

集创赛校赛设计报告

集创赛校赛设计报告

赛题四

设计使用软件版本

第一部分 FPGA神经网络加速器原理

- 1.基本原理:
- 2.卷积核结构设计
- 3.数据输入结构设计
- 4.流水线结构
- 5.通道间并行结构

第二部分 三通道3x3卷积核代码设计与仿真

- 1.题目分析:
- 2.半精度浮点数原理:
- 3.乘法器设计:
- 4.加法器设计:
- 5.单通道卷积核设计:
- 6.数据输入模块设计:
- 7.三通道卷积核总体设计:

第三部分 项目文件介绍

- 1.Verilog文件:
- 2.Testbench文件:

第四部分 参考文献与博客

赛题四

使用硬件描述语言如 Verilog、VHDL等，实现三通道3x3卷积内核，具体要求如下：

1. 学习FPGA神经网络加速器的加速原理；
2. 实现一个简单的三通道3x3卷积内核，并且使用modelsim、iverilog、verilator等其中一种仿真工具，生成波形图等仿真文件，以展示相应效果；
3. 提高部分:对所设计的内核进行并行加速、切割流水线，应用winograd算法（选做）；

设计使用软件版本

- Quartus Prime 17.1
- Modelsim SE-64 10.1c

第一部分 FPGA神经网络加速器原理

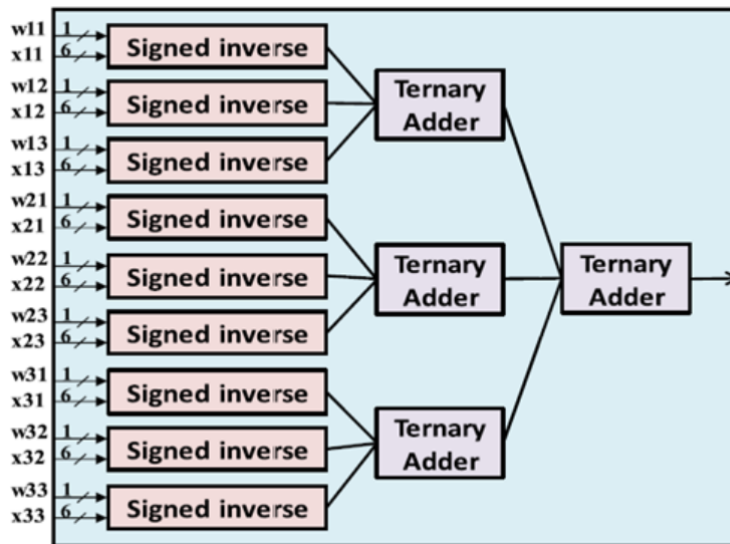
1.基本原理:

神经网络是一种基于大脑神经网络的机器学习模型。一系列节点排列在“层”中，通过操作和权重 相互连接。该模型已证明在图像分类任务中取得了成功，这些任务如今具有许多应用，从自动驾驶汽车到面部识别。标准 CNN 可以具有浮点权重和特征图——这些需要大量的内存资源和计算能力来实现必要的乘法器等。

目前卷积神经网络的FPGA加速研究主要集中在**并行计算**和**内存带宽优化**两方面，其中并行计算主要通过设计卷积层间并行、卷积内计算并行和输出通道并行3种方式来实现加速，此类单纯的硬件并行加速方法资源占用较多、带宽需求较大，实际应用中仍需做相应的改进。内存带宽优化通常采用一些优化算法定量分析计算吞吐量和所需内存的带宽，确定最佳性能进而解决资源占用量大的问题，此类方案在不同层间需要重新配置，灵活性稍显不足。

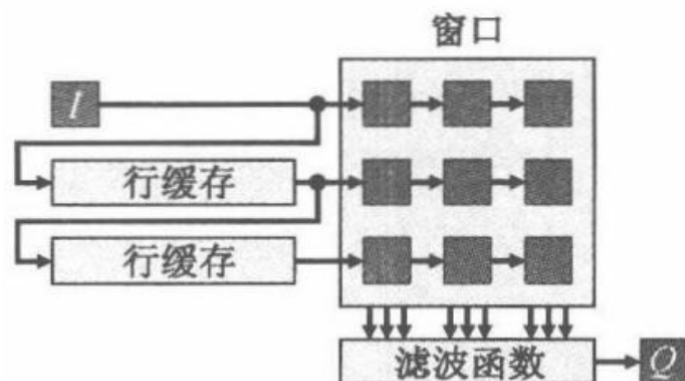
2.卷积核结构设计

基于FPGA的卷积并行加速其实有很多方法，例如**脉动阵列**、**加法树**等操作。本设计中使用基于加法树的并行化设计。其实总体原理也是很简单的。如下图所示，九个叶子节点是乘法器节点，分别代表九次乘法运算（卷积核是 3×3 的）。在得到乘法运算结果之后，将结果传送给加法节点。为了进一步增加并行性，加法树结构采用二叉树。即，对每三个子节点进行求和。最终得到一个部分和。



3.数据输入结构设计

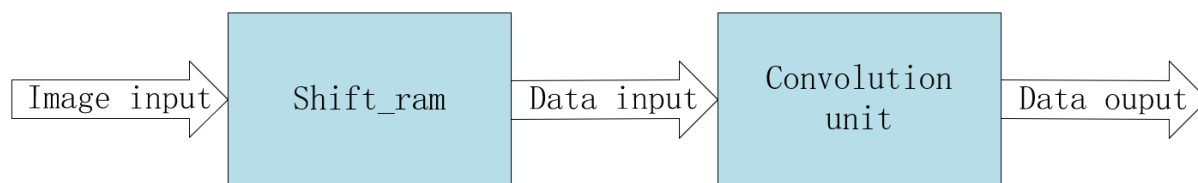
卷积是图像处理中很常见的一种操作， 3×3 是最常见的窗口大小。如果像素是一个一个来的（一个一个输入，可以节省传输带宽），要想实现 3×3 卷积，就得同时获取一个像素和它周围的8个像素，将输入像素缓存2行，这样就能同时获取3行的像素输入，此时再将这3个并行输入的像素移位进 3×3 窗口，就获得了 3×3 卷积模板，如图：



这里要注意，输入像素此时作为第三行数据输入3x3窗口，最下面的行缓存输出的才是第一行像素，上图窗口的右下角是3x3卷积模板的左上角，窗口的左上角是3x3卷积模板的右下角。实现两行缓存并获取3x3卷积窗口，用shift-ram是最简单的实现方法。

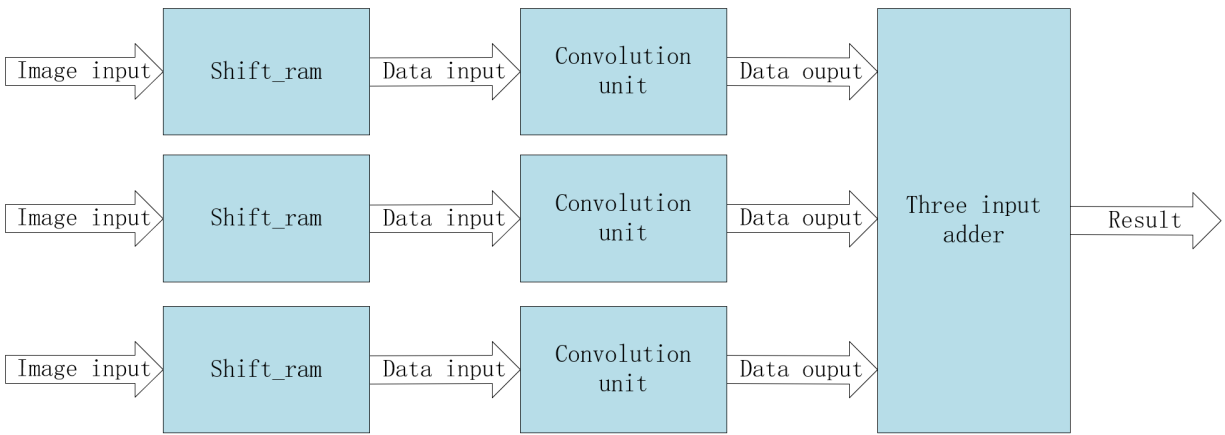
4.流水线结构

所谓流水线设计，实际上就是把规模较大、层次较多的组合逻辑电路分为几级，在每一级插入寄存器组并暂存中间数据。K级流水线就是从组合逻辑的输入到输出恰好有K个寄存器组，上一级的输出是下一级的输入而又无反馈的电路。在本设计中，可将数据输入和卷积核计算分作两级流水线，也可将卷积核内部的树形计算结构分作三级流水线。通过这种流水线的结构，可以做到**1个时钟周期**即可输出1次卷积结果，大大提高了卷积运算的效率。



5.通道间并行结构

由于三个通道的计算方式基本相同，因此考虑将三通道卷积进行并行加速，并行加速结果见下图。使用下图中的结构，可以将运算速度提升三倍（代价是占用了更多的硬件资源）。需要指出的是不仅是通道之间可以并行加速，事实上也可以在一副图像中使用多个卷积核来实现并行加速。



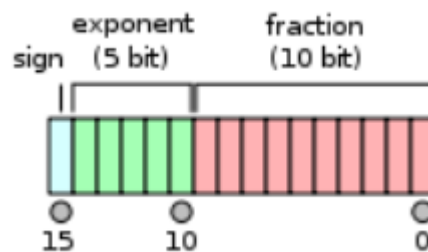
第二部分 三通道3x3卷积核代码设计与仿真

1.题目分析:

在进行卷积核设计前，首先要确定单个数据的数据格式，本设计中指定为16位浮点数（半精度浮点数）。其次考虑到卷积运算需要包含乘法和加法，因此需要设计乘法器和加法器两种电路。通常，卷积单元设计为单通道设计，考虑到题目中三通道的要求，本设计中，采用3个单通道加

2.半精度浮点数原理:

IEEE754-2008包含一种“半精度”格式，只有16位宽。故它又被称之为binary16。与单精度浮点数相比，它的优点是只需要一半的存储空间和带宽，但是缺点是精度较低。其结构图如下：



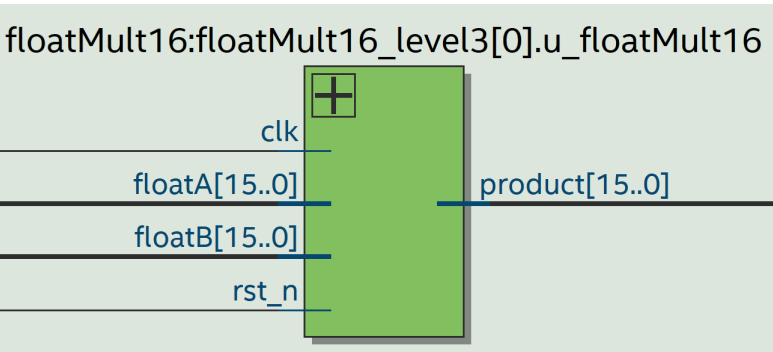
半精度的格式与单精度的格式类似，最左边的一位仍是符号位，指数有5位宽且以余-16 (excess-16) 的形式存储，尾数有10位宽，但具有隐含1。半精度浮点数的值的计算方式为 $(-1)^{\text{sign}} \times 2^{\text{(指数位的值)}} \times (1 + 0.\text{尾数位})$ 。

3.乘法器设计:

• 半精度浮点乘法的计算步骤如下:

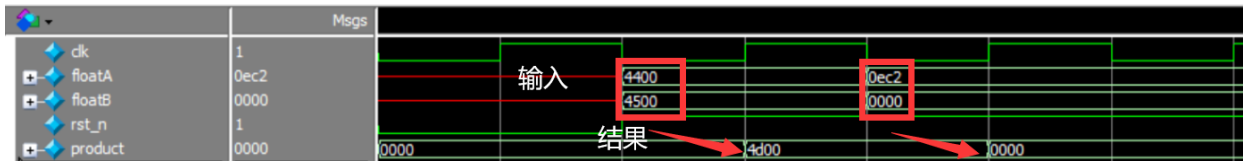
1. 首先将16位数据拆解为符号位，指数位和尾数位。之后判断是否有特殊情况，如两个乘数中有0，若有则按可直接写出结果。
2. 其次，针对一般情况下的运算，符号位为只是两乘数符号位的异或，指数位是两乘数指数的相加，尾数位则在正常的乘法之余考虑到移位（和随之而来的指数为的溢出）和进位的情况。

• 半精度浮点乘法RTL视图如下：



设计中为了提高组合电路的时序性能，在乘法器输入端加入时钟端口（clk）和异步清零端口（rst_n）。

• 半精度浮点乘法器ModelSim仿真波形如下：



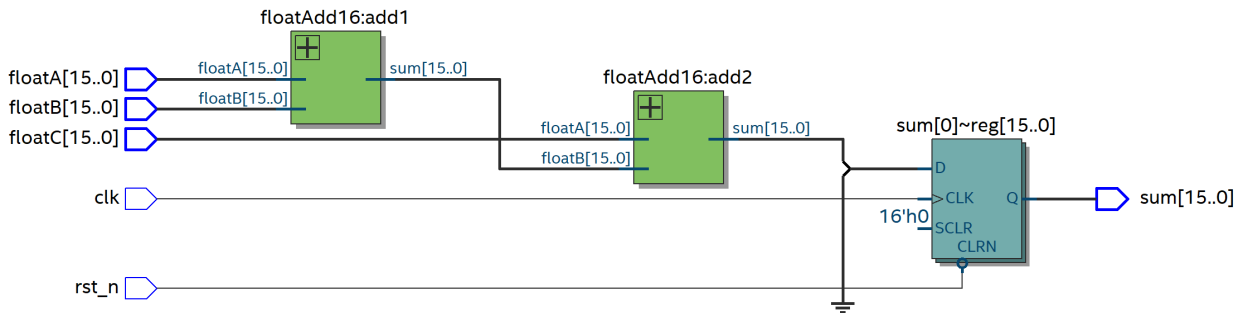
在该测试中，分别测试了4*5（即0x4400 * 0x4500）和0.0004125*0即（0x0ec2*0x0000），输出结果为20（即0x4d00）和0（即 0x0000）。测试结果正确，表明设计的电路功能正常。

4.加法器设计：

• 半精度浮点加法的计算步骤如下：

1. 首先将16位数据拆解为符号位，指数位和尾数位。之后判断是否有特殊情况，如两个加数中有0，或互为相反数，若有则按可直接写出结果。
2. 其次，针对一般情况下的运算，先进行对阶，使得两个加数的阶数相同。当阶数相同后，即可进行尾数相加，得到加法结果。此外，当出现尾数溢出的情况时，可以采用右移一位，阶码减一的方式来判断阶码是否溢出，当阶码溢出时表示两数相加的结果溢出。

• 半精度三输入加法器RTL视图如下：



设计中为了方便后续对加法的并行处理，使用两个两输入加法器级联的形式构建了三输入加法器。同时为了提高组合电路的时序性能， 添加了时钟端口和复位端口。

- 半精度浮点三输入加法器ModelSim仿真波形如下:

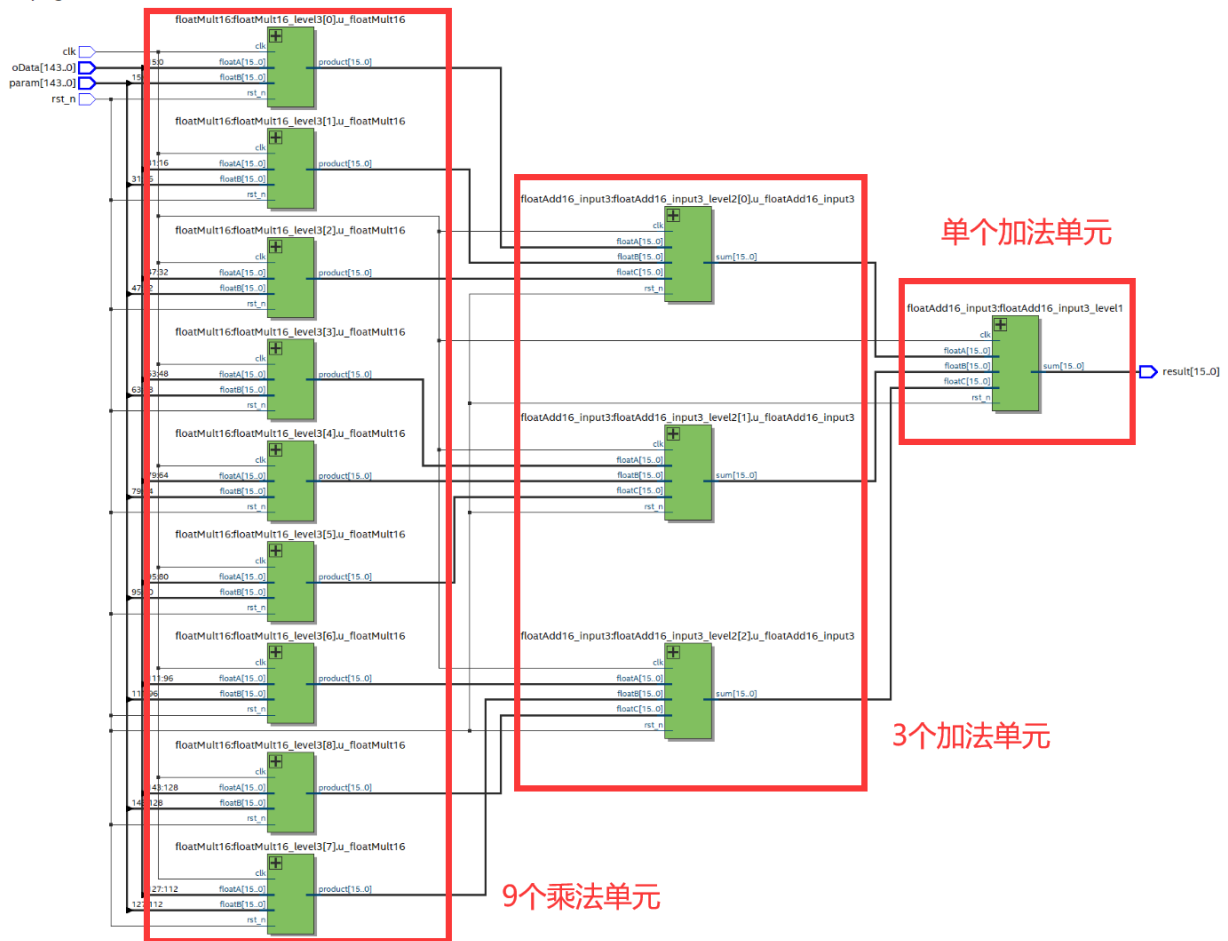


在该测试中，分别测试了0.2+0.3+0（即0x34cd+0x3266+0x0000）和0.2+0+0（即0x34cd+0x0000+0x0000=0x34cd），输出结果为0.5（即0x3800）和0.2（即0x34cd）。测试结果正确，表明设计的电路功能正常。

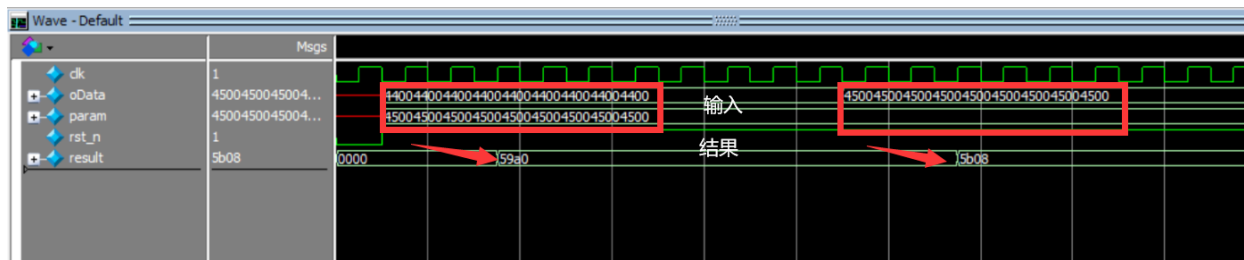
5.单通道卷积核设计:

在卷积核计算中,需要进行9次乘法和8次加法。如果全部依次执行,需要17个时钟周期,降低了计算效率。为了提高运算速度,本设计中采用了加法树进行并行加速。通过使用设计的两输入乘法器和三输入加法器,通过3个时钟周期即可得出卷积结果。考虑到流水线操作,事实上,1个时钟周期,即可得出运算结果。具体的加法树结构分为三层,第三层为9个两输入加法器,第二层为3个三输入加法器,第一层为1个三输入加法器。

- 单通道卷积核RTL视图如下:



- 单通道卷积核ModelSim仿真波形如下：



在该测试中，第一次输入数据4 4 4 4 4 4 4 4和5 5 5 5 5 5 5 5（即0x44004400440044004400440044004400和0x45004500450045004500450045004500）。经过三个时钟周期后，输出结果180即（0x59a0）。第二次输入数据5 5 5 5 5 5 5 5和5 5 5 5 5 5 5 5（即0x45004500450045004500450045004500和0x45004500450045004500450045004500）。经过三个时钟周期后，输出结果225即（0x5b08）。

6.数据输入模块设计：

要想实现3x3卷积，就得同时获取一个像素和它周围的8个像素，将输入像素缓存2行，这样就能同时获取3行的像素输入，此时再将这3个并行输入的像素移位进3x3窗口，就获得所需数据。本设计中，采用Quartus的shift_ram IP核实现移位寄存器。设置移位寄存器宽度为图像宽度（本例中设置为8），设置数据宽度为16bit，设置寄存器行数为两行。

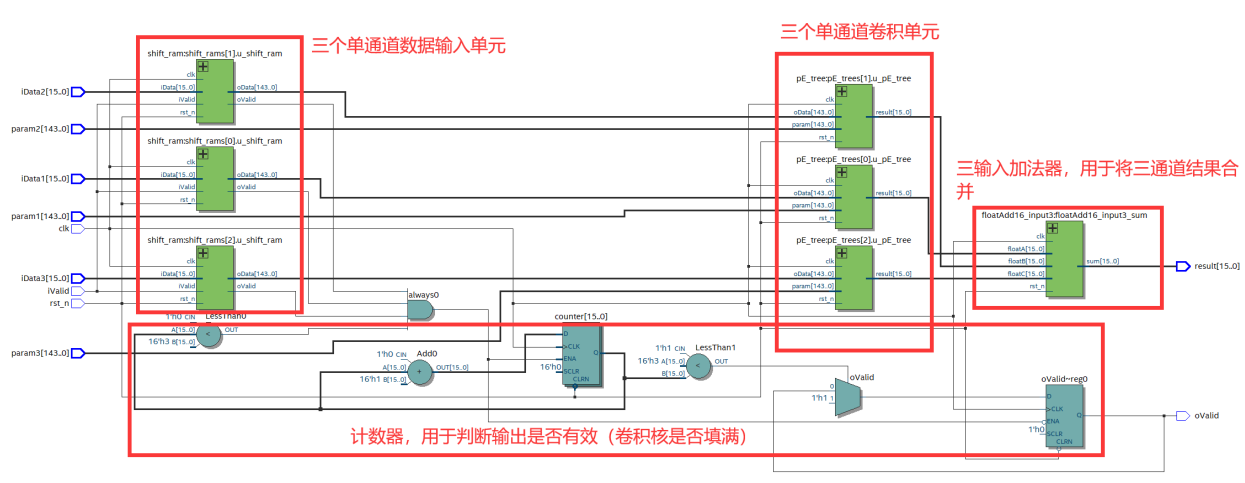
-

- 数据输入模块ModelSim仿真波形如下:



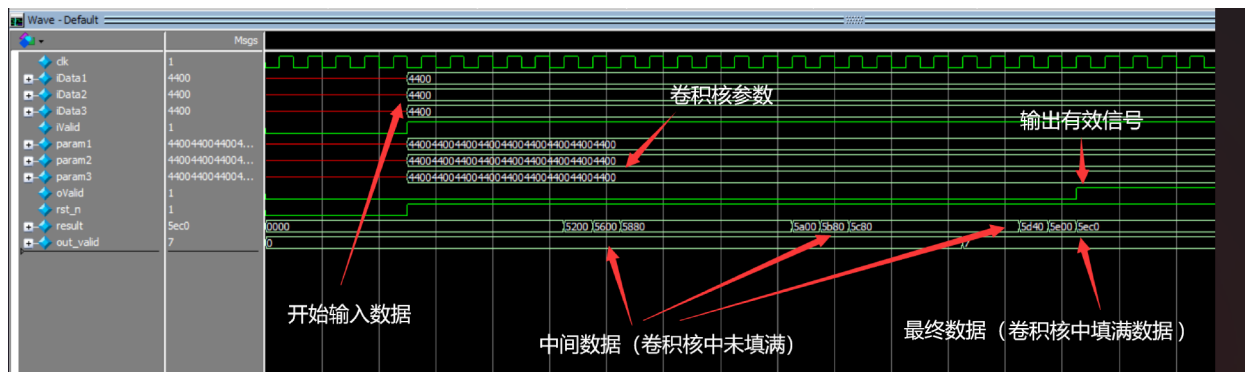
7.三通道卷积核总体设计:

- **三通道卷积核RTL视图如下:**



本设计中，三通道数据通过三个输入端口（iData1~iData3）输入，三通道卷积核参数通过三个输入端口（param1~param3）输入。result为卷积输出结果。oValid信号置1时，输出结果有效。

• 三通道卷积核ModelSim仿真图测试波形如下：

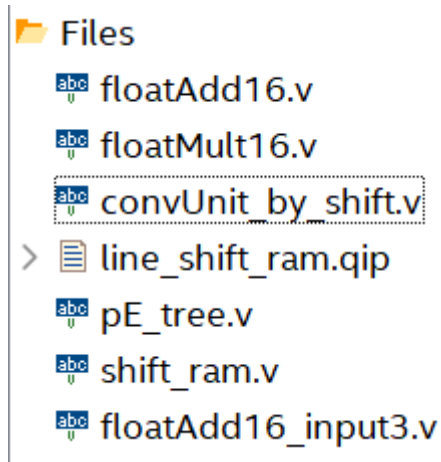


本测试中，输入三通道数据均为4（即0x4400），卷积核参数也均为4（即0x4400）。测试输出应为 $4 \times 4 \times 3 \times 3 = 144$ （即0x5ec0）。观察测试波形，在oValid置1后，result输出144（即0x5ec0），与预期一致，证明本设计功能正常，能够完成三通道3x3卷积。

第三部分 项目文件介绍

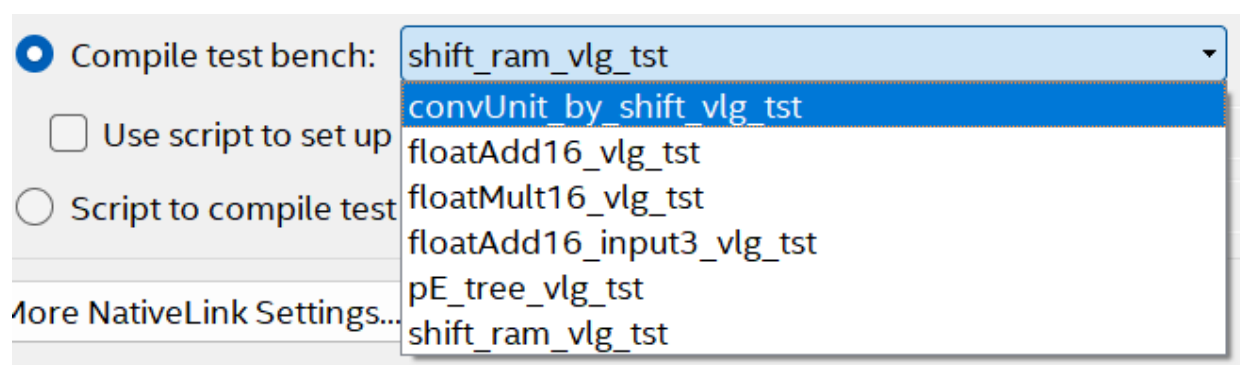
1.Verilog文件：

全部文件如下图。其中floatAdd16为两输入16位浮点加法器（组合电路），floatMult16为两输入16位浮点乘法器（时序电路），floatAdd16_input3为三输入16位浮点加法器（时序电路），pE_tree为单通道3x3卷积单元，line_shift_ram为Quartus IP核（shift_ram）文件，shift_ram为数据输入电路（由移位寄存器实现），convUnit_by_shift为三通道3x3卷积核电路（顶层电路）。



2.Testbench文件:

下面为本设计中使用到的Testbench仿真文件。



第四部分 参考文献与博客

[1]龚豪杰,周海,冯水春.基于FPGA的卷积神经网络并行加速设计[J].计算机工程与设计,2022,43(07):1872-1878.DOI:10.16208/j.issn1000-7024.2022.07.010.

[2]郑俊伟. 基于FPGA的卷积神经网络并行加速设计研究[D].西安电子科技大学,2021.DOI:10.27389/d.cnki.gxadu.2021.000844.

[3]刘志成,祝永新,汪辉等.基于FPGA的卷积神经网络并行加速结构设计[J].微电子学与计算机,2018,35(10):80-84.DOI:10.19304/j.cnki.issn1000-7180.2018.10.016.

[\(167条消息\) FPGA图像处理-3x3卷积模板_图像卷积模板学习就van事了的博客-CSDN博客](#)

[\(167条消息\) 一起学习用Verilog在FPGA上实现CNN----\(二\)卷积层设计verilog实现卷积鲁棒最小二乘支持向量的博客-CSDN博客](#)