

Robotic Inference

Naruhiko Nakanishi

Abstract—Nvidia DIGITS utilizes transfer learning from existing pretrained models such as LeNet, AlexNet, and GoogLeNet. In this Robotic Inference project, a network on the supplied data is trained, and a network using the own collected data is also chosen and trained. Upon evaluation, the model for the supplied data meets the required accuracy of above 75 percent (around 75.4 percent), and required inference time of below 10 ms (around 4.6 ms). For the own dataset, the predictions of randomly selected images are mostly good except one of car images.

Index Terms—Robot, Inference, Nvidia DIGITS, DNN, deep learning.

1 INTRODUCTION

IN 2012, AlexNet was the first DNN architecture to win the ImageNet classification challenge. Ever more complex and accurate DNNs have been developed on the ImageNet benchmark dataset ever since, including VGGNet, ResNet, Inception, GoogLeNet, and their many variations. [1]

Nvidia DIGITS (the Deep Learning GPU Training System) is a web application for training deep learning models. It utilizes transfer learning from existing pretrained models such as LeNet, AlexNet, and GoogLeNet.

In this Robotic Inference project, with DIGITS workspace, a network on the supplied data is trained, and a network using my own collected data is also chosen and trained.

2 BACKGROUND / FORMULATION

Using Digits, once the data is imported, a training model is chosen. The model has to achieve an inference time of 10 ms or less on the workspace and have an accuracy greater than 75 percent. A trained model is tested by running the command evaluate with the DIGITS server.

One of the created datasets can be selected that you want to train. The number of epochs can be changed as can the seed, batch size, solver, and learning rate.

One of the award winning standard networks, LeNet, AlexNet, or GoogLeNet can be selected. In this project, AlexNet is chosen for the supplied data set and then GoogLeNet is chosen for a robotic inference project.

3 DATA ACQUISITION

The supplied dataset are added into DIGITS and trained. In the dataset, there are photos taken from a Jetson mounted over a conveyor belt. The images of candy boxes, bottles, and nothing (empty conveyor belt) are trained for the purpose of real time sorting. This kind of design can be extrapolated to many things that require real time sorting.

For this robotic inference project, a phone is used to collect the dataset of images which are rectangular images at first. A system requires a constant input dimensionality. Therefore, the images are down-sampled to a fixed resolution of 256 x 256.

Fig.1 shows RGB type of the collected images of trains, cars, and background. In the dataset, there are 47 images of trains, 34 images of cars and 35 images of background.



Fig. 1. Collected images of trains, cars, and background

4 RESULTS

Using the supplied data, a network is created and is trained. AlexNet is chosen for the dataset. Fig.2 shows the result of training of 10 epochs and a learning rate of 0.001 with SGD as solver type.

Fig.3 shows a snapshot of the terminal output after running the evaluate command. Upon evaluation, the model for the supplied data meets the required accuracy of above 75 percent (around 75.4 percent), and required inference time of below 10 ms (around 4.6 ms).

GoogLeNet is chosen for the own dataset. Fig.4 shows the result of training of 200 epochs and a learning rate of 0.001 with Adam as solver type.

Fig.5, Fig.6 and Fig.7 show the predictions of randomly selected images. The predictions of them are mostly good except the prediction of bottom car image in Fig.6 that is 20.79 percent for car and is 79.16 percent for train.

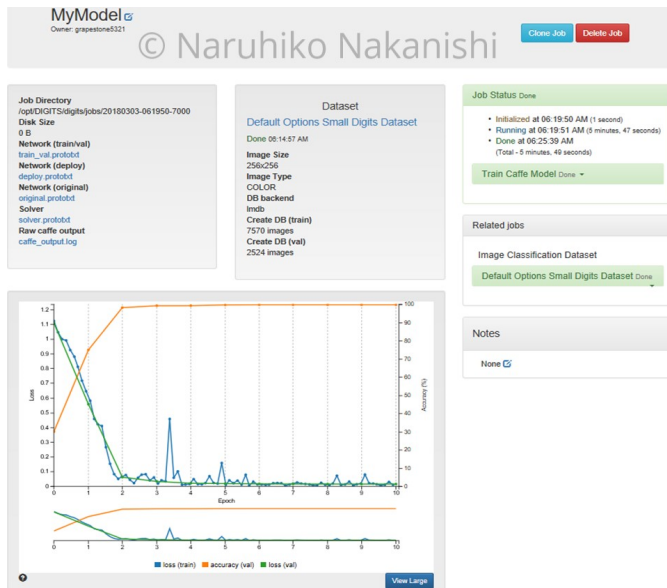


Fig. 2. Training for supplied data

```
Please enter the Job ID: 20180303-061950-7000

Calculating average inference time over 10 samples...
deploy: /opt/DIGITS/digits/jobs/20180303-061950-7000/deploy.prototxt
model: /opt/DIGITS/digits/jobs/20180303-061950-7000/snapshot_iter_600.caffemodel
output: softmax
iterations: 5
avgRuns: 10
Input "data": 3x227x227
Output "softmax": 3x1x1
name=data, bindingIndex=0, buffers.size()=2
name=softmax, bindingIndex=1, buffers.size()=2
Average over 10 runs is 4.64803 ms.
Average over 10 runs is 4.61309 ms.
Average over 10 runs is 4.60273 ms.
Average over 10 runs is 4.59012 ms.
Average over 10 runs is 4.14079 ms.

Calculating model accuracy...

% Total % Received % Xferd Average Speed Time Time Time Current
 100 14650 100 12334 100 2316 1062 199 0:00:11 0:00:11 --:--:-- 2178

Your model accuracy is 75.4098360656 %
root@58447fc2352e:/home/workspace#
```

© Naruhiko Nakanishi

Fig. 3. Evaluation

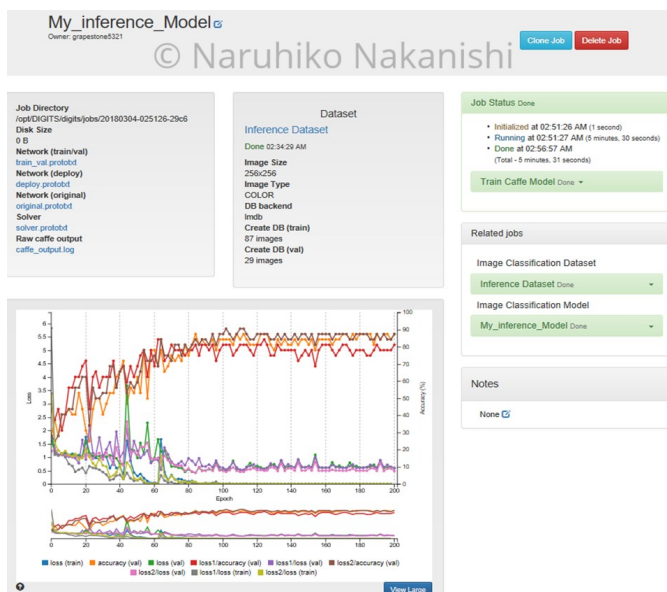


Fig. 4. Training for the own data

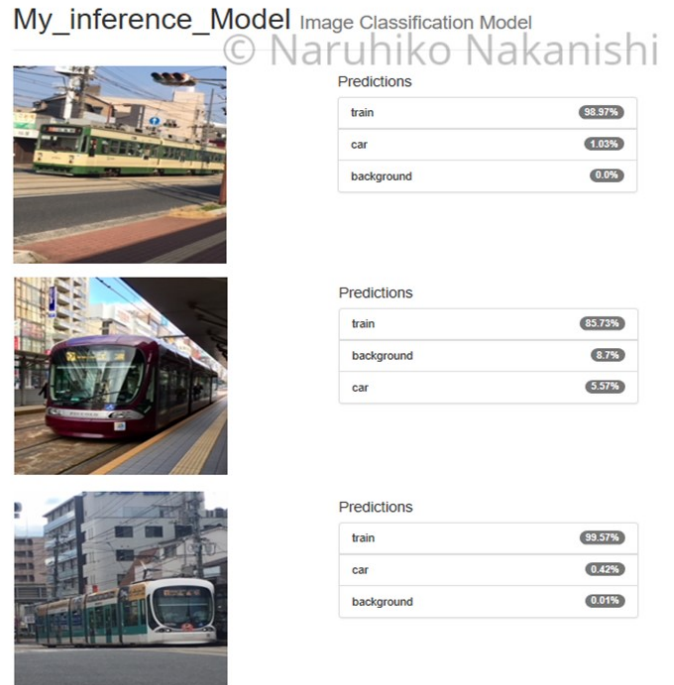


Fig. 5. train

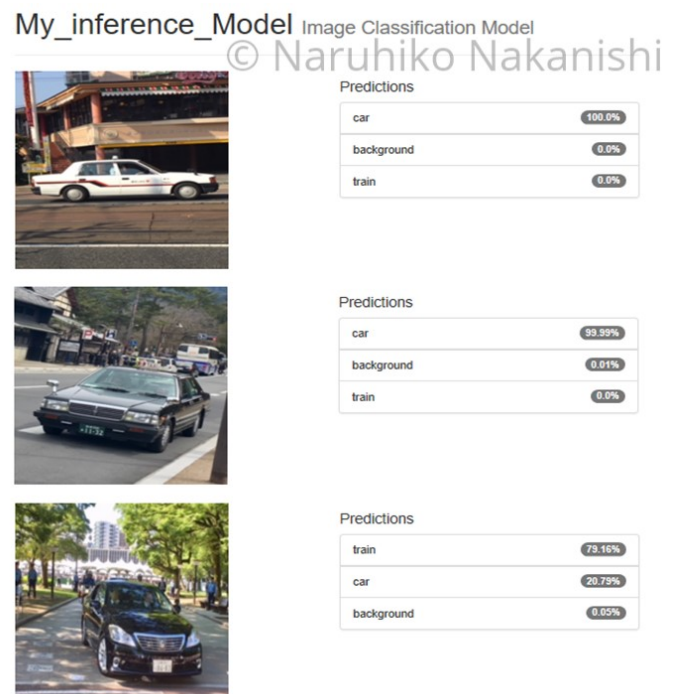


Fig. 6. car

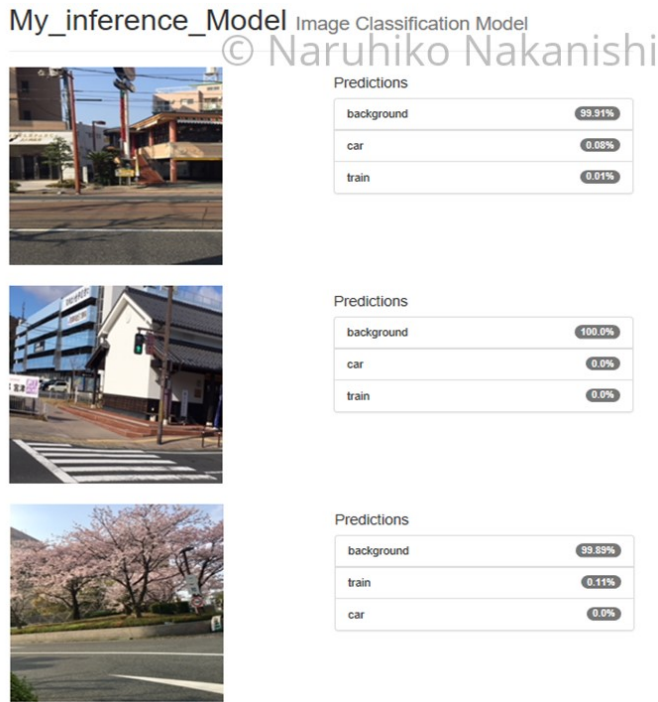


Fig. 7. background

5 DISCUSSION

For the supplied data, AlexNet is chosen. After the training of 10 epochs and a learning rate of 0.001 with SGD as solver type, the model for the supplied data meets the required accuracy of above 75 percent (around 75.4 percent), and required inference time of below 10 ms (around 4.6 ms).

To meet the requirements, parameters should be optimized. For example, the accuracy become 0 percent in the case of a learning rate of 0.01.

For the own dataset, GoogLeNet is chosen. After the training of 200 epochs and a learning rate of 0.001 with Adam as solver type, the predictions of randomly selected images are mostly good.

However the prediction of bottom car image in Fig.6 is not sufficient. A possible reason is that the bottom car image in Fig.6 looks similar to train. The countermeasure for this problem is to increase the number of data.

6 CONCLUSION / FUTURE WORK

Nvidia DIGITS utilizes transfer learning from existing pre-trained models such as LeNet, AlexNet, and GoogLeNet. In this Robotic Inference project, with DIGITS workspace, a network on the supplied data is trained, and a network using my own collected data is also chosen and trained.

The supplied dataset are added into DIGITS and trained. The images of candy boxes, bottles, and nothing (empty conveyor belt) are trained for the purpose of real time sorting. Upon evaluation, the model for the supplied data meets the required accuracy of above 75 percent (around 75.4 percent), and required inference time of below 10 ms (around 4.6 ms).

For this robotic inference project, a phone is used to collect the dataset of RGB images of trains, cars, and background. The images are down-sampled to a fixed resolution of 256 x 256.

For the own dataset, the predictions of randomly selected images are mostly good. However the prediction of one of car images is not sufficient. A possible reason is that the car image looks similar to train.

In the own dataset, there are 47 images of trains, 34 images of cars and 35 images of background. The countermeasure for this problem is to increase the number of data.

REFERENCES

- [1] Udacity, *Robotics Software Engineer Nanodegree Program Term2 Lesson7: Inference Application in Robotics*. Udacity, 2018.