

Fault-Tolerant Systolic Array Based Accelerators for Deep Neural Network Execution

Introduction:

With the rapid development of AI, the demand for AI-accelerated chips is getting higher and higher.

In order to achieve faster and more efficient computing AI models, custom AI or SOC chips have become a trend.

The following figure shows the basic architecture of the AI accelerator, or Neural-network Processing Unit (NPU) which is mainly composed of Systolic Array (SA), whose advantage is that the matrix operation is fast, and because it does not need to frequently access the memory, so it is much more energy efficient than CPU and GPU.

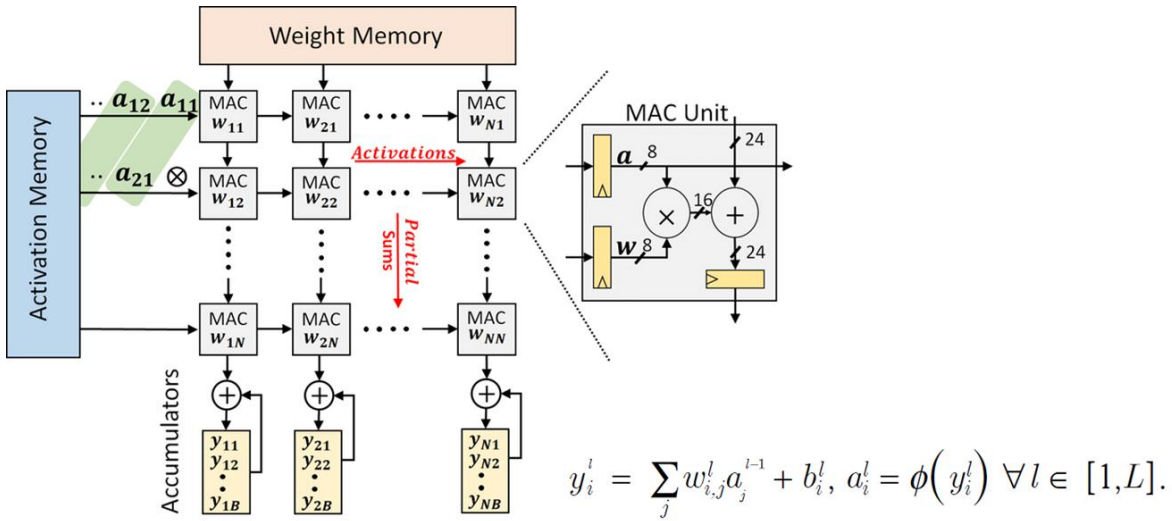


Fig. 1 Basic architecture of the AI accelerator [1]

However, because the characteristic of Systolic array is passed one by one, and the number of PEs is large, the operation error of a PE often becomes a fatal error of neural network. Fig. 2 show the classification accuracy drop due to stuck-at fault MACs. Therefore, we have adopted the following design to try to avoid operation errors caused by hardware errors, which is showed as Fig. 3.

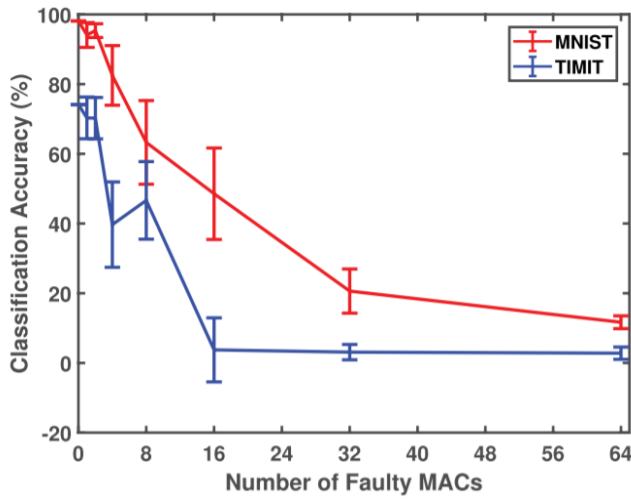


Fig. 2 Impact of classification accuracy drop due to stuck-at fault MACs

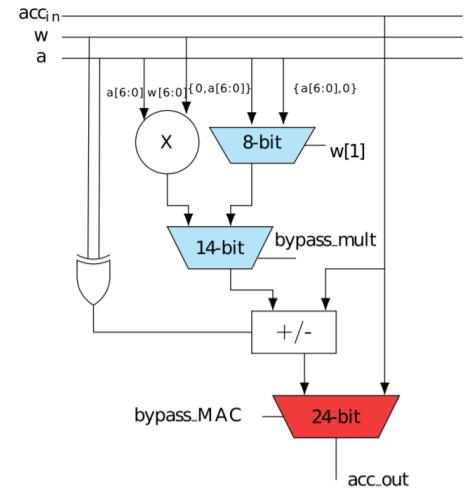
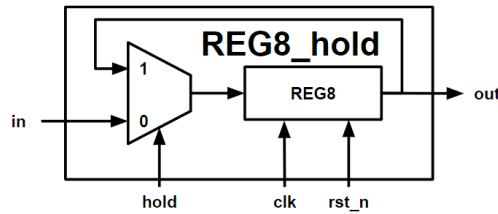


Fig. 3 Illustration of Fault-Tolerant MAC unit

Experimental and programming:

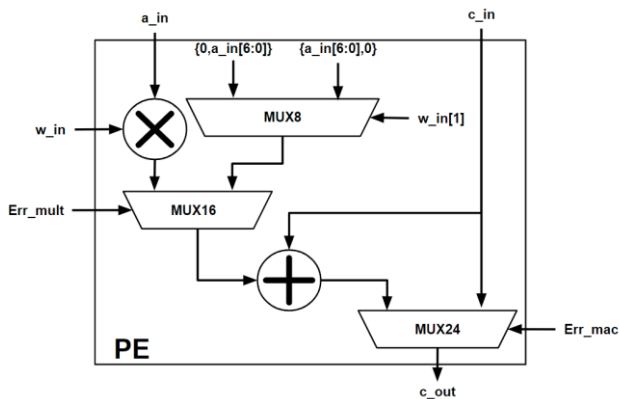


```

1 module REG8_hold(clk,rst_n,in,out,s);
2 parameter N=8;
3 input clk,s,rst_n;
4 input signed [N-1:0] in;
5 output signed [N-1:0] out;
6 reg signed [N-1:0] in_wire;
7 reg signed [N-1:0] out;
8 always @ (*)
9 begin
10     case(s)
11         1'b0:in_wire <= in;
12         1'b1:in_wire <= out;
13     endcase
14 end
15 always @ (posedge clk or negedge rst_n)
16 begin
17     if (!rst_n) out <= {N{1'b0}};
18     else out <= in_wire;
19 end
20 endmodule

```

Fig. 4 REG diagram & Verilog code

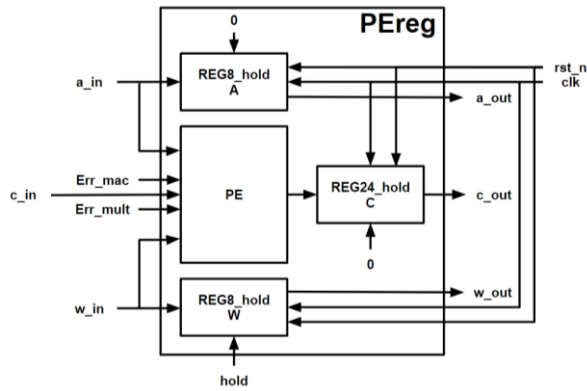


```

1 module PE(Err_mac, Err_mult, w_in, a_in, c_in, c_out);
2 input signed [7:0] w_in, a_in;
3 input signed [23:0] c_in;
4 input Err_mac, Err_mult;
5
6 output signed [23:0] c_out;
7
8 wire signed [15:0] product;
9 wire signed [7:0] shift;
10 wire signed [15:0] acc_in;
11 //reg signed [23:0] acc_in_24;
12 wire signed [23:0] acc_out;
13
14 //multiplier multi(product, w_in, a_in);
15 assign product = w_in * a_in;
16 MUX8 sft(w_in[1],1'b0,a_in[6:0]),{a_in[6:0],1'b0},shift);
17 MUX16 bypass_mult(Err_mult,product,{8'h00,shift},acc_in);
18 MUX24 bypass_MAC(Err_mac,acc_out,c_in,c_out);
19 //wire [23:0] acc_in_24 = { 8{acc_in[15]}, acc_in[15:0] };
20
21 /*
22 always @ (*)
23 begin
24     //if(acc_in[7]==1) acc_in_24 = {8'b1111_1111,acc_in};
25     //else acc_in_24 = {8'b0000_0000,acc_in};
26
27 end
28 */
29
30 assign acc_out = c_in + acc_in;
31 endmodule

```

Fig. 5 PE diagram & Verilog code



```

1 module Pereg(Err_mac, Err_mult, w_in, a_in, c_in, w_out, a_out, c_out, clk, rst_n,
2 hold);
3
4 input signed [7:0] w_in, a_in;
5 input signed [23:0] c_in;
6 input Err_mac, Err_mult;
7 input clk;
8 input rst_n;
9 input hold;
10
11 output signed [7:0] w_out, a_out;
12 output signed [23:0] c_out;
13
14 wire signed [23:0] p_sum;
15
16 PE PE_no_reg(Err_mac, Err_mult, w_out, a_out, c_in, p_sum);
17 //PE(Err_mac, Err_mult, w_in, a_in, c_in, c_out);
18
19 REG8_hold regw(clk, rst_n, w_in, w_out, hold);
20 REG8_hold rega(clk, rst_n, a_in, a_out, 1'b0);
21 REG24_hold regc(clk, rst_n, p_sum, c_out, 1'b0);
22
23 endmodule

```

Fig. 6 Pereg diagram & Verilog code

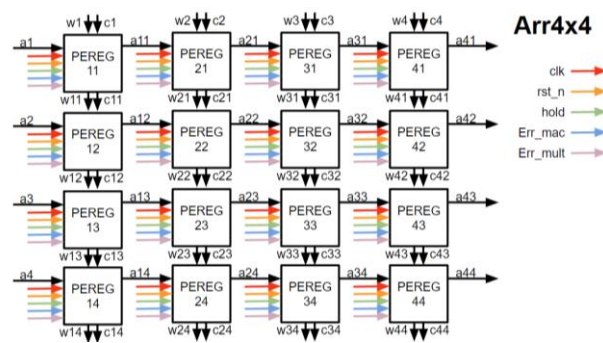


Fig. 7 Arr4x4 diagram

```

1 module Arr4x4(Err_mac, Err_mult,
2 w1_in, w2_in, w3_in, w4_in,
3 a1_in, a2_in, a3_in, a4_in,
4 w1_out, w2_out, w3_out, w4_out,
5 a1_out, a2_out, a3_out, a4_out,
6 c1_out, c2_out, c3_out, c4_out,
7 clk, rst_n, hold);
8 //Pereg(Err_mac, Err_mult, w_in, a_in, c_in, w_out, a_out, c_out, clk, rst_n, hold);
9 input [7:0] w1_in, w2_in, w3_in, w4_in;
10 input [7:0] a1_in, a2_in, a3_in, a4_in;
11 input clk, rst_n, hold;
12 input [15:0] Err_mac, Err_mult;
13
14 output [7:0] w1_out, w2_out, w3_out, w4_out;
15 output [7:0] a1_out, a2_out, a3_out, a4_out;
16 output [23:0] c1_out, c2_out, c3_out, c4_out;
17
18
19
20 wire [7:0] w11_out, w21_out, w31_out, w41_out,
21 w12_out, w22_out, w32_out, w42_out,
22 w13_out, w23_out, w33_out, w43_out,
23 w14_out, w24_out, w34_out, w44_out;
24 wire [7:0] a11_out, a21_out, a31_out, a41_out,
25 a12_out, a22_out, a32_out, a42_out,
26 a13_out, a23_out, a33_out, a43_out,
27 a14_out, a24_out, a34_out, a44_out;
28 wire [23:0] c11_out, c21_out, c31_out, c41_out,
29 c12_out, c22_out, c32_out, c42_out,
30 c13_out, c23_out, c33_out, c43_out,
31 c14_out, c24_out, c34_out, c44_out;
32
33 //Pereg(Err_mac, Err_mult, w_in, a_in, c_in, w_out, a_out, c_out, clk1, clk2);
34 Pereg PE11( Err_mac[0], Err_mult[0], w1_in, a1_in, 24'sh00_0000, w11_out, a11_out, c11_out, clk, rst_n, hold);
35 Pereg PE21( Err_mac[1], Err_mult[1], w2_in, a11_out, 24'sh00_0000, w21_out, a21_out, c21_out, clk, rst_n, hold);
36 Pereg PE31( Err_mac[2], Err_mult[2], w3_in, a21_out, 24'sh00_0000, w31_out, a31_out, c31_out, clk, rst_n, hold);
37 Pereg PE41( Err_mac[3], Err_mult[3], w4_in, a31_out, 24'sh00_0000, w41_out, a41_out, c41_out, clk, rst_n, hold);

```

```

39 Pereg PE12( Err_mac[4], Err_mult[4], w11_out, a2_in, c11_out, w12_out, a12_out, c12_out, clk, rst_n, hold);
40 Pereg PE22( Err_mac[5], Err_mult[5], w21_out, a12_out, c21_out, w22_out, a22_out, c22_out, clk, rst_n, hold);
41 Pereg PE32( Err_mac[6], Err_mult[6], w31_out, a22_out, c31_out, w32_out, a32_out, c32_out, clk, rst_n, hold);
42 Pereg PE42( Err_mac[7], Err_mult[7], w41_out, a32_out, c41_out, w42_out, a42_out, c42_out, clk, rst_n, hold);
43
44 Pereg PE13( Err_mac[8], Err_mult[8], w12_out, a3_in, c12_out, w13_out, a13_out, c13_out, clk, rst_n, hold);
45 Pereg PE23( Err_mac[9], Err_mult[9], w22_out, a13_out, c22_out, w23_out, a23_out, c23_out, clk, rst_n, hold);
46 Pereg PE33(Err_mac[10], Err_mult[10], w32_out, a23_out, c32_out, w33_out, a33_out, c33_out, clk, rst_n, hold);
47 Pereg PE43(Err_mac[11], Err_mult[11], w42_out, a33_out, c42_out, w43_out, a43_out, c43_out, clk, rst_n, hold);
48
49 Pereg PE14(Err_mac[12], Err_mult[12], w13_out, a4_in, c13_out, w14_out, a14_out, c14_out, clk, rst_n, hold);
50 Pereg PE24(Err_mac[13], Err_mult[13], w23_out, a14_out, c23_out, w24_out, a24_out, c24_out, clk, rst_n, hold);
51 Pereg PE34(Err_mac[14], Err_mult[14], w33_out, a24_out, c33_out, w34_out, a34_out, c34_out, clk, rst_n, hold);
52 Pereg PE44(Err_mac[15], Err_mult[15], w43_out, a34_out, c43_out, w44_out, a44_out, c44_out, clk, rst_n, hold);
53 //Pereg(Err_mac, Err_mult, w_in, a_in, c_in, w_out, a_out, c_out, clk, rst_n, hold);
54
55 assign c1_out = c14_out;
56 assign c2_out = c24_out;
57 assign c3_out = c34_out;
58 assign c4_out = c44_out;
59
60 assign w1_out = w14_out;
61 assign w2_out = w24_out;
62 assign w3_out = w34_out;
63 assign w4_out = w44_out;
64
65 assign a1_out = a41_out;
66 assign a2_out = a42_out;
67 assign a3_out = a43_out;
68 assign a4_out = a44_out;
69
70 endmodule

```

Fig. 8 Arr4x4 Verilog code

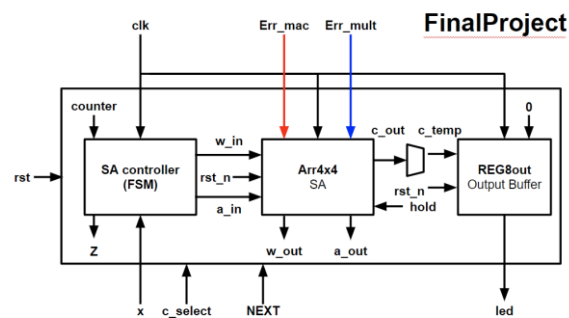


Fig. 9 FinalProject diagram

```

1 module FinalProject(clk,x,rst,led,Err_mac_in,Err_mult_in,c_select,NEXT);
2
3 input NEXT;
4 input clk;
5 input x, rst;
6 input Err_mac_in, Err_mult_in;
7 input [1:0] c_select;
8
9 output [7:0] led;
10
11 parameter S0 = 2'b00,S1 = 2'b01,S2=2'b10;
12
13 assign rst_n = !rst;
14
15 wire [1:0] z;
16 wire [3:0] counter;
17
18 reg [7:0] w1_in;
19 reg [7:0] w2_in;
20 reg [7:0] w3_in;
21 reg [7:0] w4_in;
22
23 reg [7:0] a1_in;
24 reg [7:0] a2_in;
25 reg [7:0] a3_in;
26 reg [7:0] a4_in;
27
28 reg [7:0] c_temp;

```

```

35 reg hold;
36 reg clk_s;
37 reg [28:0] clk_counter;
38 reg [15:0] Err_mac, Err_mult;
39 always @(*)
40 begin
41     if(Err_mac_in)
42         Err_mac = 16'b0000_0000_0001_0000;
43     else
44         Err_mac = 16'b0000_0000_0000_0000;
45
46     if(Err_mult_in)
47         Err_mult = 16'b0000_0010_0000_0000;
48     else
49         Err_mult = 16'b0000_0000_0000_0000;
50 end
51
52 FSM fsm(clk,x,z,rst,counter);
53 REG8out temp(clk_s, rst_n, c_temp, led, 1'b0);
54 Arr4x4 SA(Err_mac, Err_mult,
55     w1_in, w2_in, w3_in, w4_in,
56     a1_in, a2_in, a3_in, a4_in,
57     w1_out, w2_out, w3_out, w4_out,
58     a1_out, a2_out, a3_out, a4_out,
59     c1_out, c2_out, c3_out, c4_out,
60     clk, rst_n, hold);
61

```

```

62  always @(*)
63  begin
64      case(c_select)
65          2'b00 : begin c_temp = c1_out[7:0];end
66          2'b01 : begin c_temp = c2_out[7:0];end
67          2'b10 : begin c_temp = c3_out[7:0];end
68          2'b11 : begin c_temp = c4_out[7:0];end
69          default : c_temp = c1_out[7:0];
70      endcase
71  end

```

```

91  always @(*)
92  begin
93      case(z)
94          S0 : begin
95              clk_s = NEXT;
96              hold = 0;
97          end
98          S1 : begin
99              hold = 0;
100             clk_s = clk;
101             case(counter)
102                 4'h0 : begin w1_in=4; w2_in=3; w3_in=2; w4_in=1; end
103                 4'h1 : begin w1_in=8; w2_in=7; w3_in=6; w4_in=5; end
104                 4'h2 : begin w1_in=4; w2_in=3; w3_in=2; w4_in=1; end
105                 4'h3 : begin w1_in=8; w2_in=7; w3_in=6; w4_in=5; end
106             endcase
107         end
108         S2 : begin
109             hold = 1;
110             clk_s = clk;
111             case(counter)
112                 4'h0 : begin a1_in=8; a2_in=0; a3_in=0; a4_in=0; end
113                 4'h1 : begin a1_in=7; a2_in=8; a3_in=0; a4_in=0; end
114                 4'h2 : begin a1_in=6; a2_in=7; a3_in=8; a4_in=0; end
115                 4'h3 : begin a1_in=5; a2_in=6; a3_in=7; a4_in=8; end
116                 4'h4 : begin a1_in=4; a2_in=5; a3_in=6; a4_in=7; end
117                 4'h5 : begin a1_in=3; a2_in=4; a3_in=5; a4_in=6; end
118                 4'h6 : begin a1_in=2; a2_in=3; a3_in=4; a4_in=5; end
119                 4'h7 : begin a1_in=1; a2_in=2; a3_in=3; a4_in=4; end
120                 4'h8 : begin a1_in=0; a2_in=1; a3_in=2; a4_in=3; end
121                 4'h9 : begin a1_in=0; a2_in=0; a3_in=1; a4_in=2; end
122                 4'ha : begin a1_in=0; a2_in=0; a3_in=0; a4_in=1; end
123                 4'hb : begin a1_in=0; a2_in=0; a3_in=0; a4_in=0; end
124                 default: begin a1_in=0; a2_in=0; a3_in=0; a4_in=0; end
125             endcase
126         end
127     endcase
128 end
129 endmodule

```

Fig. 10 FinalProject Verilog code

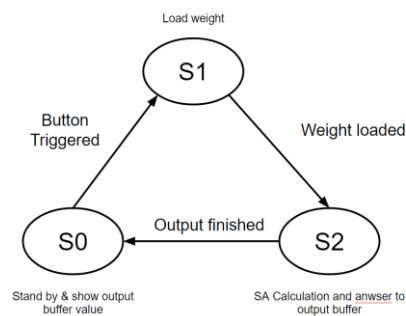


Fig. 11 FSM diagram

```

1  module FSM(clk,x,z,rst,counter);
2
3  input clk;
4  input x, rst;
5
6  output [1:0] z;
7  output reg [3:0] counter;
8
9  reg [1:0] Z;
10 reg [2:0] counter1;
11 reg [3:0] counter2;
12 reg [1:0] next_state;
13 reg [1:0] current_state;
14 //state encoding
15 parameter S0 = 2'b00,S1 = 2'b01,S2=2'b10;

```

```

17 //counter1
18 always @(posedge clk , posedge rst)
19 begin
20     if(rst==1)
21         counter1 <= 3'b000;
22     else if(counter1 < 3'b011 && current_state==S1)
23         counter1 <= counter1 +1;
24     else
25         counter1 <= 3'b000;
26 end

```

```

28 //counter2
29 always @(posedge clk , posedge rst)
30 begin
31     if(rst==1)
32         counter2 <= 4'b0000;
33     else if(counter2 < 4'b1110 && current_state==S2)
34         counter2 <= counter2 +1;
35     else
36         counter2 <= 4'b0000;
37 end
38
39 //CS
40 always @(posedge clk , posedge rst)
41 begin
42     if(rst==1)
43     begin
44         current_state <= S0;
45     end
46
47     end
48     else
49         current_state <= next_state;
50 end
51
52 //Z <= 2'b00;
53 //next_state <= S0;
54 //OL
55 assign z = Z;
56 always @(*)
57 begin
58     if(current_state == S1)
59         counter <= counter1;
60     else
61         counter <= counter2;
62 end

```

```

64 //NL
65 always @(*)
66 begin
67     case(current_state)
68     S0 : begin
69         Z = 2'b00;
70         if(x == 1)begin
71             next_state <= S1;
72         end
73         else begin
74             next_state <= S0;
75         end
76     end
77     end
78
79     S1 : begin
80         Z = 2'b01;
81         if(counter1 < 3'b011)begin
82             next_state <= S1;
83         end
84         else begin
85             next_state <= S2;
86         end
87     end
88
89     S2 : begin
90         Z = 2'b10;
91         if(counter2 < 4'b1110)begin
92             next_state <= S2;
93         end
94         else begin
95             next_state <= S0;
96         end
97     end
98     end
99     default: next_state <= current_state;
100 end
101 endcase
102 end
103
104
105

```

Fig. 12 FSM Verilog code

Results:

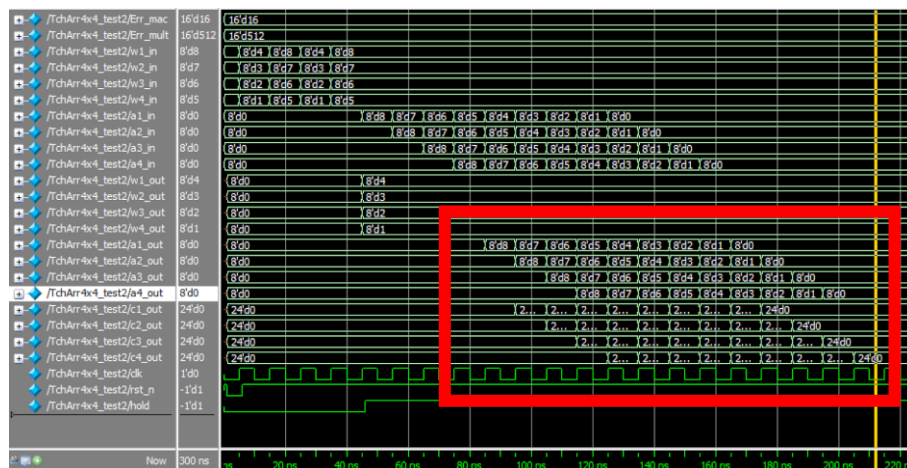


Fig. 13 simulation result

Reference:

- [1] J. Zhang, Z. Ghodsi, S. Garg and K. Rangineni, "Enabling Timing Error Resilience for

Low-Power Systolic-Array Based Deep Learning Accelerators," in IEEE Design & Test, vol. 37, no. 2, pp. 93-102, April 2020, doi: 10.1109/MDAT.2019.2947271.

[2] J. J. Zhang, K. Basu, and S. Garg. Fault-tolerant systolic array based accelerators for deep neural network execution. IEEE Design Test, 2019.