

Solutions - Midterm Exam

(February 18th @ 7:30 pm)

Presentation and clarity are very important! Show your procedure!

PROBLEM 1 (17 PTS)

- Compute the result of the following operations. The operands are signed fixed-point numbers. The result must be a signed fixed-point number. For the division, use $x = 4$ fractional bits.

10.011 + 0.11111	1.010101 - 01.0101
1.10101 × 10.111	10.0110 ÷ 0.0111

$$\begin{array}{ccccccc}
 c_7 & = & 0 & & & & \\
 c_6 & = & 0 & & & & \\
 c_5 & = & 1 & & & & \\
 c_4 & = & 1 & & & & \\
 c_3 & = & 0 & & & & \\
 c_2 & = & 0 & & & & \\
 c_1 & = & 0 & & & & \\
 c_0 & = & 0 & & & & \\
 \\
 1 & 0 & 0 & 1 & 1 & 0 & 0 \\
 0 & 1 & 1 & 1 & 1 & 1 & + \\
 \hline
 1 & 1 & 0 & 1 & 0 & 1 & 1
 \end{array}$$

$$\begin{array}{ccccccc}
 1 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\
 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\
 \hline
 1 & 1 & 0 & 1 & 0 & 1 & 0 & 1
 \end{array} - \rightarrow$$

$$\begin{array}{ccccccc}
 c_8 & = & 1 & & & & \\
 c_7 & = & 1 & & & & \\
 c_6 & = & 1 & & & & \\
 c_5 & = & 1 & & & & \\
 c_4 & = & 1 & & & & \\
 c_3 & = & 0 & & & & \\
 c_2 & = & 0 & & & & \\
 c_1 & = & 0 & & & & \\
 c_0 & = & 0 & & & & \\
 \\
 1 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\
 1 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \\
 \hline
 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1
 \end{array} +$$

$$\begin{array}{l}
 1.10101 \times \rightarrow 0.01011 \times \rightarrow 1.0110 \times \\
 10.111 \quad 01.001 \quad 1001 \\
 \hline
 1 & 0 & 1 & 1 \\
 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 \\
 1 & 0 & 1 & 1 \\
 \hline
 0 & 1 & 1 & 0 & 0 & 0 & 1 & 1
 \end{array}$$

$$\begin{array}{l}
 \downarrow \\
 0.01100011
 \end{array}$$

✓ $\frac{10.0110}{0.0111}$: To unsigned and then alignment, $a = 4$: $\frac{01.1010}{0.0111} = \frac{01.1010}{0.0111} \equiv \frac{11010}{111}$

$$\begin{array}{r}
 000111011 \\
 111 \overline{)110100000} \\
 \underline{111} \\
 \underline{1100} \\
 \underline{111} \\
 \underline{1010} \\
 \underline{111} \\
 \underline{1100} \\
 \underline{111} \\
 \underline{1010} \\
 \underline{111} \\
 \hline
 11
 \end{array}$$

Append $x = 4$ zeros: $\frac{110100000}{111}$

Integer Division:

$$\begin{array}{l}
 Q = 111011, R = 11 \\
 \rightarrow Qf = 11.1011 (x = 4)
 \end{array}$$

$$\text{Final result (2C): } \frac{10.0110}{0.0111} = 2C(011.1011) = 100.0101$$

PROBLEM 2 (8 PTS)

- Represent these numbers in Fixed Point Arithmetic (signed numbers). Use the FX format [12.4].

✓ -17.125 ✓ 13.75

✓ -17.125:

$$17.125 = 010001.001 \rightarrow -17.125 = 101110.111 = 0xEE.E$$

✓ 13.75:

$$13.75 = 01101.11 = 0x0D.C$$

PROBLEM 3 (40 PTS)

- Perform the following 32-bit floating point operations. For fixed-point division, use 4 fractional bits. Truncate the result when required. Show your work: how you got the significand and the biased exponent bits of the result. Provide the 32-bit result.

<input checked="" type="checkbox"/> C3FA8000 - C1E00000	<input checked="" type="checkbox"/> D0D80000 + D0FA0000	<input checked="" type="checkbox"/> 80C00000xFAD00000	<input checked="" type="checkbox"/> 7B380000 ÷ C8A00000
---	---	---	---

✓ $X = \text{C3FA8000} - \text{C1E00000}$:

C3FA8000: 1100 0011 1111 1010 1000 0000 0000 0000
 $e + bias = 10000111 = 135 \rightarrow e = 135 - 127 = 8$
 $\text{C3FA8000} = -1.11110101 \times 2^8$

Significand = 1.11110101

C1E00000: 1100 0001 1110 0000 0000 0000 0000 0000
 $e + bias = 10000011 = 131 \rightarrow e = 131 - 127 = 4$
 $\text{C1E00000} = -1.11 \times 2^4$

Significand = 1.11

$X = -1.11110101 \times 2^8 + 1.11 \times 2^4 = -1.11110101 \times 2^8 + \frac{1.11}{2^4} \times 2^8$

$X = -(1.11110101 - 0.000111) \times 2^8$

We perform unsigned subtraction: $X = -1.11011001 \times 2^8$

$X = -1.11011001 \times 2^8, e + bias = 8 + 127 = 135 = 10000111$

$X = 1100 0011 1110 1100 1000 0000 0000 0000 = \text{C3EC8000}$

$$\begin{array}{cccccccccc}
 b_9 &= 0 & b_8 &= 0 & b_7 &= 0 & b_6 &= 0 & b_5 &= 1 & b_4 &= 1 \\
 &&&&&&&&&& b_3 &= 0 & b_2 &= 0 & b_1 &= 0 & b_0 &= 0 \\
 &&&&&&&&&& 1 & . & 1 & 1 & 1 & 0 & 1 & 0 & 1 \\
 &&&&&&&&&& 0 & . & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
 \hline &&&&&&&&&& 1 & . & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1
 \end{array}$$

✓ $X = \text{D0D80000} + \text{D0FA0000}$:

D0D80000: 1101 0000 1101 1000 0000 0000 0000 0000
 $e + bias = 10100001 = 161 \rightarrow e = 161 - 127 = 34$
 $\text{D0D80000} = -1.1011 \times 2^{34}$

Significand = 1.1011

D0FA0000: 1101 0000 1111 1010 0000 0000 0000 0000
 $e + bias = 10100001 = 161 \rightarrow e = 161 - 127 = 34$
 $\text{D0FA0000} = -1.1111010 \times 2^{34}$

Significand = 1.1111010

$X = -1.1011 \times 2^{34} - 1.1111010 \times 2^{34}$ (unsigned addition)

$X = -11.101001 \times 2^{34} = -1.1101001 \times 2^{35}$

$e + bias = 35 + 127 = 162 = 10100010$

$X = 1101 0001 0110 1001 0000 0000 0000 0000 = \text{D1690000}$

$$\begin{array}{cccccccccc}
 c_7 &= 1 & c_6 &= 1 & c_5 &= 1 & c_4 &= 1 & c_3 &= 1 & c_2 &= 0 & c_1 &= 0 & c_0 &= 0 \\
 &&&&&&&&&& 1 & . & 1 & 0 & 1 & 1 & 0 & 0 & 0 \\
 &&&&&&&&&& 1 & . & 1 & 1 & 1 & 1 & 0 & 1 \\
 \hline &&&&&&&&&& 1 & . & 1 & 1 & 0 & 1 & 0 & 0 & 1
 \end{array}$$

✓ $X = 80C00000 \times \text{FAD00000}$:

80C00000: 1000 0000 1100 0000 0000 0000 0000 0000
 $e + bias = 00000001 = 1 \rightarrow e = 1 - 127 = -126$
 $\text{80C00000} = -1.1 \times 2^{-126}$

Significand = 1.1

FAD00000: 1111 1010 1101 0000 0000 0000 0000 0000
 $e + bias = 11110101 = 245 \rightarrow e = 245 - 127 = 118$
 $\text{FAD00000} = -1.101 \times 2^{118}$

Significand = 1.101

$X = (-1.1 \times 2^{-126}) \times (-1.101 \times 2^{118}) = 10.0111 \times 2^{-8} = 1.00111 \times 2^{-7}$

$e + bias = -7 + 127 = 120 = 01111000$

$X = 0011 1100 0001 1100 0000 0000 0000 0000 = \text{3C1C0000}$

$$\begin{array}{cccccccccc}
 1 & . & 1 & 0 & 1 & x & & & & \\
 & & & & & 1 & . & & & \\
 & & & & & 1 & 1 & 0 & 1 & \\
 & & & & & 1 & 1 & 0 & 1 & \\
 \hline & & & & & 1 & 0 & 0 & 1 & 1
 \end{array}$$

✓ $X = 7B380000 \div \text{C8A00000}$:

7B380000: 0111 1011 0011 1000 0000 0000 0000 0000
 $e + bias = 11110110 = 246 \rightarrow e = 246 - 127 = 119$
 $\text{7B380000} = 1.0111 \times 2^{119}$

Significand = 1.0111

C8A00000: 1100 1000 1010 0000 0000 0000 0000 0000
 $e + bias = 10010001 = 145 \rightarrow e = 145 - 127 = 18$
 $\text{C8A00000} = -1.01 \times 2^{18}$

Significand = 1.01

$X = -\frac{1.0111 \times 2^{119}}{1.01 \times 2^{18}} = -\frac{1.0111}{1.01} \times 2^{101}$

$$\begin{array}{r}
 000010010 \\
 10100 \overline{\sqrt{101110000}} \\
 10100 \\
 \hline
 11000 \\
 10100 \\
 \hline
 1000
 \end{array}
 \quad \text{Alignment: } \frac{1.0111}{1.01} = \frac{1.0111}{1.0100} = \frac{10111}{10100} \\
 \text{Append } x = 4 \text{ zeros: } \frac{101110000}{10100} \\
 \text{Integer division: } Q = 10010 \rightarrow Qf = 1.001$$

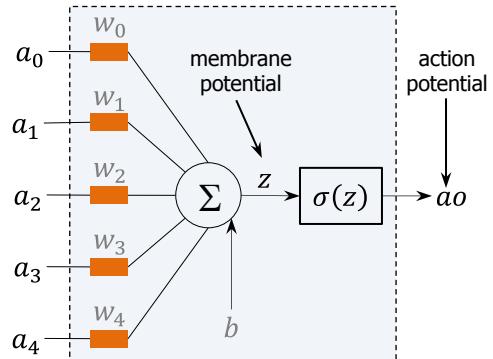
Thus: $X = -1.001 \times 2^{101}$, $e + bias = 101 + 127 = 228 = 11100100$
 $X = 1\textcolor{red}{111} \ 0010 \ 0001 \ 0000 \ 0000 \ 0000 \ 0000 = \text{F2100000}$

PROBLEM 4 (3.5 PTS)

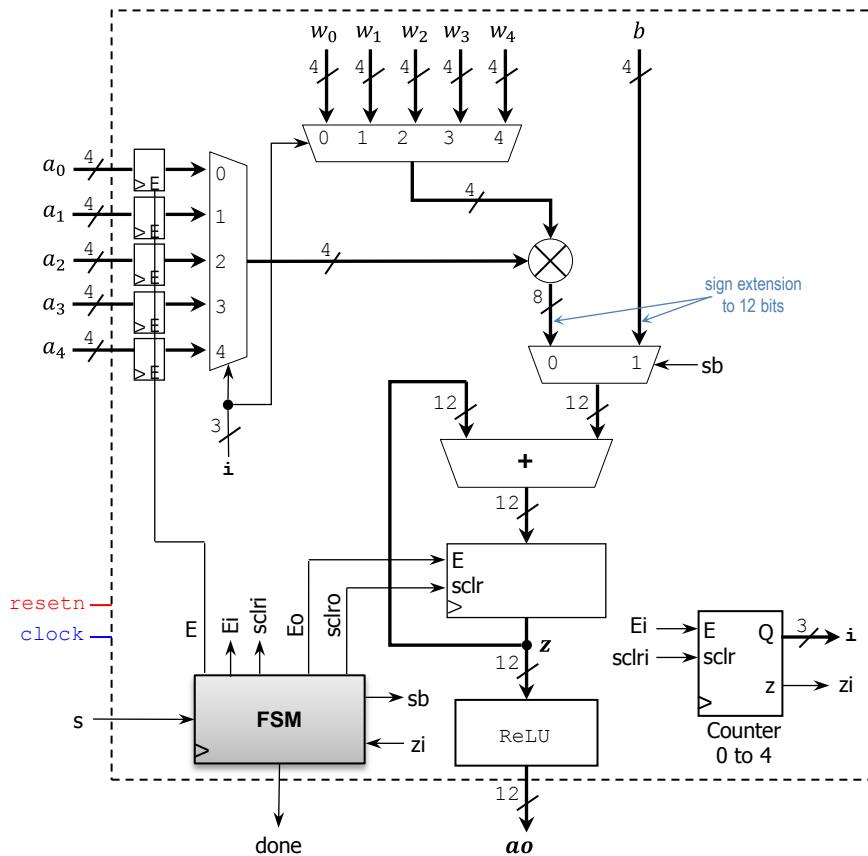
- Artificial neuron model.** The membrane potential z is a sum of products (input activations a_i by weights w_i) to which a bias term b is added. The action potential ao is modeled as a scalar function of z . The figure depicts a neuron with 5 inputs. The bias and the weights are constant values.

$$ao = \sigma \left(\sum_i a_i \times w_i + b \right)$$

- ✓ A popular and simple scalar function is the Rectified Linear Unit (ReLU):
 - $\sigma(z) = z$ if $z \geq 0$, otherwise $\sigma(z) = 0$.



- Digital System (FSM + Datapath): An iterative architecture for a 5-input neuron is depicted in next page. The circuit captures the input data (a_0, a_1, a_2, a_3, a_4) and then computes z using a multiply-and-accumulate approach (see iterative algorithm). The output ao is computed by applying the ReLU function to z .
 - ✓ All data is represented as signed integers:
 - Input activations (a_0, a_1, a_2, a_3, a_4), weights (w_0, w_1, w_2, w_3, w_4), bias (b): 4-bits wide.
 - Membrane potential (z) and action potential (ao): 12-bits wide (11 bits suffice, we select 12 for simplicity's sake).
 - ✓ Weights and biases: These are constant values (signed numbers represented as hexadecimals).
 - $w_0 = 0\times 4, w_1 = 0\times E, w_2 = 0\times C, w_3 = 0\times 5, w_4 = 0\times A$. These values appear in the timing diagram (next page).
 - Then $z = 4\times 4 + -2\times 1 + -4\times 2 + 5\times (-8) + -6\times 2 = -40 = 0\times FD8$. Finally, $ao = 0\times 000$
 - ✓ Counter 0 to 4: If $E=1, sclr=1$, then $Q \leftarrow 0$. If $E=1, sclr=0$, then $Q \leftarrow Q+1$. Also: $z=1$ if $Q=4$, else $z=0$.
 - ✓ Register: If $E=1, sclr=1 \rightarrow$ Clear. If $E=1, sclr=0 \rightarrow$ Load data.
 - ✓ 4×4 Signed Multiplier: This is a combinational circuit, whose result is 8-bits wide (it should be sign-extended to 12 bits).
 - ✓ ReLU: Combinational Block that implements the ReLU operation. For example, if $z = 0\times FD8$, then $ao = 0\times 000$
- Sketch the Finite State Machine diagram (in ASM form) given the algorithm. (20 pts.)
 - ✓ The process begins when s is asserted, at this moment we capture a_0, a_1, a_2, a_3, a_4 on the input registers. Then z is updated until the counter reaches its maximum value (4). The signal done is asserted when the final result is computed.
- Complete the timing diagram (z and ao are in hexadecimal format). (15 pts.)



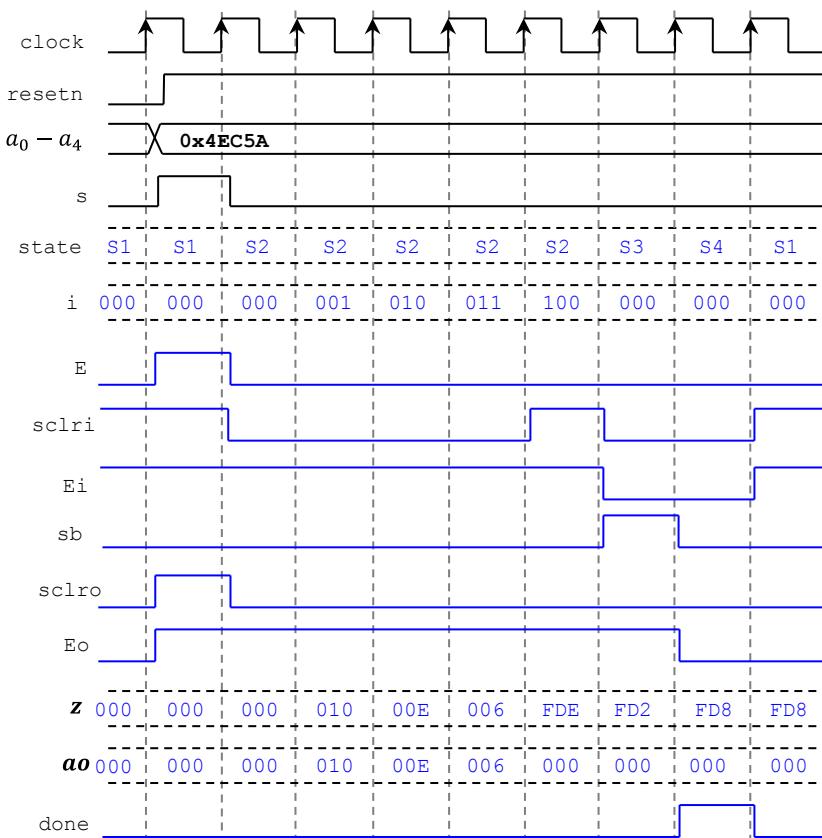
$$ao = \sigma_{ReLU} \left(\sum_{i=0}^4 a_i \times w_i + b \right)$$

ALGORITHM

```

 $z \leftarrow 0$ 
for  $i = 0$  to  $4$ 
     $z \leftarrow z + a_i \times w_i$ 
end
 $z \leftarrow z + b$ 
 $ao \leftarrow z$  if  $z \geq 0$ , else  $0$ 

```



FSM:

