

# Final Project (xml: NoC\_FP.xml)

## Part I. Architecture and Algorithm Choice

Network topology	4x4 Torus mesh
Routing Algorithm	X-Y
Switching	Wormhole switching
Buffer Size	1 ("out_flit" slot is as a buffer)
Virtual Channel	No/1 Channel per input port
Flow Control	ACK-REQ protocol

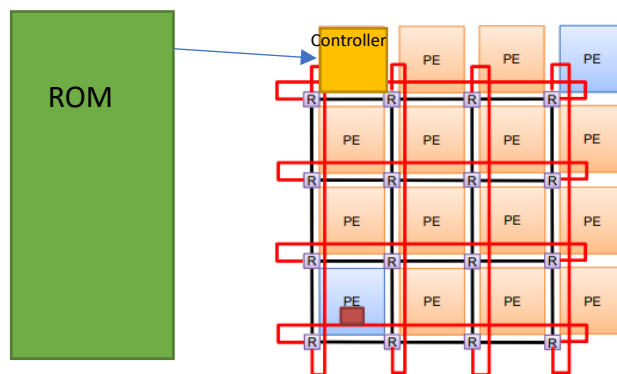


Fig.1 4x4 Torus mesh

To simplify the implementation, I choose X-Y routing algorithm. This algorithm guarantee there will no dead lock or living lock during the system running. Regarding switching method, I use wormhole switching, which only allow the flits from the same packet to fill in a channel. Below shows the design of the flit. The first flit in the packet contains the BoP (the 33th bit), source id[31:28], destination id[27:24], Is\_parameter[23], Is\_bias[22], and OP\_id[21:18]. The body flits begin with 2 zero and follow with data value. Tail flit consists of EoP[32] and the last data value.

Header	BoP [33]	EoP [32]	Src [31:28]	Dst [27:24]	Is_parameter [23]	Is_bias [22]	OP id [21:18]	0 [17:0]
Body	BoP [33]	EoP [32]	Data [31:0]					
Tail	BoP [33]	EoP [32]	Data [31:0]					

Fig2. Flit format

### A. Controller

Controller can access ROM to get the parameters and data, and send them to the assigned place. At first, Controller will send weights and bias for each conv and linear layer. To weights, the header of “Is\_parameter” signal will be high to indicate the packet is for parameter, while “Is\_bias” signal will be low. To bias, instead, “Is\_bias” will be high.

All the weights and bias will be transmitted to the core, to be more precious, they are stored in the PEs.

After finishing sending out all the parameters, Controller will send the image to the first PE (PE\_1) to compute Conv1+MaxPooling1. Then, the following data propagating will be ruled by each PE because they will know their next destination for their output feature map.

Finally, Controller will receive the output feature map from linear3 (fc8). Controller will perform the softmax function, calculate the probability of each class, and print the results on screen.

### B. PE design

PEs are included in the cores. While initializing the cores object, PE will also be initialized and got to know what layer and functions shall be executed in this PE. Also, the attributes of the layer, such as stride and kernel size, will be set then. The below lists the PE id and the corresponding operations. This figure also shows how the Vgg-16 be partitioned and how the data propagate during inference. To be noticed, PE\_0 is replaced by controller.

```
// PE_v2:
// 0: Controller
// 1: Conv1 + ReLU + Max_pooling1
// 2: Conv2 + ReLU + Max_pooling2
// 3: Conv3 + ReLU
// 7: Conv4 + ReLU
// 6: Conv5 + ReLU + Max_pooling3
// 5: Linear1 + ReLU
// 4: Linear2 + ReLU
// 12: Linear3
// 0: + Softmax + Sort
```

Fig3. PE and the corresponding operations

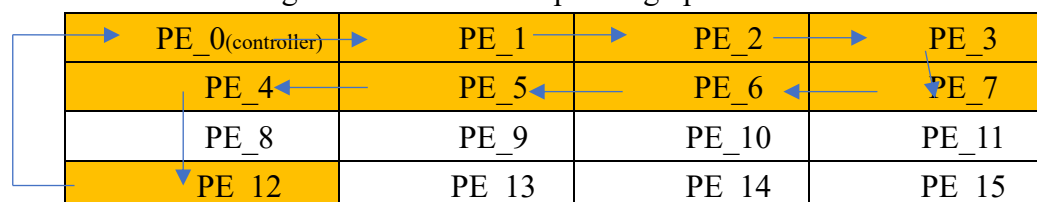


Fig4. PE id and their place  
(Orange block means some ops will be executed there.)

While initializing PEs, PE will also know what next Op will be done in the next destination PE. Then, when sending the output feature map, PE will config the OP\_id in the header of the packet to indicate what op will be executed in the next PE. The below shows the Op\_id and the corresponding operations.

Op_id (in decimal) (use [21:18] in header)	Op type
1	Conv+ReLU+MaxPooling
2	Conv+ReLU
3	Linear+ReLU
4	Linear

#### C. Core design:

It contains a PE and NI in each core. First, the core will always check if there is packet needed to send to somewhere. After fetching the packet, I split the source id, destination id, data and op\_id. Then, I reformat them into flit format, and store in a “std::queue<sc\_lv <34> >”. With ACK-REQ (ack\_tx, req\_tx) handshaking, the flits will send to the router.

Conversely, when “req\_rx” is on, which means the router want to send the flit to this core, the core will let “ack\_rx”, and prepare to format the packet.

After collecting all flit received from the router, the core will execute “check\_packet()” method to send the packet to PE.

#### D. Router design:

It contains only 1 slot buffer for 5 input ports. In the router, I separate it into 2 parts: input port control and output port control.

In the former part, a routing unit is implemented to determine the destination of the input flit. After getting the destination, it will lock the destination output port and limit the user (here means the packet). Then, it will start transmission and get all the receiving flit into the input buffer (handshake at input gate). After handshaking at input gate, the flit will pass to its output gate and wait the next node handshake with it. As for the input gate, the last node will pass new flit

until the tail flit was passed.

In the output port control, it need to check if current flit is a valid data. After handshaking at output gate, output gate will need to wait the flit from input gate handshake. Therefore, it is necessary to control “out\_req” signal.

About the routing unit, which implemented X-Y routing algorithm, it will calculate the next step in each router using its own router id. Therefore, when router get a source id from a header flit, it can calculate the next step according current router id.

## Part II. Simulation results (on PA)

```

-----
[33m***** [0m
[33;5m          Top-5 Results [0m
[33m***** [0m
[32mIndex      Val      Possibility      ClassName      [0m
-----
285      20.206686      96.381480%      Egyptian cat
281      16.136835      1.646189%      tabby
282      15.733852      1.100187%      tiger cat
287      14.790856      0.428478%      lynx
728      14.411864      0.293315%      plastic bag
SystemC: simulation stopped by user.

```

Fig4. account: mlchip007, cat

```

-----
[33m***** [0m
[33;5m          Top-5 Results [0m
[33m***** [0m
[32mIndex      Val      Possibility      ClassName      [0m
-----
207      16.594538      38.627541%      golden retriever
175      15.569661      13.861136%      otterhound
220      15.361864      11.260396%      Sussex spaniel
163      15.002670      7.862449%      bloodhound
219      14.593221      5.220789%      cocker spaniel
SystemC: simulation stopped by user.

```

Fig5. account: mlchip007, dog

## Part III. Challenges and Observations

### A. Test function by unit

Because loading weights , bias and image takes a lot of time, I try to split the whole design into two parts while implement it. One is loading all weights, bias and image. The other is inference the image. Because weights and bias will be stored in the corresponding PE, I write another “pe\_load.h” , which also declare the PE class and point out how PE work, to trace the process and transmission of weights and bias. It will directly read the weights/bias as golden\_weights and golden\_bias. After PE got its weights and bias from ROM, it will check them with golden\_weights and golden\_bias. If there exists error, it will print a fatal message in the screen.

To preload all parameters, please follow Fig6. , Fig7., Fig8..

```

C pe.h      C pe_load.h 2 X
C pe_load.h > ...
91 void PE::init(int pe_id){
270     std::string weight_addr;
271     std::string bias_addr;
272     if(this->id==1){
273         weight_addr = "../data/conv1_weight.txt";
274         bias_addr = "../data/conv1_bias.txt";
275     }
276     else if(this->id==2){
277         weight_addr = "../data/conv2_weight.txt";
278         bias_addr = "../data/conv2_bias.txt";
279     }
280     else if(this->id==3){
281         weight_addr = "../data/conv3_weight.txt";
282         bias_addr = "../data/conv3_bias.txt";
283     }
284     else if(this->id==7){
285         weight_addr = "../data/conv4_weight.txt";
286         bias_addr = "../data/conv4_bias.txt";
287     }
288     else if(this->id==6){
289         weight_addr = "../data/conv5_weight.txt";
290         bias_addr = "../data/conv5_bias.txt";
291     }
292     else if(this->id==5){
293         weight_addr = "../data/fc6_weight.txt";
294         bias_addr = "../data/fc6_bias.txt";
295     }
296     else if(this->id==4){
297         weight_addr = "../data/fc7_weight.txt";
298         bias_addr = "../data/fc7_bias.txt";
299     }
300     else if(this->id==12){
301         weight_addr = "../data/fc8_weight.txt";
302         bias_addr = "../data/fc8_bias.txt";
303     }
304     else{
305         cerr<<"[FATAL] Don't have this id!"<<endl;
306     }
}

```

Fig6. In “pe\_load.h”, after modify line273~302 to access data

```

C pe.h      C core.h 3 X
C core.h > ...

1  #ifndef CORE_H
2  #define CORE_H
3  #define SC_INCLUDE_FX
4
5  #include "systemc.h"
6  // #include "pe.h"
7  #include "pe_load.h"
8

```

Fig7. In “core.h”, comment line6 and uncomment line7.

```

C pe.h      C controller.h 2 X
C controller.h > ...
209 void Controller::run(){
265
266     if(current_step==0){
267         if(((!data_valid.read())&&(!layer_id_valid))){
268             if(next_load==0) current_step = 17;
269             else if(next_load==1) current_step = 2;

```

Fig8. In “controller.h”, modify line 267 as above (original: current\_step will be 1)

<b>data folder</b>	Include weight, bias, and input image matrix
<b>run</b>	Executable files for SystemC program
<b>controller.h</b>	Implement controller
<b>core.h</b>	Implement core and Ni
<b>pe.h</b>	Implement the pe
<b>pe_load.h</b>	Debug use, for load parameters
<b>router.v</b>	Implement the router
<b>Makefile</b>	Makefile script for compile systemC program.
<b>main.cpp</b>	Declare the main function, create the module instances, mapping the signals. It includes all operation units, pattern module, clockreset modules, and monitor module.
<b>ROM.h/ROM.o/ROM.cpp</b>	ROM stores image and parameters
<b>clockreset.h</b>	Declare the clock module and reset modules.
<b>clockreset.cpp</b>	Implement the clock module and reset modules.
<b>NoC_FP.xml</b>	PA xml file

