

Part 4a1. DeePhi ML design flow



The flow: Input and Output files

- > DECENT tool needs floating point NN network model, Pre-training weights and calibration dataset of your DNN in Caffe format to work on
 - >> **float.prototxt**: Input floating point model
 - >> **float.caffemodel**: Pre-training weights in floating point
 - >> **calibration** dataset: calibration dataset is usually a subset of training set containing about 100~1000 pictures
- > Once the quantization is done, it will generated two model files named **deploy.prototxt** and **deploy.caffemodel**, these are input files for DNNC compiler:
 - >> **deploy.prototxt**: fixed point quantized network description file
 - >> **deploy.caffemodel**: fixed point quantized weights Caffe model file (note that this is not a standard Caffe format)

The Quantization flow in 5 steps

- > 1. Compress neural network model (from host side) with **decent**
 - >> command example for miniVggNet case: `source decent_miniVggNet.sh`
- > 2. Compile neural network model (from host side) with **dnnc**
 - >> command example for miniVggNet case: `source dnnc_miniVggNet.sh`
- > 3. Edit the **main.cc** program application (from host side)
- > 4. Compile hybrid application (from target side) with **make** utility
 - >> command examples for miniVggNet case: `make clean; make build`
- > 5. Run hybrid application (from target side)
 - >> command example for miniVggNet case: `./miniVggNet`

Some terminology for step2

- > **Kernel id:** id of current kernel. Every kernel has a unique id assigned by DNNC.
- > **Kernel name:** name of current kernel. Each kernel supported by DPU will have a corresponding ELF object file with name that is same with kernel name prefixed by “dpu_” with extension “.elf”. The kernel name will be used in your application code, thus N2Cube can identify different kernels correctly.
- > **Type:** the kernel type. Three types of kernel are supported by DNNC, see Table 6 for detail.
- > **Nodes:** All nodes included in current kernel. For kernel supported by DPU, we use “NA” instead to avoid printing all names.
- > **Input nodes:** All input nodes of current kernel. For kernel not supported by DPU, programmer must get output of preceding kernel through output nodes and feed them into input nodes of current nodes using APIs provided by N2Cube.
- > **Output nodes:** All output nodes of current kernel. The output nodes’ address and size can be extracted using APIs provided by N2Cube.

Quantization of miniVggNet



miniVggNet project files

For the calibration images data use a soft link
(to save HD space):

- `cd <home>/ML/cifar10/deephi/miniVggNet`
- `mkdir quantiz`
- `mkdir quantiz/data`
- `ln -s`
`/home/danieleb/ML/cifar10/input/cifar10_jpg/calib`
`./quantiz/data/calib`
- `cd zcu102`
- `ln -s test_images ../baseline/test_images`
- `ln -s test_images ../pruned/test_images`

```
(caffe_python_2) danieleb@CentOS63-x86-64:~$ cd ML/cifar10/deephi/
(caffe_python_2) danieleb@CentOS63-x86-64:~/ML/cifar10/deephi$ tree -d
.
├── miniVggNet
│   ├── pruning
│   │   ├── quantiz
│   │   │   ├── data
│   │   │   │   └── calib -> /home/danieleb/ML/cifar10/input/cifar10_jpg/calib
│   │   │   ├── decent_output
│   │   │   ├── dnnc_output
│   │   │   └── rpt
│   │   ├── regular_rate_0
│   │   │   └── snapshots
│   │   ├── regular_rate_0.3
│   │   │   └── snapshots
│   │   ├── regular_rate_0.7
│   │   │   └── snapshots
│   │   └── rpt
│   ├── quantiz
│   │   ├── data
│   │   │   └── calib -> /home/danieleb/ML/cifar10/input/cifar10_jpg/calib
│   │   ├── decent_output
│   │   ├── dnnc_output
│   │   └── rpt
│   └── zcu102
│       ├── baseline
│       │   ├── model
│       │   │   └── arm64_4096
│       │   ├── rpt
│       │   ├── src
│       │   └── test_images -> ../test_images
│       ├── pruned
│       │   ├── model
│       │   │   └── arm64_4096
│       │   ├── rpt
│       │   ├── src
│       │   └── test_images -> ../test_images
│       └── test_images
35 directories
(caffe_python_2) danieleb@CentOS63-x86-64:~/ML/cifar10/deephi$
```

miniVggNet project files

- > There are three main folders:
- > **quantiz**: it contains the scripts and related results to quantize the CNN (which is baseline, that is, not pruned). Logfiles are stored in **rpt**
- > **pruning**: it contains the scripts and related files to prune the CNN first, and then to quantize it. Logfiles are stored in **rpt**
- > **zcu102**: it contains whatever needed to compile and run the application on ZCU102
 - >> **baseline**: application files of the baseline quantized CNN. Logfiles are stored in **rpt**
 - >> **pruned**: application files of the pruned and quantized CNN. Logfiles are stored in **rpt**
 - >> **test_images**: images used at run time to compute the top-5 predictions

Step1: decent_miniVggNet.sh script

```
1  #!/usr/bin/env bash
2
3  DNNDK_ROOT=/home/ML/DNNDK/tools
4
5  #working directory
6  work_dir=$(pwd)
7  #path of float model
8  model_dir=${work_dir}
9  #output directory
10 output_dir=${work_dir}/decent_output
11
12 #soft link to the calibration data
13 ln -s /home/ML/cifar10/input/cifar10_jpg/calib /home/ML/cifar10/deephi/miniVggNet/quantiz/data/calib
14
15 # copy input files from miniVggNet Caffe project via soft links
16 ln -s /home/ML/cifar10/caffe/models/miniVggNet/m3/deephi_train_val_3_miniVggNet.prototxt /home/ML/cifar10/deephi/miniVggNet/quantiz/float.prototxt
17
18 ln -s /home/ML/cifar10/caffe/models/miniVggNet/m3/snapshot_3_miniVggNet__iter_40000.caffemodel /home/ML/cifar10/deephi/miniVggNet/quantiz/float.caffemodel
19
20
21 # run DECENT
22 $DNNDK_ROOT/decent quantize \
23     -model ${model_dir}/float.prototxt \
24     -weights ${model_dir}/float.caffemodel \
25     -output_dir ${output_dir} \
26     -method 1 \
27
28
```

Command:

source decent_miniVggNet.sh 2>&1 | tee logfile_decent_miniVggNet.txt

Step1: run decent (with auto_test)

1/2

Fragments of decent logfile:

- check CNN description and then run calibration

```
I0831 11:28:32.822849 28315 convert_proto.cpp:160] Opening file ./data/calib/calibration.txt
I0831 11:28:32.824931 28315 convert_proto.cpp:171] A total of 1000 images.
I0831 11:28:32.826131 28315 convert_proto.cpp:2286] Merge InnerProductBatchNorm -> InnerProduct: fc1 + bn5
I0831 11:28:32.846849 28315 convert_proto.cpp:2286] Merge InnerProductBatchNorm -> InnerProduct: fc1 + bn5
I0831 11:28:32.881103 28315 net.cpp:323] The NetState phase (0) differed from the phase (1) specified by a rule in layer data
I0831 11:28:32.881122 28315 net.cpp:323] The NetState phase (0) differed from the phase (1) specified by a rule in layer accuracy-top1
I0831 11:28:32.881125 28315 net.cpp:323] The NetState phase (0) differed from the phase (1) specified by a rule in layer accuracy-top5
I0831 11:28:32.881278 28315 net.cpp:52] Initializing net from parameters:
state {
  phase: TRAIN
}
layer {
  name: "data"
  type: "ImageData"
  top: "data"
  top: "label"
  include {
    phase: TRAIN
  }
  transform_param {
    crop_size: 32
    mean_value: 125
    mean_value: 123
    mean_value: 114
  }
  image_data_param {
    source: "./data/calib/calibration.txt"
    batch_size: 10
    shuffle: true
    root_folder: "./data/calib/"
  }
}
layer {
  name: "data_fixed"
  type: "FixedNeuron"
  bottom: "data"
  top: "data"
  param {
    lr_mult: 0
  }
}
```

--:--- logfile_decent_miniVggNet_autotest.txt 1% L39 (Text)

Step1: run decent (with auto_test)

2/2

Fragments of decent logfile:

- Check accuracy on the quantized model: 0.87 (top-1) and 0.99 (top-5)
- Generate output files:
 - deploy.caffemodel
 - deploy.prototxt

```
2065 I0109 17:33:16.115645 14383 net_test.cpp:339] Test iter: 43/50, loss = 0.516615
2066 I0109 17:33:16.115649 14383 net_test.cpp:339] Test iter: 43/50, top-1 = 0.82
2067 I0109 17:33:16.115653 14383 net_test.cpp:339] Test iter: 43/50, top-5 = 0.98
2068 I0109 17:33:16.122541 14383 net_test.cpp:339] Test iter: 44/50, accuracy = 0.84
2069 I0109 17:33:16.122555 14383 net_test.cpp:339] Test iter: 44/50, loss = 0.443076
2070 I0109 17:33:16.122558 14383 net_test.cpp:339] Test iter: 44/50, top-1 = 0.84
2071 I0109 17:33:16.122561 14383 net_test.cpp:339] Test iter: 44/50, top-5 = 1
2072 I0109 17:33:16.128968 14383 net_test.cpp:339] Test iter: 45/50, accuracy = 0.94
2073 I0109 17:33:16.128983 14383 net_test.cpp:339] Test iter: 45/50, loss = 0.227701
2074 I0109 17:33:16.128986 14383 net_test.cpp:339] Test iter: 45/50, top-1 = 0.94
2075 I0109 17:33:16.128988 14383 net_test.cpp:339] Test iter: 45/50, top-5 = 1
2076 I0109 17:33:16.135323 14383 net_test.cpp:339] Test iter: 46/50, accuracy = 0.9
2077 I0109 17:33:16.135336 14383 net_test.cpp:339] Test iter: 46/50, loss = 0.289338
2078 I0109 17:33:16.135340 14383 net_test.cpp:339] Test iter: 46/50, top-1 = 0.9
2079 I0109 17:33:16.135344 14383 net_test.cpp:339] Test iter: 46/50, top-5 = 1
2080 I0109 17:33:16.141512 14383 net_test.cpp:339] Test iter: 47/50, accuracy = 0.96
2081 I0109 17:33:16.141526 14383 net_test.cpp:339] Test iter: 47/50, loss = 0.128769
2082 I0109 17:33:16.141530 14383 net_test.cpp:339] Test iter: 47/50, top-1 = 0.96
2083 I0109 17:33:16.141533 14383 net_test.cpp:339] Test iter: 47/50, top-5 = 1
2084 I0109 17:33:16.147711 14383 net_test.cpp:339] Test iter: 48/50, accuracy = 0.96
2085 I0109 17:33:16.147724 14383 net_test.cpp:339] Test iter: 48/50, loss = 0.175614
2086 I0109 17:33:16.147728 14383 net_test.cpp:339] Test iter: 48/50, top-1 = 0.96
2087 I0109 17:33:16.147730 14383 net_test.cpp:339] Test iter: 48/50, top-5 = 1
2088 I0109 17:33:16.154798 14383 net_test.cpp:339] Test iter: 49/50, accuracy = 0.94
2089 I0109 17:33:16.154812 14383 net_test.cpp:339] Test iter: 49/50, loss = 0.224946
2090 I0109 17:33:16.154816 14383 net_test.cpp:339] Test iter: 49/50, top-1 = 0.94
2091 I0109 17:33:16.154819 14383 net_test.cpp:339] Test iter: 49/50, top-5 = 1
2092 I0109 17:33:16.161334 14383 net_test.cpp:339] Test iter: 50/50, accuracy = 0.86
2093 I0109 17:33:16.161348 14383 net_test.cpp:339] Test iter: 50/50, loss = 0.50688
2094 I0109 17:33:16.161351 14383 net_test.cpp:339] Test iter: 50/50, top-1 = 0.86
2095 I0109 17:33:16.161355 14383 net_test.cpp:339] Test iter: 50/50, top-5 = 0.98
2096 I0109 17:33:16.161357 14383 net_test.cpp:346] Test Results:
2097 I0109 17:33:16.161360 14383 net_test.cpp:347] Loss: 0.383396
2098 I0109 17:33:16.161363 14383 net_test.cpp:361] accuracy = 0.8744
2099 I0109 17:33:16.161370 14383 net_test.cpp:361] loss = 0.383396 (* 1 = 0.383396 loss)
2100 I0109 17:33:16.161373 14383 net_test.cpp:361] top-1 = 0.8744
2101 I0109 17:33:16.161376 14383 net_test.cpp:361] top-5 = 0.9924
2102 I0109 17:33:16.161379 14383 net_test.cpp:387] Test Done!
2103 I0109 17:33:16.312860 14383 decent.cpp:333] Start Deploy
2104 I0109 17:33:16.341810 14383 decent.cpp:341] Deploy Done!
2105 -----
2106 Output Deploy Weights: "/home/ML/cifar10/deephi/miniVggNet/quantiz/decent_output/deploy.caffemodel"
2107 Output Deploy Model:   "/home/ML/cifar10/deephi/miniVggNet/quantiz/decent_output/deploy.prototxt"
2108
```

Step1: prepare the input files 1/2

> This and next pages are only to explain how the input files should be prepared. Such changes were already saved in the file named `deephi_train_val_3_miniVggNet.prototxt`

>> From the trained model with the highest prediction accuracy copy the caffemodel and training description files and rename them. For example, assuming the best Caffe-trained original model is placed in `<wkr>/models/miniVggNet/m3` do the following:

```
cd <wkr>/models/miniVggNet/m3
ln -s <wkr>/models/miniVggNet/m3/train_val_3_miniVggNet.prototxt float.prototxt
ln -s <wkr>/models/miniVggNet/m3/snapshot_3_miniVggNet__iter_40000.caffemodel float.caffemodel
```

>> `float.prototxt` and `float.caffemodel` are the 2 input files to the quantization process (just renamed from the `m3` models originally trained in Caffe)

>> Now you need to do some changes to `float.prototxt` file (see next page):

- Remove the **Data** type layers for the original TRAIN phase
- Add an **ImageData** type layer with the calibration images for the new TRAIN phase
- Add on the bottom two **Accuracy** type layers to compute **top-1** and **top-5** accuracies
- Remove the mean file and put separate values (DPU does not support reading a mean file)

Step1: prepare the input files

2/2

add this
new layer
on top

comment this
layer

```
# add below layer for Quantization with DeePhi' DECENT
#####
# CIFAR10 miniVggNet m3
layer {
  name: "data"
  type: "ImageData"
  top: "data"
  top: "label"
  include {
    phase: TRAIN
  }
  transform_param {
    crop_size: 32
    mean_value: 125
    mean_value: 123
    mean_value: 114
    #mean file is not supported by DeePhi DPU
    #mean_file: "/home/danieleb/ML/cifar10/input/mean.binaryproto"
  }
  image_data_param {
    source: "./data/calib/calibration.txt"
    root_folder: "./data/calib/"
    batch_size: 10
    shuffle: true
  }
}
# comments below layer for the original TRAIN
#####
#layer {
#  name: "data"
#  type: "Data"
#  top: "data"
#  top: "label"
#  include {
#    phase: TRAIN
#  }
#  transform_param {
#    mirror: true
#    crop_size: 32
#    mean_file: "/home/danieleb/ML/cifar10/input/mean.binaryproto"
#  }
#  data_param {
#    source: "/home/danieleb/ML/cifar10/input/lmdb/train_lmdb"
#    batch_size: 128
#    backend: LMDB
#  }
#}
layer {
  name: "data"
  type: "Data"
  top: "data"
  top: "label"
  include {
    phase: TEST
  }
  transform_param {
    mirror: false
    crop_size: 32
    mean_value: 125
    mean_value: 123
    mean_value: 114
  }
  data_param {
    source: "/home/danieleb/ML/cifar10/input/lmdb/valid_lmdb"
```

```
name: "fc2"
type: "InnerProduct"
bottom: "drop3"
top: "fc2"
param {
  lr_mult: 1
  decay_mult: 1
}
param {
  lr_mult: 2
  decay_mult: 0
}
inner_product_param {
  num_output: 10
  weight_filler {
    #type: "gaussian"
    type: "xavier"
    #std: 0.001
  }
  bias_filler {
    type: "constant"
    value: 1
  }
}
}
layer {
  name: "loss"
  type: "SoftmaxWithLoss"
  bottom: "fc2"
  bottom: "label"
  top: "loss"
}
#layer {
#  name: "accuracy"
#  type: "Accuracy"
#  bottom: "fc2"
#  bottom: "label"
#  top: "accuracy"
#  include {
#    phase: TEST
#  }
#}
layer {
  name: "accuracy-top1"
  type: "Accuracy"
  bottom: "fc2"
  bottom: "label"
  top: "top-1"
  include {
    phase: TEST
  }
}
layer {
  name: "accuracy-top5"
  type: "Accuracy"
  bottom: "fc2"
  bottom: "label"
  top: "top-5"
  include {
    phase: TEST
  }
  accuracy_param {
    top_k: 5
  }
}
```

add these 2
new layers
at the bottom

Step2: dnnc_miniVggNet.sh script

```
1  #!/bin/bash
2
3  DNNDK_ROOT=/home/ML/DNNDK/tools
4
5  net=miniVggNet
6  model_dir=decent_output
7  output_dir=dnnc_output
8
9  echo "Compiling network: ${net}"
10
11  $DNNDK_ROOT/dnnc --prototxt=${model_dir}/deploy.prototxt \
12    --caffemodel=${model_dir}/deploy.caffemodel \
13    --output_dir=${output_dir} \
14    --net_name=${net} \
15    --dpu=4096FA \
16    --cpu_arch=arm64 \
17    --mode=debug \
18    --save_kernel
19
20
21
22  echo " copying dpu elf file into ../zcu102/baseline/model/arm64_4096 "
23  cp ${output_dir}/dpu_${net}\*.elf ${output_dir}/../../zcu102/baseline/model/arm64_4096
24
25  echo " copying the test images to be used by the ZCU102"
26  cp -r /home/ML/cifar10/input/cifar10_jpg/test ${output_dir}/../../zcu102/test_images
27
```

Command:

```
source dnnc_miniVggNet.sh 2>&1 | tee logfile_dnnc_miniVggNet.txt
```

Step2: run dnnc

```
Compiling network: miniVggNet  
[DNNC][Warning] Layer [loss] is not supported in DPU, deploy it in CPU instead.
```

DNNC Kernel Information

1. Overview

```
kernel numbers : 2  
kernel topology : 0 -> 1
```

2. Kernel Description in Detail

```
kernel id      : 0  
kernel name    : miniVggNet_0  
type           : DPUKernel (Supported, Running on DPU)  
nodes          : NA  
input node(s)  : conv1  
output node(s) : fc2  
  
kernel id      : 1  
kernel name    : miniVggNet_1  
type           : CPUKernel (Not-Supported, Running on CPU)  
nodes          : loss  
input node(s)  : loss  
output node(s) : loss
```

```
-:***- logfile_dnnc_miniVggNet.txt All L3 (Text)
```

dnnc logfile

Steps 1 and 2: what we have done so far

- > Running `decent_miniVggNet.sh` and `dnnc_miniVggNet.sh` scripts compile the miniVggNet model into DPU kernel, which is an ELF file containing the DPU instructions and parameters for the network model.
- > The miniVggNet neural network model was divided as 2 different Kernels:
 - >> Kernel 0 : `miniVggNet_0` (running on DPU with ELF file `dpu_miniVggNet_0.elf`)
 - >> Kernel 1 : `miniVggNet_1` (running on ARM CPU)
- > Kernel 0 runs all the layers of the CNN -i.e. CONV, BN, POOL, RELU, FC- on DPU in the PL of ZCU102 platform
- > Kernel 1 run the SoftMax last layer on the PS of ZCU102 platform
- > The quantized CNN has an accuracy unchanged from the original model (at least for this CNN case), as estimated by **decent**

Step3: main.cc operations

> The main operations are:

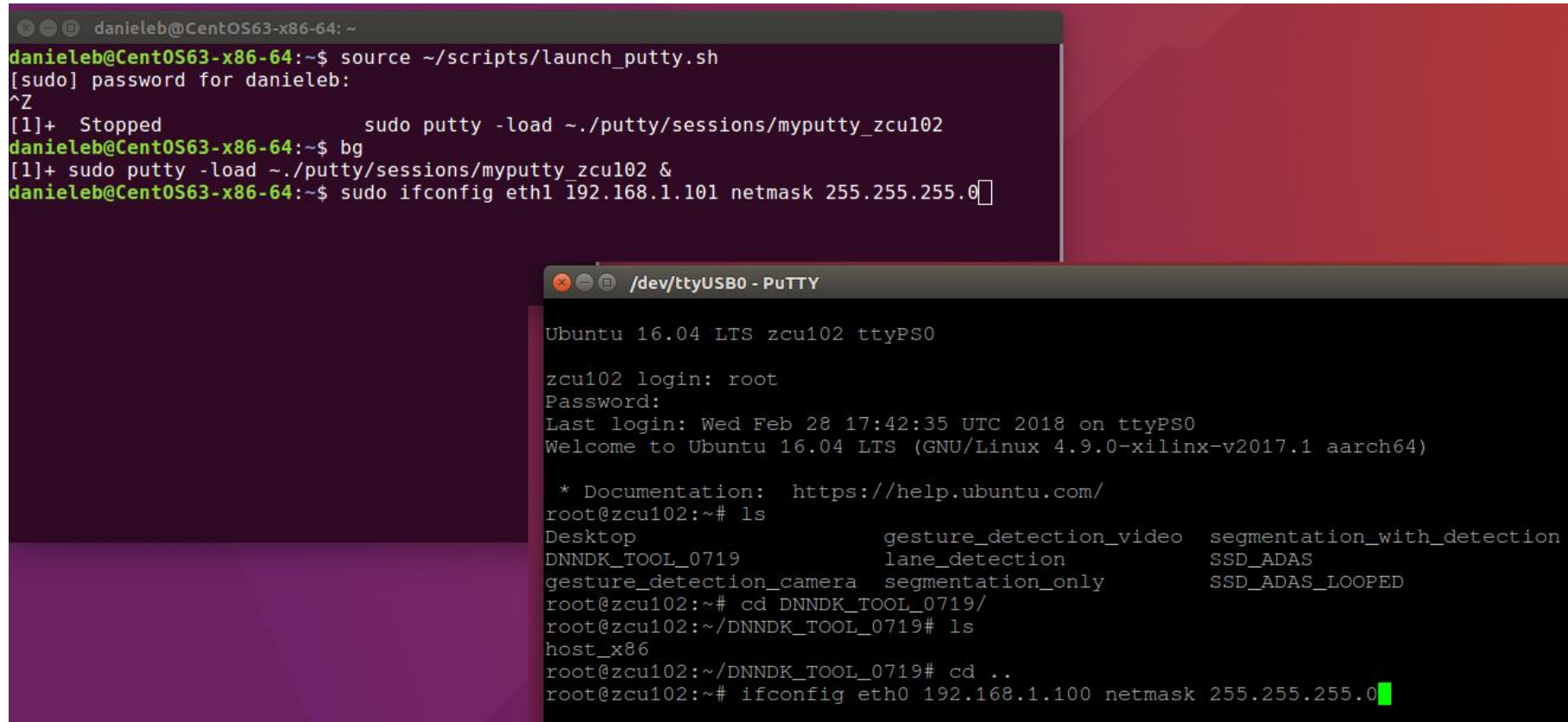
- >> 1. Call **dpuOpen ()** to open DPU device
- >> 2. Call **dpuLoadKernel ()** to load DPU kernel miniVggNet_0 from miniVggNet model
- >> 3. Call **dpuCreateTask ()** to create task for each DPU kernel
- >> 4. Coding “Softmax” CPU kernel to work with DPU kernel in **run_miniVggNet ()** to do image classification
- >> 5. Call **dpuDestroyKernel ()** and **dpuDestroyTask ()** to destroy Kernel and Task
- >> 6. Call **dpuClose ()** to close DPU device

Step3: run_miniVggNet() routine

> **run_miniVggNet()** does the following tasks:

- >> 1. Read picture and set it as the input to DPU kernel miniVggNet_0 by calling **dpuSetInputImage2()** API
- >> 2. Call **dpuRunTask()** to run taskConv convolution operation in minVggNet network model.
- >> 3. Call **dpuRunTask()** to do full connection operation of taskFC on DPU.
- >> 4. Do **Softmax** operation on CPU using the output of full connection operation as input.
- >> 5. Output the **top-5 classification** category and the corresponding probability

Step4: compile hybrid application on ZCU102 target



```
danieleb@CentOS63-x86-64: ~  
danieleb@CentOS63-x86-64:~$ source ~/scripts/launch_putty.sh  
[sudo] password for danieleb:  
^Z  
[1]+  Stopped                  sudo putty -load ~/.putty/sessions/myputty_zcu102  
danieleb@CentOS63-x86-64:~$ bg  
[1]+ sudo putty -load ~/.putty/sessions/myputty_zcu102 &  
danieleb@CentOS63-x86-64:~$ sudo ifconfig eth1 192.168.1.101 netmask 255.255.255.0  
  
/dev/ttyUSB0 - PuTTY  
Ubuntu 16.04 LTS zcu102 ttyPS0  
zcu102 login: root  
Password:  
Last login: Wed Feb 28 17:42:35 UTC 2018 on ttyPS0  
Welcome to Ubuntu 16.04 LTS (GNU/Linux 4.9.0-xilinx-v2017.1 aarch64)  
  
* Documentation:  https://help.ubuntu.com/  
root@zcu102:~# ls  
Desktop                                gesture_detection_video  segmentation_with_detection  
DNNDK_TOOL_0719                       lane_detection           SSD_ADAS  
gesture_detection_camera               segmentation_only        SSD_ADAS_LOOPED  
root@zcu102:~# cd DNNDK_TOOL_0719/  
root@zcu102:~/DNNDK_TOOL_0719# ls  
host_x86  
root@zcu102:~/DNNDK_TOOL_0719# cd ..  
root@zcu102:~# ifconfig eth0 192.168.1.100 netmask 255.255.255.0
```

Copy the elf files related to the DPU kernels from Host PC to the Target board file system by **scp** command, having the Target ZCU102 board EP2P connected to Host PC:

- from Host, run `sudo ifconfig eth1 192.168.1.101 netmask 255.255.255.0`
- from Target, run `ifconfig eth0 192.168.1.100 netmask 255.255.255.0`

Steps 4 and 5: working on the ZCU102 board

```
zcu102
├── baseline
│   ├── model
│   │   └── arm64_4096
│   ├── rpt
│   ├── src
│   └── test_images -> ../test_images
├── pruned
│   ├── model
│   │   └── arm64_4096
│   ├── rpt
│   ├── src
│   └── test_images -> ../test_images
└── test_images
```

I have created a tar file and copied it from/to host to/from target via `scp` command (see next pages)

Booth the ZCU102 board with DeePhi' SD-Card

I have created a folder named `~/cifar10/miniVggNet/zcu102` organized as shown on the left

The folder `test_images` has the same 1000 images of the CIFAR10 test folder (remember that during the training I used 50000 images of the train and 9000 images of the validation dataset)

The folder `model/arm64_4096` contains the same `dpu_miniVggNet_0.elf` kernel that was created during step2

Step5: run hybrid application on ZCU102 target

- > To check the fps performance with multithreading I do the following commands:

```
source run_fps_miniVggNet.sh 2>&1 | tee ./rpt/logfile_fps_miniVggNet.txt
```

- > To check the top-5 average accuracy I do the following commands:

```
source top5_miniVggNet 2>&1 | tee ./rpt/logfile_top5_miniVggNet.txt
```

Step5: best performance in fps

- > Comment all unnecessary code with printf or file I/O
- > Comment image resizing (the images are already 32x32)
- > Comment top-5 computation
- > **Best performance is 3799 with 4 threads**

```
12 make: warning: Clock skew detected. Your build may be incomplete.
13
14 ./miniVggNet 1
15 total image : 1000
16 [Time]695304us
17 [FPS]1438.22
18
19 ./miniVggNet 2
20 total image : 1000
21 [Time]390204us
22 [FPS]2562.76
23
24 ./miniVggNet 3
25 total image : 1000
26 [Time]271455us
27 [FPS]3683.85
28
29 ./miniVggNet 4
30 total image : 1000
31 [Time]263226us
32 [FPS]3799.02
33
34 ./miniVggNet 5
35 total image : 1000
36 [Time]276489us
37 [FPS]3616.78
38
39 ./miniVggNet 6
40 total image : 1000
41 [Time]281639us
42 [FPS]3550.64
```

Step5: compute top-5 accuracy

- > Uncomment top-5 computation lines and add a line to print the image currently read just before the top-5 computation.
- > See on the right the logfile captured at runtime
- > You need a Python script to post process such logfile and compute the overall average accuracy at runtime (otherwise add some C code lines in the main) in order to compare it with your predictions from the original Caffe model

```
logfile_run_top5_miniVggNet.txt
1 total image : 1000
2 DBG imread ./images/frog_00001.jpg
3 [Top]0 prob = 0.940882 name = frog
4 [Top]1 prob = 0.028412 name = deer
5 [Top]2 prob = 0.028412 name = cat
6 [Top]3 prob = 0.001415 name = bird
7 [Top]4 prob = 0.000858 name = dog
8 DBG imread ./images/frog_00017.jpg
9 [Top]0 prob = 0.995822 name = frog
10 [Top]1 prob = 0.003169 name = cat
11 [Top]2 prob = 0.000908 name = dog
12 [Top]3 prob = 0.000045 name = horse
13 [Top]4 prob = 0.000035 name = ship
14 DBG imread ./images/horse_00007.jpg
15 [Top]0 prob = 0.923117 name = horse
16 [Top]1 prob = 0.045959 name = dog
17 [Top]2 prob = 0.027876 name = deer
18 [Top]3 prob = 0.001782 name = bird
19 [Top]4 prob = 0.000842 name = cat
20 DBG imread ./images/truck_00074.jpg
21 [Top]0 prob = 0.997527 name = truck
22 [Top]1 prob = 0.002473 name = automobile
23 [Top]2 prob = 0.000000 name = ship
24 [Top]3 prob = 0.000000 name = frog
25 [Top]4 prob = 0.000000 name = airplane
26 DBG imread ./images/frog_00092.jpg
27 [Top]0 prob = 0.999935 name = frog
28 [Top]1 prob = 0.000035 name = cat
29 [Top]2 prob = 0.000010 name = dog
30 [Top]3 prob = 0.000008 name = deer
31 [Top]4 prob = 0.000006 name = bird
32 DBG imread ./images/dog_00070.jpg
33 [Top]0 prob = 0.879771 name = dog
34 [Top]1 prob = 0.119064 name = cat
35 [Top]2 prob = 0.000625 name = frog
36 [Top]3 prob = 0.000295 name = bird
37 [Top]4 prob = 0.000230 name = deer
38 DBG imread ./images/bird_00057.jpg
39 [Top]0 prob = 0.603890 name = bird
```

Step5: compute top-5 accuracy

```
void TopK(const float *d, int size, int k, vector<string> &vkind) {
    assert(d && size > 0 && k > 0);
    priority_queue<pair<float, int>> q;

    for (auto i = 0; i < size; ++i) {
        q.push(pair<float, int>(d[i], i));
    }

    for (auto i = 0; i < k; ++i)
    {
        pair<float, int> ki = q.top();
        printf("[Top]%d prob = %-8f  name = %s\n", i, d[ki.second], vkind[ki.second].c_str());
        q.pop();
    }
}

-:--- main.cc      57% L174   (C++/l Abbrev)
    DPUTask *taskconv = dpuCreateTask(kernelconv, DPU_MODE_NORMAL); // profiling not enabled
    //DPUTask *taskconv = dpuCreateTask(kernelconv, DPU_MODE_PROF); // profiling enabled
    //enable profiling
    //int res1 = dpuEnableTaskProfile(taskconv);
    //if (res1!=0) printf("ERROR IN ENABLING TASK PROFILING FOR CONV KERNEL\n");

    for(unsigned int ind = i ;ind < images.size();ind+=threadnum)
    {
        Mat img = imread(baseImagePath + images.at(ind));
        cout << "DBG imread " << baseImagePath + images.at(ind) << endl;
        //Size sz(32,32);
        //Mat img2; resize(img, img2, sz); //DB
        //run_miniVggNet(taskconv,img2); //DB: images are already 32x32 and do not need any resize
        run_miniVggNet(taskconv,img);
    }

    // Destroy DPU Tasks & free resources
    dpuDestroyTask(taskconv);
});
-:--- main.cc      81% L258   (C++/l Abbrev)
```


Check the top-5 accuracy with a python script

- > By capturing the output of DPU at runtime on a logfile, I can post-process it with a python script to check the effective top-1 and top-5 accuracies of quantized net:
- > **top-1: 87% average accuracy measured at runtime, as also estimated by decent during quantization step1**

```
1421 LINE:  airplane_00093.jpg
1422
1423 PREDICTED:  [ 3.] cat
1424 EXPECTED :  [ 0.] airplane
1425 [Top]0 prob = 0.925463  name = cat
1426 [Top]1 prob = 0.035884  name = deer
1427 [Top]2 prob = 0.021765  name = airplane
1428 [Top]3 prob = 0.006236  name = bird
1429 [Top]4 prob = 0.003782  name = horse
1430
1431
1432 LINE:  truck_00028.jpg
1433
1434 PREDICTED:  [ 8.] ship
1435 EXPECTED :  [ 9.] truck
1436 [Top]0 prob = 0.999815  name = ship
1437 [Top]1 prob = 0.000075  name = truck
1438 [Top]2 prob = 0.000058  name = airplane
1439 [Top]3 prob = 0.000021  name = frog
1440 [Top]4 prob = 0.000021  name = automobile
1441
1442
1443 LINE:  dog_00100.jpg
1444
1445 PREDICTED:  [ 2.] bird
1446 EXPECTED :  [ 5.] dog
1447 [Top]0 prob = 0.752371  name = bird
1448 [Top]1 prob = 0.130742  name = dog
1449 [Top]2 prob = 0.101822  name = cat
1450 [Top]3 prob = 0.010732  name = frog
1451 [Top]4 prob = 0.003075  name = deer
1452
1453
1454 number of total images predicted  999
1455 number of top1 false predictions  132
1456 number of top1 right predictions  867
1457 number of top5 false predictions  2
1458 number of top5 right predictions  997
1459 top1 accuracy = 0.87
1460 top5 accuracy = 1.00

top1_false = 0
top5_true = 0
top5_false = 0
img_count = 0
false_pred = 0

test_ids = np.zeros((NUMEL,1))
preds = np.zeros((NUMEL, 1))
idx = 0

for ln in range(0, tot_lines):

    if "DBG" in lines[ln]:

        top5_lines = lines[ln:ln+6]
        s2 = top5_lines[0].index(" ")
        class_name = top5_lines[0][20:s2].strip()
        #print 'DBG: found class ', class_name, ' in line ', ln, ': ', lines[ln]

        predicted = top5_lines[1][30 : ].strip()

        if class_name in top5_lines[1]:
            top1_true += 1
            top5_true += 1
        elif class_name in top5_lines[2]:
            top5_true += 1
            top1_false +=1
        elif class_name in top5_lines[3]:
            top5_true += 1
            top1_false +=1
        elif class_name in top5_lines[4]:
            top5_true += 1
            top1_false +=1
        elif class_name in top5_lines[5]:
            top5_true += 1
            top1_false +=1
        else:
            top5_false += 1
            top1_false +=1

        test_ids[idx] = labelNames[class_name] # ground truth
        preds[idx] = labelNames[predicted] # actual prediction

        if (predicted != class_name) :
            print "LINE: ", top5_lines[0].split(".")[1].strip()
            print "PREDICTED: ", preds[idx], predicted
            print "EXPECTED : ", test_ids[idx], class_name
            for k in range(1, 6):
                print top5_lines[k].strip()
            print "\n"

        img_count +=1
        idx += 1

    if ( idx == (NUMEL-1) ):
        break

-:--- logfile_check_dpu_top1_10Kimages.txt  Bot L14365 (Text) -:--- check_dpu_runtime_accuracy.py  40% L109 (Python)
```