

# 1 Exercises: Fixed Point Numbers

## 1.1 Integral Numbers

There are two ways to interpret a bit vector as integral number: *unsigned* and *signed*, corresponding to the *IEEE* VHDL libraries *std\_logic\_unsigned* and *std\_logic\_signed*, resp.

**Unsigned interpretation:** A bit vector of  $w$  bits represents the integer range .....

**Signed interpretation:** A bit vector of  $w$  bits represents the int. range .....

## 1.2 Fixed Point Numerical Representation: The Q Number Format

**Unsigned:** **UQg.f** with  $g$  integral (deutsch: ganze) and  $f$  fractional bits. Width  $w=g+f$ .

**Signed:** **Qg.f** with 1 sign bit plus  $g$  integral and  $f$  fraction bits. Width  $w=1+g+f$ .

**Exercise:** The bit string **110.1011** can be interpreted...

... as UQ3.4 format representing

.....

... as Q2.4 delivering

.....

**Unsigned:** Range: Resolution:  $r =$

**Singed** Range: Resolution:  $r$

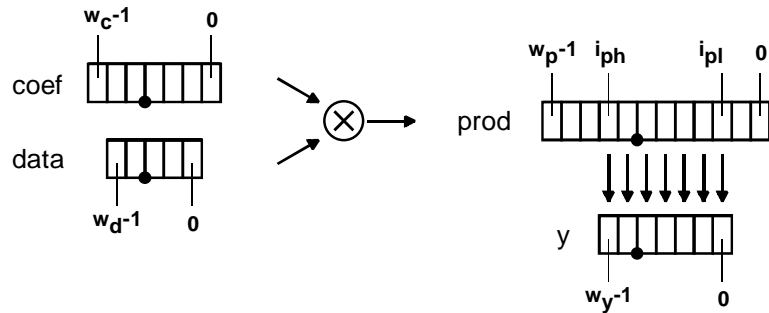
- You can append an arbitrary number of ..... after the point.
- 
- You can precede an arbitrary number of ..... before an unsigned number.
- 
- You can precede an arbitrary multiple of ..... before a signed number.

**Summation an subtraction** of fixed-point numbers is easy as they can be treated like integer numbers when they are written such that the points are over each other. Example:

Given numbers	Unsigned treatment	Signed treatment
11011011.11011...	11011011.11011...	11011011.11011...
± .....101.11101101	± .....101.11101101	± .....101.11101101

### 1.3 Multiplication of Fixed-Point Numbers

**Fig. 2.3:**  
Reducing the product length to the length of its factors with indices  $i_{ph}$  and  $i_{pl}$ .



We compute  $\mathbf{prod} = \mathbf{coef} * \mathbf{data}$  with  $\mathbf{coef}$  and  $\mathbf{data}$  having  $w_c$  and  $w_d$  binary places, respectively,  $f_c$  and  $f_d$  of them fractional.

Then  $\mathbf{prod}$  has ..... binary places, ..... of them fractional.

**Reducing the length of products:**  $\mathbf{y} = \mathbf{prod}(i_{ph} : i_{pl})$

Let  $\mathbf{coef}$  have  $w_c$  binary places,  $f_c$  of them fractional. Signal  $\mathbf{data}$  has  $w_d$  binary places,  $f_d$  of them fractional. The product has

$w_p = \dots$  binary places,  $f_p = \dots$  of them fractional.

Fig. 2.3 illustrates the multiplication of the coefficient  $\mathbf{coef}$  with  $w_c = \dots$ ,  $f_c = \dots$

and the data sample  $\mathbf{data}$  with  $w_d = \dots$ ,  $f_d = \dots$ . The product  $\mathbf{prod}$  has

$w_p = \dots$  binary places,

$f_p = \dots$  of them fractional.

We want to take result vector  $\mathbf{y}$  out of  $\mathbf{prod}$  preserving the point. For all bit vectors the LSB has index 0.

In Fig. 2.3  $\mathbf{y}$  has  $w_y = \dots$ , binary places  $f_y = \dots$  of them fractional.

To apply the VHDL command `y<=prod(iph DOWNTO ip1)` we have to compute

`ip1` = .....

`iph` = .....

### 1.4 Real → Binary Conversion

**Exercises** (for solutions see → chapter 8) :

Convert  $\pi=3.14159$  into a signed bit vector with 8 binary places, 4 of them fractional.

.....  
.....

Convert  $-\pi=-3.14159$  into a signed bit vector with 8 binary places, 4 of them fractional.

.....  
.....

## 2 Exercise Based on Executable VHDL

Listing 4: Code with gaps

```

(1)  LIBRARY ieee; USE ieee.std_logic_1164.ALL;
(2)  PACKAGE pk_filter IS
(3)      CONSTANT cDataInWidth:POSITIVE:=4;    -- Input-Data BitWidth
(4)      CONSTANT cDataInFract:POSITIVE:=2;    -- No of Input-Data fract. Bits
(5)      CONSTANT cDataOutWidth:POSITIVE:=5;   -- Output-Data BitWidth
(6)      CONSTANT cDataOutFract:POSITIVE:=3;   -- No of Output-Data fract Bits
(7)      CONSTANT cCoefWidth:POSITIVE:=4;     -- Coefficient's BitWidth
(8)      CONSTANT cCoefFract:POSITIVE:=2;     -- No of Coef's fractional Bits
(9)      SUBTYPE t_DataIn IS std_logic_vector(cDataInWidth-1 DOWNT0 0);
(10)     SUBTYPE t_DataOut IS std_logic_vector(cDataOutWidth-1 DOWNT0 0);
(11)     SUBTYPE t_coef IS std_logic_vector(cCoefWidth-1 DOWNT0 0);
(12) END PACKAGE pk_filter;
(13)
(14) LIBRARY ieee; USE ieee.std_logic_1164.ALL,
(15)     ieee.std_logic_signed."+", ieee.std_logic_signed."*";
(16) USE WORK.pk_filter.ALL;
(17) ENTITY TestBitslice IS
(18) END ENTITY TestBitslice;
(19)
(20) ARCHITECTURE rtl_TestBitslice OF TestBitslice IS
(21)     SIGNAL DataIn :t_DataIn;
(22)     SIGNAL coef   :t_coef;
(23)     SIGNAL DataOut:t_DataOut;

(24)     SIGNAL product:std_logic_vector(.....
.....

(25)     CONSTANT iPl:NATURAL:= .....
.....

(26)     CONSTANT iPh:NATURAL:= .....
.....

(27) BEGIN
(28)     DataIn  <= "0101", "0100" AFTER 10 ns; -- 1.25,   1.00 AFTER 10 ns
(29)     coef    <= "0101";                    -- 1.25
(30)     product <= coef * DataIn;              -- 1.5625, 1.25 AFTER 10 ns

(31)     DataOut <= product(iPh DOWNT0 iPl) .....
.....

(32) END ARCHITECTURE rtl_TestBitslice;

```

Correspondences with chapter 2.3:  $f_c=cCoefFract$ ,  $f_d=cDataInFract$ ,  $f_y=cDataOutFract$ ,  $w_c, w_d, w_p, w_y$ :  $cCoefWidth$ ,  $cDataInWidth$ ,  $cProdWidth$ ,  $cDataOutWidth$ , respectively.

### Exercises:

- Complete line (24) to get a *product* signal that fits to the multiplication of line (30).
- Compute *iPl* und *iPh* in lines (25), (26) to fit the bit-slice operation of line (31).
- Extend line (31) to get the bit-slice by bit-vector easy rounding.
- Verify the product, bit-slice and rounding operation of lines (39), (49) by hand.