



**HAL**  
open science

# Floating-point exponential functions for DSP-enabled FPGAs

Florent de Dinechin, Bogdan Pasca

► **To cite this version:**

Florent de Dinechin, Bogdan Pasca. Floating-point exponential functions for DSP-enabled FPGAs. International Conference on Field-Programmable Technology, Dec 2010, Beijing, China. pp.110-117. ensl-00506125

**HAL Id: ensl-00506125**

**<https://hal-ens-lyon.archives-ouvertes.fr/ensl-00506125>**

Submitted on 27 Jul 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Floating-point exponential functions for DSP-enabled FPGAs

*LIP research report RR2010-23*

Florent de Dinechin, Bogdan Pasca

*LIP (ENSL-CNRS-Inria-UCBL), École Normale Supérieure de Lyon*

*46 allée d'Italie, F-69364 Lyon, France*

{Florent.de.Dinechin, Bogdan.Pasca}@ens-lyon.fr

**Abstract**—This article presents a floating-point exponential operator generator targeting recent FPGAs with embedded memories and DSP blocks. A single-precision operator consumes just one DSP block, 18Kbits of dual-port memory, and 392 slices on Virtex-4. For larger precisions, a generic approach based on polynomial approximation is used and proves more resource-efficient than the literature. For instance a double-precision operator consumes 5 BlockRAM and 12 DSP48 blocks on Virtex-5, or 10 M9k and 22 18x18 multipliers on Stratix III. This approach is flexible, scales well beyond double-precision, and enables frequencies close to the FPGA's nominal frequency. All the proposed architectures are last-bit accurate for all the floating-point range. They are available in the open-source FloPoCo framework.

## I. INTRODUCTION

The exponential function is, after the basic arithmetic operators, one of the next most useful building block for floating-point applications. On FPGAs, it has been used for scientific or financial Monte-Carlo simulations, in phylogenetic tree reconstruction, in quantum chemistry, in the implementation of the power function among others.

### A. Previous works

Several publications have described exponential implementations. We list them here, and will discuss in more details the choices they made and their performance impact in IV-B.

For single precision, Doss and Riley [1] adapted to FPGAs a software algorithm based on floating-point operations. However, building a specific fixed-point architecture [2] proved more efficient. This architecture was later improved [3]. However, the table-based method used there doesn't scale up to double-precision.

As FPGAs are increasingly being used for double-precision applications, iterative architectures that scale better [4], [5], [6] were adapted for FPGAs [7]. This architecture was designed with 5-input LUTs in mind but is actually poorly suited to DSP-enabled FPGAs, as IV-B will show. It was parameterized in precision, but to our knowledge was never pipelined. Another pipelined, but double-precision only implementation was proposed in [8], [9].

In [10], a CORDIC-based approach using several parallel CORDIC cores was proposed. It has a complex control including input and output FIFOs. Being radix-2 CORDIC, it computes one digit per iteration and thus has a very long latency. Moreover, it is based on a floating-point adder,

whereas CORDIC is inherently a fixed-point computation, so there is probably room for improvement there.

From a user point of view, the current state of the art is probably the floating-point exponential function `ALTFP_EXP` provided with Altera Megawizard since 2008 [11]. This implementation is parameterized in exponent and mantissa size and fully pipelined. Being included in the standard Quartus releases, it is widely available, although only for Altera targets.

Many other publications have addressed the computation of exponential function in ASIC, e.g. [4], [5], [12], [6]. However, it is difficult to evaluate the relevance of such works on FPGAs.

### B. Contributions

In the present article, we propose yet another architecture for the floating-point evaluation of the exponential function, and its implementation in the open-source FloPoCo project<sup>1</sup>. Its main specificities are the following.

- The algorithm, based on the usual multiplicative range reduction followed by a polynomial approximation, was chosen with DSP blocks and embedded memories in mind, so it makes efficient use of these resources. For instance, the single-precision version now involves just one 17x17-bit multiplier and 18Kbits of dual-port memory, and runs at 375MHz on a Virtex-4, which is a large improvement in all respects over the state of the art [3].
- As we believe that floating-point on FPGA should exploit the flexibility of the target and therefore not be limited to IEEE single and double precision, the algorithm and implementation we propose are fully parametrized in exponent and mantissa size.
- They scale to double-precision and beyond.
- They are last-bit accurate for all supported mantissa sizes.
- The implementation is pipelined to a user-specified frequency. We are able to generate operators working at a frequency close to the DSP nominal frequency, as well as lower-frequency, lower-resource versions.
- The architectures are generated as synthesizable VHDL portable to any FPGA target, although many target-specific optimizations are
- A novel variation of the KCM algorithm, for multiplying a real constant by an integer, is used.

<sup>1</sup><http://www.ens-lyon.fr/LIP/Arenaire/Ware/FloPoCo/>

- All this work is freely available from the FloPoCo project. It is included in the open SVN repository [https://gforge.inria.fr/scm/?group\\_id=1030](https://gforge.inria.fr/scm/?group_id=1030) and will be part of releases 2.1.0 and following. It comes with the test vector generation framework of FloPoCo [13]. In general, it should be immediately usable for application designers.

Section II gives an overview of the algorithm used, and Section III discusses some implementation choices. Section IV compares implementation results with the literature, and Section V concludes.

## II. A FLOATING-POINT EXPONENTIAL

### A. Special cases

The exponential function is defined on the set of the reals. However, in this floating-point format, the smallest representable number is

$$X_{\min} = 2^{-E_0}$$

and the largest is

$$X_{\max} = (2 - 2^{-w_F}) \cdot 2^{2^{w_E} - 1 - E_0}.$$

The exponential should return zero for all input numbers smaller than  $\log(X_{\min})$ , and should return  $+\infty$  for all input numbers larger than  $\log(X_{\max})$ . In single precision ( $w_E = 8$ ,  $w_F = 23$ ), for instance, the set of input numbers on which a computation will take place is  $[-88.03, 88.72]$ . In addition, as for small  $x$  we have  $e^x \approx 1 + x + x^2/2$ , the exponential will return 1 for all the input  $x$  smaller than  $2^{-w_F - 2}$ .

One consequence is that the testing of a FP exponential operator should concentrate on numbers between  $X_{\min}$  and  $X_{\max}$ . In FloPoCo, we use random inputs with exponents restricted to  $[-w_F - 3, w_E - 2]$ .

### B. Algorithm overview

The algorithm used is similar to what is typically used in software [14].

The main idea is to reduce  $X$  to an integer  $E$  and a fixed-point number  $Y$  such as

$$X \approx E \cdot \log 2 + Y \quad (1)$$

where  $Y \in [-1/2, 1/2)$  – we will see below in II-C how to ensure this enclosure.

We may then use the identity

$$e^X \approx 2^E \cdot e^Y \quad (2)$$

so  $E$  is almost the exponent of the result. Indeed, if  $Y \in [-1/2, 1/2)$ , we have  $e^Y \in [0.6, 1.7]$ . Thus the exponent and mantissa of the result are

$$\begin{cases} R = 2^E \cdot e^Y & \text{if } e^Y \geq 1 \\ R = 2^{E-1} \cdot (2e^Y) & \text{if } e^Y \leq 1 \end{cases} \quad (3)$$

This test boils-down in testing the most significant bit of  $e^Y$ , and the division by 2 is just a shift.

The architecture of this operator is given on Figure 1. This figure also explicits the alignment of the fixed-point data.

### C. Range reduction

To implement equation (1), we have to implement an approximation of

$$E = \left\lfloor \frac{X}{\log 2} \right\rfloor \quad (4)$$

then

$$Y = X - E \times \log 2. \quad (5)$$

If computed infinitely accurate, this would ensure  $Y \in [-\frac{\log 2}{2}, \frac{\log 2}{2}]$ . On one hand, this is not ideal from an architectural point of view, as  $Y$  will be input to a table and  $\frac{\log 2}{2}$  is not a power of two (as  $\log 2 \approx 0.34$ , the next power of 2 is  $1/2$ , so only 69% of the table would be used). On the other hand, implementing (4) and (5) accurately enough would be expensive. A solution to both problems is therefore a relaxed implementation of (4) that will save on the computation of (4) and (5) while ensuring  $Y \in [-1/2, 1/2)$ . The idea is that the computation of  $E$  can be grossly approximate, as long as (5) is accurately implemented. A normalization process (see below (3)) will take care of the cases where  $E$  was not directly computed as the exact result exponent.

As (4) and (5) are inherently fixed-point computations, the first task is to build a fixed-point representation  $X_{\text{fix}}$  of the input  $X$ . The most significant bit (MSB) of this representation is provided by the condition  $X > \log(X_{\max}) \rightarrow \exp(X) = +\infty$ , from which we deduce  $X > 2^{w_E + 1} \rightarrow \exp(X) = +\infty$ . The MSB of  $X_{\text{fix}}$  should therefore have weight  $w_E$ . The least significant bit is provided by the condition  $X < 2^{-w_F - 2} \rightarrow \exp(x) = 1$ , which defines a LSB of weight  $-w_F - 2$ . Actually, we will improve this accuracy to  $-w_F - g$  with  $g = 3$  (see below in III-A) to allow for rounding error accumulation in these  $g$  guard bits.

Thus the *shift to fixed point* box on Figure 1 shifts the mantissa by the value of the exponent. More specifically, if the exponent is positive, it shifts to the left by up to  $w_E$  positions (more means overflow). If the exponent is negative, it shifts to the right by up to  $w_F + g$  positions. This box also generates out-of-range signals.

Let us now turn to the relaxed computation of  $E$ . Since  $E$  is almost the final exponent (of size  $w_E$ ), its size in bits will be  $w_E + 1$ , including one sign bit, the  $+1$  preventing overflow in the second case of (3). Expliciting all the roundings and truncations, (4) becomes

$$E = \left\lfloor [X_{\text{fix}}]_{w_E + g'} \times \left\lfloor \frac{1}{\log 2} \right\rfloor_{w_E + g'} \right\rfloor_{w_E + 1}. \quad (6)$$

The value  $g' = 3$  ensures that the product is computed with a relative accuracy of  $2 \times 2^{-w_E - 3}$  with respect to  $\frac{X}{\log 2}$ . Other terms  $E$  may off by 1 with respect to the ideal  $\left\lfloor \frac{X}{\log 2} \right\rfloor$  when  $X$  is within  $1/4$  of the middle between two multiples of  $\log 2$ .

Then, (5) may be implemented as

$$Y = X_{\text{fix}} - E \times \log 2. \quad (7)$$

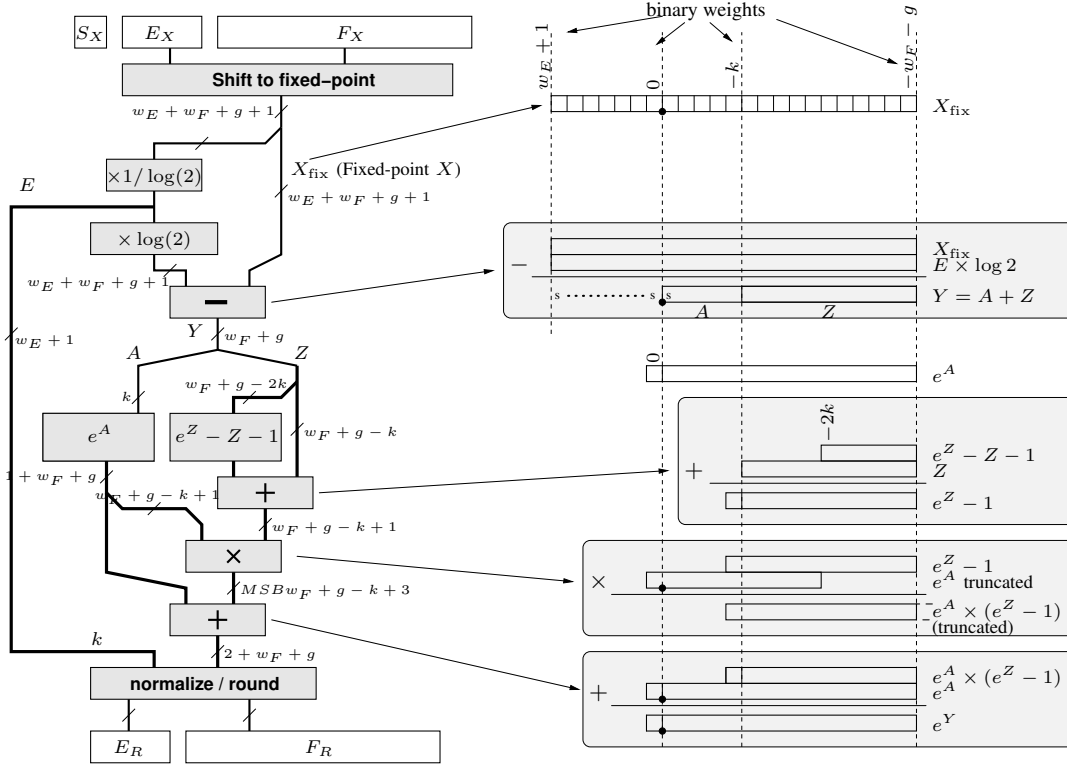


Fig. 1. Architecture and fixed-point data alignment

The previous error analysis now ensures that the bound on  $Y$  is now  $[-\frac{\log 2}{2} - \frac{\log 2}{8}, \frac{\log 2}{2} + \frac{\log 2}{8}] \approx [-0.44, 0.44]$ . This computation will cancel the integer part and the first bit of the fractional part.

In this work, we have also considered reducing to  $Y \in [0, 1)$  instead of  $Y \in [-1/2, 1/2)$ . It turns out that guaranteeing this enclosure, especially  $Y \geq 0$ , is more expensive.

#### D. Constant multiplications

As both constant multiplications (by  $1/\log 2$  and  $\log 2$ ) multiply a large constant by a small input, we use the KCM algorithm [15]. For the larger multiplication by the real value  $\log 2$ , we actually use a variation that is original to our knowledge and that we briefly present now.

Let  $\alpha$  be the LUT input size of the target FPGA. The input (here  $E$ ) is split into chunks of size  $\alpha$ :

$$E = \sum_{i=0}^p 2^{i\alpha} E_i$$

therefore

$$E \log 2 = \sum_{i=0}^p 2^{i\alpha} E_i \log 2.$$

We tabulate in LUTs the products  $2^{i\alpha} E_i \log 2$  on just the required precision, so that its LSB has value  $2^{-w_F - g - \gamma}$  where  $\gamma$  is again a number of guard bits. Each table may hold the correctly rounded value of the product of  $E_i$  by the real value  $\log 2$  to this precision, so entails an error of  $2^{-w_F - g - \gamma - 1}$ .

Finally, the first table actually holds  $E_0 \log 2 + 2^{-w_F - g - 1}$  (i.e., it adds one half-ulp of the multiplier result), so that the truncation of the sum will correspond to a rounding of the product: this provides one half-ulp accuracy at no cost.

As  $\alpha \in \{4..6\}$  for current FPGAs, and practical values of  $E$  are smaller than 15, the value  $\gamma = 2$  is usually enough to ensure that this multiplier returns a faithful multiplication by  $\log 2$ .

With respect to the technique used by all the previous works (rounding  $\log 2$  to some fixed-point value, then using an integer constant multiplier, then truncating its result), this approach saves a few LUTs, but also one half ulp in the error analysis.

FloPoCo provides an implementation of this KCM with frequency-directed pipeline.

#### E. Computation of $e^Y$

Let us now turn to the computation of  $e^Y$ . We use a second range reduction, splitting  $Y$  as

$$Y = A + Z \quad (8)$$

where  $A$  consists of the  $k$  most significant bits of  $Y$ , and  $Z$  consists of the  $w_F + g - k$  least significant bits. Then we have

$$e^Y = e^{A+Z} = e^A \cdot e^Z. \quad (9)$$

Here  $e^A$  will be tabulated in a table indexed by  $A$ , and  $Z$  is small enough to enable us to use the Taylor formula

$$e^Z \approx 1 + Z + Z^2/2 + \dots \quad (10)$$

This formula has the advantage that the three first coefficients are powers of two, therefore the corresponding multiplications can be mere shifts. Actually we define

$$f(Z) = e^Z - Z - 1 \quad (11)$$

From  $0 \leq Z < 2^{-k}$  and  $e^Z - Z - 1 \approx Z^2/2 + \dots$ , we know that the MSB of  $f(Z)$  has weight  $-2k - 1$ . As  $f(Z)$  will be added to  $Z$ , its LSB should have the same weight  $-w_F - g$ . The useful size of  $f(Z)$  is therefore  $w_F + g - 2k$ . As a consequence, we do not need to compute it out of all the bits of  $Z$ . Truncating  $Z$  to its  $w_F + g - 2k$  MSBs will entail an error of roughly the same weight as the error entailed by the fixed-point format of  $f(Z)$ .

Out of  $Z$  and  $f(Z)$ , we compute  $e^Z - 1 = f(Z) + Z$ . This addition may overflow, so the result is on  $w_F + g - k + 1$  bits, one more bit than  $Z$ .

If  $1 + w_F + g < 17$ , the final multiplication  $e^Y = e^A \cdot e^Z$  may be computed directly as a single DSP block. For larger precisions, the cost of this multiplication is reduced by implementing it as

$$\begin{aligned} e^A \cdot (1 + Z + f(Z)) \\ = e^A + e^A \cdot (Z + f(Z)) \end{aligned} \quad (12)$$

Again, the two addends have LSB weight  $-w_F - g$ . Again, the multiplier inputs need not be more accurate than their output, so we truncate  $e^A$  to its LSB  $w_F + g - k + 1$  bits.

As we need to truncate the result of this multiplier, we may as well use, for large precisions, truncated multipliers [16], to save DSP and latency.

A final normalization step possibly shifts left the mantissa by one bit, then performs the final rounding. The rounding consists in possibly adding one bit, then truncating. The IEEE-754 format has the nice property that we may use an adder of size  $w_E + w_F + 1$  to add the rounding bit to the concatenated exponent and mantissa: carry propagation from mantissa to exponent will handle the possible exponent change due to rounding up.

### III. IMPLEMENTATION ISSUES

#### A. Error analysis

This computation involves several approximation and rounding errors. The purpose of this section is to guarantee faithful rounding, *ie.* an error of less than one *unit in the last place* (ulp) of the result. Here the ulp has the value  $2^{-w_F+g}$ .

In the following, all the errors will be expressed in terms of unit in the last place of  $Y$ . Thus errors expressed this way can be made as small as required by increasing  $g$ .

First, note that the argument reduction is not exact. As already stated, numerical errors in the computation of  $E$  (6) mostly impact the range of  $Y$ . Concerning the computation of  $Y$  (1), there are two exclusive cases:

- If  $X$  is large (its exponent is larger than  $-2$ ), its mantissa is shifted without loss of information, then the computation of  $E \times \log 2$  introduces at most one ulp of error in  $Y$  as seen in II-D.

- Or,  $X$  is small, its mantissa is shifted right beyond the ulp, so its LSBs are lost, which also entails one error of one ulp in  $Y$ . However, in this case  $E = 0$ , so the computation of  $E \times \log 2$  is exact.

In both cases we may thus have an error of at most one ulp on  $Y$ . Let us now see how it propagates to  $e^Y$ .

$e^A$  is tabulated rounded to the nearest, thus with an error of 1/2 ulp.

$e^Z - Z - 1$  is either tabulated (1/2 ulp) or evaluated through polynomial approximation (1 ulp). As the higher order bits of  $Z$  are used, the error on  $Y$  (which is the error on  $Z$ ) is scaled down and becomes negligible.

Then  $e^Y - 1$  adds the error on  $Z$  and the error on  $e^Z - Z - 1$ , and thus holds an error of 1.5 or 2 ulps.

The error on the other input to the multiplier ( $e^A$  truncated) is of one ulp. The product adds these error as  $(a+\epsilon) \times (b+\epsilon') = ab + b\epsilon + a\epsilon' + \epsilon\epsilon'$ . Here is another subtlety. This formula shows that the error on  $e^Z - Z - 1$  is scaled by the value of  $e^A$ . Fortunately, the worst case error will occur for  $e^A < 1$ , since in this case the result will be shifted left by one bit. In the case  $e^A > 1$  the error on  $e^Z - Z - 1$  may be scaled up (by up to 1.6) but we will have in this case the extra bit of precision needed for the other case, so it doesn't matter.

Truncating the multiplier result would yields another error of one ulp, however we may instead round it (1/2 ulp only) at very little cost by adding its round bit to the right of  $e^A$ , so the addition of  $e^A$  will also compute the rounding of the product.

Finally the product holds an error of 3 or 3.5 ulps.

Adding the error on  $e^A$ , we deduce that the error on  $e^Y$  may be up to 3.5 ulp in the dual table case, and 4 ulp in the polynomial case.

If  $e^Y < 1$  the final 1-bit shift will multiply this error by 2, so we need 3 guard bits.

Previous works need more guard bits (5 guard bits in [2], 8 in [10] for instance), hence a wider datapath. This improvement in the present work is partly due to a finer error analysis, partly to a refined implementation, in particular of the multiplication by  $\log 2$ . It is proportionally more important for lower precisions.

More guard bits will mean a larger percentage of correctly rounded results. As  $g$  is a parameter in our implementation, it is possible to use any value larger than 3.

#### B. The case study of single precision

Setting  $w_F = 23$  and  $g = 3$  in the previous architecture, it turns out that  $k = 9$  allows for a highly efficient architecture on recent FPGAs.

Firstly, we need altogether  $2^9 \times 27$  bits of RAM for  $e^A$  and  $2^9 \times 9$  bits for  $e^Z - Z - 1$ . We can group both tables in a single  $2^9 \times 36$  table with dual-port access. This perfectly matches one Xilinx BlockRAM, or two Altera M9K.

Secondly, the multiplication is now 18x18 bits, unsigned. This perfectly matches the DSP blocks of Altera chips. On Xilinx chips up to Virtex-4, the multipliers are able of 17x17

unsigned, so the cost is one DSP block plus two 18-bit additions. On Virtex-5 the DSP block is able of 17x24 unsigned, so we only need one addition. One more trick allows us to hide the latency of this addition. We choose to input  $e^A$  on 17 bits only instead of 18. To keep the same error bound of one ulp, we now need to round it to 17bits. This rounding requires an addition (so there is no saving compared to extending the multiplier input to 18 bit), but this addition is now before the multiplier, in parallel to the addition of  $Z$  to  $e^Z - Z - 1$ .

### C. Polynomial approximation for large precisions

For larger values of  $w_F$ , a generic polynomial evaluator [17] is used as a black box. It inputs a function of  $[0, 1] \rightarrow [0, 1]$  (here  $e^{2^{-k}x} - 2^{-k}x - 1$ ) with its input and output precisions (given on Figure 1) and a degree, and implements a piecewise polynomial approximation. The input interval is decomposed into smaller intervals, and the number of such intervals is computed so that the generated architecture returns a faithfully rounded result. The architectures are optimized for the target FPGA (currently Xilinx Virtex-4, Virtex-5 and Virtex-6, and Altera Stratix II to IV), making efficient use of the DSP blocks to attain high frequencies.

One advantage of this approach is that it is DSP- and memory- based. Another one is its genericity, as future improvements to the polynomial evaluator will immediately benefit to the exponential. This includes the adaptation of the polynomial evaluator to newer FPGAs.

More specifically, the function evaluated here is easy to approximate by a low-degree polynomial approximations. It turns out that degree 2 is enough for precision up to double-extended precision.

We now have two parameters to set:  $k$ , that fixes the input to the  $e^A$  table, and the degree  $d$  of the polynomial, that fixes the trade-off between area of the coefficient table and DSP count/latency. We have varied these parameters to obtain the best trade-offs, that is a an architecture well balanced between DSP and memory consumption, with memories as full as possible and multipliers used as fully as possible. For instance, for double precision, on all targets the best choice is  $k = 9$  and a degree-2 approximation on 512 intervals.

Figure 2 details one instance of this architecture for Virtex-5.

### D. Pipeline tuning

We have designed a component generator framework that allows us to finely tune the pipelines, and this exponential operator was also a case study for this framework. This is not the subject of this article, but it explains in particular the relatively short latency we are able to obtain. For illustration, Figure 3 shows an example of the obtained component hierarchy, with the pipeline information. It also details the sizes of the various multipliers on this example.

## IV. RESULTS

### A. Synthesis results

Table I provides synthesis results for several precisions and several FPGA targets, and compares with results from

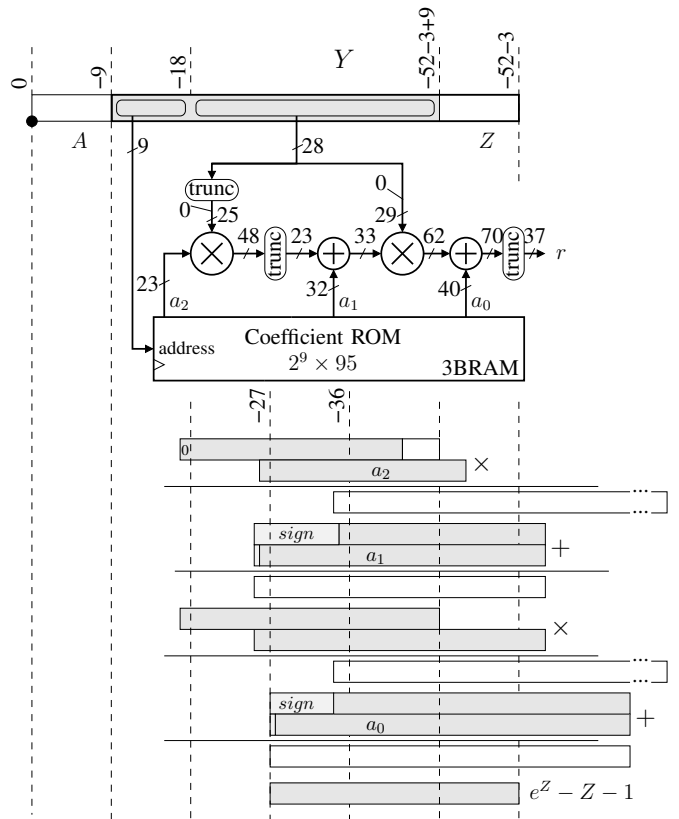


Fig. 2. The architecture evaluating  $e^Z - Z - 1$  for Virtex-5/Virtex-6

previous papers. Our approach is clearly the most efficient of the literature for all the precisions. It combines very high frequency (close to the nominal DSP block frequency), the lowest DSP and memory consumption, portability to both Xilinx and Altera targets, last-bit accuracy, flexibility in precision, and also flexibility in terms of latency versus frequency.

Note that the synthesis on Stratix III reports 2 DSP blocks for single precision. One is actually unused. The coarse-grain DSP block structure of Altera chips since Stratix III prevent using the  $18 \times 18$ -bit multipliers completely independently.

Of special interest is the last line of this table, which shows that even a quadruple-precision exponential function will consume only one tenth of the resources of a high-end FPGA while still running at a very high frequency.

This work is actually open-source and already available on the Internet, and the curious reviewer should have no difficulty to find it to reproduce these results.

### B. Comparison with other works

In [7], a double-precision combinatorial operator consumes, on VirtexII, 2045 slices for a delay of 229 ns. To our knowledge, it was never pipelined, but we estimate that a high-frequency pipelined would require a doubling of the area and roughly 40 cycles.

In addition, this architecture was based on tables inputting  $\alpha$  bits and rectangular multipliers where one dimension was also  $\alpha$  (an integer parameter) and the other dimension varied

TABLE I  
SYNTHESIS RESULTS OF THE VARIOUS INSTANCES OF THE FLOATING-POINT EXPONENTIAL OPERATOR. RESULTS OF ALTERA MEGAWIZARD FOR STRATIXIII WHERE OBTAINED FROM [18]. WE USED QUARTUSII v9.0 FOR STRATIXIII EP5K10K100-3-FF1152, VIRTEX5 XC5VFX100T-3-FF1738 AND VIRTEX6 XC6VHX380T-3-FF1923

Precision	FPGA	Tool	Performance		Resource Usage				
			$f$ (MHz)	Latency	Logic Usage			DSPs	Memory
					(A)LUTs	Reg.	Slice		
(8,23)	StratixIII	Altera MegaWizard	274	17	527	900	-	19 18-bit elem.	0
		<i>ours</i>	391	6	832	374	-	2 18-bit elem.	0
			405	7	519	382	-		2 M9K
	VirtexII 1000	[7]	1/123ns	0			728	0	0
	VirtexIV	<i>ours</i>	375	17	739	510	392	1 DSP48	1 BRAM
			178	8	680	244	367		
	Virtex5	<i>ours</i>	441	20	560	572	-	1 DSP48E	1 BRAM
		188	9	524	256	-			
Virtex6	<i>ours</i>	561	17	602	485	-	1 DSP48E1	1 BRAM	
			241	8	493	229	-		
(10,40)	Virtex5	<i>ours</i> (k=5,d=2)	310	30	1377	1141	-	10 DSP48E	4 BRAM
	Virtex6	<i>ours</i> (k=5,d=2)	488	32	1469	1344	-	10 DSP48E1	3 BRAM
(11,52)	StratixIII	Altera MegaWizard	205	25	2905	2285	-	58 18-bit elem.	0
		<i>ours</i>	327	29	1307	3757	-	22 18-bit elem.	10 M9K
			256	15	1437	1984	-		
	VirtexII 1000	[7]	1/229ns	0			2045	0	0
		[8]	?	0	1293	105		71 DSP48	6 BRAM
	VirtexIV	[9]	200	30	13614	19704		0	29 BRAM
		[10] (CORDIC)	5.25 cycles@100Mhz	>61			23455	36 DSP48	
		<i>ours</i>	319	38	2249	1964	1393		
			178	24	2128	1361	1154	17 DSP48	5 BRAM
			97	14	2034	926	1096		
	Virtex5	<i>ours</i>	310	35	1867	1456	-	12 DSP48E	5 BRAM
			204	18	1604	1018	-		
			119	12	1601	806	-		
	Virtex6	<i>ours</i>	488	38	1928	1791	-	12 DSP48E1	5 BRAM
		221	22	1642	1184	-			
		125	10	1547	629	-			
(15,64)	Virtex6	<i>ours</i> (k=11, d=2)	486	41	2894	2539	-	20 DSP48E1	11 BRAM
(15,112)	Virtex6	<i>ours</i> (k=14, d=3)	395	69	8071	7725	-	71 DSP48E1	123 BRAM

from  $\alpha$  to the mantissa size. This was a good design choice for LUT-based FPGAs, but it poorly matches the capabilities of the DSP blocks and embedded memories of modern FPGAs. For a short latency, and to use the DSP blocks optimally, one should choose  $\alpha = 17$ , but then the tables would be much too large ( $2^{17}$  entries). Or, one should chose  $\alpha \approx 10$ , but then the DSPs would be underutilized.

As Altera Megawizard produces readable source files, we could analyse the algorithm used. The range reduction is the usual one, and the architecture diverges only for the computation of  $e^Y$ . For double precision, Altera’s architecture is based on a decomposition of the input as  $Y = Y_0 + Y_1 + Y_2 + Y_L$  where  $Y_0$  consists of the 9 leading bits,  $Y_1$  and  $Y_2$  consist of the two following 9-bit chunks, and  $Y_L$  consists of the remaining lower bits. The exponential is computed as  $e^Y = (e^{y_0} \times e^{y_1}) \times (e^{y_2} e^{y_L})$ , where the three first terms are simply read from tables with  $2^9$  entries, and  $e^{y_L}$  is approximated as the Taylor polynomial  $e^{y_L} \approx 1 + Y_L$ . This is very similar to the method proposed by Wielgosz et al [8], [9], and both were probably designed independently. However the latter is not generic in precision.

We weren’t able to synthesize these Altera ALTFP\_EXP operator (our Quartus 9 hangs on it), so we report results from the documentation [18]. They do not use 9Kbit embedded

memories, although this design would be a perfect match for them (it should consume  $(61 + 51 + 42)/18 = 9$  of them, with a corresponding huge reduction in logic resources).

This approach has a potential of lower latency, as the multipliers are organized in tree and not in sequence as in our proposal. Its drawback is that it doesn’t exploit the structure of the numbers. The three multiplications are of size roughly  $60 \times 60$  bits. However,  $e^{y_1}$ ,  $e^{y_2}$ , and  $e^{y_L}$  are all of the form  $1 + \epsilon$ , so at the bit level, we have a lot of predictable multiplications by 0, for which the hardware could be saved. Table I will show the advantage of our approach in terms of multiplier count and performance.

Finally, we remark that the two references by Wielgosz et al. [8], [9] seem to use the same architecture, however the first one reports results using DSP blocks, while the second one replaces all the DSPs with logic. This actually makes sense, since in this case the parts of the large multipliers that multiply by zero will indeed be optimized out by the synthesizer.

### C. Comparison with microprocessors

This table allows us to compare the theoretical peak performance, in terms of floating-point exponentials, of a large FPGA and a high-end processor. These numbers, of course, should be taken with due care as they ignore the issue of data movements which are a limiting factor [9].

```

|---Entity LeftShifter_53_by_max_65:
|   Pipeline depth = 1
|---Entity IntAdder_67_f500_slice_SRL_noBUFFER_uid1_Classical:
|   Pipeline depth = 2
|   |---Entity KCMTable_6_11818_unsigned:
|   |   Not pipelined
|   |---Entity KCMTable_6_11818_signed:
|   |   Not pipelined
|   |   |---Entity IntAdder_22_f500_slice_SRL_noBUFFER_uid2_C
|   |   |   Pipeline depth = 1
|   |---Entity IntCompressorTree_22_3:
|   |   Pipeline depth = 2
|---Entity IntIntKCM_14_11818_signed:
|   Pipeline depth = 2
|   |---Entity KCMTable_6_51145234580810622639_unsigned:
|   |   Not pipelined
|   |---Entity KCMTable_6_51145234580810622639_signed:
|   |   Not pipelined
|   |   |---Entity IntAdder_72_f500_slice_SRL_noBUFFER_uid3_C
|   |   |   Pipeline depth = 2
|   |---Entity IntCompressorTree_72_2:
|   |   Pipeline depth = 2
|---Entity IntIntKCM_12_51145234580810622639_signed:
|   Pipeline depth = 3
|---Entity IntAdder_67_f500_slice_SRL_noBUFFER_uid4_Classical:
|   Pipeline depth = 2
|---Entity firstExpTable_11_56:
|   Not pipelined
|   |---Entity TableGenerator_7_89:
|   |   Not pipelined
|   |   |---Entity SignedIntMultiplier_25_23:
|   |   |   Pipeline depth = 2
|   |   |---Entity IntAdder_31_f500_slice_SRL_noBUFFER_uid5_C
|   |   |   Pipeline depth = 1
|   |   |   |---Entity IntAdder_60_f500_slice_SRL_noBUFFE
|   |   |   |   Pipeline depth = 1
|   |   |---Entity IntCompressorTree_60_2:
|   |   |   Pipeline depth = 1
|   |   |---Entity SignedIntMultiplier_27_31:
|   |   |   Pipeline depth = 4
|   |   |---Entity IntAdder_64_f500_slice_SRL_noBUFFER_uid7_C
|   |   |   Pipeline depth = 2
|   |---Entity PolynomialEvaluator_d2:
|   |   Pipeline depth = 10
|   |---Entity IntAdder_35_f500_slice_SRL_noBUFFER_uid8_Class
|   |   Pipeline depth = 1
|---Entity FunctionEvaluator:
|   Pipeline depth = 14
|   |---Entity IntAdder_75_f500_slice_SRL_noBUFFER_uid9_C
|   |   Pipeline depth = 2
|   |---Entity IntCompressorTree_75_2:
|   |   Pipeline depth = 2
|---Entity IntMultiplier_45_45:
|   Pipeline depth = 6
|---Entity IntAdder_57_f500_slice_SRL_noBUFFER_uid10_Classica
|   Pipeline depth = 1
|---Entity IntAdder_65_f500_slice_SRL_noBUFFER_uid11_Classica
|   Pipeline depth = 2
Entity FPExp_11_52:
  Pipeline depth = 38

```

Fig. 3. Component hierarchy for DP exponential on Virtex6

The largest Virtex-6 FPGA (XC6VVSX475T) could accommodate 168 double-precision exponential cores running above 400 MHz, so the theoretical peak performance of the FPGA is now over 60 giga FP exponentials per second (GFPEXP/s).

For a fair comparison, we have to compare to the highest performance software implementation currently available, one which was tuned with comparable effort. To our knowledge, it is the Intel Vector Math Library (VML), which can achieve a peak of 6 cycles/DPEXP on Itanium-2 or Core i7. On an 8-core processor running at 3GHz, we obtain a peak performance of 4 GFPEXP/s, with a speed-up of 15 in favor of the FPGA. On single precision, the numbers are in excess of 400GSPEXP/s

for the FPGA while the performance of VML is only improved to 6GSPEXP/s. The FPGA speed-up is now above 60.

## V. CONCLUSION AND FUTURE WORK

We have presented a state-of-the-art floating-point exponential operator generator. It produces last-bit accurate architectures for a wide range of FPGA targets, for a wide range of precisions up to IEEE-754-2008 quadruple precision, and for a wide range of latency/frequency trade-offs. It is designed to make good use of the DSP blocks and embedded memories of high-end FPGAs, and outperforms previous works in performance and resources consumption.

Two functions already constitute a library: After the logarithm completed in 2009, we intend to extend FloPoCo further to include a complete open-source and portable mathematical library (libm) for FPGAs. This is an important enabling step for the success of C-to-hardware compilers for reconfigurable computing.

## REFERENCES

- [1] C. Doss and R. L. Riley, Jr., "FPGA-based implementation of a robust IEEE-754 exponential unit," in *Field-Programmable Custom Computing Machines*. IEEE, 2004, pp. 229–238.
- [2] J. Detrey and F. de Dinechin, "A parameterized floating-point exponential function for FPGAs," in *Field-Programmable Technology*. IEEE, 2005.
- [3] P. Echeverra, D. Thomas, M. Lpez-Vallejo, and W. Luk, "An FPGA run-time parameterisable log-normal random number generator," in *Reconfigurable computing: architectures, tools and applications*, 2008, pp. 221–232.
- [4] M. Ercegovac, "Radix-16 evaluation of certain elementary functions," *IEEE Transactions on Computers*, vol. C-22, no. 6, pp. 561–566, 1973.
- [5] W. F. Wong and E. Goto, "Fast hardware-based algorithms for elementary function computations using rectangular multipliers," *IEEE Transactions on Computers*, vol. 43, no. 3, pp. 278–294, 1994.
- [6] J. A. Pineiro, M. D. Ercegovac, and J. D. Bruguera, "Algorithm and architecture for logarithm, exponential, and powering computation," *IEEE Transactions on Computers*, vol. 53, no. 9, pp. 1085–1096, Sep. 2004.
- [7] J. Detrey, F. de Dinechin, and X. Pujol, "Return of the hardware floating-point elementary function," in *18th Symposium on Computer Arithmetic*. IEEE, 2007, pp. 161–168.
- [8] M. WIELGOSZ, E. JAMRO, and K. WIATR, "Implementacja w ukadach fpga operacji eksponenty dla liczb w standardzie IEEE-754 o podwójnej precyzji," *PAK*, vol. 53, no. 5, pp. 126–128, 2007.
- [9] M. Wielgosz, E. Jamro, and K. Wiatr, "Accelerating calculations on the rasc platform: A case study of the exponential function," in *ARC '09: Proceedings of the 5th International Workshop on Reconfigurable Computing: Architectures, Tools and Applications*. Berlin, Heidelberg: Springer-Verlag, 2009, pp. 306–311.
- [10] R. Pottathuparambil and R. Sass, "A parallel/vectorized double-precision exponential core to accelerate computational science applications," in *FPGA '09: Proceeding of the ACM/SIGDA international symposium on Field programmable gate arrays*. New York, NY, USA: ACM, 2009, pp. 285–285.
- [11] M. Langhammer, "Foundation for FPGA acceleration," in *Fourth Annual Reconfigurable Systems Summer Institute*, 2008.
- [12] A. Vázquez and E. Antelo, "Implementation of the exponential function in a floating-point unit," *Journal of VLSI Signal Processing*, vol. 33, no. 1-2, pp. 125–145, Jan. 2003.
- [13] F. de Dinechin, C. Klein, and B. Pasca, "Generating high-performance custom floating-point pipelines," in *Field Programmable Logic and Applications*. IEEE, Aug. 2009, pp. 59–64.
- [14] R.-C. Li, P. Markstein, J. P. Okada, and J. W. Thomas, "The libm library and floating-point arithmetic for HP-UX on Itanium," Hewlett-Packard company, Tech. Rep., 2001.
- [15] K. Chapman, "Fast integer multipliers fit in FPGAs (EDN 1993 design idea winner)," *EDN magazine*, May 1994.



- [16] M. Schulte and E. Swartzlander, "Truncated multiplication with correction constant," in *Workshop on VLSI Signal Processing*, 1993, pp. 388–396.
- [17] F. de Dinechin, M. Joldes, and B. Pasca, "Automatic generation of polynomial-based hardware architectures for function evaluation," in *Application-specific Systems, Architectures and Processors*. IEEE, 2010.
- [18] *Floating Point Exponent (ALTFP\_EXP)*, Altera Corporation, 2008.