



**POLITECNICO**  
MILANO 1863

---

# Optimize the Floating Point Unit shipped with the Mor1kx

[ Coding Project ]

---

**Student** Lorenzo Fumagalli  
**ID** 898398

**Student** Luca Guzla  
**ID** 898601

**Course** Embedded System 1  
**Academic Year** 2018-2019

**Advisor** Davide Zoni  
**Professor** William Fornaciari

March 10, 2019

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Problem statement . . . . .	3
1.2	Summary of the work . . . . .	3
<b>2</b>	<b>Design and implementation</b>	<b>5</b>
<b>3</b>	<b>Experimental evaluation</b>	<b>10</b>
3.1	Experimental setup . . . . .	10
3.2	Results . . . . .	10
<b>4</b>	<b>Conclusions and Future Works</b>	<b>10</b>

Table 1: Summary of clock cycles required for each operation in the original implementation

Operation	Clock cycles
Compare	1
Add/sub	6
Multiplication	7
Division	19
ITOF	4
FTOI	4
Rounding(already included in the count above except compare)	2

# 1 Introduction

In this project we have analyzed the behaviour of the implementation of an FPU [1] provided by our professor with the aim of optimize the number of cycles and the operating frequency.

## 1.1 Problem statement

The goal of the project was to optimize the FPU shipped with the mor1kx from the point of view of the number of clock cycles required by each operation. In particular, we were required to reduce the number of clock cycles needed to compute addition, subtraction, multiplication and conversion from integer format (with two's complement) to floating point IEEE 754 (32 bit), and viceversa.

In order to do this, we were given a working implementation and, starting from it, we tried to move pieces of computation in different stages, with the constraint of a minimum clock frequency of 140MHz (clock period of 7 ns).

The original implementation of the FPU was composed of 7 modules: add/sub, mul/div, i2f conversion, f2i conversion, compare, rounding and the top module. In the top module there were some computations, which were useful to prepare the computation of all the other modules. In particular, the top module checks which operation should be performed, splits the input signals in sign, exponent and mantissa, and computes some special cases (infinite, invalid, zero input). In particular, all the outputs of the modules related to the operations were given in input to the rounding module, which was the final module generating the three output signals of the FPU (fpu\_arith\_valid, fpu\_result, and fpcsr\_o). The initial operating frequency of the FPU was 100MHz, i.e. a clock period was 10ns.

Table 1.

## 1.2 Summary of the work

We considered 4 operations: floating point addition/subtraction, multiplication (splitting multiplication from division, which was not required by the project), and conversion, from integer to floating point IEEE754 representation and viceversa. We optimized these operations, reducing the number of needed cycles (from the rising edge of the new\_fpu\_data to the rising edge of the fpu\_arith\_valid\_o), having as a constraint the maximum clock period of 7ns. In order to do this, we eliminated some stages, anticipating computation in the previous stages, or delaying it in the following ones.

In particular:

- for the addition/subtraction we optimized from 5 to 4 clock cycles;
- for multiplication, we moved from 7 to 4 c.c;
- for the conversions, we moved from 4 to 3 cycles.

It was possible to reduce both the number of clock cycles and the clock period (from 10ns to 7ns) thanks to the fact that the implementation at 10ns had very large slack in each stage.

## 2 Design and implementation

In the following description we will use the following convention/notation, first we will provide a theoretic description of our implementation and then we will provide between () the name of the signals regarding the description. In the top module most of the computation are performed for both the operands and so we use && to separate the signal of input a and the signal of input b. The algorithm described can be found here [2].

### TOP

The top module is the module containing all the other modules, in the original version this module was used just to execute some common calculations to all the operations such as the detection of infinite or nan cases. In our modified version we moved part of the computation of the single operations inside it in order to met timing constraints. The top module receives as input the clock signal (clk), the reset (rst), the flush, which is used to cancel the current operation just like the reset (flush\_i), a signal that says if the operation can procede to the next stage of the computation in the top (padv\_decode\_i), a signal that says if the operation can procede to the next stage of the computation in the single module (padv\_execute\_i), a signal that is used to determine which type of operation the FPU must execute (op\_fpu\_i), a signal that describes the rounding mode that must be used (round\_mode\_i) and two 32 bits operand (opa\_i && opb\_i). As first operation we have a registering of top primary input, in this way we are able to avoid timing problems on the input, than we perform an analysis of input values: the 32 bits input are then splitted into sign, exponent and mantissa (for input a in\_signa, in\_expa, in\_fracta && in\_signb, in\_expb, in\_fractb) and this values are used to determine if we have a special case like infinte, Nan, denormalized or zero input. At this point we have some computation regarding the multiplication (and the future division), the integer to float conversion and the addition/subtraction. For the multiplication, after some precomputation on the exponent and on the mantissa in order to take in consideration the denormalized case, we calculate the number of leading zeros in the mantissa for both the operands (in\_nlza && in\_nlzb): these values are then used to shift the mantissa to obtain a mantissa with a 1 in the first position (in\_fract24a\_shl && in\_fract24b\_shl). We also calculate the value of the exponent of the multiplication (in\_exp10mux).

For the integer to float conversion we just calculate the module of the first input (in\_module\_a), since the second input is not used in this case, and we check if the input 'a' is zero (in\_opa\_0).

For the addition and subtraction we determine which operand has the greatest exponent and the greater mantissa (exp\_gt, fract\_gt) and in a similar way we determine if the exponent and the mantissa are equal (exp\_eq, fract\_eq); this values are used to determine which is the greatest input value (addsub\_agtb) in order to establish which is the first operand (in\_fract24\_nsh) and which is the second one (in\_fract24\_fsh). We also use this values to determine the difference between the exponent (in\_exp\_diff).

The output are then set taking in consideration which is the operation performed by the FPU and in particular we set if the operation is valid (fpu\_arith\_valid\_o), the actual result (fpu\_result\_o) and some useful flags (fpcsr\_o).

### I2F

The integer to float module receives as input the number to be converted into IEEE 754 floating point representation (32 bit opa\_i), the sign of the number (sign\_i) and a bit that signal if the input is zero (zero\_i), all the input are registered. At this point we computed the position of the first one in order to calculate the correct exponent representation (slt\_exp\_in). The value of the exponent is stored in

a 8 bit variable: if the input number is zero, the exponent is set to zero, otherwise it is set as the sum of the position of the first one plus the bias of the IEEE standard ( $s1t\_exp8$ ). At this point we calculate the value of the possible right and left shift; if  $s1t\_exp\_in$  is greater than 23 it means that the most significant one is in one of the first 8 bit of the input (from the left) and the right shift could be at most 8 and the exact value is calculated as  $(s1t\_exp\_in - 23)$  and the left shift is zero; if  $s1t\_exp\_in$  is less than 23, the most significant bit is in one of the last 23 bit of the input (from the left) the left shift could be at most 23 in this case and the exact value is calculated as  $(23 - s1t\_exp\_in)$  and the right shift is zero ( $s1t\_shr$  and  $s1t\_shl$ ). We extend the input adding 3 bit in order to allow the calculation of the guard, round and sticky bit ( $s1t\_fract35$ ).

At this point, using the value of the right shift we determine the sticky bit ( $s2r\_sticky$ ); since the sticky bit in case of left shift equals to zero the final value of the sticky bit is exactly the sticky for the right ( $s1t\_sticky$ ). The sticky bit is calculated as the or bit-a-bit of those bits that are "lost" with right shift. The last operation of this first stage is the calculation of the shifted input accordingly to the shifts previously calculated; the shifts are applied giving the priority to the right shift ( $s1t\_fract35sh$ ).

The following values are then registered to the next stage: the sign of the input, the exponent, the input shifted, the sticky bit and the rounding mode.

We determine some flags depending on the rounding mode received ( $rm\_nearest$ ,  $rm\_to\_zero$ ,  $rm\_to\_infp$ ,  $rm\_to\_infn$ ) and we assign the value received from the previous stage. Using the shifted input we calculate the guard, round, sticky and lost bit ( $s2t\_g$ ,  $s2t\_r$ ,  $s2t\_s$ ,  $s2t\_lost$ ) which are used to establish if the shifted input must be rounded up ( $s2t\_rnd\_up$ ).

We round up in one of the following situation:

- we round to the nearest and the sticky bit and the round bit are one or if the guard and the round and not the sticky;
- we round to the nearest and the sticky bit and the round bit are one or if the guard and the round and not the sticky;
- we round to the infinite negative and the lost bit is one and the sign of the input is negative.

The round is then applied to the first 32 bits of the shifted input ignoring the round bits ( $s2t\_fract32\_rnd$ ). The rounding phase could produce a variation of the exponent previously calculated, for this reason we update it taking in consideration the round ( $s2t\_f32\_exp10$ ); the same reason apply to the mantissa ( $s2t\_f32\_fract24$ ). We then compose the result as: 1 bit for the sign, 8 bit for the exponent and 23 bit for the mantissa ( $s2t\_result$ ) checking if the result is zero ( $s2t\_zero$ ). The result and the flags are then registered.

## F2I

The float to integer module receives as input the sign of the number ( $sign\_i$ ), the exponent ( $exp10a\_i$ ), the mantissa ( $fract24a\_i$ ), a bit that signals if we have a sNaN ( $snan\_i$ ) and a bit that signals if the input is a qNaN ( $qnan\_i$ ), all the input are registered.

As first operation we subtract to the exponent received as input the constant 150, in this way we remove the bias and move the binary point at the end of the mantissa ( $s1t\_exp10m$ ). Then we calculate the possible right shift ( $s1t\_shr\_t$ ) and we limit the possible shift up to 31 bits ( $s1t\_shr$ ). In a similar way we determine if a left shift for the mantissa is required and we limit it to 15 bits ( $s1t\_shl$ ). We check if the left shift is greater than ( $s1t\_is\_shl\_gt8$ ) or equal to 8 ( $s1t\_is\_shl\_eq8$ ); a left shift greater or equal than 8 will result in an overflow ( $s1t\_is\_shl\_ovf$ ). The number received as input is then extended to 35 bits ( $f2i\_int35\_t$ ) to take in consideration the shifts and, after a check on the right one to determine if it is zero ( $s1\_shr$ ), we apply the shifts to it; the shifts are applied giving the priority to the right shift

(s1t\_fract35sh). At this point, using the value of the right shift we determine the sticky bit (s2r\_sticky); the sticky bit is calculated as the or bit-a-bit of those bits that are "lost" with right shift. In an equivalent way we calculated the possible left shift (s1l\_sticky) and once both the two sticky bit are available we determine which one we will use in the rounding phase. Just like in the shifted input calculation we give the priority to the right one depending on the presence of the right shift (s1t\_sticky).

The following values are then registered to the next stage: the sign of the input verifying if we have a qnan or snan situation, the shifted number, two bit that represent the round and sticky bit, a bit that signals if the number is not valid, a bit representing the sNaN case, a bit representing the overflow and the rounding mode.

We determine some flags depending on the rounding mode received (rm\_nearest, rm\_to\_zero, rm\_to\_infp, rm\_to\_infm) and using the shifted input, the round and sticky bit we calculate the lost bit (s2t\_lost), which is used with the guard, the round and the sticky bit to establish if the shifted input must be rounded up (s2t\_rnd\_up).

We round up in one of the following situation:

- we round to the nearest and the sticky bit and the round bit are one or if the guard and the round and not the sticky;
- we round to the infinite positive and the lost bit is one and the sign of the input is positive;
- we round to the infinite negative and the lost bit is one and the sign of the input is negative.

The round is then applied to the first 32 bit of the shifted input ignoring the round bits (s2t\_fract32\_rnd).

At this point we calculate some useful flags used to determine if an invalid case occurred (s2t\_i32\_inv) and we calculate the two's complement in case of a negative number (s2t\_i32\_int32), further that we check if the output number is zero (s2t\_i32\_int32\_00). The output is then assigned checking if we are in an invalid situation (s2t\_i32\_opc).

The result and the flags are then registered.

## MULT

The mul\_fast module receives as inputs the mantissa shifted left for both operands 'a' and 'b' (fract24ash\_ii and fract24bsh\_ii), such that the '1.' is in the most significant bit, the sign of the operands (signa\_ii and signb\_ii), some flags to detect invalid cases (snan\_ii, qnan\_ii, anan\_sign\_ii, inv\_ii, inf\_ii), a bit that says if at least one of the two operands is zero (opc0\_ii) and the exponent of the result (exp10mux\_ii). All these inputs are immediately registered.

In the first stage, we set to 0 the operands for actual multiplication (s0t\_fract24a and s0t\_fract24b) and the exponent of the result (s0t\_exp10c) if at least one of the operands is 0. Then we decompose the two operands in high and low part (s0t\_mul16\_al, s0t\_mul16\_ah, s0t\_mul16\_bl, s0t\_mul16\_bh), which will be used to compute the partial products. At this point, we compute possible right shift for the exponent and the corrected exponent (s0t\_shr\_t and s0t\_exp10rx), considering different cases. The right shift is limited by 31 (s0t\_shrx). At this point we register some flags related to exceptions on input values (s0o\_inv, s0o\_inf\_i, s0o\_snan\_i, s0o\_qnan\_i, s0o\_anan\_sign\_i) and some signals related to the multiplication computation: the sign of the result (s0o\_signc), the exponent of the result (s0o\_exp10c), the exponent for right shift (s0o\_exp10rx), the right shift (s0o\_shrx) and the partial products obtained by multiplying the two parts of each operand (s0o\_fract32\_albl, s0o\_fract32\_albh, s0o\_fract32\_ahbl, s0o\_fract32\_ahbh).

In the next stage, we have the computation of the resulting mantissa on 48 bits (fract48), obtained by summing up the four partial products obtained from the previous stage. Starting from the mantissa on 48 bits, we reduce it to 28 bits, taking the 27 most significant bits and the or bit-a-bit of all the others

(s1t\_fract28). Furthermore, we check if we have a carry (s1t\_carry). At this point we compute actual right shift and actual exponent of the result (s1t\_shr, s1t\_exp10). Then, we apply a possible right shift, while we cannot have a left shift (s1t\_fract28sh), and we compute the sticky for no shift (s1l\_sticky). We determine some flags depending on the rounding mode received (rm\_nearest, rm\_to\_zero, rm\_to\_infp, rm\_to\_infm) and we assign the value received from the previous stage.

At this point we extend the fractional part to 32 bits (s1t\_fract32) for later calculations. We compute now the round (s1t\_r) and the guard bits (s1t\_g). Since the sticky computation is in the following stage, we compute some signals for both cases, i.e. sticky equal to 1 and equal to 0. Following this idea, we compute the lost bit in case of sticky equal to 0 (s1t\_lostS0), since in case of sticky equal to 1, the lost bit is the sticky itself. Other signals that are computed in the two versions are the 'round up', detecting if we have to round up.

For the case of sticky equal to 1 (s1t\_rnd\_upS1), we round up in the following cases:

- we round to the nearest and the round bit is 1;
- we round to the infinite positive and the sign of result is positive;
- we round to the infinite negative and the sign of result is negative.

For the case of sticky equal to 0 (s1t\_rnd\_upS0), we round up in the following situations:

- we round to the nearest and both the guard and round bits are 1;
- we round to the infinite positive, the sign of the result is positive and the lost bit in case of sticky equal to 0 is 1;
- we round to the infinite negative, the sign of the result is negative and the lost bit in case of sticky equal to 0 is 1.

Following the computation of the rnd\_up bits in both cases, we decide if we have to add 1 to the fractional part, computing the values needed for rounding (s1t\_rnd\_v32S1, s1t\_rnd\_v32S0). At this point, we perform a partial computation of the sticky bit, splitting the computation in two parts (s1t\_stickyh, s1t\_stickyl) and a signal telling us which part of the sticky we need (stickyfix): in case of a shift greater than or equal to 13, we will need both parts of the sticky, otherwise we will use only the first one.

Now we register some signals: flags related to exceptions on the input, the sign of the result, a signal telling us if a shift has been performed (s1o\_is\_shifted), the partial stickies, the exponent, the fractional part on 32 bit, the lost for sticky equal to 0 and the values needed for rounding

In the last stage, we compute the final sticky bit (s2r\_sticky) and the rounded fractional for both values of the sticky (s2t\_fract32\_rndS1, s2t\_fract32\_rndS0). At this point we can choose the actual fractional part on 32 bits of the result (s2t\_fract32\_rnd) depending on the sticky and the lost bit. Then we compute the final exponent on 10 bits (s2t\_f32\_exp10) and the fractional on 24 bits (s2t\_f32\_fract24). After this, we compute some useful flags and we compose the final result (s2t\_opc), giving priority to the exceptions. At the end of the module, these outputs are registered.

## ADD/SUB

This module receives as inputs the exponent of both the operands (exp10a\_i and exp10b\_i), some flags related to invalid inputs (snan\_i, qnan\_i, anan\_sign\_i, inv\_i, inf\_i), two signals which say if 'a' is greater than or equal to 'b' (addsub\_agtb\_i and addsub\_aeqb\_i), the mantissa of the two operands (fract24nsh\_i for the operand that doesn't need shift, and fract24fsh\_i for the operand that needs to



be shifted), the difference of the exponents (`exp_diff_i`). All these input are registered, forcing the two mantissa to zero in case of an infinite input.

In the first stage, we compute the needed right shift (`s2t_shr`), in order to have the same exponent for both the operands; then this shift is applied to the second operand after having extended it from 24 to 26 bits (`s2t_fract26_shr`). Using the extended operand (not shifted) and the right shift, we compute the sticky bit (`s2t_sticky`): we use this bit to build a 28 bit mantissa for the second operand (`s2t_fract28_shr`) starting from the shifted one on 26 bits. At this point, we can compute the addition (`s2t_fract28_add`): in case of a subtraction, we just add the first operand and the two's complement of the second one. We discard the least significant bit, taking only the first 27 bits (`s2t_fract27`). We register some signals: exceptions on the input, the sign of the result, the exponent, the mantissa on 27 bit, a flag saying if we have a subtraction and the result is 0 (`s2o_sub_0`), the sticky (`s2o_sticky`) and the round mode. Furthermore, we anticipate in this stage the final sticky for later computation in case of carry equal to 0 (`s2o_stickyFinalC0`).

In the next stage, the first thing we compute is the left shift needed to normalize the result and the exponent for left shift (`s3t_shl` and `s3t_exp10shl`). The sign of the result is "fixed" in case we have a subtraction with 0 result and we are rounding to minus infinite (`s3t_add_sign`), and we extend the mantissa from 27 to 35 bits adding some zeroes and the sticky computed from the previous stage. We detect if there is a carry by taking the most significant bit of the mantissa on 27 bit (`s3t_add_carry`): at this point, we can have right shift at most equal to 1 in case of a carry equal to one (`s3t_shr`), and we compute the corrected exponent for this case (`s3t_expPlusCarry`). Now we can compute the final sticky (`sticky`), using the previous computed one (`s2o_stickyFinalC0`) and the carry. We register the sign of the result, the fractional part on 35 bits (`s4o_fract35`), two conditions needed for next stage (`s4o_and_condition_4` and `s4o_and_condition_3`), right and left shift, the carry, the two possible exponents, the final sticky and the rounding mode.

In the last stage, we determine some flags depending on the rounding mode received (`rm_nearest`, `rm_to_zero`, `rm_to_infp`, `rm_to_infm`). We compute the final exponent (`s5t_exp10`) and the exponent plus one (`s5t_exp10PlusOne`) for later use. Using this calculated exponents, we compute two signals telling us if they are greater than 254 testing the bits (`condition_notPlusOne` and `condition_plusOne`). Now we apply the shift to the 35 bit fractional part, giving priority to the right shifts due to the carry (`s5t_fract35sh`). After this, we compute the guard, round, sticky and lost bits (`s5t_g`, `s5t_r`, `s5t_s` and `s5t_lost`) needed for rounding.

Given the flags for the rounding mode, we determine if we have to round up (`s5t_rnd_up`) in this cases:

- we round to the nearest and both round and sticky are 1;
- we round to the nearest, guard and round bits are 1 while the sticky is 0;
- we round to the infinite positive, the sign of the result is positive and lost is 1; we round to the infinite positive, the sign of the result is positive and lost is 1;
- we round to the infinite negative, the sign of the result is negative and lost is 1.

At this point we compute the rounded fractional on 32 bit (`s5t_fract32_rnd`). We compute the final exponent (`s5t_f32_exp10`), choosing between the two previously computed values, and the condition telling us if the exponent is greater than 254, again choosing between the two previously computed (`exponent_condition`). Then we take the final mantissa on 24 bit (`s5t_f32_fract24`) and check if we have a denormalized result (`s5t_f32_fract24dn`).

We calculate some flags regarding the result, we compose the final result (`s5t_opc`) and at the end we register the result

Table 2: Number of cycles of our implemented design

Operation	Clock cycles	Slack
Top, all operation together		+0.036
Add/sub	4	+0.009
Multiplication	4	+0.194
ITOF	3	+0.331
FTOI	3	+0.899

### 3 Experimental evaluation

In order to check the functional correctness we tested, through a behavioural simulation using Xilinx Vivado, all the operations with 30 test cases for each rounding mode, comparing the results and corresponding flags (fpu\_result\_o and fpcsr\_o) given by our new implementation of the FPU, with the results given by the original implementation, verifying that the two results coincide. Results can be found here [3]. Through Vivado, we used Synthesizer to run synthesis and implementation, in order to analyze the timing of the various operations. The following results refer to the post-implementation slack provided by Vivado. In particular, we analyzed the timing for the single operations and for our whole implementation of the Floating Point Unit.

#### 3.1 Experimental setup

The software used to develop the project was Xilinx Vivado2018a. The project has been implemented using an FPGA taken from the Artix-7 family, the xc7a100tcsg324-1, which is mounted on the Nexys4-DDR.

#### 3.2 Results

Here we provide the result of our design, in particular we will provide both the number of clock cycles for each implemented operation **Table 2.** and the slack obtained in post-implementation setting the top on the top module.

### 4 Conclusions and Future Works

We were able to optimize the four operations, reducing the number of clock cycles and the clock period. Future works will consist in a reimplementation of all the modules, the optimization of the floating point division and its integration in the FPU, together with the module for the compare operation.

## References

- [1] Rudolf Usselmann Andrey Bacherov, Jidan Al-eryani. Fpu shipped with mor1kx. 2014.
- [2] Luca Guzla Lorenzo Fumagalli. mor1kx code repository, <https://github.com/zeta0omega/embedded-system—p06—optimize-the-floating-point-unit-shipped-with-the-mor1kx>. 2019.
- [3] Luca Guzla Lorenzo Fumagalli. mor1kx code test, <https://github.com/zeta0omega/embedded-system—p06—optimize-the-floating-point-unit-shipped-with-the-mor1kx>. 2019.