


# Is Sharing GPU to Multiple Containers feasible?

Samed Güner, SAP Artificial Intelligence  
August 18, 2020

## Is sharing GPU to multiple containers feasible? #52757

 Open tianshapjq opened this issue on 20 Sep 2017 · 67 comments



tianshapjq commented on 20 Sep 2017

Contributor + 😊 ...

Is this a **BUG REPORT** or **FEATURE REQUEST**?: feature request  
/kind feature


### What happened:

As far, we do not support sharing GPU to multiple containers, one GPU can only be assigned to one container at a time. But we do have some requirements on achieving this, is it feasible that we manage GPU just like CPU or memory?

### What you expected to happen:

sharing GPU to multiple containers just like CPU and memory.

 137  7  34  4  2

 k8s-github-robot added the **needs-sig** label on 20 Sep 2017



tianshapjq commented on 20 Sep 2017

Contributor Author + 😊 ...

@vishh @cmluciano is it workable?



tbchj commented on 20 Sep 2017

+ 😊 ...

+1



jianzhangbjz commented on 20 Sep 2017

Contributor + 😊 ...

/cc

# About me

- Developer at SAP Artificial Intelligence
- Contributor to Terraform and Cloud Foundry
- Living in Germany, having Turkish roots



# Open Source Projects

Contributions and Commitments from SAP



**CLOUD NATIVE**  
COMPUTING FOUNDATION

Platinum Member



**ECLIPSE**<sup>™</sup>  
FOUNDATION

Strategic Member



**CLOUD FOUNDRY**

Platinum Member



**kubernetes**



Gardener



Kyma



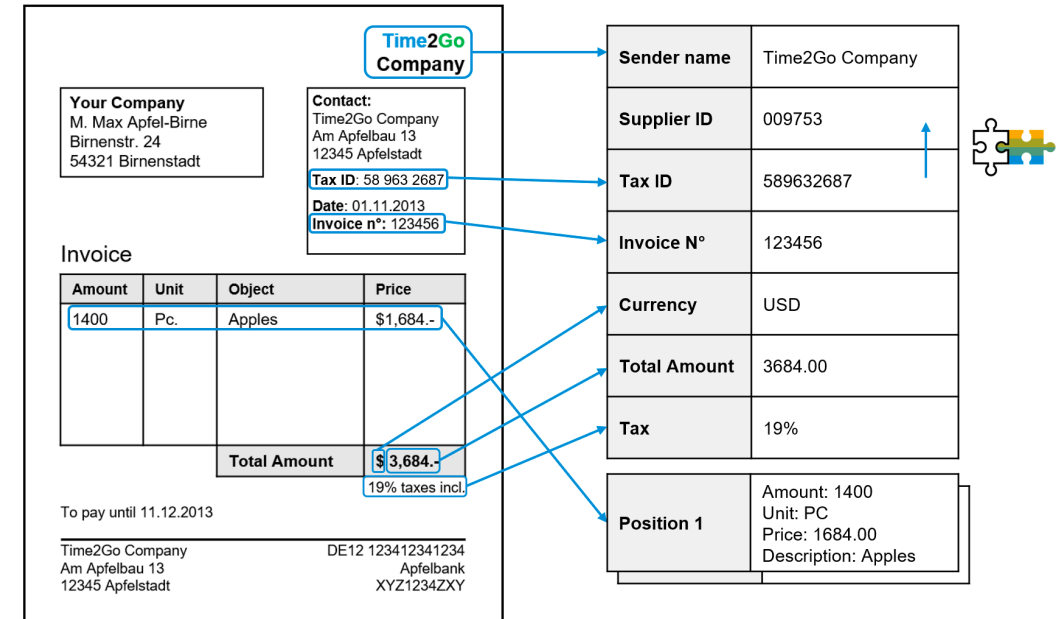
Luigi



...

# Artificial Intelligence at SAP

- Providing AI platform and services for enterprises and SAP applications
- Embedding intelligence into enterprises
- Large number of models in production



# ML on GPUs in K8s

- Expensive workload
- GPU utilization as challenge
- Sparse and limited availability

Fictional Calculation:

$$u = 70\% = 0,7 \quad n = 500 \quad p = \$3,06 \frac{1}{h}^1$$

$$C(n) = (1 - u) * n * p$$

$$= \left( 0,3 * 500 * \$3,06 \frac{1}{h} \right) * 768 h$$

$$= \$352.512 \text{ loss per month}$$



<sup>1</sup> averaged price from AWS, Azure and GCP, smallest setup with one Nvidia Tesla V100

**So we asked ourselves ..**

Is sharing GPU to multiple containers feasible? #52757

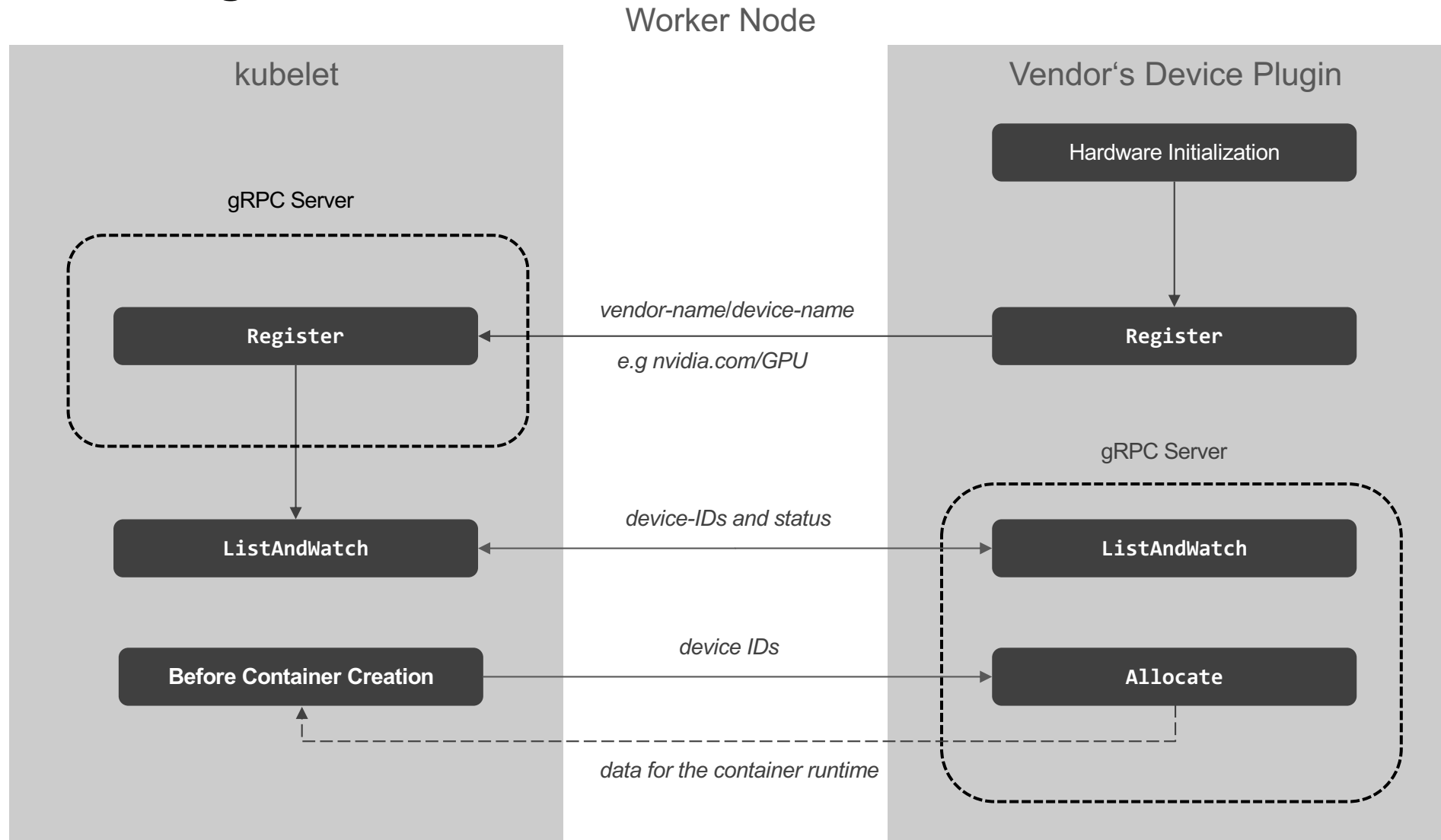


tianshapjq opened this issue on 20 Sep 2017 · 67 comments

**.. too<sup>1</sup>.**

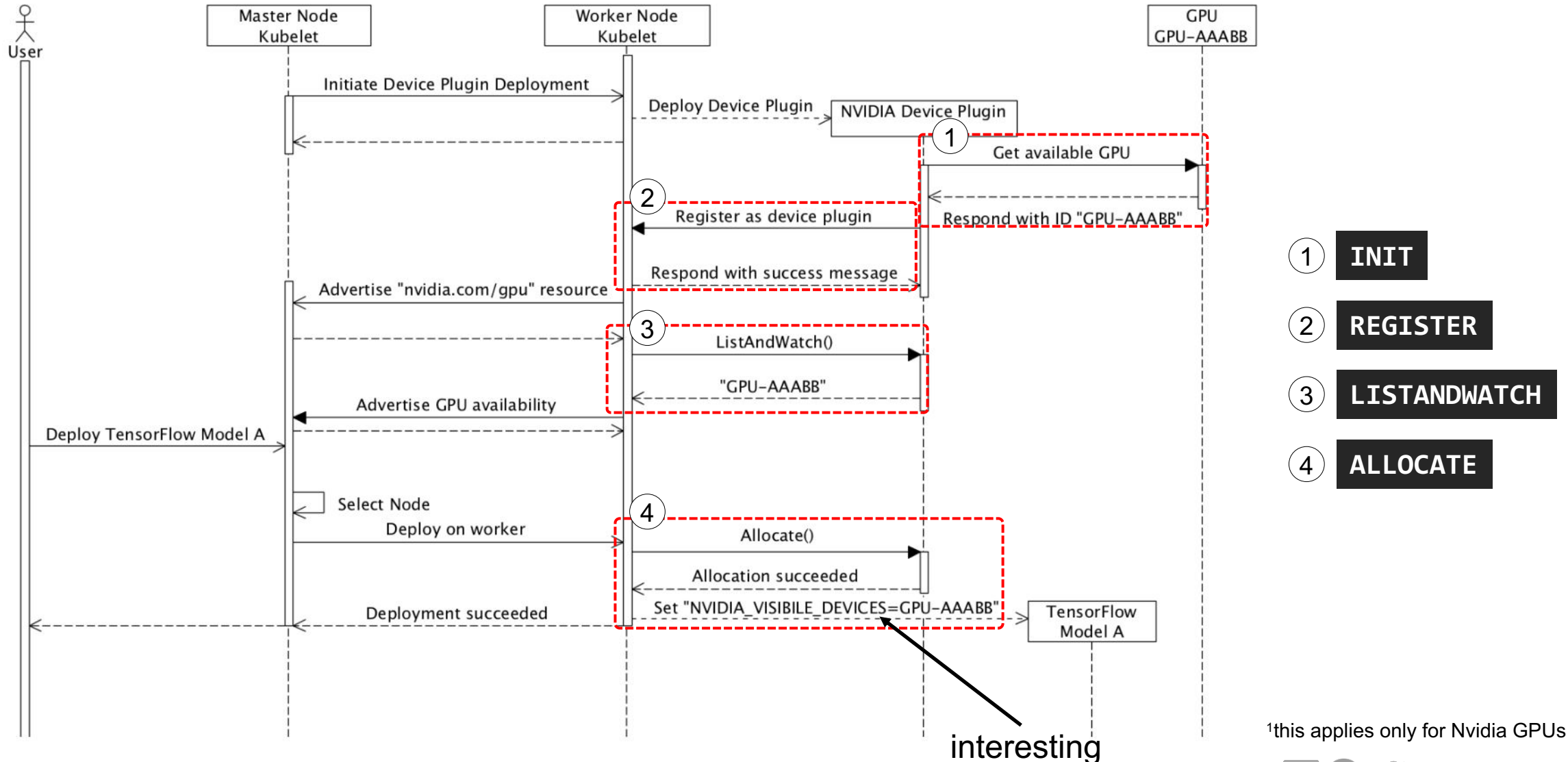
<sup>1</sup> target: Nvidia GPUs

# K8s Device Plugins<sup>1</sup>



<sup>1</sup> as of K8s 1.10

# So how are GPUs provisioned in K8s?<sup>1</sup>





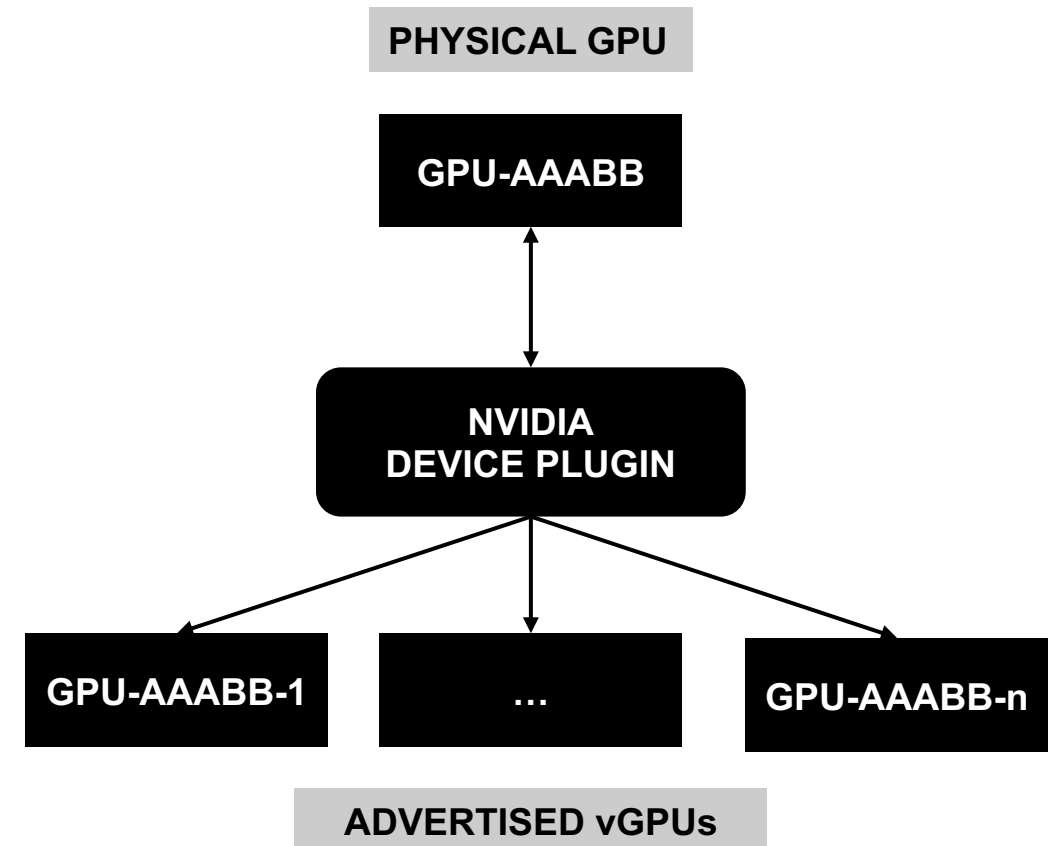
# So how did we share a GPU?

- Variant 1: Model Stuffing in Application Logic
  - no CGROUPS for cpu/memory
  - no ability to do request rate limiting per model
- Variant 2: Advertising more GPUs to K8s
  - extending the Nvidia Device Plugin for vGPUs
- Variant 3: NodeSelector Hacking
  - Privileged Container consume GPU directly
  - No Kubernetes Scheduling at all

# So how do we advertise more GPUs?

problems already start here

```
1  for k := 0; k < NUM_VGPU; k++ {
2      // generate vGPU ID
3      vGPUId := GPUtoVGPU(d.UUID, k)
4      // generate the vGPU
5      tmpVGPU := vGPU{vGPUId, 0, 0, 0,
6          pluginapi.Healthy}
7
8      // assign the vGPU to the physical GPU
9      gpus[d.UUID].vGPUs[vGPUId] = tmpVGPU
10     // make it public in ListAndWatch
11     devs = append(devs, &pluginapi.Device{
12         ID: tmpVGPU.UUID,
13         Health: pluginapi.Healthy,
14     })
15 }
```



## We asked ourselves ..

- How many models can we fit on a Nvidia K80s?
- How does the whole system behave?
- What are the trade-offs in doing so?

**.. so we did our experiments and collected data.**

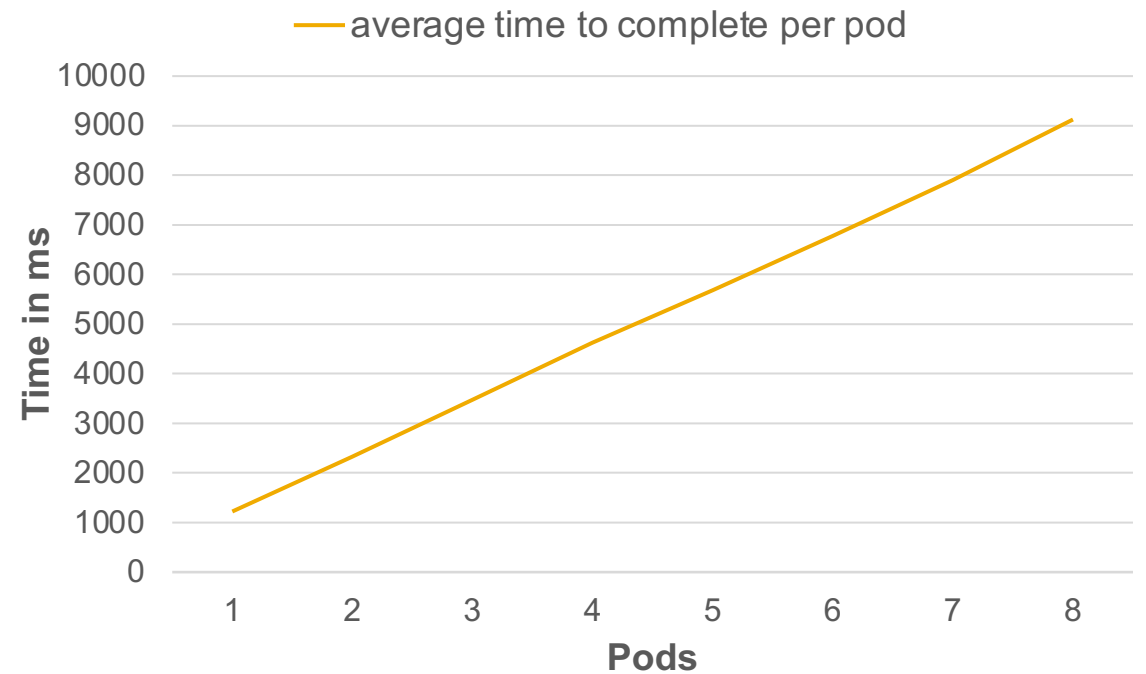
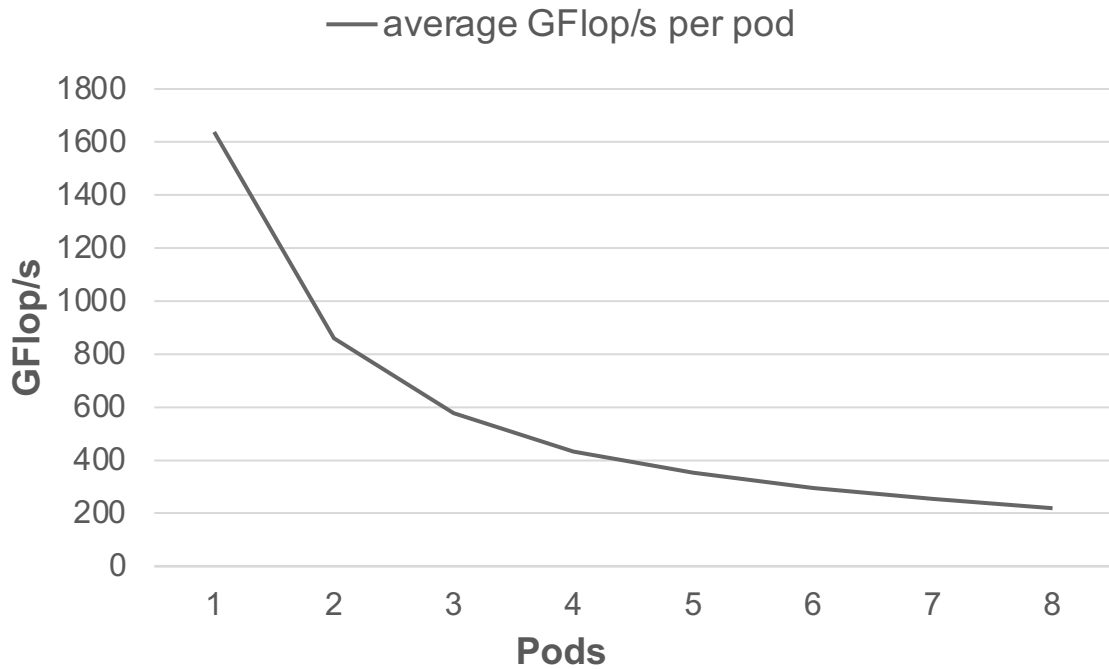
# Collecting Data from NVIDIA Device Plugin

- Device Plugins are K8s Resources!
- Collect Data with Nvidia Management Library (NVML)
  - Low Level Access to GPU
  - Go-C Bindings
  - Used by `nvidia-smi`
- Expose via Prometheus and Grafana



# Experiments<sup>1</sup>

nbody Simulation (n= 100096)

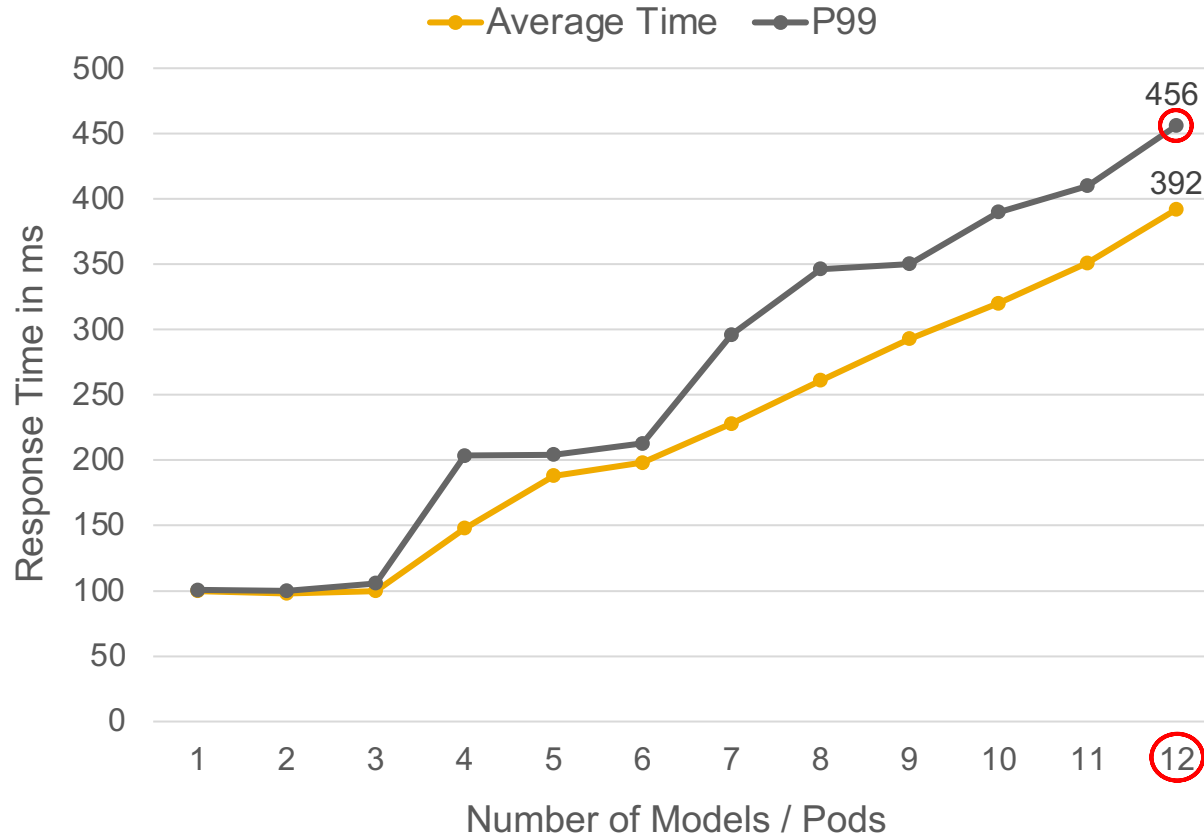


- Missing aspects for ML inference (but comparable to ML training)
- Fair-Scheduling (?)
- No GPU RAM Limiting & Fair Share of CPU

<sup>1</sup> for all experiments we use p2.xlarge with 1x Nvidia Tesla K80 (12 GB)

# Experiments<sup>1</sup>

ML Inference InceptionV3 (sequential requests)



- Low Throughput: 1 Request per Model
- 10,000 requests per Pod
- CPU Limit: 350m per Pod

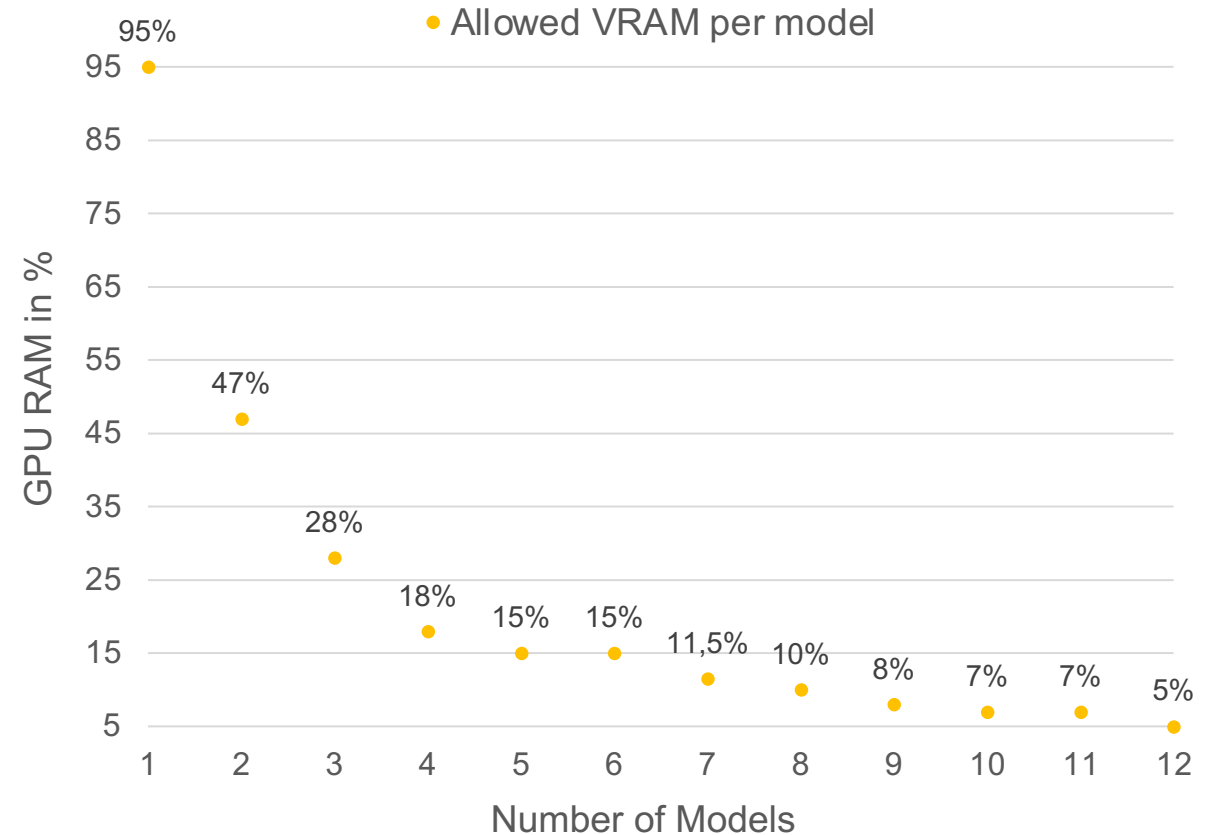
<sup>1</sup> for all experiments we use p2.xlarge with 1x Nvidia Tesla K80 (12 GB)

# Experiments<sup>1</sup>

## Dealing with GPU RAM for InceptionV3<sup>2</sup>

- Limiting Tensorflow Serving (TFS) Ram:
  - `tf.GPUOptions`
- Dynamic Scheduling:
  - Offer vGPUs until memory is full

- Can we go deeper?

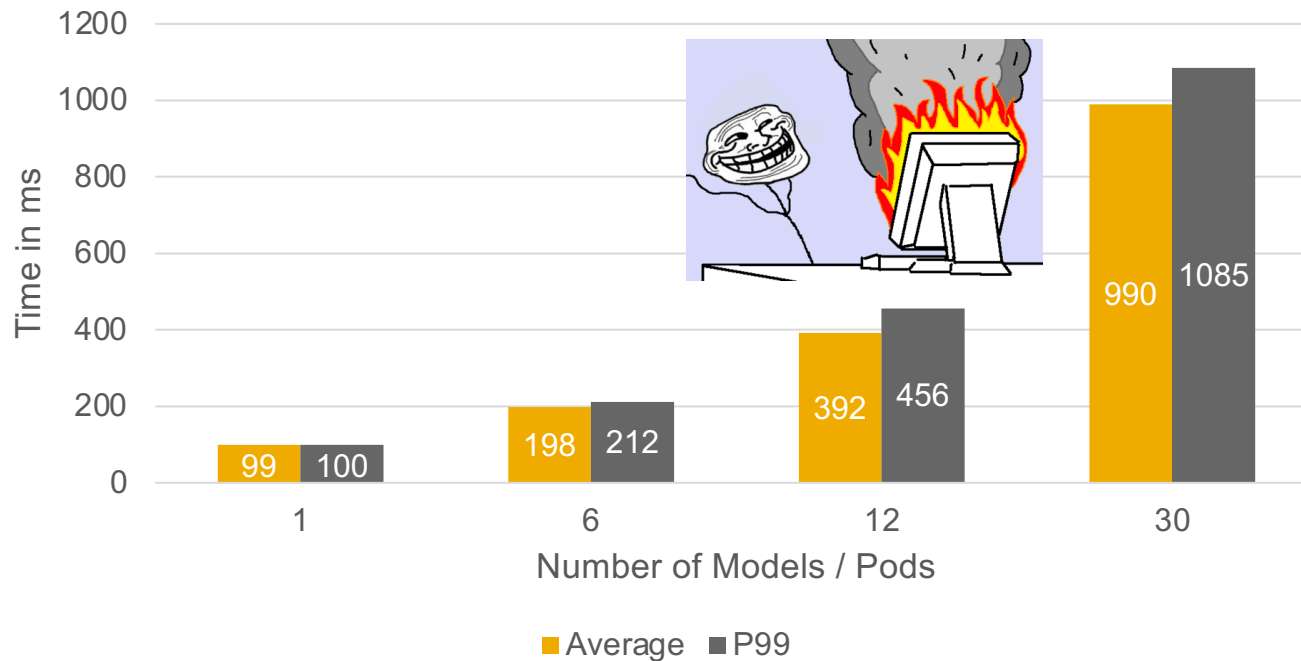


<sup>1</sup> for all experiments we use p2.xlarge with 1x Nvidia Tesla K80 (12 GB)

<sup>2</sup> only applicable for TensorFlow 1.x

# Experiments<sup>1</sup>

Dealing with GPU RAM for InceptionV3<sup>2</sup> (sequential requests)



- Minimum 228 mb per Model
  - up to 50 models on one K80
  - TFS will crash for less!
- 30 Models with each having:
  - 100m CPU Limit
  - 3% of VRAM
  - 1 Request at time

<sup>1</sup> for all experiments we use p2.xlarge with 1x Nvidia Tesla K80 (12 GB)

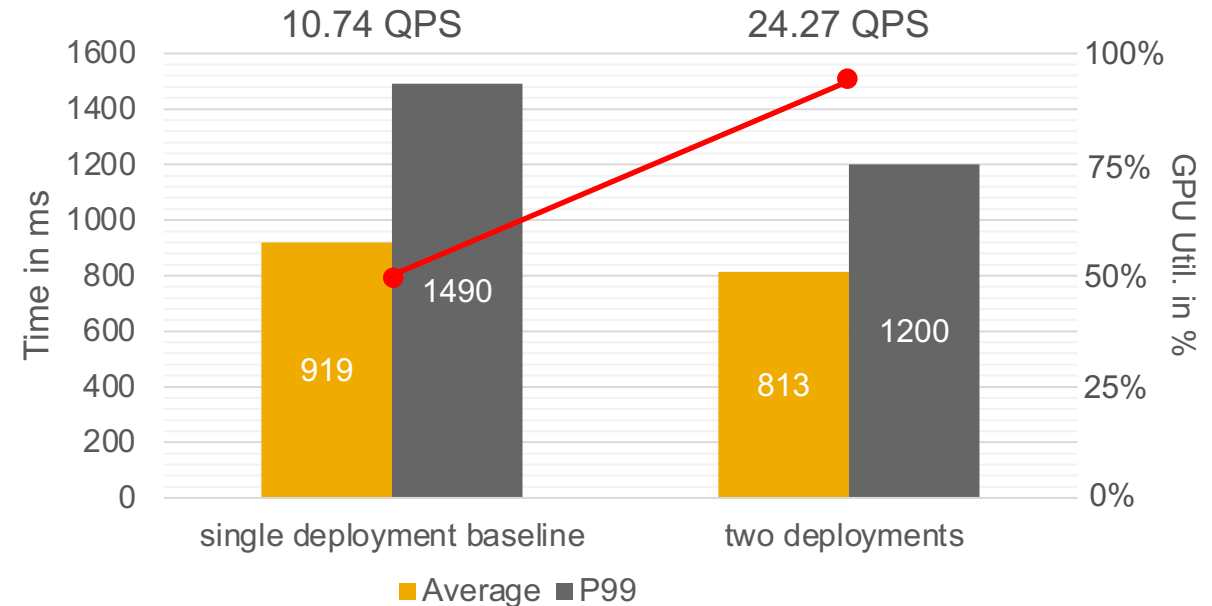
<sup>2</sup> only applicable for TensorFlow 1.x



# Experiments<sup>1</sup>

ML Inference InceptionV3 (parallel requests, no batching)

- Limits:
  - 10 parallel requests per Model
  - 350m CPU Limit
  - 3% of VRAM
- 4x worse than sequential request pattern
- Increase in throughput leads to increase in latency and GPU util.  
→ GPU has more work to do!



With batching enabled 2x latency and throughput improvement for 6.6x more VRAM.

<sup>1</sup> for all experiments we use p2.xlarge with 1x Nvidia Tesla K80 (12 GB)

<sup>2</sup> only applicable for TensorFlow 1.x

# Germans would say ..

**Jein<sup>1</sup>**

(yes, but with trade-offs and limitations)

<sup>1</sup>an answer between yes and no

# What does that all mean?<sup>1</sup>

## ML Inference

- Our solution is far away from production-ready
- GPU sharing is possible, if you try hard.
  - up to 30x cost saving<sup>2</sup>
  - Trade-Off between Throughput and Latency
- CPU limits as key factor in our setup
  - can be used to in-/decrease GPU Utilization, thus latency and throughput
  - allows deployment of multiple models with similar latencies while doubling throughput
  - GPU Util. overcommitment will increase latency

<sup>1</sup> for all experiments we use p2.xlarge with 1x Nvidia Tesla K80 (12 GB)

<sup>2</sup> K80 in combination with InveptionV3 capped at 228 MB VRAM and 100m CPU

# Limitations with our GPU Sharing

- No isolation guarantee on hardware level by Nvidia<sup>1</sup>
  - Multi-Tenancy setup is impossible
  - Still some sort of fair-scheduling of GPU duties
- No Low-Level API to specify GPU and VRAM
  - CPU Limits and TFS the only current way<sup>2</sup>
- “Scheduling” only by offering vGPUs from Device Plugin
- Resource Fragmentation

<sup>1</sup> for Nvidia Tesla K80

<sup>2</sup> without major customizations on CUDA level

# What do we need for GPU Sharing?

- Scheduling / isolation on GPU-level
  - CGROUP-like isolation of GPU resources
- Resource Defragmentation / Locality-Awareness
- Low-overhead during scheduling and processing
- Per device sub-resources (GPU Util. and VRAM) exposed at device plugin

# What does the community do?

- Deepomatic ([github.com/Deepomatic/shared-gpu-nvidia-k8s-device-plugin](https://github.com/Deepomatic/shared-gpu-nvidia-k8s-device-plugin))
  - Multiple vGPUs per GPU (cf. our implementation)
- Tencent's GPU-Manager (<https://github.com/tkestack/gpu-manager>)
  - Cuda-Core / VRAM Requests & Limits
  - Custom Wrapper around Nvidia Device Library
- KubeShare ([github.com/NTHU-LSALAB/KubeShare](https://github.com/NTHU-LSALAB/KubeShare))
  - Released June 2020
  - Set of custom controller / CRs
  - Sharing through CUDA Calls API interception
  - Powerful: Fine-grained allocation, isolation, locality-aware, low overhead

```
apiVersion: v1
kind: Pod
metadata:
  name: vcuda
  annotations:
    tencent.com/vcuda-core-limit: 50
spec:
  restartPolicy: Never
  containers:
  - image: <test-image>
    name: nvidia
    command:
    - /usr/local/nvidia/bin/nvidia-smi
    - pmon
    - -d
    - 10
    resources:
      requests:
        tencent.com/vcuda-core: 50
        tencent.com/vcuda-memory: 30
      limits:
        tencent.com/vcuda-core: 50
        tencent.com/vcuda-memory: 30
```

## KubeShare: A Framework to Manage GPUs as First-Class and Shared Resources in Container Cloud

Ting-An Yeh

National Tsing Hua University  
Hsinchu, Taiwan R.O.C.

Hung-Hsin Chen

National Tsing Hua University  
Hsinchu, Taiwan R.O.C.

Jerry Chou

National Tsing Hua University  
Hsinchu, Taiwan R.O.C.

### ABSTRACT

Container has emerged as a new technology in clouds to replace virtual machines (VM) for distributed applications deployment and

practice of DevOps [16] to achieve continuous software delivery, productivity, automation, cost saving, in the software delivery process. Therefore, containerization has been increasingly adapted and

# Thank you.

Contact information:

samed.guener@sap.com

   / samedguener