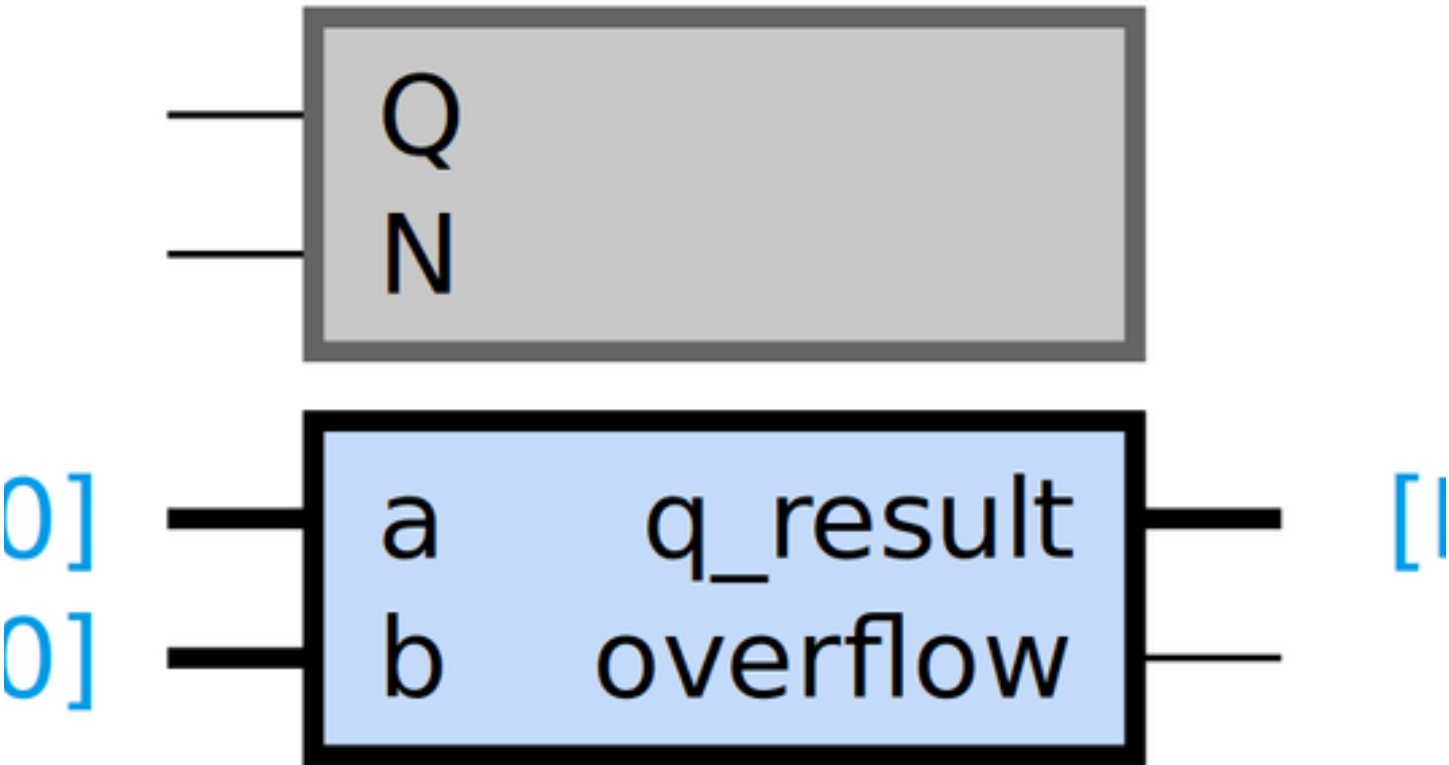


nt arithmetic to your design



ere we implement an FPGA based Convolutional Neural Network accelerator. However, the content of this article is useful for understanding the 'Fixed-po
u can understand the general principle and use the code even if you aren't interested in the entire series of articles.

ithmetic to your design

imating the Verification Process

neural network in hardware.

iber representation system?

In real life inputs and data, we need it to be able to interpret the various formats humans use to interpret this data. Real world data has everything from signals. However, when we work with hardware, all we see is registers and wires which are simply variables capable of holding a sequence of bits. All we get to is a parameter involved in defining these variables. As the author of this excellently written [whitepaper on fixed point arithmetic](#) puts it

“The meaning inherent in a binary word, although most people are tempted to think of them (at first glance, anyway) as positive integers. **However, the meaning is not inherent in the binary word, i.e., on the representation set and the mapping we choose to use.**”

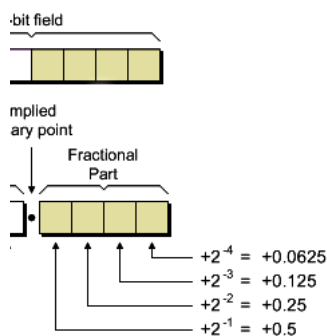
Fixed-point are simply that, they're representation sets and mapping schemes that have been widely used and standardized by the entire community. You could use it to interpret these binary words.

But in the open that can teach you about the basics of representing numbers on computers and the differences between fixed-point and floating-point numbers.

Fixed-point of floating-point numbers. It just recognizes fields with certain bit widths. This means that we need a layer above the hardware to help us convert at the same time help us interpret the results that the hardware is giving us by converting the raw binary words into the chosen format.

Using python mainly because of its simplicity and code readability.

Which the position of the decimal point remains fixed independent of the value that number is representing. This is what makes fixed point numbers easier to use than the floating point numbers. Fixed point arithmetic also uses much less resources in comparison to floating point arithmetic. Of course all of this comes at a cost of precision for a particular bit-width than fixed-point arithmetic. We take a hit in terms of the quantization noise when we use numbers in which the binary value is not exactly represented by that number.



Representing a fixed-point number

Implement fixed-point representation wherein, as shown in the above image, the first bit represents the sign bit and if the number is negative (sign bit = 1), be mindful of this when dealing with negative numbers.

Commonly denote the various parameters of a fixed point number. The most popular one is probably the **A(a,b)** format wherein **a** is the number of bits used to represent the integral portion of the number. Which means that the total number of bits used to store an A(a,b) fixed point number is **N = a + b**

and based on the *range* and *precision* that you need for the problem at hand. For example, if all or most of the numbers you're trying to represent have a very small range, you can use less bits to represent the integral portion and use more bits to represent the fractional portion thus achieving greater precision in your arithmetic. This is the application where your hardware will be used.

For fixed-point numbers, go through [this](#) document.

he fixed point arithmetic in Verilog, let's create a golden model i.e a model of the target functionality that gives use the correct target outputs which we can and fix any bugs based on the wrong outputs.

we use a function that allows us to convert numbers between the standard floating-point representation that python uses and the fixed point representation that we want our hardware to use. The function `float_to_fixed` converts a number in the standard *float* format to a user specified format. The variable *integer_precision* represents the number of bits to be used to represent the integer part of the number and the variable *fraction_precision* denotes the number of bits used for the fractional part (after the decimal point).

From hereon after, it represents the 'float' data type in python. Apparently, it is stored in the form of a 32-bit floating point number by the CPU under the hood.

```
def float_to_fixed(num, integer_precision, fraction_precision):
    """Convert a float to a fixed-point representation.

    Args:
        num: The float number to convert.
        integer_precision: The number of bits for the integer part.
        fraction_precision: The number of bits for the fractional part.

    Returns:
        A binary string representing the fixed-point number.
    """
    #sign bit is 1 for negative numbers in 2's complement representation
    sign = 1 if num < 0 else 0

    # Convert the absolute value of the number to a binary string
    abs_num = abs(num)
    integral_part = int(abs_num)
    fractional_part = abs_num - integral_part

    # Convert the integer part to binary
    integral_bin = bin(integral_part)[2:]

    # Convert the fractional part to binary
    fractional_bin = ''
    for i in range(fraction_precision):
        fractional_part *= 2
        bit = int(fractional_part)
        fractional_part = fractional_part - bit
        fractional_bin += str(bit)

    # Combine the integer and fractional parts
    fractional_bin = fractional_bin.ljust(fraction_precision, '0')
    bin_str = sign + integral_bin + fractional_bin

    # Pad the integer part to the specified precision
    integral_bin = integral_bin.zfill(integer_precision)

    return bin_str
```

Along with this code is another function that we use to generate the 2's complement of any binary number. Here is the code:

```
def twos_complement(bin_str, integer_precision, fraction_precision):
    """Generate the 2's complement of a binary string.

    Args:
        bin_str: The binary string to invert.
        integer_precision: The number of bits for the integer part.
        fraction_precision: The number of bits for the fractional part.

    Returns:
        The 2's complement of the binary string.
    """
    # Invert the bits
    inverted = ''
    for x in bin_str:
        inverted += str(1 - int(x))

    # Invert the sign bit
    inverted = inverted[1:] + str(1 - int(inverted[0]))

    # Pad the integer part to the specified precision
    inverted = inverted.zfill(integer_precision + fraction_precision)

    return inverted
```

Now that we've built:

```
float_to_fixed(110110110, 3, 12))

# FLOAT NUMBER TO FIXED-POINT
float_to_fixed(110110110, 3, 12))

# FLOAT NUMBER TO FIXED-POINT
float_to_fixed(110110110, 3, 12))
```

Now we can convert data from our fixed-point representation to a human readable *float* variable:

```
def fixed_to_float(bin_str, integer_precision, fraction_precision):
    """Convert a fixed-point binary string to a float.

    Args:
        bin_str: The binary string to convert.
        integer_precision: The number of bits for the integer part.
        fraction_precision: The number of bits for the fractional part.

    Returns:
        The float value of the binary string.
    """
    #s = input binary string
```

```

n - 1

twos_comp((s[1:]), integer_precision, fraction_precision)

s[1:]
precision + fraction_precision - 1):
complemented[j])*(2**i)

er

```

working:

```

01001001110', 3, 12))

1011111011', 3, 12))

```

tion, let's see what happens when we nest these functions together!

```

to_fp(-2.729, 3, 12), 3, 12))

```

How come the value is different when converted back and forth between fixed-point and *float* representation?

sion of the fixed-point representation in action. When we store a *float* variable in python, the CPU actually stores it in a 32-bit floating point number in its native format. This introduces a small loss of precision as the fixed-point representation does not have enough bits to achieve the same precision as the 32-bit floating point number.

Fixed-point arithmetic

Model ready, let's begin coding in Verilog!

Modified versions of the [fixed-point arithmetic library](#) from [OpenCores](#).

```

sign-bit + 3 integer-bits + 12 fractional-bits = 16 total-bits

```

```
|III|FFFFFFFFF|
,F) format would be A(3,12)
```

negative number in it's 2's complement form by default, all we
e two numbers together (note that to subtract a binary number
its two's complement)

your system (the software/testbench feeding this hardware with
y negative numbers in their 2's complement form,(some people
nitude as it is and make the sign bit '1' to represent negatives)
a look at the fixed point arithmetic modules at opencores linked

sign-bit + 3 integer-bits + 12 fractional-bits = 16 total-bits

```
|III|FFFFFFFFF|
,F) format would be A(3,12)
```

es

```
:0] a,
:0] b,
:0] q_result, //output quantized to same number of bits as the input
flow //signal to indicate output greater than the range of our format
```

assumption, here, is that both fixed-point values are of the same length (N,Q)
the results will be of length N+N = 2N bits
fies the hand-back of results, as the binimal point
n the same location

```
sult; // Multiplication by 2 values of N bits requires a
// register that is N+N = 2N deep
```

```
plicand;
plier;
mp, b_2cmp;
tized_result,quantized_result_2cmp;
```

```
-1},{(N-1){1'b1}} - a[N-2:0]+ 1'b1}; //2's complement of a
-1},{(N-1){1'b1}} - b[N-2:0]+ 1'b1}; //2's complement of b
```

```
= (a[N-1]) ? a_2cmp : a;
= (b[N-1]) ? b_2cmp : b;
```

```
= a[N-1]^b[N-1]; //Sign bit of output would be XOR of input sign bits
ltiplicand[N-2:0] * multiplier[N-2:0]; //We remove the sign bit for multiplication
ult = f_result[N-2+Q:Q]; //Quantization of output to required number of bits
ult_2cmp = {(N-1){1'b1}} - quantized_result[N-2:0] + 1'b1; //2's complement of quantized_result
0] = (q_result[N-1]) ? quantized_result_2cmp : quantized_result; //If the result is negative, we return a 2's co
//of the output value
_result[2*N-2:N-1+Q] > 0) ? 1'b1 : 1'b0;
```

hold the result of the multiplication between two **N - bit** numbers right? Right.

width as the multiplier and multiplicand?

that needs to be done at some point in the process in order to preserve the width of our data-path. If we keep increasing the size of the output at every stage datapath. This quantization adds something called the *quantization-noise* to our data since we're truncating the $2*N$ bit output to N bits. We should keep the truncation process. This is true most of the time if the numbers are pretty small and within the dynamic range of the fixed-point representation.

test benches can be found at the [Github Repo](#)

Source on 23 July, 2020.



Batman

I'm **Batman**, a silent nerd and a watchful engineer obsessed with great Technology. Get in touch via the [Discord community](#) for this site



Like what you are reading? Let me send the latest posts right to your inbox!

Subscribe

Free. No spam. Unsubscribe anytime you wish.