

文章编号 1004-924X(2023)04-0543-09

基于 ZYNQ 的 Yolo v3-SPP 实时目标检测系统

张丽丽, 陈真*, 刘雨轩, 屈乐乐

(沈阳航空航天大学 电子信息工程学院, 辽宁 沈阳 110000)

摘要: 基于卷积神经网络的目标检测算法发展迅速, 随着计算复杂度增加, 对设备的性能及功耗要求越来越高。为了使目标检测算法能够部署在嵌入式设备上, 本文采用软硬件协同设计方法, 使用 FPGA 对算法进行硬件加速, 提出了 ZYNQ 平台下的 Yolo v3-SPP 目标检测系统。本文将该系统部署在 XCZU15EG 芯片上, 并对系统所需的功耗、硬件资源及性能进行了分析。首先对要部署的网络模型进行优化, 并在 Pascal VOC 2007 数据集上进行训练, 最后使用 Vitis AI 工具对训练后的模型进行量化、编译, 使其适用于 ZYNQ 端的部署。为了选取最佳的配置方案, 探究了各配置对硬件资源及系统性能的影响, 从系统功耗(W)、检测速度(FPS)、各类别平均精度的平均值(mAP)、输出误差等方面对系统进行了分析。结果表明: 在 300 M 时钟频率下, 输入图片大小为(416, 416)时, 针对 Yolo V3-SPP 和 Yolo V3-Tiny 网络结构, 检测速度分别为 38.44 FPS 和 177FPS, mAP 分别为 80.35% 和 68.55%, 片上芯片功耗为 21.583 W, 整板功耗 23.02 W。满足嵌入式设备部署神经网络模型的低功耗、实时性、高检测精度等要求。

关键词: 目标检测; 硬件加速; ZYNQ; Yolo v3-SPP; Yolo v3-Tiny

中图分类号: TP394.1; TH691.9 **文献标识码:** A **doi:** 10.37188/OPE.20233104.0543

Yolo v3-SPP real-time target detection system based on ZYNQ

ZHANG Lili, CHEN Zhen*, LIU Yuxuan, QU Lele

(College of Electronic and Information Engineering, Shenyang Aerospace University,
Shenyang 110000, China)

* Corresponding author, E-mail: chenzen_1996@qq.com

Abstract: The target detection algorithm based on the convolutional neural network is developing rapidly, and with the increase in computational complexity, requirements for device performance and power consumption are increasing. To enable the target detection algorithm to be deployed on embedded devices, this study proposes a Yolo v3-SPP target detection system based on the ZYNQ platform by using a hardware and software co-design approach and hardware acceleration of the algorithm through FPGA. The system is deployed on the XCZU15EG chip, and the required power consumption, hardware resources, and performance of the system are analyzed. The network model to be deployed is first optimized and trained on the Pascal VOC 2007 dataset, and finally, the trained model is quantified and compiled using the Vitis AI tool to make it suitable for deployment on the ZYNQ platform. To select the best configuration scheme, the impact of each configuration on hardware resources and system performance is explored. The system power consumption (W), detection speed (FPS), mean value of average precision (mAP) for

收稿日期: 2022-06-02; 修订日期: 2022-07-14.

基金项目: 国家自然科学基金资助项目(No. 61671310); 辽宁省兴辽英才计划项目基金资助项目(No. XLYC1907134); 辽宁省教育厅项目资助(No. LJKZ0174)

each category, output error, etc. are also analyzed. The experimental results show that the detection speed is 38.44 FPS and 177 FPS for Yolo V3-SPP and Yolo V3-Tiny network structures, respectively, with mAPs of 80.35% and 68.55%, on-chip power consumption of 21.583 W, and board power consumption of 23.02 W at 300 M clock frequency and input image size of (416, 416). This shows that the proposed target detection system meets the requirements of embedded devices for deploying neural network models with low power consumption, real-time, and high detection accuracy.

Key words: object detection; hardware acceleration; ZYNQ; Yolo v3-SPP; Yolo v3-Tiny

1 引言

目标检测是计算机视觉领域的一个重要研究方向,近年来,基于卷积神经网络的目标检测算法取得了巨大突破,在自动驾驶、人脸识别、行人检测等领域都获得了广泛应用,文献[1]针对 Yolo v3 提出使用 Generalized Intersection over Union(GIOU)计算损失、密集连接等方法,实现了 2.11% 的性能提升,文献[2]指出传统的目标检测算法对于小尺寸目标的检测效果较差,针对此问题,文献[3]通过不同通道特征图的叠加,引入空间注意力机制,增强了模型对于小目标的检测能力,文献[4]使用 Mish 激活函数替换 ReLu 激活函数,使用 Complete Intersection over Union(CIoU)计算损失,实现 1.88% 的性能提升,然而随着检测性能的提升,目标检测算法计算冗余、网络参数繁多、计算复杂度大幅增加,使得其只能在高性能计算机上运行,在实时性要求更高的应用场景中,传统的中央处理器的计算架构无法满足实时计算的需求,需要硬件加速器进行加速计算,降低延时,主流的解决方法之一是采用现场可编程逻辑门阵列(Field Programmable Gate Array, FPGA)提高运算速度,FPGA 具有高实时性、低功耗以及并行处理等特点,使其能够完成多种情形的工作,具有较好的实用性和灵活性,适用于目标检测算法的各种应用场景。

文献[5]中,Wei 等人实现了 Yolo 网络的加速,并将网络中的 Leaky ReLu 激活函数替换为 ReLu 激活函数以减少资源消耗,在 ZYNQ7035 中实现了 19 FPS 的性能;文献[6]提出了一种轻量级 Yolo v2 的实现方法,该方法通过使用二值化的特征提取网络减少了计算量与内存消耗,并使用支持向量机回归对物体进行分类,在 XCZU9EG 器件上实现了 40.81 FPS 的性能,该设

计通过降低算法的复杂度成功地提高了速度;文献[7]对 Yolo 网络进行了优化,针对优化后的网络,使用 AXI4 总线封装了应用程序接口,并且使用 ReLu 激活函数替换 Leaky ReLu 激活函数,参数模型全部存储在片上存储器中以减少外部存储器的访问,实现了 19.6 FPS 的性能;文献[8]采用流水线架构,所部署的神经网络中的每一层均映射到专门的硬件模块,在 Virtex XC7VX486 T 器件上实现了 109 FPS 的性能,该方案需要相当大的片上存储空间,对于一些中低端芯片,部署难度大。相对于 Yolo v2 而言,Yolo v3 在检测精度方面有了巨大提升,同时也带来了更多的计算量,于是部分工作中采用 Yolo v3-Tiny 网络进行部署,文献[9]中通过将 Yolo v3-Tiny 网络中的特征图映射为矩阵并且将归一化层与卷积层合并以降低计算复杂度,在 XCZU7EV 器件上实现了 8.3 FPS 的性能;文献[10]使用 Yolo v3-Tiny 网络在 XCZU9EG 器件上实现了 104 FPS 的性能,但其未提到图片大小、可检测物体种类数

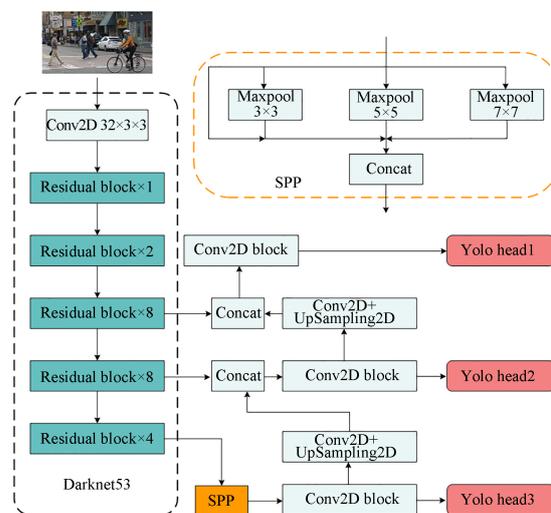


图 1 Yolo v3-SPP 网络结构

Fig. 1 Yolo v3-SPP network structure

以及资源消耗,以上工作需要针对网络进行特定的优化,灵活性低,对开发者硬件知识储备要求较高,对于非专业硬件开发人员而言上手难度大,且开发周期较长,难以适应快速迭代神经网络模型,因此需要一种普适性强,开发周期短,开发流程简洁的神经网络模型加速方法。

针对上述需求,为达到嵌入式设备部署神经网络时所需的低功耗、高检测准确度、实时性以及方便移植的目的,本文提出了一种目标检测网络模型在ZYNQ平台上的实现方法,该方法采用软硬件协同设计,使用ZYNQ芯片中的FPGA部分对算法进行硬件加速处理,实现了两个目标检测模型,分别是改进后的Yolo v3-Tiny与Yolo v3-SPP。首先优化Yolo v3-Tiny与Yolo v3-SPP模型结构使其适用于ZYNQ端的部署,并对其进行了训练,然后将训练好的模型进行量化,再对量化后的模型根据构建的硬件信息进行编译,得到可以在ZYNQ端执行的模型文件,最后编写程序调用该模型文件,达到硬件加速的目的。

2 目标检测算法介绍及优化

Yolo v1-v3^[11-13]是一种被广泛使用的one stage目标检测架构,Yolo v3网络模型在保证检测精度的同时兼顾了检测速度,由于该模型参数量较大,在使用ZYNQ芯片中的FPGA部分对算法进行加速时,无法直接将模型的全部参数存储于FPGA上有限的片上存储器中,且其参数类型为float32,不适于FPGA等硬件设备进行计算加速,因此需要对其进行量化、编译等操作以使其适用于FPGA的部署,而Yolo v3模型在量化、编译之后会有较大的精度损失,为了弥补这个损失,本文采用Yolo v3-SPP网络模型进行部署,该模型相较于Yolo v3模型,主要区别在于特征提取网络之后加入了Spatial Pyramid Pooling(SPP)模块,在该结构中分别使用大小为 1×1 , 5×5 , 9×9 , 13×13 的池化核进行最大池化处理,再将不同尺度的特征图在第一维度进行concatenate操作,得到SPP结构的输出。

SPP模块通过使用不同大小的池化核进行池化操作,增加了特征提取网络提取全局信息的能力,可以提高特征图的表达能力,提高模型的检测性能,而上述SPP结构无法直接应用于

ZYNQ端的部署,因为该结构中包含大小为 9×9 与 13×13 的池化核,由Xilinx公司的产品指南^[14]可知,FPGA上支持的最大池化核尺寸为 8×8 ,一种解决方案是保持SPP结构不变,将模型的中间结果从FPGA传回至CPU处理,待CPU计算完成后再将结果传至FPGA完成之后的计算,这样的缺点是会造成额外的数据搬移,此时数据在CPU与FPGA之间的传输会成为系统整体性能的瓶颈,另一种方法是更改网络结构,使其适用于FPGA的部署从而避免不必要的的数据搬移,提升系统性能。本文选择更改SPP结构,然后对模型进行训练,使用训练完成的模型进行部署,整体网络结构如图1所示,SPP结构对运行速度的影响及模型性能测试见本文4.2节。

3 系统设计

3.1 系统结构设计

ZYNQ为Xilinx公司推出的一种异构加速平台,该平台由两部分构成,为Processing System(PS)端与Programmable Logic(PL)端,即FPGA部分。本文所用芯片属于ZYNQ UltraScale+MPSoCs系列芯片,型号为XCZU15EG-ffvb1156-2-i,所用编译环境为Vivado 2021.2, Petalinux 2021.2, Vitis 2021.2以及Vitis AI 1.4.0。系统采用软硬协同实现目标检测功能,系统框图如图2所示,其中PS端挂载了一张SD卡,卡中烧录有Linux操作系统,系统中包含有opencv、numpy等常用python库,还包含VART

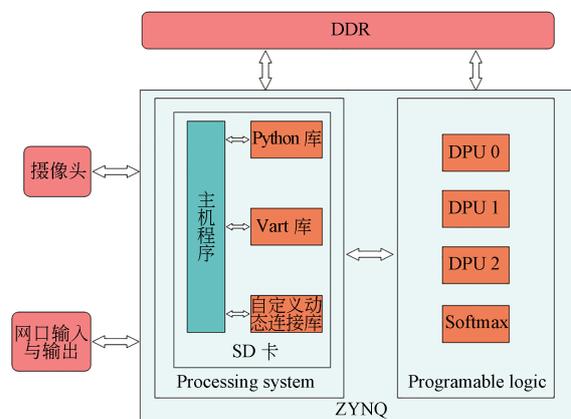


图2 系统框图

Fig. 2 System block diagram

以及自定义动态链接库,摄像头负责获取图像数据,获取到的图像数据传入主机程序中进行处理,处理完成的数据通过 AXI 总线传输至 PL 端的 DPU IP 核中进行运算,PL 端再将运算后的结果返回至 PS 端,PS 端再对结果进行后处理,得到模型的输出结果,另外 PS 端还负责结果的储存与输出。

3.2 DPU IP 核的设计

Deeplearning Process Unity (DPU) 是 Xilinx 公司开发的一款 IP 核。该 IP 核在 ZYQN 中的 PL 端实现,可由 Vitis AI 编译器生成的指令驱动。在使用该 IP 核时需要对其进行配置,部分可配置参数如下,Arch,该参数用于配置 DPU 架构并行度,如 B512, B4096,数字 512 或 4 096 代表每秒最大计算量 (Ops); RAM Usage,该参数用于配置 Block RAM (BRAM) 的使用模式,可配置为低 BRAM 使用与高 BRAM 使用; Number of DPU Cores,该参数用于配置一个 DPU IP 核中的内核个数,可配置为 1~4 个; DSP Usage,该参数用于配置 DSP 的使用模式,可配置为低 DSP 使用与高 DSP 使用; URAM Use per DPU,该参数用于配置 Ultra RAM (URAM) 的使用模式,可配置 DPU 中每个内核使用的 URAM 数量,各参数的配置以及对系统性能影响的测试见 4.1 节。

3.3 ZYNQ 端可执行模型文件生成

ZYNQ 端的可执行模型文件,即 xmodel 文件,基于 Xilinx 公司的 Vitis AI 工具生成,该工具中包含有 Vitis AI 量化器与 Vitis AI 编译器,其中, Vitis AI 量化器负责将浮点型数据量化为定点数据,在量化过程中,需要输入校准图像数据; Vitis AI 编译器负责根据 DPU 的配置参数将经 Vitis AI 量化器量化后的模型文件编译为 xmodel 文件。当 DPU 不支持模型中的某个运算时, Vitis AI 编译器仍可编译该模型。在这种情况下,模型被分成几个部分,每个部分称为子图, DPU 上无法执行的子图将被放在 CPU 上执行,此时, DPU 执行完子图后需要将结果回传至 CPU 中进行相应子图的计算, DPU 需要等待 CPU 计算完成,才可以进行下一子图的计算,这样会增加数据传输以及等待时间,影响计算效率。

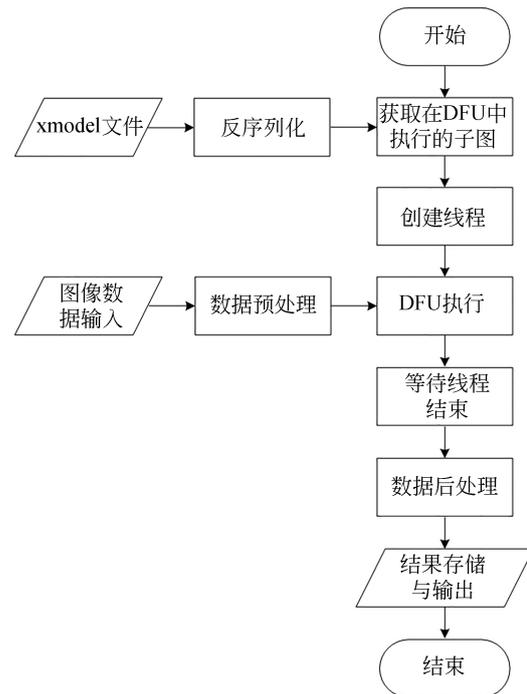


图3 主程序结构

Fig. 3 Main program structure

3.4 ZYNQ 端执行程序设计

执行程序基于 Python 语言编写,采用多线程设计,程序结构如图 3 所示,通过 ctypes 库实现对 c/c++ 的兼容以及动态链接库的调用,程序执行时首先加载 xmodel 文件并对 xmodel 文件进行反序列化,得到所有子图,然后获取各子图执行的位置,取出需要在 DPU 中执行的子图;然后通过 opencv 库中的函数接口获取摄像头输出的图片数据,对图片数据进行预处理,预处理包括数据归一化、数据缩放以及数据类型转换,数据缩放倍数由模型编译时确定。再将预处理后的数据送入 DPU 或 CPU 中完成对应子图的计算,得到网络输出结果,结果经处理后可以存储在 SD 卡中或传输给下一级。

4 性能测试与分析

4.1 DPU 性能测试

为探究 DPU IP 核的内核数、架构并行度以及系统时钟频率对系统性能的影响,本文针对 7 种 DPU 配置方案进行了测试,测试使用改进后的 Yolo v3-SPP 模型,每种方案的配置方式如表 1 所示,其中方案 1、方案 3 架构并行度相同,时钟

频率相同,内核数不同,用于探究时内核数对系统性能的影响;方案 2、方案 3 架构并行度相同,内核数相同,时钟频率不同,用于探究时钟频率对系统性能的影响;方案 4、方案 5 架构并行度相同,内核数相同,时钟频率不同,均使用相同数量的 URAM 代替部分 BRAM,用于探究使用 BRAM 情况下时钟频率对系统性能的影响;方案 6、方案 7 内核数相同,时钟频率相同,架构并行度不同,用于探究架构并行度对系统性能的影响。除上述参数外,每种方案的激活函数开启 Leaky-ReLU, ReLU, ReLU6 支持,开启 Softmax 支持、Depthwise Conv 支持、ElementWise Multiply 支持、AveragePool 支持,开启 channel augmentation,配置为低 RAM 使用率,高 DSP 使用率。

针对以上 7 种方案使用 1 000 张图片对系统进行性能测试,测试结果包括各类型资源占用率(%)、系统功耗(W)以及运行速度(FPS),结果如图 4 所示,对比方案 1 和方案 3,在架构并行度为 B4096,时钟频率为 300 MHz 的情况下,增加一个 DPU 内核,FPS 提升 51.68%,同时会增加约 96.68% 的功耗,各类别资源占用平均增加 17.25%;对比方案 2 和方案 3,在架构并行度为 B4096,内核数为 2 的情况下,时钟频率增加 50%,FPS 可提升 31.30%,同时会增加约

表 1 DPU 配置方式

Tab. 1 DPU configuration method

方案	Arch	Number of Cores	Ultra RAM/ Cores	Clock /MHz
1	B4096	1	0	300
2	B4096	2	0	200
3	B4096	2	0	300
4	B4096	3	28	200
5	B4096	3	28	300
6	B1152	4	0	300
7	B512	4	0	300

31.71% 的功耗,各类别资源占用平均减少 0.33%;对比方案 6 和方案 7,在内核数为 4,时钟频率为 300 MHz 的情况下,增加 1.25 倍的架构并行度,FPS 可提升 44.84%,同时会增加约 29.95% 的功耗,各类别资源占用平均增加 10.89%;对比方案 4 和方案 5,在架构并行度为 B4096,内核数为 3 且使用 URAM 的情况下,时钟频率增加 50%,FPS 可提升 35.59%,同时会增加约 22.33% 的功耗,各类别资源占用平均增加 5.34%;综上所述,考虑到板卡功耗,资源占用以及系统性能等因素,本文使用方案 5 进行部署。

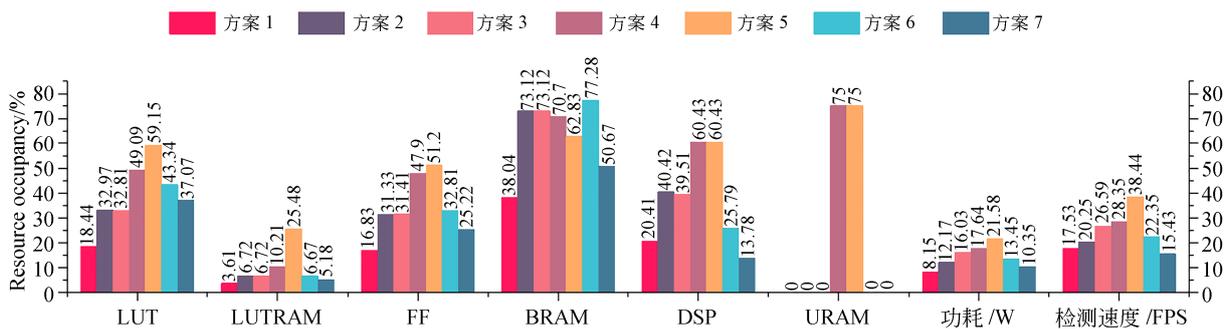


图 4 不同方案资源占用率及性能分析

Fig. 4 Analysis of resource occupancy of different schemes

4.2 SPP 结构性能测试

本文采用更改 SPP 结构的方式使得 Yolo v3-SPP 模型适用于 ZYNQ 端的部署,将原 SPP 结构的池化核大小由 $1 \times 1, 5 \times 5, 9 \times 9, 13 \times 13$ 改为 $1 \times 1, 3 \times 3, 5 \times 5, 7 \times 7$,为探究更改后的 SPP 结构与原 SPP 结构对模型性能以及部署后模型

运行速度的影响,使用不同的 SPP 结构进行了测试,分别对包含二者的模型采用相同的训练方案进行训练,使用 Pascal VOC 2007 测试集对二者进行性能测试,测试指标为 mAP,阈值取 0.5;另外在 ZYNQ 端使用 1 000 张图片分别对各模型进行推断速度的测试,结果如表 2 所示,其中 SPP-o

表示使用原 SPP 结构, SPP-n 表示使用更改后的 SPP 结构, 可知, 更改后的 SPP 结构模型性能提升效果与原 SPP 结构相当; 表中还列出了每个 DPU 内核推理一张图片所需的平均时间 (avg) 以及 3 个 DPU 内核推理一张图片所需的平均时间 (total avg), 该时间由每个 DPU 内核所需的平均时间计算得到, 从表中可以看出, 增加了 SPP 结构以后, ZYNQ 端的 DPU 内核平均推断速度由 76.225 ms 分别增加到 76.391 ms 与 76.331 ms, 平均推断耗时分别增加了 0.166 ms 与 0.106 ms, 增加的 SPP 结构对系统推断速度影响较小, 因此, 本文采用包含 SPP-n 结构的模型进行后续测试与部署。

4.3 模型性能测试

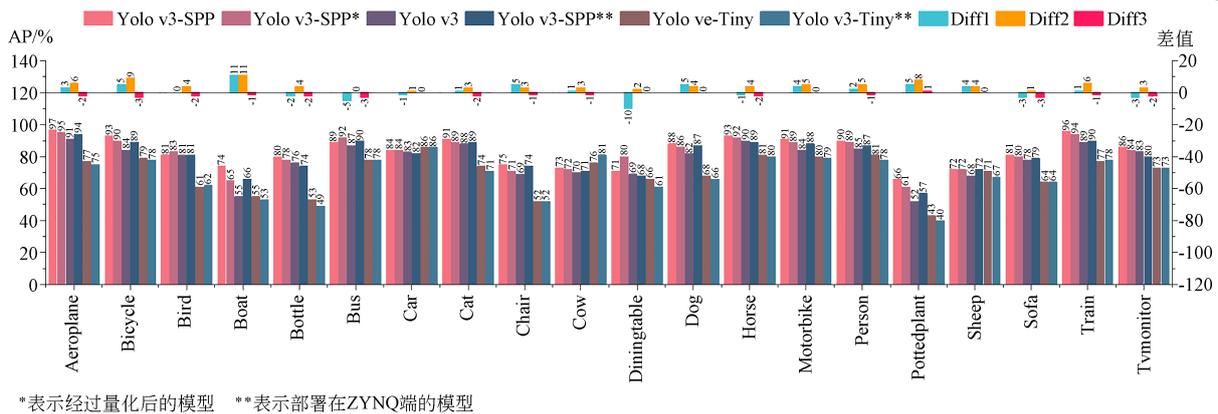
为探究各模型之间的性能差异, 本文分别在 PC 端与 ZYNQ 端共部署了 6 个神经网络模型, 分别是部署在 PC 端的 Yolo v3-SPP, Yolo v3, Yolo v3-Tiny 以及量化后的 Yolo v3-SPP 模型, 部署在 ZYNQ 端的 Yolo v3-SPP 与 Yolo v3-Tiny 模型, 针对各类别的 AP 值对各模型进行了测试, 阈值取 0.5, 使用 Pascal VOC 2007 数据集进行测试, 共 20 个物体类别, 测试结果如图 5 所示, 图中展示出了不同模型各类别 AP 值以及不同模型

表 2 SPP 结构对运行速度的影响

Tab. 2 Influence of SPP structure on running speed

SPP 结构	内核	Avg/ms	Total avg/ms	mAP/%
NA	1	76.941		
	2	74.056	76.225	79.25
	3	77.679		
SPP-o	1	76.658		
	2	74.866	76.391	83.61
	3	77.651		
SPP-n	1	77.232		
	2	74.228	76.331	83.55
	3	77.534		

的结果在该类别处的差值; 其中, Diff 1 为部署在 ZYNQ 端的 Yolo v3-SPP 模型与部署在 PC 端的 Yolo v3 模型所得结果的差值, Diff 2 为部署在 PC 端的 Yolo v3-SPP 模型与 Yolo v3 模型所得结果的差值, Diff 3 为部署在 PC 端的 Yolo v3-SPP 模型量化后与量化前所得结果的差值, 整体来看, 模型经量化后其性能并未发生明显下降, 符合量化要求。



*表示经过量化后的模型 **表示部署在ZYNQ端的模型

图 5 各模型性能对比及差值
Fig. 5 Analysis of model performance after quantification

部署在 PC 端的 Yolo v3-SPP 模型的 mAP 为 83.55%, 量化后的 Yolo v3-SPP 模型的 mAP 为 82.62%, 部署在 PC 端的 Yolo v3 模型的 mAP 为 79.25%, Yolo v3-Tiny 模型的 mAP 为 69.75%; 部署在 ZYNQ 端的 Yolo v3-SPP

模型的 mAP 为 80.35%, 部署在 ZYNQ 端的 Yolo v3-Tiny 模型的 mAP 为 68.55%, 部署在 ZYNQ 端的 Yolo v3-SPP 模型的 mAP 相较于 PC 端有 3.2% 的下降, 但仍高于 PC 端的 Yolo v3 模型, 整体来看, ZYNQ 端的模型检测性

能发生了少许下降,在可接受范围之内,符合部署要求。

图 6 为 ZYNQ 端运行的 Yolo v3-SPP 模型输出结果与 PC 端运行的 Yolo v3-SPP 模型输出对比,其中(a)为 ZYNQ 端执行模型时输出的结果,

(b)为 PC 端执行模型时输出的结果,在简单场景下,ZYNQ 端的检测结果与 PC 端结果无明显区别,如图片 A,复杂场景下,ZYNQ 端对图片中小尺寸目标的检测能力略差于 PC 端,如图片 B、图片 C、图片 D。

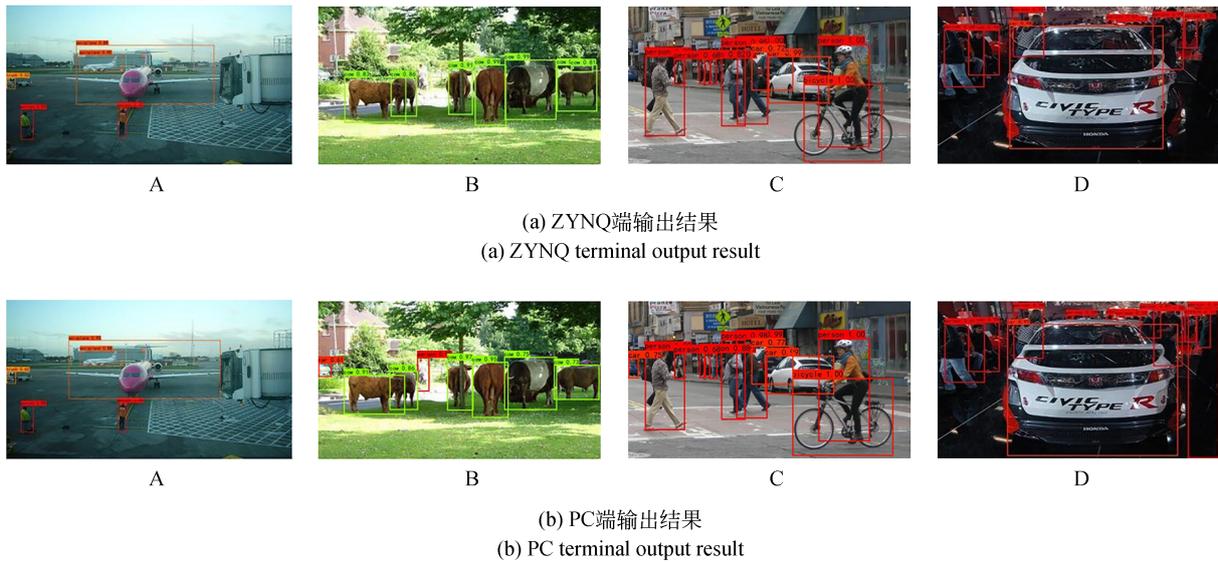


图 6 运行结果对比

Fig. 6 Comparison of running results

本文实现的 Yolo v3-SPP 模型可检测类别个数为 80,模型的最终输出包含有 3 个输出特征层,每个输出特征层有 255 个维度,其计算方式为 $(80 + 4 + 1) \times 3$,其中 80 为各类别对应概率,模型每个输出特征层中的每个特征点存在 3 个先验框,故进行乘 3 操作,每个先验框含有 4 个调整参数,另外还有 1 个参数表示该先验框内是否含有物体,图 7 展示了使用该模型对同一张图片进行推断时,每个输出特征层上 ZYNQ 端的输出结果、PC 端的输出结果以及二者的差值,横坐标表示维度,纵坐标表示输出值的大小,其中黄色折线表示 ZYNQ 端的输出结果,蓝色折线表示 PC 端的输出结果,红色折线表示二者的差值;可以看出,体现二者差值的红色折线稳定在 0 值附近,表明 ZYNQ 端的输出结果与 PC 端的输出结果基本保持一致,满足实际部署要求(彩图见期刊电子版)。

另外,将本文的实验结果与其他文献进行了比较,主要比较指标为检测速度(FPS)、功耗

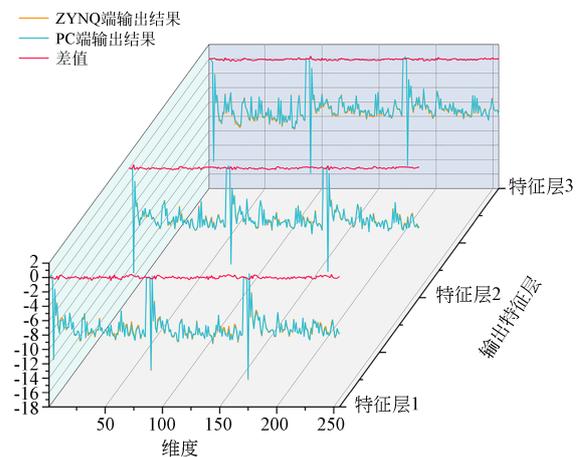


图 7 特征层输出对比

Fig. 7 Feature layer output comparison

(W)以及 mAP(阈值取 0.5),测试使用 Pascal VOC 2007 数据集;结果如表 3 所示,可知,对于 Yolo v3-Tiny 模型,本文实现了最高的 FPS,为 177FPS,分别是文献[15]、文献[16]的 57.65 倍和 51 倍,mAP 仅次于 GPU 和文献[9],为 68.55%,实现的 YOLOv3-SPP 模型达到了除 GPU

外了最高的 mAP, 为 80.35%, 较文献[15]、文献[9]、文献[6]分别有 21.95%, 5.35% 和 12.75%

的提升, 检测速度仅次于 GPU 和文献[6], 为 38.44 FPS。

表 3 与其他硬件平台的性能对比

Tab. 3 Performance comparison with other hardware platforms

平台	GPU	文献[15]	文献[9]	文献[6]	文献[16]	本文			
年份	2021	2021	2021	2018	2020	2021			
数据集	PascalVOC 2007								
器件	GTX1660Ti	Zedboard	Ultra96V2	XCZU9EG	XCZU9EG	XCZU15EG			
网络类型	Yolo v3-SPP	Yolo v3-Tiny	Yolo v3	Yolo v3-Tiny	Yolo v3-Tiny	Lightweight Yolo v2			
数据类型	float32	fixed-16	fixed-8	binary	fixed-16	fixed-8			
图像大小	(416, 416)	(416, 416)	(416, 416)	(224, 224)	(416, 416)	(416, 416)			
FPS	48	258	52	3.07	8.3	40.81	3.47	177	38.44
功耗(W)	118	24.32	4.26	4.5	11.8	23.02			
mAP(%)	83.55	69.75	79.25	58.4	75	67.6	\	68.55	80.35

5 结 论

本文通过对 Yolo v3-SPP 网络结构进行优化, 使其目标检测性能进一步提升的同时适用于 ZYNQ 端的部署, 弥补了 ZYNQ 端部署模型时带来的性能损失, 使用 Vitis AI 工具对预训练网络模型进行量化、编译, 生成板卡运行所需的模型文件, 采用多线程思想编写程序, 充分利用板卡资源, 完成了 Yolo v3-Tiny 和 Yolo v3-SPP 网络

模型在 ZYNQ 端的部署, 然后对系统进行了性能测试, 测试结果表明, 本文所实现的系统在保证检测速度和网络性能的前提下完成了模型在 ZYNQ 端的部署, 对于 Pascal VOC 2007 数据集, 图片输入大小为 (416, 416), 可检测类别数为 80, 分别实现了 177 FPS 和 38.44 FPS 的性能, mAP 分别为 68.55% 和 80.35%, 优于以往部署在嵌入式设备的目标检测系统的设计, 本文所提出的设计有望促进边缘设备实现实时目标检测。

参考文献:

- [1] 唐悦, 吴戈, 朴燕. 改进的 GDT-YOLOV3 目标检测算法[J]. 液晶与显示, 2020, 35(8): 852-860.
TANG Y, WU G, PIAO Y. Improved algorithm of GDT-YOLOV3 image target detection[J]. *Chinese Journal of Liquid Crystals and Displays*, 2020, 35(8): 852-860. (in Chinese)
- [2] 范丽丽, 赵宏伟, 赵浩宇, 等. 基于深度卷积神经网络的目标检测研究综述[J]. 光学精密工程, 2020, 28(5): 1152-1164.
FAN L L, ZHAO H W, ZHAO H Y, *et al.* Survey of target detection based on deep convolutional neural networks[J]. *Opt. Precision Eng.*, 2020, 28(5): 1152-1164. (in Chinese)
- [3] 鞠默然, 罗海波, 刘广琦, 等. 采用空间注意力机制的红外弱小目标检测网络[J]. 光学精密工程, 2021, 29(4): 843-853.

- JU M R, LUO H B, LIU G Q, *et al.* Infrared dim and small target detection network based on spatial attention mechanism [J]. *Opt. Precision Eng.*, 2021, 29(4): 843-853. (in Chinese)
- [4] 王宸, 张秀峰, 刘超, 等. 改进 YOLOv3 的轮毂焊缝缺陷检测[J]. 光学精密工程, 2021, 29(8): 1942-1954.
WANG CH, ZHANG X F, LIU CH, *et al.* Detection method of wheel hub weld defects based on the improved YOLOv3 [J]. *Opt. Precision Eng.*, 2021, 29(8): 1942-1954. (in Chinese)
- [5] WEI G J, HOU Y Z, CUI Q M, *et al.* YOLO acceleration using FPGA architecture [C]. 2018 *IEEE/CIC International Conference on Communications in China (ICCC)*. August 16-18, 2018, Beijing, China. IEEE, 2019: 734-735.
- [6] NAKAHARA H, YONEKAWA H, FUJII T, *et al.* A lightweight YOLOv2: a binarized CNN with

- A parallel support vector regression for an FPGA [C]. *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays. February 25 - 27, 2018, Monterey, CALIFORNIA, USA. New York: ACM, 2018: 31-40.*
- [7] LI Z G, WANG J T. An improved algorithm for deep learning YOLO network based on Xilinx ZYNQ FPGA [C]. *2020 International Conference on Culture-oriented Science & Technology (ICCTST). October 28-31, 2020, Beijing, China. IEEE, 2020: 447-451.*
- [8] NGUYEN D T, NGUYEN T N, KIM H, *et al.* A high-throughput and power-efficient FPGA implementation of YOLO CNN for object detection [J]. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems, 2019, 27(8): 1861-1873.*
- [9] ADIONO T, PUTRA A, SUTISNA N, *et al.* Low latency YOLOv3-tiny accelerator for low-cost FPGA using general matrix multiplication principle [C]. *IEEE Access. October 15, 2021, IEEE, 2021: 141890-141913.*
- [10] OH S, YOU J H, KIM Y K. Implementation of compressed YOLOv3-tiny on FPGA-SoC [C]. *2020 IEEE International Conference on Consumer Electronics-Asia (ICCE-Asia). November 1-3, 2020, Seoul, Korea (South). IEEE, 2020: 1-4.*
- [11] REDMON J, DIVVALA S, GIRSHICK R, *et al.* You only look once: unified, real-time object detection [C]. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). June 27-30, 2016, Las Vegas, NV, USA. IEEE, 2016: 779-788.*
- [12] REDMON J, FARHADI A. YOLO9000: better, faster, stronger [C]. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). July 21-26, 2017, Honolulu, HI, USA. IEEE, 2017: 6517-6525.*
- [13] REDMON J, FARHADI A. YOLOv3: an incremental improvement [EB/OL]. 2018: *arXiv: 1804.02767*. <https://arxiv.org/abs/1804.02767>
- [14] DPUCZDX8G for zynq ultraScale+ MpSoCs product guide PG338 (v3.4) [EB/OL]. Xilinx, [2022-01-20]. https://docs.xilinx.com/r/en-US/pg338-dpu? tocId=3xsG16y_QFTWvAJKHbisEw
- [15] ZHANG H B, JIANG J Q, FU Y H, *et al.* Yolo v3-tiny object detection SoC based on FPGA platform [C]. *2021 6th International Conference on Integrated Circuits and Microsystems (ICICM). October 22-24, 2021, Nanjing, China. IEEE, 2021: 291-294.*
- [16] ZHANG S, CAO J, ZHANG Q, *et al.* An FPGA-based reconfigurable CNN accelerator for YOLO [C]. *2020 IEEE 3rd International Conference on Electronics Technology (ICET). IEEE, 2020: 74-78.*

作者简介:



张丽丽(1979—),女,黑龙江省讷河人,博士,副教授,硕士生导师,2002年、2005年、2012年于吉林大学分别获得学士、硕士、博士学位,主要从事FPGA系统设计及深度学习算法的研究。E-mail: 20052727@sau.edu.cn

通讯作者:



陈真(1996—),男,河南商丘人,硕士研究生,2020年于沈阳航空航天大学获得学士学位,主要从事深度学习以及FPGA的算法研究。E-mail: chen zhen_1996@qq.com