Convolutional Neural Network(CNN)

AI RTL

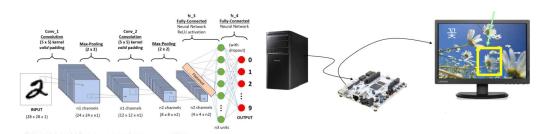
이재성

목차

- ■1. Project 목표
- ■2. 설계 내용
- ■3. 결과

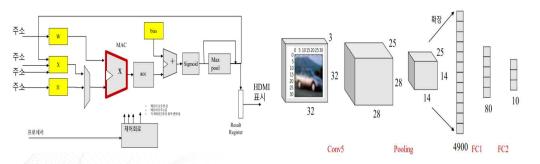
Project 목표

- ■Verilog HDL로 Convolutional Neural Network 구현
- ■HDMI로 모니터와 연결후 인식 결과 확인



Project 목표

CNN 연산 구조



MAC 하나만 사용 연산 중 MAC은 무휴로 운용 4개의 32x32x3 이미지, 25개의 5x5x3 weight, maxpooling, bias, sigmoid, fully-connected 연산

■1. 입력 이미지 i1 , weight address value

```
always @(posedge clk) begin
   case(s_fsm)
       2'bill: begin s mireal: lastcall: weakel: may nool on officel: encel://opens/ution operation
                if (addr. i>12287) addr. i<=0:
                if (n==4) begin
                     n<=0:
                     if (1==4) begin
                         1<=0;
                         if (n=-2) begin
                             nceff:
                             last <=1:
                             if (kee27) begin
                                 ks=0:
                                 if (i==27) hegin
                                     i<=0:
                                     if (i==24) begin
                                         i<=0;
                                         s_fse<-2'b01;
                                         addr_w<=addr_w+1;
                                         addr_bias<=addr_bias+1;
                                         if (image==3) hegin
                                              inage«-fi:
                                             addr is=0:
                                         end etse begin
                                         image<=image+1;
                                         addr_i<-addr_i+1;
```

3개의 state(convolution, fully-connect 1,2) Image와 weight의 행과 열 값에 따라 address 계산

```
i<=i+1:
                    addr_w<=75+i;
                    addr i<=addr i-3071:
                    addr bias<=addr bias+1;
                    addr_x_Write<=addr_x_Write+1:
                    wea<=1;
                    end
                end else hegin
                i<=i+1;
                addr_i<=addr_i=2175; addr_w<=75*i; //addr_i<=image <3072+32*j
            end else begin
            k<=k+1;
            addr_i<=addr_i-2179; addr_w<=75+i; //addr_i<=image+3072+32+i+kc
        end etse begin
        ncen+1:
        addr is=addr i+892; addr ws=addr w+1; //addr is=image+3072+32+i+k+1024+n
        and
    end else hegin
    1<=1+1;
    addr_i<=addr_w<=addr_w+1; //image +3072+32+j+k+1024+n+32+1;
end else begin
m<=m+1:
addr i<=addr i+1; addr w<=addr w+1;
```

THE PERSON IS NOT THE OWNERS OF THE PERSON NAMED IN COLUMN 2 IN CO

주소 값 직접 연산 Bias, x memory write 주소 설정

■1. 입력 이미지 i1, weight address value

> • ex10/image[1:0]	1						1				
> W ex10/i[4:0]	21						21				
> W ex10/j[4:0]	3						3				
> ₩ ex10/k[4:0]	18		18								
→ ex10/n[1:0]	1		1						2		
> ₩ ex10/I[2:0]	4	3			4					0 '	
₩ ex10/m[2:0]	3	2	3	4	0	1	2	3	4	0	1
₩ ex10/addr_i[13:0]	4341	4308	4309	4310	4338	4339	4340	4341	4342	5234	523
₩ ex10/out_i[15:0]	265	206	190	135	119	159	307	265	225	191	19
₩ ex10/addr_w[18:0]	1623	1617	1618	1619	1620	1621	1622	1623	1624	1625	162
₩ ex10/out_w_1clock[9:0]	41	21	8	1009	12	30	62	41	23	3	100
■ ex10/mux_out[15:0]	307	333	206	190	135	119	159	307	265	225	19
₩ ex10/mul_out[25:0]	9858	263160	6993	1648	191710	1620	3570	9858	12587	6095	675

i1 address = addr_i 첫번째 이미지에서 i=21, j=3, k=18, n=1, l=4, m=0일때 addr_i=4338 weight address = addr_w 첫번째 이미지에서 i=21, n=1, l=4, m=0일때 addr_w=1620

2. MAC operation

multiple

```
always @(posedge clk) begin
  out_w_lclock<=out_w;
  mul_out<=mux_out+out_w_lclock;
end</pre>
```

mux의 output과 클락을 맞추기 위해 read한 weight값 한클락 지연시킨 후 multiple

accumulate

```
always @(posedge clk) begin
   acc_out<=mul_out+acc_out;
   if(last) acc_out<=0;
end</pre>
```

마지막 데이터일때 last=1로 만들어 데이터 초기화



3. bias

```
if (i==24) begin
    i<=0:
    s_fsm<=2'b01;
    addr w<=addr w+1;
    addr_bias<=addr_bias+1;
    if(image==3) begin
        image<=0:
       addr_i<=0;
    end else begin
    image<=image+1:
    addr_i<=addr_i+1;
    addr x Read<=0:
    end
end else begin
i<=i+1:
addr_w<=75*i;
addr i<=addr i-3071;
addr_bias<=addr_bias+1;
```

bias address 값 연산 i값이 바뀔때마다 1씩 증가

```
mem bias Bias (
  .clka(clk), // input wire olka
  .addra(addr_bias), // input wire [6 : 0] addra
  .douta(out_bias) // output wire [15: 0] doute
);
  bias memory read (16bits)
 (* mark_debug = "true" *) reg [31:0] add_bias;
 always @(posedge clk) begin
    add bias<= out bias conv + acc out;
 end
```

Overflow 방지와 sigmoid 연산을 위해 출력인 add bias는 32비트로 설정

4. sigmoid

$$sigmoid(x) = \frac{1}{1 + e^{-x}}$$

bias output floating-point convert

floating-point division

floating-point exponential

floating-point ip(exp, add, divide) 사용하여 연산

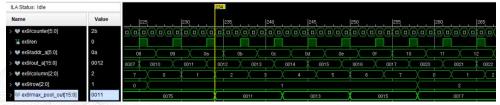
5. maxpooling

```
always @(posedge clk) begin
if(max_pool_no_of) begin
if(en) begin
column <= column +1:
    reg_array(0]<=sigmoid_out_fload:
    reg_array(0]<=sigmoid_out_fload:
    reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0]<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=reg_array(0)<=
```

```
if(column=5'bl1100) besin row-crow+1; column <=5'b000001: end
if(row=5'b11000) row-5'b000001;
if(column[0]==88rou[0]==1) besin
mex!=max(rea_array[14], rea_array[15]);
mex2-max(rea_array[1], max[1];
mex2-max(rea_array[1], max[2]);
end
end
end
else mex_pool_out<=signoid_out_tixed:
-en 값 설정하여 신호제어
-167# Q! register array 사용하여 row column
```

-16개의 register array 사용하여 row, column 에 따라 연산

-max_pool_on_off 값에 따라 maxpool 연산 수행여부 결정



6. fully-connect operation - 1

```
if (fc1==79) begin
   fc1<=0:
   addr_x_Read<=addr_x_Read+1;
   addr_w<= addr_w+1;
   addr_bias<=addr_bias+1;
   addr x Write<=0:
   s fsm<=2'b10:
etse begin
   fc1<=fc1+1:
   addr_x_Read <= 0;
   addr_w<= addr_w+1;
   addr_bias<=addr_bias+1;
   addr_x_Write<=addr_x_Write+1;
   wea<=1:
end
```

```
-input: 4900개의 maxpooling 출력 데이터
-output: 80개
-weight: 4900x80개의 데이터
-x memory write address: 출력값 나올때
마다 1씩 증가
-x memory read address, weight address:
데이터 읽을때마다 1씩 증가
```

7. fully-connect operation - 2

```
if (r==79) begin
    r<=0:
    if(fc2==9) begin
        fc2<=0:
        addr × Read<=0:
        addr bias<=0;
        addr w<=0;
        s fsm<=2'b00;
    else begin
        fc2<=fc2+1:
        addr_bias<=addr_bias+1;
        addr_x_Read<=4900;
        addr w<=addr w+1:
        en2<=1:
    end
end else begin
    r<=r+1;
    addr x Read<=addr x Read+1;
    addr w<=addr w+1;
end
end
```

-input: 80개의 첫번째 fully-connect 출력 데이터 -output: 10개 -weight: 80x10개의 데이터 -x memory write address: 0으로 설정 -x memory read address, weight address: 데이터 읽을때마다 1씩 증가

8. output

```
always @(posedge clk) begin
    if (en2) begin
        resultReg[0]<=max_pool_out;
        resultReg[1]<=resultReg[0];
        resultReg[2]<=resultReg[1];
        resultReg[3]<=resultReg[2];
        resultReg[4]<=resultReg[3];
        resultReg[5]<=resultReg[4];
        resultReg[6]<=resultReg[5];
        resultReg[7]<=resultReg[6];
        resultReg[8]<=resultReg[7];
        resultReg[9]<=resultReg[8];
        result count <= result count +1:
        if (result count == 4 b 1010) begin
            maxla<=max(resultReg[0],resultReg[1]); max2a<=max(maxla,resultReg[2]);
            max3a<=max(max2a,resultReg[3]); max4a<=max(max3a,resultReg[4]);
            max5a<=max(max4a,resultReg[5]); max6a<=max(max5a,resultReg[6]);
            max7a<=max(max6a,resuitReg[7]); max8a<=max(max7a,resuitReg[8]);
            max_result<=max(max8a,resultReg[9]);
    end
```

- -2번째 fully-connect 연산의 결과를 저장하기 위해 10개의 register array 사용
- -register array에 데이터가 모두 들어 오면 최대값 구한 후 출력

결과

- ■최종 미완성
- ■완성하지 못한 것
- ■1. sigmoid에서 exponential 값 에러



- ■2. multiple 연산 이후의 delay 등을 고려한 타이밍 조절
- 3. max pooling 후 x memory write address 할당
- ■4. accumulator overflow 조절
- ■5. 연산 결과 및 HDMI 케이블로 연결하여 화면에서 인식 결과 확인