# Yolov5 pytorch mlu移植教程

官方Yolov5文档: https://github.com/ultralytics/yolov5/releases/tag/v4.0

官方模型下载地址: https://github.com/ultralytics/yolov5/releases/download/v4.0/yolov5s.pt

移植好后代码: http://gitlab.software.cambricon.com/neuware/software/ae/solution/270\_sample/yolov5\_demo

MLU SDK版本: 1.7.0

MLU PYTORCH版本: 1.3.0

【以下版本在v4版本上完成(v5和v4类似)】

### 高版本PyTorch降级到低版本PyTorch

为什么要做这一步?

由于官方yolov5模型是在PyTorch1.8.x版本上进行,而MLU的PyTorch还是1.3.0版本。

高版本的PyTorch具有zip压缩模型功能,而在1.3.0上不支持,如果在1.3.0的版本上直接打开高版本的pt,会出现报错。

# Traceback (most recent call last): File "detect.py", line 282, in <module> detect() File "detect.py", line 84, in detect model = attempt\_load(weights, map\_location=device) # load FP32 model File "/home/Cambricon-Test/zhangkun\_yolov5/yolov5\_/models/experimental.py", line 141, in attempt\_load state\_dict = torch.load(weights[0], map\_location='cpu') File "/torch/venv3/pytorch/lib/python3.6/site-packages/torch/serialization.py", line 507, in load return \_load(f, map\_location, pickle\_module, \*\*pickle\_load\_args) File "/torch/venv3/pytorch/lib/python3.6/site-packages/torch/serialization.py", line 680, in \_load raise RuntimeError("{} is a zip archive (did you mean to use torch.jit.load()?)".format(f.name)) RuntimeError: weights/yolov5s.pt is a zip archive (did you mean to use torch.jit.load()?)

### 搭建 PyTorch1.8.0环境

可以先用conda或者virtualenv创建一个python虚拟环境,本教程用的是python3.6.10,然后执行如下命令安装pytorch1.8.0 cpu版。

pip install torch==1.8.0+cpu torchvision==0.9.0+cpu torchaudio==0.8.0 -f https://download.pytorch.org/whl
/torch\_stable.html

### 降低pt模型的版本

下载官方yolov5代码:

```
git clonehttps://github.com/ultralytics/yolov5.git
```

```
git checkout v4.0
```

修改detect.py

添加代码: torch.save(model.state\_dict(), "./weights/unzip.pt", \_use\_new\_zipfile\_serialization=False)

```
detect(save img=False)
        20
21
22
23
24
25
26
27
         # Directories
        save_dir = Path(increment_path(Path(opt.project) / opt.name, exist_ok=opt.exist_ok)) # increment run
(save_dir / 'labels' if save_txt else save_dir).mkdir(parents=True, exist_ok=True) # make dir
        # Intration
set_logging()
device = select_device(opt.device)
half = device.type != 'cpu' # half precision only supported on CUDA
29
30
31
32
33
34
        # Load model
        model = attempt load(weights, map location=device) # load FP32 model
        torch.save(model.state_dict(), "./weights/unzip.pt", _use_new_zipfile_serialization=False)
imgsz = check img_size(imgsz, s=model.stride.max())  # check img_size
        imgsz = check_img_size(imgsz, s=model.stride.max())
if half:
    model.half() # to FP16
35
                                                                                 check img_size
36
37
```

wgethttps://github.com/ultralytics/yolov5/releases/download/v4.0/yolov5s.pt

运行python detect.py --device cpu --weights yolov5s.pt 生成unzip.pt文件,这个pt可以放到MLU PyTorch使用。 到此, PT降版本工作完成。

### 增加算子

修改/workspace/neuware\_sdk\_ubuntu\_prebuild/venv/lib/python3.6/site-packages/torch/nn/modules/activation.py,添加SiLU,Hardswish激活函数

class Hardswish(Module): # export-friendly version of nn.Hardswish()

@staticmethod

def forward(x):

# return x \* F.hardsigmoid(x) # for torchscript and CoreML

return x \* F. hardtanh(x + 3, 0., 6.) / 6. # for torchscript, CoreML and ONNX

class SiLU(Module): # export-friendly version of nn.SiLU()

@staticmethod

def forward(x):

return x \* torch.sigmoid(x)

修改/workspace/neuware\_sdk\_ubuntu\_prebuild/venv/lib/python3.6/site-packages/torch/nn/modules/\_\_init\_\_.py,注册SiLU,Hardswish激活函数

1, from .activation import 中添加SiLU, Hardswish,

2, \_\_all\_\_ = 中添加'SiLU', 'Hardswish',

### MLU移植工作

首先请确认您安装的是MLU SDK 1.7.0版本,具体的移植步骤如下:

1. 环境搭建(云平台已搭好)

使用镜像: yellow.hub.cambricon.com/pytorch/pytorch:0.15.0-ubuntu16.04

#/bin/bash

export MY\_CONTAINER="Cambricon-MLU270-v1.7.0-pytorch"

num=`docker ps -a|grep "\$MY\_CONTAINER"|wc -l`

echo \$num

echo \$MY\_CONTAINER

if [ O -eq \$num ];then

xhost +

docker run -e DISPLAY=unix\$DISPLAY --device /dev/cambricon\_dev0 --net=host --pid=host -v /sys/kernel/debug:/sys/kernel/debug -v /tmp/.X11 -unix:/tmp/.X11-unix -it --privileged --name \$MY\_CONTAINER -v \$PWD/Cambricon-MLU270/:/home/Cambricon-MLU270 \

-v \$PWD/../tar\_package/datasets:/home/Cambricon-MLU270/datasets \

-v \$PWD/../tar\_package/models:/home/Cambricon-MLU270/models \

yellow.hub.cambricon.com/pytorch/pytorch:0.15.0-ubuntu16.04 /bin/bash

else

docker start \$MY\_CONTAINER

#sudo docker attach \$MY\_CONTAINER

docker exec -ti \$MY\_CONTAINER /bin/bash

fi

-v部分自定义

### 2. 激活MLU PyTorch

source /torch/venv3/pytorch/bin/activate

### 3. 修改部分代码

在yolov5的目录下,需要修改models/experimental.py部分代码

### 1) 引入库

(pytorch) root@localhost:/home/Cambricon-MLU270/test/yolov5s-src# git diff

diff --git a/models/experimental.py b/models/experimental.py

index 2dbbf7f..c83c685 100644

--- a/models/experimental.py

+++ b/models/experimental.py

@@ -1,4 +1,6 @@

# This file contains experimental modules

+import matplotlib

+matplotlib.use('Agg')

import numpy as np

import torch

在 models/experimental.py, 头部增加 import matplotlib 和 matplotlib.use('Agg') 两行

2). 修改读取model代码

 ${\tt def attempt\_load} \ ({\tt weights, map\_location=None}):$ 

from models.yolo import Model

model = Model('./models/yolov5s.yaml')

state\_dict = torch.load(weights[0], map\_location='cpu')

model.float().fuse().eval()

model.load\_state\_dict(state\_dict, strict=False)

model.float().fuse().eval()

return model

这里默认使用yolov5s, 所以固定载入 ./models/yolov5s.yaml

3)验证

python detect.py --device cpu --weights unzip.pt

### :/home/Cambricon-Test/yolov5\_demo# python detect.py --device cpu --weights weights/unzip.pt

(pytorch) rost@localhost:/home/Cambricon-Test/yolov5\_demo# python detect.py --device cpu --weights weights/unzip.pt CMML: 7.10.2 08592c0ee CMRT: 4.10.1 0884a9a Namespace(agnostic\_nms=False, augment=False, cfg='cpu', classes=None, conf\_thres=0.25, device='cpu', exist\_ok=False, img\_size=640, iou\_thres=0.45, jit=False, name='exp', project='result', save\_co nf=False, save\_txt=False, source='data/images', update=False, view\_img=False, weights/unzip.pt']) Using torch 1.3.0a0 CPU

	trom	n	params	module	arguments
Θ			3520	models.common.Focus	[3, 32, 3]
1			18560	models.common.Conv	[32, 64, 3, 2]
2			18816	models.common.C3	[64, 64, 1]
3			73984	models.common.Conv	[64, 128, 3, 2]
4			156928	models.common.C3	[128, 128, 3]
5			295424	models.common.Conv	[128, 256, 3, 2]
6			625152	models.common.C3	[256, 256, 3]
7			180672	models.common.Conv	[256, 512, 3, 2]
8			656896	models.common.SPP	[512, 512, [5, 9, 13]]
9			182720	models.common.C3	[512, 512, 1, False]
10			131584	models.common.Conv	[512, 256, 1, 1]
11				torch.nn.modules.upsampling.Upsample	[None, 2, 'nearest']
12	[-1, 6]			models.common.Concat	[1]
13			361984	models.common.C3	[512, 256, 1, False]
14			33024	models.common.Conv	[256, 128, 1, 1]
15				torch.nn.modules.upsampling.Upsample	[None, 2, 'nearest']
16	[-1, 4]			models.common.Concat	
17			90880	models.common.C3	[256, 128, 1, False]
18			147712	models.common.Conv	[128, 128, 3, 2]
19	[-1, 14]		Θ	models.common.Concat	[1]
20			296448	models.common.C3	[256, 256, 1, False]
21			590336	models.common.Conv	[256, 256, 3, 2]
22	[-1, 10]		Θ	models.common.Concat	[1]
23	-1	1 1	182720	models.common.C3	[512, 512, 1, False]
24 [1	7, 20, 23]	1	229245	models.yolo.Detect	[80, [[10, 13, 16, 30, 33, 23], [30, 61, 62, 45, 59, 119], [116, 90, 156, 198, 373, 326]], [128
Model Summa	ry: 283 lay	ers, 7	276605	parameters, 7276605 gradients	
Fucing lave					
Modol Summo	ry: 224 low	0.50 7	266072	parameters 220245 gradients	
Fusing lave	ry. 224 (ay	515, /	2009/3	parameters, 229245 gradients	
Model Summa	rv: 224 lav	are 7	266072	parameters 220245 gradients	
image 1/2 /	homo/Cambri	COD-TO	et/volo	v5 domo/data/imagos/bus_ing: run_cou	
640×490 4 0		on-le	420e)	vs_demo/data/images/bus.jpg. run cpu	
image 2/2 /	home (Cambri	con-Te	et/volo	v5 demo/data/images/zidane_ing: run_cnu	
284×640 2 0	<ul> <li>1 27</li> </ul>	nne (	0 27201	to_demo, data, imageo, i dane. Jpg. Tun epu	
Recults save	ed to result	t/evn	0.3135/		
Done (A 87	/el	er exp			
001101-001-071	131				



到此,官方的pt模型已经可以在MLU pytorch上执行cpu推理。

# 量化

1. 在detect.py中,在 Run inference上增加代码

import torch\_mlu

import torch\_mlu.core.mlu\_quantize as mlu\_quantize

import torch\_mlu.core.mlu\_model as ct

```
. . .
```

global quantized\_model

if opt.cfg == 'qua':

qconfig = {'iteration':2, 'firstconv':False}

quantized\_model = mlu\_quantize.quantize\_dynamic\_mlu(model, qconfig, dtype='int8', gen\_quant=True)

```
# Run inference
t0 = time.time()
img = torch.zeros((1, 3, imgsz, imgsz), device=device) # init img
....
    # Inference
t1 = time_synchronized()
if opt.cfg == 'cpu':
pred = model(img, augment=opt.augment)[0]
print('run cpu')
elif opt.cfg == 'qua':
pred = quantized_model(img)[0]
torch.save(quantized_model.state_dict(), 'yolov5s_int8.pt')
print('run qua')
```

parser.add\_argument('--cfg', default='cpu', help='qua and off')
parser.add\_argument('--jit', type=bool, default=False)

### 2. 执行

if \_\_name\_\_ == '\_\_main\_\_':

python detect.py --weights weights/unzip.pt --source data/images/ --cfg qua

F <mark>alse</mark> , s Using to	ave_txt=False, rch 1.3.0a0 CP	soui U	rce='data	/images/', update=False, view_img≐False	, weights=['weights/unzip.pt'])
A	- 1	1	2520	models common Focus	arguments [3 32 3]
ĩ	-1	î	18560	models.common.Conv	[32, 64, 3, 2]
2	-1	1	18816	models.common.C3	[64, 64, 1]
			73984	models.common.Conv	[64, 128, 3, 2]
			156928	models.common.C3	[128, 128, 3]
5	-1	1	295424	models.common.Conv	[128, 256, 3, 2]
6	-1		625152	models.common.C3	[256, 256, 3]
2	-1		1180672	models.common.conv	[25b, 512, 3, 2]
0	-1		1102720	models.common.SPP	[512, 512, [5, 9, 13]]
10			131594	models.common.Conv	[512, 514, 1, Pates]
11	-1	î	0	torch.nn.modules.upsampling.Upsample	[None 2. 'nearest']
12	[-1, 6]	1	ē	models.common.Concat	
			361984	models.common.C3	[512, 256, 1, False]
14			33024	models.common.Conv	[256, 128, 1, 1]
15	-1	1	Θ	torch.nn.modules.upsampling.Upsample	[None, 2, 'nearest']
16	[-1, 4]		٥	models.common.Concat	
1/	-1		90880	models.common.C3	[256, 128, 1, False]
10	[1 14]		14//12	models.common.conv	
20	(-1, 14)		296//18	models common C3	[256 256 ] False]
21	-1	î	590336	models.common.Conv	[256, 256, 3, 2]
22	[-1, 10]	ĩ	Θ	models.common.Concat	
23			1182720	models.common.C3	[512, 512, 1, False]
24	[17, 20, 23]		229245	models.yolo.Detect	[80, [[10, 13, 16, 30, 33, 23], [30, 61, 62, 45, 59, 119], [116, 90, 156, 198, 373, 326]], [128, 256, 512]]
Model Su	mmary: 283 lay	ers,	7276605	parameters, 7276605 gradients	
Fusing l Model Su Fusing l Model Su The 'ite set this image 1/ image 2/ Results	ayers mmary: 224 lay ayers mmary: 224 lay ration' key in key. We will 2 /home/Cambri saved to resul	ers, qcor calcu con- con- t/exp	7266973 7266973 nfig has ulate qua Test/yolo Test/yolo D2	parameters, 229245 gradients parameters, 229245 gradients been deprecated in this version. Users ntization parameters by real forward ti VS_demo/data/images/bus.jpg; run qua vS_demo/data/images/zidane.jpg; run qua	needn't mes controlled by users.

torch.save(quantized\_model.state\_dict(), 'weights/yolov5s\_int8.pt')

```
这句代码会在weights目录下保存一个量化好的yolov5s_int8.pt文件
```

## 逐层

由于每一层都需要跑在mlu上,通过打印可以看到其他层目前是能支持的,最后一层有mlu算子完成。

修改models/yolo.py的部分代码。

```
class Detect(nn.Module):
    stride = None # strides computed during build
    export = False # onnx export
```

```
def __init__(self, nc=80, anchors=(), ch=()): # detection layer
    super(Detect, self). __init__()
self.anchors_list = list(np.array(anchors).flatten())
    self.num_anchors = len(self.anchors_list)
    self.nc = nc # number of classes
    self.no = nc + 5 # number of outputs per anchor
    self.nl = len(anchors) # number of detection layers
    self.na = len(anchors[0]) // 2 # number of anchors
self.grid = [torch.zeros(1)] * self.nl # init grid
    a = torch.tensor(anchors).float().view(self.nl, -1, 2)
    self.register_buffer('anchors', a) # shape(n1, na, 2)
self.register_buffer('anchor_grid', a.clone().view(self.n1, 1, -1, 1, 1, 2)) # shape(n1, 1, na, 1, 1, 2)
    self.m = nn. ModuleList(nn. Conv2d(x, self.no * self.na, 1) for x in ch) # output conv
    self.img_h = 640
    self.img w = 640
    self. conf_thres = 0.25
    self. iou thres = 0.45
    self.maxBoxNum = 1024
  def forward(self, x):
    # x = x.copy() # for profiling
    z = [] # inference output
    output = []
    self.training |= self.export
    if x[0].device.type == 'mlu':
      for i in range(self.nl):
        x[i] = self.m[i](x[i]) # conv
        y = x[i].sigmoid()
        # print('y.shape: ', y.shape)
        output.append(y)
      detect_out = torch.ops.torch_mlu.yolov5_detection_output(output[0], output[1], output[2],
                self. anchors_list, self. nc, self. num_anchors,
                self.img_h, self.img_w, self.conf_thres, self.iou_thres, self.maxBoxNum)
               # [10, 13, 16, 30, 33, 23, 30, 61, 62, 45, 59, 119, 116, 90, 156, 198, 373, 326]
      return detect_out
    for i in range(self.nl):
      x[i] = self.m[i](x[i]) # conv
      bs, _, ny, nx = x[i]. shape # x(bs, 255, 20, 20) to x(bs, 3, 20, 20, 85)
可以看到,这里将Detect类,修改为,我们mlu支持的算子torch.ops.torch_mlu.yolov5_detection_output
```

detect\_out = torch.ops.torch\_mlu.yolov5\_detection\_output

这里算子是把img宽高固定为 640 640作为输入。

detect.py需要继续增加逐层运行的方式

增加一个get\_boxes的函数,用来从torch.ops.torch\_mlu.yolov5\_detection\_output获取结果后画框,不再依赖cpu的获取框的后处理方式了。

import numpy as np

def get\_boxes(prediction, batch\_size=1, img\_size=640): .... Returns detections with shape: (x1, y1, x2, y2, object\_conf, class\_score, class\_pred) .... reshape\_value = torch.reshape(prediction, (-1, 1)) num\_boxes\_final = reshape\_value[0].item() print('num\_boxes\_final: ', num\_boxes\_final) all\_list = [[] for \_ in range(batch\_size)] for i in range(int(num\_boxes\_final)): batch\_idx = int(reshape\_value[64 + i \* 7 + 0].item()) if batch\_idx  $\geq 0$  and batch\_idx  $\leq$  batch\_size: bl = reshape\_value[64 + i \* 7 + 3].item() br = reshape\_value[64 + i \* 7 + 4].item() bt = reshape\_value[64 + i \* 7 + 5].item() bb = reshape\_value[64 + i \* 7 + 6].item() if bt - bl > 0 and bb - br > 0: all\_list[batch\_idx].append(bl) all\_list[batch\_idx].append(br) all\_list[batch\_idx].append(bt) all\_list[batch\_idx]. append (bb) all\_list[batch\_idx].append(reshape\_value[64 + i \* 7 + 2].item()) all\_list[batch\_idx].append(reshape\_value[64 + i \* 7 + 1].item()) output = [np. array(all\_list[i]).reshape(-1, 6) for i in range(batch\_size)] return output 增加 mlu 的配置方式 def detect(save img=False): source, weights, view\_img, save\_txt, imgsz = opt.source, opt.weights, opt.view\_img, opt.save\_txt, opt.img\_size webcam = source.isnumeric() or source.endswith('.txt') or source.lower().startswith( ('rtsp://', 'rtmp://', 'http://')) # Directories save\_dir = Path(increment\_path(Path(opt.project) / opt.name, exist\_ok=opt.exist\_ok)) # increment run (save\_dir / 'labels' if save\_txt else save\_dir).mkdir(parents=True, exist\_ok=True) # make dir # Initialize set\_logging() device = select\_device(opt.device) half = device.type != 'cpu' # half precision only supported on CUDA # Load model model = attempt\_load(weights, map\_location=device) # load FP32 model imgsz = check\_img\_size(imgsz, s=model.stride.max()) # check img\_size if half: model.half() # to FP16 # Second-stage classifier classify = False if classify: modelc = load\_classifier(name='resnet101', n=2) # initialize modelc.load\_state\_dict(torch.load('weights/resnet101.pt', map\_location=device)['model']).to(device).eval() # Set Dataloader vid\_path, vid\_writer = None, None if webcam: view\_img = True cudnn.benchmark = True # set True to speed up constant image size inference dataset = LoadStreams(source, img\_size=imgsz) else: save\_img = True dataset = LoadImages(source, img\_size=imgsz) # Get names and colors names = model.module.names if hasattr(model, 'module') else model.names colors = [[random.randint(0, 255) for \_ in range(3)] for \_ in names] global quantized\_model global quantized\_net if opt.cfg == 'qua':
 qconfig = {'iteration':2, 'firstconv':False} quantized\_model = mlu\_quantize.quantize\_dynamic\_mlu(model, qconfig, dtype='int8', gen\_quant=True) elif opt.cfg == 'mlu':

from models.yolo import Model model = Model('./models/yolov5s.yaml').to(torch.device('cpu')) model.float().fuse().eval() quantized\_net = torch\_mlu.core.mlu\_quantize.quantize\_dynamic\_mlu(model) #state\_dict = torch. load("./yolov5s\_int8.pt") state\_dict = torch.load(weights[0]) quantized\_net.load\_state\_dict(state\_dict, strict=False) quantized\_net.eval() quantized\_net.to(ct.mlu\_device()) if opt. jit: print("### jit") ct.save\_as\_cambricon('weights/yolov5s\_int8\_1\_4') torch.set\_grad\_enabled(False) ct.set\_core\_number(4) trace\_input = torch.randn(1, 3, 640, 640, dtype=torch.float) input\_mlu\_data = trace\_input.type(torch.HalfTensor).to(ct.mlu\_device()) quantized\_net = torch.jit.trace(quantized\_net, input\_mlu\_data, check\_trace=False) 增加mlu的推理方式 # Run inference t0 = time.time() img = torch.zeros((1, 3, imgsz, imgsz), device=device) # init img
\_ = model(img.half() if half else img) if device.type != 'cpu' else None # run once for path, img, imOs, vid\_cap in dataset: img\_letter = letterbox0(imOs, new\_shape=imgsz)[0] img = torch. from\_numpy(img). to(device) img = img.half() if half else img.float() # uint8 to fp16/32 img /= 255.0 # 0 - 255 to 0.0 - 1.0 if img.ndimension() == 3: img = img.unsqueeze(0) # Inference t1 = time\_synchronized() #pred = model(img, augment=opt.augment)[0] if opt.cfg == 'cpu': pred = model(img, augment=opt.augment)[0] print('run cpu') elif opt.cfg == 'qua': pred = quantized\_model(img)[0] torch.save(quantized\_model.state\_dict(), 'weights/yolov5s\_int8.pt') print('run qua') elif opt.cfg == 'mlu': img = img.type(torch.HalfTensor).to(ct.mlu\_device()) img = img. to(ct.mlu\_device()) pred = quantized\_net(img)[0] pred=pred.data.cpu().type(torch.FloatTensor) box\_result = get\_boxes(pred) det = box\_result[0].tolist() with open("result/yolov5s\_mlu\_output.txt", "w+") as f: for \*xyxy, conf, cls in reversed(det): f.write("{}\n{}\n{}\n{}\n".format(\*xyxy)) label = f' {names[int(cls)]} {conf:.2f} plot\_one\_box(xyxy, img\_letter, label=label, color=colors[int(cls)], line\_thickness=1) import os cv2. imwrite("result/mlu\_out\_{}.jpg".format(os.path.basename(path).split('.')[0]), img\_letter) img = img \* 255 print('run mlu') elif opt.cfg == 'cpu': pred = model(img, augment=opt.augment)[0] print('run cpu') if opt.cfg != 'cpu': continue # Apply NMS

执行

python detect.py --device cpu --weights weights/yolov5s\_int8.pt --source data/images/ --cfg mlu

结果:

Fusing	using layers									
Model S	Summary: 224 layers, 7266973 parameters, 229245 gradients									
Fusing	Jsing Lavers									
Model S	Summary: 224 lay	ayers, 7266973 parameters, 229245 gradients								
	from		params	module	arguments					
Θ			3520	models.common.Focus	[3, 32, 3]					
1			18560	models.common.Conv	[32, 64, 3, 2]					
2			18816	models.common.C3	[64, 64, 1]					
3	- 1	1	73984	models.common.Conv	[64, 128, 3, 2]					
4	- 1	1	156928	models.common.C3	[128, 128, 3]					
5	- 1	1	295424	models.common.Conv	[128, 256, 3, 2]					
6	-1	1	625152	models.common.C3	[256, 256, 3]					
7	-1	1	1180672	models.common.Conv	256, 512, 3, 2					
8	-1	1	656896	models.common.SPP	[512, 512, [5, 9, 13]]					
9	-1	1	1182720	models.common.C3	[512, 512, 1, False]					
10	-1	1	131584	models.common.Conv	[512, 256, 1, 1]					
11	-1	1	Θ	torch.nn.modules.upsampling.Upsample	[None, 2, 'nearest']					
12	[-1, 6]		Θ	models.common.Concat						
13	-1		361984	models.common.C3	[512, 256, 1, False]					
14	-1		33024	models.common.Conv						
15	-1		Θ	torch.nn.modules.upsampling.Upsample	[None, 2, 'nearest']					
16	[-1, 4]		0	models.common.Concat						
1/	-1		90880	models.common.C3	[256, 128, 1, False]					
18			14//12	models.common.Conv	[128, 128, 3, 2]					
19	[-1, 14]		0	models.common.Concat						
20	-1		296448	models.common.us	[256, 256, 1, False]					
21			590336	models.common.conv	[256, 256, 3, 2]					
22	[-1, 10]		0	models.common.concat						
23	1- 00 001		1182720	models.common.cs	[312, 312, 1, False] [40. [10. 10. 16. 20. 22. 22] [20. 61. 62. 45. 50. 120] [116. 00. 156. 100. 270. 226]] [120. 256. 510]]					
24 Madal 6	[17, 20, 23]	1	229245	models.yoto.belect	[80, [[10, 13, 10, 30, 33, 23], [30, 01, 02, 45, 54, 114], [110, 40, 150, 148, 373, 320]], [128, 250, 512]]					
modet s	Summary: 283 tay	ers,	/2/0005	parameters, 7276665 gradients						
FUSING Layers										
Would Summary: 224 Cayers, 7260973 parameters, 229245 gradients										
Datchmount 1 Jamaa 1/ Jama/Cambricon-Test/valav5.dama/data/imagas/hus jag: pum bayas final: 5.0										
in a guine and a convest of yours _uemo/ data / images/ous.jpg. tum_boxes_lintet										
hard Mine 1										
image 2	image 2/2 /home/Cambricon-Test/volov5 demo/data/images/zidane.jpg: num boxes final: 4.0									
anage are from ready for ready for a second and a second and a second and a second a second a second a second a										

run mlu Results saved to result/exp3 Done. (11.260s)



```
elif opt.cfg == 'mlu':
    from models.yolo import Model
    model = Model('./models/yolov5s.yaml').to(torch.device('cpu'))
    model.float().fuse().eval()
    quantized_net = torch_mlu.core.mlu_quantize.quantize_dynamic_mlu(model)
    #state_dict = torch.load("./yolov5s_int8.pt")
    state_dict = torch.load("./yolov5s_int8.pt")
    state_dict = torch.load(weights[0])
    quantized_net.load_state_dict(state_dict, strict=False)
    quantized_net.eval()
    quantized_net.to(ct.mlu_device())
```

```
if opt.jit:
    print("### jit")
    ct.save_as_cambricon('weights/yolov5s_int8_1_4')
    torch.set_grad_enabled(False)
    ct.set_core_number(4)
    trace_input = torch.randn(1, 3, 640, 640, dtype=torch.float)
    input_mlu_data = trace_input.type(torch.HalfTensor).to(ct.mlu_device())
    quantized_net = torch.jit.trace(quantized_net, input_mlu_data, check_trace=False)
```

执行

python detect.py --weights weights/yolov5s\_int8.pt --cfg mlu --jit True

最后得到 yolov5s\_int8\_1\_4. cambricon 离线模型

结果:



完整代码参考: http://gitlab.software.cambricon.com/neuware/software/ae/solution/270\_sample/yolov5\_demo