

## CS M152A Lab 1

---

### Floating Point Conversion

*In this lab, you will learn how to use the Xilinx ISE program to design and test a floating point converter.*

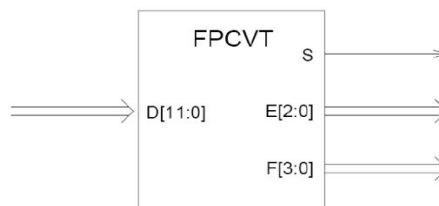
### Introduction

For this lab, you will use the Xilinx ISE software to design and test a combinational circuit that converts a 12-bit linear encoding of an analog signal into a compounded 8-bit Floating Point (FP) Representation.

This lab will be based on simulation only; no FPGA use will be involved. you are going to implement your design in Verilog HDL. At the end of the lab, you are expected to present design project with source code and test bench, and the design will be tested against a test bench that runs through all possible input patterns.

### Overview

The module for the floating-point conversion (called FPCVT) that you will be implementing is shown below:



The inputs and outputs of the FPCVT logic block are described in the following table:

| FPCVT Pin Descriptions |  |
|------------------------|--|
| D [11 : 0]             | Input data in Two's Complement Representation. D0 is the Least Significant Bit (LSB). D11 is the Most Significant Bit (MSB). |
| S                      | Sign bit of the Floating Point Representation.   |
| E [2 : 0]              | 3-Bit Exponent of the Floating Point Representation.   |
| F [3 : 0]              | 4-Bit Significand of the Floating Point Representation.  |

## Background

Analog signals are often converted to digital form for storage or transmission. A linear encoding using 8 bits can represent the unsigned number within the range 0 – 255 or a signed number within the range -128 – 127 using Two's Complement representation. Seven or eight bits of precision are adequate for intelligible speech or music almost good enough to listen. However, seven or eight bits do not provide sufficient dynamic range to capture both loud and soft sounds. Therefore, nonlinear encodings are used in most commercial systems.

## Output Format:

For this laboratory assignment, we will use a simplified Floating Point Representation consisting of a 1-Bit Sign Representation, a 3-Bit Exponent, and a 4-Bit Significand (the significand is sometimes called the *mantissa*).

|   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|
| 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
| S | E |   |   | F |   |   |   |

The value represented by an 8-Bit Byte in this format is:

$$V = (-1)^S \times F \times 2^E$$

The **S**-Bit signifies the **Sign** of the number. The 4-Bit **Significand**, **F**, ranges from [0000] = 0 to [1111] = 15, and the 3-Bit **Exponent**, **E**, ranges from [000] = 0 to [111] = 7. The following table shows the values corresponding to several Floating Point Representations.

| Floating Point Representation Examples |                  |       |
|--|------------------|-------|
| Floating Point Representation          | Formula          | Value |
| [0 000 0000]                           | $0 \times 2^0$   | 0     |
| [1 010 1010]                           | $-10 \times 2^2$ | -40   |
| [0 011 0111]                           | $7 \times 2^3$   | 56    |
| [0 010 1110]                           | $14 \times 2^2$  | 56    |

The last two rows of the above table demonstrate that some numbers have multiple Floating Point Representations. The preferred representation are the ones in which the Most Significant Bit of the Significand is 1. This representation is said to be the **normalized** representation. It is quite straightforward to produce the linear encoding corresponding to a floating-point representation; this operation is called *expansion*.

---

The goal of this laboratory assignment is to build a combinational circuit for the inverse operation, called *compression*. A device that performs both expansion and compression is called a compounder. The compression half of a compounder is more challenging because there are more input bits than output bits. As a result, many different linear encodings must be mapped to the same Floating Point Representation. Values that do not have Floating Point Representations should be mapped to the closest Floating Point encoding; this process is called *rounding*.

### Input Format:

| Leading Zeroes | Exponent |
|----------------|----------|
| 1              | 7        |
| 2              | 6        |
| 3              | 5        |
| 4              | 4        |
| 5              | 3        |
| 6              | 2        |
| 7              | 1        |
| $\geq 8$       | 0        |

The Significand consists of the 4 bits immediately following the last leading 0. When the exponent is 0, the Significand is the Least Significant 4 bits. For example,  $422 = [000110100110]$  has 3 leading zeroes, including the sign bit, thus its Exponent is 5 (according to the table above), and its Significand is  $[1101]$ . In case of a negative number which is in the 2's complement format, you should first negate it, and then count the number of leading zeroes. For example,  $-422 = [111001011010]$ , after negation it will become  $422 = [000110100110]$ , thus it also has 3 leading zeroes.

This FP representation expands to  $13 \times 2^5 = 416$ . The number 422 cannot be represented exactly, so it is represented with an error of about 1.5%.

### Rounding

The procedure presented above produces the correct Floating Point Representation for about half the possible linear encodings. However, it does not guarantee the most accurate representation. The circuit that you will design is required to round the linear encoding to the nearest Floating Point encoding. You should use the simple rounding rule that depends only on the fifth bit following the last leading 0. Recall that the first 4 bits following the last leading 0 make up the Significand. The next (fifth) bit then tells us whether to round up or down. If that bit is 0, the nearest number is obtained by truncation – simply using the first four bits. If, on the other hand, the fifth bit is 1, the representation is obtained by rounding the first four bits up – by adding 1.

The following table gives examples of rounding:

| Rounding Examples |                         |          |
|-------------------|-------------------------|----------|
| Linear Encoding   | Floating Point Encoding | Rounding |
| 000000101100      | [0 010 1011]            | Down     |
| 000000101101      | [0 010 1011]            | Down     |
| 000000101110      | [0 010 1100]            | Up       |
| 000000101111      | [0 010 1100]            | Up       |

The rounding stage of Floating Point conversion can lead to a complication. When the maximum Significand [1111] is rounded up, adding one causes an overflow. The result, 10000, does not fit in the 4-Bit Significand field. This problem is solved by dividing the Significand by 2, or shifting right, to obtain 1000, and increasing the Exponent by 1 to compensate.

For example:

000001111101 → 

|   |   |       |
|---|---|-------|
| 0 | 3 | 10000 |
|---|---|-------|

 OOPS! → 

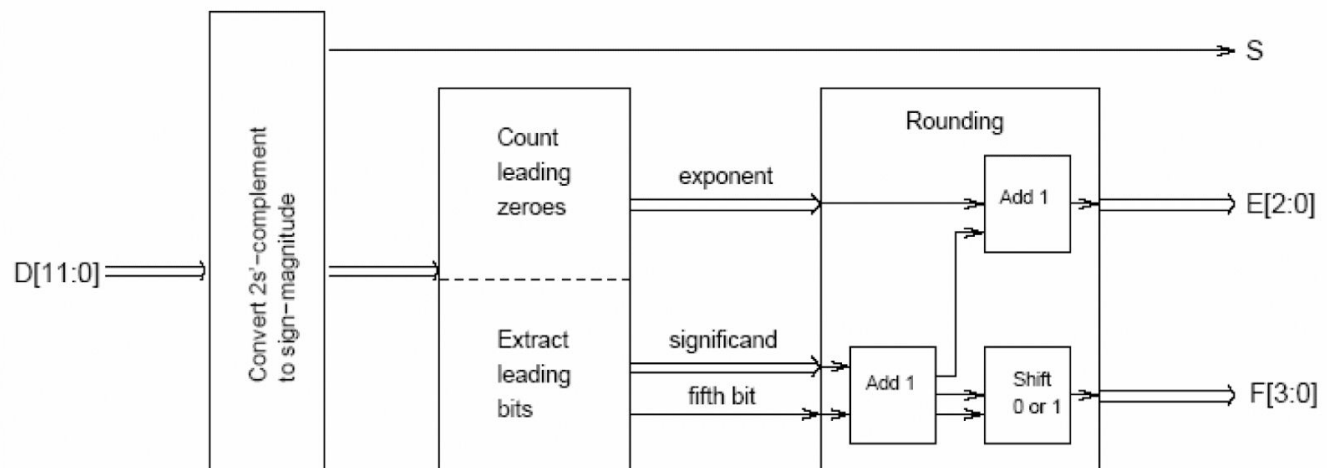
|   |   |      |
|---|---|------|
| 0 | 4 | 1000 |
|---|---|------|

In this example, 125 is converted to  $8 \times 2^4 = 128$ , which is indeed the closest Floating Point number. Note that the overflow possibility can be detected either before or after the addition of the rounding bit. **Which method is easier?**

When rounding very large linear encodings, such as  $2047 = [0111111111]$ , the exponent may be incremented beyond 7, to 8, which cannot be stored in the exponent field. Our solution to this problem is to use the largest possible Floating Point Representation.

## Overall Design

An overall block diagram for the floating-point conversion circuit is shown below:



The first block converts the 12-bit two's-complement input to sign-magnitude representation. Nonnegative numbers (sign bit 0) are unchanged, while negative numbers are replaced by their absolute value. As you should know, the negative of a number in twos'-complement representation can be found by complementing (inverting) all bits, then incrementing (adding 1) to this intermediate result. One problem case is the most negative number,  $-2048 = (100000000000)$ ; when complement-increment is applied, the result is  $100000000000$ , which looks like  $-2048$  instead of  $+2048$ .

The second block performs the basic linear to Floating Point conversion. The Exponent output encodes the number of leading zeroes of the linear encoding, as shown in table A above. To count the leading zeroes, we recommend you implement a **priority encoder**. The Significand output is obtained by right shifting the most significant input bits from bit positions 0 to 7. What this means is that each bit of the Significand comes from one of 8 possible magnitude bits.

The third block performs rounding of the Floating Point Representation. If the fifth bit following the last leading 0 of the intermediate Floating Point Representation is 1, the Significand is incremented by 1. If the Significand over flows, then we shift the Significand right one bit and increase the Exponent by 1.

## Deliverables

When you finish, the following should be submitted for this lab:

1. Project Code: the Xilinx ISE project folder should be cleaned up (*Project > Cleanup Project Files*), zipped and uploaded in the corresponding assignment page on the course website. **The TAs will run your design against a golden test case to test the correctness of your design**

2. Lab Report (Electronic Version): the lab report should be uploaded in the corresponding assignment page
3. Lab Report (Paper Version): the paper version of the lab report should be printed out on both sides and handed in on assigned date



