラウドのベストプラクティス



#### MASAFUMI OHTA

A 'STACKER' LOOKING
INTO GPGPU USE.
VOLUNTEER FOR THE
RASPBERRY PI
FOUNDATION:)



#### PRESENTATION RENEWED

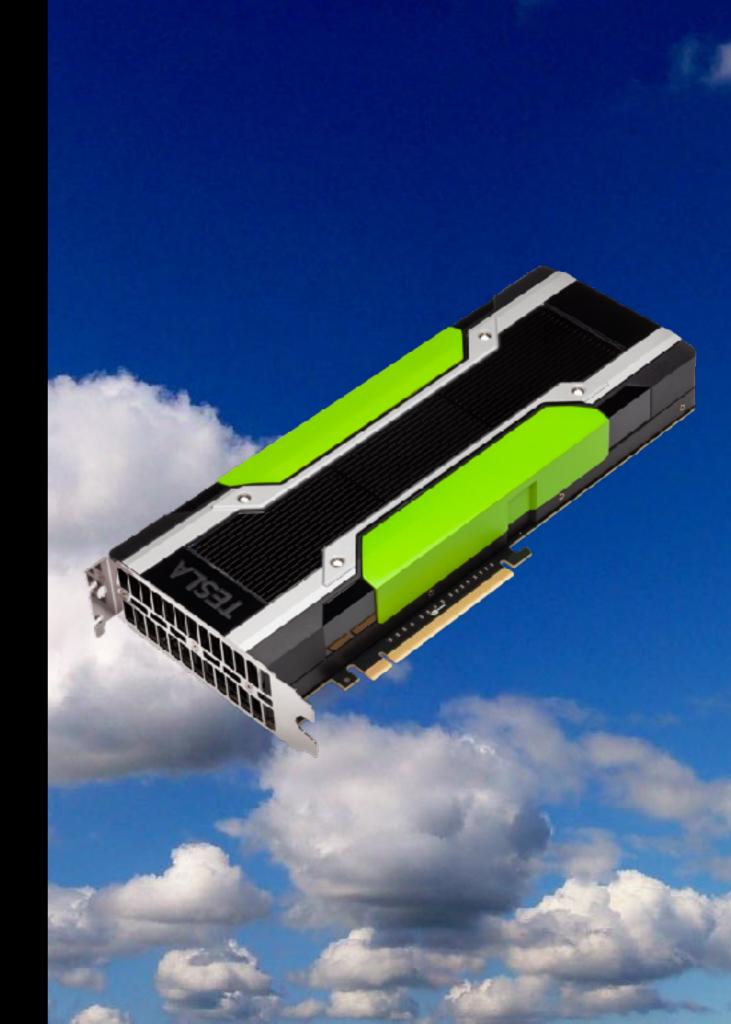
THIS PRESENTATION IS
RENEWED FOR
OPENSTACK DAYS TOKYO
2017.

IT IS INCLUDED SOME
FEEDBACKS FROM THE
SESSION AT #LC3 CHINA
IN BEIJING, 2017



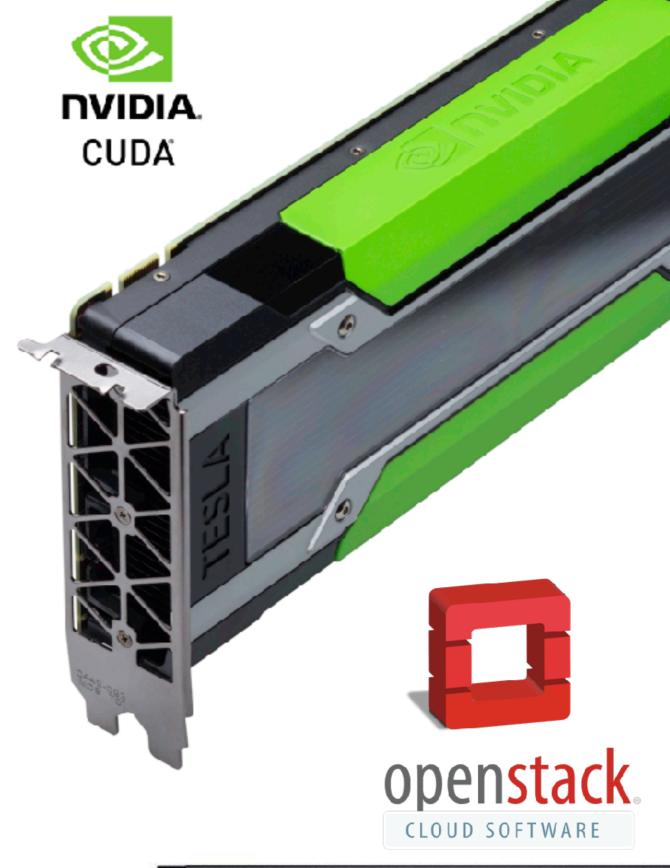
#### THIS PROJECT IS...

ORGIZNED BY VTJ AND
GATHER
ME, NVIDIA (ASK), DELL EMC,
NESIC AND SOME
COMPANIES WHO
INTERESTED GPU USE ON
OPENSTACK ENVIRONMENT
FOR OUR CUSTOMERS



#### NOW EVALUATING...

TESLA M60+DELLEMC
POWEREDGE
C4130+OPENSTACK
TO EVALUATE GPU ON
OPENSTACK ENVIRONMENT
THANKS TWO HELPING
OUT AND PROVIDING
THOSE SERVERS+CARDS







# 日本仮想化技術、OpenStack環境でのGPU技術の活用「GPU on OpenStack」を推進

報道関係者各位 プレスリリース 2017年7月18日 日本仮想化技術株式会社

オープンソースのクラウド技術である"OpenStack"の導入支援やコンサルティングを手掛ける日本仮想化技術株式会社 (本社:東京都渋谷区、代表取締役社長:宮原 徹) は、OpenStack環境でのGPU技術の活用「GPU on OpenStack」を 推進いたします。

「GPU on OpenStack」は、OpenStack環境の仮想マシン/コンテナに対してGPUの計算リソースの割り当てを行い、ハイパフォーマンスコンピューティングやデスクトップ仮想化などの様々な用途で利用できるインフラ技術です。「GPU on OpenStack」を使用したOpenStack環境を用意することで、GPUを業務で活用する利用者に対して利便性の向上を実現し、インフラ管理者に対して一元管理による運用の効率化を支援いたします。「GPU on OpenStack」の環境構築や利用方法についての情報は十分でなく、日本仮想化技術ではエヌビディア合同会社(以下、NVIDIA)と株式会社アスク(以下、アスク)の協力により環境構築と技術評価を行い、利用者やインフラ管理者にとって役立つ情報を発信していきます。日本仮想化技術およびこの取り組みに賛同いただける企業と共に「GPU on OpenStack」の導入支援やコンサルティングサービスを推進していきます。

国内最大規模のオープンクラウドイベントである「OpenStack Days Tokyo」(7月20日~7月21日、虎ノ門ヒルズフォーラム)ではNVIDIA/アスクの展示ブースでデモンストレーションを用意しております。弊社主催イベントである「OpenStack最新情報セミナー」(7月26日、渋谷)ではOpenStack環境でのGPU技術の活用をテーマとした技術セミナーを企画いたします。

OpenStack最新情報セミナー https://openstack-update.compass.com

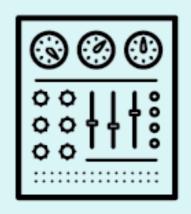
はじめに

#### なぜ GPU ON OPENSTACKなのか?



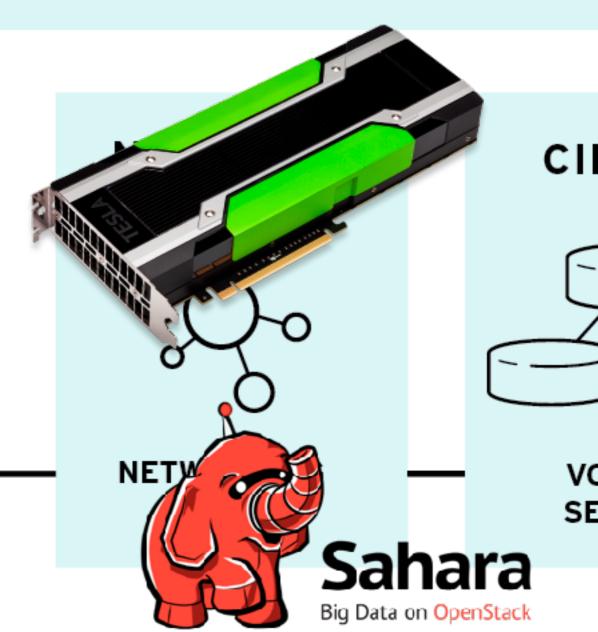
## OPENSTACK 特殊利用の需要

- OpenStackの特殊な利用の需要が 増えてきている
  - Hadoop(Sahara),HPC,などなど
- ほとんどのものはまとめたのが なく、ググって調べるしかない...
  - 曰く「ドキュメントロスト状態」
- これらのものはOpenStackの
  Docsにまとめられるべき



#### HORIZON

DASHBOARD



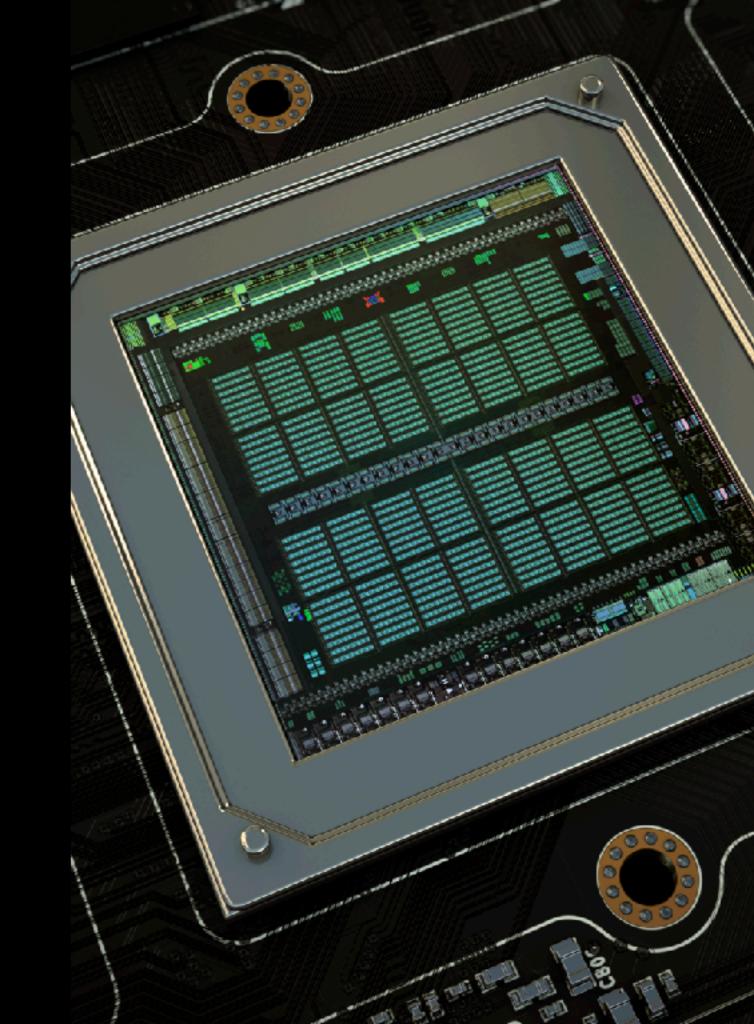
'GPU ON OPENSTACK'とは何か?

どうやってGPUはOPENSTACK環境で動作するのか?



## 'GPU'トレンド のおさらい

- 多くのGPUコアを利用する
  - 一部の計算は単一MPUユニットが 小さくスピードが遅かろうとも多く のMPUコアを使うことはいくつか の計算処理には有効である
  - サーバ各々は高パフォーマンスでコンパクトにできる
- 低消費電力化はHPCエンドユーザにとって大事なこと
  - 多くの省電力・ハイパフォーマンス サーバを所有しやすくなる



## どうやってGPU は動くのか?

- 現在のところPCIパススルーあるいは NVIDIA(GPU) dockerでの利用となる
- 'PCIパススルー' は土台となるベアメタル仮 想環境に依存する
  - VSphereとXenは各々VMにGPUコアを任意の単位 で分割して割り当てすることが可能
  - OpenStackの標準的なベアメタル仮想環境として 使われるKVMではコア分割ができず、GPUユニットごとをVMに割り当てることとなる
- 昨今使われるコンテナ仮想環境でのデプロイとなるNVIDIA(GPU) DockerはDockerコンテナ同士で一つのGPUユニットを共有しているが、明示的な分割はしない
  - ウインドウズは'docker vm'としては動作しない



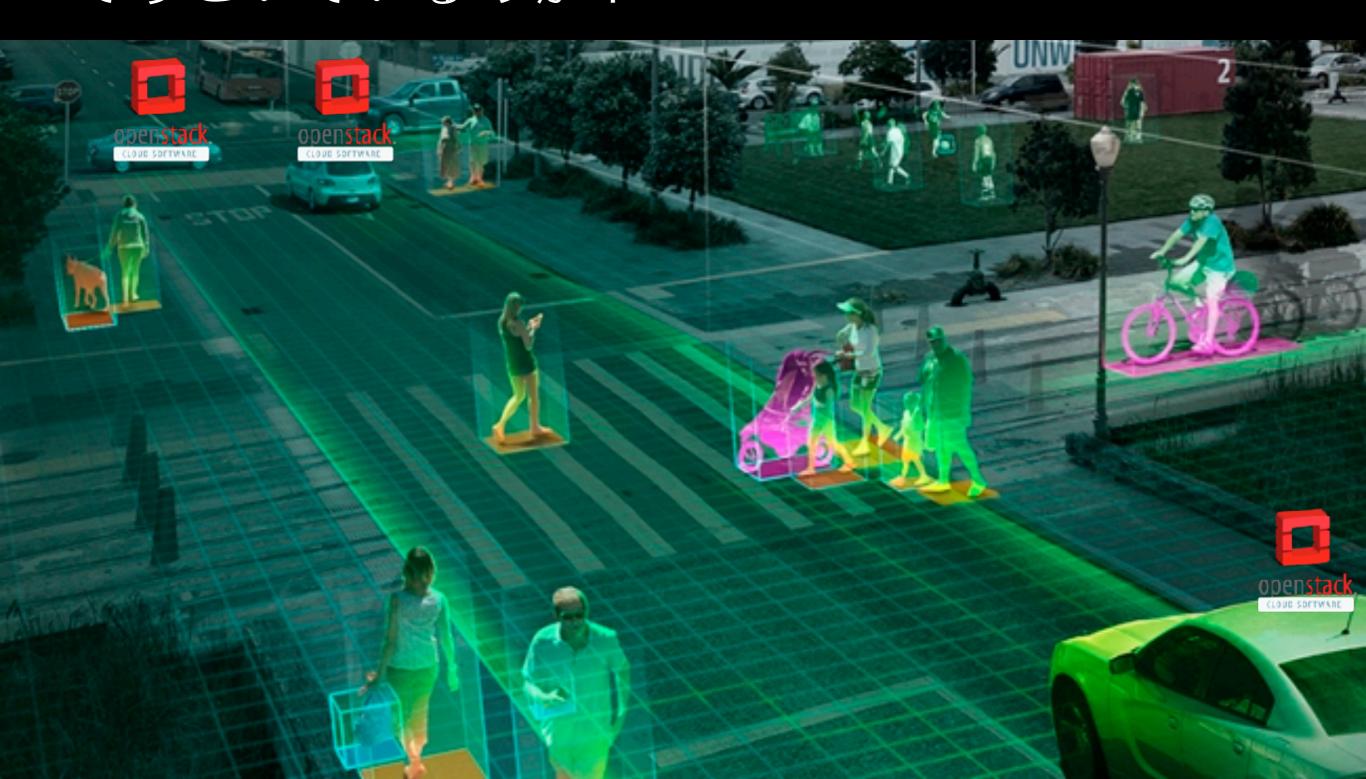
## GPU OPENSTACKとは

- インスタントなHPC利用
  - いくつかの計算をしおえたらVMその ものを壊す。
  - ちょっとしたお試しでいくつかのVM を使ってHPCグリッド試してみる。終 わったらすぐ壊す。
- GPUインターナルクラウドとしてのGPU利用
  - 主に製造業において、いくつかのシス テムは情報管理上、パブリッククラウ ドに外だしができない。



SETUP: GPU ON OPENSTACK

どんなGPUのメカニズムがOPENSTACK環境でうごいているのか?



#### OpenStackでGPUを動かす方法 (現在)

- PCIパススルー
  - PCIデバイスをダイレクトに接続する
  - ComputeNodeのハイパーバイザー依存、OpenStack依存ではない
  - Xen利用でのGPUコア分割、KVMはNVIDIA/AMDともコア分割不可
  - Intel GVT-g(Xen)/GVT-d(KVM)によるIntel GPUのコア分割

#### コンテナ

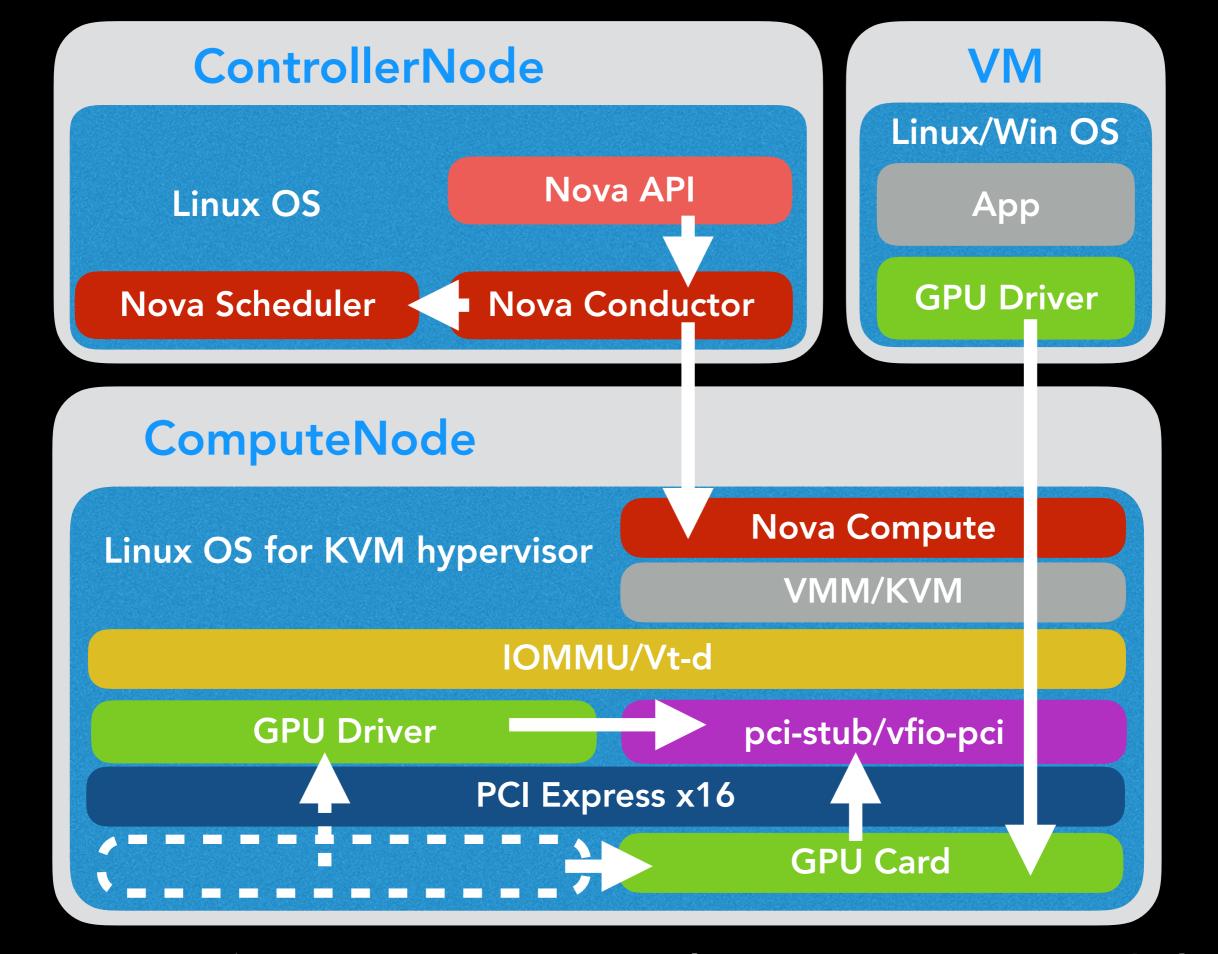
- NVIDIA Dockerの利用
- 複数のコンテナによってGPUを利用できるが、明示的なGPUコア分割はできない タスク次第でのGPU利用
- Kubernetes/Mesos/Docker Swarmなどとの組み合わせ管理

#### PCIパススルーはVM上でどう動作するのか?

- PCIデバイスをダイレクトにLinuxホストを通じて接続する
  - 物理ホストよりデバイスを切り離す必要あり
    - (NVIDIA) Dockerと違いホスト切り離しとなるため、監視管理等は各VMに必要
  - OpenStack依存のことではなく、Linuxのベアメタル環境に依存
- KVM上では一つのGPUデバイスにひとつのVM
  - Docker/Xen/VSphereと違いGPUは分割もできなければ共有もできない
  - これはKVMの制限であってOpenStackによるものではない

### PCI Passthrough on OpenStack

- Redhatは公式サポート
  - ただ、あまり推しでない気が...
- Ubuntuはドキュメントすらまともにない...
  - ググって探すしかない(さすがUbuntu..orz)
  - 目下小職/NVIDIA JAPAN/DELLEMC/VTJで再度細く検証と OpenStackコミュニティのためにドキュメントをまとめ上げる予 定(趣旨に同意いただき参画いただける企業様歓迎いたします)



#### Step1:ホストのGPUの状態をまずは調べる

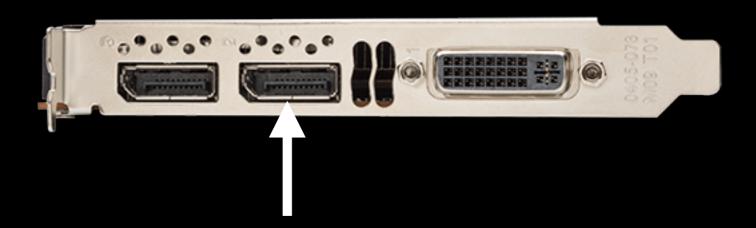
● Ispci -nn | grep -i nvidia でGPUの状態をまずはチェック

Ispci -nn | grep -i nvidia

88:00.0 VGA compatible controller [0300]: NVIDIA Corporation Device [10de:11b4] (rev a1) 88:00.1 Audio device [0403]: NVIDIA Corporation GK104 HDMI Audio Controller [10de:0e0a] (rev a1)

- 全てのGPUユニットがパススルーされている必要がある
  - GPUだけでなく、ユニットとして認識されるGPUデバイス自体も パススルーする必要あり
  - でないとVMは動いてくれない.. (完全にパススルーにならない)

#### GPUポートがパススルーされているかチェック



- QuadroなどHDMIはビデオポートだけでなくオーディ オポートももっているため注意
  - Ispciでチェックされたものについては全てパススルーさせる必要がある。

#### STEP2:IOMMUセットアップ

- IOMMU(Input/Output Memory Management Unit) は物理デ バイスを仮想化システムで使う上で必要なもの
- もちろんvt-dはオンにする必要あり (EFI/BIOSののデフォルトはON)
- intel\_iommuとvfio\_iommu\_type1.allow\_unsafe\_interruptsの 設定を/etc/default/grubに行う必要あり

GRUB\_CMDLINE\_LINUX\_DEFAULT="quiet splash intel\_iommu=on vfio\_iommu\_type1.allow\_unsafe\_interrupts=1"

#### STEP3:pci-stubとVFIO

- pci-stubは物理デバイスをLinuxホスト側が利用できないようにする
- VFIO(Virtual Function IO) は pci-stubと同様の働きをする。kernel 4.1以 降のサポート
  - 未使用時はデバイスをD3ステート(低消費電力モード)に変更する
- It is not used by default デフォルトでは使われないため /etc/module を編集して、これらと関連するコンポーネンツを追記する (kvm,kvm\_intel)

```
pci_stub
vfio
vfio_iommu_type1
vfio_pci
kvm
kvm_intel
```

### STEP4-1:ブラックリスト(1)

• ramfsでGPUデバイスを認識できないようにする

/etc/initramfs-tools/modules to initramfs (ubuntuの場合)

echo 'pci\_stub ids=10de:11b4,10de:0e0a' >> /etc/initramfs-tools/modules sudo update-initramfs -u && sudo reboot

#### STEP4-2:Blacklist(2)

ブート時にGPUデバイスを認識できないようにする。

/etc/modprobe.d/blacklist.conf に次を追加:

blacklist nvidia blacklist nvidia-uvm

ドライバについてもブラックリストに入れる必要あり

blacklist nouveau

#### STEP5:物理からのアンバインド

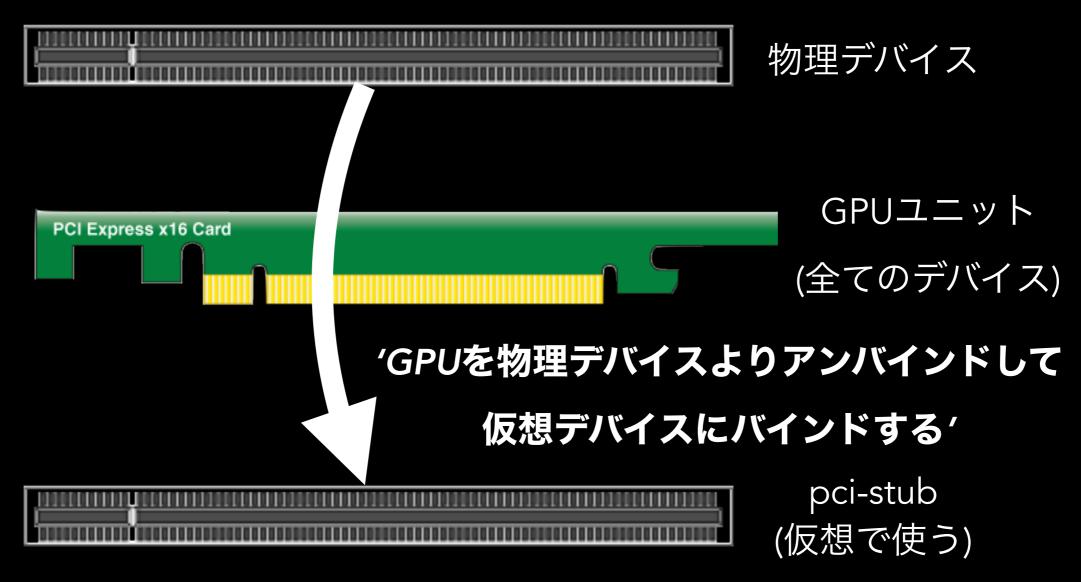
- 『物理ホストから切り離しVMへのバインドするため』にpci-stubを チェックする
  - 1. pci-stub/new\_idにパススルーするドライバのPCIIDをエントリする。
  - 2. 関連するPCI識別子を物理ホストから切り離す
  - 3. pci-stubにこのPCI識別子をバインドする。

```
echo 11de 11b4 > /sys/bus/pci/drivers/pci-stub/new_id
echo 11de 0e0a > /sys/bus/pci/drivers/pci-stub/new_id
echo 0000:88:00.0 > /sys/bus/pci/devices/0000:88:00.0/driver/unbind
echo 0000:88:00.1 > /sys/bus/pci/devices/0000:88:00.1/driver/unbind
echo 0000:88:00.0 > /sys/bus/pci/drivers/pci-stub/bind
echo 0000:88:00.1 > /sys/bus/pci/drivers/pci-stub/bind
```

物理ホストのdmesgに以下の'claimed'がブートプロセスにあるかどう か確認する。

pci-stub 0000:88:00.1: claimed by stub

echo 0000:88:00.0 > /sys/bus/pci/devices/0000:88:00.0/driver/unbind echo 0000:88:00.1 > /sys/bus/pci/devices/0000:88:00.1/driver/unbind



echo 11de 11b4 > /sys/bus/pci/drivers/pci-stub/new\_id echo 11de 0e0a > /sys/bus/pci/drivers/pci-stub/new\_id echo 0000:88:00.0 > /sys/bus/pci/drivers/pci-stub/bind echo 0000:88:00.1 > /sys/bus/pci/drivers/pci-stub/bind

図2:物理にあるGPUを切り離し仮想につなぐ

#### **UEFI/BIOS** Vt-d



**GRUB** /etc/default/grub

**IOMMU** 

ramfs /etc/initramfs-tools/modules **BLACK** LIST

modules /etc/modules **BLACK LIST** 

**IOMMU** 

modprobe /etc/modprobe.d/blacklist.conf **BLACK** LIST

**BLACK** 

**LIST** 

pci-stub /sys/bus/pci/drivers/pci-stub/ /sys/bus/pci/devices/\$(Identifier)/driver/unbind

図3:ブート時のGPUブラックリストのプロセス(Ubuntuの場合)

### GPUを追加する(1)

Ispci の結果を確認-2つの関連するPCI識別子が確認で きるはず(識別子番号は使用するシステムに依存)

```
Ispci -nn | grep -i nvidia
88:00.0 VGA compatible controller [0300]: NVIDIA Corporation Device [10de:11b4] (rev a1)
88:00.1 Audio device [0403]: NVIDIA Corporation GK104 HDMI Audio Controller [10de:0e0a] (rev a1)
84:00.0 VGA compatible controller [0300]: NVIDIA Corporation Device [10de:11b4] (rev a1)
84:00.1 Audio device [0403]: NVIDIA Corporation GK104 HDMI Audio Controller [10de:0e0a] (rev a1)
```

● pci-stubにさらにパススルーするGPUの追記をする。

```
echo 0000:84:00.0 > /sys/bus/pci/devices/0000:84:00.0/driver/unbind echo 0000:84:00.1 > /sys/bus/pci/devices/0000:84:00.1/driver/unbind echo 0000:84:00.0 > /sys/bus/pci/drivers/pci-stub/bind echo 0000:84:00.1 > /sys/bus/pci/drivers/pci-stub/bind
```

#### GPUを追加する(2)

(rev a1)

追加のGPUがうまく動いたかVMで確認

```
ubuntu@guestos$ Ispci -nn | grep -i nvidia 00:07.0 VGA compatible controller [0300]: NVIDIA Corporation GK104GL [Quadro K4200] [10de:11b4]
```

00:08.0 VGA compatible controller [0300]: NVIDIA Corporation GK104GL [Quadro K4200] [10de:11b4] (rev a1)

• アプリによっては両方同時に利用する際など、同一のGPUである必要がある可能性がある。

/nbody-benchmark-numdevices=2-num bodies=65536

## Openstack:nova-apiの設定

- ControllerNodeの/etc/nova/nova.confを編集しnovaapiを再起動する
  - pci\_aliasにPCIデバイスの情報、エイリアス名を記述

pci\_alias={"name":"K4200","vendor\_id":"10de","product\_id":"11b4"}

## OpenStack:nova-computeの設定

- ComputeNodeにある/etc/nova/nova.confを編集し、 nova-computeを再起動する。
  - pci\_passthrough\_whitelistにPCIデバイスの情報、エイリアス名を記述

```
pci_passthrough_whitelist={"name":"K4200","vendor_id":"10de","product_id":"11b4"}
```

\*このケースの場合、ベンダーIDとプロダクトIDが適合したデバイスは全てVMにパススルーする。

• pci\_aliasににPCIデバイスの情報、エイリアス名を同様に追記 \*Neuton以降

```
pci_alias={"name":"K4200","vendor_id":"10de","product_id":"11b4"}
```

#### OpenStack:nova-schedulerの設定

- ControllerNodeの/etc/nova/nova.confを設定しnovaschedulerを再起動する。
  - PciPassthroughFilterを利用できるようにするために
     PciPassthroughFilterをscheduler\_default\_filtersに追記する
  - 同様にPciPassthroughFilterをscheduler\_available\_filtersに記述 する

scheduler\_available\_filters=nova.scheduler.filters.all\_filters
scheduler\_available\_filters=nova.scheduler.filters.pci\_passthrough\_filter.<u>PciPassthroughFilter</u>
scheduler\_default\_filters=DifferentHostFilter,RetryFilter,AvailabilityZoneFilter,RamFilter,CoreFilter,DiskFilter,ComputeFilter,ComputeCapabilitiesFilter,ImagePropertiesFilter,ServerGroupAntiAffinityFilter,
ServerGroupAffinityFilter,AggregateInstanceExtraSpecsFilter,<u>PciPassthroughFilter</u>

#### ControllerNode

Linux OS

scheduler\_default\_filters

scheduler\_available\_filters

Nova Scheduler

pcipassthroughfilterを利用してGPUパススルー

されたComputeNodeを選ぶ

pci\_alias

**Nova API** 

PCIデバイスを利用可にしてPCI利用のリクエスト を送れるようにする

**Nova Conductor** 

VI

Linux/W

Ap

**GPU D** 

#### ComputeNode

**Linux OS for KVM hypervisor** 

pci\_passthrough\_whitelist pci\_alias

**Nova Compute** 

GPUパススルーしたインスタンスをpci\_aliasと pci\_passthrough\_whitelistによって発生させる

図4:NovaがGPU(PCI)パススルーされているComputeNodeで動作するプロセス

## OpenStack:flavor-keyの設定

- flavor-keyを設定しGPUインスタンスで利用できるよう にPCIパススルーの設定をflavorに追記する
  - pci\_passthrough:alias=\$(pci\_alias\_name):\$(the number of GPUs we would like to use)

nova flavor-key \$flavor\_name set "pci\_passthrough:alias"="K4200:\$(the number\_of\_gpus)"

既知の問題

GPU ON OPENSTACKを使うにあたって注意すべき問題点



#### Cloudイメージの問題

- CloudイメージはGPUを使う上ではとても小さくqemu-imgでリサイズ する必要がある
- CUDAドライバはperl-packages(dev packages)がインストール時に必要
  - それが.debあるいは.rpmパッケージであろうとインストールが必要になる。なぜなら CUDAパッケージ自体がバイナリパッケージでなく、ソースコードよりバイナリをビ ルドしていて、makeをパッケージインストールの際に実行している
- NVIDIA日くCUDAドライバの次期リリースでFIX予定とのこと
  - CUDA 7.6以降でfixの予定だったのに..まだfixされてない..orz

#### VDIとしてのWindows利用

- CUDA on Windows は思ったより早いけどカクカクしてしまう
- 多分DISKのスピード、ネットワークなどなどいろいろなものが予想される。多 分工フェメラルモードやその他爆速系のSSD/NVMeの利用やら10g以上のネット ワーク環境が必要となるだろう
  - まだ改善対応をやってみたことはないが、なぜ発生するのかを調べるべく近々検証予定。
  - VMは基本メモリ/NW転送などコンテキストスイッチでの動作をするため、CUDAの重い ワークロードはカクカクする可能性があるかもしれない。
  - 詳しくはデモビデオにて...
- GPU on OpenStack上でのWindowsの動作はもっと調査が必要..時間がほしい...

#### VDIとしてのWindows利用(feedback)

- Thanks giving some feedbacks to my session at LC3 China!
- Should checked and investigate the Windows-related issues below, will update later.
  - Windows 10 on KVM issue
    - http://bart.vanhauwaert.org/hints/installing-win10-on-KVM.html
      - Windows 10 deployment is succeeded from my KVM but I still have failed the same deployment from my OpenStack.I should investigate more what happened in detail.
  - Nvidia Driver issue (version 337.88 or later)
    - https://wiki.archlinux.org/index.php/PCI\_passthrough\_via\_OVMF#.
       22Error\_43:\_Driver\_failed\_to\_load.22\_on\_Nvidia\_GPUs\_passed\_to\_Windows\_VMs
      - '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. 'I haven't found this issue on my
        Windows 7 VMs so I should check more in detail
      - Related links libvirt for adding the driver, should be checked.
        - https://github.com/openstack/nova/blob/master/nova/virt/libvirt/config.py#L2025-L2036
        - https://github.com/openstack/nova/blob/master/nova/virt/libvirt/driver.py#L4262-L4264

#### ライブマイグレーションの問題

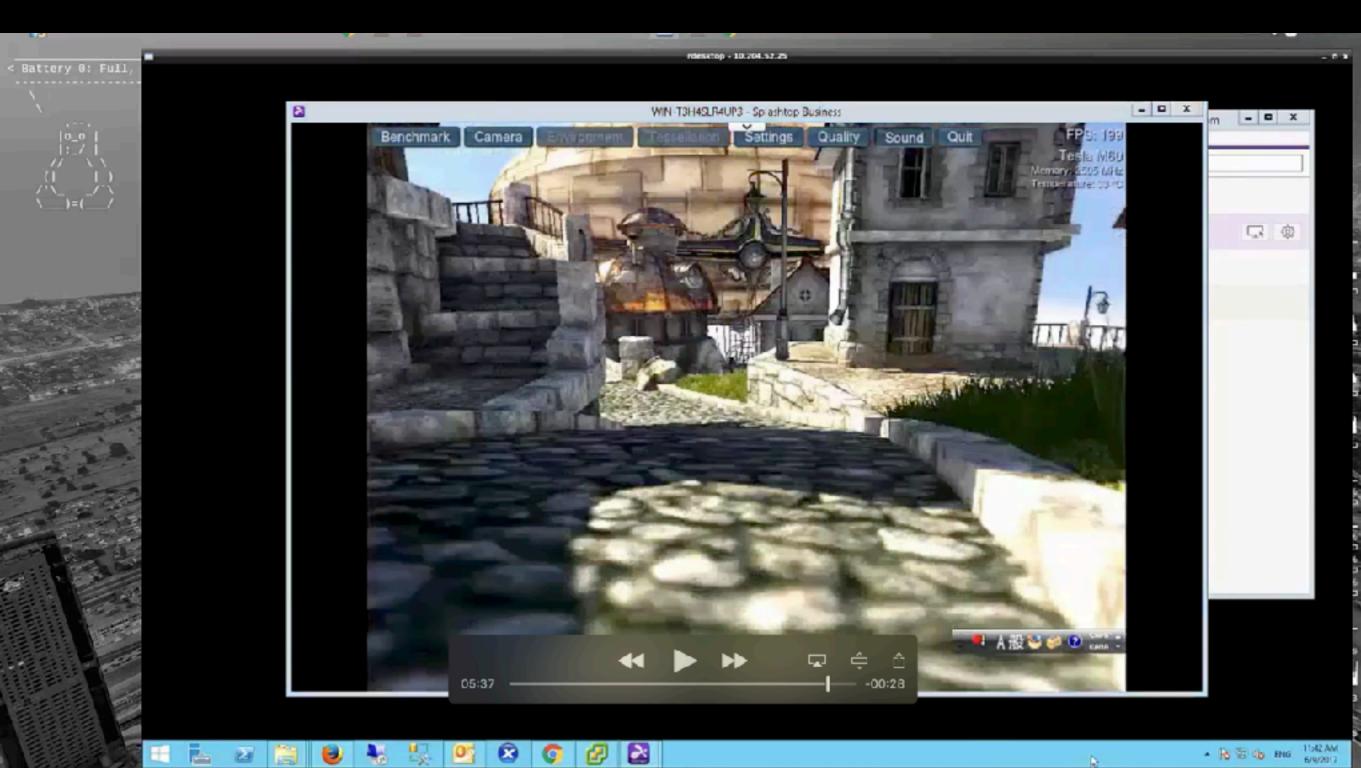
- PCIパススルーを使ったVMはライブマイグレーションができない。マイグレーション前のホストのPCIコネクションが切れないままになってしまう。
- ワークアラウンド:古いコネクションを下記にあるようなmysql DBの nova.pci\_devicesから削除し古いホストを再起動する。
  - 再起動は他のホストに影響するため、ありえない。他にmysql DBなど関連プロセスの再起動で乗り切る形もあるが、これも他のホストに影響するためありえないことになる。

2016-08-11 00:54:45   2016-08-19 04:58:01   NULL	0   45	21   0000:84:00.0   11b4	10de	type-PCI   pci_0000_84_00_0
label_10de_11b4   available   {}   NULL	NULL	1 << old-host		
2016-08-11 00:54:45   2016-08-19 04:58:01   NULL	0   48	21   0000:88:00.0   11b4	10de	type-PCI   pci_0000_88_00_0
label_10de_11b4   available   {}   NULL	NULL	1   << old-host		

デモ(ビデオ)

#### GPUがOPENSTACK環境でどう動くか

チェック!



#### インスタンスを立ち上げGPUの動作を確認

- RHEL OpenStackを利用
- パススルーしてあるインスタンスを立ち上げる
- Ubuntuインスタンスを立ち上げIspciとdevicequeryを実行、GPUの動作を確認
- Windows上でのGPUの動きをGraphicsベンチで確認、RDSでのリ モートサイトにあるものの接続のため、カクカクしてしまうがベン チの結果だけ?良いことを確認する。
  - デモビデオではSlashtopを利用

#### 関連リンク

- Attaching physical PCI devices to guests: https://docs.openstack.org/admin-guide/compute-pci-passthrough.html
- Container as a Service on GPU Cloud- Our Decision Among K8s, Mesos, Docker Swarm, and OpenStack Zun:
  - https://www.slideshare.net/secret/AiUdO4dLxNTkfl
- OVMF\_による\_PCI\_パススルー https://goo.gl/icq9mV

#### Special Thanks to:

GPU on OpenStack project members
 VirtualTech Japan
 NVIDIA
 DellEMC
 NEC Networks & System Integration

- @zgock999 at Tokaido-LUG, Nagoya, Japan Teach me some hints how to use GPGPU on VM!
- Matthew Treinish of IBM attended my session at LC3 china and figure out and feedback some point!
- Our customers! give the chance to evaluate!

