

# Design and Implementation of Binary to IEEE754 64-bit Floating Point Converter using Verilog

Srinivas v<sup>1</sup> Harshitha B<sup>2</sup>

<sup>1</sup>M Tech. Student <sup>2</sup>Assistant Professor

<sup>1</sup>Electronics <sup>2</sup>Dept. of E&C

<sup>1,2</sup>B M S College of Engg, Bangalore, India.

**Abstract**— Floating point numbers are used in many applications due to their dynamic representation capabilities in digital signal processing. The main objective of the work is to design and implement a binary to IEEE 754 floating point converter for representing 64 bit double precision floating point values. The converter at the input side of the floating point unit module helps to improve the overall design and reduces the complexity of manually converting the binary numbers into IEEE 754 format numbers. Hence the method solves the problems in occurrence of human errors to an extent in direct feeding of IEEE 754 format numbers to the FPU module.

The modules are written using Verilog and are then synthesized for Xilinx Virtex 5 FPGA using Xilinx Integrated Software Environment(ISE) design suite 14.2

**Key words:** Floating point, IEEE754, Double Precision.

## I. INTRODUCTION

The demand for floating point arithmetic operations in most of the commercial, financial and internet based applications is increasing day by day. Hence it becomes essential to find out an option to feed binary numbers directly as input for these applications. This helps in saving time and is much easier. In the current scenario, this is not possible, because, in the floating point unit, inputs should be given in IEEE 754 format i.e[3]. the binary inputs cannot be given as such, but it needs to be converted to the sign, exponent and mantissa form, about which, will be described in detail later. Hence the implemented binary to floating point converter for double precision bits will solve this issue to an extent.

The fundamental difference between fixed and floating point digital signal processors (DSPs) is their respective numeric representation of data. While fixed point hardware performs strictly integer arithmetic, floating point DSPs support integer or real arithmetic, the latter normalized in the form of scientific notation. A 64-bit, binary floating point DSP, supporting industry-standard, double precision operations, provides greater accuracy and greater precision than fixed point devices due to its wider word width, exponentiation and exact internal representations of data. Fixed point devices had to implement real arithmetic indirectly through software routines which add algorithmic instructions and development time, while with floating point format, real arithmetic could be coded directly into hardware operations. So, this thesis emphasizes on utilizing the capabilities of floating point format. The binary input given will range from 0-2048 bits, which is the maximum input range that can be provided which will satisfy the exponent range in the 64 bit IEEE 754 double precision format.

The document is organized as follows: Section II summarizes the important aspects of IEEE 754 single precision format. Section III describes binary to floating

point conversion. Section IV presents the results. Section V concludes the thesis.

## II. IEEE FLOATING POINT REPRESENTATION

The Institute of Electrical and Electronics Engineering (IEEE) issued 754 standards for binary floating point arithmetic in 1985[1,2,4,6]. This standardization was needed to eliminate computing industry's arithmetic vagaries. Later revisions were made to the existing standard in 2008. There are five basic formats—three binary floating point formats and two decimal formats.

The first two binary formats are called 'Single precision' and 'Double precision' respectively. IEEE double precision format alone is considered in this work. Hence it is described below.

In IEEE 754-2008, the 64-bit, base 2 format is officially referred to as binary 64. Double precision format uses 1-bit for sign bit, 11-bits for exponent and 52-bits to represent the fraction as shown in Fig 1.

Sign	Exponent	Mantissa
63	62..... 52	51..... 0

Fig 1: IEEE 754 double precision format

The double precision floating point number is calculated as  $(-1)^s \times 1.F \times 2^{(E-1023)}$ . Sign bit determines the sign of a number, which is either 0 for a non-negative number or 1 for a negative number. For IEEE double precision format, a bias of 1023 is added to the actual exponent.

The IEEE 754 standard specifies some special values like positive infinity, negative infinity, positive zero, negative zero and Not a Number (NaN). The standard also specifies the following rounding modes also like: Round to nearest, ties to even; Round to nearest, ties away from zero; Round toward positive infinity; round toward negative infinity and round toward zero. These special cases are considered in this work.

## III. BINARY TO FLOATING POINT CONVERSION

Our proposal is to convert the floating point number into a simple yet quite precise integral representation and perform the calculations on the same, followed by the final conversion of the output into its expected floating point result format.

The floating point data is inputted in two parts. The first part is a 64 bit binary value of the integer part of the floating point operand and other is a 64 bit binary value of fractional part of the floating point operand. This is done because Verilog cannot deal with floating point numbers. So we need to consolidate the two parts (integral and fractional) of the operand into a single 64 bit effective operand.



was simulated using Modelsim SE 6.1 and was targeted towards Vertex 5 FPGA.

#### REFERENCES

- [1] J. D. Bruguera and T. Lang. "Leading one prediction with concurrent position correction" IEEE Transactions on Computers, 48(10):1083-1097, Oct 1999.
- [2] A. Y. Duale, M. H. Decker, H.-G. Zipperer, M. Aharoni, and T. J. Bohizic. "Decimal floating-point in z9: An implementation and testing perspective". IBM Journal of Research and Development, 51(1/2):217-228, 2007.
- [3] W. Haller, U. Krauch, T. Ludwig, and H. Wetter. "Combined binary/decimal adder unit". U.S. Patent 5,928,319, July 1999.
- [4] E. Hokenek and R. K. Montoye. "Leading-zero anticipator (LZA) in the IBM RISC System/6000 floating-point execution unit". IBM Journal of Research and Development, 34(1):71-77, 1990.
- [5] Charles Farnum, "Compiler Support for Floating-Point Computation" Software Practices and Experience," pp. 701-9 vol. 18, July 1988.
- [6] D. Goldberg, "What every computer scientist should know about floating-point Arithmetic," pp. 5-48 in ACM Computing Surveys vol.23-1 (1991).

