

Ingenic<sup>®</sup>

## Magik Training Quantization Guide

Date: Feb. 2022

**Ingenic®**

**Magik Training Quantization Guide**

Copyright© Ingenic Semiconductor Co. Ltd 2022. All rights reserved.

**Release history**

<b>Date</b>	<b>Author</b>	<b>Revision</b>	<b>Change</b>
Feb. 2022	Heidi	1.0.1	First release

**Disclaimer**

This documentation is provided for use with Ingenic products. No license to Ingenic property rights is granted. Ingenic assumes no liability, provides no warranty either expressed or implied relating to the usage, or intellectual property right infringement except as provided for by Ingenic Terms and Conditions of Sale.

Ingenic products are not designed for and should not be used in any medical or life sustaining or supporting equipment.

All information in this document should be treated as preliminary. Ingenic may make changes to this document without notice. Anyone relying on this documentation should contact Ingenic for the current documentation and errata.

**Ingenic Semiconductor Co., Ltd.**

**Ingenic Headquarters, East Bldg. 14, Courtyard #10  
Xibeiwang East Road, Haidian District, Beijing, China,  
Tel: 86-10-56345000  
Fax:86-10-56345001  
Http: //www.ingenic.com**

## Contents

1. Abstract.....	1
2. Environmental requirements.....	1
3. Environmental installation steps.....	1
3.1 Ubuntu.....	1
3.2 GCC.....	1
3.3 Nvidia Driver.....	2
3.4 CUDA.....	3
3.5 CUDNN.....	3
3.6 Python.....	4
3.7 Pytorch.....	5
3.8 Torchvision.....	5
3.9 Magik Trainingkit.....	5
3.10 Others.....	6
4. YOLOv5s Training Process.....	6
4.1 Downloading of code and model.....	6
4.2 Data Preparation.....	7
4.3 Network Training.....	8
4.4 Model Testing.....	10
5. Model Transformation.....	11
5.1 From pt to onnx.....	11
5.2 From onnx to bin.....	12
6. Run on-board.....	13
6.1 Code compiling.....	13
6.2 Run.....	14
7. Check data.....	17
1. PC.....	17
2. Board.....	18
8. Q & A.....	19

## 1. Abstract

This Guide is mainly for novices who use the magik platform of Ingenic processing core. Here, we take the yolov5s-person under the pytorch framework as an example, at the same time, the target device is set to “T40”. The whole process is introduced in detail from environment construction, data preparation, network training, model transformation to final board operation, which aims to guide users to be familiar with the use method of our magik platform and get the board operation process clearly.

## 2. Environmental Requirements

- Linux
- GCC ( $\geq 5.4.0$ )
- Nvidia Driver
- CUDA ( $\geq 9.0$ )
- CUDNN
- Python ( $\geq 3.5$ )
- Pytorch ( $\geq 1.3$ )
- Torchvision

If all the above environments are available (the version may not be subject to fixed requirements), you can directly go to step 3.9.

## 3. Environmental Installation Steps

### 3.1 Ubuntu

- system : Ubuntu16.04

### 3.2 GCC

- Version: 5.4.0
- Specific Steps
  1. Check Version  
\$ gcc -v
  2. Download Link  
<http://ftp.gnu.org/gnu/gcc>
  3. Compile and Install

```

$ tar -zxvf gcc-5.4.0.tar.bz2
$ cd gcc-5.4.0
$ ./contrib/download_prerequisites
$ cd ..
$ mkdir gcc-build-5.4.0
$ cd gcc-build-5.4.0
$ ../gcc-5.4.0/configure --enable-checking=release --enable-languages=
c,c++ --disable-multilib
$ sudo make
$ sudo make install
4. Check again
$ gcc -v

```

### 3.3 Nvidia Driver

- Version: >= 440
- Specific Steps
  1. Download link  
<http://www.nvidia.cn/Download/index.aspx?lang=cn>
  2. Disable Nouveau third party drivers  
 open edit profile:  

```
$ sudo gedit /etc/modprobe.d/blacklist.conf
```

 add to the last line: blacklist nouveau  

```
$ sudo update-initramfs -u
```

```
$ reboot
```
  3. Installation  

```
$ lsmod | grep nouveau
```

```
$ sudo /etc/init.d/lightdm stop (or: sudo service lightdm stop)
```

```
$ sudo chmod a+x NVIDIA-Linux-x86_64-440.64.run
```

```
$ sudo ./NVIDIA-Linux-x86_64-440.64.run --no-opengl-files
```

```
$ sudo /etc/init.d/lightdm start (or: sudo service lightdm start) --no-opengl-files --no-x-check --no-nouveau-check
```
  4. Check Version  

```
$ reboot
```

```
$ nvidia-smi
```

If the required driver version appears, the installation is over.

NVIDIA-SMI 440.31				Driver Version: 440.31				CUDA Version: 10.2			
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 108...	Off	00000000:02:00.0	Off	0%	N/A					
0%	35C	P8	10W / 250W	0MiB / 11178MiB	0%	Default					
1	GeForce GTX 108...	Off	00000000:03:00.0	Off	0%	N/A					
0%	50C	P2	60W / 250W	829MiB / 11178MiB	0%	Default					
2	GeForce GTX 108...	Off	00000000:83:00.0	Off	0%	N/A					
0%	39C	P8	9W / 250W	10MiB / 11177MiB	0%	Default					

### 3.4 CUDA

- Version: 10.0
- Specific Steps
  1. Download link  
<https://developer.nvidia.com/cuda-10.0-download-archive?>  
Choose: linux—x86\_64—Ubuntu—16.04—runfile(local)  
Download: cuda\_10.0.130\_410.48\_linux.run

#### CUDA Toolkit 10.0 Archive

Select Target Platform ⓘ

Click on the green buttons that describe your target platform. Only supported platforms will be shown.

Operating System: Windows, Linux, Mac OSX

Architecture ⓘ: x86\_64, ppc64le

Distribution: Fedora, OpenSUSE, RHEL, CentOS, SLES, Ubuntu

Version: 18.04, 16.04, 14.04

Installer Type ⓘ: runfile (local), deb (local), deb (network), cluster (local)

Download Installers for Linux Ubuntu 16.04 x86\_64

The base installer is available for download below.  
There is 1 patch available. This patch requires the base installer to be installed first.

Base Installer [Download (2.0 GB) ⬇]

Installation Instructions:

1. Run `sudo sh cuda_10.0.130_410.48_linux.run`
2. Follow the command-line prompts

Patch 1 (Released May 10, 2019) [Download (3.3 MB) ⬇]

In this patch we introduce new APIs for JPEG stream parsing and device and pinned memory control as well as a new hybrid decode API that decouples decoding process into pure host and device stages enabling more flexible control flow. The new APIs also support ROI decoding and 4 channel jpeg bitstreams.

2. Installation  
`$ sudo sh cuda_10.0.130_410.48_linux.run`
3. Check Version  
`$ cat /usr/local/cuda/version.txt`
4. CUDA Setting  
`$ export PATH=/usr/local/cuda/bin:$PATH`  
`$ export LD_LIBRARY_PATH=/usr/local/cuda/lib64:$LD_LIBRARY_PATH`

### 3.5 CUDNN

- Version: 7.6.5
- Specific Steps
  1. Download link  
<https://developer.nvidia.com/rdp/cudnn-download>  
cuDNN Library for Linux: cudnn-10.0-linux-x64-v7.6.5.32.tgz

2. Unzip
 

```
$ tar -zxvf cudnn-10.0-linux-x64-v7.6.5.32.tgz
cuda/include/cudnn.h
cuda/NVIDIA_SLA_cuDNN_Support.txt
cuda/lib64/libcudnn.so
cuda/lib64/libcudnn.so.7
cuda/lib64/libcudnn.so.7.6.5
cuda/lib64/libcudnn_static.a
```
3. Copy
 

```
$ sudo cp cuda/lib64/* /usr/local/cuda-10.0/lib64/
$ sudo cp cuda/include/* /usr/local/cuda-10.0/include/
```
4. Check version
 

```
$ cat /usr/local/cuda/include/cudnn.h | grep CUDNN_MAJOR -A 2
#define CUDNN_MAJOR 7
#define CUDNN_MINOR 6
#define CUDNN_PATCHLEVEL 5
#define CUDNN_VERSION (CUDNN_MAJOR * 1000 +
CUDNN_MINOR * 100 + CUDNN_PATCHLEVEL)
#include "driver_types.h"
```

### 3.6 Python

- Version: 3.7 (Please note that do not uninstall or delete the system comes with python2.7, so as not to affect the operation of the system.)
- Specific Steps
  1. Download Link
 

```
https://www.python.org/ftp/python/3.7.3/Python-3.7.3.tar.xz
```
  2. Compile & Install
 

```
$ tar -xvJf Python-3.7.3.tar.xz
$ cd Python-3.7.3
$ ./configure --prefix=/usr/local/bin/python3
$ sudo make
$ sudo make install
```
  3. Create Soft Link
 

```
$ ln -s /usr/local/bin/python3/bin/python3 /usr/bin/python3
$ ln -s /usr/local/bin/python3/bin/pip3 /usr/bin/pip3
```
  4. Check Version
 

```
$ python3 --version
```

### 3.7 Pytorch

- Version: 1.3.0
- Specific Steps
  1. Installation  
`$ pip3 install torch==1.3.0`
  2. Check Version  
`$ python3`  
`>> import torch`  
`>> print(torch.__version__)`

### 3.8 Torchvision

- Version: 0.4.2
- Specific Steps
  1. Installation  
`$ pip3 install torchvision==0.4.2`
  2. Check Version  
`$ python3`  
`>> import torchvision`  
`>> print(torchvision.__version__)`

### 3.9 Magik Trainingkit

- Version: 1.1.1
- Specific Steps
  1. Check Version  
`$ python`  
`>> import torch`  
`>> print(torch.__version__)` #version of torch  
`>> print(torch.version.cuda)` #version of cuda  
`>> print(torch.backends.cudnn.version())` #version of cudnn  
You can also set the other version of cuda and cudnn to meet your own demands.
  2. Get trainingkit  
You can get your trainingkit installation package(whl file) from our plugin directory “magik-toolkit/TrainingKit/pytorch/magik\_whl” according to the above environment and then you should install it by pip, for example, “pip3 install magik\_trainingkit\_torch\_130-1.1.1-py3-none-any.whl”.  
Please note that training installation package here is provided by us



according to the customer's environment. If there is any change and you cannot find the package in our directory, please inform the relevant technicians of the magik platform to compile the corresponding plug-ins which adapt to your new environment.

If you only update the previous installation package, you need to uninstall the old installation package first.

### 3. Check

```
>>> from ingenic_magik_trainingkit.QuantizationTrainingPlugin.python import ops
INFO(magik): trainingkit version:1.1.1(00010101_84f712d) built:20220121-1849(5.4.0 pytorch)
>>>
```

If the above figure appears, it indicates that the import is successful. The green font contains the version of the trainingkit, the commit number and the compilation date.

## 3.10 Others

Other packages:

pandas (\$ pip install pandas)

requests (\$ pip install requests)

cv2 (\$ pip install opencv-python)

yaml (\$ pip install pyyaml)

tqdm (\$ pip install tqdm)

matplotlib (\$ pip install matplotlib)

seaborn (\$ pip install seaborn)

If other installation packages are missing at runtime, please install them by “pip install” according to the prompt until normal operation.

## 4. YOLOv5s Training Process

### 4.1 Downloading of code and model

1. Download link of yolov5(official):

<https://github.com/ultralytics/yolov5.git>

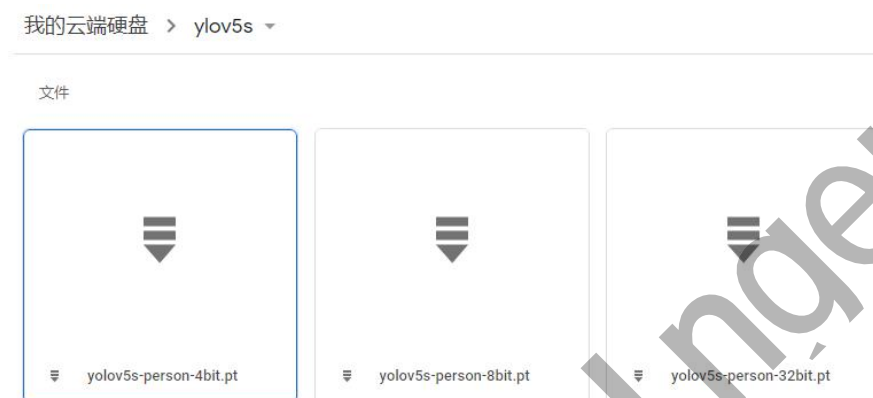
2. Our yolov5 code(modified based on the official)

Since the training plug-in re encapsulates the relevant operators (such as convolution, depthwise-convolution and etc), we will provide a whole set of training code yolov5s-person based on the original yolov5. Our yolov5 code is in the directory “magik-toolkit/Models/training/pytorch/T40/yolov5s-person/”, which the process of training, conversion and operation on board described below are completed in this directory.

### 3. Download address of the model

Link: [https://drive.google.com/drive/folders/1J5KqIN91vpoM8g7xg6\\_7CIWbkQCXR1qe?usp=sharing](https://drive.google.com/drive/folders/1J5KqIN91vpoM8g7xg6_7CIWbkQCXR1qe?usp=sharing)

As shown in the figure below, these three models are 32bit, 8bit and 4bit models trained with 15000 person data of coco2017 step by step, and yolov5s-person-4bit.pt is used for later PC end testing and board end testing, please download and put it under “magik-toolkit/Models/training/pytorch/T40/yolov5s-person/runs/train”.



## 4.2 Data Preparation

Take coco2017 as an example, We need to extract the person (id:0) and convert it into the data format required for Yolo training.

Firstly, please download and decompress the coco2017 data set, the format of the annotation file is json. We are supposed to use wget to download the pictures and labels of coco2017 from the link below.

Pictures:

<http://images.cocodataset.org/zips/train2017.zip>

<http://images.cocodataset.org/zips/test2017.zip>

<http://images.cocodataset.org/zips/val2017.zip>

Labels:

[http://images.cocodataset.org/annotations/stuff\\_annotations\\_trainval2017.zip](http://images.cocodataset.org/annotations/stuff_annotations_trainval2017.zip)

[http://images.cocodataset.org/annotations/image\\_info\\_test2017.zip](http://images.cocodataset.org/annotations/image_info_test2017.zip)

[http://images.cocodataset.org/annotations/annotations\\_trainval2017.zip](http://images.cocodataset.org/annotations/annotations_trainval2017.zip)

Secondly, The script for generating data is in COCO\_forYOLO folder.

```
$ python batch_split_annotation_foryolo.py
```

Pay attention to modify the coco path by setting 'coco\_data\_dir=' in the file.

After running, three folders (person/data/images, person/data/ImageSets and

person/data/labels) will be generated in the coco path, then please put train2017.txt, val2017.txt and test2017.txt under Imagesets into the folder “persondet/data/coco” and using absolute path of the pictures in the txt file.

## 4.3 Network Training

### 1. Training configuration

Please check the default parameters like is\_quantize, bitw, bita and etc in “models/commom.py” for reference.

```
bita = 32
if bita==32:
    bitw = 32
    is_quantize = 0
    clip_max_value = 6.0
    shortcut_clip_max_value = 2.0
elif bita==8:
    bitw = 8
    is_quantize = 1
    clip_max_value = 6.0
    shortcut_clip_max_value = 2.0
elif bita==4:
    bitw = 4
    is_quantize = 1
    clip_max_value = 4.0
    shortcut_clip_max_value = 1.5
weight_factor = 3.0
target_device = "T40"
```

When you want to train model of 32-bit, please set bita to 32. Similarly, you can set bita to 8 if you are going to train model of 8-bit. Of course, 4-bit is in the same way.

Parameter interpretation:

is\_quantize -- quantize or not, note that set it to 0 when bita is 32.

bitw - bitwidth of weight

bita - bitwidth of feature

clip\_max\_value - clip value of feature when quantization

shortut\_clip\_max\_value -- clip value of feature in shortcut when quantization

### 2. Training script

```
$ sh yolov5s-person/train.sh
export NCCL_IB_DISABLE=1
export NCCL_DEBUG=info
GPUS=6
```

```
python3 -m torch.distributed.launch --nproc_per_node=$GPUS
--master_port=60051 train.py \
--data data/coco-person.yaml \
--cfg models/yolov5s.yaml \
--weights '' \
```

```
--batch-size 132 \
--hyp data/hyp.scratch.yaml \
--project ./runs/train/yolov5s-person-32bit \
--epochs 300 \
--device 0,1,2,3,4,5
```

#### (1) About pre training model

There is no pre training model in floating-point training, so -weights is '', 32bit model with sufficient accuracy is loaded in 8bit training, and 8bit trained is loaded in 4bit as the pre training model. In this way, the effect is better step by step.

#### (2) About multi-GPU training

We added “torch.distributed.launch” during training to support multi-GPU training. The value of GPUs corresponds to the total number of devices below, which can be modified according to the actual situation. If it is single gpu training, directly use “python3 train.py” plus the following parameters. Batch size is the total batch number of all GPUs, which can be set according to the actual size of the GPUs.

#### (3) About learning rate

The super parameters are set in file data/hyp.scratch.yaml, lr0 is the initial learning rate, and the rest adopts the default value. We’ve removed the part about loading optimizer in file train.py, the low bit model re training based on 32-bit model instead of subsequent training, so that we can avoid affecting the training effect of low bit. We need to know that “ema” and “epoch” also have similar problems, so the corresponding modifications have been made as shown in the figure below.

```
# Resume
start_epoch, best_fitness = 0, 0.0
if pretrained:
    # Optimizer
    #if ckpt['optimizer'] is not None:
    #    optimizer.load_state_dict(ckpt['optimizer'])
    #    best_fitness = ckpt['best_fitness']

    # EMA
    #if use_ema and ema and ckpt.get('ema'):
    #    ema.ema.load_state_dict(ckpt['model'].float().state_dict())
    #    ema.updates = ckpt['updates']

    # Epochs
    # start_epoch = ckpt['epoch'] + 1

    if resume:
        start_epoch = ckpt['epoch'] + 1
        assert start_epoch > 0, '%s training to %g epochs is finished, nothing to resume.' % (weights, epochs)
```

#### (4) About saving model

We can use “--project” in train.py to set the saved directory about training model, there will be two models under runs/train/project/weights: best.pt and

last.pt, obviously, best.pt is the best training so far, and last.pt the latest training so far. The results of each epoch are saved in result.txt, which can be viewed easily.

#### (5) Some experiences about training

About yov5s-person, our experiences at present are that 32-bit uses sgd and lr0.01, 8bit uses sgd and lr0.01 for pre training based on 32bit, 4bit uses adam and lr0.001 for pre training based on 8bit. Please see train.sh to get specific commands.

Since 32-bit has no pre training model, the convergence is slightly slow; 8bit has pre training and can express, and the convergence is fast; Although 4bit has pre training, its expression ability is slightly weak and its convergence is slightly slow.

## 4.4 Model Testing

### 1. Test picture(s)

Through run “python detect.py -h” to view and select the required parameters. Set “--source” to detect pictures or videos or picture folders. And the detection results can be set to display (- - view-img) or save (- - save-img)

- Image: `--source file.jpg`
- Video: `--source file.mp4`
- Directory: `--source dir/`

The usage is as follows:

```
$ sh detect.sh(python detect.py --source data/images/bus.jpg \
--weights ./runs/train/yolov5s-person-4bit.pt \
--imgs 640 --device 0 --view-img)
```

Parameter interpretation:

- source - picture(s) or video to be tested
- weights - trained model for testing
- imgs - dest size when testing
- device - test with which GPU
- view-img - display image or not

Please make sure the detected model configurations are consistent with the training configuration.

### 2. Test accuracy

```
$ sh test.sh(python test.py --data data/coco-person.yaml \
--weights ./runs/train/yolov5s-person-4bit.pt \
--imgs 640 --device 0 --batch-size 40)
```

The model to be tested is specified by “--weights”, and the verification set is specified by val2017.txt in data/coco-person.yaml, and other parameters are

given according to actual needs.

Our testing accuracy are as follows, with resolution is 640x640.

	Class	Images	Targets	P	R	mAP@0.5	mAP@.5:.95
32bit:	all	5000	11004	0.771	0.615	0.700	0.422
8bit:	all	5000	11004	0.751	0.638	0.706	0.430
4bit:	all	5000	11004	0.786	0.602	0.698	0.419

As you can see, we've also provided yolov5m.yaml and yolov5l.yaml under yolov5s-person/models, so that you can train your own yolov5m or yolov5s like the process of yolov5s we introduced above to meet your needs.

## 5. Model Transformation

### 5.1 From pt to onnx

Please pay attention to that, the onnx file must be generated in the above training environment.

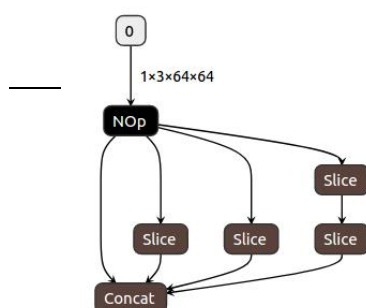
```
$ sh convert_onnx.sh(python convert_onnx.py \
--weights ./runs/train/yolov5s-person-4bit.pt)
```

In this script, input is the path of pt -- the model we are looking forward to convert, and output is the onnx file that we want.

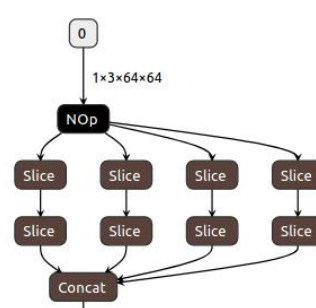
We are supposed to use opset9 when convert pt to onnx, so that we can use transform tool successfully later. If there are some errors occur when you set opset9, Please refer to the question answer in step 8. Different versions of torch will have some different problems.

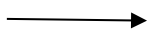
```
def _slice(g, input, axes, starts, ends):
    assert len(starts) == len(ends)
    # if len(starts) == 1 and starts[0] == 0 and ends[0] == 9223372036854775807:
    #     return input
    return g.op("Slice", input, axes_i=axes, starts_i=starts, ends_i=ends)
```

In particular, for yolov5, we need to modify function \_slice() in torch/onnx/symbolic\_opset9.py like above picture before using the focus operator to onnx. Or we will get wrong focus operation, we can do visual view with netron as follows.



iconductor Co. Ltc





before

after

## 5.2 From onnx to bin

This step does not depend on python, torch, plug-ins or other environments, but only our conversion tool “magik-transform-tools”(which is under magik-toolkit/TransformKit) and run.sh.

```
$ cd transform
```

```
$ sh run.sh
```

Specifically,

```
path=../../../../../TransformKit
```

```
$ path/magik-transform-tools \
```

```
--framework onnx \
```

```
--target_device T40 \
```

```
--outputpath yolov5s-person-4bit.mk.h \
```

```
--inputpath ../runs/train/yolov5s-person-4bit.onnx \
```

```
--mean 0,0,0 \
```

```
--var 255,255,255 \
```

```
--img_width 416 \
```

```
--img_height 416 \
```

```
--img_channel 3
```

Parameter interpretation:

inputpath - onnx file we got before

outputpath - output of transformation, we will get yolov5s-person-4bit.bin in this directory after running

target\_device - set same as target\_device during training

mean - mean value of preprocessing during training

var - variance of preprocessing during training

img\_width - input's width of running on-board

img\_height - input's height of running on-board

img\_channel - input's channels of running on-board



```

INFO(magik): magik-transform-tools version:1.0.0(00010000 bbb8835) cuda version:9.0.176 built:20220214-1120 GPU
2022-02-15 14:53:37.819065: W clip_model_by_nop_optimizer.cc:48-collect_remove_nodes] Multiple output NOP exists on the current pathway!
2022-02-15 14:53:37.820572: W clip_model_by_nop_optimizer.cc:48-collect_remove_nodes] Multiple output NOP exists on the current pathway!
2022-02-15 14:53:37.820625: W clip_model_by_nop_optimizer.cc:48-collect_remove_nodes] Multiple output NOP exists on the current pathway!
2022-02-15 14:53:37.820665: W clip_model_by_nop_optimizer.cc:48-collect_remove_nodes] Multiple output NOP exists on the current pathway!
2022-02-15 14:54:00.861688: W quantize_operation.cc:305-search] Not found QuantizeOperation function for Unpool2D
2022-02-15 14:54:00.875195: W quantize_operation.cc:305-search] Not found QuantizeOperation function for Unpool2D
*****
*               ^ ^ Convert successfully, Enjoy it ^ ^               *
*****

```

## 6. Run on-board

### 6.1 Code compiling

#### 1. Preparation

First, download address of opencv library compiled by our mips

Link:<https://drive.google.com/drive/folders/1Gdj4vF2TpaPzbPNT0H1hIEKD TXEj92Jg?usp=sharing>

我的云端硬盘 > opencv

文件



Please download and decompress it into magik-toolkit/ThirdParty.

```
$ tar -xvf opencv_mips_7.2.0_2.29.tar.gz
```

Secondly, we also provide inference.cpp, Makfile and test data, which are necessary when run model on board. In addition, we also need venus library and mips compilation tool. You may need to know that venus library is under the library “magik-toolkit/InferenceKit/nnal”, and mips compilation tool is provided by our solutions’ colleagues.

#### 2. Input of network

Here, in order to all of you can realize the process quickly, we’ve added opencv library functions here to facilitate calling “imread()” and other functions, so we can test jpg diagram directly on board.



### 3. Model loading

The model in inference.cpp is passed in through parameters. So please pay attention to synchronously copying it to the corresponding directory on the board end and passing it in at run time.

### 4. Setting of super parameters

```
void generateBBox(std::vector<venus::Tensor> out_res, std::vector<magik::venus::ObjBbox_t>& candidate_boxes)
{
    float person_threshold = 0.3;
    int classes = 1;
    float nms_threshold = 0.6;
    std::vector<float> strides = {8.0, 16.0, 32.0};
    int box_num = 3;
    std::vector<float> anchor = {10,13, 16,30, 33,23, 30,61, 62,45, 59,119, 116,90, 156,198, 373,326};
}
```

These settings should be consistent with the test code on the PC end.

### 5. Compile

TOPDIR -- directory of venus

opencv\_path -- opencv library compiled by our mips

libtype - muclibc or other(mglibc)

build\_type -- release (Default setting.)

-- profile (Visualization of network structure at runtime and statistics of running time and GOPs at each layer of the network.)

-- debug (Save the results of quantization features of each layer.)

-- nmem (Count the memory usage of nmem when the model is running, and it's in /tmp/nmem\_memory.txt.)

Directly compile inference.cpp can generate venus\_yolov5s\_bin\_uclibc\_\*, and that is the executable file we need on the board.

We also provide the code inference\_nv12.cpp and input data 10\_w1024\_h714\_nv12 when you want to use nv12 type input. If necessary, just modify Makefile for compilation and use test.

## 6.2 Run

Opencv library:

magik-toolkit/ThirdParty/opencv\_mips\_7.2.0\_2.29/uclibc/lib

Venus library:

magik-toolkit/InferenceKit/nnal/mips720-glibc229/lib/uclibc/

#### 1. release(default)

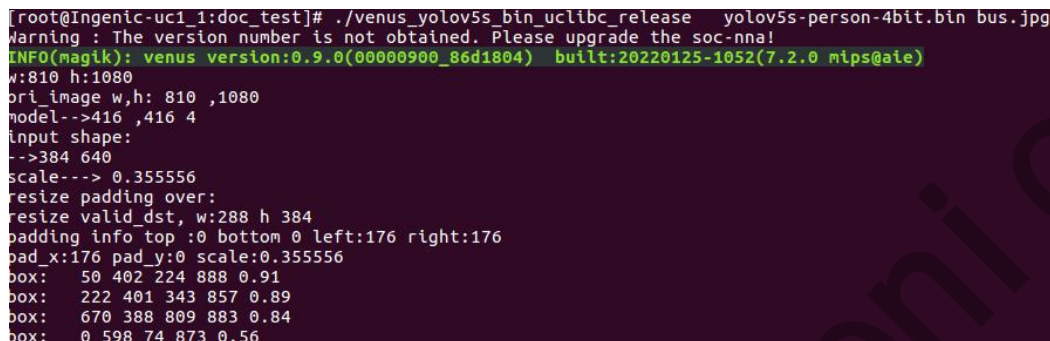
\$ make build\_type=release

(use "make build\_type=release clean" if you want to clean)

We can get executable file venus\_yolov5s\_bin\_uclibc\_release after compile, then copy libvenus.so, opencv library(which consists of libopencv\_core.so.3.3, libopencv\_highgui.so.3.3, libopencv\_imgcodecs.so.3.3, libopencv\_imgproc.so.3.3, libopencv\_videoio.so.3.3), venus\_yolov5s\_bin\_uclibc\_release, yolov5s-person-

4bit.bin and bus.jpg to board.

When you enter the end of the board, please add library path first(`export LD_LIBRARY_PATH=$lib_path:$LD_LIBRARY_PATH`), after that , you just need run on board by command “`./venus_yolov5s_bin_uclibc_release yolov5s-person-4bit.bin bus.jpg`”. The operation results are shown in the figure below.



```
[root@Ingenic-uc1_1:doc_test]# ./venus_yolov5s_bin_uclibc_release yolov5s-person-4bit.bin bus.jpg
Warning : The version number is not obtained. Please upgrade the soc-nna!
INFO(magik): venus version:0.9.0(00000900_86d1804) built:20220125-1052(7.2.0 mips@ale)
w:810 h:1080
ori_image w,h: 810 ,1080
model-->416 ,416 4
input shape:
-->384 640
scale--> 0.355556
resize padding over:
resize valid_dst, w:288 h 384
padding info top :0 bottom 0 left:176 right:176
pad_x:176 pad_y:0 scale:0.355556
box: 50 402 224 888 0.91
box: 222 401 343 857 0.89
box: 670 388 809 883 0.84
box: 0 598 74 873 0.56
```

## 2. debug(check data)

See step 7.2 for details. The input processing is a little different, and we do not need opencv library this step.

## 3. profile(network visualization)

`$ make build_type=profile`

(use “`make build_type=profile clean`” if you want to clean)

We can get executable file `venus_yolov5s_bin_uclibc_profile` after compile, then copy `libvenus.p.so`, `opencv` library(which consists of `libopencv_core.so.3.3`, `libopencv_highgui.so.3.3`, `libopencv_imgcodecs.so.3.3`, `libopencv_imgproc.so.3.3`, `libopencv_videoio.so.3.3`), `venus_yolov5s_bin_uclibc_profile`, `yolov5s-person-4bit.bin` and `bus.jpg` to board.

When you enter the end of the board, please add library path first(`export LD_LIBRARY_PATH=$lib_path:$LD_LIBRARY_PATH`), after that , you just need run on board by command “`./venus_yolov5s_bin_uclibc_profile yolov5s-person-4bit.bin bus.jpg`”. The operation results are shown in the figure below.

```

INFO(magik): venus version:0.9.0(00000900_86d1804) built:20220125-1057(7.2.0 mips@ale)
w:810 h:1080
ori_image w,h: 810 ,1080
model-->416 ,416 4
input shape:
-->384 640
scale---> 0.355556
resize padding over:
resize valid_dst, w:288 h 384
padding info top :0 bottom 0 left:176 right:176
pad_x:176 pad_y:0 scale:0.355556
box: 50 402 224 888 0.91
box: 222 401 343 857 0.89
box: 670 388 809 883 0.84
box: 0 598 74 873 0.56
[Peak:venus]
Timing cycle = 10
===== Detailed DispatchProfiler Summary: N/A, Exclude Swarn-ups =====
LayerName OperatorType InputShape OutputShape FilterShape Stride Pad Dila Avg(ns) Max(ns) Min(ns) Last(ns) Avg(S) Max(S) Min(S) Bandwidth(M)
layer_2.QuantizeFeature Convolution (1,384,640,4) (1,192,320,32) (3,3,4,32) (1,1) (1,1,1,1) (1,1) 0.594 0.612 0.567 0.5070 0.56% 0.142 2.817
layer_4.QuantizeFeature Convolution (1,192,320,32) (1,96,160,64) (3,3,32,64) (2,2) (0,1,0,1) (1,1) 2.132 2.171 2.121 2.1210 0.26% 0.566 2.831
layer_6.QuantizeFeature Convolution (1,96,160,64) (1,96,160,32) (1,1,64,32) (1,1) (0,0,0,0) (1,1) 0.081 0.114 0.071 0.1140 0.64% 0.063 1.173
layer_12.QuantizeFeature Convolution (1,96,160,32) (1,96,160,32) (1,1,32,32) (1,1) (0,0,0,0) (1,1) 0.253 0.254 0.252 0.2540 0.98% 0.031 0.469
layer_14.QuantizeFeature Convolution (1,96,160,32) (1,96,160,32) (3,3,32,32) (1,1) (1,1,1,1) (1,1) 0.190 0.191 0.190 0.1900 0.74% 0.283 0.473
layer_15.QuantizeFeature Add (1,96,160,32) (1,96,160,32) N/A N/A N/A 0.194 0.195 0.194 0.1940 0.73% 0.000 0.703
layer_8.QuantizeFeature Convolution (1,96,160,64) (1,96,160,32) (1,1,64,32) (1,1) (0,0,0,0) (1,1) 0.065 0.067 0.064 0.0670 2.58% 0.063 1.173
layer_10.QuantizeFeature Convolution (1,96,160,64) (1,96,160,64) (1,1,64,64) (1,1) (0,0,0,0) (1,1) 0.521 0.522 0.520 0.5210 2.02% 0.126 0.940
layer_17.QuantizeFeature Convolution (1,96,160,64) (1,48,80,128) (3,3,64,128) (2,2) (0,1,0,1) (1,1) 0.587 0.615 0.579 0.5790 2.27% 0.566 0.739
layer_19.QuantizeFeature Convolution (1,48,80,128) (1,48,80,64) (1,1,128,64) (1,1) (0,0,0,0) (1,1) 0.186 0.170 0.165 0.1660 0.65% 0.063 0.356
layer_25.QuantizeFeature Convolution (1,48,80,64) (1,48,80,64) (1,1,64,64) (1,1) (0,0,0,0) (1,1) 0.146 0.147 0.146 0.1460 0.57% 0.031 0.237
layer_27.QuantizeFeature Convolution (1,48,80,64) (1,48,80,64) (3,3,64,64) (1,1) (1,1,1,1) (1,1) 0.172 0.172 0.171 0.1710 0.67% 0.283 0.257
layer_28.QuantizeFeature Add (1,48,80,64) (1,48,80,64) N/A N/A N/A 0.105 0.105 0.105 0.1050 0.41% 0.000 0.352
layer_30.QuantizeFeature Convolution (1,48,80,64) (1,48,80,64) (1,1,64,64) (1,1) (0,0,0,0) (1,1) 0.152 0.154 0.150 0.1500 0.59% 0.031 0.237
layer_32.QuantizeFeature Convolution (1,48,80,64) (1,48,80,64) (3,3,64,64) (1,1) (1,1,1,1) (1,1) 0.162 0.162 0.161 0.1620 0.63% 0.283 0.252
layer_33.QuantizeFeature Add (1,48,80,64) (1,48,80,64) N/A N/A N/A 0.104 0.105 0.104 0.1040 0.40% 0.000 0.352
layer_35.QuantizeFeature Convolution (1,48,80,64) (1,48,80,64) (1,1,64,64) (1,1) (0,0,0,0) (1,1) 0.151 0.154 0.148 0.1480 0.59% 0.031 0.237
layer_37.QuantizeFeature Convolution (1,48,80,64) (1,48,80,64) (3,3,64,64) (1,1) (1,1,1,1) (1,1) 0.161 0.162 0.161 0.1610 0.63% 0.283 0.252
layer_38.QuantizeFeature Add (1,48,80,64) (1,48,80,64) N/A N/A N/A 0.104 0.105 0.104 0.1040 0.40% 0.000 0.352
layer_21.QuantizeFeature Convolution (1,48,80,128) (1,48,80,64) (1,1,128,64) (1,1) (0,0,0,0) (1,1) 0.163 0.163 0.162 0.1630 0.63% 0.063 0.356
layer_23.QuantizeFeature Convolution (1,48,80,128) (1,48,80,128) (1,1,128,128) (1,1) (0,0,0,0) (1,1) 0.268 0.262 0.257 0.2610 1.01% 0.126 0.478
layer_40.QuantizeFeature Convolution (1,48,80,128) (1,24,40,256) (3,3,128,256) (2,2) (0,1,0,1) (1,1) 0.547 0.578 0.539 0.5400 2.12% 0.566 0.494
layer_42.QuantizeFeature Convolution (1,24,40,256) (1,24,40,128) (1,1,256,128) (1,1) (0,0,0,0) (1,1) 0.093 0.100 0.090 0.0900 0.36% 0.063 0.192
layer_48.QuantizeFeature Convolution (1,24,40,128) (1,24,40,128) (1,1,128,128) (1,1) (0,0,0,0) (1,1) 0.080 0.082 0.078 0.0780 0.31% 0.031 0.126
layer_50.QuantizeFeature Convolution (1,24,40,128) (1,24,40,128) (3,3,128,128) (1,1) (1,1,1,1) (1,1) 0.174 0.176 0.173 0.1740 0.67% 0.283 0.188
layer_51.QuantizeFeature Add (1,24,40,128) (1,24,40,128) N/A N/A N/A 0.066 0.067 0.066 0.0660 0.26% 0.000 0.176
layer_53.QuantizeFeature Convolution (1,24,40,128) (1,24,40,128) (1,1,128,128) (1,1) (0,0,0,0) (1,1) 0.083 0.084 0.083 0.0830 0.32% 0.031 0.126
layer_55.QuantizeFeature Convolution (1,24,40,128) (1,24,40,128) (3,3,128,128) (1,1) (1,1,1,1) (1,1) 0.170 0.171 0.170 0.1700 0.66% 0.283 0.188
layer_56.QuantizeFeature Add (1,24,40,128) (1,24,40,128) N/A N/A N/A 0.066 0.066 0.065 0.0660 0.25% 0.000 0.176
layer_58.QuantizeFeature Convolution (1,24,40,128) (1,24,40,128) (1,1,128,128) (1,1) (0,0,0,0) (1,1) 0.082 0.082 0.081 0.0820 0.32% 0.031 0.126
layer_60.QuantizeFeature Convolution (1,24,40,128) (1,24,40,128) (3,3,128,128) (1,1) (1,1,1,1) (1,1) 0.170 0.171 0.169 0.1690 0.66% 0.283 0.188
layer_61.QuantizeFeature Add (1,24,40,128) (1,24,40,128) N/A N/A N/A 0.073 0.066 0.065 0.0650 0.28% 0.000 0.176
layer_44.QuantizeFeature Convolution (1,24,40,256) (1,24,40,128) (1,1,256,128) (1,1) (0,0,0,0) (1,1) 0.091 0.098 0.089 0.0900 0.35% 0.063 0.192
layer_46.QuantizeFeature Convolution (1,24,40,256) (1,24,40,256) (1,1,256,256) (1,1) (0,0,0,0) (1,1) 0.139 0.140 0.138 0.1390 0.54% 0.126 0.266
layer_63.QuantizeFeature Convolution (1,24,40,256) (1,12,20,512) (3,3,256,512) (2,2) (0,1,0,1) (1,1) 0.559 0.553 0.549 0.5490 2.13% 0.566 0.742
layer_65.QuantizeFeature Convolution (1,12,20,512) (1,12,20,256) (1,1,512,256) (1,1) (0,0,0,0) (1,1) 0.076 0.077 0.076 0.0760 0.30% 0.063 0.152
518 Pooling (1,12,20,256) (1,12,20,256) (5,5) (1,1) (1,1,1,1) (1,1) 0.391 0.391 0.390 0.3910 1.17% 0.032 0.652

```

#### 4. nmem(nmem occupancy)

\$ make build\_type=nmem

(use “make build\_type=nmemclean” if you want to clean)

We can get executable file venus\_yolov5s\_bin\_uclibc\_release after compile, then copy libvenus.m.so, opencv library(which consists of libopencv\_core.so.3.3, libopencv\_highgui.so.3.3, libopencv\_imgcodecs.so.3.3, libopencv\_imgproc.so.3.3, libopencv\_videoio.so.3.3), venus\_yolov5s\_bin\_uclibc\_nmem, yolov5s-person-4bit.bin and bus.jpg to board.

When you enter the end of the board, please add library path first(export LD\_LIBRARY\_PATH=\$lib\_path:\$LD\_LIBRARY\_PATH), after that , you just need run on board by command “./venus\_yolov5s\_bin\_uclibc\_nmem yolov5s-person-4bit.bin bus.jpg”. The operation results are shown in the figure below.

```

[root@Ingenic-uc1 1:doc_test]# ./venus_yolov5s_bin_uclibc_nmem yolov5s-person-4bit.bin bus.jpg
Warning : The version number is not obtained. Please upgrade the soc-nna!
INFO(magik): venus version:0.9.0(00000900_86d1804) built:20220125-1045(7.2.0 mips@ale)
w:810 h:1080
ori_image w,h: 810 ,1080
model-->416 ,416 4
input shape:
-->384 640
scale---> 0.355556
resize padding over:
resize valid_dst, w:288 h 384
padding info top :0 bottom 0 left:176 right:176
pad_x:176 pad_y:0 scale:0.355556
box: 50 402 224 888 0.91
box: 222 401 343 857 0.89
box: 670 388 809 883 0.84
box: 0 598 74 873 0.56
[root@Ingenic-uc1 1:doc_test]# cat /tmp/nmem_memory.txt
nmem total = 131072K
nmem used = 13024K
nmem free = 118048K
### %9.94 ###

```

## 7. Check Data

To verify whether the data of PC end and board end can be aligned, you can do the following:

### 1. PC

When we test single picture like Step 4.4-1, we can add the environment variable “MAGIK\_TRAININGKIT\_DUMP=1” to save the quantization results of each layer to “/tmp/trainingkit\_data/feature”, what’s more, you can also set the environment variable “MAGIK\_TRAININGKIT\_PATH” specify the directory to save. The specific running command is as below.

```
$ MAGIK_TRAININGKIT_DUMP=1  
MAGIK_TRAININGKIT_PATH="/tmp" python detect.py \  
--source data/images/bus.jpg \  
--weights ./runs/train/yolov5s-person-4bit.pt \  
--imgs 640 --device 0
```

The original resolution of bus.jpg is 810x1080. During the test, the target size “--imgs” is set to 640. According to the scaling principle of yolov5 code (the long side is 640, the short side is scaled in equal proportion, and then filled to a multiple of 32), so the resolution of the final test is 480x640. The operation results of this network are finally saved under “./trainingkit\_data/feature” (see the figure below for details), it can be seen that the input layer is saved as input\_data\_shape\_1\_640\_480\_3.bin, among them, height is 640, width is 480, channel is 3. As you can see in the picture, we also save the last three output layers(out\_input\_index\_\*.bin), which can be used to check the network output.

```
input_data_shape_1_640_480_3.bin layer_44 QuantizeFeature.bin  
layer_101 QuantizeFeature.bin layer_46 QuantizeFeature.bin  
layer_103 QuantizeFeature.bin layer_48 QuantizeFeature.bin  
layer_106 QuantizeFeature.bin layer_4 QuantizeFeature.bin  
layer_108 QuantizeFeature.bin layer_50 QuantizeFeature.bin  
layer_10 QuantizeFeature.bin layer_51 QuantizeFeature.bin  
layer_110 QuantizeFeature.bin layer_53 QuantizeFeature.bin  
layer_112 QuantizeFeature.bin layer_55 QuantizeFeature.bin  
layer_114 QuantizeFeature.bin layer_56 QuantizeFeature.bin  
layer_116 QuantizeFeature.bin layer_58 QuantizeFeature.bin  
layer_119 QuantizeFeature.bin layer_60 QuantizeFeature.bin  
layer_121 QuantizeFeature.bin layer_61 QuantizeFeature.bin  
layer_123 QuantizeFeature.bin layer_63 QuantizeFeature.bin  
layer_125 QuantizeFeature.bin layer_65 QuantizeFeature.bin  
layer_127 QuantizeFeature.bin layer_67 QuantizeFeature.bin  
layer_129 QuantizeFeature.bin layer_69 QuantizeFeature.bin  
layer_12 QuantizeFeature.bin layer_6 QuantizeFeature.bin  
layer_14 QuantizeFeature.bin layer_71 QuantizeFeature.bin  
layer_15 QuantizeFeature.bin layer_73 QuantizeFeature.bin  
layer_17 QuantizeFeature.bin layer_75 QuantizeFeature.bin  
layer_19 QuantizeFeature.bin layer_77 QuantizeFeature.bin  
layer_21 QuantizeFeature.bin layer_80 QuantizeFeature.bin  
layer_23 QuantizeFeature.bin layer_82 QuantizeFeature.bin  
layer_25 QuantizeFeature.bin layer_84 QuantizeFeature.bin  
layer_27 QuantizeFeature.bin layer_86 QuantizeFeature.bin  
layer_28 QuantizeFeature.bin layer_88 QuantizeFeature.bin  
layer_2 QuantizeFeature.bin layer_8 QuantizeFeature.bin  
layer_30 QuantizeFeature.bin layer_90 QuantizeFeature.bin  
layer_32 QuantizeFeature.bin layer_93 QuantizeFeature.bin  
layer_33 QuantizeFeature.bin layer_95 QuantizeFeature.bin  
layer_35 QuantizeFeature.bin layer_97 QuantizeFeature.bin  
layer_37 QuantizeFeature.bin layer_99 QuantizeFeature.bin  
layer_38 QuantizeFeature.bin output_index_1_shape_1_80_60_18.bin  
layer_40 QuantizeFeature.bin output_index_2_shape_1_40_30_18.bin  
layer_42 QuantizeFeature.bin output_index_3_shape_1_20_15_18.bin
```



## 2. Board

The input data required at the board end is 4 channels. For the bus.jpg tested at the PC end above, we process input\_data\_shape\_1\_640\_480\_3.bin(which saved at run time) to header file(\*.h), and then use it to test the board end to ensure that the inputs of the two are completely consistent. You can do offline processing first, you can get the processing script is in the “yolov5s-person/venus\_sample\_yolov5s/generate\_img\_input.py” .

Usage:

```
$ python generate_img_input.py bin_path w h c bin
```

```
($ python generate_img_input.py input_data_shape_1_640_480_3.bin \
  480 640 3 bin)
```

bin\_path -- the input bin file to be converted

w, h -- width and height of the bin file, you can get them from the name

bin -- flag, which indicates that the bin file turns to the header file, and another one is “img” means image to header file

You'd better clear the previously compiled mode by “make build\_type=\* clean” before compiling, otherwise, the link library of OpenCV may fail. Then compile it by command “make build\_type=debug”, We can get executable file venus\_yolov5s\_bin\_uclibc\_debug after compile, then copy libvenus.d.so, venus\_yolov5s\_bin\_uclibc\_debug, yolov5s-person-4bit.bin to board.

When you enter the end of the board, please add library path first(export LD\_LIBRARY\_PATH=\$lib\_path:\$LD\_LIBRARY\_PATH), after that , you just need run on board by command “./venus\_yolov5s\_bin\_uclibc\_debug yolov5s-person-4bit.bin 480 640”, among them, 480 is width, 640 is height. The operation results are shown in the figure below

```

layer_116_QuantizeFeature_bt.bin layer_27_QuantizeFeature_out.bin layer_53_QuantizeFeature_weight.bin
layer_116_QuantizeFeature_out.bin layer_27_QuantizeFeature_weight.bin layer_55_QuantizeFeature_bt.bin
layer_116_QuantizeFeature_weight.bin layer_28_QuantizeFeature_out.bin layer_55_QuantizeFeature_out.bin
layer_119_QuantizeFeature_bt.bin layer_2_QuantizeFeature_bt.bin layer_55_QuantizeFeature_weight.bin
layer_119_QuantizeFeature_out.bin layer_2_QuantizeFeature_out.bin layer_56_QuantizeFeature_out.bin
layer_119_QuantizeFeature_weight.bin layer_2_QuantizeFeature_weight.bin layer_58_QuantizeFeature_bt.bin
layer_121_QuantizeFeature_bt.bin layer_30_QuantizeFeature_bt.bin layer_58_QuantizeFeature_out.bin
layer_121_QuantizeFeature_out.bin layer_30_QuantizeFeature_out.bin layer_58_QuantizeFeature_weight.bin
layer_121_QuantizeFeature_weight.bin layer_30_QuantizeFeature_weight.bin layer_60_QuantizeFeature_bt.bin
layer_123_QuantizeFeature_bt.bin layer_32_QuantizeFeature_bt.bin layer_60_QuantizeFeature_out.bin
layer_123_QuantizeFeature_out.bin layer_32_QuantizeFeature_out.bin layer_60_QuantizeFeature_weight.bin
layer_123_QuantizeFeature_weight.bin layer_32_QuantizeFeature_weight.bin layer_61_QuantizeFeature_out.bin
layer_125_QuantizeFeature_bt.bin layer_33_QuantizeFeature_out.bin layer_63_QuantizeFeature_bt.bin
layer_125_QuantizeFeature_out.bin layer_35_QuantizeFeature_bt.bin layer_63_QuantizeFeature_out.bin
layer_125_QuantizeFeature_weight.bin layer_35_QuantizeFeature_out.bin layer_63_QuantizeFeature_weight.bin
layer_127_QuantizeFeature_out.bin layer_35_QuantizeFeature_weight.bin layer_65_QuantizeFeature_bt.bin
layer_127_QuantizeFeature_bt.bin layer_37_QuantizeFeature_bt.bin layer_65_QuantizeFeature_out.bin
layer_129_QuantizeFeature_bt.bin layer_37_QuantizeFeature_out.bin layer_65_QuantizeFeature_weight.bin
layer_129_QuantizeFeature_out.bin layer_37_QuantizeFeature_weight.bin layer_67_QuantizeFeature_bt.bin
layer_129_QuantizeFeature_weight.bin layer_38_QuantizeFeature_out.bin layer_67_QuantizeFeature_out.bin
layer_129_QuantizeFeature_weight.bin layer_40_QuantizeFeature_bt.bin layer_67_QuantizeFeature_weight.bin
layer_12_QuantizeFeature_bt.bin layer_40_QuantizeFeature_out.bin layer_69_QuantizeFeature_bt.bin
layer_12_QuantizeFeature_out.bin layer_40_QuantizeFeature_weight.bin layer_69_QuantizeFeature_out.bin
layer_12_QuantizeFeature_weight.bin layer_42_QuantizeFeature_bt.bin layer_69_QuantizeFeature_weight.bin
layer_14_QuantizeFeature_bt.bin layer_42_QuantizeFeature_out.bin layer_6_QuantizeFeature_bt.bin
layer_14_QuantizeFeature_weight.bin layer_42_QuantizeFeature_weight.bin layer_6_QuantizeFeature_out.bin
layer_15_QuantizeFeature_out.bin layer_44_QuantizeFeature_bt.bin layer_6_QuantizeFeature_weight.bin
layer_17_QuantizeFeature_bt.bin layer_44_QuantizeFeature_out.bin layer_71_QuantizeFeature_bt.bin
layer_17_QuantizeFeature_out.bin layer_44_QuantizeFeature_weight.bin layer_71_QuantizeFeature_out.bin
layer_17_QuantizeFeature_weight.bin layer_46_QuantizeFeature_bt.bin layer_71_QuantizeFeature_weight.bin
layer_19_QuantizeFeature_bt.bin layer_46_QuantizeFeature_out.bin layer_73_QuantizeFeature_bt.bin
layer_19_QuantizeFeature_weight.bin layer_46_QuantizeFeature_weight.bin layer_73_QuantizeFeature_out.bin
layer_21_QuantizeFeature_bt.bin layer_48_QuantizeFeature_bt.bin layer_73_QuantizeFeature_weight.bin
layer_21_QuantizeFeature_out.bin layer_48_QuantizeFeature_out.bin layer_75_QuantizeFeature_bt.bin
layer_21_QuantizeFeature_weight.bin layer_48_QuantizeFeature_weight.bin layer_75_QuantizeFeature_out.bin
layer_23_QuantizeFeature_bt.bin layer_4_QuantizeFeature_bt.bin layer_75_QuantizeFeature_weight.bin
layer_23_QuantizeFeature_out.bin layer_4_QuantizeFeature_out.bin layer_77_QuantizeFeature_bt.bin
layer_23_QuantizeFeature_weight.bin layer_4_QuantizeFeature_weight.bin layer_77_QuantizeFeature_out.bin
layer_25_QuantizeFeature_bt.bin layer_50_QuantizeFeature_bt.bin layer_77_QuantizeFeature_weight.bin
layer_25_QuantizeFeature_out.bin layer_50_QuantizeFeature_out.bin layer_80_QuantizeFeature_bt.bin
layer_25_QuantizeFeature_weight.bin layer_50_QuantizeFeature_weight.bin layer_80_QuantizeFeature_out.bin
layer_27_QuantizeFeature_bt.bin layer_51_QuantizeFeature_out.bin layer_80_QuantizeFeature_weight.bin
layer_27_QuantizeFeature_out.bin layer_53_QuantizeFeature_bt.bin layer_82_QuantizeFeature_bt.bin
layer_27_QuantizeFeature_weight.bin layer_53_QuantizeFeature_out.bin layer_82_QuantizeFeature_out.bin

```

Where layer\*\_QuantizeFeature\_out.bin on board and layer\*\_QuantizeFeature.bin saved on PC , They correspond to each other one by one, You just need to check whether the md5 value of the two files is consistent.

If the input is completely consistent, but there are two layers that do not correspond, please give feedback in time.

```

86859cb6efb8cfd932dbb57e32a260e3 /tmp/trainingkit_data/feature/layer_2_QuantizeFeature.bin
6e2c04cf0e560cbd3b76860a23132053 /tmp/trainingkit_data/feature/layer_129_QuantizeFeature.bin

[root@Ingenic-uc1_1:doc_test]# md5sum layer_2_QuantizeFeature_out.bin
86859cb6efb8cfd932dbb57e32a260e3 layer_2_QuantizeFeature_out.bin
[root@Ingenic-uc1_1:doc_test]# md5sum layer_129_QuantizeFeature_out.bin
6e2c04cf0e560cbd3b76860a23132053 layer_129_QuantizeFeature_out.bin

```

In the debug mode, we only check the network results here, and there is no post-processing part, so there is not any results.

## 8. Q & A

1. Q: Why is the original resolution 810\*1080 of bus.jpg 640\*384 in release and 480\*640 in debug?

A: When we set the target resolution to 640, the scaling rule at the board end is to fill width to 640 after scaling in equal proportion (of course that can be modified), where 640 is width and 384 is height. But the rule in yolov5 code is to scale the long side by 640, and then fill the short side to a multiple of 32, so here 640 is height and 480 is width, and debug is consistent with the PC end. It can be

set as needed in actual use.

2. Q: Similar errors occurred like “undefined symbol:

`_ZN6caffe28TypeMeta21_typeMetaInlineDataIN3c108BFloat16EEEEPKNS_6detail12TypeMetaInlineDataEv`” while importing the magik package?

A: The version of torch in the current environment does not correspond to that of magik package.

3. Q: “ImportError: version ‘libcudart.so.10.1’ not found” ?

A: The version of cuda in the current environment does not correspond to that of magik package.

4. Q: Can the network with magik plug-ins load the native torch model for pre training?

A: The theory is OK as long as the parameters and their names can correspond one by one. However, the 32bit model based on the native float still needs to be fine tuned and trained. It is not recommended to directly load the native 8-bit model.

5. Q: It appears an error “RuntimeError: Address already in use” when training?

A: Specify an unused port such as “--master\_port=60053” after “python -m torch.distributed.launch”.

6. Q: When training, there is “expected scalar type Float but found Half” ?

A: You may use “amp.autocast” in torch to accelerates semi-precision training, but quantization does not support semi-precision training temporarily. Just comment out this. In addition, the original yolov5 saves the model according to half(). At present, it is found that the accuracy of the last saved model will be lost by quantifying, so it is recommended to remove it. In fact, we’ve done this in the code we provide.

7. Q: A error “step!=1 is currently not supported” when convert pt to onnx?

A: When we use Focus in yolov5, we may get this error after torch1.7, Please comment out the line corresponding to the error report(if + raise RuntimeError) in torch/onnx/symbolic\_opset9.py.

8. Q: Error like “RuntimeError:input\_shape\_value==

`reshape_value || input_shape_value==1 || reshape_value == 1 INTERNAL ASSERT FAILED .....`” when convert pt to onnx?

A: You may get this error after torch1.9 in yolov5, please modify the last parameter “onnx\_shape\_inference=True” in function `_export()` in file torch/onnx/utils.py to “onnx\_shape\_inference=False”. Different versions may be on different lines, you just need to find out that.

9. Q: Warning like “Warning: Unsupported operator NOp. No schema registered for this operator.” when convert pt to onnx after torch1.8?

A: It doesn't matter.

10. Q: An error occurs “After checking the input dimension, the input of concat cannot be concatenated!” when convert model from onnx to bin?

A: The problem about slice, please see step 5.1 to modify that (if+return).

11. Q: Errors like “error: 'constexpr bool std::isinf(double)' conflicts with a previous declaration” occurred while compiling C code?

A: The compiler of mips conflicts with the built-in compiler of X86, please make sure the environment variable `CPLUS_INCLUDE_PATH` about “/usr/include/x86\_64-linux-gnu” is removed.

12. Q: “make: mips-linux-gnu-g++: Command not found” while compiling C code?

A: Mips compiler is not specified, please set it by export.

13. Q: What is the loss of board end accuracy and PC end accuracy?

A: Theoretically, the loss is very small(within one percentage point). If there is a large gap, you can check whether the inputs are consistent, the training and verification data is BGR or RGB, the threshold setting is consistent or not. You can also compare the network output results to see whether the output is consistent. See Step 7 for details.

14. Q: Does the resolution of onnx, conversion tool and board end test need to be same?

A: It is better to keep consistent, especially the latter two. The reason why they are not consistent in this example is that they have no impact on this. For other non full convolution networks, it may affect the converted model. Therefore, when testing the accuracy of the verification set, it is recommended that the input data be processed offline to the same size, and then the PC end and board end are compared.