

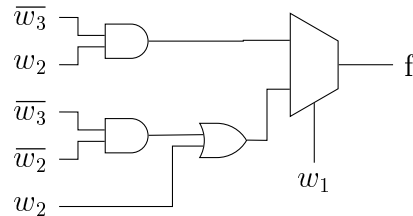
1. *Book Problems: 4.3, 4.21*

(4.3) Consider  $f = \overline{w_1}\overline{w_3} + w_2\overline{w_3} + \overline{w_1}w_2$ . Use the truth table to derive a circuit for  $f$  that uses a 2-to-1 multiplexer.

Note that the truth table for the above can be found as:

$w_1$	$w_2$	$w_3$	$f$
0	0	0	1
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	0
1	1	0	1
1	1	1	0

Which is equivalent to the below circuit, using  $w_1$  as the select in a multiplexer.



(4.21) Write Verilog code for an 8-to-3 binary encoder.

```

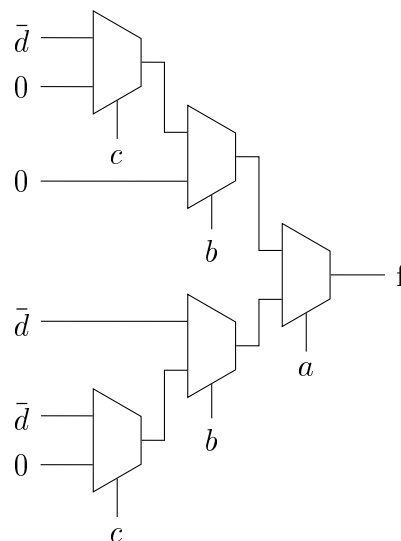
module Encode (b, f);
    input [7:0] b;
    output reg [2:0] f;
    always@(*)
    begin
        case(b)
            8'b00000001: f = 3'b000;
            8'b00000010: f = 3'b001;
            8'b00000100: f = 3'b010;
            8'b00001000: f = 3'b011;
            8'b00010000: f = 3'b100;
            8'b00100000: f = 3'b101;
            8'b01000000: f = 3'b110;
            8'b10000000: f = 3'b111;
        endcase
    end
endmodule
    
```

2. Implement the following circuits using only 2-to-1 multiplexers.

(a)  $f = \sum m(2, 5, 6, 14)$

Note that the truth table, and circuit visually derived from it, are as shown below:

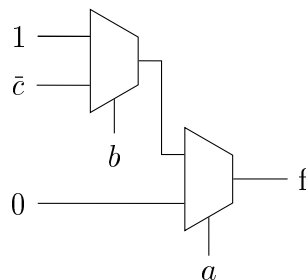
a	b	c	d	f
0	0	0	0	0
0	0	0	1	0
0	0	1	0	1
0	0	1	1	0
0	1	0	0	1
0	1	0	1	0
0	1	1	0	1
0	1	1	1	0
1	0	0	0	0
1	0	0	1	0
1	0	1	0	0
1	0	1	1	0
1	1	0	0	0
1	1	0	1	0
1	1	1	0	1
1	1	1	1	0



(b)  $f = \prod M(3, 4, 5, 6, 7)$

Note that the truth table, and circuit derived from it, are as shown below:

a	b	c	f
0	0	0	1
0	0	1	1
0	1	0	1
0	1	1	0
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	0



3. Convert the following decimal numbers to 32-bit floating point format.

(a) 33554430

First, note that  $(\frac{33554430}{2^{24}})_{10} = 1.99\dots$  and that

$$\begin{aligned} 33554430 &= 2^{24} + 2^{23} + 2^{22} + \dots + 2^3 + 2^2 \\ &= (111111111111111111111110)_2 \\ &= (1.1111111111111111111111)_2 * 2^{24} \end{aligned}$$

Which implies that  $Exponent = 24 \implies E = 24 - 127 = 151 = (10010111)_2$  for the floating point form below.

$$\begin{aligned} &S, E, E, E, E, E, E, E, E, 23bitsM \\ &(0\ 10010111\ 111111111111111111111111)_2 \end{aligned}$$

(b) 33554431

First, note that  $(\frac{33554431}{2^{24}})_{10} = 1.99\dots$  and that

$$\begin{aligned} 33554430 &= 2^{24} + 2^{23} + 2^{22} + \dots + 2^3 + 2^2 + 2^1 \\ &= (111111111111111111111111)_2 \\ &= (1.1111111111111111111111)_2 * 2^{24} \end{aligned}$$

Which implies that  $Exponent = 24 \implies E = 24 - 127 = 151 = (10010111)_2$  for the floating point form below.

$$\begin{aligned} &S, E, E, E, E, E, E, E, E, 23bitsM \\ &(0\ 10010111\ 111111111111111111111111\textcolor{red}{1})_2 \end{aligned}$$

However, there needs to be 24 bits in the Mantissa to accurately represent 33554431. Therefore, there is no way to accurately represent the given number without double precision or rounding.

4. Convert the following decimal numbers to fixed point unsigned binary with at least 8-bits of binary precision

(a) 12.45897

Note that  $(12)_{10} = (1100)_2$ , and  $(.45897)_{10} \approx \frac{1}{4} + \frac{1}{8} + \frac{1}{16} = (.0111)_2$ , which implies that  $(12.45897)_{10} \approx \boxed{(1100.0111)_2}$

(b) 0.333333

Note that  $.333333 \approx \frac{1}{4} + \frac{1}{16} + \frac{1}{64} + \frac{1}{256} = \boxed{(.01010101)_2}$

5. *For 32-bit Precision Floating point numbers, E=0x00 and E=0xFF are used for special numbers (like 0 and  $\infty$ ). What are the decimal values of the floating point numbers (32-bit) of smallest (non-zero) and largest (non-infinity) magnitude*

The smallest number would have E=0x1 and M=0x1, which results in a number of value  $(1 + 2^{-24}) * 2^{1-127} \approx \boxed{1.1754945 * 10^{-38}}$

Similarly, the largest number would have E=0xFE and M=b111...111, which results in  $(1 + 2^{-2} + 2^{-4} + .. + 2^{-24}) * 2^{127} \approx \boxed{3.4028235 * 10^{38}}$