# INTRO

Eufrat Tsaqib

- Software Engineer at Nodeflux
- In-charge of Nodeflux Snapshot (Cloud)
- Just graduated from Physics major
- Latest research on Hyperspectral Imaging

Nodeflux

- Perusahaan berbasis kecerdasan buatan (AI) untuk melakukan analisis video dan gambar.
- Perusahaan pertama Indonesia di bidang Intelligent Video Analytics (IVA).
- Official NVIDIA Global Partner.

# CONTENTS

# WHAT IS A GPU

- Used to accelerate image or graphics computation
- Highly parallel mechanism
- Good to compute large matrix/tensor
- As an additional extension to computer
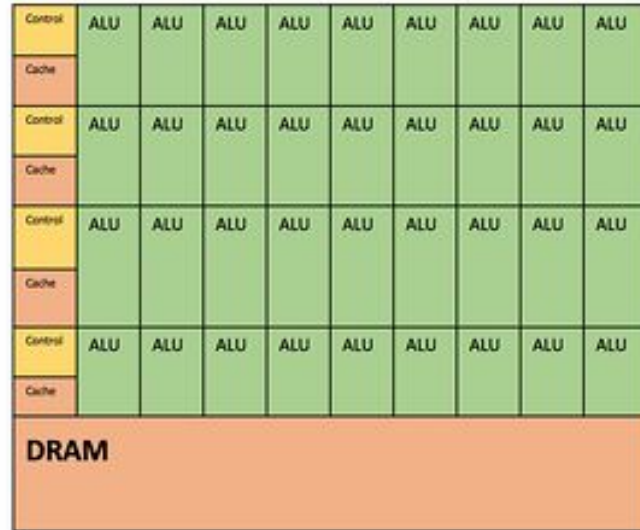- Scarce resource, especially one designed for ML (eg. NVIDIA Tesla, AMD Instinct)
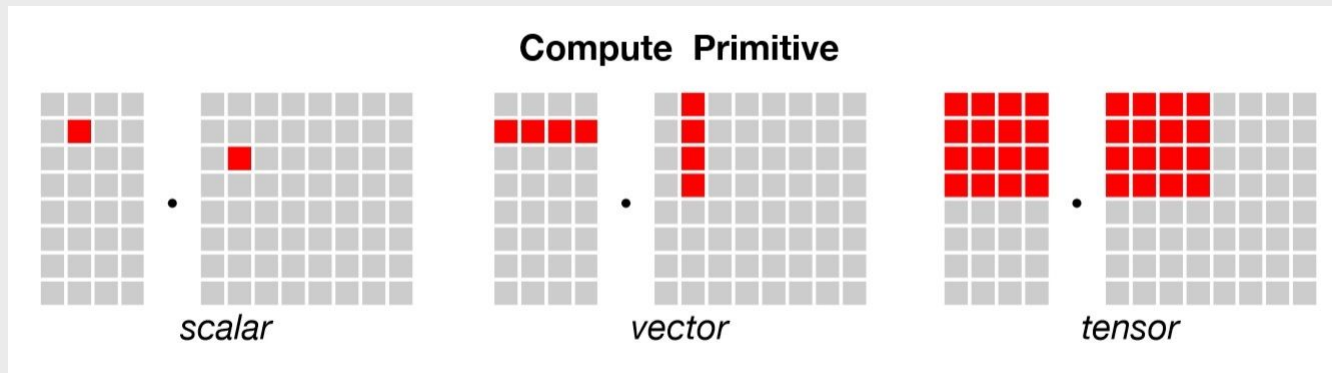
# WHAT IS A GPU

Highly parallel mechanism



| Control | ALU | ALU |
|---------|-----|-----|
|         | ALU | ALU |
| Cache   |     |     |
| DRAM    |     |     |

**CPU**

| Control / Cache | ALU | ALU | ALU | ALU | ALU | ALU | ALU | ALU | ALU |
|-----------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Control / Cache | ALU | ALU | ALU | ALU | ALU | ALU | ALU | ALU | ALU |
| Control / Cache | ALU | ALU | ALU | ALU | ALU | ALU | ALU | ALU | ALU |
| Control / Cache | ALU | ALU | ALU | ALU | ALU | ALU | ALU | ALU | ALU |
| DRAM            |     |     |     |     |     |     |     |     |     |

**GPU**

A GPU has more Arithmetic Logic Units (ALU) than a typical CPU.
- Increased ability to process simple operations in parallel

# WHAT IS A GPU

Good to compute large matrix/tensor



**Compute Primitive**

*scalar*     *vector*     *tensor*

CPU: 1 X 1 data unit
GPU: 1 X N data unit
TPU: N X N data unit

```
with tf.device('/gpu:0'):
  a = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], shape=[2, 3], name='a')
  b = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], shape=[3, 2], name='b')
c = tf.matmul(a, b)
```

# WHAT IS A GPU

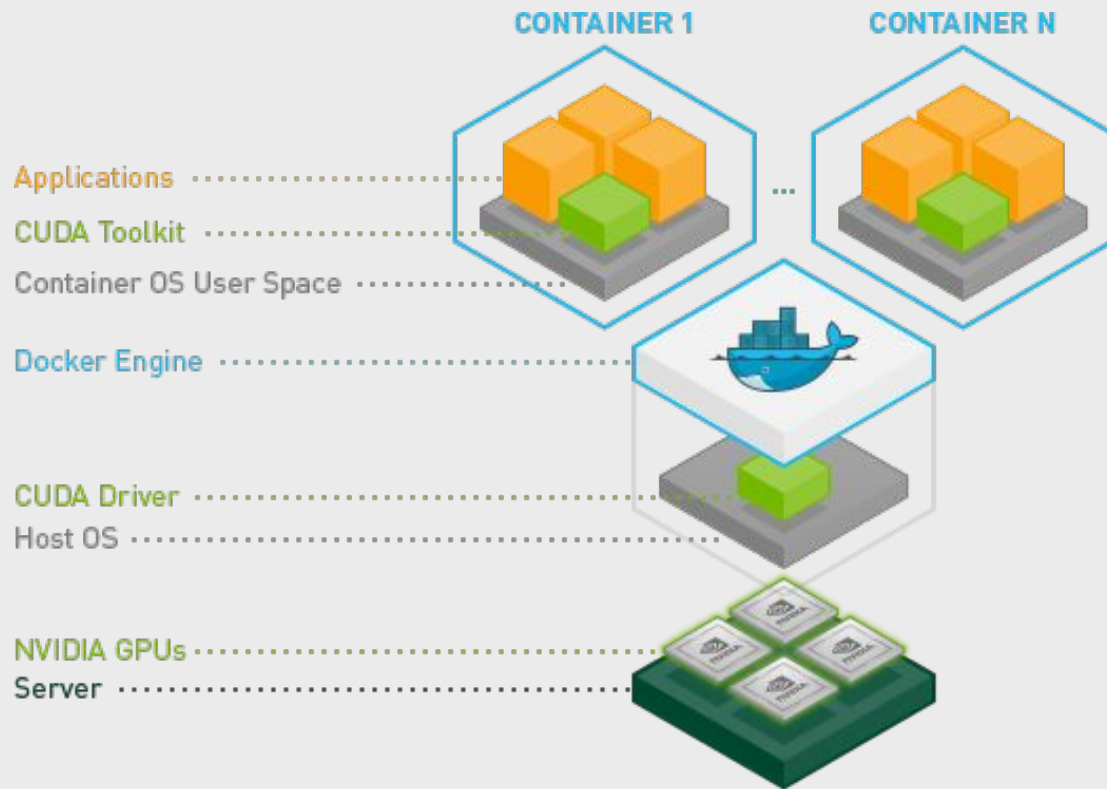**Scarce resource**, especially one designed for ML



NVIDIA Tesla V100 SXM2

=





NVIDIA Tesla T4 PCIe

=

# GPU ON CONTAINER



https://github.com/NVIDIA/nvidia-docker

# GPU ON CONTAINER

```
$ docker run --gpus all nvidia/cuda:9.0-base nvidia-smi
$ docker run --runtime nvidia nvidia/cuda:9.0-base nvidia-smi
$ nvidia-docker run nvidia/cuda:9.0-base nvidia-smi
```

```
Mon Sep  2 14:57:28 2019
+-----------------------------------------------------------------------------+
| NVIDIA-SMI 418.56       Driver Version: 418.56       CUDA Version: 10.1      |
|-------------------------------+----------------------+----------------------+
| GPU  Name        Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap|         Memory-Usage | GPU-Util  Compute M. |
|===============================+======================+======================|
|   0  GeForce GTX 105...  Off  | 00000000:01:00.0 Off |                  N/A |
| N/A   48C    P8    N/A /  N/A |      2MiB /  4040MiB |      0%      Default |
+-------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------+
| Processes:                                                       GPU Memory |
|  GPU       PID   Type   Process name                             Usage      |
|=============================================================================|
|  No running processes found                                                 |
+-----------------------------------------------------------------------------+
```

# DEVICE PLUGINS API HISTORY

alpha.kubernetes.io/gpu

- Deprecated in v1.10, removed 1.11
- Using kernel module and library on host, expose with hostPath
- K8s expose GPU as scheduleable resources

Device Plugins (e.g. nvidia.com/gpu)

- Introduced in v1.8, beta 1.10
- Advertise hardware resources to Kubelet
- Deploy manually or with DaemonSet
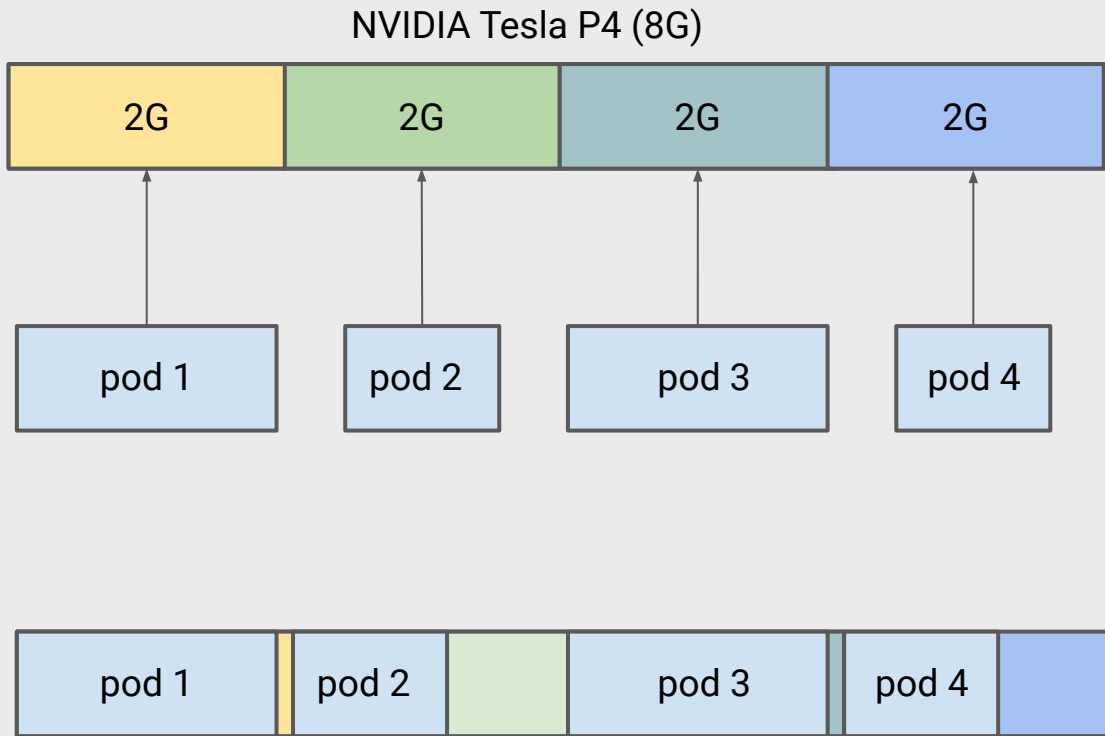- GPUs, high-performance NICs, FPGAs and InfiniBand adapters
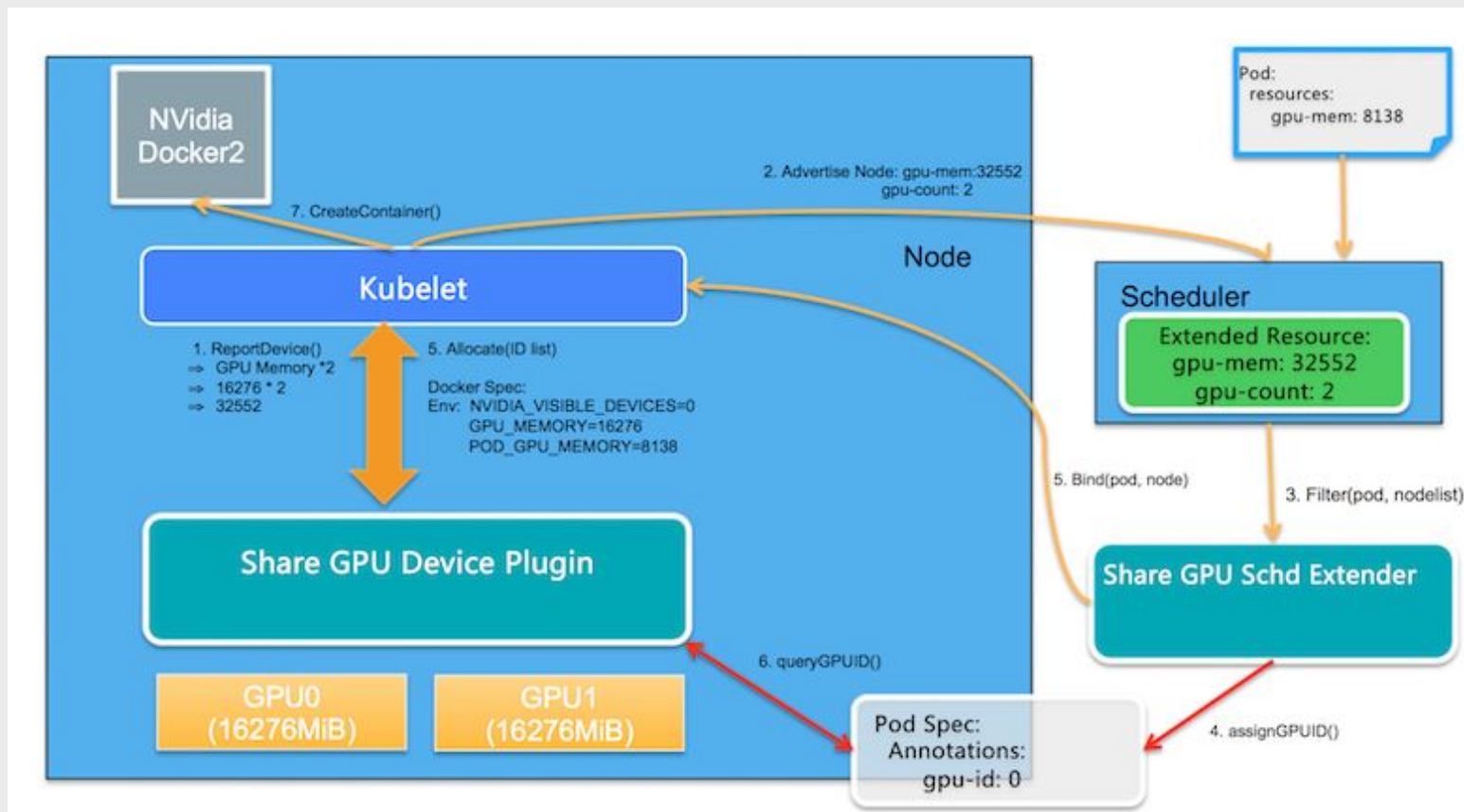
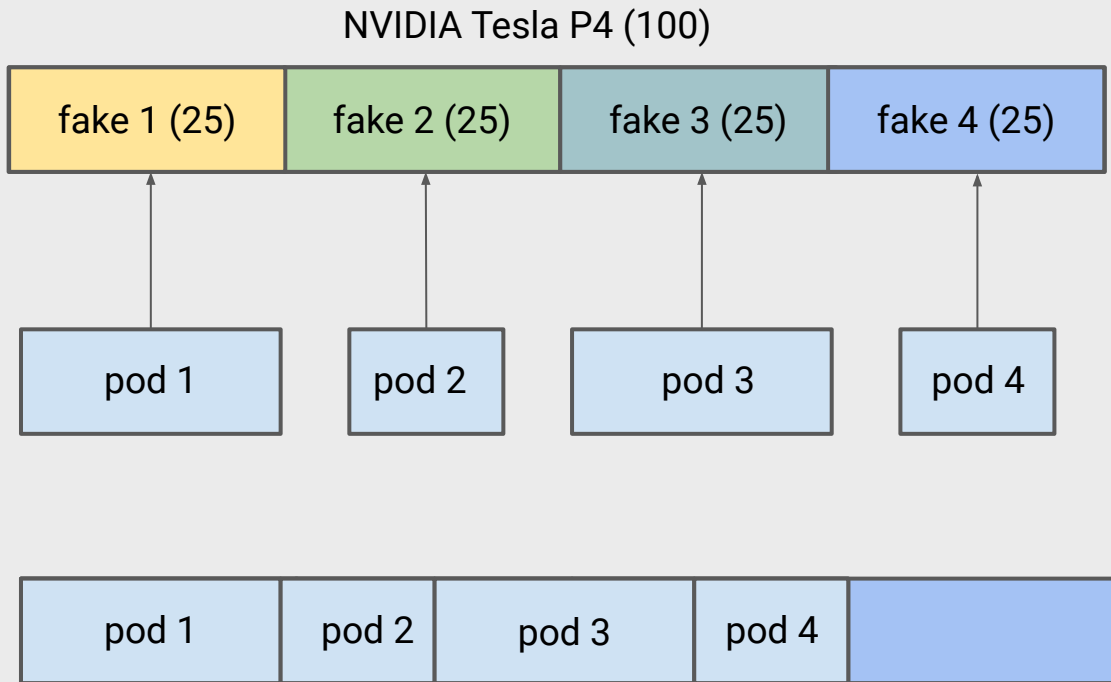# BIN–PACKING ON GPU

# BIN-PACKING ON GPU

Bin-pack based on GPU memory limit

NVIDIA Tesla P4 (8G)

| 2G | 2G | 2G | 2G |
|---|---|---|---|

pod 1  pod 2  pod 3  pod 4

| pod 1 | | pod 2 | | pod 3 | | pod 4 | |
|---|---|---|---|---|---|---|---|

# BIN–PACKING ON GPU

https://github.com/AliyunContainerService/gpushare-scheduler-extender

# BIN–PACKING ON GPU

Bin-pack based on fake GPU device

NVIDIA Tesla P4 (100)

| fake 1 (25) | fake 2 (25) | fake 3 (25) | fake 4 (25) |
|---|---|---|---|

| pod 1 | pod 2 | pod 3 | pod 4 |

| pod 1 | pod 2 | pod 3 | pod 4 | |

# BIN-PACKING ON GPU

| nvidia-driver daemonset | custom-device-plugin daemonset |
|---|---|

```
...
  spec:
    tolerations:
      - key: nvidia.com/gpu
        operator: Exists
        effect: NoSchedule
    containers:
      - name: cat-detection
        image:
        resources:
          limits:
            cpu: "4"
            memory: "8G"
            nodeflux.io/shared-gpu: 25
...
```

**gpu-node**

| Label | nodeflux.io/shared-gpu=true |
|---|---|
| Taints | NoSchedule nodeflux.io/shared-gpu=present |
| | NoSchedule nvidia.com/gpu=present |

**non-gpu-node**

https://github.com/Deepomatic/container-engine-accelerators/

# CURRENT HURDLES

- GPU sharing
  (Open K8s Issue #52757)
- Autoscaling support
- A general solution for GPU
- No Minikube support
- No local managed provider



Fractional GPUs: Software-based Compute and Memory Bandwidth Reservation for GPUs

https://ieeexplore.ieee.org/document/8743200

# EXAMPLE ON GKE

https://github.com/eufat/presentasi-k8s

# Thank you

Psst, we're hiring!

www.nodeflux.io