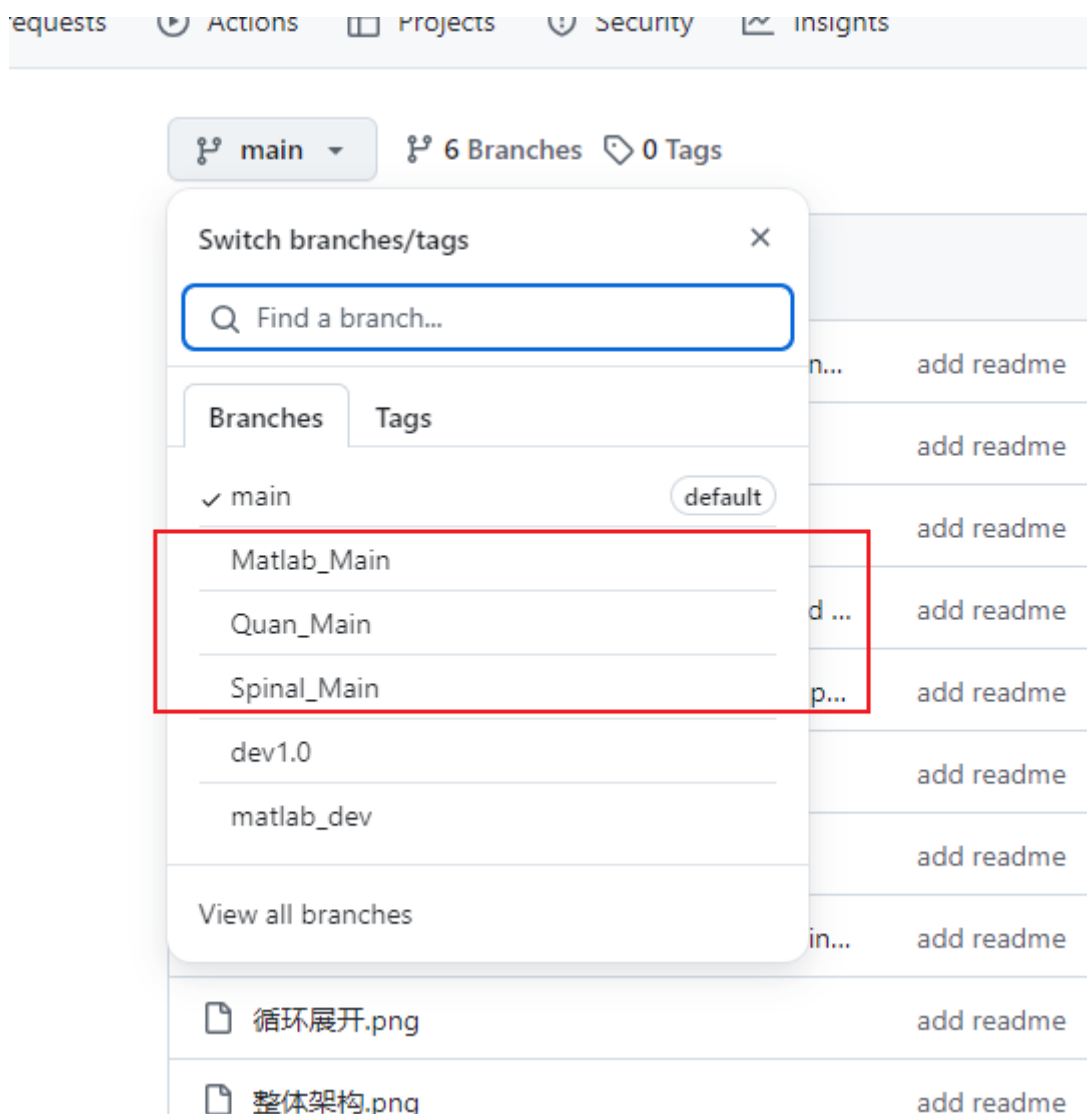


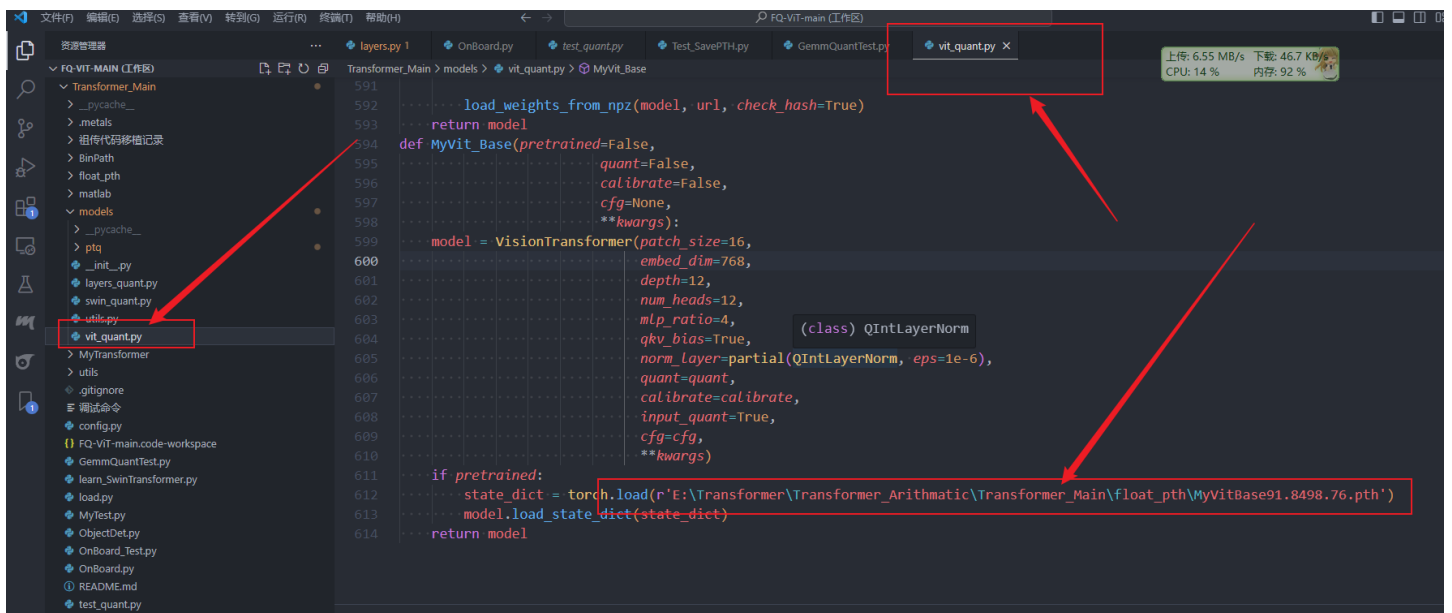
# 上板测试流程（仿真流程也差不多）

git仓库: <https://github.com/Simon-K1/DETR/tree/main>(代尧+翟锐凯)

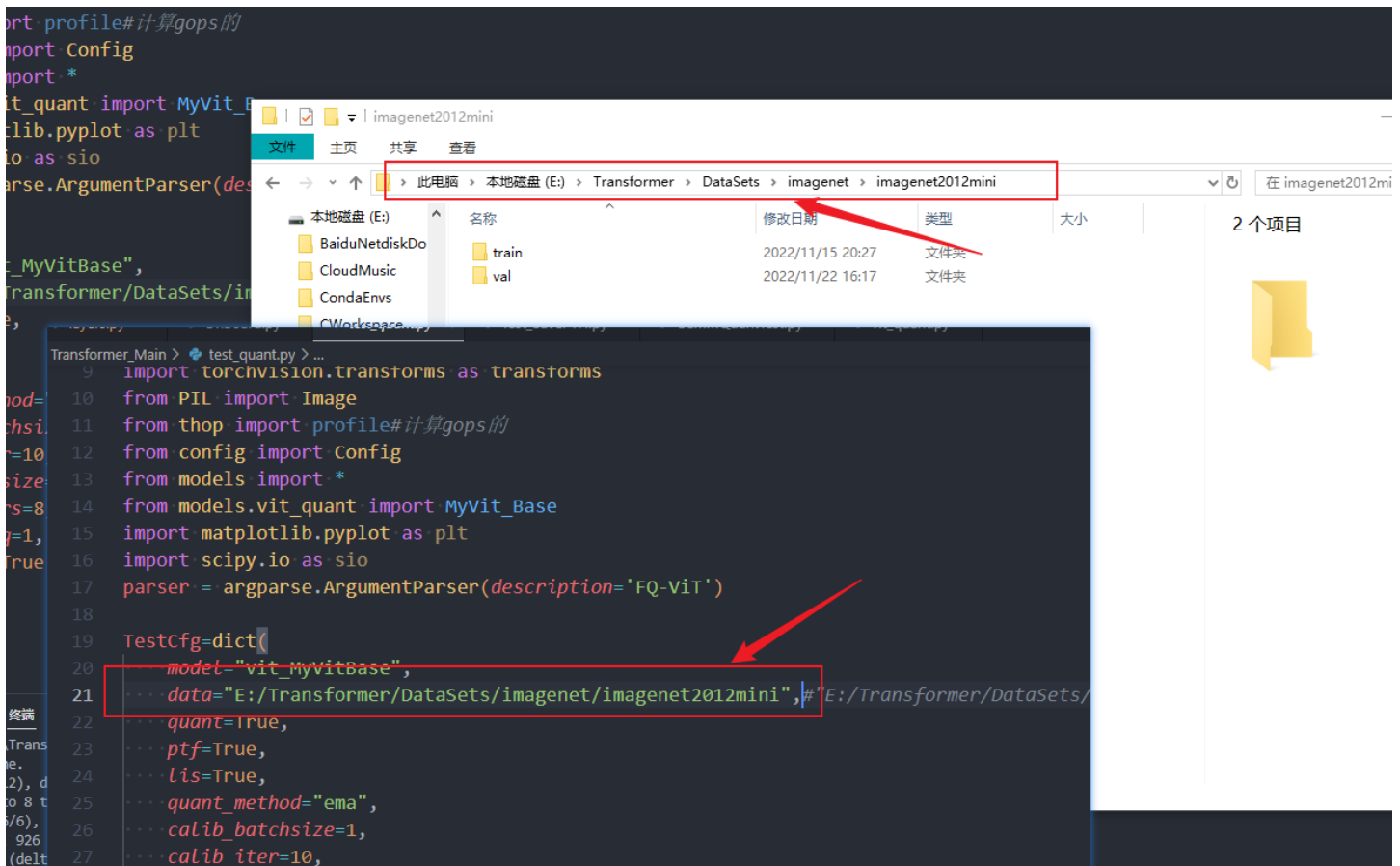
需要Clone下来的分支分别为Spinal\_Main,Quan\_Main和Matlab\_Main



(创建三个文件夹，分别把这三个分支Clone下来，最新代码都是能跑的，Quan\_Main跑不通是因为需要数据集和训练好的模型来跑，数据集和模型放在下面链接，这里面的模型是我自己随便训的，只能分类三个类，需要自己训练就需要跑train.py的代码)



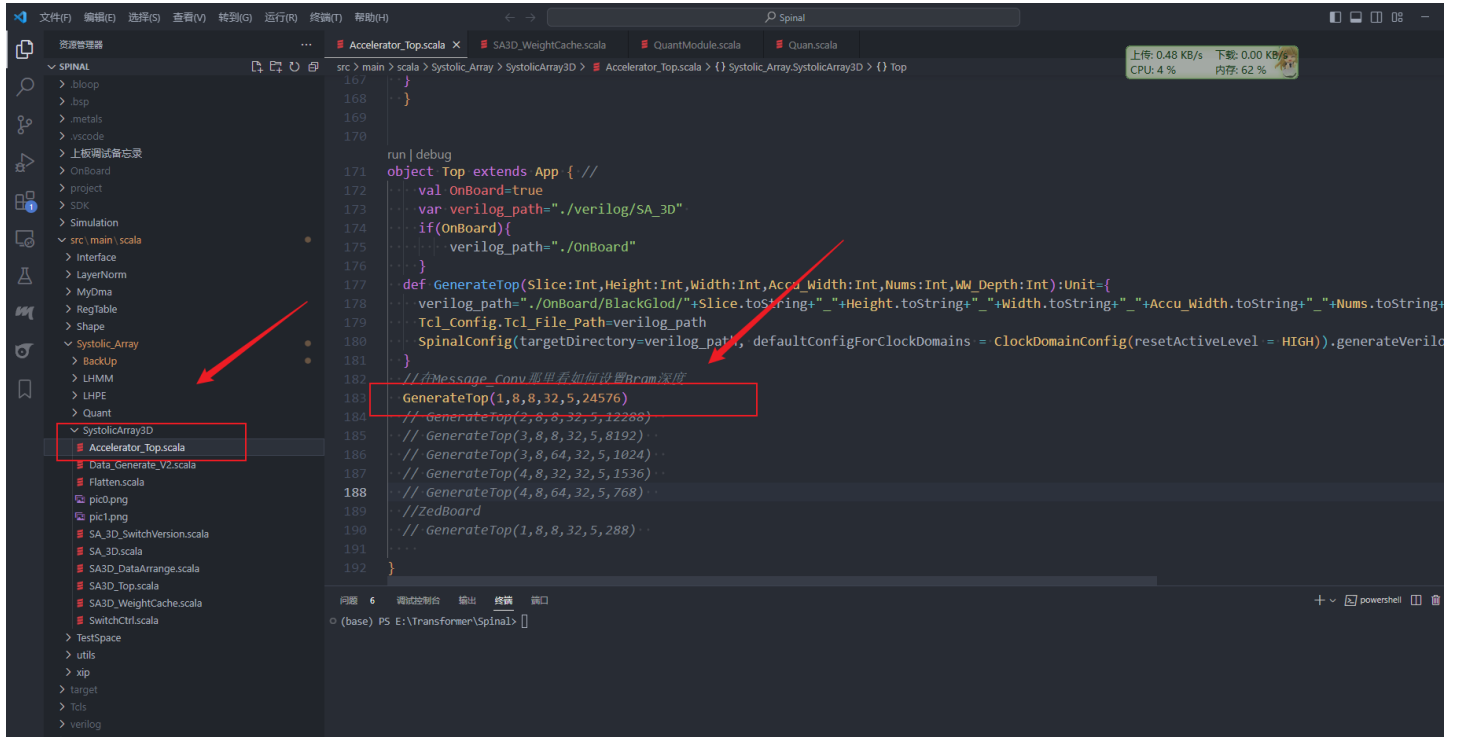
数据集：解压即用，在test\_quant的cfg中改一下路径就行



## 第一步：生成RTL

分支：Spinal\_V3--目前进度：卷积量化和矩阵量化均实现了，剩下的就是无穷无尽的测试和调试了。。。

SpinalHDL现在可以一键生成各种配置的加速器，现在3\*8\*64的上面验证通过了卷积+完整量化，1\*8\*8上面验证通过了矩阵计算+完整量化。



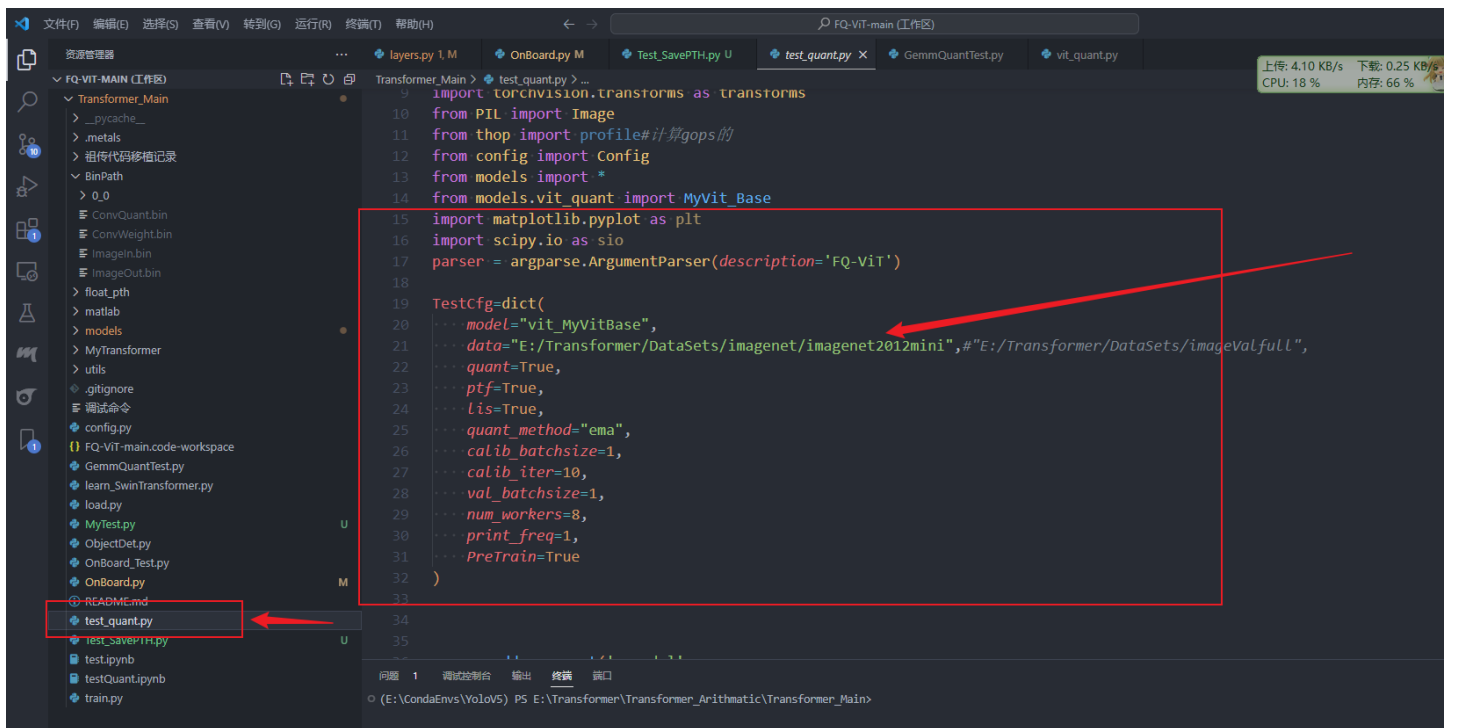
[Tile, Height, Width, 32(累加和位宽, 不要改) , 5 (不要改) , 权重缓存Bram的深度(位宽为64bit)]

拿到RTL代码后自己搭一下BD, 导入TCL (这部分待更新)

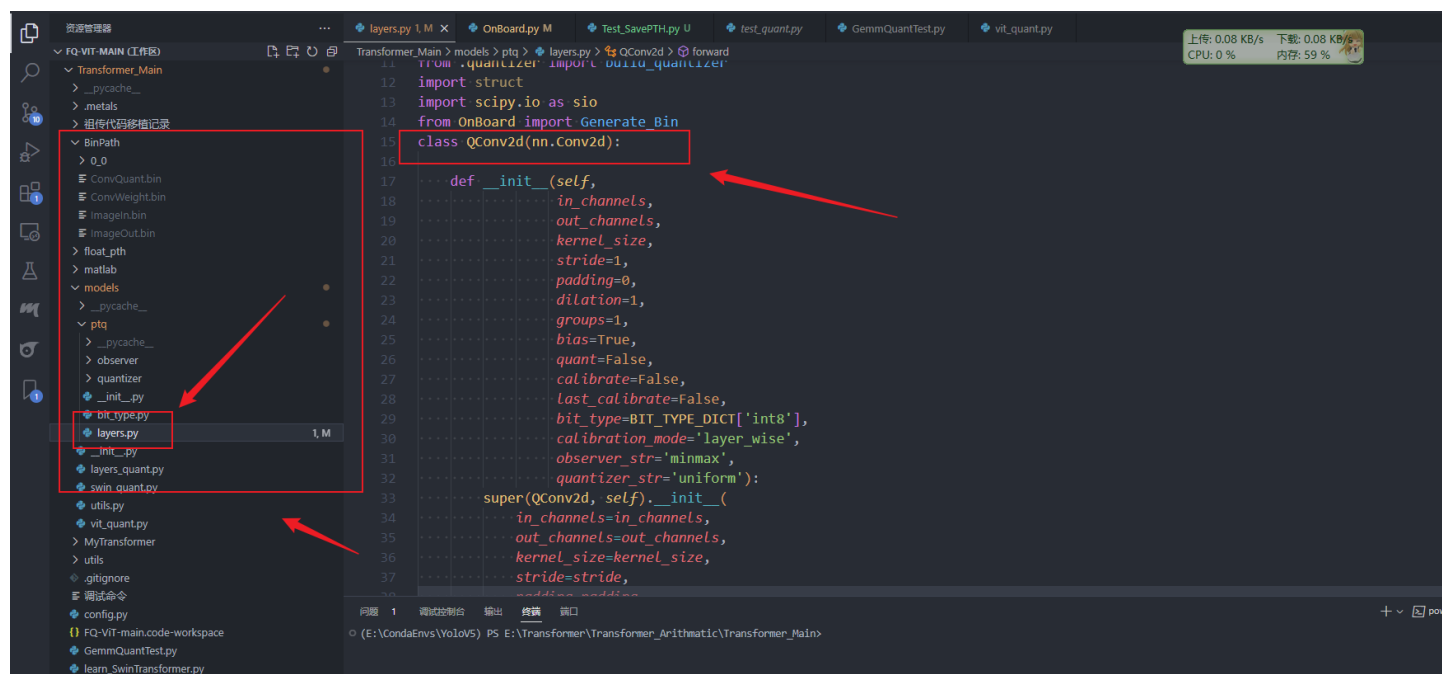
## 第二步：Pytorch生成真实的权重，图片和量化参数

分支：Quan\_v6

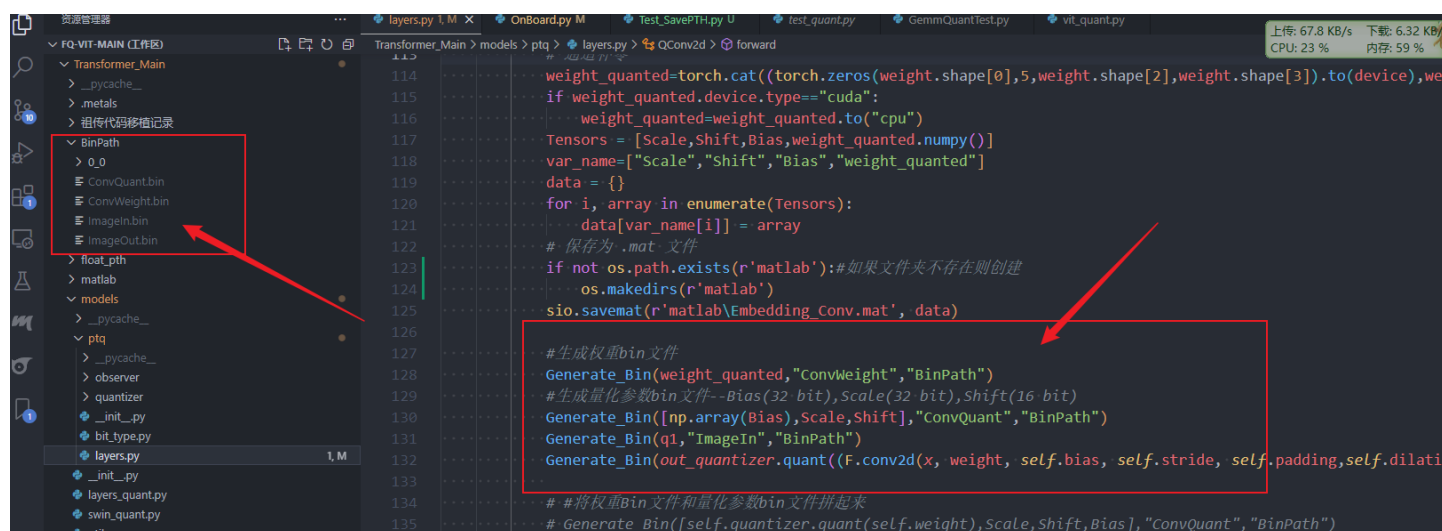
在test\_Quant.py中配置好图片路径等参数



然后找到QConv2这个类，断点调试打到这里，祖传量化已经被移植到这里了，校准完生成Scale, Zeropoint后，程序运行到这里会自动生成权重，图片和量化参数和输出图片共四个bin文件。



下面这四行代码用于生成bin文件，生成完bin文件后就可以退出调试

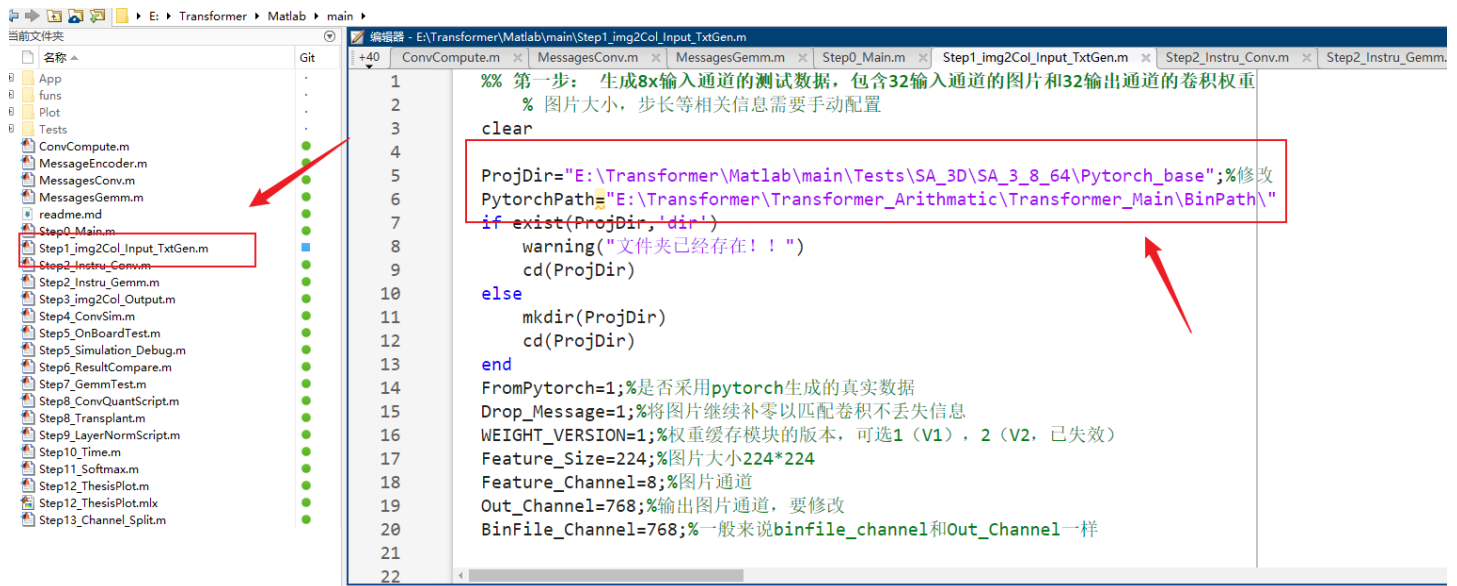


## 第三步：打开Matlab，用Matlab读取Pytorch生成的真实数据并生成指令

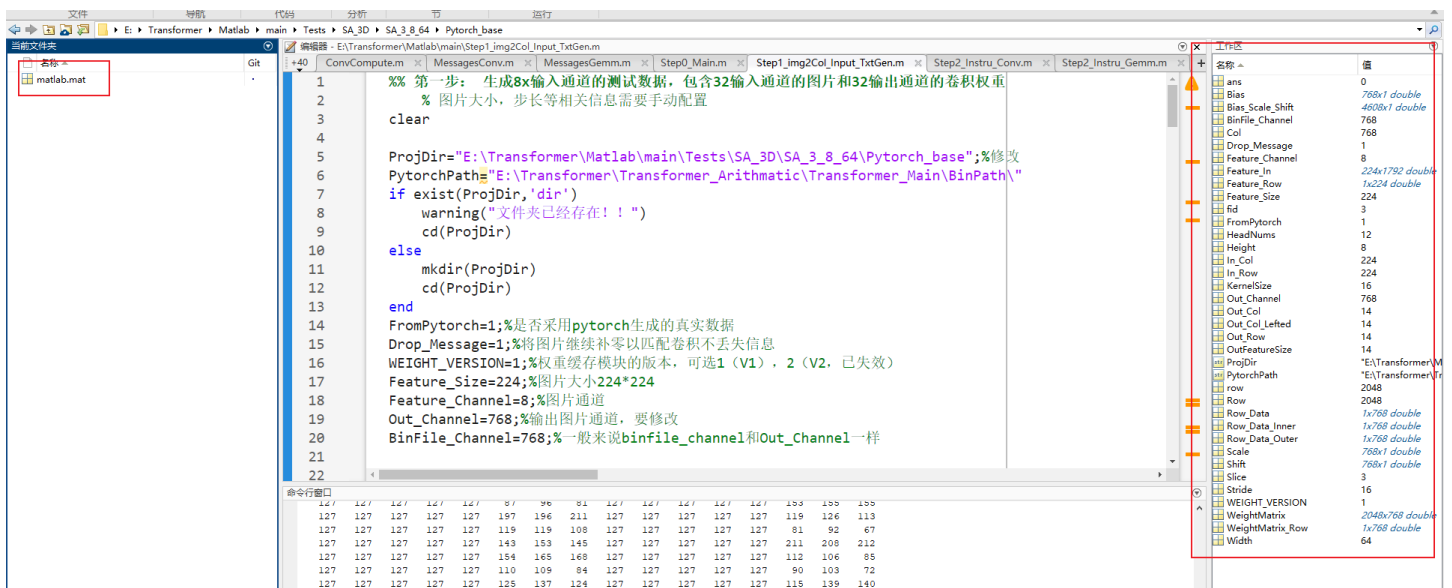
如果要上板测的话，就几步而已

## 进入matlab的Step1，

- 配置好自己的ProjDir，以及Pytorch生成的bin文件目录
- FromPytorch选择1
- Drop\_Message选择1，不用管他，用来padding的，现在还没实现padding，不理他
- WEIGHT\_VERSION:默认1，因为权重有两种数据排列格式：行优先和列优先，现在我们用的全是列优先
- Feature\_Size;与Pytorch对应，Pytorch是224的图片
- Feature\_Channel:图片通道，默认8通道，Pytorch那虽然图片是3通道，但是我在生成bin文件的时候补成了8通道
- Out\_Channel：这个得和Vit-Base对应，768输出通道
- BinFile\_Channel：Pytorch生成的BinFile中卷积是多少通道就填多少通道。
- 剩下的待补充



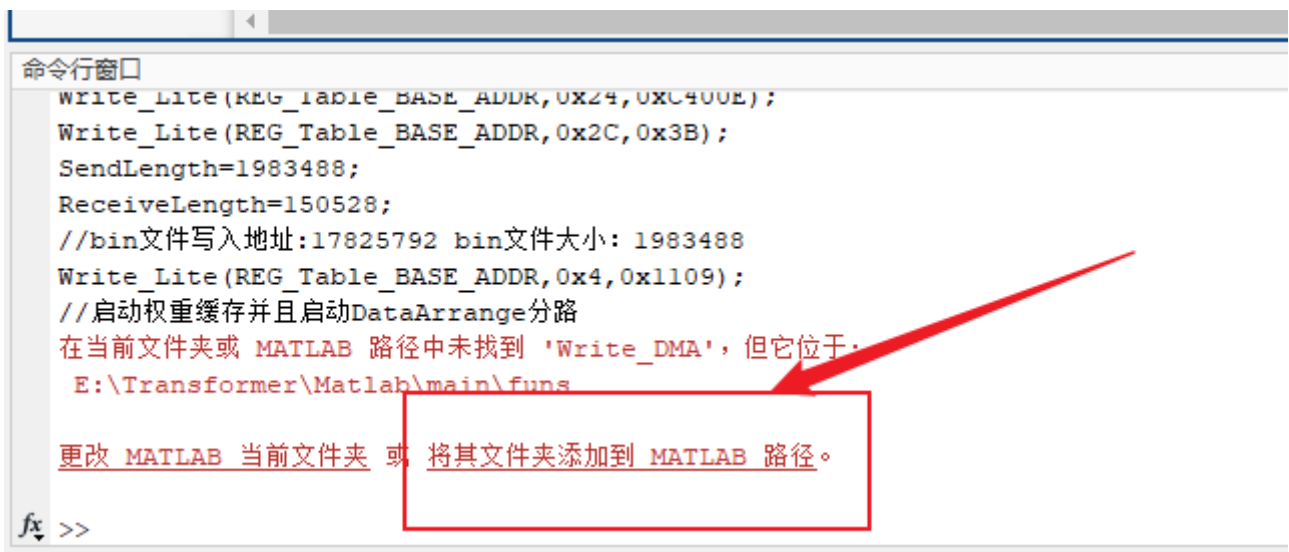
全部配好后直接运行Step1，自动创建包含了真实数据的全部变量，保存在matlab.mat文件中，以后只需要读这个matlab.mat文件即可



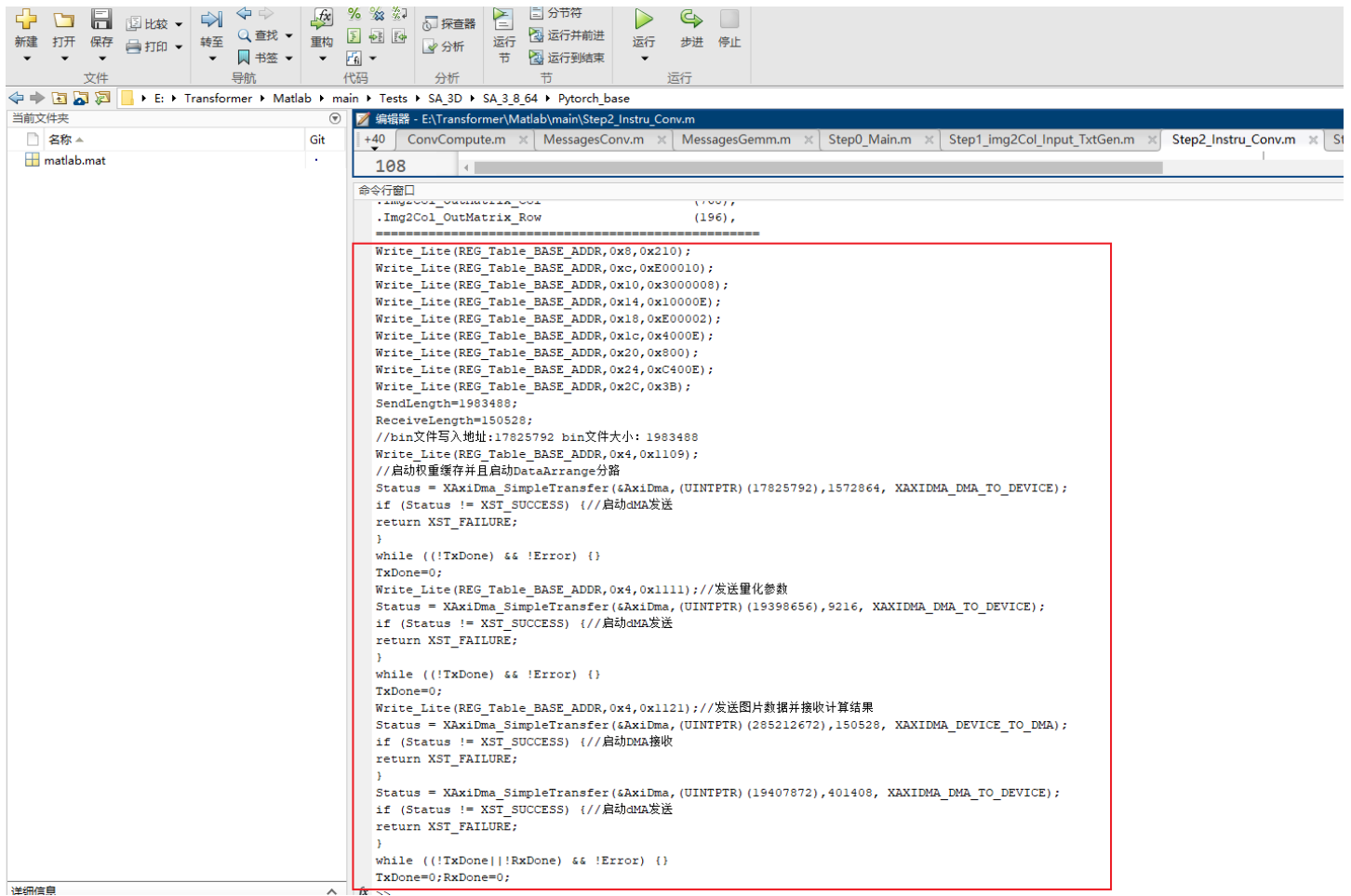
## 进入Step2\_Instru\_Conv

Step2用来生成指令，应该没啥问题，测了好多遍了

如果报这个错，是因为funs这个文件夹没有被添加到搜索路径，将其添加一下就好再重新运行即可

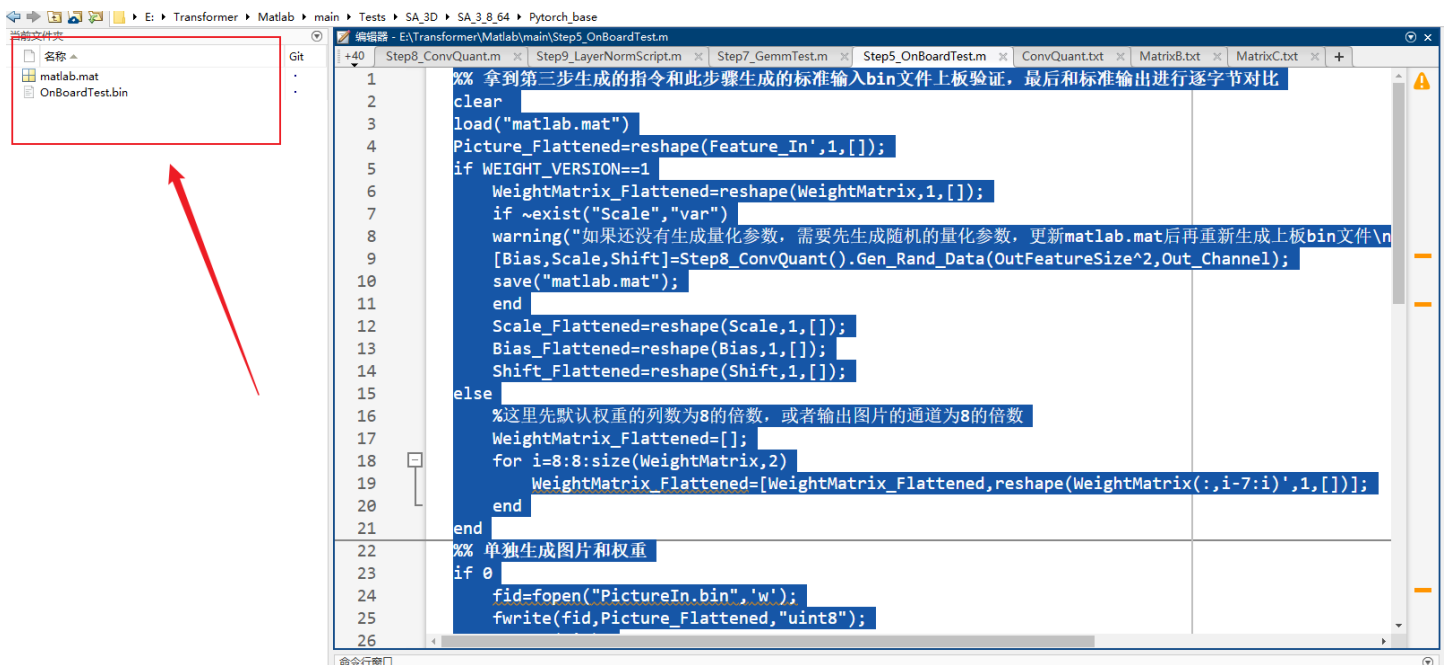


跑完后生成上板测试指令和仿真指令，到时候复制粘贴就行



## 进入Step5-OnBoadeTest

运行即可一键生成上板测试bin文件，然后将其导入到Zynq的ddr里即可

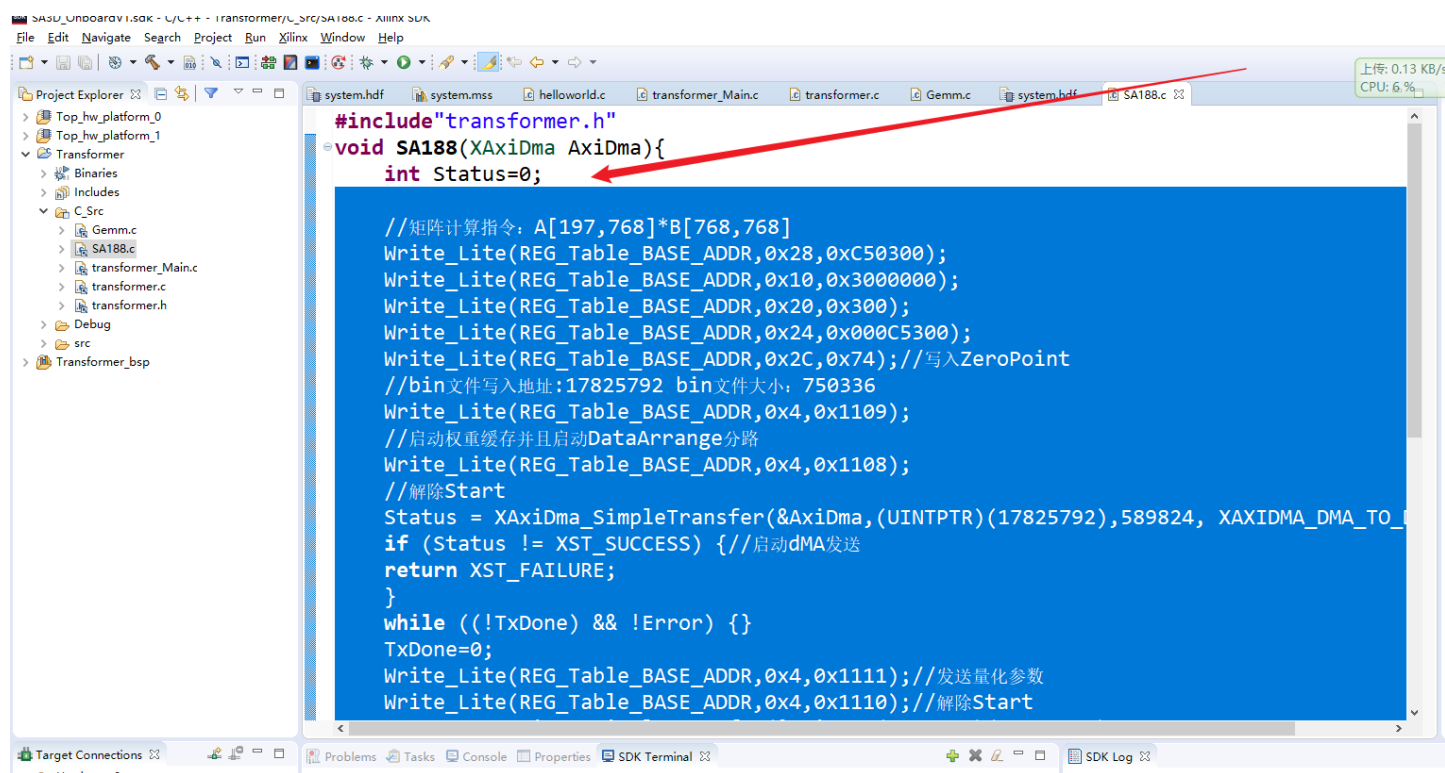


## 第四步：打开SDK

SDK那边就4文件，需要修改的文件就是SA188.c或者Gemm.c（这种文件可以自己创建多个，再在transforme\_main中调用即可。



把刚刚Matlab生成的C代码复制到SA188.c中（SA188.c只是一个源文件，可以自己创建任意名称的文件并include"transformer.h"），然后再在transformer\_Main中调用即可。



SA188.c

transformer\_Main.c

transformer.h

transformer.c

## 第五步：写入Bin文件

写入生成好的OnBoardTest.bin文件，然后跑一把看看结果。

Bin文件的写入数据量和写入地址在刚刚复制的C代码中已经标出来了



```

#include"transformer.h"
void SA188(XAxiDma AxiDma){
    int Status=0;

    //矩阵计算指令: A[197,768]*B[768,768]
    Write_Lite(REG_Table_BASE_ADDR,0x28,0xC50300);
    Write_Lite(REG_Table_BASE_ADDR,0x10,0x3000000);
    Write_Lite(REG_Table_BASE_ADDR,0x20,0x300);
    Write_Lite(REG_Table_BASE_ADDR,0x24,0x000C5300);
    Write_Lite(REG_Table_BASE_ADDR,0x2C,0x74);//写入ZeroPoint
    //bin文件写入地址:17825792 bin文件大小: 750336
    Write_Lite(REG_Table_BASE_ADDR,0x4,0x1109);
    //启动权重缓存并且启动DataArrange分路
    Write_Lite(REG_Table_BASE_ADDR,0x4,0x1108);

```

## 第六步：FPGA算完后导出bin文件

从哪导出导出多少数据量，具体需要查看最后的DMA接收的那行代码中的数据：

```

ain.c
    if (Status != XST_SUCCESS) { // 启动DMA发送
        return XST_FAILURE;
    }
    while ((!TxDone) && !Error) {
        TxDone=0;
        Write_Lite(REG_Table_BASE_ADDR,0x4,0x1181); //发送图片数据并接收计算结果
        Write_Lite(REG_Table_BASE_ADDR,0x4,0x1180); //解除Start
        Status = XAxiDma_SimpleTransfer(&AxiDma, (UINTPTR)(285212672), 151296, XAXIDMA_DEVICE);
        if (Status != XST_SUCCESS) { //启动DMA接收
            return XST_FAILURE;
        }
    }
}

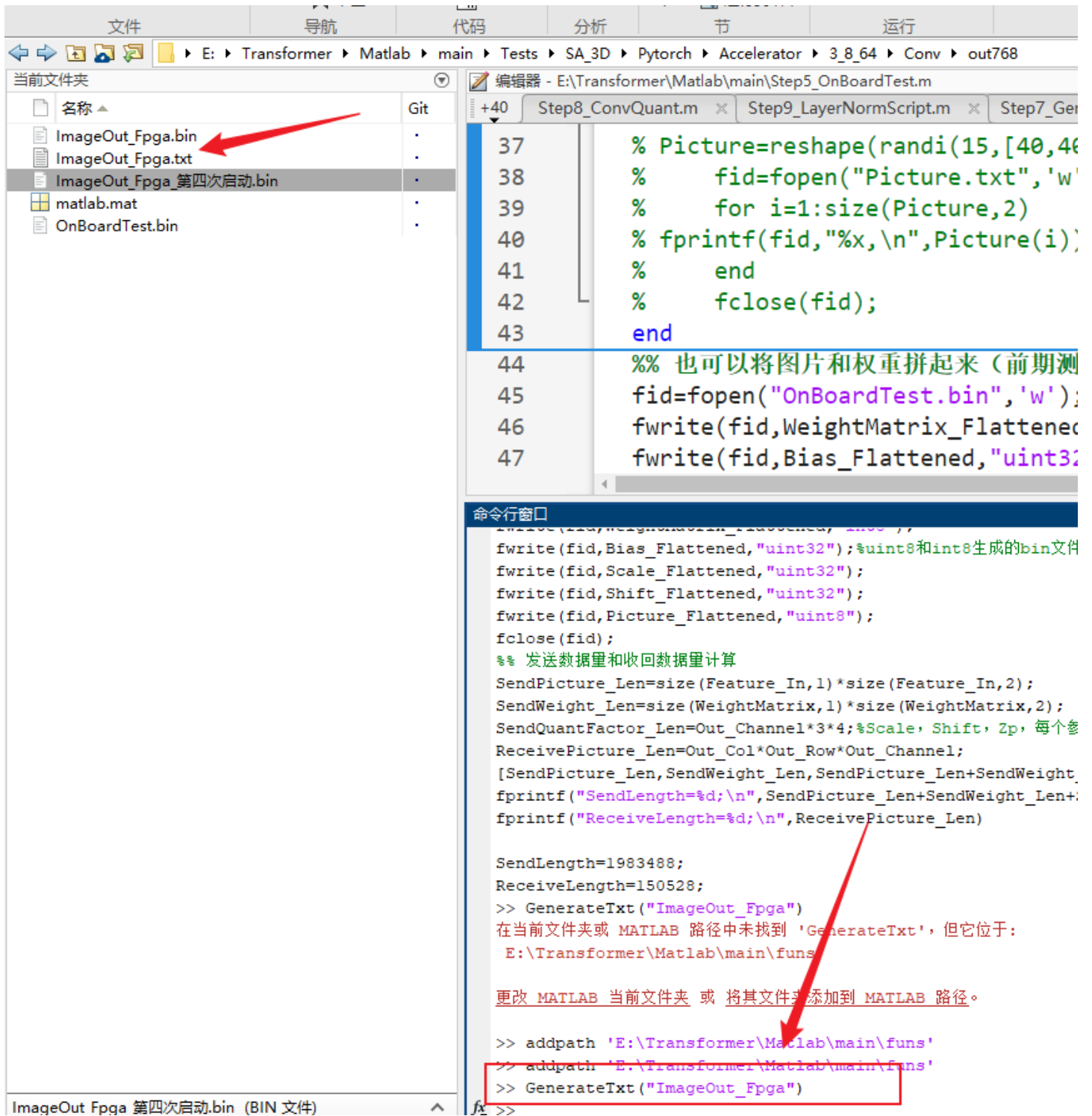
```

接收地址      接收数据量

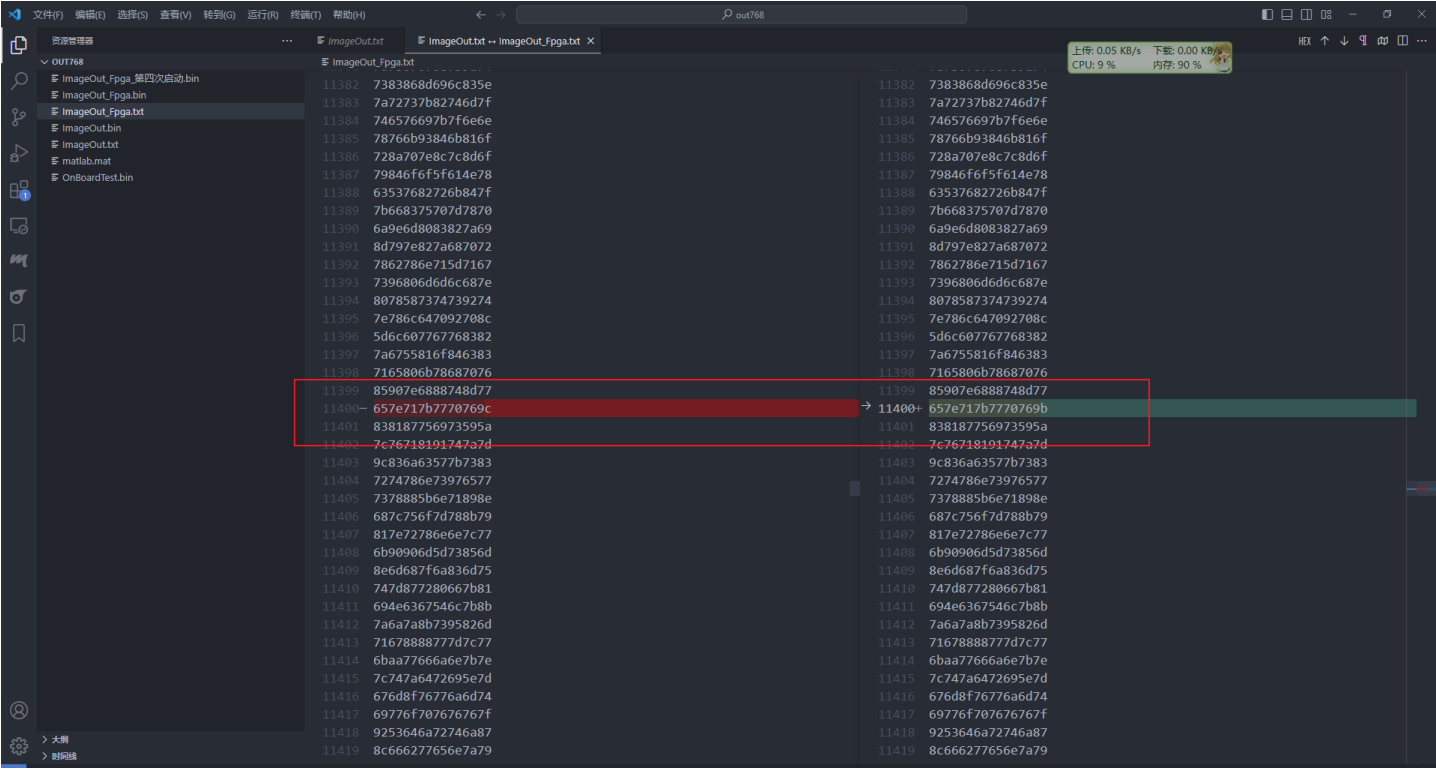
## 第七步：拿导出的bin文件和Pytorch的输出进行逐一字节的对比

很多方法，这里可以用matlab生成txt然后对比：

命令行中输入GenerateTxt("ImageOut\_Fpga")即可自动生成对应bin文件的txt文件



同样地，将Pytorch生成的输出图片的bin文件也转成txt然后用vscode对比即可：



最后就1bit没对上。