



KubeCon



CloudNativeCon

Europe 2019



KubeCon



CloudNativeCon

Europe 2019

GPU Machine Learning From Laptop to Cloud

Mark Puddick - Pivotal

@mpuddick



KubeCon



CloudNativeCon

Europe 2019

- What and Why
- Running on Laptop
- Moving to cloud challenges
- Running on Cloud
- Conclusion/Question



KubeCon



CloudNativeCon

Europe 2019

- Get started on laptop
- Keras, Tensorflow and Python
- Dockerising you Keras app
- Minikube
- Jupyter notebooks?

Laptop - Install



KubeCon



CloudNativeCon

Europe 2019

- Nvidia drivers - Latest Proprietary vs Open Source
- Cuda Toolkit and cuDNN SDK - Linux packages
 - Tensorflow GPU Install guide
 - Secure Boot
- Tensorflow GPU - pip

- Versions

- Now lets run some training!

Laptop - Docker



KubeCon



CloudNativeCon

Europe 2019

Docerkfile

```
FROM
nvidia/cuda:10.0-cudnn7-devel-ubuntu18.04

RUN apt-get update
RUN apt-get install -y python3.6 python3-pip

COPY requirements.txt /
RUN pip3 install -r /requirements.txt
```

REQUIREMENTS.txt

```
tensorflow-gpu
keras
pandas
numpy
s3fs
requests
```

Laptop - Docker



KubeCon



CloudNativeCon

Europe 2019

Dockerfile

```
FROM gcr.io/ml-gpu/base-gpu:v1  
  
COPY . /app  
  
CMD [ "python3", "/app/legos.py" ]
```

Laptop - Enabling GPU in Docker



KubeCon



CloudNativeCon

Europe 2019

```
docker run --runtime=nvidia --rm nvidia/cuda nvidia-smi  
sudo apt-get install nvidia-docker2
```

/etc/docker/daemon.json

```
{  
  "runtimes": {  
    "nvidia": {  
      "path": "nvidia-container-runtime",  
      "runtimeArgs": []  
    }  
  },  
  "default-runtime" : "nvidia"  
}
```


Laptop - Enabling GPU in Docker



KubeCon



CloudNativeCon

Europe 2019

- Mount local directory as volume

```
docker run --runtime=nvidia --mount  
type=bind,source="/home/markpudd/mldemo/data/testml",target=/app/data testml
```

Laptop - Enabling GPU in Minikube



KubeCon



CloudNativeCon

Europe 2019

- Driver kvm2 or None
- Use none to use underlying docker

```
minikube start --vm-driver=none --apiserver-ips 127.0.0.1 --apiserver-name
localhost

kubectl create -f
https://raw.githubusercontent.com/NVIDIA/k8s-device-plugin/v1.10/nvidia-dev
ice-plugin.yml
```

Laptop - Enabling GPU in Minikube



KubeCon



CloudNativeCon

Europe 2019

```
apiVersion: batch/v1
kind: Job
metadata:
  name: ml-job
spec:
  template:
    metadata:
      labels:
        app: testml
    spec:
      containers:
      - name: testml
        image: testml
        resources:
          limits:
            nvidia.com/gpu: 1
        imagePullPolicy: Never
        volumeMounts:
        - mountPath: /app/data
          name: test-volume
      restartPolicy: Never
```

```
apiVersion: v1
kind: PersistentVolume
metadata:
  name: pv0001
spec:
  accessModes:
    - ReadWriteOnce
  capacity:
    storage: 5Gi
  hostPath:
    path: /data/pv0001/
```

Running on Laptop



KubeCon



CloudNativeCon

Europe 2019

- Close browser windows
- `nvidia -smi`
- Decrease Batch size so we get a bit more memory

Moving to cloud challenges



KubeCon



CloudNativeCon

Europe 2019

So now we can just drop our container onto our Kubernetes cluster and run it?

Moving to cloud challenges



KubeCon



CloudNativeCon

Europe 2019

Not quite....

So all of the the previous demo loads from disk and send to stdout and displays graph using pyplot.....

Need to be a bit more cloud native

And it would be nice if we could spin up multiple instance with different hyper parameters

Where do we get our data?



KubeCon



CloudNativeCon

Europe 2019

- Object Store
- Database
- Disk

Where do we get our data?



KubeCon



CloudNativeCon

Europe 2019

- PersistentVolume
- One pod is going to copy data and effectively cache
- Other pods will load data from claim, so we don't need to change our application
- Other option is pods just pull the data when they need it....

Flow

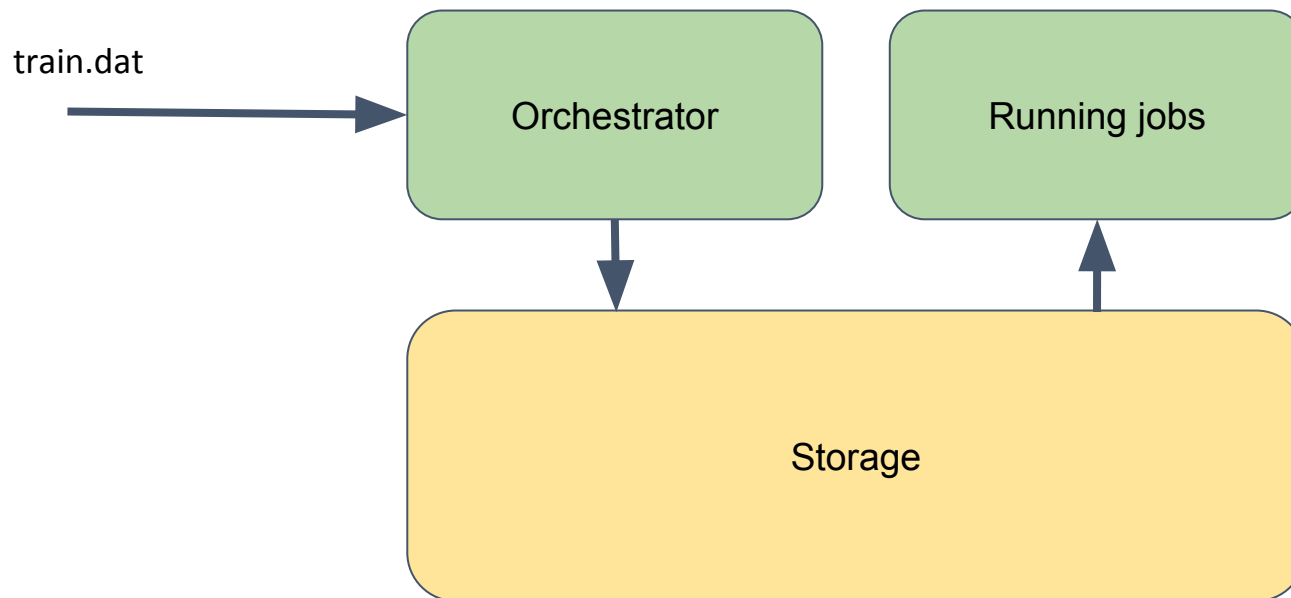


KubeCon



CloudNativeCon

Europe 2019



Moving to cloud challenges



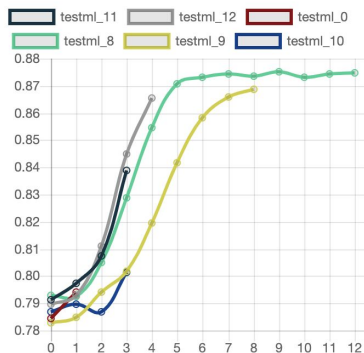
KubeCon



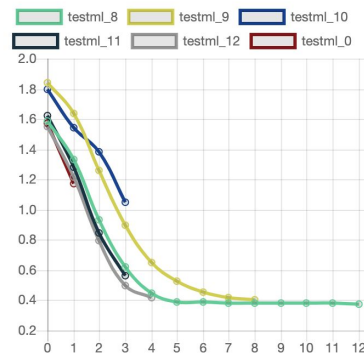
CloudNativeCon

Europe 2019

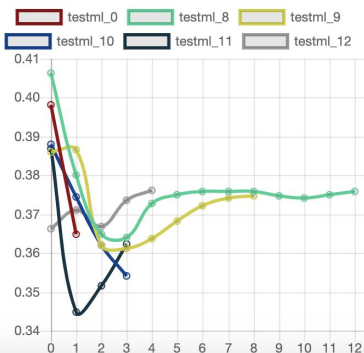
Accuracy



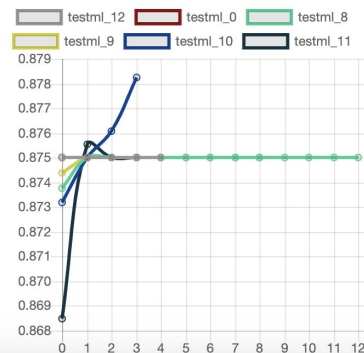
Loss



Cross Validation Accuracy



Cross Validation Loss



Moving to cloud challenges



KubeCon



CloudNativeCon

Europe 2019

Lets use a keras callback to send the data some were....

```
rm = callbacks.RemoteMonitor(root=server_root, path=endpoint_path,  
field='data', headers=None, send_as_json=True)  
  
tm =model.fit(X_train, Y_train, epochs = 75, batch_size = 64,  
validation_data=(X_test, Y_test), callbacks=[rm])
```

Flow

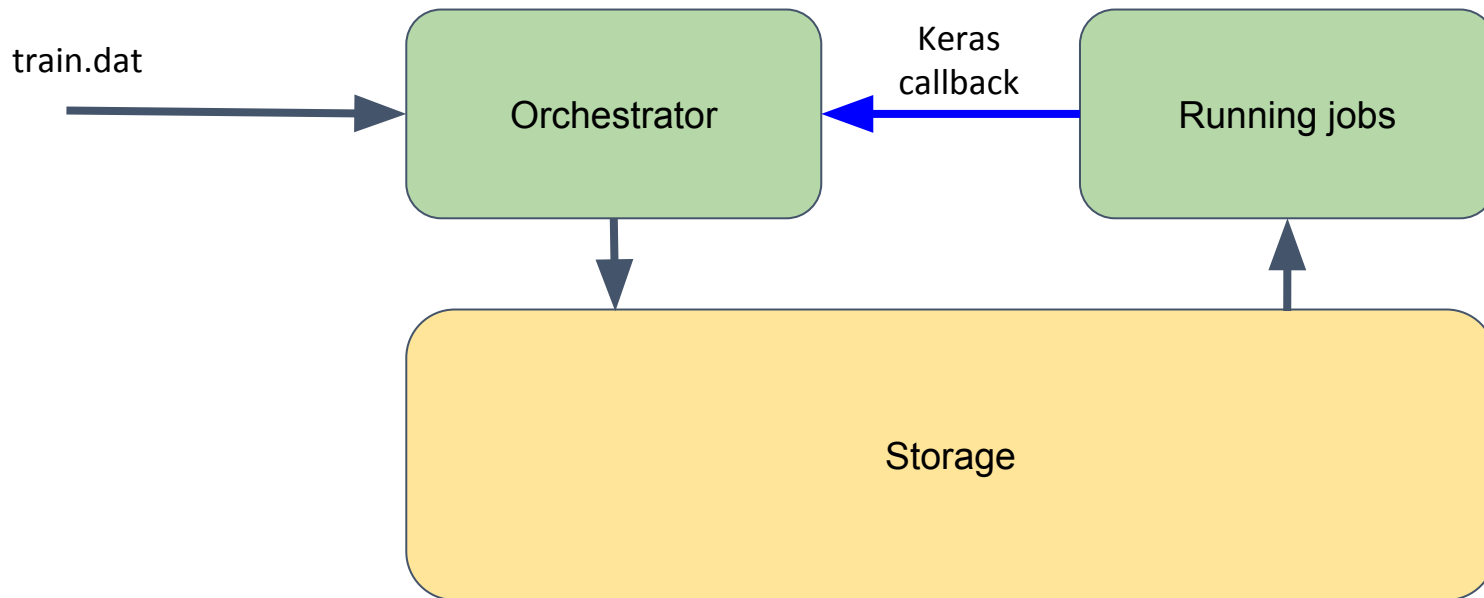


KubeCon



CloudNativeCon

Europe 2019



Where do we get our data?



KubeCon



CloudNativeCon

Europe 2019

- Job changes
 - Image
 - imagePullPolicy
- Still works on minikube

```
apiVersion: batch/v1
kind: Job
metadata:
  name: ml-job
spec:
  template:
    metadata:
      labels:
        app: testml
    spec:
      containers:
        - name: testml
          image: gcr.io/karttech/testml:v1
          resources:
            limits:
              nvidia.com/gpu: 1
            imagePullPolicy: Never
          volumeMounts:
            - mountPath: /app/data
              name: test-volume
      restartPolicy: Never
```

Moving to cloud challenges



KubeCon



CloudNativeCon

Europe 2019

- Tuning hyperparameters
- Hyperparameters set as environment

```
apiVersion: batch/v1
kind: Job
metadata:
  name: ml-job
spec:
  template:
    metadata:
      labels:
        app: testml
    spec:
      containers:
      - name: testml
        image: gcr.io/karttech/testml:v1
        env:
        - name: LEARNING_RATE
          value: "0.00008"
        resources:
          limits:
            nvidia.com/gpu: 1
          imagePullPolicy: Never
          volumeMounts:
          - mountPath: /app/data
            name: test-volume
        restartPolicy: Never
```

Flow

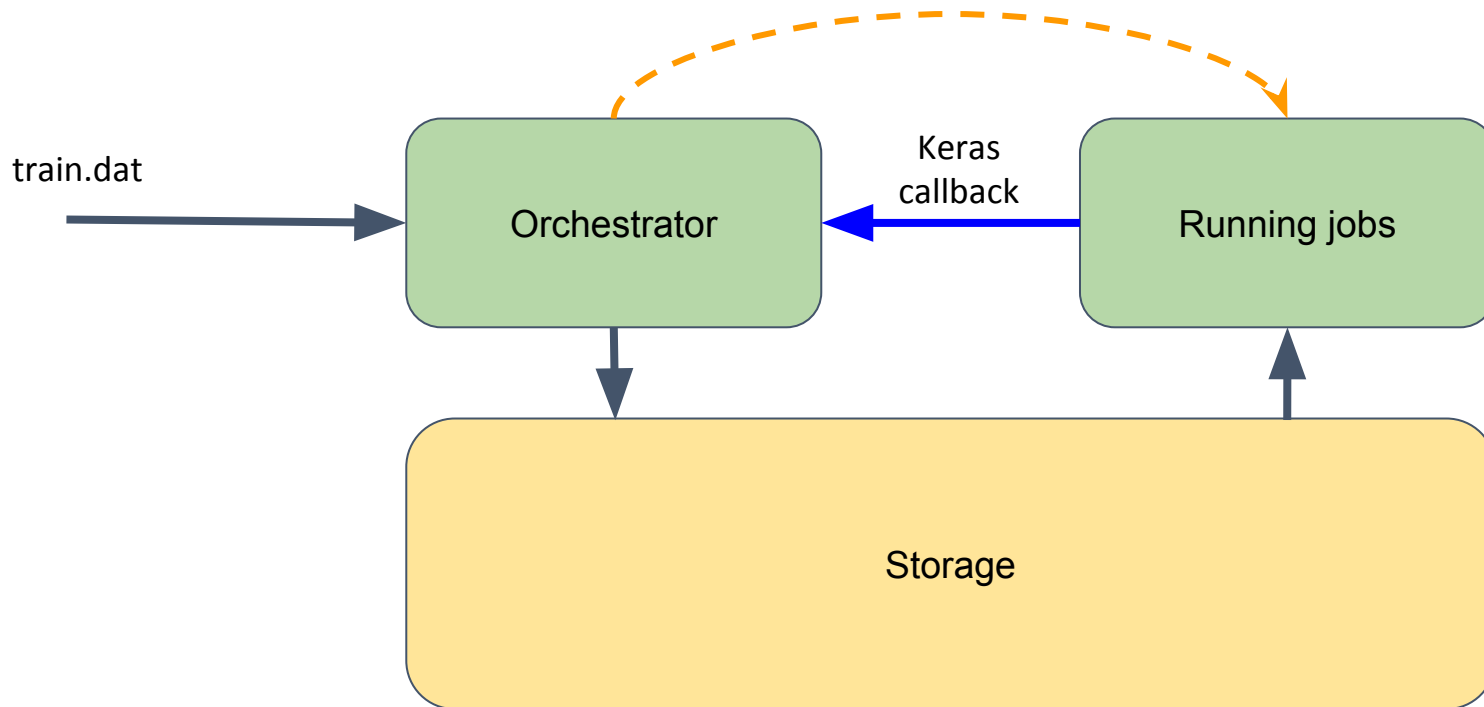


KubeCon



CloudNativeCon

Europe 2019



Running on Cloud



KubeCon



CloudNativeCon

Europe 2019

- Running multiple versions of the
- We've done a few tricks in the go app so we're only running one GPU pod per node

Running on Cloud



KubeCon

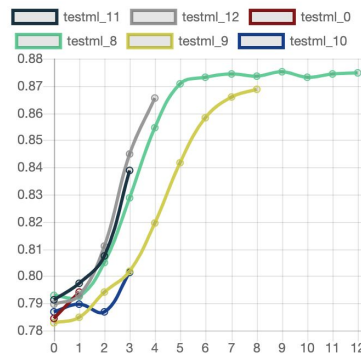


CloudNativeCon

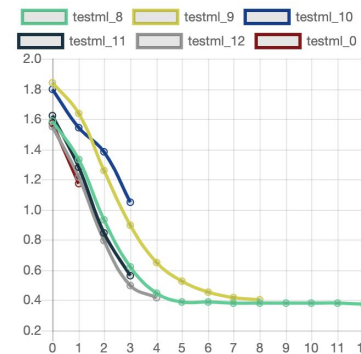
Europe 2019

- Here is our result!
- Weights stored on disk

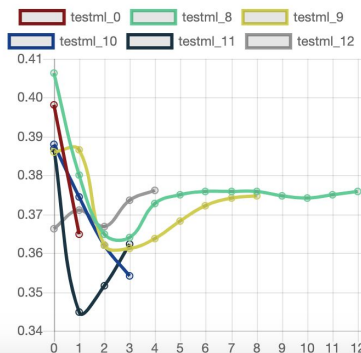
Accuracy



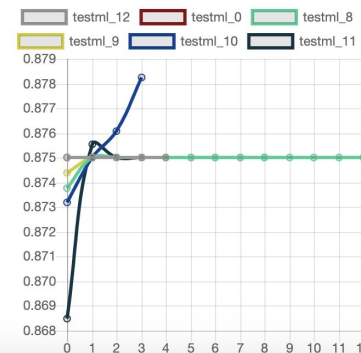
Loss



Cross Validation Accuracy



Cross Validation Loss



Conclusion



KubeCon



CloudNativeCon

Europe 2019

- Let point at some bigger data set
 - Running parallel TF jobs
 - Is it worth containerisation?
-
- Github - github.com/markpudd/mlorc



KubeCon



CloudNativeCon

Europe 2019

Thank you!

Mark Puddick - Pivotal
@mpuddick