

## 第7章 综合实验

前述章节实验完成了图像处理领域中的风格迁移应用在智能处理器上的完整开发和优化。随着人工智能在更广泛领域取得了良好效果，本章的综合实验将涉及在多个不同领域（如目标检测、文本识别、自然语言处理和语音合成等）的人工智能应用在智能处理器上的开发和优化。具体而言，第7.1节介绍基于YOLOv3网络模型<sup>[20]</sup>实现目标检测应用，并借助TensorFlow框架及流处理框架在智能处理器上进行部署。第7.2节介绍基于EAST网络模型<sup>[21]</sup>实现文字识别应用在智能处理器上的部署。第7.3介绍基于BERT网络模型<sup>[22]</sup>在智能处理器上的部署。

### 7.1 目标检测-YOLOv3

#### 7.1.1 实验目的

本实验主要完成将目标检测的代表性算法——YOLOv3网络移植到智能处理器DLP上，并进行性能优化与离线部署，使读者可以借助DLP处理器完成完整的目标检测任务。

因此，本实验的主要目的包括：

1. 通过用智能编程语言实现YOLOv3的部分后处理操作，深入掌握智能编程语言的算子开发和优化方法；
2. 通过在TensorFlow中添加大算子，掌握在TensorFlow框架中添加融合算子的方法；
3. 通过使用TensorFlow进行在线推理，掌握使用TensorFlow编写目标检测应用并在典型DLP上进行优化的方法；

实验时间预计为两周。

#### 7.1.2 背景介绍

##### 7.1.2.1 目标检测

目标检测是指针对一个或多个特定类别在数字图像中进行视觉对象的分类和定位。目标检测属于多任务学习，也就是说一个任务分支需要通过分类算法将物体根据类别划分，另一个任务分支是在判断类别后标记出每一个类别的具体位置信息，即中心坐标位置或坐标位置（以目标图片的左上角为原点 $(0,0)$ ）。随着技术的发展，目标检测逐渐成为其他计算机视觉任务的基础，如图像分割、图像字幕、物体追踪等；目前，目标检测已广泛应用于自动驾驶、机器人视觉、安防监控等领域。目标检测的两个任务分支如下图7.2所示：

在介绍目标检测算法原理前，先列举一些目标检测中的常见概念如表7.1所示，包括平均精度均值（mean Average Precision, mAP）、滑动窗口（Sliding Window）、选择搜索算法

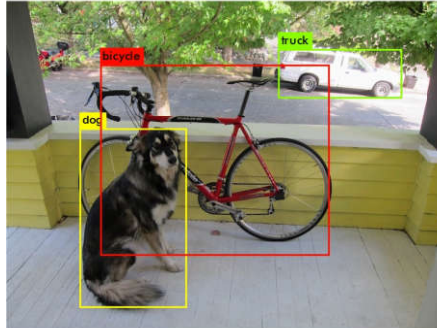


图 7.1 目标检测效果

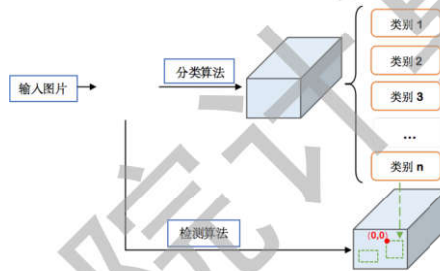


图 7.2 目标检测的多个任务分支

表 7.1 目标检测常见知识

名称	定义
平均精度均值 mAP	AP 定义为 Precision-recall 曲线下的面积，mAP 为多个检测目标类别的 AP 均值，区间在 [0,1] 之间，越接近 1 证明精确度越高
滑动窗口 Sliding Window	采用不同大小和比例（宽高比）的窗口在全图片上以一定的步长进行滑动，以便于对应的区域做图像分类
选择搜索算法 Selective Search	将图片分为多个子区域，根据子区域的相似性不断迭代合并，在此过程中对合并区域实现外切矩形，得到候选框
交并比 IOU	用于预测 Bounding-box 和人工标注 Bounding-box 的重叠度，两个矩形框的 IOU 计算公式： $IOU=(A \cap B)/(A \cup B)$
兴趣区域 (RoI)	从被处理的图像以方框或不规则形状等方式勾勒出需要处理的区域
非极大值抑制 NMS	依据候选区域和 score 矩阵，判断置信度最高的 bounding box，若预测框重叠，则选择置信度得分高的边框

(Selective Search, SS)、交并比 (Intersection over Union, IoU)、兴趣区域 (Region of Interest, RoI)、非极大值抑制 (Non-Maximum Suppression, NMS) 等。

依据深度学习算法阶段不同, 目标检测方法主要可以分为两类: 一种是基于 Region Proposal 的两阶段检测, 它首先由算法计算出目标候选框, 然后再通过卷积神经网络对目标候选框进行分类与回归, 主流的算法有 R-CNN、Fast R-CNN、Faster R-CNN 等; 另一种是单阶段检测, 不再需要候选区域, 仅使用卷积神经网络直接将目标定位框转为回归问题, 预测不同类别的目标位置, 主流的算法有 YOLO、SSD 等。

### • 两阶段检测

我们以典型的 Faster R-CNN 算法介绍两阶段检测算法的主要流程。

Faster R-CNN 针对 FAST RCNN 检测速度的瓶颈做了进一步改进, 摆脱了传统意义上的 SS 算法的限制, 取而代之的是使用 Region Proposal Network (RPN) 生成 RoI, 计算步骤如下:

- 通过共享卷积层提取特征图
- 采用 RPN 提取特征图中的目标候选框, 判断 RoI 后进行首次边框修正
- 利用 RoI Pooling Layer 在 RPN 输出的特征图上计算 RoI 对应的固定维度特征
- 将目标框使用全连接层进行分类, 并对边框进行二次修正

Faster R-CNN 是第一个真正意义上的端到端的算法, 将 RPN 替代 SS 使得检测速度大幅提升, 同时在 VOC07 数据集上的 mAP 达到 73.2%<sup>[23]</sup>。Faster R-CNN 有一个主要问题在于后续的检测阶段候选框区域提取上耗时较长。

### • 单阶段检测

我们以典型的 YOLO 算法介绍单阶段检测算法的主要流程。

针对两阶段检测算法存在的计算冗余、检测速度较慢等问题, 在 YOLO (You Only Look Once) 算法中提出“一步完成”的单级检测器, 这意味着将目标分类和目标检测在同一个步骤中实现, 其算法的主要步骤如下:

- 图片缩放: 输入图片格式调整为 448\*448;
- 网络预测: 采用单个卷积神经网络对整个图像进行 Region Proposal 的选取以及目标置信度和位置输出;
- 预测框筛选: 利用 IOU 筛选 Bounding Box 得到目标检测的预测框。

和两阶段检测网络相比, 单阶段检测网络的速度得到大幅度提升, 但由于它对输入图像的分割较粗糙, 导致后续目标定位精度有所下降, 对于小目标的检测结果不尽如人意。

### 7.1.2.2 YOLOv3

本节实验所使用的YOLOv3网络是经典的单阶段检测方法,它在延续之前版本YOLOv1和YOLOv2主要特点(包括:通过划分单元格来做检测,使用Leaky ReLU<sup>[24]</sup>作为激活函数,实现端到端训练,使用Batch Normalization<sup>[7]</sup>避免过拟合、使用多尺度训练等)的基础上,通过修改Backbone网络来提升精度与性能。

下面主要介绍YOLOv3的网络结构和实验中需要采用编程语言实现的采用NMS来进行候选框筛选的部分。

#### • 网络结构

图7.3描述了YOLOv3的结构图,包含了从图片输入经过Darknet-53,Concat拼接到输出3个不同尺度的Feature Map的过程。

其中DBL是Yolov3的基本计算单元,它由一个DarknetConv2d、一个BatchNorm和一个LeakyRelu操作线性排列组成,BatchNorm和LeakyRelu是conv2d卷积层后必不可少的部分(不包括最后一个卷积层),每个算子的功能如下表7.2所示:

表 7.2 DBL 的基本计算单元

算子名称	功能
DarknetConv2d	Darknet是Yolov3的基础网络,1个Darknet的2维卷积Conv2D层,即DarknetConv2D():操作流程:通过kernel_regularizer将卷积核的参数进行L2正则化,当Padding的strides=(2,2)时采用valid模式,其他Padding一般采用same模式;其余操作与常规二维卷积Conv2D()保持一致。
BN	BatchNormalization(),将输入数据进行正则化,把每层神经网络的输入值的分布强制拉回到均值为0方差为1,的标准正态分布,得到的输出值作为激活输入数据,分布在非线性函数的敏感区,加快收敛速度。
Leaky()	LeakyReLU激活函数是ReLU的变换,公式为: 合 $y = \max(0, x) + leak * \min(0, x)$ ; (斜率为0.1)

Res\_unit是Resblock\_body的基本组件之一,它是由两个DBL和一个ADD组合,借鉴了Resnet的残差结构,即将每个Res\_unit的输入经过两次DBL处理后生成的特征图与输入自身进行叠加(Add),Add操作仅是直接相加并不会导致Tensor维度的改变。引入残差结构使Yolov网络深度从Darknet-19(Yolov2)增加至Darknet-53(Yolov3),残差层可以减小网络过深而引起的梯度爆炸风险,提高网络的训练效率。

Resblock\_body在网络结构中作为Yolov3的大组件,用Res\_n表示,n为在此结构中Res\_unit的个数,它是由一个ZeroPadding2D、一个DBL、N个Res\_unit线性排列组成。ZeroPadding2D根据下一步操作的需要,将输入矩阵的边界进行0充填,以达到信息补齐的效果。

Concat是将不同层的Tensor进行拼接,在Yolov3中主要是将darknet中间层的输出与上采样结果进行维度扩张,以获得不同的输出。

Darknet.output输出的底层信息包含全局特征信息,DarkNet中间层包含局部特征信息,通过Concat拼接,可提高检测精度。

值得注意的是 YOLOv3 的对象分类预测并没有延续 YOLOv2 中的 Softmax，而是采用了 Logistic Regression 对输出的 bbox 进行预测类别，这是因为 Softmax 要求每个目标对象具有相互独立的标签，而数据集中每个对象需要支持多个标签，将 Softmax 用 Logistic 可以预测每个类别的得分。

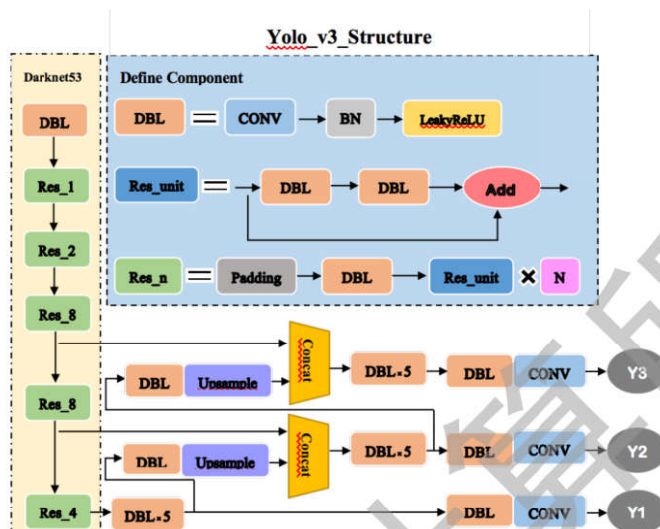


图 7.3 YOLOv3 结构图

对于网络输出，YOLOv3 中的获得 anchorbox 的方案依旧使用了 YOLOv2 的聚类法，首先计算出 9 个聚类中心，再按照大小均匀分给 3 个尺度，经过卷积层对输出采用多尺度检测技术，也就是说有个三个不同尺度的输出特征图，分别为  $13 \times 13$ 、 $26 \times 26$ 、 $52 \times 52$ ，其深度都为 255。由于每个网络单元都有 3 个检测框，每个框包含检测框中心坐标  $(x,y)$ 、检测框的长度  $(h)$ 、检测框的宽度  $(w)$  以及检测框的置信度  $(confidence)$  五个基本参数，80 个检测类别都有对应的概率，所以特征图的深度为 255  $(3 \times (5 + 80) = 255)$ 。三个不同尺度的输出如下表 7.3 所示，对应与图 7.3 中的三个输出 (Y1、Y2 和 Y3)：

表 7.3 YOLOv3 的网络输出

尺度	检测类型	实现方式
13*13	大型目标	Darknet 的网络最终输出结构为 13131024，通过设置 512 通道个数和 $num\_anchors * (num\_classes + 5)$ ，得到预测矩阵 Y1。
26*26	中型目标	通过 DBL 将通道数由 512 调整至 256，再经过上采样 UpSampling2D，将 $13 \times 13$ 的网络结构转换为 $26 \times 26$ ，得到的网络结构与 DarkNet 的中间层输出结果进行拼接，最终结果即是第 2 个尺度特征图 Y2。
52*52	小型目标	通过 DBL 将通道数由 256 调整至 128，再经过上采样 UpSampling2D，将 $26 \times 26$ 的网络结构转换为 $52 \times 52$ ，得到的网络结构与 DarkNet 的中间层输出结果进行拼接，最终结果即是第 3 个尺度特征图 Y3。

YOLOv3 中使用 Darknet-53 作为 backbone 网络用以提取图片的特征信息，相较于 YOLOv2 中的 Darknet-19，其使用步长为 2 的卷积来替代之前的池化操作，同时采用 ResNet<sup>[2]</sup> 的残差结构 (resblock) 提升精度。与 ResNet-152 和 ResNet-101 相比，Darknet-53 网络层数较少，不仅分类精度没有下降，检测速度也有所提升。图 7.4 展示了 Darknet-53 的结构，通过 5 次

步长为 2 的卷积进行尺寸变换，最终输出的 Feature Map 为输入的 1/32。

Type	Filters	Size	Output
Convolutional	32	3 × 3	256 × 256
Convolutional	64	3 × 3 / 2	128 × 128
Convolutional	32	1 × 1	128 × 128
Convolutional	64	3 × 3	
1x Residual			128 × 128
Convolutional	128	3 × 3 / 2	64 × 64
2x Convolutional	64	1 × 1	64 × 64
Convolutional	128	3 × 3	
2x Residual			64 × 64
Convolutional	256	3 × 3 / 2	32 × 32
8x Convolutional	128	1 × 1	32 × 32
Convolutional	256	3 × 3	
8x Residual			32 × 32
Convolutional	512	3 × 3 / 2	16 × 16
8x Convolutional	256	1 × 1	16 × 16
Convolutional	512	3 × 3	
8x Residual			16 × 16
Convolutional	1024	3 × 3 / 2	8 × 8
4x Convolutional	512	1 × 1	8 × 8
Convolutional	1024	3 × 3	
4x Residual			8 × 8
Avgpool		Global	
Connected		1000	
Softmax			

图 7.4 Darknet-53 结构图

Yolov3 中的 Darknet-53 网络结构有以下两个特点：

1. 将全卷积层应用于整个网路: Darknet-53 是 53 个卷积层和池化层的组合，但在 Yolov3 中并没有包含 darknet-53 最后一个全局平均池化层，所以减小张量维度是通过改变 5 次卷积核的步长来实现的。全卷积结构也应用于预测分支，最后一个卷积核的规模为  $1 \times 1 \times 1024 \times 255$ ，其中 255 是针对 COCO 数据集的 80 类别确定的卷积核个数。
2. 借鉴 Resnet 思想，引入残差 (Residual) 结构: 网络深度越深，特征信息表达越好，但随着网路深度的加深，训练误差趋势会先下降再上升，也就意味着网络深度越深，模型训练难度越大；所以采用残差模块不仅可以保证在网络深度加深的情況仍能保持收敛状态，还在一定程度上大大减少了计算量。

• 非极大值抑制: NMS

得到 3 种不同尺度的 Feature Map 后需要进行后处理操作来筛选出合适的框以完成检测，NMS 非极大值抑制是其中重要的操作。NMS<sup>[25]</sup> (非极大值抑制) 是一种针对目标框的多重性而提出的检测技术，是目标检测网络中一个重要环节，由于目标检测相邻的交叉窗口经常出现分数相近的情况，使用 NMS 可以去掉重复的检测框。下图 7.5 展示了从输出预选框到生成最终检测框的过程 (待增加更多 NMS 原理流程部分介绍)。

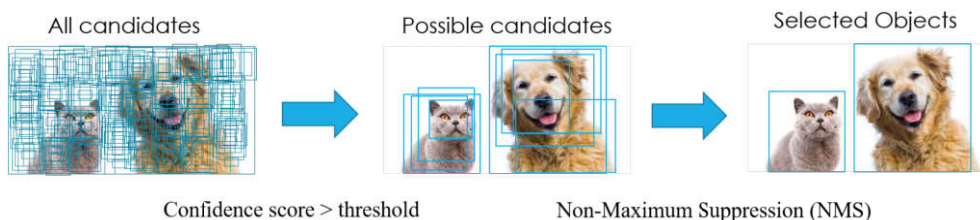


图 7.5 NMS 示意图

NMS 的计算过程可概括为以下四个步骤：

1. 对所有预测框的置信度降序排序；
2. 选出置信度最高的预测框，确认其为正确预测，并计算它与其他预测框的 IOU；
3. 根据 (2) 中计算的 IOU 去除重叠度高的框，即  $IOU > thresh\_iou$  就删除；
4. 剩下的预测框返回第 1 步，直到没有置信度  $> thresh\_score$  的框为止。

由于目标检测类算法一般都包含 NMS 计算环节，之外许多其它场景下也可以用到 NMS 计算过程（例如文本识别场景）。在这样的背景下，有必要将 NMS 模块在 DLP 硬件上的实现整体抽象剥离出来，单独封装成一个 NMS 通用模块，不仅适用于任何检测算法，还可以缩短其它检测算法的开发时间，增加代码的可维护性与稳定性。因此检测网络综合实验中设计了一个 NMS 通用模块任务，以.h 文件提供给调用程序，将 NMS 计算过程封装成 `nms_detection` 函数，供 YOLOv3 后处理直接调用。

### 7.1.3 流处理框架

为了方便在 DLP 上的离线部署，DLP 硬件在运行时库的基础上进一步提供了流处理框架 CNStream，用于完成从数据处理到离线推理等任务的流水线处理。CNStream 是面向 DLP 的数据流处理 SDK。用户可以根据 CNStream 提供的接口，开发实现自己的组件。还可以通过组件之间的互连，灵活地实现自己的业务需求。CNStream 能够大大简化 DLP 提供的推理和其他处理，如视频解码、神经网络图像前处理的集成。也能够在兼顾灵活性的同时，充分发挥 DLP 的硬件解码和机器学习算法的运算性能。

CNStream 基于模块化和流水线的思想，提供了一套基于 C++ 语言的接口来支持流处理多路并发的 Pipeline 框架。为用户提供可自定义的模块机制以及通过高度抽象的 CNFrame-Info 类型进行模块间的数据传输，满足用户对性能和可伸缩性的需求。

### 7.1.4 实验环境

本节实验所涉及的硬件平台和软件环境如下：

- 硬件平台：硬件平台使用第 ?? 节介绍的云平台中的 DLP 硬件。
- 软件环境：所涉及的 DLP 软件开发模块包括编程框架 TensorFlow、高性能库 CNML、运行时库 CNRT、编程语言及编译器、运行时库 CNRT 以及流式编程框架 CNStream。

### 7.1.5 实验内容

本节实验主要是在 DLP 硬件上实现并优化以 YOLOv3 网络模型为核心的目标检测应用。针对用户从网络上下载的标准 YOLOv3 模型，经过模型量化并配置 DLP 相关参数，则可以采用 DLP 上的定制 TensorFlow 版本来运行。为了充分发挥 DLP 的计算能力，进一步采用智能编程语言 BCL 实现 YOLOv3 中的 NMS 部分，并将其作为一个大算子集成到 TensorFlow 框架中。考虑推理应用的实际需求，采用 DLP 的流式框架编写完整目标检测应用，基于离线模型进行部署。具体实验内容包括：

1. **YOLOv3 DLP 运行**: 将标准的 YOLOv3 网络模型进行量化后, 通过定制的 TensorFlow 版本在 DLP 硬件上运行;
2. **NMS 的 BCL 实现**: 采用 BCL 实现 NMS 后处理运算;
3. **框架算子集成**: 将用 BCL 实现的 NMS 算子集成到 TensorFlow 框架中, 进行在线推理。

## 7.1.6 实验步骤

如前所述, 详细的实验步骤主要包括: YOLOv3 DLP 运行、NMS 的 BCL 实现、框架算子集成等。

### 7.1.6.1 YOLOv3 DLP 运行

为了使得标准的 YOLOv3 网络模型在 DLP 上运行, 主要需要完成: YOLOv3 CPU 运行、模型量化和 DLP 运行等流程。

#### • YOLOv3 CPU 运行

通过在 CPU 上正确执行 YOLOv3 标准网络模型完成推理, 确保相关依赖都已安装完成。具体步骤如下:

1. 获取开源 YOLOv3 工程: `git clone https://github.com/YunYang1994/tensorflow-yolov3.git`
2. 安装推理所需软件包: 一些必备的 Python 安装包在所提供的云平台软件环境中已经安装完成, 但对于每个综合实验的特有依赖还需进行单独安装, 使用图 7.6 中的命令完成依赖的安装与升级。

```
1 $ cd tensorflow-yolov3
2 $ pip install -r ./docs/requirements.txt
```

图 7.6 YOLOv3 依赖安装

3. 模型下载与转换: 执行图 7.7 中的命令完成模型的下载与转换。要执行这一步的原因是模型文件和 ckpt 文件之间有一些版本上的差异, 需要将节点名等细节对齐。

```
1 $ cd checkpoint
2 $ wget https://github.com/YunYang1994/tensorflow-yolov3/releases/download/v1.0/
   yolov3_coco.tar.gz
3 $ tar -xvf yolov3_coco.tar.gz
4 $ cd ..
5 $ python convert_weight.py
6 $ python freeze_graph.py
```

图 7.7 YOLOv3 模型下载与转换

4. 模型推理与验证: 执行图 7.8 中的命令完成在 CPU 上的推理。



```
1 $ python image_demo.py
```

图 7.8 YOLOv3 CPU 推理

### • 模型量化

下面介绍如何将开源的 pb 模型量化为 INT8 pb 模型。DLP 软件环境提供 fppb\_to\_intpb 工具来进行模型量化。该工具在 TensorFlow 根目录的/tensorflow/cambricon\_examples/tools/fppb\_to\_intpb 目录下。具体的量化分为以下几个步骤：

1. **生成数据文件**:按照模型所需要的数据集以及数据集的具体路径来修改 generate\_image\_list.py 文件中的 image\_libs\_map, 来生成相应的数据文件。YOLOv3 需要 coco 数据集, 因此将 image\_libs\_map 进行如图 7.9所示的修改。完成修改后执行 generate\_image\_list.py 来生成数据文件。

```
1 import os
2 from os.path import isfile, join
3
4 max_image_count = 100
5 image_libs_map = {
6     "coco":join(os.environ.get( "COCO_DATASET_HOME" ), "COCO/test2017")
7 }
8 if __name__ == "__main__":
9     for image_lib, image_path in image_libs_map.items():
10         n = 0
11         with open("image_list_{}".format(image_lib), "w") as image_list:
12             for root, dirs, files in os.walk(image_path):
13                 for f in files:
14                     if n >= max_image_count:
15                         break
16                     if f.endswith(".jpg"):
17                         image_list.write(join(image_path, f) + "\n")
18                         n = n + 1
```

图 7.9

2. **生成量化配置文件**: 如第 ?? 所述, 在进行量化前需设置相应的配置文件。以指定量化相关的参数。DLP 软件环境提供了 generate\_ini.py 来生成指定模型的量化配置文件。针对 YOLOv3, 需要将其中的 model\_name 进行如图 7.10所示的修改, 即指定要生成的配置文件是针对 YOLOv3 的。

```
1 models_name = ['yolov3']
```

图 7.10

执行 generate\_ini.py 脚本后会生成 config 文件夹, 其中包含了 yolov3\_naive\_int8.ini 参数配置文件。该配置文件中包含如图 7.11所示的信息。

其中, 可以通过 activation\_quantization\_alg 和 weight\_quantization\_alg 来设置具体的量化模式。包括 naive 和 threshold\_search 两种方式。

```

1 [preprocess]
2 mean = 0, 0, 0      #均值, 顺序依次为 mean_r、 mean_g、 mean_b
3 std = 255.0        #方差
4 color_mode = rgb   #网络的输入图片是 rgb 还是 bgr
5 crop = 416, 416    #前处理最终将图片处理为 416 * 416 大小
6 calibration = yolov3_preprocess_cali #校准数据读取及前处理的方式, 可以根据需求进行自
   定义, [preprocess] 和 [data] 中定义参数均为 calibration 的输入参数
7
8 [config]
9 activation_quantization_alg = naive #输入量化模式, 可选 naive 和 threshold_search, naive
   为基础模式, threshold_search 为阈值搜索模式
10 device_mode = clean #可选 clean、mlu 和 origin, 建议使用 clean, 使用 clean
   生成的模型在运行时会自动选择可运行的设备
11 use_convfirfirst = False #是否使用 convfirfirst
12 quantization_type = int8 #量化位宽, 目前可选 int8 和 int16
13 debug = False #是否为debug模式
14 weight_quantization_alg = naive #权重量化模式, 可选 naive 和 threshold_search, naive 为
   基础模式, threshold_search 为阈值搜索模式
15 int_op_list = Conv, FC, LRN #要量化的 layer 的类型, 目前只能量化 Conv、FC 和 LRN
16 channel_quantization = False #是否使用分通道量化, 目前不支持为 True
17
18 [model]
19 output_tensor_names = pred_sbbox/concat_2:0, pred_mbbox/concat_2:0, pred_lbbox/concat_2
   :0 #输出 Tensor 的名字, 可以是多个, 以逗号隔开
20 original_models_path = ./realpath-of-yolov3_coco.pb #输入 pb
21 save_model_path = ./pbs/yolov3/yolov3_int8.pb #输出 pb
22 input_tensor_names = input/input_data:0 #输入 Tensor 的名字, 可以是多个, 以逗号隔开
23
24 [data]
25 num_runs = 2 #运行次数
26 data_path = ./image_list_coco #校准数据文件路径
27 batch_size = 10 #每次运行的 batch_size

```

图 7.11

naive 模式对应于第??节的量化方式, 直接统计数据的最大值和最小值完成量化。

threshold\_search 阈值搜索模式用于处理存在异常值的待量化数据集, 该模式能够过滤部分异常值, 重新计算出数据集的最值, 用新最值来计算数据集的量化参数, 从而提高数据集整体的量化质量。在这种模式下, 数据量化所使用的最大值和最小值不再由简单的统计得出, 而是通过搜索的方式得到。但是对于不存在异常值, 数据分布紧凑的情况下, 不建议使用该算法, 比如网络权值的量化。

其中 device\_mode 可以设置输出 pb 所有节点的 device, 有三种配置方式: clean、mlu 和 origin。其中, mlu 将输出 pb 的所有节点的 device 都设置为 MLU, 即都在 MLU 上运行; clean 将输出 pb 的所有节点的 device 清除, 运行时根据算子注册情况自动选择可运行的设备; origin 则使用和输入 pb 一样的设备指定。

3. **完成模型量化:** 运行量化工具 python fppb\_to\_intpb.py config/yolov3\_naive\_int8.ini 完成模型的量化。

#### • DLP 运行

完成 CPU 侧的推理后, 接下来需要将其移植在 DLP 上, 以加快其推理速度。此时我

们只需仿照第7章编程语言算子实验中的DLP推理移植部分，通过配置session config以及进行INT8量化即可完成DLP移植。下面将详细介绍这两个步骤：

修改上述执行CPU推理时的代码image\_demo.py，增加如下几个部分：

```

1 pb_file = "../realpath-to-int8-pb"
2 .....
3 config = tf.ConfigProto(allow_soft_placement=True,
4                         inter_op_parallelism_threads=1,
5                         intra_op_parallelism_threads=1)
6 config.mlu_options.data_parallelism = 1
7 config.mlu_options.model_parallelism = 1
8 config.mlu_options.core_num = 4
9 config.mlu_options.core_version = "MLU270"
10 config.mlu_options.precision = "int8"
11 config.mlu_options.save_offline_model = False
12 config.mlu_options.offline_model_name = "yolov3_int8.cambricon"
13 with tf.Session(config = config, graph=graph) as sess:
14     .....

```

图 7.12 YOLOv3 DLP 移植 (1)

### 7.1.6.2 NMS BangC 代码实现

在使用BCL设计NMS模块时，需要考虑以下涉及原则：

1. 能够处理的数据来源：GDRAM, NRAM, SRAM中的一种
2. 能够将数据存放到：GDRAM, NRAM, SRAM中的一种
3. 能够以不同的储存格式存放数据
4. 能够在BLOCK和U1的模式下工作
5. 能够处理任意数据规模

使用BangC实现时，如果按照背景介绍章节中NMS原理进行计算，则会显得较为复杂，比如第(3)步中删除交并比大于阈值的框，使用BangC语言实现这个操作需要来回拷贝数据，导致效率不高。为了解决此问题，标准程序中采用另外一种方式，即不改变之前的数据，将被删除框对应的分数置为0，下一轮筛查的时，删除的框必然不会对此轮筛查造成影响。整体设计流程图如图7.13所示：

下面以图7.14中的代码为例进行说明。NMS通用模块接口介绍：

- output\_box\_num: NMS筛选出的框的总个数；
- output\_data: NMS计算结果存放的数据地址；
- dst: NMS计算结果存放的数据地址类型：GDRAM/SRAM/NRAM；
- input\_data\_score: 输入框score的数据地址；

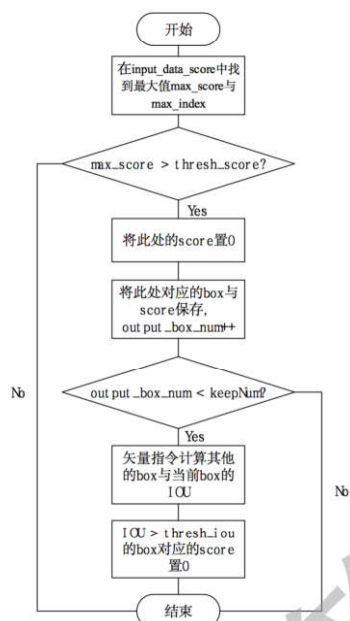


图 7.13 BangC NMS 设计流程图

```

1 __mli_func__ void nms_detection(int& output_box_num,
2                               NMS_DT* output_data,
3                               Addr dst,
4                               NMS_DT* input_data_score,
5                               NMS_DT* input_data_box,
6                               Addr src,
7                               NMS_DT *buffer,
8                               int buffer_size,
9                               NMS_DT* sram,
10                              SplitMode split_mode,
11                              int input_box_num,
12                              int input_stride,
13                              int output_stride,
14                              int keepNum,
15                              NMS_DT thresh_iou,
16                              NMS_DT thresh_score,
17                              int save_method);

```

图 7.14 接口设计

- `input_data_box`: 输入框坐标的数据地址, 储存顺序是: `x1, y1, x2, y2`, 同一个类型的数据储存在一起;
- `src`: NMS 输入数据存放的数据地址类型: GDRAM/SRAM/NRAM;
- `buffer`: NMS 计算使用的 NRAM 空间首地址;
- `buffer_size`: NMS 计算使用的 NRAM 空间大小, 单位为字节, 当 `dst` 为 GDRAM/SRAM, 则 `buffer_size` 至少需要:  $(64 * 4 + 64 + 256 * 5) * \text{sizeof}(\text{NMS\_DT})$  字节; 当 `dst` 为 NRAM, 则 `buffer_size` 至少需要:  $(64 * 4 + 64) * \text{sizeof}(\text{NMS\_DT})$  字节。注意: 如果 `src == NRAM` 并且 `input_box_num` 是对齐的, NMS 计算会更加高效, 因为可以省去数据搬运工作。

- sram: SRAM 的地址空间, 用来在 U1 模式下通过核间通讯找 score 的最大值, 大小至少为  $30 * \text{sizeof}(\text{NMS\_DT})$  字节;
- split\_mode: 拆分模式: NMS\_BLOCK/NMS\_U1;
- input\_box\_num: 输入框的数量;
- input\_stride: 输入数据的 stride, 即 x1, y1, x2, y2 之间的 stride;
- output\_stride: 输出数据的 stride, 即 x1, y1, x2, y2 之间的 stride;
- keepNum: 最多保留筛选框的个数;
- thresh\_iou: 交并比阈值;
- thresh\_score: score 阈值;
- save\_method: 储存模式, 分为 0/1/2

具体实现主要分为两个步骤:

#### 1. 准备阶段: 准备所需的变量, 防呆检测, 空间划分, 多核拆分

##### (a) 主要变量有:

- core\_limit: 启用的核数, BLOCK 模式对应 1, U1 模式对应 4;
- loop\_end\_flag: U1 模式结束标识;
- nram\_save\_limit\_count: NRAM 上临时储存筛选框的数量;
- MODE: 当 src=NRAM 时, 0 代表数据需要先加载到制定的 NRAM 上进行计算, 1 代表数据直接在 NRAM 上可以直接运算, 保证高效, 但需要满足两个条件: buffer 空间足够, input\_box\_num 是对齐的;
- input\_score\_ptr: 输入框 score 的数据指针;
- input\_x1\_ptr: 输入框 box 的 x1 坐标数据指针;
- input\_y1\_ptr: 输入框 box 的 y1 坐标数据指针;
- input\_x2\_ptr: 输入框 box 的 x2 坐标数据指针;
- input\_y2\_ptr: 输入框 box 的 y2 坐标数据指针;
- x1: buffer 空间, 存放 x1;
- y1: buffer 空间, 存放 y1;
- x2: buffer 空间, 存放 x2;
- y2: buffer 空间, 存放 y2;
- score: buffer 空间, 存放 score;
- inter\_x1: buffer 空间, IOU 筛选临时空间;
- inter\_y1: buffer 空间, IOU 筛选临时空间;
- inter\_x2: buffer 空间, IOU 筛选临时空间;

- `inter_y2`: buffer 空间, IOU 筛选临时空间;
  - `max_box`: buffer 空间, 筛选框的信息储存空间, 顺序为 `max score, x1, y1, x2, y2`;
  - `nram_save`: buffer 空间, 筛选框的临时储存空间;
  - `limit = 0`: 根据 buffer 进行 `findlimit`, 一次最多能处理的输入框的个数;
  - `len_core = 0`: 每个核处理的框的个数;
  - `max_seg_pad = 0`: 每次处理输入框的个数, 根据 `limit` 进行下补齐, 满足硬件限制;
  - `repeat`: 整数段, 需要处理几次 `max_seg_pad`;
  - `remain`: 余数段, 剩下得需要处理的框的个数;
  - `remain_pad`: 余数段进行补齐后的框个数, 满足硬件限制;
  - `input_offset`: 当前核处理的输入数据的起始地址偏移;
  - `nram_save_count`: 临时储存空间储存多少个框后再整体拷贝到实际的目标地址;
- (b) 空间划分 `MODE` 为 0 的时候, `x1, y2, x2, y2, score` 指向的 buffer 空间用来加载数据, `MODE` 为 1 的时候, 不再需要加载数据的空间, `inter_x1, inter_y1, inter_x2, inter_y2` 是 buffer 上的临时空间, 用来计算 NMS, 4 块临时空间是考虑了空间复用后最少需要的数量, 以节省 buffer, 可以一次处理更多的数据。
- (c) 多核拆分 `src==NRAM` 模块内只支持单核模式, 用户可根据需求在模块外按类拆分, 即每个核负责计算一个类的 NMS 过程, 因为模块内无法在 U1 模式时对某一个核上的 `NRAM` 上的数据进行访问; `src==GDRAM/SRAM`, 模块内支持按数据块进行多核拆分, 即一个类的 NMS 过程拆分到多个核上进行计算。多核拆分使用的拆分方式是: 每个核上分到框的数量之间相差不超过 1, 即尽可能的将数据平均分配到多个核上。
2. 执行阶段: 主要操作都放在 `for (int keep = 0; keep < keepNum; keep++)` 循环中, 一次找一个框, 这个循环有两个退出条件: 找到的框数量达到 `keepNum`, 最大的 `score` 小于 `thresh_score`。
- (a) 搜索最大值模块: 通过 `__bang_max` 指令找到 `score` 存放数据的最大值与最大索引, 如果是 `BLOCK` 模式, 则直接即可找到 `score` 的最大值与最大索引, 如果是 `U1` 模式, 则需要每个核上计算所分到数据的最大值, 然后通过 `SRAM` 进行通信, 进一步找到四个核上的最大值与最大值索引。
- (b) 保存模块: 保存方式主要包括两种, 一种是按 `score, x1, y1, x2, y2` 的顺序一组一组进行储存, 一种是按 `score, x1, y1, x2, y2` 的顺序以此存放。两者区别如下图 7.15 所示:
- (c) 根据交并比筛选模块: 以图 7.16 中的代码为例进行说明。

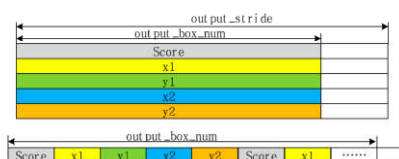


图 7.15 BangC NMS 两种保存方式

```

1  *IOU计算部分*
2  // 得到相交部分的面积:
3  __nramset(inter_y1, seg_len, max_box[1]); // max_x1
4  __svmax_relu(inter_x1, x1, inter_y1, seg_len); // inter_x1
5  __nramset(inter_y2, seg_len, max_box[3]); // max_x2
6  __svmin_relu(inter_x2, x2, inter_y2, seg_len); // inter_x2
7  __bang_sub(inter_x1, inter_x2, inter_x1, seg_len);
8  __bang_active_relu(inter_x1, inter_x1, seg_len); // inter_w
9  __nramset(inter_x2, seg_len, max_box[2]); // max_y1
10 __svmax_relu(inter_y1, y1, inter_x2, seg_len); // inter_y1
11 __nramset(inter_x2, seg_len, max_box[4]); // max_y2
12 __svmin_relu(inter_y2, y2, inter_x2, seg_len); // inter_y2
13 __bang_sub(inter_y1, inter_y2, inter_y1, seg_len);
14 __bang_active_relu(inter_y1, inter_y1, seg_len); // inter_h
15 __bang_mul(inter_x1, inter_x1, inter_y1, seg_len); //area_l
16 // 得到每个框的面积: area = (x2 - x1) * (y2 - y1):
17 __bang_sub(inter_y1, x2, x1, seg_len);
18 __bang_sub(inter_y2, y2, y1, seg_len);
19 __bang_mul(inter_x2, inter_y1, inter_y2, seg_len); // area
20 // 得到相并部分的面积area_U: area + max_area - area_l:
21 __nramset(inter_y1, seg_len, max_area);
22 __bang_add(inter_x2, inter_x2, inter_y1, seg_len);
23 __bang_sub(inter_x2, inter_x2, inter_x1, seg_len); //area_U
24 *筛选*
25 // 如果IOU超过阈值, 则将相应的框的score置0: area_U * thresh > area_l?
26 __bang_mul_const(inter_x2, inter_x2, thresh_iou, seg_len);
27 __bang_gt(inter_x1, inter_x2, inter_x1, seg_len);
28 __bang_mul(score, score, inter_x1, seg_len);

```

图 7.16 NMS 交并比筛选模块

### 7.1.6.3 TensorFlow 算子添加

TensorFlow 中添加算子可分为两种情况, 一种为 TensorFlow 官方定义的标准流程, 可以从官方文档中获取; 另一种为添加 MLU 算子的方式, 本节将对如何添加 MLU YOLOv3 后处理算子作出介绍。

添加 MLU 算子可分为以下几个步骤:

1. MLULib 层, 该部分代码在 tensorflow/stream\_executor/mlu/mlu\_api/lib\_ops/目录下, 主要作用是进一步对 SDK 层进行 API 封装。需要注意的是 MLU LIB 层禁止引用 TensorFlow 的头文件, 目前只使用了 tensorflow::Status 类方便做返回值处理。对于该算子来说主要对该目录下的 mlu\_lib\_ops.cc 和 mlu\_lib\_ops.h 文件进行修改, 示例代码如下所示。

2. MLUOps 层实现, MLUOps 层的职责是使用 LIB 层的 API 完成算子的实现。该部分代码在 tensorflow/stream\_executor/mlu/mlu\_api/ops 下, 注意每个算子首先需要在该目录下的 mlu\_ops.h 文件中添加算子类的声明。下图 7.18 为添加算子类声明的程序。

```

1 tensorflow::Status CreateYolov3DetectionOutputOp(
2     MLUBaseOp** op, MLUTensor** input_tensors, MLUTensor** output_tensors,
3     cnmlPluginYolov3DetectionOutputOpParam_t param) {
4     CNML_RETURN_STATUS(cnmlCreatePluginYolov3DetectionOutputOp(
5         op, param, input_tensors, output_tensors));
6 }
7
8 tensorflow::Status ComputeYolov3DetectionOutputOp(MLUBaseOp* op,
9                                                     MLUCnrtQueue* queue,
10                                                    void* inputs[], int input_num,
11                                                    void* outputs[],
12                                                    int output_num) {
13     int dp = 1;
14     cnrtInvokeFuncParam_t compute_forw_param;
15     u32_t affinity = 0x01;
16     compute_forw_param.data_parallelism = &dp;
17     compute_forw_param.affinity = &affinity;
18     compute_forw_param.end = CNRT_PARAM_END;
19
20     cnmlComputePluginYolov3DetectionOutputOpForward(
21         op, inputs, input_num, outputs, output_num, &compute_forw_param, queue);
22 }

```

(a)

```

1 tensorflow::Status CreateYolov3DetectionOutputOp(
2     MLUBaseOp** op, MLUTensor** input_tensors, MLUTensor** output_tensors,
3     cnmlPluginYolov3DetectionOutputOpParam_t param);
4
5 tensorflow::Status ComputeYolov3DetectionOutputOp(MLUBaseOp* op,
6                                                     MLUCnrtQueue* queue,
7                                                     void* inputs[], int input_num,
8                                                     void* outputs[],
9                                                     int output_num);

```

(b)

图 7.17 MLU.lib 层实现

添加完算子声明后还需同一个目录.cc 文件中添加算子的实现，每个算子都需实现 Create 和 Compute 两个方法。下图 7.19 和 7.20 程序显示了 yolov3detectionoutput.cc 的实现。

Create 方法的职责是将创建好的 baseop 指针调用 base\_ops.\_push\_back(op\_ptr); 储存起来；Compute 方法的职责在于调用 lib 层的 compute 函数进行计算，目前算子实现为同步方式，因此需要调用 SyncQueue。

3. MLUStream 实现, MLUStream 负责 MLUOps 算子类的实例化, 其接口都定义在 tensorflow/stream\_executor/mlu/mlu\_stream.h 中。以 MLUOpKernel 和 MLUStream 层的接口特点, 可以把 MLUStream 层的算子分为 2 类: 通用模版算子与特例化算子。对于 Yolov3DetectionOutput 算子来说, 其属于通用模版算子。

通用模版算子满足以下三个条件: MLU 算子的所有输入均来自 OpKernelContext; MLU 算子的所有输出顺序须与 OpkernelContext 的输出顺序一致; MLUTensor 可以被 CreateMLUTensorFromTensor 创建, 即 MLU Tensor 的形状、数据类型与 TensorFlow tensor 一致。

4. MLUOpKernel 实现, OpKernel 为 TensorFlow 对算子抽象定义。MLUOpKernel 继承了 OpKernel 类, 其使用方法与 OpKernel 基本一致。对于每个 MLU 算子均需要实现其



```

1 struct MLUYolov3DetectionOutputOpParam {
2     int batchNum_;
3     int inputNum_;
4     int classNum_;
5     int maskGroupNum_;
6     int maxBoxNum_;
7     int netw_;
8     int neth_;
9     float confidence_thresh_;
10    float nms_thresh_;
11    int* inputWs_;
12    int* inputHs_;
13    float* biases_;
14    MLUYolov3DetectionOutputOpParam(int batchNum, int inputNum, int classNum,
15                                    int maskGroupNum, int maxBoxNum, int netw,
16                                    int neth, float confidence_thresh,
17                                    float nms_thresh, int* inputWs, int* inputHs,
18                                    float* biases)
19    : batchNum_(batchNum),
20      inputNum_(inputNum),
21      classNum_(classNum),
22      maskGroupNum_(maskGroupNum),
23      maxBoxNum_(maxBoxNum),
24      netw_(netw),
25      neth_(neth),
26      confidence_thresh_(confidence_thresh),
27      nms_thresh_(nms_thresh),
28      inputWs_(inputWs),
29      inputHs_(inputHs),
30      biases_(biases) {}
31
32 };

```

(a)

```

1 DECLARE_OP_CLASS(MLU...);
2 DECLARE_OP_CLASS(MLUYolov3DetectionOutput);
3 DECLARE_OP_CLASS(MLU...);

```

(b)

图 7.18 Yolov3DetectionOutput 算子 MLUOps 声明

构造函数与 ComputeOnMLU 方法。MLUOpKernel 实现的主要功能有参数检查、参数处理、输出形状推断及输出内存分配、调用 MLUStream 层接口完成算子计算。下图 7.22、7.23、7.24 中的代码分别显示了这四个不同的步骤。input0、input1、input2 为 YOLOv3 结构图 7.3 中的三个 feature map 输出。

5. 开始注册 MLU Op Kernel, 在 tensorflow/core/kernels 中添加 yolov3\_detection\_output\_op.cc 注册文件和添加 yolov3\_detection\_output\_op\_mlu.h 算子实现文件, Op Kernel 注册使用的均为 REGISTER\_KERNEL\_BUILDER 宏, 如下图 7.25 注意: 当前只支持注册 MLU 支持的数据类型, 不支持的无法注册; 可以通过 HostMemory 对特定输入进行限制, 如需要在 CPU 上做处理的输入数据。

6. 算子注册在 tensorflow/core/ops 目录下找到对应的 Op 注册文件, 对于 Yolov3DetectionOutput 为 image\_ops.cc 文件。通过这个文件可知 Yolov3DetectionOutput 算子有三个输入和一个输

```

1 Status MLUYolov3DetectionOutput::CreateMLUOp(std::vector<MLUTensor*> &inputs, \
2     std::vector<MLUTensor*> &outputs, void *param) {
3     TF_PARAMS_CHECK(inputs.size() > 0, "Missing input");
4     TF_PARAMS_CHECK(outputs.size() > 0, "Missing output");
5     MLUBaseOp *op_ptr = nullptr;
6     MLUTensor* input0 = inputs.at(0);
7     MLUTensor* input1 = inputs.at(1);
8     MLUTensor* input2 = inputs.at(2);
9     MLUTensor* output = outputs.at(0);
10    MLUTensor* buffer = outputs.at(1);
11
12    MLULOG(3) << "CreateYolov3DetectionOutputOp"
13        << ", input0: " << lib::MLUTensorUtil(input0).DebugString()
14        << ", input1: " << lib::MLUTensorUtil(input1).DebugString()
15        << ", input2: " << lib::MLUTensorUtil(input2).DebugString()
16        << ", output: " << lib::MLUTensorUtil(output).DebugString()
17        << ", buffer: " << lib::MLUTensorUtil(buffer).DebugString();
18    int batchSize = ((ops::MLUYolov3DetectionOutputOpParam*)param)->batchNum_;
19    int inputNum = ((ops::MLUYolov3DetectionOutputOpParam*)param)->inputNum_;
20    int classNum = ((ops::MLUYolov3DetectionOutputOpParam*)param)->classNum_;
21    int maskGroupNum = ((ops::MLUYolov3DetectionOutputOpParam*)param)->maskGroupNum_;
22    int maxBoxNum = ((ops::MLUYolov3DetectionOutputOpParam*)param)->maxBoxNum_;
23    int netw = ((ops::MLUYolov3DetectionOutputOpParam*)param)->netw_;
24    int neth = ((ops::MLUYolov3DetectionOutputOpParam*)param)->neth_;
25    float confidence_thresh = ((ops::MLUYolov3DetectionOutputOpParam*)param)->
26        confidence_thresh_;
27    float nms_thresh = ((ops::MLUYolov3DetectionOutputOpParam*)param)->nms_thresh_;
28    int* inputWs = ((ops::MLUYolov3DetectionOutputOpParam*)param)->inputWs_;
29    int* inputHs = ((ops::MLUYolov3DetectionOutputOpParam*)param)->inputHs_;
30    float* biases = ((ops::MLUYolov3DetectionOutputOpParam*)param)->biases_;
31    cnmlPluginYolov3DetectionOutputOpParam_t mlu_param;
32    const int num_anchors = 3;
33    cnmlCoreVersion_t core_version = CNML_MLU270;
34    cnmlCreatePluginYolov3DetectionOutputOpParam(
35        &mlu_param,
36        batchSize,
37        inputNum,
38        classNum,
39        num_anchors,
40        maxBoxNum,
41        netw,
42        neth,
43        confidence_thresh,
44        nms_thresh,
45        core_version,
46        inputWs,
47        inputHs,
48        biases);
49    std::vector<MLUTensor*> input_tensors = {input0, input1, input2};
50    std::vector<MLUTensor*> output_tensors = {output, buffer};
51    TF_STATUS_CHECK(lib::CreateYolov3DetectionOutputOp(
52        &op_ptr, input_tensors.data(), output_tensors.data(), mlu_param));
53
54    base_ops_.push_back(op_ptr);
55
56    return Status::OK();
57 }

```

图 7.19 Create 函数实现

```

1 Status MLUYolov3DetectionOutput::Compute(const std::vector<void *> &inputs,
2     const std::vector<void *> &outputs, cnrtQueue_t queue) {
3     int num_input = inputs.size();
4     int num_output = outputs.size();
5     assert(num_input == 3);
6     assert(num_output == 2);
7     TF_STATUS_CHECK(lib::ComputeYolov3DetectionOutputOp(
8         base_ops_.at(0), queue,
9         const_cast<void**>(inputs.data()),
10        num_input,
11        const_cast<void**>(outputs.data()),
12        num_output));
13
14     // todo, delete sync queue after delay copy
15     cnrtSyncQueue(queue);
16     return Status::OK();
17 }

```

(a)

图 7.20 Compute 函数实现

出，三个输入分别用 input0、input1、input2 标识，输出用 predicts 标识，输入输出的数据类型均用 T 表示即该算子的输出与输入数据类型一致。Attr("T: type") 表示 T 允许的数据类型为 type，也就是 TensorFlow 支持的所有数据类型，其余为各超参数配置包括 confidence 阈值以及 NMS 阈值等。

以上 6 个步骤为添加 cnplugin 中 Yolov3DetectionOutput 算子的基本流程，除此之外，为了将该算子最终集成在 TensorFlow 中，编译时还需要在 tensorflow/core/kernels/BUILD 中添加以下信息。

#### 7.1.6.4 DLP 执行

在将 BCL 实现的 Yolov3DetectionOutput 算子成功集成到框架后，需要修改之前转换完成的 int8 pb 模型，在模型中增加 Yolov3DetectionOutput 后处理算子，这一步可以使用编程框架机理章节中介绍的 pb 与 ptxt 相互转换工具实现。具体在 pb 中添加如下图 7.28 所示的节点：注意，由于后处理节点已经将三种类型的 feature map 做过一定处理，与 pb 原模型中的操作有一定冗余，为了确保精度的准确性，选择"conv\_lbbox/BiasAdd"，"conv\_mbbox/BiasAdd"，"conv\_sbbox/BiasAdd" 作为后处理算子的输入节点。至此，只需修改推理代码，将后处理部分注释掉即可完成整个 YOLOv3 模型的在线加速推理工作。

#### 7.1.7 实验评估

相比其它两个综合实验，YOLOv3 实验难度中等，因此难度系数分为 1.1。学生最终得分为 1.1 乘以相应的得分。

- 60 分标准：补全 nms\_detection.h 文件，cnplugin 可正常编过。
- 70 分标准：在 60 分的基础上完成 TensorFlow 的集成编译，完成 pb 模型添加后处理大算子的操作，执行测试脚本时 map 值高于 30%，单 batch 延时（包含后处理）低于

```

1  Status Yolov3DetectionOutput(OpKernelContext* ctx ,
2      Tensor* tensor_input0 ,
3      Tensor* tensor_input1 ,
4      Tensor* tensor_input2 ,
5      int batchSize ,
6      int inputNum ,
7      int classNum ,
8      int maskGroupNum ,
9      int maxBoxNum ,
10     int netw ,
11     int neth ,
12     float confidence_thresh ,
13     float nms_thresh ,
14     int* inputWs ,
15     int* inputHs ,
16     float* biases ,
17     Tensor* output1 ,
18     Tensor* output2){
19     ops :: MLUYolov3DetectionOutputOpParam op_param(
20         batchSize ,
21         inputNum ,
22         classNum ,
23         maskGroupNum ,
24         maxBoxNum ,
25         netw ,
26         neth ,
27         confidence_thresh ,
28         nms_thresh ,
29         inputWs ,
30         inputHs ,
31         biases);
32     return CommonOpImpl<ops :: MLUYolov3DetectionOutput>(
33         ctx ,
34         {tensor_input0 , tensor_input1 , tensor_input2} ,
35         {output1 , output2} ,
36         static_cast<void*>(&op_param));
37 }

```

图 7.21 Yolov3DetectionOutput 算子 MLUOpKernel 实现

300ms。

- 80 分标准：在 70 分的基础上，执行测试脚本时 map 值高于 50%，单 batch 延时（包含后处理）低于 100ms。
- 90 分标准：在 80 分的基础上，执行测试脚本时 map 值高于 54%，单 batch 延时（包含后处理）低于 50ms。
- 100 分标准：在 90 分的基础上，执行测试脚本时 map 值高于 56%，单 batch 延时（包含后处理）低于 25ms。

### 7.1.8 实验思考

1.NMS BCL 相比 CPU 实现有何优势？ 2. 有哪些方法可以提升 BCL 算子的性能？

```

1 namespace tensorflow {
2 template<typename T>
3 class MLUYolov3DetectionOutputOp: public MLUOpKernel{
4     public:
5         explicit MLUYolov3DetectionOutputOp(OpKernelConstruction* context):MLUOpKernel(
6             context){
7             OP_REQUIRES_OK(context, context->GetAttr("batchNum",&batchNum_));
8             OP_REQUIRES_OK(context, context->GetAttr("inputNum",&inputNum_));
9             OP_REQUIRES_OK(context, context->GetAttr("classNum",&classNum_));
10            OP_REQUIRES_OK(context, context->GetAttr("maskGroupNum",&maskGroupNum_));
11            OP_REQUIRES_OK(context, context->GetAttr("maxBoxNum",&maxBoxNum_));
12            OP_REQUIRES_OK(context, context->GetAttr("netw",&netw_));
13            OP_REQUIRES_OK(context, context->GetAttr("neth",&neth_));
14            OP_REQUIRES_OK(context, context->GetAttr("confidence_thresh",&
15                confidence_thresh_));
16            OP_REQUIRES_OK(context, context->GetAttr("nms_thresh",&nms_thresh_));
17            OP_REQUIRES_OK(context, context->GetAttr("inputWs",&inputWs_));
18            OP_REQUIRES_OK(context, context->GetAttr("inputHs",&inputHs_));
19            OP_REQUIRES_OK(context, context->GetAttr("biases",&biases_));
20        }
21
22        void ComputeOnMLU(OpKernelContext* context) override {
23            auto* stream = context->op_device_context()->mlu_stream();
24            auto* mstream_exec =
25                context->op_device_context()->mlu_stream()->parent();
26            // 参数检查与处理
27            Tensor* input0 = const_cast<Tensor*>(&context->input(0));
28            Tensor* input1 = const_cast<Tensor*>(&context->input(1));
29            Tensor* input2 = const_cast<Tensor*>(&context->input(2));
30            string op_parameter = context->op_kernel().type_string();
31            MLU_OP_CHECK_UNSUPPORTED(mstream_exec, op_parameter, context);
32            int c_arr_data[3] = {255, 255, 255};
33            std::vector<int> input0_shape(4, 1);
34            std::vector<int> input1_shape(4, 1);
35            std::vector<int> input2_shape(4, 1);
36            input0_shape[0] = batchNum_;
37            input0_shape[1] = c_arr_data[0];
38            input0_shape[2] = inputHs_[0];
39            input0_shape[3] = inputWs_[0];
40
41            input1_shape[0] = batchNum_;
42            input1_shape[1] = c_arr_data[1];
43            input1_shape[2] = inputHs_[1];
44            input1_shape[3] = inputWs_[1];
45
46            input2_shape[0] = batchNum_;
47            input2_shape[1] = c_arr_data[2];
48            input2_shape[2] = inputHs_[2];
49            input2_shape[3] = inputWs_[2];

```

图 7.22 Yolov3DetectionOutput 算子 MLUOpKernel 实现

```

1 // 输出形状推断及输出内存分配
2     int buffer_size = 255 * (inputHs_[0] * inputWs_[0] +
3                             inputHs_[1] * inputWs_[1] +
4                             inputHs_[2] * inputWs_[2]);
5     std::vector<int> buffer_shape = {batchNum_, buffer_size, 1, 1};
6     std::vector<int> output_shape(4, 1);
7     output_shape[0] = batchNum_;
8     output_shape[1] = 7 * maxBoxNum_ + 64;
9     Tensor* output;
10    Tensor* buffer;
11    TensorShape tf_output_shape {output_shape[0], output_shape[1], output_shape[2],
output_shape[3]};
12    TensorShape tf_buffer_shape {buffer_shape[0], buffer_shape[1], buffer_shape[2],
buffer_shape[3]};
13    OP_REQUIRES_OK(context, context->allocate_output(0, tf_output_shape, &output))
;
14    OP_REQUIRES_OK(context, context->allocate_output(0, tf_buffer_shape, &buffer))
;

```

图 7.23 Yolov3DetectionOutput 算子 MLUOpKernel 实现

```

1 // 调用MLUstream层接口完成算子计算
2     if (output->NumElements() > 0 && buffer->NumElements() > 0 ) {
3         OP_REQUIRES_OK(
4             context,
5             stream->Yolov3DetectionOutput(
6                 context, input0, input1, input2, batchNum_, inputNum_,
7                 classNum_, maskGroupNum_, maxBoxNum_, netw_, neth_,
8                 confidence_thresh_, nms_thresh_, inputWs_.data(), inputHs_.data(),
9                 biases_.data(), output, buffer));
10        } else {
11            mlustream_exec->insert_unsupported_op(context, op_parameter);
12        }
13    }
14    private:
15    int batchNum_;
16    .....
17 };
18 }

```

图 7.24 Yolov3DetectionOutput 算子 MLUOpKernel 实现

```

1 #if CAMBRICON_MLU
2 #include "tensorflow/core/kernels/yolov3_detection_output_op_mlu.h"
3 namespace tensorflow {
4 #define REGISTER_MLU(T) \
5     REGISTER_KERNEL_BUILDER( \
6         Name("Yolov3DetectionOutput") \
7         .Device(DEVICE_MLU) \
8         .TypeConstraint<T>("T"), \
9         MLUYolov3DetectionOutputOp<T>);
10 TF_CALL_MLU_FLOAT_TYPES(REGISTER_MLU);
11 #undef REGISTER_MLU
12 #endif // CAMBRICON_MLU
13 }

```

图 7.25 Yolov3DetectionOutput 算子注册

```

1 REGISTER_OP("Yolov3DetectionOutput")
2   .Output("predicts: T")
3   .Input("input0: T")
4   .Input("input1: T")
5   .Input("input2: T")
6   .Attr("batchNum: int")
7   .Attr("inputNum: int")
8   .Attr("classNum: int")
9   .Attr("maskGroupNum: int")
10  .Attr("maxBoxNum: int")
11  .Attr("netw: int")
12  .Attr("neth: int")
13  .Attr("confidence_thresh: float")
14  .Attr("nms_thresh: float")
15  .Attr("inputWs: list(int) = [13, 26,52]")
16  .Attr("inputHs: list(int) = [13, 26,52]")
17  .Attr("biases: list(float) = [116, 90, 156, 198, 373, 326, 30, 61, 62, 45, 59, 119,
18     10, 13, 16, 30, 33, 23]")
19  .Attr("T: type")
20  .SetShapeFn([](InferenceContext *c){
21      return SetOutputForYolov3DetectionOutput(c);
22  });

```

图 7.26 Yolov3DetectionOutput 算子注册

```

1 .....
2 ce_library(
3   name = "image",
4   deps = [
5     ":adjust_contrast_op",
6     .....
7     ":yolov3_detection_output_op",
8   ],
9 )
10 .....
11 tf_kernel_library(
12   name = "yolov3_detection_output_op",
13   prefix = "yolov3_detection_output_op",
14   deps = [
15     "//tensorflow/core:core_cpu",
16     "//tensorflow/core:framework",
17     "//tensorflow/core:lib",
18     "//tensorflow/core:lib_internal",
19     "//third_party/eigen3",
20     ] + if_mlu([
21     "//tensorflow/stream_executor:mlu_stream_executor"]);
22 )

```

图 7.27 TensorFlow 编译

```
1 node {
2   name: "Yolov3DetectionOutput"
3   op: "Yolov3DetectionOutput"
4   input: "conv_lbbox/BiasAdd"
5   input: "conv_mbbox/BiasAdd"
6   input: "conv_sbbox/BiasAdd"
7   attr {
8     key: "T"
9     value {
10      type: DT_FLOAT
11    }
12  }
13  ...
14  attr {
15    key: "nms_thresh"
16    value {
17      f: 0.449999988079
18    }
19  }
20 }
21 versions {
22   producer: 38
23 }
```

图 7.28 YOLOv3 增加后处理节点



```
1 def predict_bak(self, images):
2
3     org_h = [0 for i in range(self.batch_size)]
4     org_w = [0 for i in range(self.batch_size)]
5     for i in range(self.batch_size):
6         org_h[i], org_w[i], _ = images[i].shape
7
8     image_data = utils.images_preporcess(images, [self.input_size, self.input_size])
9
10    start = time.time()
11    pred_sbbox, pred_mbbox, pred_lbbox = self.sess.run(
12        [self.pred_sbbox, self.pred_mbbox, self.pred_lbbox],
13        feed_dict={
14            self.input_data: image_data,
15            #self.trainable: False
16        }
17    )
18    np.savetxt("pred_sbbox.txt", pred_sbbox.flatten())
19    np.savetxt("pred_mbbox.txt", pred_mbbox.flatten())
20    np.savetxt("pred_lbbox.txt", pred_lbbox.flatten())
21    end = time.time()
22
23    batch_bboxes = []
24    for idx in range(self.batch_size):
25        pred_bbox = np.concatenate([np.reshape(pred_sbbox[idx], (-1, 5 + self.
26            num_classes)),
27                                   np.reshape(pred_mbbox[idx], (-1, 5 + self.
28            num_classes)),
29                                   np.reshape(pred_lbbox[idx], (-1, 5 + self.
30            num_classes))], axis=0)
31        bboxes = utils.postprocess_boxes(pred_bbox, (org_h[idx], org_w[idx]), self.
32            input_size, self.score_threshold)
33        batch_bboxes.append(utils.nms(bboxes, self.iou_threshold))
34        print("bbox num : ", len(batch_bboxes))
35        for i in range(len(batch_bboxes)):
36            print(batch_bboxes[i])
37    exit(0)
38    return batch_bboxes, (end - start)
```

图 7.29 原始推理

```

1  def predict(self, images):
2      .....
3      start = time.time()
4      bbox_raw = self.sess.run(
5          self.bbox_raw,
6          feed_dict={
7              self.input_data: image_data,
8              #self.trainable: False
9          }
10     )
11     end = time.time()
12     batch_bboxes = []
13     num_batches = 1
14     num_boxes = 1024 * 2
15     predicts_mlu = bbox_raw.flatten()
16     print("org_h[0], ", org_h[0])
17     print("org_w[0], ", org_w[0])
18     for batchIdx in range(num_batches):
19         result_boxes = int(predicts_mlu[batchIdx * (64 + num_boxes * 7)])
20         current_bboxes = []
21         print("result_boxes :", result_boxes)
22         print("x1, y1, x2, y2, score, classId")
23         for i in range(result_boxes):
24             # batchId, classId, score, x1, y1, x2, y2
25             batchId = predicts_mlu[i * 7 + 0 + 64 + batchIdx * (64 + num_boxes * 7)]
26             classId = predicts_mlu[i * 7 + 1 + 64 + batchIdx * (64 + num_boxes * 7)]
27             score = predicts_mlu[i * 7 + 2 + 64 + batchIdx * (64 + num_boxes * 7)]
28             x1 = predicts_mlu[i * 7 + 3 + 64 + batchIdx * (64 + num_boxes * 7)]
29             y1 = predicts_mlu[i * 7 + 4 + 64 + batchIdx * (64 + num_boxes * 7)]
30             x2 = predicts_mlu[i * 7 + 5 + 64 + batchIdx * (64 + num_boxes * 7)]
31             y2 = predicts_mlu[i * 7 + 6 + 64 + batchIdx * (64 + num_boxes * 7)]
32             print(x1, y1, x2, y2, score, classId)
33             # bbox = [xmin, ymin, xmax, ymax, score, class]
34             bbox = [x1, y1, x2, y2, score, classId]
35             current_bboxes.append(np.array(bbox))
36             batch_bboxes.append(current_bboxes)
37         print("bbox num : ", len(batch_bboxes))
38         for i in range(len(batch_bboxes)):

```

图 7.30 添加后处理大算子

## 7.2 文本识别 OCR-EAST

### 7.2.1 实验目的

本实验主要完成光学字符识别的代表性算法——EAST 网络移植到智能处理器 DLP 上, 并进行性能优化与离线部署, 使读者可以借助 DLP 处理器完成完整的学字符识别任务。

因此, 本实验的主要目的包括:

1. 通过使用 Tensorflow 在 CPU 和智能处理器 DLP 上进行在线推理, 掌握使用 Tensorflow 编写光学字符识别应用;
2. 通过用智能编程语言实现 EAST 网络中的 Split+sub+concat 合并算子, 深入掌握智能编程语言的算子开发和优化方法;
3. 通过在 Tensorflow 中添加合并算子, 代替单算子, 掌握在 Tensorflow 框架中添加合并算子的方法并优化性能;

实验时间预计为两周。

### 7.2.2 背景介绍

#### 7.2.2.1 OCR



图 7.31 OCR 应用举例

从传说中的仓颉造字以来, 文字或文本就为人类的发展的进步带来了重要贡献。到现在为止, 文本依旧是包含丰富而准确信息的一个大类场景, 因此自然场景中的文本检测和识别已成为计算机视觉和文档分析中重要且活跃的研究主题。尤其是近年来, 尽管仍然存在各种挑战 (例如, 噪声, 模糊, 失真, 遮挡和变化), 但社区在这些领域的研究工作激增并取得了实质性进展。

在实际应用中, 在新闻出版, 邮政快递, 金融服务等领域都广泛使用了文本识别算法。例如对于如图 7.31 所示的针对会计领域的发票文本识别就大大减轻了人工比对的负担。

光学字符识别 (Optical Character Recognition, OCR) 是指对文本资料的图像文本进行文字识别, 获取文本资料版面信息的过程。OCR 算法的研究有很长的历史, 在近一些年基于各种深度学习的算法逐渐成为重要的研究方向。

OCR 过程大致可以分为: 图像获取-> 图像预处理-> 文字检测-> 文字识别-> 输出。

### 7.2.2.2 典型算法

文本检测和识别旨在检测和定位场景图像或视频中的文本，并自动对其进行识别。目前学术界已经开发了许多技术来检测和识别场景图像和视频中的文本。这些技术大致可以分为三大类：文本检测与定位，文本识别和端到端文本识别系统。文本检测与定位的目的是确定给定图像或视频中是否存在文本并将其定位。文本识别旨在将图像中定位出的文本转换为字符编码。而端到端文本识别系统将检测和识别结合为一个完整的框架。

#### • 文本检测

现有的基于 CNN 的方法可以大致分为几类：基于 region proposal 的方法，基于 segmentation 的方法以及使用多任务学习的混合方法。

基于 region proposal 的方法一般都利用了基于 region 的方法和深度神经网络的优势。Zhang 等人<sup>[26]</sup>提出了一种可以检测多种方向，语言和字体文本的框架。它主要由两个完全卷积网络（FCN）组成：一个用于文本区域的 salient map，另一个 FCN 用于预测每个字符的中心。但是，在某些情况下可能会出现误报和字符丢失的情况，例如对比度极低，曲率大，反射光强，文本行太封闭或字符之间的巨大间隙。另一个限制是计算效率低。

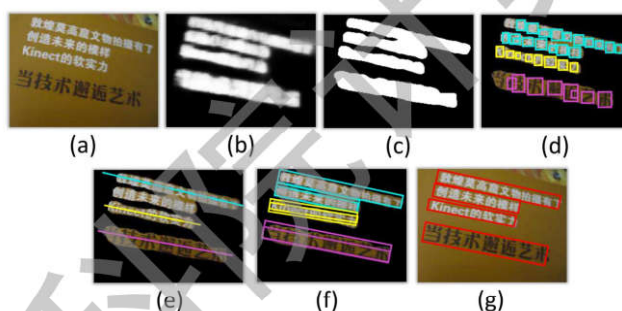


图 7.32 典型文本检测算法流程

如图7.32所示，这个算法主要分为以下几个步骤：

- (a) 输入图像；
- (b) TextBlock FCN 预测的文本区域的 salient map；
- (c) 文本块生成；
- (d) 候选字符成分提取；
- (e) 按组成部分预测的方向估计；
- (f) 提取候选文本行；
- (g) 算法的检测结果。

#### • 文本识别

对于 OCR 的文本识别，基于图像的序列识别问题，Shi<sup>[27]</sup> 等人提出了一种端到端的文本识别系统，称为卷积递归神经网络（CRNN），它由卷积层，递归层和转录组成。通过将 CNN 与 RNN 结合使用，CRNN 可以处理任意长度的序列，并且不依赖于其他词典。如图 7.33 所

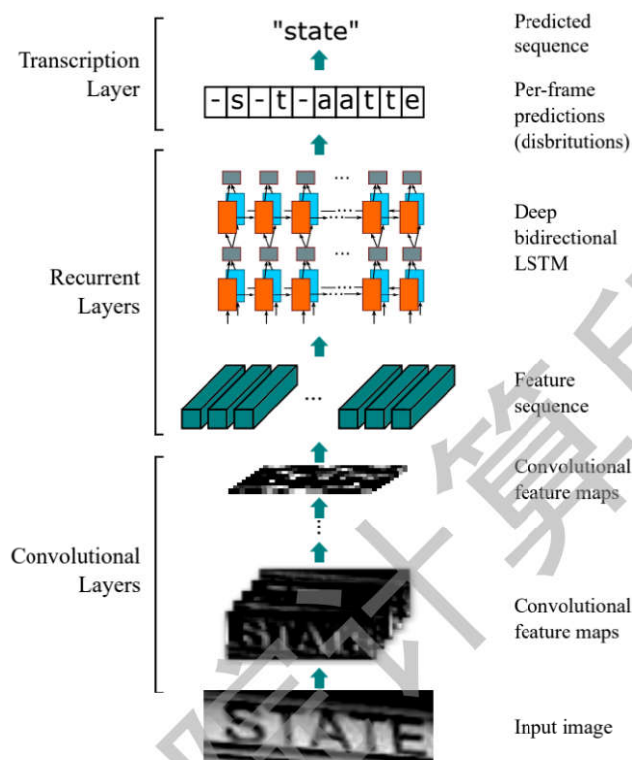


图 7.33 典型文本识别网络结构

示网络结构主要包括三个部分：

- 卷积层，从输入图像中提取特征序列；
- 循环层，预测每个帧的标签分布；
- 转录层，它将每帧的预测翻译成最终的标记序列。

#### • 端到端方法

EAST 算法全名是 Efficient and Accuracy Scene Text, 它提出了一个快速而准确的轻量级场景文本检测 pipeline, 解决大多数 OCR 算法准确性和效率缺陷. 该 Pipeline 只有两个阶段: 第一阶段是基于 FCN 模型的文本框检测; 第二阶段是对生成的文本框 (旋转或矩形) 经过非极大值抑制得到最终结果, 从而免去其他中间步骤, 实现端到端训练.

#### • 网络结构

EAST 的网络结构分为三大部分, 如图 7.34 所示:

### 第一部分:Feature extractor stem(PVANet)

该部分是主干特征检测网络,引入多尺度检测的思想,提取不同尺寸卷积核下的特征并用于后期的特征组合,以适应文本行尺度的变化.作者采用了 PVANet 模型作为主干网络,实际使用的时候也可以采用 VGG16 或者 Resnet;

### 第二部分:Feature-merging branch

该部分运用了 Unet 的思想,主要是合并第一部分提取的特征,通过池化和 Concat 来恢复尺寸.

其中,特征合并主要通过逐层合并的方式,首先将第一部分对应的特征图输入到一个池化层恢复尺度,然后与当前层特征图进行 Concat,再通过  $1 \times 1$  卷积来减少通道数,最后利用  $3 \times 3$  卷积将局部信息融合以最终产生该合并阶段的输出;

### 第三部分:Output Layer

该部分输出分为 3 类:score\_map,RBOX geometry,QUAD geometry.

Score\_map 代表此处是否有文字的可能性;

RBOX geometry 中的五个通道,分别代表每个像素点到文本框边线的距离,加上一个文本框的旋转角;

QUAD geometry 中的八个通道,分别代表着每个像素点到文本框任意四边形的 xy 距离.

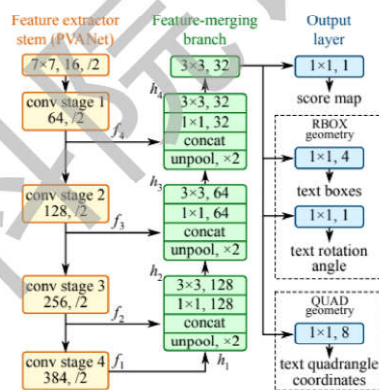


图 7.34 east 网络结构

- 后处理——局部感知 NMS

EAST 的后处理与通用目标检测类似,需要将网络的结果经过非极大值抑制 (NMS) 来取得最终结果.由于 EAST 的网络输出是成千上万个几何文本框,需要根据阈值过滤,然后将文本框的置信度得分加权平均,得到合并后的文本框坐标.

将经过合并后的文本框经过通用的 NMS 处理,详情可参考 7.1 章节,得到最终结果.

### 7.2.3 实验环境

本实验所涉及的硬件平台和软件环境如下:

- 硬件平台: 硬件平台使用第??节介绍的云平台中的 DLP 硬件.
- 软件环境: 所涉及的 DLP 软件开发模块包括编程框架 TensorFlow、高性能库 CNML、运行时库 CNRT、编程语言及编译器、运行时库 CNRT 以及流式编程框架 CNStream.

## 7.2.4 实验内容

本实验主要是在 DLP 硬件上实现并优化以 EAST 网络模型为核心的目标检测应用. 针对用户从网络上下载的标准 EAST 模型, 经过模型量化并配置 DLP 相关参数, 则可以采用 DLP 上的定制 TensorFlow 版本来运行. 为了充分发挥 DLP 的计算能力, 进一步采用智能编程语言 BCL 实现 EAST 网络模型中的 Split+Sub+Concat 合并算子, 并将其添加至 TensorFlow 框架中, 代替原有的 Split、Sub、Concat 单算子. 考虑推理应用的实际需求, 采用 DLP 的流式框架编写完整目标检测应用, 基于离线模型进行部署. 具体实验内容包括:

1. EAST CPU 运行: 利用标准的与训练 CKPT 模型, 通过定制的 TensorFlow 版本在 CPU 上运行 EAST, 得到 Float16 精度的 PB 模型;
2. EAST DLP 运行: 将步骤 1 得到的 Float16 精度的 PB 模型进行 INT8 量化后, 通过定制的 TensorFlow 版本在 DLP 硬件上运行;
3. Split+Sub+Concat 合并算子的 BCL 实现: 采用 BCL 实现的 Split+Sub+Concat 合并运算;
4. 框架算子集成: 将用 BCL 实现的 Split+Sub+Concat 合并算子集成到 TensorFlow 框架中。

## 7.2.5 实验步骤

如前所述, 详细的实验步骤主要包括: EAST CPU 运行、EAST DLP 运行、Split+Sub+Concat 合并算子的 BCL 实现、框架算子集成等。

### 7.2.5.1 EAST CPU 运行

为了使得标准的 EAST 网络模型在 DLP 上运行, 需要通过运行 CPU 实现, 将 CKPT 预训练模型转换成 Float16 精度的 PB 模型. 具体步骤如下:

1. 获取开源数据集和 CKPT 模型:<https://github.com/argman/EAST>
2. 按照以下目录结构配置数据集和模型:

数据集目录结构:

```

1 /home/Cambricon-MLU270/datasets/tensorflow_models/east
2     |-- 2015Challenge4_Test_Task1_GT/
3     |-- 2015ch4_test_images/
4     |-- 2015ch4_training_images/
5     |-- 2015ch4_training_localization_transcription_gt/

```

图 7.35 EAST 数据集目录结构

```

1 /home/Cambricon-MLU270/models/tensorflow_models/east/east_icdar2015_resnet_v1_50_rbox
2     |-- checkpoint
3     |-- model.ckpt-49491.data-00000-of-00001
4     |-- model.ckpt-49491.index
5     |-- model.ckpt-49491.meta

```

图 7.36 EAST 模型目录结构

模型目录结构:

CPU 运行之前,需要 MODEL\_PATH 指定到 east\_icdar2015\_resnet\_v1\_50\_rbox 下,DATASET\_PATH 指定到 2015ch4\_test\_images 下,GT\_PATH 指定到 2015Challenge4\_Test\_Task1\_GT 下。

3. DLP 软件环境中已经适配了 EAST CPU 运行环境和 DLP 运行环境,需要进入 tensorflow\_models/cambricon\_examples/EAST/ 目录下;
4. 安装推理所需软件包:一些必备的 Python 安装包在所提供的云平台软件环境中已经安装完成,但对于每个综合实验的特有依赖还需进行单独安装,使用 pip 命令完成依赖的安装与升级: pip install -r requirements.txt
5. 模型推理和转换:执行如下命令完成在 CPU 上的推理,并将 CKPT 模型转换为 Float16 精度的 PB 模型: ./run\_cpu.sh
6. 整理模型路径:生成的 PB 模型在 ./cpu\_pb 中,需要将得到的 east.pb 拷贝至 /home/Cambricon-MLU270/models/tensorflow\_models/EAST/east.pb

### 7.2.5.2 EAST DLP 运行

下面介绍如何将 7.2.5.1 小节得到的 Float16 精度 PB 模型量化为 INT8 精度的 PB 模型。DLP 软件环境提供 fppb\_to\_intpb 工具来进行模型量化。

1. 进入 tensorflow/cambricon\_examples/tools/fppb\_to\_intpb 目录;
2. 生成数据文件:按照模型所需要的数据集以及数据集的具体路径来修改 generate\_image\_list.py 文件中的 image\_libs\_map,具体如图 7.37 所示,来生成相应的数据文件。
3. 生成量化配置文件:如第 ?? 所述,在进行量化前需设置相应的配置文件。以指定量化相关的参数。DLP 软件环境提供了 generate\_ini.py 来生成指定模型的量化配置文件。针对 EAST,需要将 generate\_ini.py 中的 model\_name 进行修改,即指定要生成的配置文件是针对 EAST 的: models\_name = [ 'EAST' ]

执行 DLP 定制化的 PB 模型量化脚本 generate\_ini.py 后会生成 config 文件夹,其中包含了 EAST\_naive\_int8.ini 参数配置文件,该配置文件包含如图 7.38 所示的信息。同时为了规范模型路径,需要修改 save\_model\_path;

而提高数据集整体的量化质量。但是对于不存在异常值,数据分布紧凑的情况下,不建议使用该算法,比如网络权值的量化。其中 device\_mode 可以设置输出 pb 所有节点的 device,有三种配置方式: clean、mlu 和 origin。其中,mlu 将输出 pb 的所有节点的 device



```

1 import os
2 from os.path import isfile, join
3
4 max_image_count = 100
5 image_libs_map = {
6     "east": join(os.environ.get("TENSORFLOW_MODELS_DATA_HOME"), "east/2015
7         ch4_training_images")
8     }
9 if __name__ == "__main__":
10     for image_lib, image_path in image_libs_map.items():
11         n = 0
12         with open("image_list_{}".format(image_lib), "w") as image_list:
13             for root, dirs, files in os.walk(image_path):
14                 for f in files:
15                     if n >= max_image_count:
16                         break
17                     if f.endswith(".jpg"):
18                         image_list.write(join(image_path, f) + "\n")
19                         n = n + 1

```

图 7.37 EAST 量化生成数据文件

都设置为 MLU, 即都在 MLU 上运行;clean 将输出 pb 的所有节点的 device 清除, 运行时根据算子注册情况自动选择可运行的设备;origin 则使用和输入 pb 一样的设备指定.

4. 执行命令:python fppb\_to\_intpb.py 完成模型量化最终得到 east\_int8.pb.

完成 CPU 侧的推理后, 接下来需要将其移植在 DLP 上, 以加快其推理速度. 此时我们只需仿照第 7 章编程语言算子实验中的 DLP 推理移植部分.DLP 软件栈中已经完成了 EAST 的移植, 只需要将 MODEL\_PATH 指定到 east\_int8.pb, 执行得到精度和性能数据:

```
./run.sh 16 MLU270 int8 100 1
```

### 7.2.5.3 BANGC 代码实现

#### 1. 设计原则

在使用 BCL 设计 Split+Sub+Concat 合并算子的时候, 需要考虑以下 BCL 设计原则:

- 能够处理的数据来源:GDRAM,NRAM,SRAM 中的一种
- 能够将数据存放到:GDRAM,NRAM,SRAM 中的一种
- 每个 Core 都有一块片上存储空间 NRAM, 其大小约为 512KB
- 每个 Cluster 上分配一块 SRAM, 其大小 2MB, 并且单核任务不能 持 SRAM 的使, UNION1 任务和单核任务不 持 SRAM 通信 UNION2 以上任务才 持 SRAM 通信

保证算子正确性的情况下, 需要考虑以下四种优化方向:

- 对访存进 优化以充分利 上存储, 比如 NRAM 复用
- 对计算逻辑进 优化以充分削减计算量, 充分利 张量运算器

```

1 [preprocess]
2 mean = 0, 0, 0 #均值,顺序依次为 mean_r、mean_g、mean_b
3 std = 1.0 #方差
4 color_mode = rgb #网络的输入图片是 rgb 还是 bgr
5 crop = 544, 544 #前处理最终将图片处理为 416 x 416 大小
6 calibration = east_preprocess_cali #校准数据读取及前处理的方式,可以根据需求进行自定义, [preprocess] 和 [data] 中定义的参数均为 calibration 的输入参数
7
8 [config]
9 activation_quantization_alg = naive #输入量化模式,可选 naive 和 threshold_search, naive 为基础模式, threshold_search 为阈值搜索模式
10
11 device_mode = clean #可选 clean、mlu 和 origin,建议使用 clean,使用 clean 生成的模型在运行时会自动选择可运行的设备
12
13 use_convfirst = False #是否使用 convfirst
14 quantization_type = int8 #量化位宽,目前可选 int8 和 int16
15 debug = False #是否为 debug 模式
16 weight_quantization_alg = naive #权值量化模式,可选 naive 和 threshold_search, naive 为基础模式, threshold_search 为阈值搜索模式
17
18 int_op_list = Conv, FC, LRN #要量化的 layer 的类型,目前只能量化 Conv、FC 和 LRN
19 channel_quantization = False #是否使用分通道量化,目前不支持为 True
20
21 [model]
22 output_tensor_names = feature_fusion/Conv_7/Sigmoid:0, feature_fusion/concat_3:0 #输出 Tensor 的名字,可以是多个,以逗号隔开
23 original_models_path = /home/Cambricon-MLU270/models/tensorflow_models/east/east.pb #输入 pb
24 save_model_path = /home/Cambricon-MLU270/models/tensorflow_models/east/east_int8.pb #输出 pb
25 input_tensor_names = input_images:0 #输入 Tensor 的名字,可以是多个,以逗号隔开
26
27 [data]
28 num_runs = 2 #运行次数

```

图 7.38 EAST 量化配置文件

- (c) 充分利用多核并行（计算任务拆分）来提升并行度以及隐藏访存延迟
- (d) 充分利用编译器提供的各种编译优化选项,比如 O2 或者 O3 编译优化

## 2. 设计原理

首先,我们需要分析输入数据的规模,通过开源模型可视化软件 Netron 查看 PB 模型,可以确定整个网络模型的结构,并在 DLP MLU 执行过程中,打印出 input\_images 的输出维度为 (1,672,1280,3),即是 Split 的输入维度;

由于 BCL 使用的数据类型是 Half 类型,占据 2 个字节,所以计算出输入规模是  $1 \times 672 \times 1280 \times 3 \times 2 / 1024$ . 在充分利用多核并行的情况下,将数据分块到每个 core,每块 NRAM 为  $5040 / 16 = 315 \text{KB}$ ,在 512KB 的限制内;

使用 BangC 实现 Split 的时候,需要注意,输入数据的摆数是 NHWC,C 方向的数据是连续的. 如果在 C 方向进行 Split,则需要先转置为 NCHW,或者使用 \_\_bang\_\_maskmove 来进行 Split,再去做 Sub,这样会增加计算耗时或者 NRAM 的使用. 而最简洁的方法是

使用 `__bang_cyclesub`, 对 C 方向进行 `cyclesub`, 省略掉了 `Split` 和 `Concat` 部分, 大大提升效率, 减少 NRAM 的使用. 整体设计流程图如图 7.39 所示:



图 7.39 SBC 运行流程

下面以图7.40中的代码为例子进行说明, 介绍 `Split+Sub+Concat` 合并算子的接口介绍:

```

1  __mlu_entry__ void SBCKernel(half* input_data_,
2  half* output_data_,
3  int batch_num_)
4

```

图 7.40 `Split+Sub+Concat` 合并算子接口

`input_data_`: 输入数据存放的数据地址;

`output_data_`: 输出数存放的数据地址;

`batch_num_`: Batchsize 的规模.

### 3. 实现解析

代码实现分为两部分: Kernel 函数实现的 `.mlu` 和 `cnrt` 实现的 `.cc`

#### 7.2.5.4 Tensorflow 算子添加

本节介绍如何将实现的 BANGC kernel 函数集成至 Tensorflow, 为 EAST 添加 `Split+Sub+Concat` 合并算子. 添加 MLU 算子主要分为以下几个步骤, 步骤之间有承接关系:

##### 1. CNPlugin 编译

如前所述, CNML 通过 `PluginOp` 相关接口提供了用户自定义算子和高性能库已有算子协同工作 (如算子融合) 的机制. 因此, 在完成 `Split+Sub+Concat` 算子的开发后, 可以利用 `PluginOp` 相关接口封装出方便用户使用的接口 (主要包括 `PluginOp` 的创建、计算和销毁等接口), 使得用户自定义算子和高性能库已有算子有一致的编程接口和模式.

图7.43给出了 `PluginOp` 接口封装的部分示例代码, 主要包括 `Create`、`Compute` 和 `Destroy` 三类函数的具体实现.

-creat 函数:

```

1  #define HWC_SPLIT (((HEIGHT*WIDTH/16) - 1) / ALIGN_SIZE + 1) * ALIGN_SIZE*CHANNELS
2  #define CHANNELS 3
3  #define HEIGHT 672
4  #define WIDTH 1280
5  #define ALIGN_SIZE 64
6  #define DATA_COUNT ((CHANNELS) * (WIDTH) * (HEIGHT))
7
8  #include "mlu.h"
9
10 __mlu_entry__ void SBCKernel(half* input_data_, half* output_data_, int batch_num){
11
12     int batch_num = batch_num_;
13
14     // struct timeval start;
15     // struct timeval end;
16     // gettimeofday(&start, NULL);
17
18     __nram__ half split_sub_concat[HWC_SPLIT];
19     __nram__ half tmp0[192];
20
21     // 多核拆分
22     int core_loop = 16 / taskDim;
23
24     // 循环创建 cycle_sub mask
25     for (int i = 0; i < 192; i++){
26         tmp0[i*3] = 123.68;
27         tmp0[i*3+1] = 116.78;
28         tmp0[i*3+2] = 103.94;
29     }
30
31     for (int i = 0; i < batch_num; i++){
32         for (int j = 0; j < core_loop; j++){
33             // 数据拆分至每个 Core 的 NRAM
34             __memcpy(split_sub_concat, input_data_ + (j * taskDim + taskId) *
HWC_SPLIT + i * DATA_COUNT, HWC_SPLIT * sizeof(half), GDRAM2NRAM);
35             // cycle_sub 代替 split+sub
36             __bang_cycle_sub(split_sub_concat, split_sub_concat, tmp0, HWC_SPLIT, 192);
37             // 按顺序拷贝回 GDRAM
38             __memcpy(output_data_ + (j * taskDim + taskId) * HWC_SPLIT + i *
DATA_COUNT, split_sub_concat, HWC_SPLIT * sizeof(half), NRAM2GDRAM);
39         }
40         __sync_all();
41     }
42
43     // 计算耗时
44     // gettimeofday(&end, NULL);
45     // uint32_t time_usec = (uint32_t)end.tv_usec - (uint32_t)start.tv_usec;
46     // printf("Hardware Total Time: %u us\n", time_usec);
47     // printf("batch_size: %d us\n", batch_num);
48     // printf("core_num: %d us\n", taskDim);
49
50 }

```

图 7.41 EAST BANGC 算子实现

-Compute 函数:

-Destroy 函数:

## 2. 算子注册

算子注册在 `tensorflow/core/ops` 目录下找到对应的 Op 注册文件, 对于 Split+Sub+Concat 合并算子, 在 `math_ops.cc` 中注册. 如图7.45, Split+Sub+Concat 合并算子包含一个输入节点 `input` 和一个输出节点 `output`, 输入输出的数据类型均用 `T` 表示即该算子的输出与输入数据类型一致. 其中, `Attr(“T:type”)` 表示 `T` 允许的数据类型为 `type`, 也就是 TensorFlow 支持的所有数据类型.

### 3. MLULib

MLULib 层主要是对 CNML 和 CNRT 接口的封装, 该部分代码在 `tensorflow/stream_executor/mlu/mlu_` 目录下, 需要注意的是 MLULIB 层禁止引用 TensorFlow 的头文件, 目前只使用了 `tensorflow::Status` 类方便做返回值处理. 添加 Split+Sub+Concat 合并算子需要对该目录下的 `mlu_lib_ops.cc` 和 `mlu_lib_ops.h` 文件进行修改, 示例代码如图7.46所示.

### 4. MLUOps

MLUOps 层主要是使用 MLULib 层封装好的 API 完成算子的实现, 代码文件在 `tensorflow/stream_executor/mlu/mlu_api/ops` 下, 需要在目录下的 `mlu_ops.h` 文件中添加新算子类的声明, 如图7.47所示.

同时, 需要新建一个 `sbc.cc` 文件来添加 Split+Sub+Concat 合并算子的实现, 其中包含 `Creat` 方法和 `Compute` 方法, 如图7.48所示.

`Create` 方法的职责是将创建好的 `baseop` 指针调用 `base_ops_.push_back(op_ptr)`; 储存起来: `Compute` 方法的职责在于调用 `lib` 层的 `compute` 函数进行计算, 目前算子实现为同步方式, 因此需要调用 `SyncQueue`.

### 5. MLUStream

MLUStream 层负责 MLUOps 算子类的实例化, 其接口都定义在 `tensorflow/stream_executor/mlu/mlu_` 中. 可以把 MLUStream 层的算子分为 2 类: 通用模版算子与特例化算子. 对于 Split+Sub+Concat 合并算子来说, 其属于通用模版算子.

通用模版算子满足以下三个条件: MLU 算子的所有输入均来自 `OpKernelContext`; MLU 算子的所有输出顺序须与 `OpkernelContext` 的输出顺序一致; `MLUTensor` 可以被 `CreateMLUTensorFromTensor` 创建, 即 `MLUTensor` 的形状、数据类型与 `TensorFlowTensor` 一致. 如图7.49

### 6. MLUOpKernel

MLUOpKernel 层负责对算子抽象定义, 其继承了 `OpKernel` 类, 其使用方法与 `OpKernel` 基本一致. 对于每个 MLU 算子均需要实现其构造函数与 `ComputeOnMLU` 方法. `MLUOpKernel` 实现的主要功能有参数检查、参数处理、输出形状推断及输出内存分配、调用 `MLUStream` 层接口完成算子计算. 需要在 `cwise_op_sbc_mlu.cc` 文件中进行注册和 `cwise_op_sbc_mlu.h` 实现. 在下图??和7.51中的代码分别显示了这四个不同的步骤

### 7. BUILD 为了将该算子最终集成在 TensorFlow 中, 编译时还需要在 tensorflow/core/kernels/BUILD 中添加以下信息.

### 7.2.5.5 DLP 执行

在将 BCL 实现的 Split+Sub+Concat 合并算子成功集成到框架后, 需要修改之前转换完成的 int8pb 模型, 在模型中增加 Split+Sub+Concat 合并算子, 这一步可以使用编程框架机理章节中介绍的 pb 与 ptxt 相互转换工具实现. 具体在 ptxt 中用 Split+Sub+Concat 节点 7.52 替换原生 Split、Sub、Concat 节点 7.53

### 7.2.6 实验评估

相比其它两个综合实验, EAST 实验困难度较低, 因此难度系数分为 1.0。学生最终得分为 1.0 乘以相应的得分。

- 60 分标准: 完成 BANGC 算子的实现与 CNRT 测试, CNRT 测试时延时低于 30ms。
- 70 分标准: 在 60 分的基础上完成 TensorFlow 的算子集成, 包括 cnplugin 集成与 TensorFlow 的编译, 完成 pb 模型的修改操作。执行测试脚本时, 修改后的模型精度与修改前误差在 5% 以内, 且延时不高于原始模型。
- 80 分标准: 在 70 分的基础上, 执行测试脚本时, 修改后的模型精度与修改前误差在 1% 以内, 且延时至少优于原始模型 10ms。
- 90 分标准: 在 80 分的基础上, 执行测试脚本时, 修改后的模型精度与修改前误差在 0.1% 以内, 且延时至少优于原始模型 25ms。
- 100 分标准: 在 90 分的基础上, 执行测试脚本时, 修改后的模型精度于修改前完全一致, 且延时至少优于原始模型 40ms。

### 7.2.7 实验思考

1. 为什么 BCL 算子相比 Split+Sub+Concat CNML 单算子有显著优势?

```

1  ...
2  #include "macro.h"
3  #include "cnrt.h"
4  #include "utils.h"
5
6  #define CHANNELS 3
7  #define HEIGHT 672
8  #define WIDTH 1280
9  #define BATCH_SIZE 1
10 #define DATA_COUNT ((CHANNELS) * (WIDTH) * (HEIGHT))
11
12 using namespace std;
13 typedef unsigned short half;
14
15 extern "C" {
16     void SBCKernel(half* input_data_, half* output_data_, int batch_num_, int core_num_
17 );
18 }
19 int main() {
20     const int data_count = DATA_COUNT*BATCH_SIZE;
21     int batch_num_ = BATCH_SIZE;
22     int core_num_ = NUM_MULTICORE;
23     const int channels_ = CHANNELS;
24     const int height_ = HEIGHT;
25     const int width_ = WIDTH;
26
27     //开辟CPU 内存
28     ...
29
30     //读取数据文件
31     ...
32
33     //初始化设备
34     ...
35
36     //选择 UNION 模式
37     switch (core_num_) {
38         ...
39     }
40
41     //float2half
42     ...
43
44     // Passing param
45     ...
46
47     // create cnrt Notifier
48     ...
49
50     // hardware time
51     cnrtPlaceNotifier(Notifier_start, pQueue);
52     CNRT_CHECK(cnrtInvokeKernel_V2 (...));
53     cnrtPlaceNotifier(Notifier_end, pQueue);
54     CNRT_CHECK(cnrtSyncQueue(pQueue));
55
56     gettimeofday(&end, NULL);
57     float time_use = ((end.tv_sec - start.tv_sec) * 1000000 + (end.tv_usec - start.
58 tv_usec))/1000.0;
59     printf("time use: %.3f ms\n", time_use);
60
61     float* output_tmp = (float*)malloc(data_count * sizeof(float));
62     ...
63
64     // save data
65     FILE* mluOutputFile = fopen("./mluoutput.txt", "w");
66     for (int i = 0; i < data count; i++) {

```

```

1  cnmlStatus_t cnmlCreatPluginSBCOpParam(
2  cnmlPluginSBCOpParam_t *param,
3  int batch_num_
4  ){
5  *param = new cnmlPluginSBCOpParam();
6  (*param)->batch_num_ = batch_num_;
7
8  return CNML_STATUS_SUCCESS;
9  }
10
11 cnmlStatus_t cnmlDestroyPluginSBCOpParam(
12 cnmlPluginSBCOpParam_t *param
13 ){
14 delete (*param);
15 *param = nullptr;
16
17 return CNML_STATUS_SUCCESS;
18 }
19

```

图 7.43 CNPlugin 添加头文件接口声明

```

1  cnmlStatus_t cnmlCreatePluginSBCOp(
2  ...
3  ){
4
5  int input_num = 1;
6  int output_num = 1;
7
8  void** InterfacePtr;
9  InterfacePtr = reinterpret_cast<void**>(&SBCKernel);
10
11 cnrtKernelParamsBuffer_t params;
12 ...
13
14 cnmlCreatePluginOp(...);
15
16 cnrtDestroyKernelParamsBuffer(params);
17 return CNML_STATUS_SUCCESS;
18 }
19
20 cnmlStatus_t cnmlComputePluginSBCOpForward(
21 ...
22 ){
23
24 cnrtInvokeFuncParam_t compute_forw_param;
25 ...
26 cnmlComputePluginOpForward_V3(...);
27
28 return CNML_STATUS_SUCCESS;
29 }
30
31

```

图 7.44 CNPlugin 算子注册



```

1 REGISTER_OP("SBC")
2 .Input("input: T")
3 .Output("output: T")
4 .Attr("T:type")
5 .SetShapeFn(shape_inference::UnchangedShape);
6

```

图 7.45 EAST 算子注册

```

1 tensorflow::Status CreateSBCOp(MLUBaseOp** op, MLUTensor* input,
2                               MLUTensor* output, int batch_num_) {
3
4     MLUTensor* inputs_ptr[1] = {input};
5     MLUTensor* outputs_ptr[1] = {output};
6
7     CNML_RETURN_STATUS(cnmlCreatePluginSBCOp(op, inputs_ptr, outputs_ptr, batch_num_));
8 }
9
10 tensorflow::Status ComputeSBCOp(
11     MLUBaseOp* op,
12     void* input,
13     void* output,
14     MLUCnrtQueue* queue) {
15
16     void* inputs_ptr[1] = {input};
17     void* outputs_ptr[1] = {output};
18
19     CNML_RETURN_STATUS(cnmlComputePluginSBCOpForward(
20         op, inputs_ptr, 1, outputs_ptr, 1, queue));
21 }
22
23 tensorflow::Status CreateSBCOp(MLUBaseOp** op, MLUTensor* input,
24                               MLUTensor* output, int batch_num_);
25
26 tensorflow::Status ComputeSBCOp(
27     MLUBaseOp* op,
28     void* input1,
29     void* output,
30     MLUCnrtQueue* queue);
31

```

图 7.46 EAST MLULib 层

```

1 DECLARE_OP_CLASS(MLU...) ;
2 DECLARE_OP_CLASS(MLUSBC);
3 DECLARE_OP_CLASS(MLU...) ;
4

```

图 7.47 east-mluops1

```

1  /* Copyright 2018 Cambricon*/
2  #if CAMBRICON_MLU
3
4  namespace stream_executor {
5  namespace mlu {
6  namespace ops {
7
8  Status MLUSBC::CreateMLUOp(std::vector<MLUTensor *> &inputs ,
9                          std::vector<MLUTensor *> &outputs , void *param) {
10     TF_PARAMS_CHECK(inputs.size() > 0, "Missing input");
11     TF_PARAMS_CHECK(outputs.size() > 0, "Missing output");
12
13     MLUBaseOp *op_ptr = nullptr;
14     MLUTensor *input = inputs.at(0);
15     MLUTensor *output = outputs.at(0);
16
17     int batch_num_ = *((int*)param);
18
19     MLULOG(3) << "CreateSBCOp"
20              << ", input: " << lib::MLUTensorUtil(input).DebugString()
21              << ", output: " << lib::MLUTensorUtil(output).DebugString();
22
23     TF_STATUS_CHECK(lib::CreateSBCOp(&op_ptr, input, output, batch_num_));
24
25     base_ops_.push_back(op_ptr);
26
27     return Status::OK();
28 }
29
30 Status MLUSBC::Compute(const std::vector<void *> &inputs ,
31                      const std::vector<void *> &outputs , cnrtQueue_t queue) {
32     void *input = inputs.at(0);
33     void *output = outputs.at(0);
34     TF_STATUS_CHECK(lib::ComputeSBCOp(base_ops_.at(0), input, output, queue));
35
36     TF_CNRT_CHECK(cnrtSyncQueue(queue));
37
38     return Status::OK();
39 }
40
41 } // namespace ops
42 } // namespace mlu
43 } // namespace stream_executor
44
45 #endif // CAMBRICON_MLU
46
47

```

图 7.48 east-mluops2

```

1  Status SBC(OpKernelContext* ctx , Tensor* input , Tensor* output , int batch_size) {
2      return CommonOpImpl<ops::MLUSBC>(ctx , {input} , {output} , static_cast<void*>(&
3      batch_size));
4  }

```

图 7.49 east-mlustream

```

1  #ifndef TENSORFLOW_CORE_KERNELS_CWISE_OP_POWER_DIFFERENCE_MLU_H_
2  #define TENSORFLOW_CORE_KERNELS_CWISE_OP_POWER_DIFFERENCE_MLU_H_
3  #if CAMBRICON_MLU
4
5  namespace tensorflow {
6  template <typename T>
7  class MLUSBCOp : public MLUOpKernel {
8  public:
9      explicit MLUSBCOp(OpKernelConstruction* ctx) :
10         MLUOpKernel(ctx) {}
11
12     void ComputeOnMLU(OpKernelContext* ctx) override {
13
14         if (!ctx->ValidateInputsAreSameShape(this)) return;
15         auto* stream = ctx->op_device_context()->mlu_stream();
16         auto* mlustream_exec = ctx->op_device_context()->mlu_stream()->parent();
17         Tensor input = ctx->input(0);
18
19         // 参数检查与处理
20         const Tensor& a = ctx->input(0);
21         int batch_size = a.dim_size(0);
22
23         // MLU_OP_CHECK_DIM(input, 4, mlustream_exec, ctx);
24
25         string op_parameter = ctx->op_kernel().type_string()
26             + "/" + input.shape().DebugString();
27         MLU_OP_CHECK_UNSUPPORTED(mlustream_exec, op_parameter, ctx);
28
29         TensorShape shape = TensorShape(input.shape());
30         // std::cout << shape << std::endl;
31
32         // 输出形状推断及输出内存分配
33         Tensor* output;
34         OP_REQUIRES_OK(ctx, ctx->allocate_output(0, shape, &output));
35
36         // 调用MLUStream层接口完成算子计算
37         OP_REQUIRES_OK(ctx, stream->SBC(ctx, &input, output, batch_size));
38     }
39 };
40
41
42 } // namespace tensorflow
43
44 #endif // CAMBRICON_MLU
45 #endif // TENSORFLOW_CORE_KERNELS_CWISE_OP_SQUARED_DIFFERENCE_MLU_H_
46

```

图 7.50 east-MLUOpKernel

```
1 namespace tensorflow {
2 #if CAMBRICON_MLU
3 #define REGISTER_MLU(T) \
4     REGISTER_KERNEL_BUILDER(Name("SBC") \
5                             .Device(DEVICE_MLU) \
6                             .TypeConstraint<T>("T"), \
7                             MLUSBCOp<T>);
8 TF_CALL_MLU_FLOAT_TYPES(REGISTER_MLU);
9 #undef REGISTER_MLU
10 #endif
11 }
12
```

图 7.51 east-MLUOpkernel2

```
1 node {
2   name: "concat"
3   op: "SBC"
4   input: "input_images"
5   batch_num_: "1"
6   attr {
7     key: "T"
8     value {
9       type: DT_FLOAT
10    }
11  }
12 }
```

图 7.52 east-pbtxt1

```
1 node {
2   name: "split/split_dim"
3   op: "Const"
4   attr {
5     key: "dtype"
6     value {
7       type: DT_INT32
8     }
9   }
10  attr {
11    key: "value"
12    value {
13      tensor {
14        dtype: DT_INT32
15        tensor_shape {
16        }
17        int_val: 3
18      }
19    }
20  }
21 }
22 .....
23 node {
24   name: "concat"
25   op: "ConcatV2"
26   input: "sub"
27   input: "sub_1"
28   input: "sub_2"
29   input: "concat/axis"
30   attr {
31     key: "N"
32     value {
33       i: 3
34     }
35   }
36   attr {
37     key: "T"
38     value {
39       type: DT_FLOAT
40     }
41   }
42   attr {
43     key: "Tidx"
44     value {
45       type: DT_INT32
46     }
47   }
48 }
```

图 7.53 east-pbtxt2

## 7.3 自然语言处理-BERT

### 7.3.1 实验目的

本实验主要完成将自然语言处理的代表性算法——BERT 网络移植到智能处理器 DLP 上，并进行性能优化与比较，使读者可以借助 DLP 处理器完成完整的目标检测任务。

因此，本实验的主要目的包括：

1. 通过用智能编程语言实现 BERT 模型中的 BatchMatMulV2 算子，深入掌握智能编程语言的算子开发和优化方法；
2. 通过在 TensorFlow 中添加大算子，掌握在 TensorFlow 框架中添加融合算子的方法；
3. 通过使用 TensorFlow 进行在线推理，掌握使用 TensorFlow 编写目标检测应用并在典型 DLP 上进行优化的方法；
4. 通过与 BANGC 大算子实现的比较，体会 DLP 相比 GPU 和 CPU 的优势。

实验时间预计为两周。

### 7.3.2 背景介绍

#### 7.3.2.1 NLP 介绍

总的来说，NLP 的目标是要让计算机正确处理人类的自然语言。常见的 NLP 问题一般包括阅读理解，问答系统，对话系统，文本摘要和文本分类等。与一般的分类或者检测网络不同，NLP 相关的深度学习技术一般都关注于字符之间的序列关系。

斯坦福问答数据集 (Stanford Question Answering Dataset, SQuAD) 是一种阅读理解数据集，由工作人员在一组 Wikipedia 文章上提出的问题组成，其中每个问题的答案是对应阅读文章或问题的一段文本，而问题本身可能是无法回答的。例如对于文本：

James Bryant Conant led the university through the Great Depression and World War II and began to reform the curriculum and liberalize admissions after the war.

有对应问题：

What was the name of the leader through the Great Depression and World War II?

算法应当正确回答：

James Bryant Conant

#### 7.3.2.2 NLP 常见算法、原理

1. **LSTM** 关于 LSTM 的相关原理，在智能计算系统理论书中已经有所介绍。其是在循环神经网络 (RNN) 的基础上，将 RNN 的隐层单元用 LSTM(Long Short Term Memory) 单元来代替，使其能保留更长的时间信息，减缓梯度消失的现象。相较于卷积神经网络，循环神经网络有存储信息的能力，可有效处理序列数据。此外由于其权值共享的特性，循环神经网络的参数一般要少许多卷积神经网络，其结构图如下图 7.54 所示。

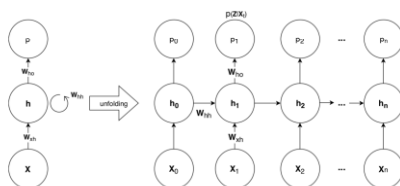


图 7.54 循环神经网络结构图

2. **Attention** Attention 是一种用于提升基于 RNN (LSTM 或 GRU) 的 Encoder + Decoder 模型效果的的机制，一般也称为注意力机制。注意力机制现已广泛应用于机器翻译、语音识别、图像标注 (Image Caption) 等领域。Attention 之所以受欢迎的原因在于其给模型赋予了区分辨别的能力，例如，在机器翻译、语音识别应用中，为句子中的每个词赋予不同的权重，使神经网络模型的学习变得更加灵活 (soft)。具体来说，注意力机制可以帮助模型对输入的每个部分赋予不同的权重，抽取出更加关键及重要的信息，使模型做出更加准确的判断，同时不会对模型的计算和存储带来更大的开销。

对于一般的 Seq2Seq 模型，其通常在编码阶段 (Encoder) 把输入编码成一个固定的长度，这样做会带来两个问题：当输入序列很长时，模型性能会有影响；其对输入的每个元素赋予相同的权重，无法体现重要程度。而注意力机制一定程度上可以解决此问题，其会对输入的每一个元素计算出其对输出端各个元素之间的权重，表示源端第  $i$  个元素对目标端各元素的影响程度。

常见的注意力机制可分为许多种类，常见的有 soft Attention 和 Hard Attention；Global Attention 和 Local Attention；以及 Self Attention 等。其中 Self Attention 为 Transformer 的重要部分，此处简单对齐作出介绍。

Self Attention 与传统的 Attention 机制不同：传统的 Attention 如上所述是基于输入端和输出端的隐层单元计算 Attention 的，得到的结果是源端的每个词与目标端每个词之间的依赖关系。而 Self Attention 是分别对输入端与输出端自身进行计算，捕捉输入端或输出端自身的词与词之间的依赖关系；然后再把输入端得到的 self Attention 加入到输出端得到的 Attention 中，捕捉输入端和输出端词与词之间的依赖关系。因此 Self Attention 最终得到的信息中不仅包含输入与输出的依赖关系，还包含自身词与词之间的依赖关系，其效果要好于传统的注意力机制。

3. **Transformer** 如图 7.55 所示，其编码部分输入 (*inputs*) 对应训练的数据集，解码部分输入 (*outputs-shifted right*) 对应数据集的标签。

像大多数 Seq2Seq 模型一样，Transformer 也分为 Encoder 和 Decoder 两个部分组成 (分别对应图 7.55 中的左半部分和右半部分)。

对于 Encoder 来说，其由 6 个相同的 Layer 组成 (即  $N=6$ )，每个 Layer 又由两个 Sub-Layer 组成，分别是 multi-head self-attention mechanism 和 fully connected feed-forward network。其中每个 Sub-Layer 还增加了 Residual Connection 和 Layer Normalisation。

对于 Decoder 来说，其在 Encoder 的基础上，增加了额外的 Attention Sub-Layer。该子

层同时增加了 Mask 操作，以此确保预测第  $i$  个位置时仅可以获取到  $i$  之前的输出（不会接触到未来的信息）。在训练时其输入为 Encoder 的输出，以及编码器输入对应的标签（整体向右偏移一位），由于标签是始终存在的，因此训练过程是可以做到并行的。在推理时，Decoder 没有 ground truth 作为输入，其输入来自上一个位置的输出，也就意味着需要逐个位置进行解码，无法并行。

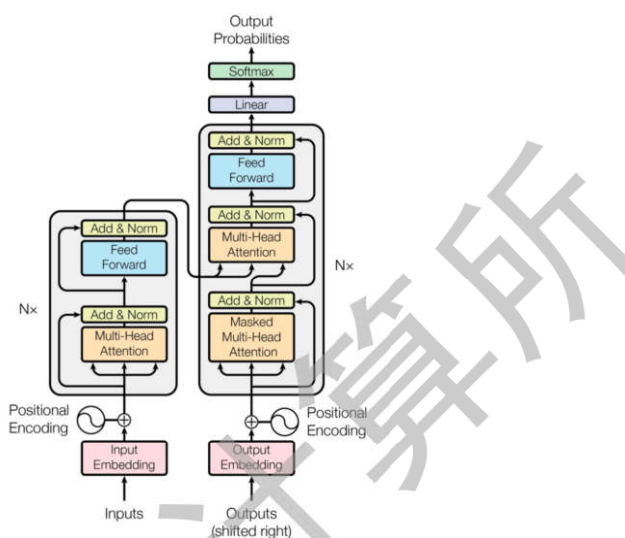


图 7.55 Transformer 结构图 (1)

从图 7.55和图 7.56的结构示意图可以发现，Encoder 和 Decoder 是层叠多个 Multi-Head Attention 单元构成，而每一个 Multi-Head Attention 单元由多个结构相似的 Scaled Dot-Product Attention 单元组成。Self Attention 也是在 Scaled Dot-Product Attention 单元里面实现的，如上图左图所示，首先把输入经过不同的线性变换分别得到  $Q$ 、 $K$ 、 $V$  ( $Q$ 、 $K$ 、 $V$  都来自于相同的输入)。然后把  $Q$  和  $K$  做 dot Product 相乘，得到输入词与词之间的依赖关系，然后经过尺度变换 (scale)、掩码 (Mask) 和 Softmax 操作，得到最终的 Self Attention 矩阵。方程式表达如下，其中  $\text{key}(K)$  意味着原始隐层单元的输出， $\text{value}(V)$  代表每一个隐层单元之间的关系， $Q$  为 outputs-shifted right 部分的输出。

$$Attention(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

4. **BERT** BERT 的全称为 Bidirectional Encoder Representation from Transformers。其采用了 Transformer 的 Encoder 部分，并且进行双向的 Block 连接，与 Transformer 相比，不仅可获取到语句之前的信息，还可获取到未来的信息。其结构图如图 7.57所示。

图中的黄色部分为输入的嵌入 (Embedding) 表示，其由三种不同的 Embedding 求和而成，分别为 Token Embeddings, Segment Embeddings 以及 Position Embeddings。Token Embeddings 是模型输入中每个元素的词向量表示；Segment Embeddings 用来区分输入的前后两个语句；Position Embeddings 含义与 Transformer 相同，但 Transformer 中使用



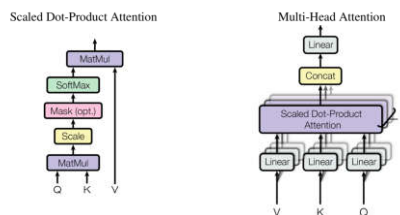


图 7.56 Transformer 结构图 (2)

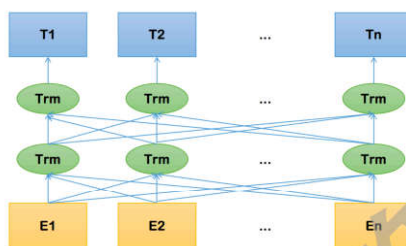


图 7.57 BERT 结构图

不同频率的三角函数直接计算得出，而 BERT 是通过训练学习得到的。图 7.57 中每一个黄色框均代表模型输入的一个元素对应的 Embedding 表示。紧接着将 Embedding 输入到 Transformer 的编码 (Encoder) 部分，由于 Encoder 没有对 Multi-Head Attention 进行 Mask 操作，这也是 BERT 能获取到双向信息的原因。经过多层串联的 Transformer 后即得到最终的输出 (图 7.57 中的蓝色框部分)。

在训练时，BERT 包含两个预训练子任务。第一个子任务为 Masked Language Model，通过随机遮挡每一个句子中 15% 的词，用其上下文来做预测，在计算损失函数时只计算被遮挡掉的输入元素。第二个子任务为 Next Sentence Prediction，可以理解简单的分类任务，用以判断两个输入语句是否为前后关系。

使用预训练模型进行微调时，BERT 针对不同的预测任务给予了不同的方法。如下图 7.58 所示，由此可以看出 BERT 不仅适合做简单的文本分类、序列标注等常见任务，对于问答系统、信息检索、聊天机器人等任务都有着重大的突破。在本实验中，我们完成的是问答系统的任务，对应与图 7.58 中的 (c) 部分。

### 7.3.2.3 BatchMatMul 算子

BatchMatMulV2 算子是 BERT 模型中重要的一个操作，其主要进行两个张量相乘的操作，其中这两个输入张量的 shape 可能不一致，因此还涉及到 Broadcast 操作，其还有两个参数可对两个输入张量进行共轭转制操作。

BatchMatMulV2 的输入输出如下：

Arguments: x: 2 维或更高维度，输入 shape 为  $[..., r_x, c_x]$  y: 2 维或更高维度，输入 shape 为  $[..., r_y, c_y]$

Attrx: adj\_x: 当为 True 时，对输入 x 进行共轭转制，默认为 False adj\_y: 当为 True 时，

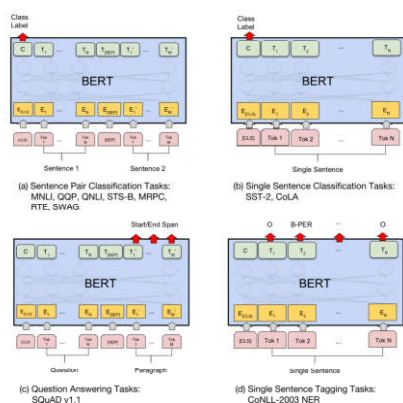


图 7.58 BERT Fine Tuning

对输入  $y$  进行共轭转制，默认为 False

Return: Output: 返回 3 维或更高维度的张量，输出 shape 为  $[...,r_o,c_o]$

其和 BatchMatMul 算子的区别在于 BatchMatMulV2 支持 batch 层面的 broadcasting 操作。

### 7.3.3 实验环境

本实验所涉及的硬件平台和软件环境与其它综合实验相同，参照附录章节 ?? 云平台的使用。

### 7.3.4 实验内容

1. BERT 重训练 (可选): 通过 Fine-Tuning 的方式，进一步提升模型的精度。目前 DLP 还不支持训练，学生可以使用 CPU 或 GPU 完成模型的重训练工作，也可直接使用教材提供的模型做下一步实验。
2. BERT DLP 移植: 将 BERT 网络移植在 DLP 中，包含模型量化，config 设置等操作。
3. BANGC BatchMatMulV2 算子实现与集成。
4. 分别比较 CPU 性能、MLU 性能 (BatchMatMulV2 放在 CPU 上计算)、MLU 性能 (BatchMatMulV2 使用 BANGC 实现)、BERT 大算子网络性能。

### 7.3.5 实验步骤

#### 7.3.5.1 BERT 重训练

由于 BERT 只提供了预训练模型，在做特点任务时，需要先做微调训练。

首先在 GitHub 中，拉取工程: <https://github.com/google-research/bert.git>

根据提示，分别下载模型和训练测试数据集。在这里，我们选择 SQuAD 1.1 任务，因此需要分别下载 Bert-Base(Uncased, 12-layer, 768-hidden, 12-heads, 110M parameters) 模型以及 train-v1.1.json、dev-v1.1.json 数据集。

执行如下所示的命令，在 CPU 或 GPU 环境中进行训练，并验证精度。

```

1 export BERT_BASE_DIR=/path/to/bert-model/uncased_L-12_H-768_A-12
2 export SQUAD_DIR=/path/to/dataset/squad
3 python run_squad.py \
4   --vocab_file=$BERT_BASE_DIR/vocab.txt \
5   --bert_config_file=$BERT_BASE_DIR/bert_config.json \
6   --init_checkpoint=$BERT_BASE_DIR/bert_model.ckpt \
7   --do_train=True \
8   --train_file=$SQUAD_DIR/train-v1.1.json \
9   --do_predict=True \
10  --predict_file=$SQUAD_DIR/dev-v1.1.json \
11  --train_batch_size=12 \
12  --learning_rate=3e-5 \
13  --num_train_epochs=2.0 \
14  --max_seq_length=384 \
15  --doc_stride=128 \
16  --output_dir=/path/to/squad_base/
17
18 python $SQUAD_DIR/evaluate-v1.1.py $SQUAD_DIR/dev-v1.1.json /path/to/squad_base/
   predictions.json

```

图 7.59 BERT Fine Tuning

执行完上述指令后，会打印相应的输出精度：“f1”: xx.xx, “exact\_match”: xx.xx。

由于微调训练实验耗时较长且对硬件环境有所要求，不作为本章的重点，学生只需了解训练的流程，在实际实验阶段可直接使用 DLP 提供的训练模型进行推理实验。

### 7.3.5.2 BERT DLP 移植

为了将 BERT 网络推理在 DLP 硬件上，需要对相关代码进行修改移植。与之前的综合实验类似，首先需要增加 DLP 专属的 config 配置，其次还需要对模型进行量化操作。为了方便推理，还在原有的 run\_squad.py 文件中增加了 checkpoint 模型转 pb 模型的操作，以此简化后续调试的流程。该部分的代码修改增加部分如下图 7.60 所示：

从上图中代码可知，在将 checkpoint 模型转 pb 模型时，由于 checkpoint 格式不包含模型的结构信息，所以需要首先在 session 实例中构造图的结构，然后加载相应的模型并作初始化。紧接着需要将计算图中的变量 (variable) 转为常量 (constant)，指定输出节点并做序列化字符的操作。最终生成的 pb 模型相比原模型仅保留了与输出节点相关的节点，相当于在原模型的基础上做了剪枝的操作，并且保存的常量在推理时不需要初始化操作，进一步提高了推理时的性能。

在配置 DLP 相关的 config 时，增加了 otype\_black\_list 设置，经过该设置会将配置的节点自动执行在 CPU 而非 DLP 设备上。可以将一些 DLP 支持但由于性能精度暂不满足要求的算子暂时执行在 CPU 上。

在进行量化时，为了兼顾性能与精度的要求，本次实验选择 INT16 作为量化精度。首先进行配置文件 bert\_int16.ini 的编写，然后执行 DLP 提供的量化脚本 fppb\_to\_intpb.py 得到量化后的模型。在该实验中，DLP 会将 MatMul、MLP、BatchMatMulV2 等操作量化为 INT16 精度，其它节点仍然以 FP32 或 FP16 的精度进行运算。该部分代码如图 7.61 所示：

### 7.3.5.3 BatchMatMulV2 BANGC 代码实现

完成 DLP 移植后，推理时会发现 DLP 暂时不支持 BatchMatMulV2 算子，由于 BERT 中存在大量的 BatchMatMulV2 操作，不仅会导致 CPU 利用率的升高，还会由于大量的分段造成性能的下降。为了解决该问题，我们设计了 BatchMatMulV2 算子的 BCL 实现，在减小分段的同时通过 INT16 定点计算减小时延提高效率。

1. 设计目标使用 BCL 实现 BatchMatMulV2 算子，并以 cnplugin 封装，供框架直接调用。此模块的参数规格设计为：

输入矩阵 1: input\_0[dim\_0, dim\_1, m, k], fp16; 输入矩阵 2: input\_1[dim\_0, dim\_1, n, k], fp16; 输出矩阵: output[dim\_0, dim\_1, m, n], fp16

2. 设计思路考虑到使用循环向量乘指令不能很好地体现 MLU 的矩阵运算性能，因此设计使用卷积指令来执行 MatMul，右矩阵将会在量化完成后拷贝至 wram 空间，使用 `__bang_conv` 指令进行计算。因为最终是要调用卷积指令来进行矩阵乘法，所以代码中 `n`, `ci`, `co` 是沿用的卷积指令符号系统，这与接口中 `m`, `n`, `k` 的对应关系是： $m \rightarrow n$ ； $n \rightarrow co$ ； $k \rightarrow ci$ 。其中  $n \times ci$  为卷积的输入， $ci \times co$  为卷积的 filter， $n \times co$  为卷积的输出。

3. 函数设计

- **\_\_mlu\_entry\_\_ void BatchMatMulV2Kernel\_MLU270\_half ():**

入口函数，主要用于开辟 NRAM, WRAM 和 SRAM 的空间，遍历 dim\_0 和 dim\_1 循环进行卷积，并调用 `int_matmul_wrap` 函数进行计算，其主要参数如下：

void\* left\_ddr: GDRAM 上的左矩阵

void\* right\_ddr: GDRAM 上的右矩阵

void\* dst\_ddr: GDRAM 上的输出矩阵

int dim\_0: 输入维度，左右矩阵的第一维

int dim\_1: 输入维度，左右矩阵的第二维

int m: 输入维度，左矩阵的第三维

int n: 输入维度，右矩阵的第四维

int k: 输入维度，左矩阵的第四维，右矩阵的第三维

float scale\_0: 左矩阵的量化 scale 参数

int pos\_0: 左矩阵的量化 pos 参数

float scale\_1: 右矩阵的量化 scale 参数

int pos\_1: 右矩阵的量化 pos 参数

- **template <typename IT, typename FT>**

**\_\_mlu\_func\_\_ void int\_matmul\_wrap():**

主要计算流程函数，将左右矩阵分别量化并拷贝至对应的空间，调用 `compute` 函数进行计算，然后将结果拷贝回 GDRAM 上，其主要参数如下（省略与上文重复的参数，下同）：

IT \*nram\_buf: NRAM 上开辟的用于计算的空间

IT \*wram\_buf: WRAM 上开辟的用于计算的空间

IT \*sram\_buf: SRAM 上开辟的用于计算的空间

- **template <typename T> \_\_mlu\_func\_\_ void load\_left():**

将左矩阵拷贝至 NRAM 上。并进行量化操作，然后拷贝至 SRAM 上，将 NRAM 的空间空余出来用于量化和摆放右矩阵。

- **\_\_mlu\_func\_\_ void SegStrategyByFactor():**

m, n 拆分策略函数。按照预设的 factor 搜索 n 值，并按照此 n 值拆分 m 值。

int factor: 预设的 n 的拆分倍数

int k\_up: k 维度进行 64 位对齐的结果

int byteit: 输入数据类型的字节数

int byteft: 输出数据类型的字节数

int &work\_m: (输出) 拆分后每次读取的 m 值

int &work\_n: (输出) 拆分后每次读取的 n 值

- **\_\_mlu\_func\_\_ void SegStrategyByMaxCo():**

最大空间的 m, n 拆分策略函数，按照最大 WRAM 空间寻找 n 值，并按照此 n 值拆分 m 值。

- **template <typename T> \_\_mlu\_func\_\_ void load\_right():**

将右矩阵拷贝至 NRAM 上。并进行量化操作，然后拷贝至 WRAM 上。

T \*right\_inchip: WRAM 上的最大空间

T \*right\_ddr: GDRAM 上的右矩阵

T \*nram\_buf: 临时 NRAM 空间，用于量化右矩阵

int n\_in\_nram: 根据 k\_up 计算出的 NRAM 上的剩余空间一次所能计算的最大 n 值

- **template <typename IT, typename FT> \_\_mlu\_func\_\_ void compute():**

计算卷积，得到 MatMul 的结果

FT \*dst: 存放结果的 NRAM 空间

IT \*left: 存放在 NRAM 上的左矩阵

IT \*right: 存放在 WRAM 上的右矩阵

FT recip\_scale: 量化 scale 参数

int pos: 量化 pos 参数

#### 4. 具体实现

- 准备所需的变量，进行空间划分。
- 最外层对  $n$  方向和  $c$  方向使用两重循环，保证多维度输入结果的准确性。
- 按照最大 SRAM 空间进行  $m$  维度拆分:  $work\_m = SRAM\_BUF\_SIZE / sizeof(TYPE) / k$

记为  $work\_m$ ，循环调用 `int_matmul_wrap` 计算维度为  $[work\_m, n, k]$  的 MatMul。如下图所示：图中蓝色阴影部分即为  $SRAM\_BUF\_SIZE$

- 在每个 `int_matmul_wrap` 运算中， $work\_m * k < SRAM\_BUF\_SIZE$ ，所以能保证完全载入到 SRAM 空间上，故先执行 `load_left`，将 `input_0` 数据 load 到 SRAM 空间上。
- 调用 `SegStrategyByFactor` 和 `SegStrategyByMaxCo` 函数寻找  $n$  和  $work\_m$  的拆分策略，其中一种拆分策略为：

```
int work_n_1 = WRAM_BUF_SIZE / byteit / k_up;
int work_n_2 = (NRAM_BUF_SIZE - byteit * k_up) / byteft;
work_n = ALIGN_DN(MIN(work_n_1, work_n_2), 64);
int work_m_1 = NRAM_BUF_SIZE / (byteit * k_up + byteft * work_n);
int work_m_2 = (NRAM_BUF_SIZE - byteit * 64 * k_up) / (byteit * k_up);
work_m' = MIN(work_m_1, work_m_2);
```

将上述的计算结果记为  $work\_n$ ， $work\_m'$ 。

通过循环 `for (int cur_n = 0; cur_n < n; cur_n += taskDim * work_n)` 计算维度为  $[work\_m, work\_n, k]$  的 MatMul，流程图如 7.64 所示：

图中黄色阴影部分即为  $work\_n * k\_up * byteit$ ，根据选取的拆分策略不同可能为  $WRAM\_BUF\_SIZE$  也可能为较小值。

- 在每个  $[work\_m, work\_n, k]$  的 MatMul 计算中，因为拆分策略中已经保证  $work\_n * k\_up * byteit < WRAM\_BUF\_SIZE$ ，所以能保证完全载入到 WRAM 空间上，故执行 `load_right`，将 `input_1` 数据 load 到 WRAM 空间上。
- 在每个  $[work\_m, work\_n, k]$  的 MatMul 计算中，还需要按照  $work\_m'$  进行拆分，以保证左矩阵和输出矩阵有足够的空间放在 NRAM 上，如下图 7.65 所示：  
上图中将 `output` 展示出来是由于 `output` 的结果与 `input_0` 共用 NRAM 空间，拆分策略中已经保证  $(work\_n * byteit + k\_up * byteft) * work\_m' < NRAM\_BUF\_SIZE$ ，所以输入输出均能完整的载入到 NRAM 空间上。
- 进行 `compute` 计算，计算每个  $[work\_m', work\_n, k]$  的 MatMul。
- 多重循环中按照输出结果应有的地址偏移将 `output` 拷出到 GDRAM 上。

#### 7.3.5.4 BatchMatMulV2 BANGC 代码集成

类似于前两个实验的算子添加过程，集成主要分为 CNPlugin 的集成和 TensorFlow 框架的集成。

1. CNPlugin 集成:分别完成 `cnmlCreatePluginBatchMatMulV2OpParam()`; `cnmlDestroyPluginBatchMatM` 以及 `cnmlStatus_t cnmlComputePluginBatchMatMulV2OpForward()` 构造函数。传入 BCL 计算程序所需的参数:  $scale_0, pos_0, scale_1, pos_1$  (代表算子两个输入的量化参数);  $dim_0, dim_1$  (代表两个输入的第 0 个和第 1 个维度的值);  $m, n, k$  (代表第一个输入的最后两个维度为  $[m, k]$ , 第一个输入的最后两个维度为  $[n, k]$ )。

此外还需在 `cnplugin.h` 头文件中增加相应的接口, 如下图 7.67 所示:

2. TensorFlow 集成:

- 算子注册:
- 定义 MLULib 层接口, 用以连接 CNPlugin 的封装接口:
- 定义 MLUOp 层接口, 实现算子类的 Create 和 Compute 方法。
- 定义 MLUStream 层接口
- 定义 MLUOpKernel 层接口: 在 `tensorflow/core/kernels/batch_matmul_v2_op_mlu.h` 文件中完成以下函数, 并调用 MLUStream 层的 `BatchMatMulV2` 函数。

分别完成 CNPlugin 和 TF 的集成后, 即可运行正常的推理代码, 完成 SQuAD 1.1 验证数据集的测试。完整的推理执行命令如图 7.72

### 7.3.5.5 BERT 大算子

为了进一步比较 BCL 的优越性, DLP 还实现了完整 BERT BCL 大算子的推理。其在 pb 模型中仅使用一个 CNPlugin 中的集成算子代替其它所有小算子, 因此避免了多余的分段。同时由于开发者可以通过 BCL 自由控制硬件计算资源和内存的使用, 进一步提高了片上内存的使用效率, 极大的提升了整体端到端的性能。

使用者在运行 BERT 大算子模型时, 仍然使用 DLP 移植后的开源代码, 仅需将模型来源指定为相应的 pb 模型。

在生成大算子模型时, 与第 7.3.5.2 章节不同的是, 其在量化配置文件中选择大算子作为输出节点。配置文件如下图所示:

运行推理代码后, 会发现此种方式下性能有着大幅度的提升。

### 7.3.6 实验评估

相比其它两个综合实验, BERT 实验困难度最高, 因此难度系数分为 1.2。学生最终得分为 1.2 乘以相应的得分。

- 60 分标准: 完成 `BatchMatMulV2` BCL 算子的实现与 `cnplugin` 的集成, 可以跑通单算子测试程序。
- 70 分标准: 在 60 分的基础上, 执行单算子测试脚本时, 错误率在 1% 以内。

- 80分标准：在70分的基础上，完成TensorFlow算子集成，执行完整的BERT模型推理验证脚本，f1值大于80、exact\_math值大于70，单Batch单次推理平均延时小于100ms。
- 90分标准：在80分的基础上，执行测试脚本时，f1值大于82、exact\_math值大于73，单Batch单次推理平均延时小于60ms。
- 100分标准：在90分的基础上，执行测试脚本时，f1值大于82.5、exact\_match值大于74，单Batch单次推理平均延时小于50ms。

### 7.3.7 实验思考

1. 本章提供的样例代码中通过使用卷积来做矩阵运算，可否使用矩阵乘等基本BCL操作来完成？2. 使用上述不同方式完成BCL算子实现有何不同？哪种效率更高？3. BatchMatMulV2 BCL算子的两个输入最后两个维度为何为[m,k]、[n,k]形式？而不是[m,k]、[k,n]形式，这对TensorFlow框架集成BCL算子时调用CNML有何影响。4. 与大算子BERT相比，有何方法可以进一步提升小算子BERT网络的性能？



```

1 from tensorflow.core.framework import graph_pb2
2 from tensorflow.python.framework import importer
3 from tensorflow.python.framework import ops
4 ...
5 flags.DEFINE_bool("export_frozen_graph", False, "Whether to export SavedModel.")
6 ...
7 def save_dlp_frozen_graph(bert_config, seq_len, init_checkpoint):
8     tf_config = tf.ConfigProto()
9     with tf.Session(config=tf_config) as tf_sess:
10        input_ids = tf.placeholder(tf.int32, [None, seq_len], 'input_ids')
11        input_mask = tf.placeholder(tf.int32, [None, seq_len], 'input_mask')
12        segment_ids = tf.placeholder(tf.int32, [None, seq_len], 'segment_ids')
13
14        # construct origin graph
15        (start_logits, end_logits) = create_model(
16            bert_config=bert_config,
17            is_training=False,
18            input_ids=input_ids,
19            input_mask=input_mask,
20            segment_ids=segment_ids,
21            use_one_hot_embeddings=False)
22
23        output_node_names = ['unstack']
24        tvars = tf.trainable_variables()
25        initialized_variable_names = {}
26        if init_checkpoint:
27            (assignment_map, initialized_variable_names
28             ) = modeling.get_assignment_map_from_checkpoint(tvars, init_checkpoint)
29
30            tf.train.init_from_checkpoint(init_checkpoint, assignment_map)
31            tf.logging.info("**** Load Variables Done ****")
32            tf_sess.run(tf.global_variables_initializer())
33        else:
34            raise ValueError("No init_checkpoint!")
35
36        tf.logging.info("**** Trainable Variables ****")
37        for var in tvars:
38            init_string = ""
39            if var.name in initialized_variable_names:
40                init_string = ", *INIT_FROM_CKPT*"
41            tf.logging.info("  name = %s, shape = %s%s", var.name, var.shape, init_string)
42
43        frozen_graph = tf.graph_util.convert_variables_to_constants(tf_sess,
44            tf_sess.graph.as_graph_def(), output_node_names)
45
46        export_file = os.path.join(FLAGS.output_dir, "frozen_model.pb")
47        with tf.gfile.GFile(export_file, "wb") as f:
48            f.write(frozen_graph.SerializeToString())
49        tf.logging.info("**** Save Frozen Graph Done ****")
50    ...
51 def main(_):
52    ...
53    config = tf.ConfigProto(allow_soft_placement=True,
54                            inter_op_parallelism_threads=1,
55                            intra_op_parallelism_threads=1,
56                            log_device_placement=False)
57    config.mlu_options.core_version = "MLU270"
58    config.mlu_options.precision = "int16"
59    config.mlu_options.save_offline_model = False
60    config.mlu_options.core_num = 16
61    config.mlu_options.optype_black_list = "OneHot"
62    run_config = tf.contrib.tpu.RunConfig(...
63        session_config=config)

```

图 7.60 BERT DLP 移植

```

1 [config]
2 quantization_type = int16
3 device_mode = clean
4 int_op_list = FC
5
6 [model]
7 activation_quantization_alg = naive
8 input_tensor_names = input_ids:0, input_mask:0, segment_ids:0
9 weight_quantization_alg = naive
10 seq_length = 128
11 original_models_path = /path/to/frozen/pb
12 output_tensor_names = unstack:0,unstack:1
13 save_model_path = /path/to/output/frozen_model_int16.pb
14
15 [data]
16 do_lower_case = False
17 doc_stride = 128
18 max_query_length = 64
19 batch_size = 8
20 num_runs = 1
21 vocab_file = bert/uncased_L-12_H-768_A-12/vocab.txt
22 data_list = bert/squad/dev-v1.1.json

```

图 7.61 BERT 量化

```

1 for (int cur_dim0 = 0; cur_dim0 < dim_0; cur_dim0++) {
2     for (int cur_dim1 = 0; cur_dim1 < dim_1; cur_dim1++) {
3         ...
4     }
5 }

```

图 7.62 两重循环

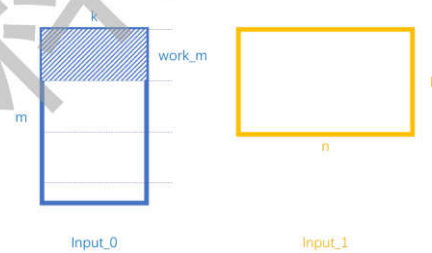


图 7.63 INPUT0 拆分

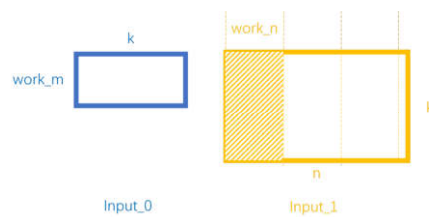


图 7.64 INPUT1 拆分

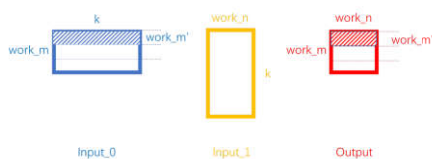


图 7.65 NRAM 拆分

```

1 struct cnmlPluginBatchMatMulV2OpParam {...};
2
3 typedef cnmlPluginBatchMatMulV2OpParam *cnmlPluginBatchMatMulV2OpParam_t;
4
5 cnmlStatus_t cnmlCreatePluginBatchMatMulV2OpParam (...);
6
7 cnmlStatus_t cnmlDestroyPluginBatchMatMulV2OpParam (...);
8
9 cnmlStatus_t cnmlCreatePluginBatchMatMulV2Op (...);
10
11 cnmlStatus_t cnmlComputePluginBatchMatMulV2OpForward (...);

```

图 7.66 BERT-CNPlugin 接口

```

1 在 tensorflow/core/kernels/batch_matmul_op_impl.h 中增加注册
2 #if CAMBRICON_MLU
3 #define REGISTER_BATCH_MATMUL_MLU(T) \
4 REGISTER_KERNEL_BUILDER( \
5     Name("MLUBatchMatMulV2").Device(DEVICE_MLU).TypeConstraint<T>("T"), \
6     MLUBatchMatMulV2<T>)
7 #endif // CAMBRICON_MLU
8
9 在 tensorflow/core/ops/math_ops.cc 中增加注册，并声明该算子只可用于DLP中。
10 REGISTER_OP("MLUBatchMatMulV2")
11     .Input("x: T")
12     .Input("y: T")
13     .Output("output: T")
14     .Attr("T: {half, float}")
15     .Attr("input1_position: int = 0")
16     .Attr("input1_scale: float = 1.0")
17     .Attr("input2_position: int = 0")
18     .Attr("input2_scale: float = 1.0")
19     .Attr("adj_x: bool = false")
20     .Attr("adj_y: bool = false")
21     .SetShapeFn(shape_inference::BatchMatMulV2Shape)
22     .Doc(R"doc(
23 MLUBatchMatMulV2 (DLP only)
24 )doc");

```

图 7.67 算子注册

```

1  ###tensorflow/stream_executor/mlu/mlu_api/lib_ops/mlu_lib_ops.cc
2  tensorflow::Status CreateBatchMatMulV2Op(MLUBaseOp** op,
3      cnmlPluginBatchMatMulV2OpParam_t param,
4      MLUTensor* input1, MLUTensor* input2,
5      MLUTensor* output) {
6      MLUTensor* inputs[2];
7      MLUTensor* outputs[1];
8      inputs[0] = input1;
9      inputs[1] = input2;
10     outputs[0] = output;
11     CNML_RETURN_STATUS(cnmlCreatePluginBatchMatMulV2Op(op, param, inputs, outputs));
12 }
13
14 tensorflow::Status ComputeBatchMatMulV2Op(MLUBaseOp* op,
15     MLUCnrtQueue* queue,
16     void** inputs, int input_num,
17     void** outputs, int output_num) {
18     CNML_RETURN_STATUS(cnmlComputePluginBatchMatMulV2OpForward(
19         op, inputs, input_num, outputs, output_num, queue));
20 }
21
22 ###tensorflow/stream_executor/mlu/mlu_api/lib_ops/mlu_lib_ops.h
23 tensorflow::Status CreateBatchMatMulV2Op(MLUBaseOp** op,
24     cnmlPluginBatchMatMulV2OpParam_t param,
25     MLUTensor* input1, MLUTensor* input2,
26     MLUTensor* output);
27 tensorflow::Status ComputeBatchMatMulV2Op(MLUBaseOp* op,
28     MLUCnrtQueue* queue,
29     void** inputs, int inputs_num,
30     void** outputs, int outputs_num);
31
32 ###tensorflow/stream_executor/mlu/mlu_api/ops/mlu_ops.h
33 struct MLUBatchMatMulV2OpParam {
34     float scale_0_;
35     int pos_0_;
36     float scale_1_;
37     int pos_1_;
38     int dim_0_;
39     int dim_1_;
40     int m;
41     int n;
42     int k;
43     MLUBatchMatMulV2OpParam(float scale_0, int pos_0,
44         float scale_1, int pos_1,
45         int dim_0, int dim_1,
46         int m, int n, int k)
47     : scale_0_(scale_0), pos_0_(pos_0), scale_1_(scale_1), pos_1_(pos_1),
48     dim_0_(dim_0), dim_1_(dim_1), m(m), n(n), k(k) {}
49 };
50 ...
51 DECLARE_OP_CLASS(MLUBatchMatMulV2);

```

图 7.68 MLULib 接口

```

1  ###tensorflow/stream_executor/mlu/mlu_api/ops/batch_matmul_v2.cc
2  #include "tensorflow/stream_executor/mlu/mlu_api/lib_ops/mlu_lib_ops.h"
3  #include "tensorflow/stream_executor/mlu/mlu_api/ops/mlu_ops.h"
4  #include "tensorflow/stream_executor/mlu/mlu_api/tf_mlu_intf.h"
5
6  namespace stream_executor {
7  namespace mlu {
8  namespace ops {
9
10 Status MLUBatchMatMulV2::CreateMLUOp(std::vector<MLUTensor*> &inputs,
11                                     std::vector<MLUTensor*> &outputs,
12                                     void *param) {
13     TF_PARAMS_CHECK(inputs.size() > 1, "Missing input");
14     TF_PARAMS_CHECK(outputs.size() > 0, "Missing output");
15
16     MLUTensor *in0 = inputs.at(0);
17     MLUTensor *in1 = inputs.at(1);
18     MLUTensor *output = outputs.at(0);
19
20     MLULOG(3) << "CreateBatchMatMulV2Op, input1: "
21               << lib::MLUTensorUtil(in0).DebugString()
22               << ", input2: " << lib::MLUTensorUtil(in1).DebugString()
23               << ", output: " << lib::MLUTensorUtil(output).DebugString();
24
25     MLUBatchMatMulV2OpParam *op_param = static_cast<MLUBatchMatMulV2OpParam*>(param);
26     ...
27
28     TF_STATUS_CHECK(lib::CreateBatchMatMulV2Op(...));
29
30     base_ops_.push_back(batch_matmul_op_ptr);
31
32     return Status::OK();
33 }
34
35 Status MLUBatchMatMulV2::Compute(const std::vector<void*> &inputs,
36                                 const std::vector<void*> &outputs,
37                                 cnrtQueue_t queue) {
38     MLULOG(3) << "ComputeMLUBatchMatMulV2";
39     ...
40     TF_STATUS_CHECK(lib::ComputeBatchMatMulV2Op(...));
41
42     TF_CNRT_CHECK(cnrtSyncQueue(queue));
43
44     return Status::OK();
45 }
46
47 } // namespace ops
48 } // namespace mlu
49 } // namespace stream_executor

```

图 7.69 MLUOp 接口

```

1 Status BatchMatMulV2(OpKernelContext* ctx,
2   Tensor* input1, Tensor* input2,
3   float scale_0, int pos_0, float scale_1,
4   int pos_1, int dim_0, int dim_1, int m,
5   int n, int k, Tensor* output) {
6   ops::MLUBatchMatMulV2OpParam op_param(scale_0, pos_0, scale_1,
7     pos_1, dim_0, dim_1, m, n, k);
8   return CommonOpImpl<ops::MLUBatchMatMulV2>(ctx, {input1, input2},
9     {output}, static_cast<void*>(&op_param));
10 }

```

图 7.70 MLUStream 接口

```

1 #ifndef TENSORFLOW_CORE_KERNELS_BATCH_MATMUL_V2_OP_MLU_H_
2 #define TENSORFLOW_CORE_KERNELS_BATCH_MATMUL_V2_OP_MLU_H_
3 #if CAMBRICON_MLU
4 ...
5 #include "tensorflow/stream_executor/mlu/mlu_stream.h"
6
7 namespace tensorflow {
8
9 template <typename T>
10 class MLUBatchMatMulV2 : public MLUOpKernel {
11 public:
12   explicit MLUBatchMatMulV2(OpKernelConstruction* context)
13     : MLUOpKernel(context) {
14     OP_REQUIRES_OK(context, context->GetAttr("adj_x", &adj_x_));
15     OP_REQUIRES_OK(context, context->GetAttr("adj_y", &adj_y_));
16     ...
17   }
18
19   void ComputeOnMLU(OpKernelContext* ctx) override {
20     se::mlu::MLUStream* stream = static_cast<se::mlu::MLUStream*>(
21       ctx->op_device_context()->stream()->implementation());
22     ...
23
24     OP_REQUIRES_OK(ctx, stream->BatchMatMulV2(ctx, ...));
25   };
26
27 private:
28   bool adj_x_;
29   bool adj_y_;
30   ...
31 };
32 } // namespace tensorflow
33 #endif // CAMBRICON_MLU
34 #endif // TENSORFLOW_CORE_KERNELS_BATCH_MATMUL_V2_OP_MLU_H_

```

图 7.71 MLUOpKernel 层接口

```
1 python run_squad.py \  
2 --vocab_file=${BERT_BASE_DIR}/vocab.txt \  
3 --bert_config_file=${BERT_BASE_DIR}/bert_config.json \  
4 --predict_batch_size=8 \  
5 --max_seq_length=128 \  
6 --hidden_size=768 \  
7 --output_dir=${OUTPUT_DIR} \  
8 --export_frozen_graph=false \  
9 --do_predict=true \  
10 --predict_file=${SQUAD_DIR}/dev-v1.1.json \  
11 python ${SQUAD_DIR}/evaluate-v1.1.py \  
12   ${SQUAD_DIR}/dev-v1.1.json \  
13   ${OUTPUT_DIR}/predictions.json
```

图 7.72 BERT 推理

```
1 [config]  
2 quantization_type = int16  
3 device_mode = clean  
4 int_op_list = FC, BatchMatMul  
5  
6 [model]  
7 input_nodes = input_ids, input_mask, segment_ids  
8 layer_num = 12  
9 seq_length = 128  
10 original_models_path = /path/to/frozen_model.pb  
11 output_nodes = unstack  
12 quantization_output_nodes = bert_plugin_op  
13 scope_names = attention/self/query/MatMul, attention/self/key/MatMul,  
14               attention/self/value/MatMul, attention/self/MatMul,  
15               attention/self/MatMul_1, attention/output/dense/MatMul,  
16               intermediate/dense/MatMul, output/dense/MatMul  
17 post_process_name = MatMul  
18 save_model_path = /path/to/output/frozen_model_int16.pb  
19  
20 [data]  
21 do_lower_case = False  
22 doc_stride = 128  
23 max_query_length = 64  
24 batch_size = 8  
25 iters = 1  
26 vocab_file = bert/uncased_L-12_H-768_A-12/vocab.txt  
27 data_list = bert/squad/dev-v1.1.json
```

图 7.73 BERT 大算子配置文件