

文章编号: 2095-2163(2023)06-0108-04

中图分类号: TP332

文献标志码: A

基于硬件描述语言的目标识别硬件加速器设计

张嘉, 金婕

(上海工程技术大学 电子电气工程学院, 上海 201620)

摘要: 本文针对深度学习方法在目标识别领域内的应用,设计了一种基于硬件描述语言的目标识别硬件加速器,运用数据流架构优化方法,设计了二值化卷积神经网络算法所对应的硬件模块单元,实现了对输入图片的识别。实验结果表明,基于多帧分辨率为224×224的图片输入,在硬件平台仿真软件中达到了不俗的识别速率以及识别准确率,为基于硬件系统的深度学习加速器的研究奠定了基础。

关键词: 硬件描述语言; 二值化卷积神经网络; 深度学习; 硬件加速器

Design of object recognition hardware accelerator based on hardware description language

ZHANG Jia, JIN Jie

(College of Electronic and Electrical Engineering, Shanghai University of Engineering Science, Shanghai 201620, China)

[Abstract] Aiming at the application of deep learning method in the field of target recognition, this paper designs a hardware accelerator for target recognition based on hardware description language. This paper uses the data flow architecture optimization method to design the hardware module unit corresponding to the binarized convolutional neural network algorithm, and realizes the recognition of the input image. The experimental results show that under the input of multi-frame images with a resolution of 224×224, excellent recognition speed and recognition accuracy are achieved in the hardware platform simulation software, which lays the foundation for the research of deep learning accelerators based on hardware systems.

[Key words] Hardware Description Language (HDL); binarized convolutional neural network; deep learning; hardware accelerator

0 引言

图像理解在现实生活中是一种非常常见的应用。但在边缘化、高精度、运算量需求巨大等的智能计算机平台的应用场景中,需要进行目标识别、姿态估计、图像增强等基于图像任务的复杂计算,计算机视觉系统的重要性逐渐突出,其稳定性、实时性以及便捷性也变得愈发重要。

近年来,卷积神经网络在计算机视觉领域取得巨大进展,神经网络处理器(Neural Network Processor, NNP)IP在片上系统中广泛应用,例如寒武纪的思元370等。然而业界进行神经网络设计时均采用寄存器传输级别(Register Transfer Level, RTL)的硬件单元模块的搭建,如神经网络卷积处理核单元、并行处理单元

以及多级静态随机存取存储器(Static Random-Access Memory, SRAM)缓存存储单元等。

本文深入研究并分析每一层的二值化卷积神经网络结构以及网络层算法的数学模型,运用硬件描述语言描述了算法的数学模型所一一对应的硬件电路模块,设计并优化了二值化深度卷积神经网络寄存器传输级别(Register Transfer Level, RTL)的硬件模型,实现了中小型尺寸输入的二值化深度学习卷积神经网络硬件加速器的设计与验证,并在该硬件加速器中通过设计验证测试平台验证了目标识别应用的准确性与速率。

1 二值化卷积神经网络

二值化卷积神经网络的网络拓扑结构类似于通用

基金项目: 国家自然科学基金(61801286)。

作者简介: 张嘉(1997-),男,硕士研究生,主要研究方向: FPGA、神经网络; 金婕(1978-),女,博士、副教授,主要研究方向: 视频编解码、数字信号处理和VLSI。

通讯作者: 金婕 Email: 02140001@sues.edu.cn

收稿日期: 2022-06-30

的单精度浮点卷积神经网络, 同样拥有卷积层、最大池化层、批归一化层以及全连接层。而与通用的单精度浮点卷积神经网络不同的是二值化卷积神经网络在训练及推断过程中所进行计算的权重阈值数据为 1 或 -1, 且随之对每一层网络层做了特定优化, 从而大幅度地减少了硬件内存资源以及计算资源的占用率^[1]。

与其他数据类型的卷积神经网络不同, 二值化卷积神经网络(Binary Neural Network, BNN)在网络训练与网络推断的运算过程中, 网络输入与网络权重阈值参数的数据类型均为 1 bit 的二进制数即 1 或 -1, 而在 PC 端进行网络训练时凭借符号函数对网络参数数据二值化, 公式(1):

$$y^a = \text{sign}(x) = \begin{cases} +1, & x \geq 0 \\ -1, & x < 0 \end{cases} \quad (1)$$

其中, x 为单精度浮点参数, y^a 为二值化参数。

因为量化的参数具有泛化性, 即可将网络输入与网络参数均二值化为数据 1 或 -1, 而硬件中最为常用的二进制比特位为 1 或 0, 故可将 0 代表上述数据中的 -1 以便在硬件中实现。由此, 卷积运算中的乘累加计算单元即可转换为对硬件较为友好的异或运算和点累加计算方式, 公式(2):

$$y_{out,m,n} = \sum_{i=0}^{k_m} \sum_{j=0}^{k_n} \text{xnor}(w_{i,j}^b, y_{m+i,n+j}^b) \quad (2)$$

其中, $w_{i,j}^b$ 为二进制的权重数据; $y_{m+i,n+j}^b$ 为二进制的输入数据; xnor 即为异或运算。

2 架构设计

数据流架构如图 1 所示。在搭载于 FPGA/ASIC 硬件平台的加速器数据流体系结构中, 通过为

神经网络拓扑结构中的每一层网络层分配成比例的计算资源来实现层间并行化与层内并行化, 而在各个搭载了逐级网络层的计算单元之间插入了数据缓冲区, 加速了各个计算单元之间的数据流动, 更好地实现硬件资源中的数据流水化。此外, 每一个计算单元均有特定的参数接口可供配置, 根据搭载于这个计算单元内的网络层进行特定优化。同时数据流架构中的参数存储分布于片上内存和片外内存两侧存储区, 可根据数据访问频率来划分, 如将部分输入图像的像素数据存储于片外内存, 将权重阈值数据存储于片上内存, 由此可降低片上内存资源占用与片外内存带宽开销。

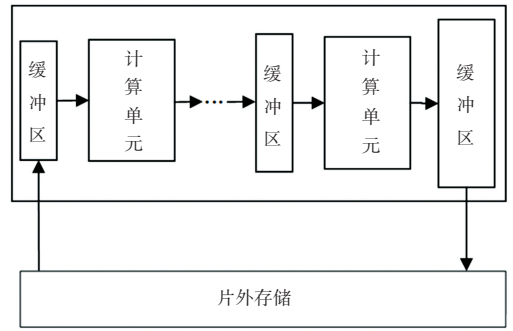


图 1 数据流架构
Fig. 1 Dataflow architecture

3 硬件单元设计

运用并行化可综合的硬件描述语言 Verilog 将神经网络软件设计转化为硬件思想并加以实现, 如网络中的乘法操作转化为硬件中的异或阵列。整体硬件电路的框架图如图 2 所示。

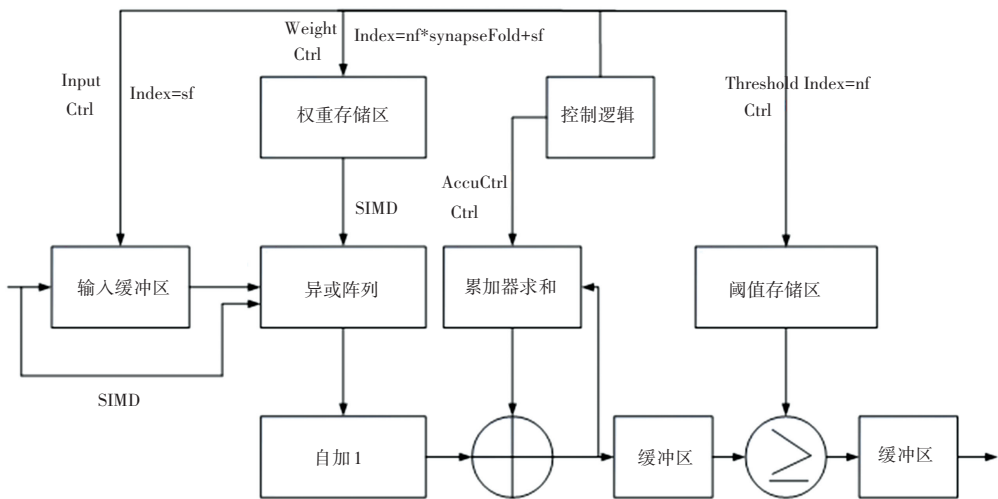


图 2 整体硬件电路框架图
Fig. 2 Overall hardware circuit frame diagram

3.1 位宽转换单元设计

依据算法特性,需将存储于文件中的 64 bit 像素数据转换成 24 bit 数据,即为串并数据的转换功能。位宽转换的时序设计如图 3 中的时序图所示, MERGE 即融合阶段标定为串转并, SLICE 即分割阶

段标定为并转串,其中 M0、M1、M2 为随时序触发的 3 段拼接操作, S0、S1、S2、S3、S4、S5、S6、S7 即为随时序触发的 8 段截取操作,其余阴影部分为无效数据,在实际实验中将阴影部分的数据设置为 0。

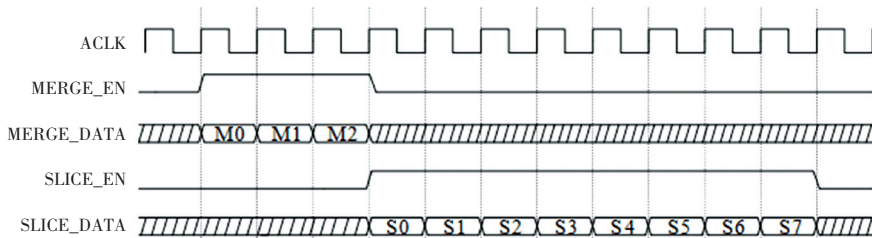


图 3 位宽转换时序图

Fig. 3 Timing diagram of bit width conversion

3.2 零填充单元设计

零填充单元维持了特征图边缘信息的完整性以及特征图形状大小的稳定性。在零值填充单元之前,已通过位宽转换单元将数据整合为每一个像素数据为 24 bit 位宽的形式,并通过设定为 1 的填充维度对整个像素矩阵的周围进行零元素的填充。因

为输入像素矩阵为特定的特征图维度,假定该维度为 224,则在该特征图周围填充维度为 1 的零元素后,整体特征图维度则变为 226。本设计通过特征图边缘检测算法,配合位宽转换单元中的细粒度有限状态机,利用计数器对该模块的输入输出数据进行控制,零填充时序图如图 4 所示。

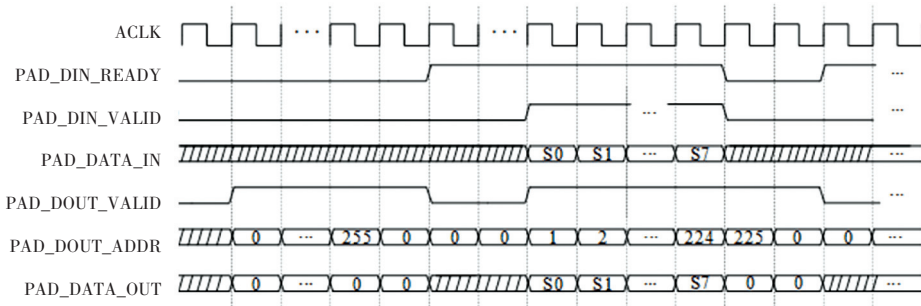


图 4 零填充时序图

Fig. 4 Timing diagram of zero padding

3.3 滑动窗口地址生成器单元设计

滑动窗口主要针对于输入特征图进行资源优化。滑动窗口地址生成器单元的设计如图 5 所示,设定滑动窗口大小为 4 乘以图像宽度,因为采用的卷积核尺寸为 3×3,故系统启动时可采用初始化 3 乘以图像宽度初始化缓冲区,类似于乒乓(Ping-pong)操作,在最大化利用缓冲区以及提升数据处理速度的同时,在 3×3 卷积核进行滑动时,可缓冲覆盖 4 行中的剩余一行,缓冲区即为双口随机存取存储器可同时输入输出,但需要非常严格的地址输入。故在内存资源优化架构中,到 T5 操作时,即可开启特征图中第一行的覆盖操作,即缓冲新数据进入特征图第一行内,将已经计算过的一行旧数据覆盖,即将数据 0,1,2,3,4,5 覆盖。

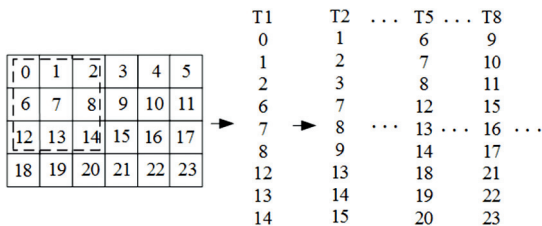


图 5 滑动窗口地址生成图

Fig. 5 Sliding window address generation

3.4 矩阵乘法单元设计

二值化卷积神经网络中的矩阵乘法所用到的硬件逻辑资源非常少,且不同的网络层需定制化矩阵乘法运算单元。本文研究的二值化神经网络在第一层网络层中并未将输入特征图向量进行二值化,在算法层面上将输入特征图向量进行强制类型转换,并

进行定点量化,所以在硬件层面的实现上,仍需要高占比的内存资源进行中间结果的存储,即整个二值化神经网络硬件加速器的内存占用瓶颈在第一层卷积层。因为第一层的权值仍可以通过取权值中单个单指令多数据 (Single Instruction Multiple Data, SIMD) 中的某个值进行权值两比特的量化,求得带符号位的 2 bit 位宽的权值,并与特征图中的定点像素数据进行乘累加运算,即取特征图数据的正数与负数进行累加运算。

而在二值化卷积神经网络中除第一层的其他网络层,均是将权值与输入向量约束在 SIMD 位宽范围内,并以相同的位宽即相同的数据长度进行按位异或并取反,实际上即为二值数据乘累加过程中的同或操作,并使用点累加算法计算乘法过后的无符号数据单比特加法,进行矩阵运算乘累加中的无符号位累加计算。

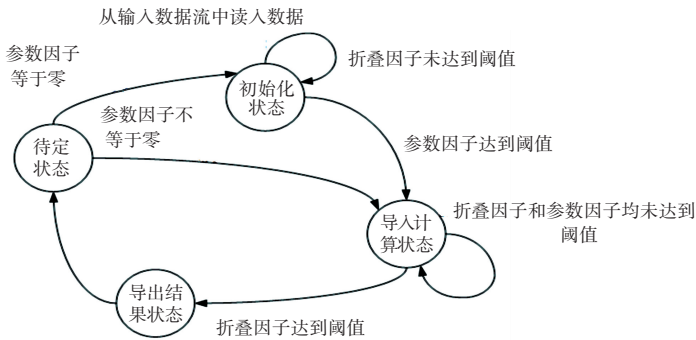


图 6 控制单元设计的状态转移图

Fig. 6 State transition of control unit

4 实现

在电子设计自动化软件 VCS 以及 Verdi 中进行仿真实验,得到仿真时序图如图 7 所示。

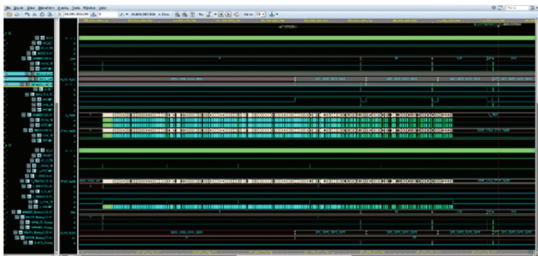


图 7 目标识别硬件加速器输出结果

Fig. 7 Output of object recognition hardware accelerator

当输入 5 张图片时,所消耗的一次性权重数据存储并计算识别得出的时间约为 467 ms,即当进行小批量的图片输入时可达到的识别速率为 10.7 FPS,且此时输入图片的数量很小,故所承担的权重阈值数据导入的比例系数较大,而当输入大批量图片时,速率效果则

3.5 控制单元设计

在整个二值化卷积神经网络硬件电路设计过程中,控制单元的设计尤为关键。数字电路设计中的两大通路,数据通路和控制通路相辅相成,缺一不可。本文深入分析了整个二值化卷积神经网络层的数据通路中数据运算方式以及针对权值与特征图数据所采用的不同的量化方法,但是在硬件底层电路方面,对于例如编码器译码器等组合逻辑电路以及触发器寄存器等时序逻辑电路的精细控制,并未进行深层次的研究。通过较为底层的贴近硬件逻辑器件的硬件描述语言 Verilog/System Verilog 进行算法的设计以及二值化神经网络结构的编写时,需要较多的考虑到数字电路时序逻辑与组合逻辑之间的交互,所以在控制单元硬件设计的过程中,需要使用可紧密联系组合逻辑与时序逻辑的状态转移图表示数字电路设计进程,如图 6 所示。

更为优胜,表明了本文所设计出的目标识别二值化卷积神经网络硬件加速器的性能有较大的提升空间。与此同时,通过脚本大批量测试识别,计算得出识别结果,与 10.7 FPS 识别速率相比较,均有所提高,进一步验证了硬件加速器目标识别的功能正确性。

5 结束语

本文深入研究了二值化卷积神经网络算法模型的数学模型与硬件实现方式,通过设计算法模型中多个模块化的硬件单元并组合,设计了一种中大型尺寸输入特征图像的可综合的目标识别二值化卷积神经网络硬件加速器,达到了较为可观的目标识别加速器性能。

参考文献

[1] COURBARIAUX M, HUBARA I, SOUDRY D, et al. Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1[J]. arXiv preprint arXiv: 1602.02830, 2016.