# How to Prepare and Run NPU System Demo (CPU + FPGA)

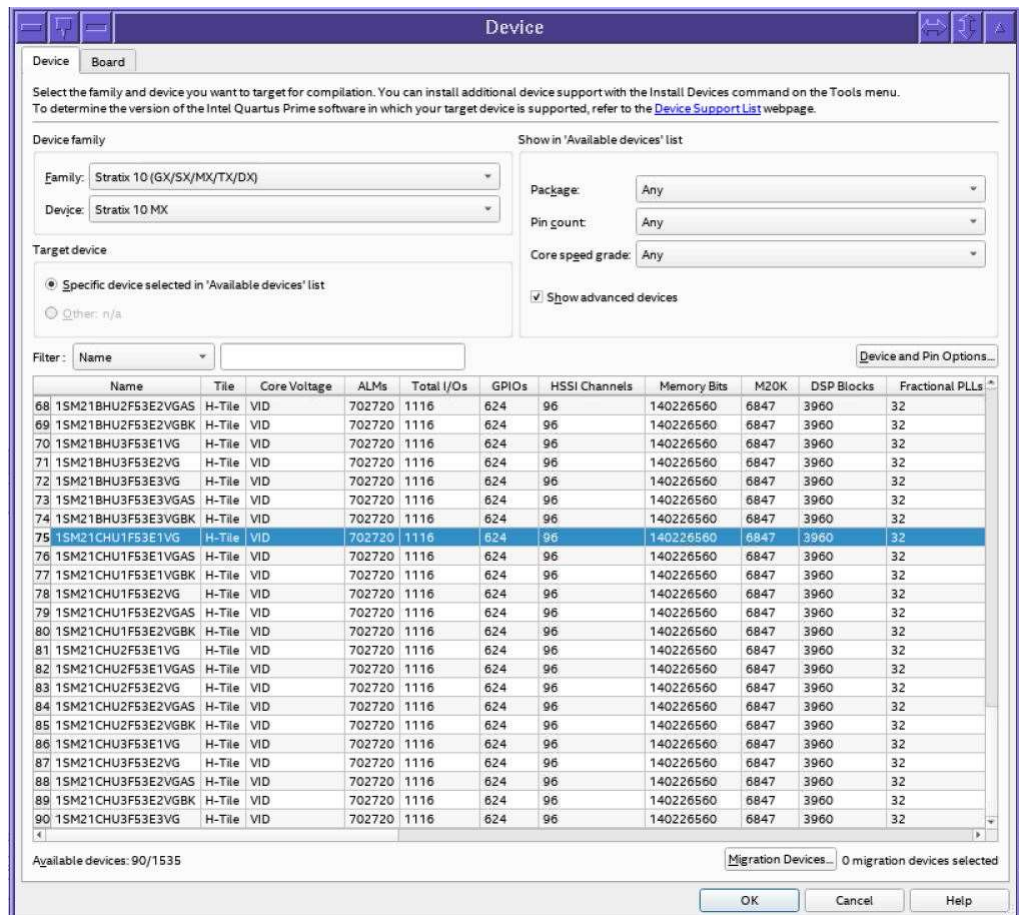## 1. Start in the directory where the zip file has been unpacked
1) We call this directory the <project home> directory
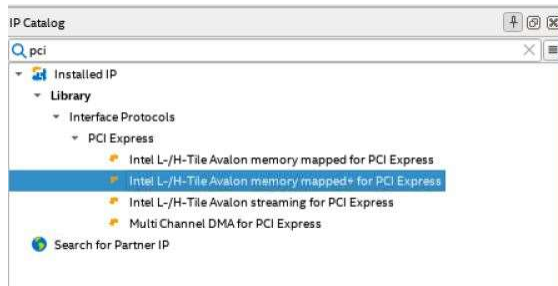
## 2. PCIe EP example setup & Integration with NPU
2.1 PCIe EP example project setup
1) Go to <project home> directory
2) Create one new project with new project wizard & choose empty project
3) Choose S10 MX device with below item at first
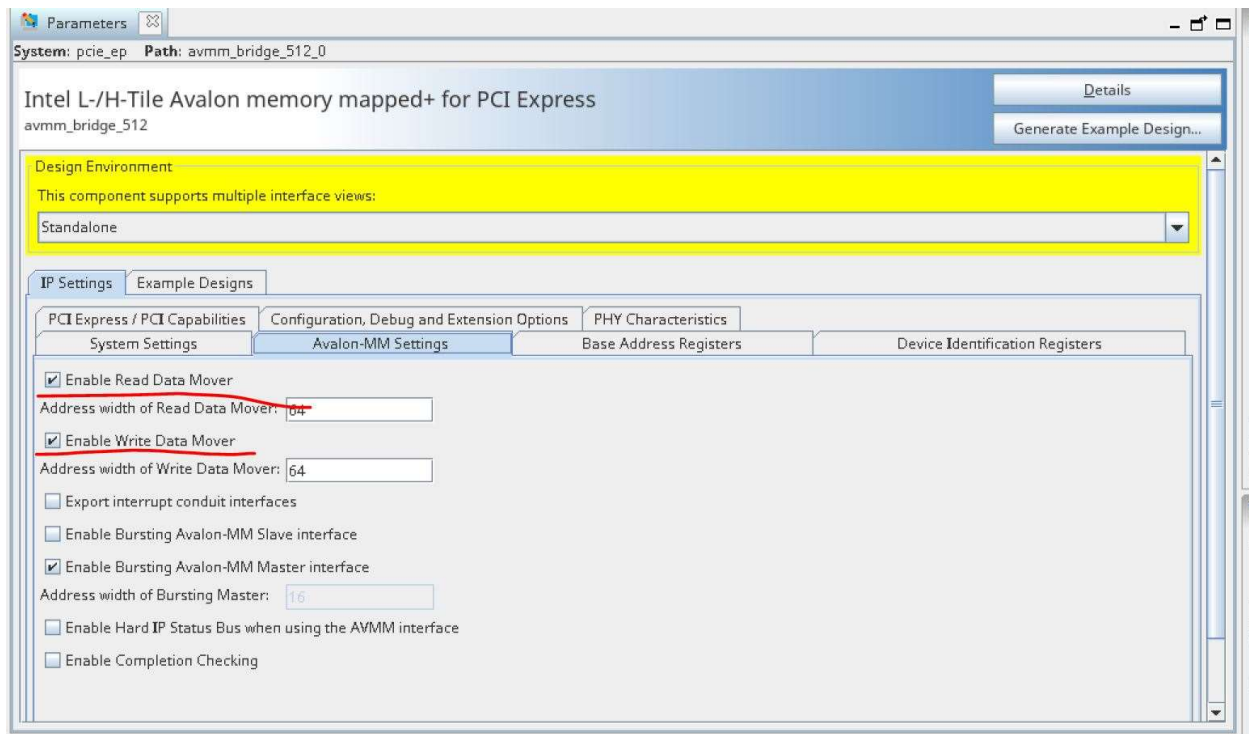


4) Directly click "finish"

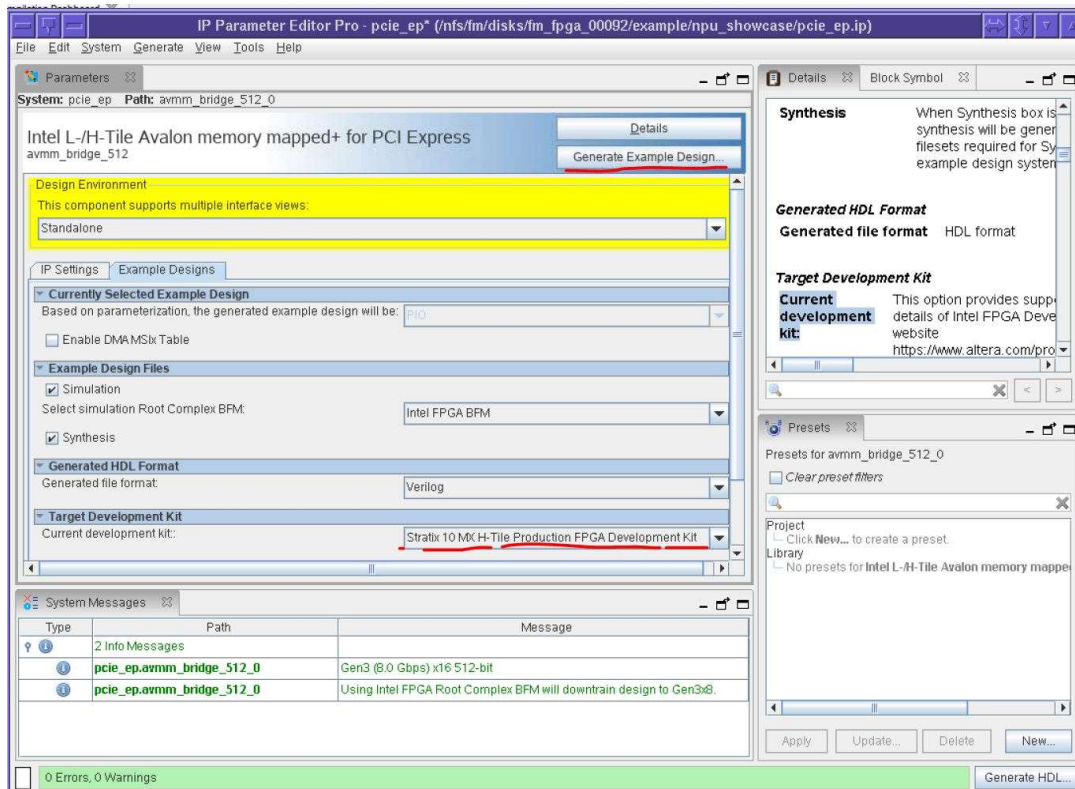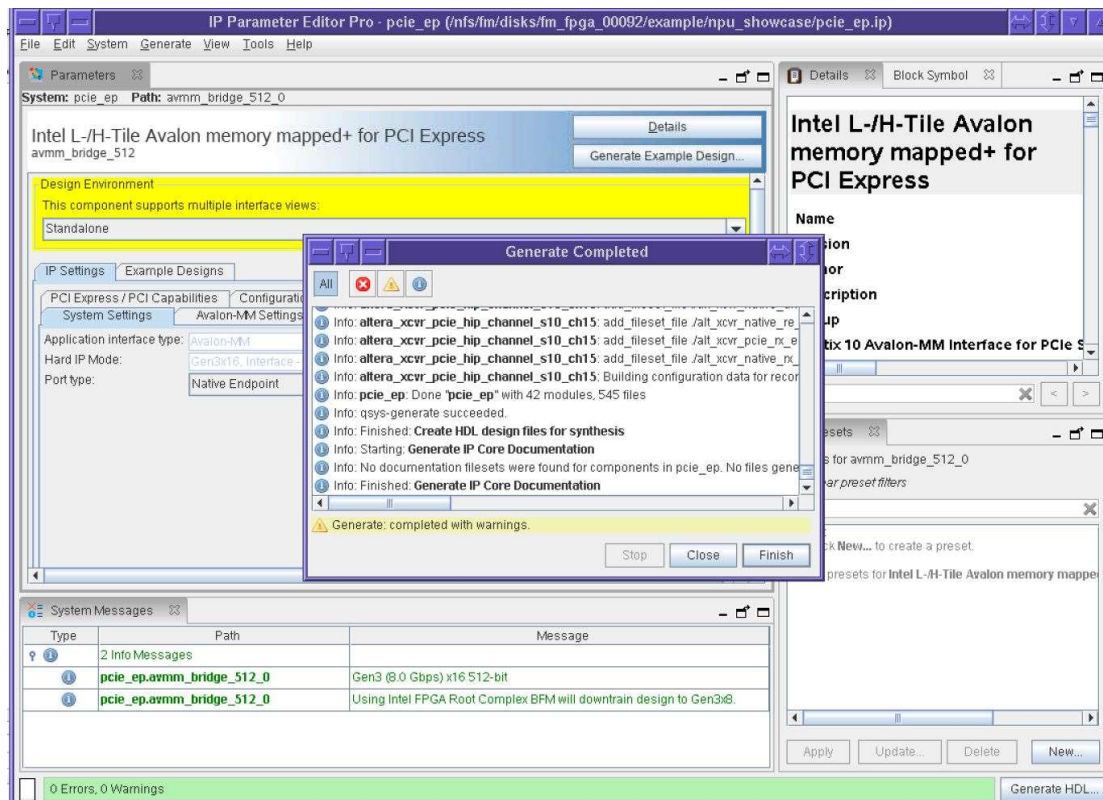5) Search "pci" in IP Catalog & choose that map + as below, then click "add.." button.

6) Create new PCIe endpoint IP
7) In PCIe EP config wizard, in "Avalon-MM settings" tab, make sure selection for "Enable Read Data Mover" & "Enable Write Data Mover". (default selected)
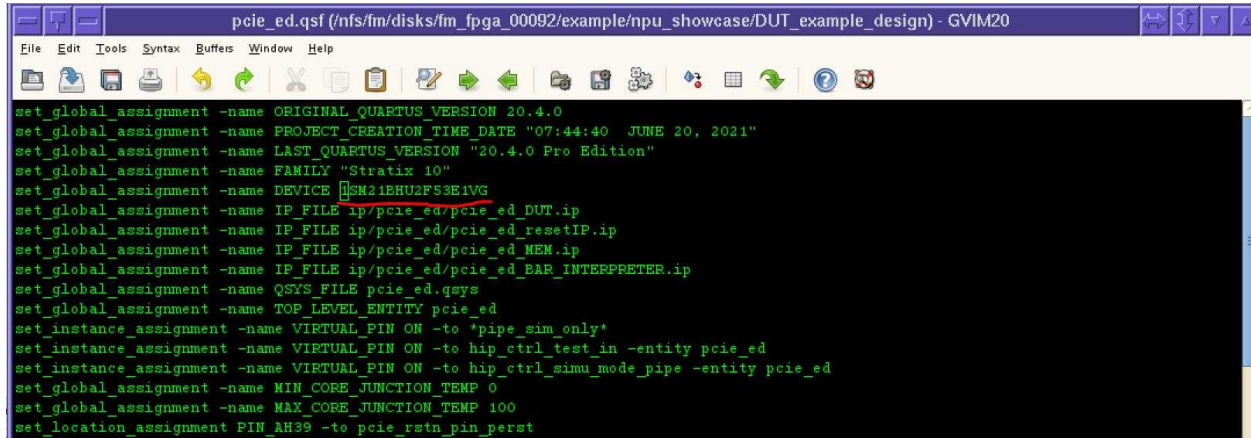


8) In "Example Design" tab, choose S10 MX development kit, then click "Generate Example Design" button.
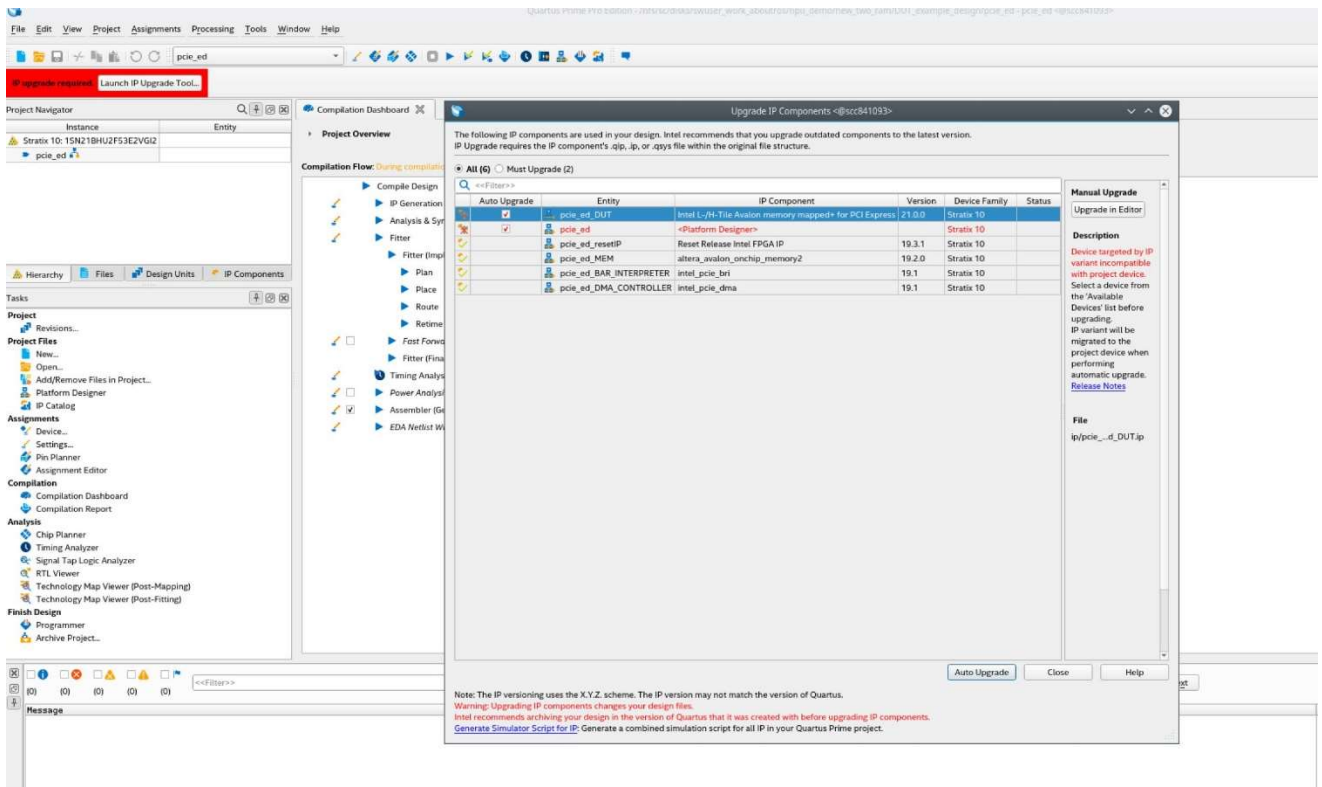
9) Save & generate.

10) Since Stratix 10 NX can't be selected from Quartus GUI, you have to manually edit *.qsf under your new PCIe example generate folder to update device name to "1SN21BHU2F53E2VGI2"
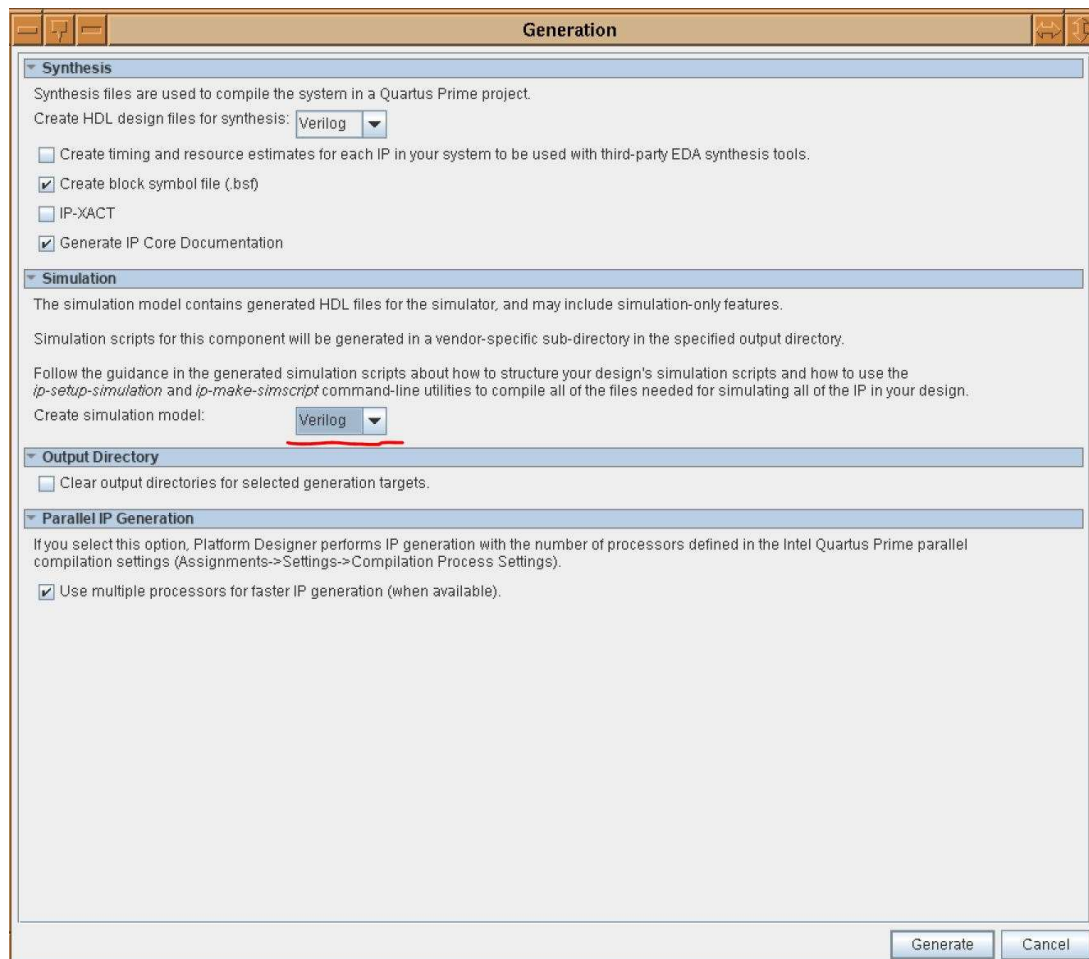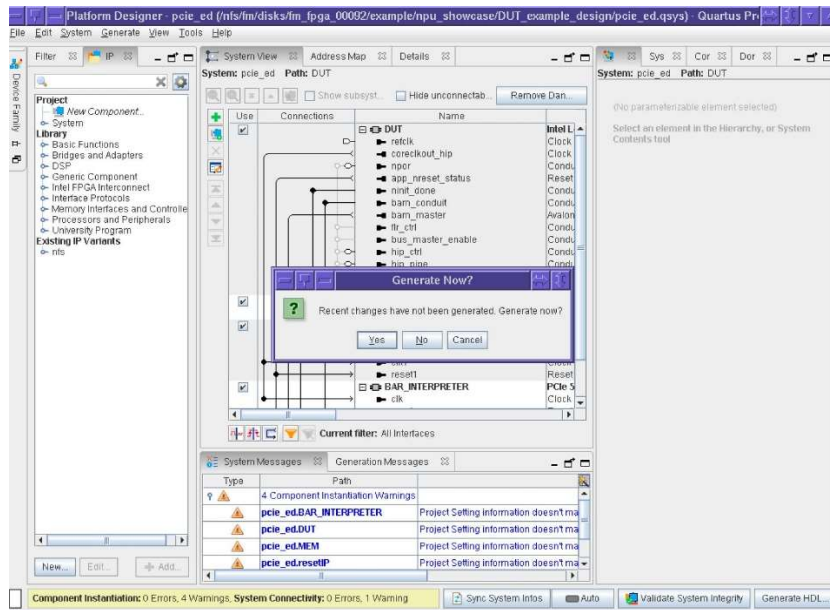


11) Re-load project (*.qpf) under your new PCIe example generate folder, then you will see below notice with RED (you have to update IP for you manually change device name)
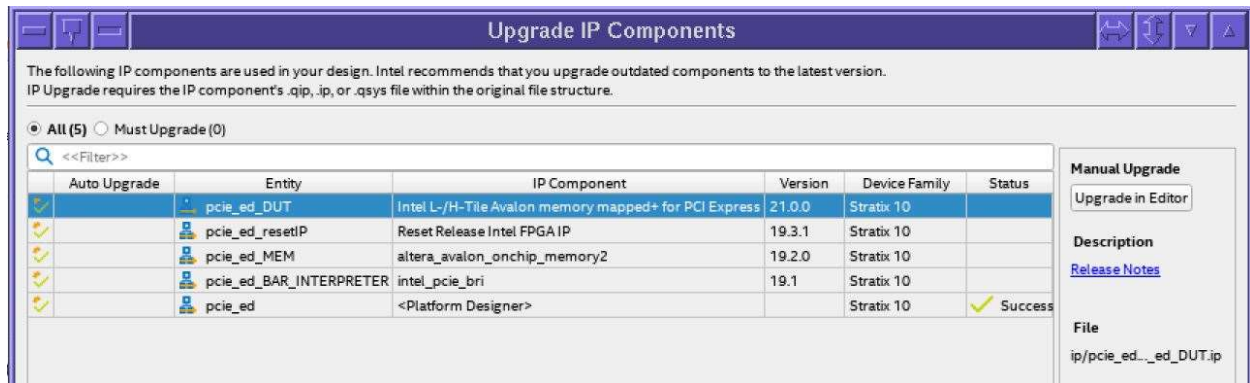
12) Click "Update in Editor" button, load platform designer, save → exit → generate.
(Note, please remember to choose Verilog for simulation)

13) Update IP in red one by one to finally all IP updated & without highlighted with red.



14) Click "run" button to generate one PCIe EP image at first before integrating with NPU IP.



## 2.2 How to get current dual-RAM FPGA module (1Mbyte DMA address)

2.2.1    Extend DMA address (connected to DUT) from 32KB to 1MB

a)    Open your project qsys file, double click that "MEM" IP, in right side property window, you need to change "Total memory size" from 32768(32KB) to 1048576(1MB)

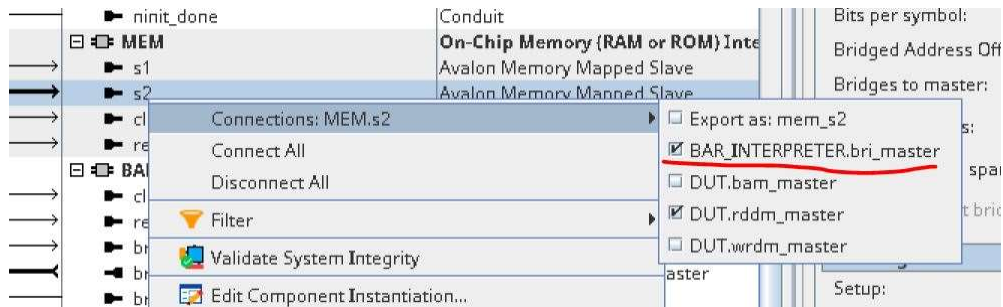b) In Qsys system view, on "mem.s1" & "mem.s2" interfaces, dis-connect  their connection with BAR_INTERPRETER.bri_master (this is for we don't need access DMA RAM content through PIO mode, we use small size DMA read to replace PIO polling,  & big buffer size easy to get conflict on BAR_INTERPRETER bus address).

c) In "Address Map" window of Qsys, you need to unblock MEM.s1 & MEM.s2 memory map setting at first, then change base address from 0x1_0000 to 0x10_0000 to resolve address map conflict.

| Slave | BAR_INTERPRETER.bri_master | DUT.wrdm_master | DUT.rddm_master | |
|---|---|---|---|---|
| BAR_INTERPRETER.b… | | | | 0x000 |
| DMA_CONTROLLER…. | 0x0000_0000 - 0x0000_0fff | | 0x0000_0000_0000_0000 - | |
| MEM.s1 | | 0x0000_0000_0001_0000 - | | |
| MEM.s2 | | | 0x0000_0000_0001_0000 - | |

| Slave | BAR_INTERPRETER.bri_master | DUT.wrdm_master | DUT.rddm_master | |
|---|---|---|---|---|
| BAR_INTERPRETER.b… | | | | 0x000 |
| DMA_CONTROLLER…. | 0x0000_0000 - 0x0000_0fff | | 0x0000_0000_0000_0000 - | |
| MEM.s1 | | 0x0000_0000_0010_0000 - | | |
| MEM.s2 | | | 0x0000_0000_0010_0000 - | |

## 2.3 Apply NPU patch to the Quartus project

1) Go to patch directory
   cd <project home>/patch
2) Apply patch
   ./setup.sh

# 3. Hardware platform setup & Run NPU demo

## 3.1 host PC selection

- X86 CPU & can support at least AVX256, better for AVX512 (can use lscpu in Linux to check CPU flags)
- Support Gen3 x16 PCIe slot

## 3.2 host OS selection

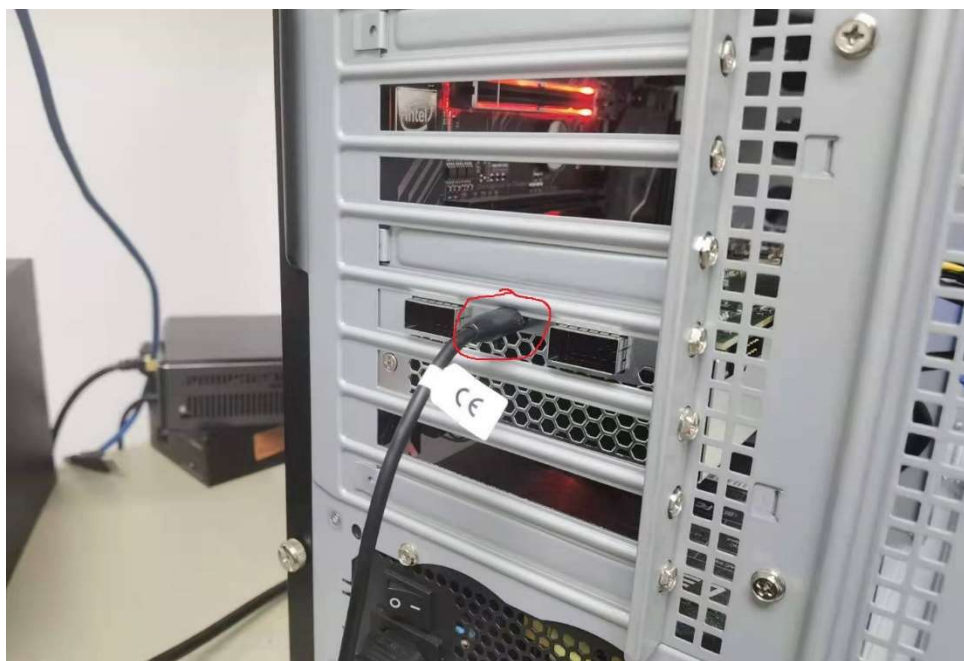Suggest Ubuntu (> 18.0), or PCIe kernel driver maybe need minor change for kernel mode function support.

## 3.3 hardware setup

- Connect S10 NX card to one Gen3 x16 PCIe slot on host PC.
- Need to connect ATX power cable to power slot on S10 NX card.

- Connect USB cable between S10 NX card & your control PC (run Quartus programmer & SigTap tool)
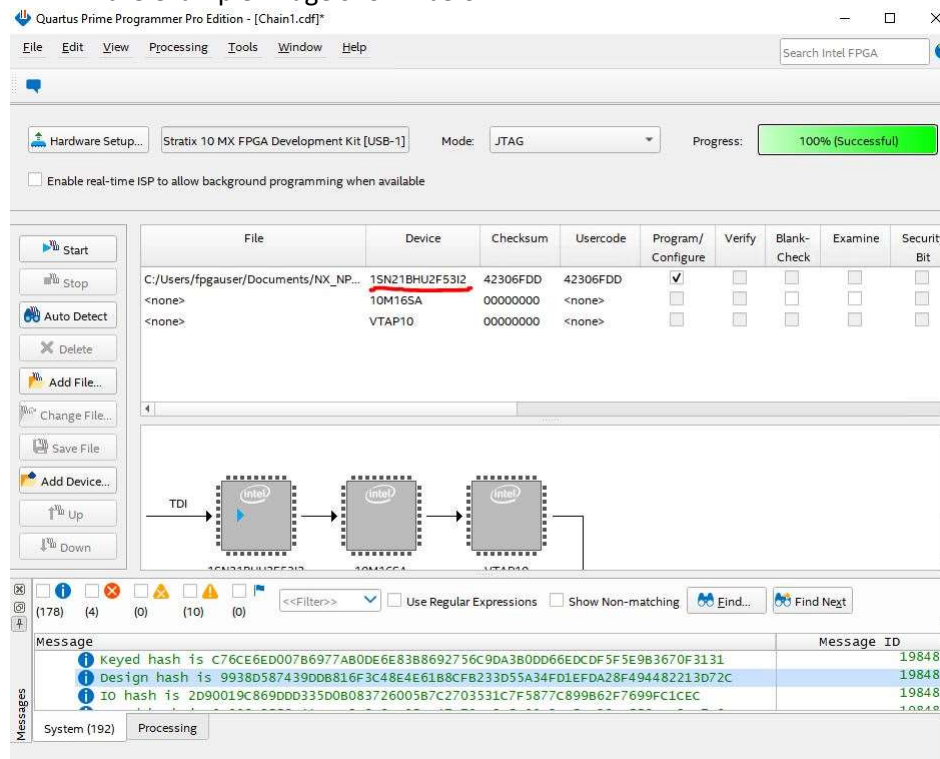
## 3.4 Control PC tool

please install Quartus program tool (not need full Quartus install, but please use 20.4 version)

## 3.5 How to run NPU on FPGA board

a) Download pcie_ed.sof generated in <project directory>/DUT_example_design/ to S10 NX (make sure image is using this device name)
   a. Click "Auto Detect"
   b. Click "Add file" to add the pcie_ed.sof bitstream
   c. Remove the redundant entry, Use "Up" button to make the sequence matching the example image shown below.



b) You need to reboot host PC (Note: it is "reboot", can't be "power down")
c) Copy <project home>/DUT_example_design/software folder to host PC
d) Go to "kernel/linux". For the first time, compile the source files by typing "make".
e) Load the driver by typing "sudo ./load"
f) Type "lspci -d 1172: -vvv" to make sure you can find that PCIe device, BAR2 enabled & driver load

g) Go to "user/npu_test". For the first time, compile the source files by typing "make".
h) Copy benchmark data. E.g., to run mlp_batch4104, type "cp ./mlp_batch4104/* ./"
i) Type "./real_npu_test" to run the test
j) Type 0 for automatic PCIe device find

```
**************************************
Intel FPGA PCIe Link Test
Version 2.0
0: Automatically select a device
1: Manually select a device
**************************************
>
```

```
intel@flexpstbj-i9:~/rui/software/user/npu_test$ ./real_npu_test

***********************************************************
Intel FPGA PCIe Link Test
Version 2.0
0: Automatically select a device
1: Manually select a device
***********************************************************
> 0
Opened a handle to BAR 0x2 of a device with BDF 0x6500
Allocate Kernel memory succesully!
Filling 280 MRFs with 40 elements x 88 words
MRF buffer size = 1835008 Bytes
Finished parsing MRF file!
Filling instruction memory with 48 elements x 34 words
Instruction buffer size = 262144 Bytes
Finished parsing Instructions file!
Filling input memory with 40 elements x 114912 words
Input buffer size = 7602176 Bytes
Finished parsing Inputs file!
Finished parsing Golden Outputs file!
Finished sending MRFs
Finished sending instructions
Latency is 3.563000 ms
Finished sending all inputs and receiving all results!
TEST PASSED!
```

3.6 [Optional] Sigtap debugging on FPGA
- new one *stp file with debug signals you want, & include it into FPGA image build
- Copy that *.stp to control PC
- Use SigTap tool to open this *.stp
- Select correct hardware & device in tool
- Set trigger condition in that "setup" tab.

- Then kick off "capture" button, it will stop & show waveform when trigger condition hit.