

君正®

Magik 训练量化使用指南

Date: Feb. 2022

君正®

Magik训练量化

使用指南

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

Release history

Date	Author	Revision	Change
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**

目录

1. 简介	1
2. 环境安装要求	1
3. 环境安装步骤	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 训练插件	5
3.10 其他	6
4. YOLOv5s 训练流程	6
4.1 代码和模型的下载	6
4.2 数据准备	6
4.3 网络训练	7
4.4 模型测试	9
5. 模型转换	10
5.1 pt 转 onnx	10
5.2 onnx 转 bin	11
6. 模型上板	13
6.1 代码编译	13
6.2 上板运行	错误！未定义书签。
7. 数据核对	16
1. PC 端	16
2. 板端	17
8. 常见疑问解答	18

1. 简介

本指南主要针对使用君正处理芯的 Magik 平台的新手，这里以 pytorch 框架下的 YOLOv5s 的 person 为例(target_devic 为 T40)，从环境搭建、数据准备、网络训练、模型转换到最终的上板运行进行全流程的详细介绍，旨在引导使用者熟悉 Magik 的使用方法，进行实现上板流程。

2. 环境安装要求

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

若上述环境均已具备(版本可不作固定要求)，可直接转到步骤 3.9。

3. 环境安装步骤

3.1 Ubuntu

- 装机要求: Ubuntu16.04

3.2 GCC

- 版本要求: 5.4.0
- 具体步骤
 1. 查看当前系统 Ubuntu16.04 的原装 GCC 版本:
\$ gcc -v
 2. 下载: <http://ftp.gnu.org/gnu/gcc>
 3. 编译安装:
\$ 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. 检查:
再次 gcc -v 查看版本
```

3.3 Nvidia Driver

- 版本要求: ≥ 440
- 具体步骤
 1. 下载驱动程序
英伟达官网地址:
<http://www.nvidia.cn/Download/index.aspx?lang=cn>
 2. 禁用 nouveau 第三方驱动
打开编辑配置文件:
`$ sudo gedit /etc/modprobe.d/blacklist.conf`
在最后一行添加: `blacklist nouveau`
改好后执行命令: `$ sudo update-initramfs -u`
重启使之生效: `$ reboot`
 3. 安装驱动
执行命令: `$ lsmod | grep nouveau`
禁用 X 服务: `$ sudo /etc/init.d/lightdm stop` (或者: `sudo service lightdm stop`)
给驱动 run 文件赋予可执行权限: `$ sudo chmod a+x NVIDIA-Linux-x86_64-440.64.run` (下载的驱动文件名)
安装: `$ sudo ./NVIDIA-Linux-x86_64-440.64.run -no-opengl-files`
开启 X 服务: `sudo /etc/init.d/lightdm start` (或者: `sudo service lightdm start`)
`-no-opengl-files` 只安装驱动文件, 不安装 OpenGL 文件。这个参数最重要
`-no-x-check` 安装驱动时不检查 X 服务, `-no-nouveau-check` 安装驱动时不检查 nouveau, 后面两个参数可不加。
 4. 检查
重启, 没有问题, 输入命令: `$ nvidia-smi`
如果出现了驱动版本就表示安装成功了。

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		N/A			
0%	35C	P8	10W / 250W	0MiB / 11178MiB	0%	Default			
1	GeForce GTX 108...	Off	00000000:03:00.0	Off		N/A			
0%	50C	P2	60W / 250W	829MiB / 11178MiB	0%	Default			
2	GeForce GTX 108...	Off	00000000:83:00.0	Off		N/A			
0%	39C	P8	9W / 250W	10MiB / 11177MiB	0%	Default			

3.4 CUDA

- 版本要求: 10.0
- 具体步骤

1. 下载:

登录官网 <https://developer.nvidia.com/cuda-10.0-download-archive?>

选择 linux—x86_64—Ubuntu—16.04—runfile(local)

下载 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. 安装:

```
$ sudo sh cuda_10.0.130_410.48_linux.run
```

3. 查看版本:

```
$ cat /usr/local/cuda/version.txt
```

4. cuda 设置(适用于多个不同版本时):

```
$ export PATH=/usr/local/cuda/bin:$PATH
```

```
$ export LD_LIBRARY_PATH=  
/usr/local/cuda/lib64:$LD_LIBRARY_PATH
```

3.5 CUDNN

- 版本要求: 7.6.5
- 具体步骤

1. 下载:

登陆官网: <https://developer.nvidia.com/rdp/cudnn-download>

选择 cuDNN Library for Linux: cudnn-10.0-linux-x64-v7.6.5.32.tgz

2. 解压:

```
$ tar -zxvf cudnn-10.0-linux-x64-v7.6.5.32.tgz
```

得到:

```
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. 拷贝:

```
$ sudo cp cuda/lib64/* /usr/local/cuda-10.0/lib64/
$ sudo cp cuda/include/* /usr/local/cuda-10.0/include/
```

4. 版本查看:

```
$ cat /usr/local/cuda/include/cudnn.h | grep CUDNN_MAJOR -A 2
```

得到:

```
#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

- 版本要求: 3.7 (注: 系统自带的 python2.7 不要卸载或删除, 以免影响系统运行)

- 具体步骤:

1. 下载:

```
$ wget https://www.python.org/ftp/python/3.7.3/Python-3.7.3.tar.xz
```

2. 编译安装:

```
$ 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. 创建软链接:

```
$ ln -s /usr/local/bin/python3/bin/python3 /usr/bin/python3
$ ln -s /usr/local/bin/python3/bin/pip3 /usr/bin/pip3
```

4. 版本查看:

```
$ python3 --version
```

3.7 Pytorch

- 版本要求：1.3.0
- 具体步骤：
 1. pip 安装：pip3 install torch==1.3.0
注：这里的 pip3 是上面的 python3.7 下的，可通过 pip3 -V 查看
 2. 版本查看：

```
$ python3
>> import torch
>> print(torch.__version__)
```

3.8 Torchvision

- 版本要求：0.4.2
- 具体步骤：
 1. pip 安装：pip3 install torchvision==0.4.2
注：这里的 pip3 是上面的 python3.7 下的，可通过 pip3 -V 查看
 2. 版本查看：

```
$ python3
>> import torchvision
>> print(torchvision.__version__)
```

3.9 Magik 训练插件

- 版本要求：1.1.1
- 具体步骤：
 1. 环境确认：

```
$ python
>> import torch
>> print(torch.__version__) #torch 的版本
>> print(torch.version.cuda) #对应的 cuda 版本
>> print(torch.backends.cudnn.version()) #对应的 cudnn 的版本
```

也可以根据自身需要设置 cuda 和 cudnn 的版本
 2. 安装前根据上述环境在插件目录
在我们提供的插件目录下 magik-toolkit/TrainingKit/pytorch/magik_whl 下找到对应的安装包。
pip 安装：pip3 install magik_trainingkit_torch_130-1.1.1-py3-none-any.whl
注：这里的 whl 安装包由我们根据客户的环境提供，如果插件目录下没有对应的版本，及时同步给相关人员以适配环境编译对应的插件。另外，

若只是更新插件，需要先卸之前的旧插件。

2. 查看是否安装成功：

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

出现上图表明导入成功，绿色的字体中有当前使用插件的版本，commit号及编译的日期。

3.10 其他

运行 yolov5 需要一些其他的安装包：

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)

运行时如还有其他，根据提示用 pip install 进行安装，直至正常运行。

4. Yolov5s 训练流程

4.1 代码和模型的下载

yolov5 的代码地址：

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

这里因为加入了插件部分，对相关算子进行了重新的封装，我们会提供一套基于原生 yolov5 改后的训练代码 yolov5s-person，其中原生的算子也有做注释保留，使用者以对比修改前后的算子了解和熟悉具体修改和使用方法。若对原生的代码感兴趣，或者想先验证一下环境是否正常，也可先跑一下原生的 yolov5 的代码。

提供的代码在 magik-toolkit/Models/training/pytorch/yolov5s-person/下，测试和上板试验用的是 yolov5s-person-4bit.pt 模型地址在 magik-toolkit/Models/training/pytorch/yolov5s-person/runs/train 下。

4.2 数据准备

1. 以 COCO2017 为例，提取出里面的 person(id: 0)部分并转换成 Yolo 训练需要的数据格式，COCO2017 数据集下载好并解压，annotation 文件的格式为 json 格式。

下载地址(用 wget 进行下载):

图片:

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

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

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

标注:

http://images.cocodataset.org/annotations/stuff_annotations_trainval2017.zip

http://images.cocodataset.org/annotations/image_info_test2017.zip

http://images.cocodataset.org/annotations/annotations_trainval2017.zip

2. 生成数据的脚本在 COCO_forYOLO 文件夹下:

(1)运行 `python batch_split_annotation_foryolo.py` (注意修改程序中的 coco 路径'coco_data_dir=')。

(2)运行完会在 coco 路径下生成 person/data/images、 person/data/ImageSets、 person/data/labels 三个文件夹, 将 ImageSets 下的 train2017.txt、 val2017.txt、 test2017.txt 放入 persondet/data/coco 文件夹下, txt 文件里存储图片的绝对路径。

4.3 网络训练

1. 训练配置:

默认配置参数可参考 models/commom.py 中 is_quantize, bitw, bita 等参数:

```
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"
```

32bit, 设置: bita = 32

8bit, 设置: bita = 8

4bit, 设置: bita = 4

其中,

is_quantize - 是否进行量化, 0-不量化, 1-量化, bita 为 32 时, 注意该值设为 0

bitw - weight 的位宽

bita - feature 的位宽

clip_max_value - feature 的截断值，建议 8bit 设为 6.0，4bit 设为 4.0

shortcut_clip_max_value 代表 shortcut 的 feature 的截断值，建议 8bit 设为 2.0，4bit 设为 1.5

2. 训练脚本 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)关于预训练

浮点训练时没有预训练模型，--weights 为 ‘ ’，8bit 训练时加载已得到的精度足够的 32bit 模型，4bit 加载训练好的 8bit 作为预训练模型，这样一步步推进，效果更佳。

(2)关于多卡训练

训练时加入了 torch.distributed.launch，为多卡的分布式训练，GPUS 的值和下面 device 的总数对应，按实际情况对应修改，如果是单卡训练，直接用 python3 train.py 加上后面的参数即可。batch-size 是所有显卡总的 batch 数目，按实际显卡大小设置即可。

(3)关于学习率

超参数的设置在 data/hyp.scratch.yaml 里，lr0 设置初始学习率，其余采用默认值。train.py 里加载预训练模型有关于加载 optimizer 的部分，这里因为低 bit 模型相对于 32bit 模型不是接着训练而是类似重训，所以这里建议去掉加载 optimizer 的部分，以免影响低 bit 的训练效果。

```

# 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)

```

ema 和 epoch 也存在类似的问题，因此如图也对应做了修改。

(4)模型的保存

train.py 里--project 可设置保存模型路径, 在 runs/train/project 下, weights 下保存的有两个模型, best.pt 本次训练到目前最好的, last.pt 本次训练到目前最新的, 每个 epoch 测试的结果保存在 result.txt 中, 可随时查看。

(5)训练经验

针对 yolov5 系列的训练, 目前总结的经验是 32bit 用 sgd 和 lr0.01, 8bit 基于 32bit 的预训练用 sgd 和 lr0.01, 4bit 基于 8bit 的预训练用 adam 和 lr0.001。具体命令参见 train.sh。32bit 没有预训练, 收敛稍慢; 8bit 有预训练且表达能力可以, 收敛较快; 4bit 虽有预训练但表达能力稍弱, 收敛稍慢。

4.4 模型测试

1.待检测图片

通过 python3 detect.py -h 查看选取所需参数, 通过设置 source 可检测图片或视频或图片文件夹, 检测结果可设置显示(--view-img)或保存(--save-img)

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

示例如下, 若需要其他操作可选择相应的参数;

```

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

```

其中,

source - 待检测图片
 weights - 训练好的用于测试的模型
 imgs - 测试图片的大小
 device - 用第几块显卡测试
 view-img - 是否画框显示检测结果

注: 检测的模型配置一定要和训练时候的配置(bits)一致。

2. 测试模型精度

```
$ 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)
```

待测试的模型通过---weights 指定，验证集通过 data/coco-person.yaml 中的 val2017.txt 确定，其余参数根据实际需要给定。

具体测试结果如下：

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	

models 下的 yolov5m.yaml 和 yolov5l.yaml 也分别对应原始的 yolov5m 和 yolov5l，训练流程也都同 yolov5s，有需要可以直接使用。

5. 模型转换

5.1 pt 转 onnx

在上述训练的环境下生成 onnx 文件：

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

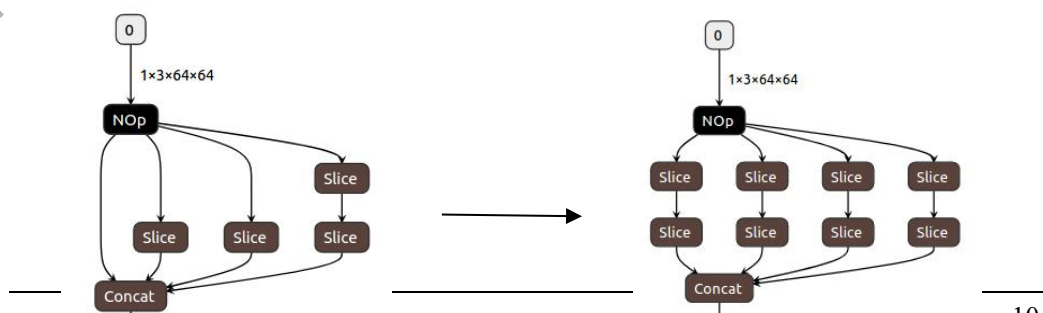
input: 指定待转模型的路径

output: .onnx，位置和.pt 在同一目录

转 onnx 时必须使用 opset9，否则会产生多余的节点导致后续转换工具不支持，如因使用 opset9 转 onnx 失败，可参见步骤 7 的疑问解答，不同版本的 torch 会有一些不同的需要注意的点。

```
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)
```

特别地，对于 yolov5 用到了 focus 算子的转 onnx 之前需要注意修改 torch/onnx/symbolic_opset9.py 中的 _slice() 函数，如下：



改前

改后

转换后的 onnx 文件可用 netron 进行可视化查看。

5.2 onnx 转 bin

这一步不依赖 python, torch 及插件等环境, 只需要我们的转换工具 magik_transformer-*py3-none-any.whl (在 magik_toolkit/TransFormKit 下) 及转换脚本 run_t40/t41/a1/x2500.sh:

magik_transformer-*py3-none-any.whl 拿到后直接使用 pip install magik_transformer-*py3-none-any.whl 安装即可, 若先前有安装过则先使用 pip uninstall magik_transformer-*py3-none-any.whl 卸载转换工具。

```
cd Models/training/pytorch/yolov5s-person/transform_sample
```

```
$ sh run_t40/t41/a1/x2500.sh
```

不同的部署芯片型号使用不同的 run_*.sh, 这里以 T40 芯片为例进行介绍, 具体的,

run_t40.sh

```
cd ../
pod -ri 's/(target_device = ")[^"]*/\1T40/' models/common.py models/yolo.py
sh detect.sh
sh convert_onnx.sh
cd -
CUDA_VISIBLE_DEVICES=0 python transform.py --model_file ../runs/train/yolov5s-person-4bit.onnx --output_file ./venus_sample_yolov5s/yolov5s_t40_magik.mk.h --config_file cfg/
magik_t40.cfg
cd venus_sample_yolov5s
cp makefile_files/Makefile_t40 Makefile
```

首先将网络中的 target_device 替换为对应的芯片型号, 在使用上一步的 detecte.py 生成 4bit 的 onnx 模型文件, 再执行最终的转换命令 transform.py, 具体的转换模型参数如下:

- model_file 转换的输入, 对应上一步得到的 onnx 模型文件路径。
- output_file 转换的输出, 转换完成后会在 venus_sample_yolov5s 下生成所需的.bin 文件。
- config_file 转换时的配置文件

不同的芯片对应的 cfg 文件不同, 主要是 SOC 设备不同, 使用时选择不同的 cfg 文件即可, 这里同样以 SOC 为 T40 为例。

magik_t40.cfg


```

"ARCH": {
  "SOC": "T40"
},
"MODEL": {
  "FORMAT": "onnx",
  "INPUT": [
    {
      "BATCH": 1,
      "WIDTH": 640,
      "HEIGHT": 640,
      "CHANNEL": 3,
      "COLOR": "RGB",
      "NORMAL": [255,255,255],
      "MEAN": [0,0,0]
    }
  ]
},
"EXTRA": {
  "BIZ_CODE": "Ingenic",
  "VERBOSE": "ON"
}
}

```

SOC	--	部署芯片型号, 不同的 SOC 芯片这里设置会不同
FORMAT	--	输入的模型格式
INPUT	--	指定模型的输入信息
BATCH	--	模型的输入个数
WIDTH	--	模型输入的宽
HEIGHT	--	模型输入的高
CHANNEL	--	模型输入的通道数
COLOR	--	模型输入源是图片, 指定图片处理方式 ("BGR" / "RGB" / "GRAY")
NORMAL	--	指定模型输入的方差信息
MEAN	--	指定模型输入的均值信息

```

2022-07-21 12:22:38.484535: W clip_model_by_nop_optimizer.cc:48-collect_remove_nodes] Multiple output NOp exists on the current pathway!
2022-07-21 12:23:46.524015: W quantize_device.cc:55-get asynchronous quantize_type] weight_quantize devices is empty
2022-07-21 12:23:46.810610: W quantize_operation.cc:521-search] Not found QuantizeOperation function for Unpool2D
2022-07-21 12:23:46.826601: W quantize_operation.cc:521-search] Not found QuantizeOperation function for Unpool2D
*****
^ ^ Convert successfully, Enjoy it ^ ^
*****

```

转换成功后会出现上图中的标志, 并在 `venus_sample_yolov5s` 下生成 `yolov5s_t40_magik.bin` 文件, 这就是我们最终上板要用到的二进制模型文件。

```

10 w1024_h714.nv12 magik_model yolov5s_t40_magik.mk.h stb
bus.jpg Makefile yolov5s_t40_magik.bin
inference.cpp makefile_files
inference_nv12.cpp readme.md

```

6. 模型上板

6.1 代码编译

1. 准备工作

transform_sample/venus_sample_yolov5s 目录下我们提供了上板运行需要的 inference.cpp、Makefile、测试数据，另外我们还需要用到 venus 库以及 mips 编译工具。其中，
venus 库在 magik-toolkit/InferenceKit/下；
Mips 编辑工具由方案同事提供。

2. 网络输入

这里为了用户快速实现流程，我们在在 venus_sample_yolov5s 下加入了 stb 一些图片读取函数，所以测试的时候直接传入 jpg 图即可。

3. 模型的加载

提供的实例 inference.cpp 中模型是通过参数传入的，运行前注意同步拷贝到板端对应的目录并在运行时传入。

4. 超参数的设置

```
void generateBBox(std::vector<venus::Tensor> out_res, std::vector<magik::venus::ObjBbox_t>& candidate_boxes, int img_w, int img_h)
{
    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};
```

strides 和 anchor 按 yolov5 的实际使用设置，person_threshold 和 nms_threshold 分别对应原始代码中的 conf-thres（置信度阈值）和 iou-thres（iou 阈值），classes 是类别数。

5. 编译

TOPDIR - venus 库的相对目录

libtype - 根据实际板子的需求确定是否 muclibc（这里以 muclibc 为例）

build_type - release 模式，运行得到结果，默认设置

- profile 模式，运行时网络结构可视化及网络每层运行时间及 GOPs 统计

- debug 模式，运行得到结果的同时保存每层量化 feature 的结果

- nmem 模式，统计模型运行时 nmem 内存占用情况，运行程序，内存使用情况保存在/tmp/nmem_memory.txt

直接 make 编译 inference.cpp 即可生成 venus_yolov5s_bin_uclibc_*, 即我们上板需要的可执行文件。

注意：我们同时提供了用于实际板端运行的输入为 nv12 的代码用例

inference_nv12.cpp 及输入数据 10_w1024_h714.nv12, 如有需要, 可修改 Makefile 进行编译及使用测试, 在 makefile_files 下面有不同部署芯片使用的 Makefile_a1/t40/t41/x2500, 不同的芯片使用的 vnues 库有所不同, 所以在使用的時候选择不同的 Makefile 文件进行编译即可, 这里还是以 t40 为例。

不同部署芯片对应的 vnues 库:

芯片型号	对应的 venus 库
T40	InferenceKit/nna1/mips720-glibc229/
T41	InferenceKit/nna2/mips720-glibc229/T41/
A1	InferenceKit/nna2/mips720-glibc229/A1/
X2500	InferenceKit/nna1/mips720-glibc229/

6.2 代码编译

注:

venus 库在 magik-toolkit/InferenceKit/nn1/mips720-glibc229/下

1.1. release(发布库)

编译: make build_type=release

在当前文件夹下生成 venus_yolov5s_bin_uclibc_release 可执行文件, 拷贝 venus 库 (libvenus.so)、可执行文件 (venus_yolov5s_bin_uclibc_release)、模型文件 (yolov5s_t40_magik.bin)、测试图片 (bus.jpg) 至开发板运行即可: ./venus_yolov5s_bin_uclibc_release yolov5s_t40_magik.bin bus.jpg

(注: 运行前添加库路径至 LD_LIBRARY_PATH:

export LD_LIBRARY_PATH=\$lib_path:\$LD_LIBRARY_PATH).

清除 make build_type=release clean

```
[root@ingenic-uc1_1:v5s-fx]# ./venus_yolov5s_bin_uclibc_release yolov5s_t40_magik.bin bus.jpg
The soc-nna version is 20220525
INFO(magik): venus memory map size: 0
INFO(magik): venus version:0.9.6.1.ALPHA(00000906_184a23e) built:20220715-1559(7.2.0 r5.1.3 glibc2.29 mips@NNA1)
INFO(magik): model version:0.9.6.NNA1_c2436c4
[I/magik:venus]: kv_size = 0
ori_image w,h: 810,1080
model-->640,640 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
test_net run time_ms:38.919000ms
pad_x:176 pad_y:0 scale:0.355556
post_net time_ms:1.046000ms
box: 51 408 239 904 0.90
box: 217 401 351 869 0.83
box: 669 409 811 891 0.67
[root@ingenic-uc1_1:v5s-fx]#
```

2. debug (用于核对数据的库)

详见步骤 7.2, 输入的处理有些不同。

3. profile (网络可视化及每层运行时间统计)

编译: make build_type=profile

在当前文件夹下生成 venus_yolov5s_bin_uclibc_profile 可执行文件, 拷

贝 venus 库 (libvenus.p.so) 、 可 执 行 文 件 (venus_yolov5s_bin_uclibc_profile) 、 模 型 文 件 (yolov5s_t40_magik.bin)、 测试 图片 (bus.jpg) 至 开发板 运行 即可 : ./venus_yolov5s_bin_uclibc_profile yolov5s_t40_magik.bin bus.jpg)

注: 运行前添加库路径至 LD_LIBRARY_PATH:

export LD_LIBRARY_PATH=\$lib_path:\$LD_LIBRARY_PATH

清除 make build_type=profile clean

```
989 Pooling (1,12,20,256) (1,12,20,256) (5,5) (1,1) (2,2,2,2) N/A 0.29 0.29 0.29 0.29 0.77 0.001 0.059 0.00
990 Pooling (1,12,20,256) (1,12,20,256) (5,5) (1,1) (2,2,2,2) N/A 0.30 0.34 0.29 0.29 0.80 0.001 0.059 0.00
991 Concat (1,12,20,1024) (1,12,20,1024) N/A N/A N/A 0.36 0.36 0.35 0.35 0.95 0.000 0.234 0.00
layer_40_Quantiz.. Convolution (1,12,20,512) (1,1,1024,512) (1,1) (0,0,0,0) (1,1) 0.35 0.35 0.35 0.35 0.94 0.222 0.430 269.00
layer_58_Quantiz.. Convolution (1,12,20,512) (1,12,20,256) (1,1,512,256) (1,1) (0,0,0,0) (1,1) 0.09 0.09 0.09 0.09 0.25 0.063 0.152 66.00
992 Upsample (1,12,20,256) (1,24,40,256) N/A N/A N/A 0.08 0.08 0.08 0.08 0.22 0.000 0.146 0.00
993 Concat (1,24,40,512) (1,24,40,512) N/A N/A N/A 0.64 0.64 0.63 0.63 1.70 0.000 0.469 0.00
layer_55_Quantiz.. Convolution (1,24,40,512) (1,24,40,128) (1,1,512,128) (1,1) (0,0,0,0) (1,1) 0.18 0.24 0.17 0.24 0.49 0.126 0.325 33.00
layer_53_Quantiz.. Convolution (1,24,40,128) (1,24,40,128) (1,1,128,128) (1,1) (0,0,0,0) (1,1) 0.11 0.13 0.11 0.13 0.30 0.031 0.126 9.00
layer_51_Quantiz.. Convolution (1,24,40,128) (1,24,40,128) (3,3,128,128) (1,1) (1,1,1,1) (1,1) 0.22 0.22 0.22 0.22 0.60 0.213 0.188 73.00
layer_49_Quantiz.. Convolution (1,24,40,128) (1,24,40,128) (1,1,512,128) (1,1) (0,0,0,0) (1,1) 0.15 0.15 0.15 0.15 0.40 0.126 0.325 33.00
994 Concat (1,24,40,256) (1,24,40,256) N/A N/A N/A 0.34 0.35 0.34 0.35 0.92 0.000 0.234 0.00
layer_46_Quantiz.. Convolution (1,24,40,256) (1,24,40,128) (1,1,256,128) (1,1) (0,0,0,0) (1,1) 0.19 0.19 0.19 0.19 0.51 0.126 0.268 24.00
995 Upsample (1,24,40,128) (1,48,80,128) N/A N/A N/A 0.12 0.12 0.12 0.12 0.32 0.000 0.293 0.00
996 Concat (1,48,80,256) (1,48,80,256) N/A N/A N/A 1.23 1.27 1.21 1.22 3.28 0.000 0.938 0.00
layer_41_Quantiz.. Convolution (1,48,80,256) (1,48,80,64) (1,1,256,64) (1,1) (0,0,0,0) (1,1) 0.25 0.26 0.25 0.25 0.68 0.126 0.594 8.50
layer_39_Quantiz.. Convolution (1,48,80,64) (1,48,80,64) (1,1,64,64) (1,1) (0,0,0,0) (1,1) 0.19 0.19 0.19 0.19 0.51 0.031 0.237 2.50
layer_37_Quantiz.. Convolution (1,48,80,64) (1,48,80,64) (3,3,64,64) (1,1) (1,1,1,1) (1,1) 0.23 0.23 0.23 0.23 0.63 0.283 0.252 18.50
layer_35_Quantiz.. Convolution (1,48,80,64) (1,48,80,64) (1,1,256,64) (1,1) (0,0,0,0) (1,1) 0.22 0.23 0.22 0.22 0.60 0.126 0.594 8.50
1022 Concat (1,48,80,128) (1,48,80,128) N/A N/A N/A 0.65 0.69 0.64 0.64 1.74 0.000 0.469 0.00
layer_32_Quantiz.. Convolution (1,48,80,128) (1,48,80,128) (1,1,128,128) (1,1) (0,0,0,0) (1,1) 0.34 0.34 0.34 0.34 0.91 0.126 0.478 9.00
1131 Convolution (1,48,80,32) (1,48,80,32) (1,1,128,32) (1,1) (0,0,0,0) (1,1) 0.33 0.33 0.33 0.33 0.89 0.031 0.707 4.25
Format_convert.. FormatConvert (1,48,80,32) (1,48,80,18) N/A N/A N/A 0.25 0.25 0.25 0.25 0.68 0.000 0.732 0.00
layer_30_Quantiz.. Convolution (1,48,80,128) (1,24,40,128) (3,3,128,128) (2,2) (1,1,1,1) (1,1) 0.42 0.43 0.42 0.43 1.14 0.283 0.364 73.00
1040 Concat (1,24,40,256) (1,24,40,256) N/A N/A N/A 0.34 0.35 0.34 0.34 0.92 0.000 0.234 0.00
layer_27_Quantiz.. Convolution (1,24,40,256) (1,24,40,128) (1,1,256,128) (1,1) (0,0,0,0) (1,1) 0.14 0.14 0.14 0.14 0.37 0.063 0.192 17.00
layer_25_Quantiz.. Convolution (1,24,40,128) (1,24,40,128) (1,1,128,128) (1,1) (0,0,0,0) (1,1) 0.12 0.10 0.10 0.10 0.31 0.031 0.126 9.00
layer_23_Quantiz.. Convolution (1,24,40,128) (1,24,40,128) (3,3,128,128) (1,1) (1,1,1,1) (1,1) 0.22 0.23 0.22 0.22 0.60 0.283 0.188 73.00
layer_21_Quantiz.. Convolution (1,24,40,256) (1,24,40,128) (1,1,256,128) (1,1) (0,0,0,0) (1,1) 0.12 0.12 0.12 0.12 0.31 0.063 0.192 17.00
1072 Concat (1,24,40,256) (1,24,40,256) N/A N/A N/A 0.34 0.34 0.34 0.34 0.92 0.000 0.234 0.00
layer_18_Quantiz.. Convolution (1,24,40,256) (1,24,40,256) (1,1,256,256) (1,1) (0,0,0,0) (1,1) 0.19 0.19 0.19 0.19 0.51 0.126 0.268 34.00
1124 Concat (1,24,40,256) (1,24,40,32) (1,1,256,32) (1,1) (0,0,0,0) (1,1) 0.13 0.13 0.13 0.13 0.36 0.016 0.242 6.25
Format_convert.. FormatConvert (1,24,40,32) (1,24,40,18) N/A N/A N/A 0.09 0.09 0.09 0.09 0.24 0.000 0.183 0.00
layer_16_Quantiz.. Convolution (1,24,40,256) (1,12,20,256) (3,3,256,256) (2,2) (1,1,1,1) (1,1) 0.52 0.52 0.51 0.52 1.38 0.283 0.430 290.00
1090 Concat (1,12,20,512) (1,12,20,512) N/A N/A N/A 0.20 0.24 0.19 0.19 0.54 0.000 0.117 0.00
layer_13_Quantiz.. Convolution (1,12,20,512) (1,12,20,256) (1,1,512,256) (1,1) (0,0,0,0) (1,1) 0.21 0.21 0.20 0.20 0.55 0.063 0.244 130.00
layer_11_Quantiz.. Convolution (1,12,20,256) (1,12,20,256) (1,1,256,256) (1,1) (0,0,0,0) (1,1) 0.13 0.13 0.13 0.13 0.36 0.031 0.182 66.00
layer_9_Quantiz.. Convolution (1,12,20,256) (1,12,20,256) (3,3,256,256) (1,1) (1,1,1,1) (1,1) 0.68 0.68 0.68 0.68 1.81 0.283 0.602 378.00
layer_7_Quantiz.. Convolution (1,12,20,512) (1,12,20,256) (1,1,512,256) (1,1) (0,0,0,0) (1,1) 0.20 0.20 0.20 0.20 0.54 0.063 0.244 130.00
1122 Concat (1,12,20,512) (1,12,20,512) N/A N/A N/A 0.24 0.24 0.23 0.24 0.63 0.000 0.234 0.00
layer_4_Quantiz.. Convolution (1,12,20,512) (1,12,20,512) (1,1,512,512) (1,1) (0,0,0,0) (1,1) 0.45 0.45 0.45 0.45 1.20 0.126 0.488 260.00
1137 Convolution (1,12,20,32) (1,12,20,32) (1,1,512,32) (1,1) (0,0,0,0) (1,1) 0.14 0.18 0.13 0.18 0.38 0.000 0.162 16.25
Format_convert.. FormatConvert (1,12,20,32) (1,12,20,18) N/A N/A N/A 0.04 0.04 0.03 0.04 0.09 0.000 0.046 0.00
[/magik:venus]:
Timing cycle = 10
===== Concise Dispatch Profiler Summary: N/A, Exclude 5 warn-ups =====
OperatorType Avg(ms) Max(ms) Min(ms) Avg(%) GOPs TotalMB(MB) CalledTimes
Add 1.924 1.936 1.918 5.15 0.001 2.021 7
Concat 6.720 6.876 6.676 18.02 0.000 1.922 13
Convolution 25.350 25.989 25.113 67.87 9.618 25.627 60
FormatConvert 0.377 0.380 0.373 1.01 0.000 0.961 3
Normalize 1.808 1.931 1.677 5.65 0.000 1.875 1
Pooling 0.884 0.926 0.872 2.37 0.004 0.176 3
Upsample 0.199 0.202 0.195 0.53 0.000 0.439 2
-----
Total: 37.351 ms
[root@ingenic-ucl1v5s-fx]#
```

4. nmem 模式 (用于统计网络运行时 nmem 占用情况, 保存在 /tmp/nmem_memory.txt)

编译: make build_type=nmem

在当前文件夹下生成 venus_yolov5s_bin_uclibc_nmem 可执行文件, 拷贝贝 venus 库 (libvenus.m.so) 、 可 执 行 文 件 (venus_yolov5s_bin_uclibc_nmem)、 模型文件(yolov5s_t40_magik.bin)、 测试 图片 (bus.jpg) 至 开发板 运行 即可: ./venus_yolov5s_bin_uclibc_nmem yolov5s_t40_magik.bin bus.jpg

注: 运行前添加库路径至 LD_LIBRARY_PATH:

export LD_LIBRARY_PATH=\$lib_path:\$LD_LIBRARY_PATH

清除 make build_type=nmem clean

```
[root@ingenic-uc1_1:v5s-fx]# ./venus_yolov5s_bin_uclibc_nmem yolov5s_t40_magik.bin bus.jpg
The soc-nna version is 20220525
INFO(magik): venus memory map size: 0
INFO(magik): venus version:0.9.6.1.ALPHA(00000906_184a23e) built:20220715-1546(7.2.0 r5.1.3 glibc2.29 mips@NNA1)
INFO(magik): model version:0.9.6.NNA1_c2436c4
[t/magik::venus]: kv_size = 0
ori_image w,h: 810 ,1080
model-->640 ,640 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
test_net run time_ms:375.395000ms
pad_x:176 pad_y:0 scale:0.355556
post_net time_ms:1.414000ms
box: 51 408 239 904 0.90
box: 217 401 351 869 0.83
box: 669 409 811 891 0.67
[root@ingenic-uc1_1:v5s-fx]#
```

7. 数据核对

验证 PC 端和板端的数据是否能对齐，可按如下操作：

1. PC 端

步骤 4.4-1 测试单张图片加入环境变量 `MAGIK_TRAININGKIT_DUMP=1` 可保存每层量化结果到 `/tmp/trainingkit_data/feature` 下，也可通过设置环境变量 `MAGIK_TRAININGKIT_PATH` 指定保存的目录，具体运行命令为：

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

这里的 `bus.jpg` 原始分辨率是 `1080*810`，测试时加入了目标 `imgs` 为 `640`，按 `yolov5` 代码的缩放原则（长边 `640`，短边等比缩放再填充至 `32` 的倍数），最后进入测试的分辨率为 `640*480`，最终保存 `./trainingkit_data/feature` 下（具体见下图），可见输入层有保存为 `input_data_shape_1_640_480_3.bin`，后面的 `shape` 命名的规则是 `n,h,w,c`，即高为 `640`，宽为 `480`，最后面是三个输出层；

参数 `imgs`、`conf-thres`、`iou-thres` 的设置也是为了和 `c` 代码(详见 6.1-3)保持一致，以保证比对条件一致；

```

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. 板端

(1) 编译 c 代码：编译之前最好先将之前编译过的模式做一下 clean(make build_type=release/profile/nmem clean)，之后 make build_type=debug，在当前文件夹下生成 venus_yolov5s_bin_uclibc_debug 可执行文件。拷贝 venus 库 (libvenus.d.so)、可执行文件 (venus_yolov5s_bin_uclibc_debug)、模型文件 (yolov5s_t40_magik.bin)、模型输入 magik_input_nhwc_1_640_480_3.bin 至开发板运行：

```
./venus_yolov5s_bin_uclibc_debug yolov5s_t40_magik.bin
```

```
magik_input_nhwc_1_640_480_3.bin
```

其他信息（如下图）：


```

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_bt.bin layer_35_QuantizeFeature_weight.bin layer_65_QuantizeFeature_bt.bin
layer_127_QuantizeFeature_out.bin layer_37_QuantizeFeature_bt.bin layer_65_QuantizeFeature_out.bin
layer_127_QuantizeFeature_weight.bin layer_37_QuantizeFeature_out.bin layer_65_QuantizeFeature_weight.bin
layer_129_QuantizeFeature_bt.bin layer_37_QuantizeFeature_weight.bin layer_67_QuantizeFeature_bt.bin
layer_129_QuantizeFeature_out.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_out.bin layer_42_QuantizeFeature_weight.bin layer_6_QuantizeFeature_out.bin
layer_14_QuantizeFeature_weight.bin layer_44_QuantizeFeature_bt.bin layer_6_QuantizeFeature_weight.bin
layer_15_QuantizeFeature_out.bin layer_44_QuantizeFeature_out.bin layer_71_QuantizeFeature_bt.bin
layer_17_QuantizeFeature_bt.bin layer_44_QuantizeFeature_weight.bin layer_71_QuantizeFeature_out.bin
layer_17_QuantizeFeature_out.bin layer_46_QuantizeFeature_bt.bin layer_71_QuantizeFeature_weight.bin
layer_17_QuantizeFeature_weight.bin layer_46_QuantizeFeature_out.bin layer_73_QuantizeFeature_bt.bin
layer_19_QuantizeFeature_bt.bin layer_46_QuantizeFeature_weight.bin layer_73_QuantizeFeature_out.bin
layer_19_QuantizeFeature_out.bin layer_48_QuantizeFeature_bt.bin layer_73_QuantizeFeature_weight.bin
layer_19_QuantizeFeature_weight.bin layer_48_QuantizeFeature_out.bin layer_75_QuantizeFeature_bt.bin
layer_21_QuantizeFeature_bt.bin layer_48_QuantizeFeature_weight.bin layer_75_QuantizeFeature_out.bin
layer_21_QuantizeFeature_out.bin layer_4_QuantizeFeature_bt.bin layer_75_QuantizeFeature_weight.bin
layer_21_QuantizeFeature_weight.bin layer_4_QuantizeFeature_out.bin layer_77_QuantizeFeature_bt.bin
layer_23_QuantizeFeature_bt.bin layer_4_QuantizeFeature_weight.bin layer_77_QuantizeFeature_out.bin
layer_23_QuantizeFeature_out.bin layer_50_QuantizeFeature_bt.bin layer_77_QuantizeFeature_weight.bin
layer_23_QuantizeFeature_weight.bin layer_50_QuantizeFeature_out.bin layer_80_QuantizeFeature_bt.bin
layer_25_QuantizeFeature_bt.bin layer_50_QuantizeFeature_weight.bin layer_80_QuantizeFeature_out.bin
layer_25_QuantizeFeature_out.bin layer_51_QuantizeFeature_out.bin layer_80_QuantizeFeature_weight.bin
layer_25_QuantizeFeature_weight.bin layer_53_QuantizeFeature_bt.bin layer_82_QuantizeFeature_bt.bin
layer_27_QuantizeFeature_bt.bin layer_53_QuantizeFeature_out.bin layer_82_QuantizeFeature_out.bin

```

其中 layer_*_QuantizeFeature_out.bin 和 PC 端运行时保存在 ./trainingkit_data/feature 下的 layer_*_QuantizeFeature.bin 是一一对应的，直接核对 md5 值是否一致即可，在保证输入完全一致的情况下中间有层不对应及时反馈。

```

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

```

debug 模式这里只核对网络结果，没有加后处理部分，因此无结果打印。

8. 常见疑问解答

1. 问：为什么 release 运行 810*1080 的 bus.jpg 的分辨率是 640*384，debug 的时候是 480*640？

答：当我们设置目标分辨率为 640 时，板端的缩放规则是按等比缩放之后 w 填充至 640（因为实际的视频流大都是 w>h），这里 640 是 w，384 是 h；而 yolov5 代码里的规则是按长边 640 等比缩放再将短边填充至 32 的倍数，所以这里 h 是 640,480 是 w，debug 是和 pc 端保持一致的，实际使用按需设置即可。

2. 问：导入 magik 包出现 “undefined symbol:

_ZN6caffe28TypeMeta21_typeMetadataInstanceIN3c108BFloat16EEEEPKNS_6

detail12TypeMetaDataEv” 类似错误？

答：当前环境下 torch 的版本和 magik 包的不对应。

3. 问：“importError: version ‘libcudart.so.10.1’ not found”？

答：当前 cuda 版本和插件编译时的版本不对应。

4. 问：加插件之后的网络是否可以加载原生 torch 模型做预训练？

答：理论可以，层对应好就行，不过基于原生 float 的 32bit 模型还是要再微调训练一下，不建议直接加载原生的跑 8bit。

5. 问：训练出现 “RuntimeError: Address already in use”？

答：在 python -m torch.distributed.launch 后指定一个未被使用的端口——master_port=60053

6. 问：训练出现 “expected scalar type Float but found Half”？

答：torch 中用了 amp.autocast 进行半精度训练加速，但量化暂不支持半精度训练，注掉这块即可。另，原生 v5 在保存模型时按 half() 存的，目前发现量化这样操作最后存的模型的精度会有损失，建议去掉。

7. 问：torch 转 onnx 出现 “step!=1 is currently not supported”？

答：yolov5 slice 转换遇到的问题，1.6 及之前版本可以可以，1.7 以后变成 error 了，torch/onnx/symbolic_opset9.py 对应报错的行 if + raise RuntimeError 注掉即可。

8. 问：torch 转 onnx 出错 “RuntimeError:input_shape_value==
reshape_value || input_shape_value==1 || reshape_value == 1 INTERNAL
ASSERT FAILED ……”？

答：1.9 版本之后的错，torch/onnx/utils.py 中 635 行 _export() 最后一个参数的 onnx_shape_inference=True 改为 onnx_shape_inference=False；
torch1.10 是在 677 行。不同版本找这个函数相应修改即可。

9. 问：torch1.8/torch1.9 版本转换 onnx 会出现 “Warning: Unsupported operator NOp. No schema registered for this operator.” 等？

答：warning 不用管，不影响结果。

10. 问：转 bin 出现 “After checking the input dimension, the input of concat cannot be concatenated!”？

答：slice 的问题（参见 5.1），将 torch/onnx/symbolic_opset9.py 中 _slice 函数中的 if+return 这两句去掉，可视化 onnx 的 slice 正常。

11. 问：c 代码编译出现 “error: ‘constexpr bool std::isinf(double)’ conflicts with a previous declaration” 等一系列冲突错误？

答：mips 的编译器和 x86 自带的出现冲突，将环境变量

CPLUS_INCLUDE_PATH 中关于/usr/include/x86_64-linux-gnu 的部分去掉。

12. 问：c 端编译出现 “make: mips-linux-gnu-g++: Command not found” ？

答：未指定 mips 编译器，通过 export 设置。

13. 问：板端精度和 PC 端精度损失有多少？

答：理论上损失很小，1 个点以内，若出现较大差距，可检查输入是否一致，训练和验证的数据是 BGR/RGB，阈值的设置，也可对比网络输出结果看看输出是否一致，具体见步骤 7。

14. 问：转 onnx 的分辨率，转换工具的分辨率以及板端测试的分辨率需要保持一致吗？

答：最好保持一致，尤其后两者，本例中没有保持一致是因为对这个没有影响，对于其他非全卷积网络可能导致转换出来的模型有影响，所以测试验证集精度时建议输入数据统一线下处理到同一大小再做 PC 端和板端的比较。