

TF serving Docker -v100	tensorflow/serving:1.13.0-gpu	tensorflow/serving	1.13.0-gpu	f46687a9404d	2 weeks ago	2.38GB
TF serving Docker -1070	tensorflow/serving:latest-gpu	tensorflow/servir latest-gpu	ad003bd9bd92		4 months ago	2.18 GB

How models are created (Ke <https://colab.research.google.com/drive/1u79vDN4MZuq6gYIOkPmWsbghjunbDq6m>)
Client code <https://gist.github.com/alexcpn/f2366cbfcd0e0dfb22f0f551f8fc0161>
Client on same machine alexcpn/tfserving-keras-retinanet-dev-gpu latest 18797ab9f0e7

```
Thu Mar 14 22:59:42 2019 v100
-----+
| NVIDIA-SMI 384.145      Driver Version: 384.145      |
|-----+-----+-----+
| GPU Name      Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap|  Memory-Usage | GPU-Util  Compute M. |
|=====+=====+=====+
| 0  Tesla V100-PCIE...  Off | 00000000:B3:00.0 Off |             0 |
|N/A  34C   P0   35W / 250W | 31175MiB / 32502MiB | 0%      Default |
|-----+-----+-----+
|
| Processes:                      GPU Memory |
| GPU   PID  Type  Process name      Usage   |
|=====+=====+=====+
| 0   218965  C   tensorflow_model_server  31165MiB |
```

```
alex@drone-OMEN:~$ nvidia-smi
Sat Mar 16 01:12:00 2019
-----+
| NVIDIA-SMI 384.130      Driver Version: 384.130      |
|-----+-----+-----+
| 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 1070    Off | 00000000:01:00.0 Off |             N/A |
|N/A  48C   P8   12W / N/A | 7851MiB / 8105MiB | 1%      Default |
|-----+-----+-----+
|
| Processes:                      GPU Memory |
| GPU   PID  Type  Process name      Usage   |
|=====+=====+=====+
| 0   1862  G   /usr/lib/xorg/Xorg  202MiB |
| 0   7302  G   ...-token=EBF4E6EC29813767B25DB5C2E5B69C9D  152MiB |
| 0   9018  G   ...-token=D5172C7A8844F7AA8608E1B4C4EF62A2  40MiB |
| 0   21994  C   tensorflow_model_server  7451MiB |
|-----+-----+-----+
```

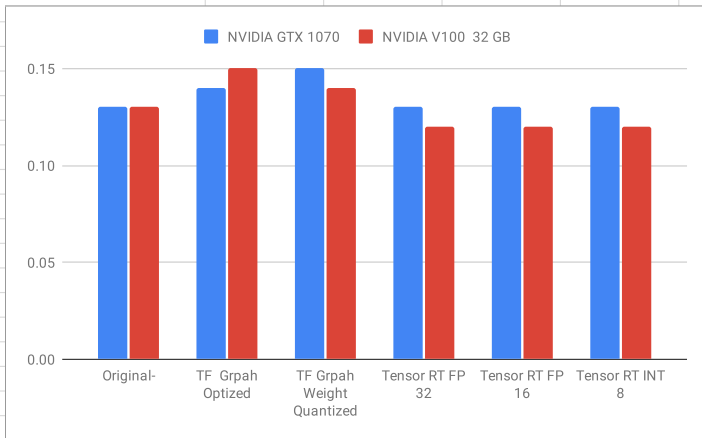
Reinnet model converted from Keras

Client - https://github.com/alexcpn/tf_serving_clients/blob/master/reinnet_client.py

Batchsize is 1 here

Model 4	Keras converted TF sering	(in image shape', (800, 1067, 3))
Model 6	TF Graph simple optimisation	(in tf shape', (1, 800, 1067, 3))
Model 7	TF Graph simple optimisation + Weight Qunatization	
Model 8	TF Graph simple optimisation + Weight + Model Qunatization	
Model 9	Based on Model 4 frozen; NVIDIA Tensor RT Optimisation FP 32	
Model 10	Based on Model 4 frozen; NVIDIA Tensor RT Optimisation FP 16	
Model 11	Based on Model 4 frozen; NVIDIA Tensor RT Optimisation INT 8	

No of Runs 1		
Model	NVIDIA GTX 1070	NVIDIA V100 32 GB
Original-	0.13	0.13
TF Grpah Optized	0.14	0.15
TF Grpah Weight Quanti	0.15	0.14
Tensor RT FP 32	0.13	0.12
Tensor RT FP 16	0.13	0.12
Tensor RT INT 8	0.13	0.12
No of runs :10		
4	1.15	0.81
6	1.34	1.16
7	1.15	1.27
9	1.23	0.82
10	1.22	0.83
11	1.22	0.85



MODEL = <https://github.com/tensorflow/models/tree/master/official/resnet>

FP32 = http://download.tensorflow.org/models/official/20181001_resnet/savedmodels/resnet_v2_fp32_savedmodel_NCHW.tar.gz

FP16 = http://download.tensorflow.org/models/official/20181001_resnet/savedmodels/resnet_v2_fp16_savedmodel_NCHW.tar.gz

Server

docker run --net=host --runtime=nvidia -it --rm -p 8900:8500 -p 8901:8501

-v /home/alex/coding/Python_neuralnet:/models/ tensorflow/serving:latesgpu

--model_config_file=/models/resnet50_fp32.json or resnet50_fp16.json

in image shape', (224, 224, 3))

('in tf shape', (2, 224, 224, 3))

('result no', 0)

('scores output', (2, 1001))

('labels output', (2,))

('Label', 535)

Client - https://github.com/alexcnpn/tf_serving_clients/blob/master/resnet50_client.py

Dev container with retinanet - alexcnpn/tfserving-keras-retinanet-dev-gpu

python resnet50_client.py -model_name=resnet_32 -batch_size=64

('Label', 535)

docker run -it --runtime=nvidia --net=host -v /home/alex/coding/Python_neuralnet:/coding --rm alexcnpn/tfserving-keras-retinanet-dev-gp

No of Runs 1/batch size 64	Model resnet_32	
	NVIDIA GTX 1070	NVIDIA V100 32 GB

Best time 1.3 1.25

per image 0.0203125 0.01953125

No of Runs 1/batch size 128 Model resnet_32

best time 3 2.5

per image 0.0234375 0.01953125

No of Runs 1/batch size 64 Model resnet_16

best time 1.3 1.25

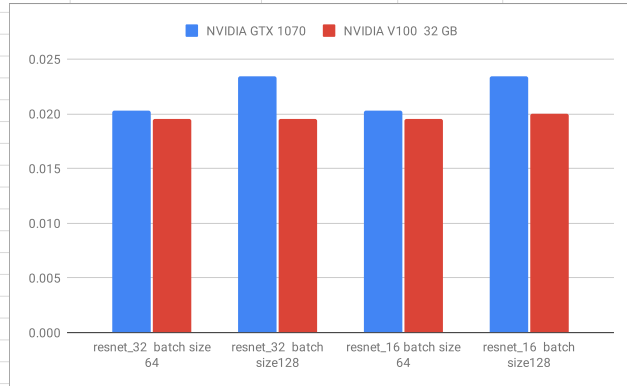
per image 0.0203125 0.01953125

No of Runs 1/batch size 128 Model resnet_16

best time 3 2.56

per image 0.0234375 0.02

Model Batch-size	NVIDIA GTX 1070	NVIDIA V100 32 GB
resnet_32 batch size 64	0.0203125	0.01953125
resnet_32 batch size128	0.0234375	0.01953125
resnet_16 batch size 64	0.0203125	0.01953125
resnet_16 batch size128	0.0234375	0.02



```
Model -ssd_resnet_50_fpn_coco -https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/detection_model_zoo.md
Server docker run --net=host --runtime=nvidia -it --rm -p 8900:8500 -p 8901:8501 -v /usr/alex:/models tensorflow/serving:1.13.0-gpu --rest_api_port=0 --enable_batching=true --model_config_file=/models/ssd_inception_v3_coco.json
Colab -https://colab.research.google.com/drive/1wQpWoc40kf_WSjTqDaReMx6FfjUn48
Client
docker run --entrypoint=/bin/bash --env http_proxy=<my proxy> --env https_proxy=<my proxy> --runtime=nvidia -it --rm -p 8900:8500 -p 8901:8501 -v /usr/alex:/coding --net=host tensorflow/tensorflow:1.13.0rc1-gpu-jupyter
pip install tensorflow-serving-api
pip install opencv-python==3.3.0.9
cd coding
python ssd_client_1.py -num_tests=1 -server=127.0.0.1:8500 -batch_size=1 -img_path='./examples/google1.jpg'
```

Original Model	0.5993821621
TensorRT FP 16	0.588560513
TensorRT INT 8	0.5967140198
TF Weights Quantized	0.6182711124

