

# BKM: A New Hardware Algorithm for Complex Elementary Functions

Jean-Claude BAJARD, Sylvanus KLA and Jean-Michel MULLER

CNRS, Laboratoire LIP-IMAG, Ecole Normale Supérieure de Lyon  
46 Allée d'Italie, 69364 Lyon Cedex 07 FRANCE

## Abstract

A new algorithm for computing complex logarithms and exponentials is proposed. This algorithm is based on shift-and-add elementary steps, and it generalizes the CORDIC algorithm. It can compute the usual real elementary functions. This algorithm is more suitable for computations in a redundant number system than CORDIC, since there is no scaling factor for computation of trigonometric functions.

## Index terms

Elementary functions, CORDIC.

## 1. Introduction

The point at stake here is the search for algorithms that rapidly compute elementary functions. Many methods have been used, e.g. approximation by polynomials, Newton's method, E-Method [4], and shift-and-add methods. The shift-and-add methods use simple elementary steps: additions, and multiplications by a power of the radix of the number system. They go back to the 17<sup>th</sup> century: Briggs used such

an algorithm for building the first tables of logarithms [9]. For instance, in radix 2, to compute  $\ln x$  with approximately  $n$  significant bits, numerous methods [2], [7] consist of finding a sequence  $d_k = \pm 1, 0$ , such that  $x \prod_{k=1}^n (1 + d_k 2^{-k}) \approx 1$ . Then  $\ln(x) \approx -\sum_{k=1}^n \ln(1 + d_k 2^{-k})$ . Another important shift-and-add method is the CORDIC algorithm, introduced by Volder [11] for computing trigonometric functions, and generalized by Walther [12]. CORDIC has been implemented in many machines (e.g. Hewlett Packard's HP 35, Intel 8087). It consists of the following iteration:

$$(1) \begin{cases} x_{n+1} = x_n - m d_n y_n 2^{-\sigma(n)} \\ y_{n+1} = y_n + d_n x_n 2^{-\sigma(n)} \\ z_{n+1} = z_n - d_n e_{\sigma(n)} \end{cases}$$

$m$  equals 0, 1 or -1, and  $d_n$  is equal to 1 or -1, so this iteration is reduced to a few additions and shifts. The results and the appropriate values of  $d_n$ ,  $m$  and  $\sigma(n)$  are given in Table 1. For a recent survey of CORDIC, see [5].

	rotation mode ( $d_n = \text{sign } z_n$ )	vectoring mode ( $d_n = -\text{sign } y_n$ )	$\sigma(n)$	$e_n$	scale factor $K$
$m = 1$ (circular)	$x_n \rightarrow K (x_0 \cos z_0 - y_0 \sin z_0)$ $y_n \rightarrow K (y_0 \cos z_0 + x_0 \sin z_0)$	$x_n \rightarrow K \sqrt{x_0^2 + y_0^2}$ $z_n \rightarrow z_0 + \tan^{-1} y_0/x_0$	$n$	$\tan^{-1} 2^{-n}$	$\prod_{n=0}^{\infty} \sqrt{1 + 2^{-2n}}$ $\approx 1.64676$
$m = 0$ (linear)	$x_n = x_0$ $y_n \rightarrow y_0 + x_0 z_0$	$x_n = x_0$ $z_n \rightarrow z_0 + y_0/x_0$	$n$	$2^{-n}$	no scale factor
$m = -1$ (hyperbolic)	$x_n \rightarrow K (x_0 \cosh z_0 + y_0 \sinh z_0)$ $y_n \rightarrow K (y_0 \cosh z_0 + x_0 \sinh z_0)$	$x_n \rightarrow K \sqrt{x_0^2 - y_0^2}$ $z_n \rightarrow z_0 + \tanh^{-1} y_0/x_0$	$\sigma(n) = n - k$ , ( $k = \text{largest integer such that } 3^{k+1} + 2^{k-1} \leq 2n$ )	$\tanh^{-1} 2^{-n}$	$\prod_{n=1}^{\infty} \sqrt{1 - 2^{-2\sigma(n)}}$ $\approx 0.82816$

Table 1. Different functions computable using CORDIC.

The major drawback of CORDIC arises when performing the iterations using a redundant number system. Such number systems are advantageous for quickly-performed arithmetic, since they make it possible to perform carry-free additions [1]. With these systems,  $d_n$  is difficult to evaluate. For instance, assume that we are in the *rotation* and *circular*

modes of CORDIC (see Table 1), and that numbers are represented in radix 2 with digits in  $\{-1, 0, 1\}$ .  $d_n$  equals the sign of the most significant non zero digit of  $z_n$ : to find its value, we may have to examine all the digits of  $z_n$ , and the advantage of the redundant representation (constant time elementary step) is lost. An alternative is to accept  $d_n = 0$ ,

but with such a method the scale factor  $K$  is no longer constant.  $K = \prod_{n=0}^{\infty} \sqrt{1 + d_n^2 2^{-2n}}$  is a constant if the  $d_i$ 's are all equal to  $\pm 1$ , but not if they can be 0. Many solutions have been suggested to solve this problem. They lead to a repetition of iterations in time [10], or in space [3]. To avoid this, we need to work out a new algorithm. Throughout this paper, we assume that we use a radix-2 usual or signed-digit number system. The main advantage of our algorithm (constant-time elementary step without scale factor) appears if the signed-digit system is used. Extension to binary carry-save representation is simple.

Consider the basic step of CORDIC in circular mode (i.e. (1) with  $m=1$  and  $\sigma(n)=n$ ), and define the complex number  $L_n = x_n + iy_n$ . We get:  $L_{n+1} = L_n (1 + id_n 2^{-n})$ . This brings us to a generalization of this algorithm: we could perform multiplications by terms of the form  $(1 + d_n 2^{-n})$ , where the  $d_n$ 's are complex numbers, chosen such that a multiplication by  $d_n$  can be reduced to a few additions. In this paper, we study the following iteration, called *BKM*:

$$(2) \quad \begin{cases} L_{n+1} = L_n (1 + d_n 2^{-n}) \\ E_{n+1} = E_n - \ln(1 + d_n 2^{-n}) \end{cases}$$

with  $d_n = -1, 0, 1, -i, i, 1-i, 1+i, -1-i, -1+i$

In  $z$  is the number  $t = a+ib$  such that  $e^t = e^a(\cos b + i \sin b) = z$ , with  $b$  lying in  $[-\pi, \pi]$ .

If we find a sequence  $d_n$  such that  $L_n$  goes to 1, then we obtain  $E_n \rightarrow E_1 + \ln(L_1)$ : we call this iteration the *L-mode* of the BKM algorithm. If we find a sequence  $d_n$  such that  $E_n$  goes to 0, then we obtain  $L_n \rightarrow L_1 e^{E_1}$ : we call this iteration the *E-mode* of BKM. Therefore, in the next sections, we focus on the problem of finding sequences  $d_n$  such that  $L_n$  goes to 1 or such that  $E_n$  goes to 0.

## 2. Computation of the complex exponential function (E-mode)

For computing  $e^{E_1}$  using BKM, one needs to find a sequence  $d_n \in D = \{-1, 0, 1, -i, i, 1-i, 1+i, -1-i, -1+i\}$  such that the sequence  $E_n$  of Eq. (2) goes to 0 as  $n$  goes to  $+\infty$ . At the outset, let us examine the numbers whose exponential can be computed. The set  $A = \left\{ \sum_{n=1}^{\infty} \ln(1 + d_n 2^{-n}) \mid d_n \in D \right\}$  of the numbers  $E_1$  such that there exists a sequence  $d_n \in D$  satisfying  $E_n \rightarrow 0$  is shown in Fig. 1.

Define  $d_n^x$  and  $d_n^y$  as the real and imaginary parts of  $d_n$  and  $E_n^x$  and  $E_n^y$  as the real and imaginary parts of  $E_n$ . We find:

$$(3) \quad \begin{cases} E_{n+1}^x = E_n^x - \frac{1}{2} \ln \left[ 1 + d_n^x 2^{-n+1} + (d_n^{x2} + d_n^{y2}) 2^{-2n} \right] \\ E_{n+1}^y = E_n^y - d_n^y \tan^{-1} \left( \frac{2^{-n}}{1 + d_n^x 2^{-n}} \right) \end{cases}$$

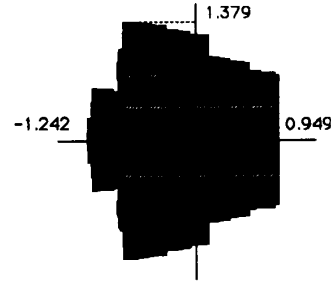


Fig.1. The set A, and the convergence area of the E-mode.

In this section, we give an algorithm which computes the sequence  $d_n$  for any  $E_1$  belonging to a rectangular set  $R_1$ . The algorithm uses a sequence  $R_n = [-s_n^x, r_n^x] + i[-r_n^y, r_n^y]$  of rectangles, whose length goes to 0 as  $n$  goes to  $+\infty$ , and such that for any  $E_n$  belonging to  $R_n$ ,  $d_n$  is such that  $E_{n+1}$  belongs to  $R_{n+1}$ .  $d_n^x$  is chosen by examining a few digits of  $E_n^x$  and  $d_n^y$  is chosen by examining a few digits of  $E_n^y$ : this allows a simple implementation of the choice of  $d_n$ .

### 2.a. Choice of $d_n^x$

Fig 2 shows the parameters involved in determining  $d_n^x$ . This figure is close to the *Robertson Diagrams* that appear in many division algorithms [6], [8]. In the following, we call such a diagram a *Robertson diagram*.

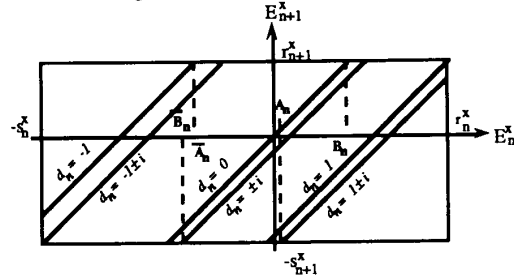


Fig. 2. The Robertson Diagram for  $E_n^x$

The diagram is constructed as follows:

□ we assume that  $E_n^x$  belongs to  $[-s_n^x, r_n^x]$ , which is the real part of  $R_n$ .  $s_n^x$  and  $r_n^x$  will be determined later.

$$\square E_{n+1}^x = E_n^x - \frac{1}{2} \ln \left[ 1 + d_n^x 2^{-n+1} + (d_n^{x2} + d_n^{y2}) 2^{-2n} \right],$$

so the value of  $E_{n+1}^x$  vs.  $E_n^x$  is given by various straight lines parametrized by  $d_n^x$  and  $d_n^y$ .

□ An adequate value of  $d_n^x$  is such that for any value of  $d_n^y$  ( $-1, 0$  or  $1$ ),  $E_{n+1}^x$  "remains in the diagram" (i.e.  $E_{n+1}^x \in [-s_{n+1}^x, r_{n+1}^x]$ ). Therefore  $r_{n+1}^x$  must be the largest value of



### BKM Algorithm - E-mode

□ Start with  $E_1 \in R_1 = [-0.829802\dots, +0.868876\dots] + i[-0.749780\dots, +0.749780\dots]$

□ Iterate: 
$$\begin{cases} L_{n+1} = L_n(1 + d_n 2^{-n}) \\ E_{n+1} = E_n - \ln(1 + d_n 2^{-n}) \end{cases}$$

with  $d_n = d_n^x + id_n^y$ , ( $d_n^x, d_n^y = -1, 0, 1$ ), chosen as follows: define  $\tilde{E}_n^x$  as the number obtained by truncating the real part of  $2^n E_n$  after its 3<sup>rd</sup> fractional digit, and  $\tilde{E}_n^y$  as the number obtained by truncating the imaginary part of  $2^n E_n$  after its 4<sup>th</sup> fractional digit.

$$\begin{cases} \text{if } \tilde{E}_n^x \leq -5/8 \text{ then } d_n^x = -1 \\ \text{if } -1/2 \leq \tilde{E}_n^x \leq 1/4 \text{ then } d_n^x = 0 \\ \text{if } \tilde{E}_n^x \geq 3/8 \text{ then } d_n^x = +1 \end{cases} \quad \begin{cases} \text{if } \tilde{E}_n^y \leq -13/16 \text{ then } d_n^y = -1 \\ \text{if } -3/4 \leq \tilde{E}_n^y \leq 3/4 \text{ then } d_n^y = 0 \\ \text{if } \tilde{E}_n^y \geq 13/16 \text{ then } d_n^y = +1 \end{cases}$$

□ Result:  $L_n \rightarrow L_1 e^{E_1}$

In practice, instead of computing  $E_n$  and examining the first digits of  $\alpha_n = 2^n E_n$ , one could directly compute the sequence  $\alpha_{n+1} = 2\alpha_n - 2^{n+1} \ln(1 + d_n 2^{-n})$ .

### 2.d. Number of iterations

Let us estimate the number of iterations required to obtain a given accuracy. We want to compute  $L_1 e^{E_1}$ . The sequence  $d_i$  satisfies:  $L_1 e^{E_1} = L_1 \prod_{i=1}^{\infty} (1 + d_i 2^{-i})$ . After  $n$  iterations, we have computed  $L_1 \prod_{i=1}^n (1 + d_i 2^{-i})$ . The relative error made by approximating  $L_1 e^{E_1}$  by this value is:

$$(12) \quad \left| \frac{L_1 e^{E_1} - L_1 \prod_{i=1}^n (1 + d_i 2^{-i})}{L_1 e^{E_1}} \right| = \left| 1 - \frac{1}{\prod_{i=n+1}^{\infty} (1 + d_i 2^{-i})} \right|$$

One can show that this value is bounded by a term equivalent to  $2^{-n}$ . Thus, after  $n$  iterations of the E-mode of BKM, we obtain a relative error approximately equal to  $2^{-n}$ . So the error behaviour of BKM is the same as that of CORDIC.

### 2.e. Number of constants stored

This algorithm requires the pre computation and storage of:

- $\ln\left(1 + d_i^x 2^{-i+1} + (d_i^{x2} + d_i^{y2}) 2^{-2i}\right)$ ,  $d_i^x, d_i^y = -1, 0, 1$
- $\tan^{-1}\left(\frac{2^{-i}}{1 + d_i^x 2^{-i}}\right)$ ,  $d_i^x = -1, 0, 1$

so, we need to store 8 terms for each value of  $i$ . From section 2.d, we deduce that, in order to obtain approximately  $n$  accuracy binary digits, we need to store  $8n$  constants.

### 3. Computation of the complex logarithm function (L-mode)

Computing  $\ln(L_1)$  using BKM requires the calculation of a sequence  $d_n \in D = \{-1, 0, 1, -i, i, i-1, i+1, -i-1, -i+1\}$ , such that:

$$(13) \quad L_{n+1} = L_n (1 + d_n 2^{-n}) \rightarrow 1$$

Fig. 4 shows the set  $B = \left\{ \prod_{n=1}^{\infty} (1 + d_n 2^{-n})^{-1} \mid d_n \in D \right\}$  of the numbers  $L_1$  such that such a sequence  $d_n$  exists.



Fig. 4 The set B, and the domain T where the convergence of the algorithm is proven.

### 3.a A Straightforward strategy

In the following, we use the norm  $\|a+ib\| = \max\{|a|, |b|\}$ . Define a sequence  $\varepsilon_n$  as:  $\varepsilon_n = 2^n (L_n - 1)$ . We obtain:

$$(14) \quad \varepsilon_{n+1} = 2(\varepsilon_n + d_n) + d_n \varepsilon_n 2^{-n+1}$$

If we find a sequence  $d_n$  such that the terms  $\varepsilon_n$  are bounded, then (13) will be satisfied. An intuitive solution is to choose  $d_n$  roughly equal to  $-\varepsilon_n$ . So, in this section, we consider the following strategy which gives  $\|\varepsilon_n\| \leq 3/2$ :

- at step  $i$ , we examine the value  $\tilde{\varepsilon}_i$  obtained by truncating the real and imaginary parts of  $\varepsilon_i$  after their  $p^{\text{th}}$  fractional digits, where  $p$  is a small integer.

- $d_i$  is obtained by rounding the real and imaginary part of  $-\tilde{\varepsilon}_i$  to the nearest integer. Since  $p$  is small, this is easily performed. If  $\|\varepsilon_i\| \leq 3/2$ , this choice will give  $d_i \in D$ .

If this algorithm actually gives  $\|\varepsilon_n\| \leq 3/2$  for any  $n$ , then the sequence  $d_n$  will fulfill (13). From  $\|\tilde{\varepsilon}_i - \varepsilon_i\| \leq 2^{-p}$  and  $\|d_i + \tilde{\varepsilon}_i\| \leq 1/2$ , using (14), we deduce:  $\|\varepsilon_{n+1}\| \leq 1 + 2^{1-p} + 2^{-n+1} \|d_n \varepsilon_n\|$ . The norm  $\|\cdot\|$  satisfies  $\|zz'\| \leq 2\|z\| \|z'\|$ , therefore:

$$(15) \quad \|\varepsilon_{n+1}\| \leq 1 + 2^{1-p} + 2^{-n+2} \|\varepsilon_n\|$$

If  $n \geq 4$ ,  $p \geq 4$ , and if  $\|\varepsilon_4\| \leq 3/2$ , then, using (15), one can prove that for any  $n \geq 3$ ,  $\|\varepsilon_n\| \leq 3/2$ . Thus, if we start the iteration (14) at step 4, from  $\varepsilon_4$  satisfying  $\|\varepsilon_4\| \leq 3/2$ , then the strategy presented above will hold. This strategy allows computation of logarithms in a very tiny domain only: we can use it to compute  $\ln(L_4)$  if  $\|\varepsilon_4\| = \|16(L_4 - 1)\| \leq 3/2$ , i.e. if  $L_4 \in [1-3/32, 1+3/32] + i[-3/32, 3/32]$ .

### 3.b Computation in a larger domain

We still study the sequence  $\varepsilon_k = 2^k (L_k - 1)$ . Our purpose is to start its evaluation with  $\varepsilon_1$  belonging to a domain that will be given later, and to obtain, after  $n$  steps ( $n \geq 3$ ), a value  $\varepsilon_{n+1}$  such that  $\|\varepsilon_{n+1}\| \leq 3/2$ . After this, the strategy of section 3.a can be used. The following algorithm was found through simulations, before being proved.

#### BKM Algorithm - L-mode

□ Start with  $L_1$  belonging to the trapezoid  $T$  delimited by the straight lines  $x = 1/2$ ,  $x = 1.3$ ,  $y = x/2$ ,  $y = -x/2$ .

□ Iterate: 
$$\begin{cases} L_{n+1} = L_n (1 + d_n 2^{-n}) \\ E_{n+1} = E_n - \ln(1 + d_n 2^{-n}) \end{cases}$$

with  $d_n = d_n^x + i d_n^y$ , ( $d_n^x, d_n^y = -1, 0, 1$ ), chosen as follows:

- define  $\varepsilon_n^x$  and  $\varepsilon_n^y$  as the real and imaginary parts of  $\varepsilon_n = 2^n (L_n - 1)$ , and  $\tilde{\varepsilon}_n^x$  and  $\tilde{\varepsilon}_n^y$  as the values obtained by truncating these numbers after their 4<sup>th</sup> fractional digits.

- At step 1:

$$\begin{cases} \text{if } \tilde{\varepsilon}_1^x \leq -7/16 \text{ and } 6/16 \leq \tilde{\varepsilon}_1^y \text{ then } d_1 = 1 - i \\ \text{if } \tilde{\varepsilon}_1^x \leq -7/16 \text{ and } \tilde{\varepsilon}_1^y \leq -6/16 \text{ then } d_1 = 1 + i \\ \text{if } -6/16 \leq \tilde{\varepsilon}_1^x \text{ and } 8/16 \leq \tilde{\varepsilon}_1^y \text{ then } d_1 = -i \\ \text{if } -6/16 \leq \tilde{\varepsilon}_1^x \text{ and } \tilde{\varepsilon}_1^y \leq -9/16 \text{ then } d_1 = i \\ \text{if } \tilde{\varepsilon}_1^x \leq -7/16 \text{ and } -5/16 \leq \tilde{\varepsilon}_1^y \leq 5/16 \text{ then } d_1 = 1 \\ \text{if } -6/16 \leq \tilde{\varepsilon}_1^x \text{ and } -1/2 \leq \tilde{\varepsilon}_1^y \leq 1/2 \text{ then } d_1 = 0 \end{cases}$$

- At step  $n$ ,  $n \geq 2$ :

$$\begin{cases} \text{if } \tilde{\varepsilon}_n^x \leq -1/2 \text{ then } d_n^x = 1 \\ \text{if } -1/2 < \tilde{\varepsilon}_n^x < 1/2 \text{ then } d_n^x = 0 \\ \text{if } 1/2 \leq \tilde{\varepsilon}_n^x \text{ then } d_n^x = -1 \end{cases} \quad \begin{cases} \text{if } \tilde{\varepsilon}_n^y \leq -1/2 \text{ then } d_n^y = 1 \\ \text{if } -1/2 < \tilde{\varepsilon}_n^y < 1/2 \text{ then } d_n^y = 0 \\ \text{if } 1/2 \leq \tilde{\varepsilon}_n^y \text{ then } d_n^y = -1 \end{cases}$$

□ Result:  $E_n \rightarrow E_1 + \ln(L_1)$

In practice, instead of computing  $L_n$ , one could directly compute  $\varepsilon_n = 2^n (L_n - 1)$  using (14).

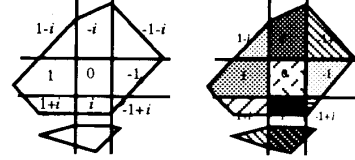
**Proof of the algorithm:** our goal is to show that if  $L_1 \in T$ , then there exists  $n \geq 4$  such that  $\|\varepsilon_n\| \leq 3/2$ . In order to do this, we build a sequence  $\beta_k$  of *bounding sets*, such that for any  $L_1 \in T$ ,  $\varepsilon_k \in \beta_k$ . Our problem is reduced to show that there exists  $n \geq 4$  such that  $\beta_n$  is included in the square  $\|z\| \leq 3/2$ . At the outset, let us explain how the sequence  $\beta_k$  is computed. Figures 5, and 6 show how  $\beta_{k+1}$  is deduced from  $\beta_k$ . The example described in these figures is imaginary: the "true" bounding sets are shown in Fig. 7.

□  $\beta_1$  is equal to  $2(T-1)$ , and  $\beta_k$  is defined as an aggregate of convex polygons, represented by their vertices.

A step of the algorithm can be represented by a splitting of the complex plane into 9 convex *d*-areas. The *d*-area associated with  $\delta \in D$  is the domain  $DA(\delta)$  such that if  $\tilde{\varepsilon}_k$  be-

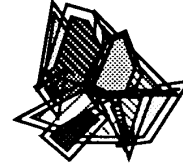
longs to  $DA(\delta)$ , then the algorithm gives  $d_k = \delta$ . For instance, if  $k \geq 2$ , then  $DA(-1-i)$  is the set of the complex whose real and imaginary parts are greater than  $1/2$ . In  $DA(\delta)$ , the transformation  $\varepsilon_{k+1} = 2(\varepsilon_k + \delta) + \delta \varepsilon_k 2^{-k+1}$  is a *similarity*, i.e. the combination of a *rotation* and a *multiplication* by a real factor.

□ Each convex polygon of  $\beta_k$  is splitted into sub-convex polygons, obtained by intersecting it with the *d*-areas. Fig. 5 shows the bounding set at step  $k$ , and the various *d*-areas ( $k \geq 2$ ), and the splitting of the polygons of  $\beta_k$ .



**Fig. 5** Left: The bounding set at step  $k$  and the different *d*-areas (for  $k \geq 2$ )  
right: The bounding set is split into convex polygons following the *d*-areas

Broadly speaking,  $\beta_{k+1}$  is obtained by computing the transformation of each sub-convex polygon generated by the splitting (the image of a polygon is obtained by computing the image of its vertices). We must take into account that  $d_k$  is deduced from  $\tilde{\varepsilon}_k^x$  and  $\tilde{\varepsilon}_k^y$ , which are obtained by truncating the real and imaginary parts of  $\varepsilon_k$  after their 4<sup>th</sup> fractional digits. For instance, if  $\tilde{\varepsilon}_k = \tilde{\varepsilon}_k^x + i \tilde{\varepsilon}_k^y$  belongs to  $DA(-1)$ , this does not prove that  $\varepsilon_k$  belongs to  $DA(-1)$ . Thus, to each sub-convex polygon, a "ribbon" of length  $2^{-4}$  is added, so that if  $\tilde{\varepsilon}_k$  belongs to the "old" sub-polygon, then  $\varepsilon_k$  belongs to the "new" one. Then, for each new sub-polygon, we compute the image of its vertices by the similarity defined by the value of  $d_k$  assigned to the polygon (Fig. 6). This gives the new bounding set  $\beta_{k+1}$ .



**Fig. 6** The iteration is applied to each of the vertices of the sub-polygons, to obtain the new bounding set

The proof that  $\varepsilon_k \in \beta_k$  for any  $L_1 \in T$  is obvious. Thus, if we find  $n \geq 4$  such that all the vertices of the sub-convex polygons of  $\beta_n$  are in the square  $\|z\| \leq 3/2$ , then the algorithm is proven. The adequate value of  $n$  is 6: this leads to a number of vertices much too large to be verified manually. We used a program written in ML for computing all the vertices of  $\beta_6$ , using exact rational arithmetic. Fig. 7 shows the bounding sets  $\beta_1$ ,  $\beta_4$  and  $\beta_6$ . Using this pro-

gram, we have verified that all the vertices of  $\beta_6$  are included in the square  $\|z\| \leq 3/2$ . Fig. 4 shows the domain  $T$  where the convergence of the algorithm is proven.

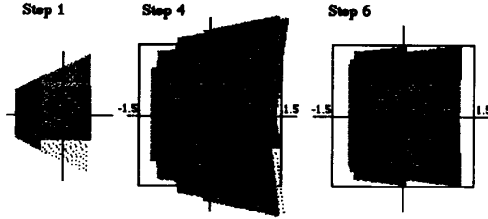


Fig. 7 The bounding sets  $\beta_1, \beta_4$  and  $\beta_6$ .

### 3.c Number of iterations

As we did in section 2.d for the E-mode, let us estimate the number of iterations required to obtain a given accuracy. The sequence  $d_k$  defined by the algorithm satisfies:

$$\ln(L_1) = -\sum_{k=1}^{\infty} \ln(1 + d_k 2^{-k})$$

after  $n$  iterations of the L-mode of BKM, we have computed  $E_1 - \sum_{k=1}^n \ln(1 + d_k 2^{-k})$ . The absolute error made by approximating  $E_1 + \ln(L_1)$  by this value is:

$$(16) \quad \text{error}(n) = \left| \sum_{k=n+1}^{\infty} \ln(1 + d_k 2^{-k}) \right|$$

Using the Taylor expansion of the logarithm, one can show that this expression is bounded by a term equivalent to  $2^{-n} \sqrt{2}$ . Therefore, in order to obtain an absolute error less than  $2^{-n}$ , one needs to perform  $n+1$  iterations.

## 4. Application: computation of elementary functions

As shown in the previous sections, the BKM algorithm makes it possible to compute the following functions:

- in E-mode,  $L_1 e^{E_1}$ , where  $E_1$  belongs to the domain  $[-0.829802, +0.868876] + i \cdot [-0.749780, +0.749780]$ .
- in L-mode,  $E_1 + \ln(L_1)$ , where  $L_1$  belongs to the trapezoid  $T$  delimited by the straight lines  $x = 1/2$ ,  $x = 1.3$ ,  $y = \pm x/2$  (the actual convergence domain looks larger, but the algorithm is proven only for  $L_1 \in T$ ).

Therefore, using BKM, one can compute the following functions of real variables:

### 4.a Functions computable using one mode of BKM.

□ **real sine and cosine functions.** In the E-mode of BKM, one can compute the exponential of  $E_1 = i\theta$  (where  $\theta$  is a real number), and obtain  $L_n = \cos \theta + i \sin \theta \pm 2^{-n}$ .

□ **real exponential function.** If  $E_1$  is a real number belonging to  $[-0.829802, +0.868876]$ , the E-mode of BKM will give a value  $L_n$  equal to  $L_1 e^{E_1 \pm 2^{-n}}$ .

□ **real logarithm.** If  $L_1$  is a real number belonging to  $T$ , the E-mode of BKM will give  $E_n = E_1 + \ln(L_1) \pm 2^{-n}$ . Furthermore, in this case, the iteration is reduced to Brigg's algorithm, and the algorithm works for  $L_1 \in \left[ \prod_{n=1}^{\infty} (1 + 2^{-n})^{-1}, \prod_{n=2}^{\infty} (1 - 2^{-n})^{-1} \right] = [0.419, 1.731]$ .

□ **2-D rotations.** As pointed out in many papers dealing with CORDIC (e.g. [5]), performing rotations is useful for Fast Fourier Transformation, Digital Filtering, and Matrix Computations. The vector  $(c \ d)^t$  obtained by rotating the 2-D vector  $(a \ b)^t$  of an angle  $\theta$  is computed using the E-mode of BKM, with  $L_1 = a + ib$  and  $E_1 = i\theta$ .

□ **real  $\tan^{-1}$  function.** From the relation:

$$\ln(x + iy) = \begin{cases} \frac{1}{2} \ln(x^2 + y^2) + i \tan^{-1} \frac{y}{x} \mod(2\pi) & \text{if } x > 0 \\ \frac{1}{2} \ln(x^2 + y^2) + i \left( \pi + \tan^{-1} \frac{y}{x} \right) \mod(2\pi) & \text{if } x < 0 \end{cases}$$

one can deduce that, if  $x + iy$  belongs to the convergence domain of the L-mode of BKM, then  $\tan^{-1} y/x$  is the imaginary part of the limit value of  $E_n$ , while  $0.5 \ln(x^2 + y^2)$  is its real part, assuming that the L-mode is used with  $E_1 = 0$  and  $L_1 = x + iy$ .

### 4.b Functions computable using two consecutive modes of BKM.

□ **Complex multiplication:** The product  $zt$  is evaluated as  $z \cdot e^{\log t}$ , see fig. 8.

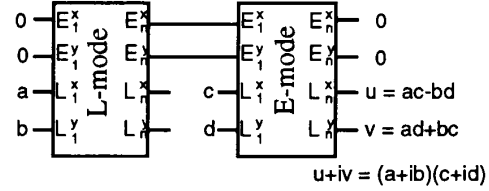


Fig. 8 Complex multiplication

In fact, one can compute  $zte^u$ , where  $z$ ,  $t$  and  $u$  are complex numbers, using the same operator, by choosing  $E_1^x$  equal to the real part of  $u$ , and  $E_1^y$  equal to its imaginary part.

□ **Computation of  $x\sqrt{a}$  and  $y\sqrt{a}$  in parallel** ( $x, y$  and  $a$  are real numbers): we use the relation  $\sqrt{a} = e^{1/2 \ln a}$ .

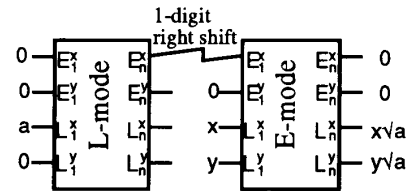


Fig. 9 Computation of  $x\sqrt{a}$  and  $y\sqrt{a}$  in parallel

□ **Computation of lengths and normalization of 2D-vectors:** The L-mode of BKM allows the computation of  $F = 1/2 \ln(a^2 + b^2) = \ln \sqrt{a^2 + b^2}$ , where  $a$  and  $b$  are real numbers. Using the E-mode of BKM, we can compute  $e^F$ , or  $e^{-F}$ . See Fig. 10. The normalization of 2D vectors (i.e. the computation of  $x / \sqrt{a^2 + b^2}$  and  $y / \sqrt{a^2 + b^2}$ , where  $x, y, a$  and  $b$  are real numbers) is a basic step of Givens' QR factorization algorithm.

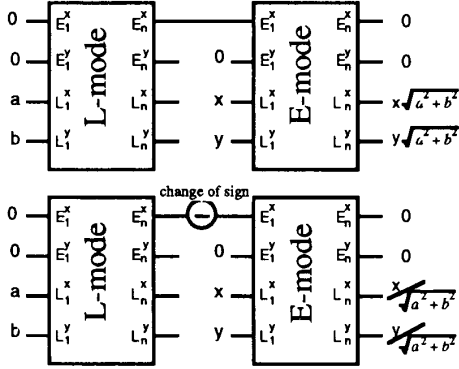


Fig. 10 Lengths and normalization of 2D-vectors

## 5. Comparison with CORDIC

To obtain  $p$  significant bits, CORDIC and BKM roughly need  $p$  iterations. BKM requires the storage of  $8p$  constants, while CORDIC requires the storage of  $p$  constants. Since these constants are represented by  $p$ -digits, both algorithms need a  $O(p^2)$  area for storage of them. Both algorithms need a shifter able to perform an  $n$ -position shift at step  $n$ . A barrel shifter makes it possible to perform a  $n$ -position shift (for any  $n \leq p$ ) in constant time, and lies in an area  $O(n^2)$ . Since the area complexity of most adders is better than  $O(n^2)$ , we deduce that the area complexity of CORDIC and BKM is  $O(n^2)$ . The computations performed during a BKM iteration are:

● **For the variable  $E_n$ :**

$$\begin{cases} E_{n+1}^x = E_n^x - \frac{1}{2} \ln \left[ 1 + d_n^x 2^{-n+1} + \left( d_n^{x^2} + d_n^{y^2} \right) 2^{-2n} \right] \\ E_{n+1}^y = E_n^y - d_n^y \tan^{-1} \left( \frac{2^{-n}}{1 + d_n^x 2^{-n}} \right) \end{cases}$$

or

$$\begin{cases} \alpha_{n+1}^x = 2 \alpha_n^x - 2^n \ln \left[ 1 + d_n^x 2^{-n+1} + \left( d_n^{x^2} + d_n^{y^2} \right) 2^{-2n} \right] \\ \alpha_{n+1}^y = 2 \alpha_n^y - 2^{n+1} d_n^y \tan^{-1} \left( \frac{2^{-n}}{1 + d_n^x 2^{-n}} \right) \end{cases}$$

if instead of computing  $E_n$  and examining the first digits of  $\alpha_n = 2^n E_n$ , we directly compute  $\alpha_n$ .

● **For the variable  $L_n$ :**

$$\begin{cases} L_{n+1}^x = L_n^x + \left( d_n^x L_n^x - d_n^y L_n^y \right) 2^{-n} \\ L_{n+1}^y = L_n^y + \left( d_n^y L_n^x + d_n^x L_n^y \right) 2^{-n} \end{cases}$$

or

$$\begin{cases} \epsilon_{n+1}^x = 2 \left( \epsilon_n^x + d_n^x \right) + \left( d_n^x \epsilon_n^x - d_n^y \epsilon_n^y \right) 2^{-n+1} \\ \epsilon_{n+1}^y = 2 \left( \epsilon_n^y + d_n^y \right) + \left( d_n^x \epsilon_n^y + d_n^y \epsilon_n^x \right) 2^{-n+1} \end{cases}$$

if instead of computing  $L_n$  and examining the first digits of  $\epsilon_n = 2^n (L_n - 1)$ , we directly compute  $\epsilon_n$ .

So the BKM iterations look more complicated than the CORDIC iterations. As a matter of fact, in order to compare CORDIC and BKM, we have to assume that we use a redundant number system. Using such a system, the time complexities of both algorithms are  $O(p)$ . As pointed out in many papers dealing with CORDIC, efficient use of CORDIC with such a number system requires a doubling of the iterations in space [3] or in time [10]. For instance, doubling the CORDIC iterations in time gives:

$$\begin{cases} x_{n+1} = x_n - d_n y_n 2^{-n} - d_n^2 x_n 2^{-2n-2} \\ y_{n+1} = y_n + d_n x_n 2^{-n} - d_n^2 y_n 2^{-2n-2} \\ z_{n+1} = z_n - 2 d_n \tan^{-1} 2^{-n-1} \end{cases}$$

which is at least as complex as the BKM iteration (because at step  $n$  one needs to perform an  $n$  position shift and a  $2n-2$  position shift: this requires a larger shifter). Doubling the iterations in space requires more control: in the *branching CORDIC* method proposed by Duprat and Muller [3], one needs to compare at each step the values given by two CORDIC modules. Furthermore, doubling the iterations makes it possible to obtain a constant scaling factor, but this factor remains different from 1, therefore, for computing many functions, one needs to perform a multiplication after the CORDIC iterations. So, although both methods have the same time and space complexities, BKM looks more interesting when using a redundant number system.

## 6. Conclusion

We have proposed a new algorithm for the computation of many elementary functions (complex exponential and logarithms, complex multiplication, real functions  $\sin$ ,  $\cos$ ,  $\tan^{-1} y/x$ ,  $\ln(x^2 + y^2)$ ,  $x \sqrt{a}$ ,  $x \sqrt{a^2 + b^2}$ ,  $x / \sqrt{a^2 + b^2}$ , 2D rotations). This algorithm matches the CORDIC algorithm, since it allows the use of a redundant number system without any scaling factor problem. Moreover, several functions (complex exponentials and logarithms, complex multiplications), are directly computable using one or two BKM operations, while this is not true using CORDIC.

## References

- [1] A. Avizienis, *Signed-digit number representations for fast parallel arithmetic*, IRE Transactions on electronic computers, 10, pp. 389-400, 1961.
- [2] T.C. Chen, *Automatic Computation of Exponentials, Logarithms, Ratios and Square Roots*, IBM J. Res. Develop, 16, pp 380-388, 1972.
- [3] J. Duprat and J.M. Muller, *The CORDIC Algorithm: new results for fast VLSI implementation*, to appear in IEEE Trans. Computers.
- [4] M.D. Ercegovac, *A general hardware-oriented method for evaluation of functions and computations in a digital computer*, IEEE Transactions on Computers, Vol. C-26 No 7, July 1977, pp 667-680.
- [5] Y.H. Hu, *CORDIC-based VLSI Architectures for Digital Signal Processing*, IEEE Signal Processing Magazine, pp 16-35, July 1992.
- [6] K. Hwang, *Computer Arithmetic principles, architecture and design*, J. Wiley&Sons Inc., New-York, 1979.
- [7] B. De Lugish, *A class of algorithms for automatic evaluation of functions and computations in a digital computer*, PhD dissertation, Department of Computer Science, Univ. of Illinois, Urbana, June 1970.
- [8] J.E. Robertson, *A new class of digital division methods*, IRE Transactions on Electronic Computers, Vol. EC-17, Sept. 1958.
- [9] H.E. Salzer, *radix tables for finding the logarithm of any number of 25 decimal places*, in *Tables of functions and of zeros of functions*, National Bureau of Standards Applied Mathematics Series No 37, 1954, pp 143-144.
- [10] N. Takagi, T. Asada and S. Yajima, *Redundant CORDIC methods with a constant scale factor*, IEEE Trans. on Computers, Vol. 40 No 9, pp 989-995, Sept. 1991.
- [11] J. Volder, *The CORDIC Computing technique*, IRE Transactions on Computers, Sept. 1959.
- [12] J. Walther, *A unified algorithm for elementary functions*, Joint Computer Conf. proc., Vol. 38, 1971.

## Appendix

### Computation of the parameters occurring in the E-mode.

We want to prove that  $\bar{B}_n < \bar{A}_n$ ,  $A_n < B_n$ ,  $C_n < D_n$ , and to find two numbers  $\bar{A}$  and  $A$  whose binary representations have only  $p_1$  fractional digits, and a  $p_2$  fractional digit number  $C$ , such that for any  $n$ :

$$\begin{aligned} 2^n \bar{B}_n &\leq \bar{A} - 2^{-p_1} \leq \bar{A} \leq 2^n \bar{A}_n \\ 2^n A_n &\leq A \leq A + 2^{-p_1} \leq 2^n B_n \\ 2^n C_n &\leq C \leq C + 2^{-p_2} \leq 2^n D_n \end{aligned}$$

with:

$$\begin{aligned} \bar{A}_n &= \ln(1 - 2^{-n}) + \sum_{k=n+1}^{\infty} \ln(1 + 2^{-k}) \\ A_n &= \frac{1}{2} \ln(1 + 2^{-n+1} + 2^{-2n+1}) + \frac{1}{2} \sum_{k=n+1}^{\infty} \ln(1 - 2^{-k+1} + 2^{-2k+1}) \\ \bar{B}_n &= \frac{1}{2} \ln(1 + 2^{-2n}) + \frac{1}{2} \sum_{k=n+1}^{\infty} \ln(1 - 2^{-k+1} + 2^{-2k+1}) \\ B_n &= \sum_{k=n+1}^{\infty} \ln(1 + 2^{-k}) \\ C_n &= \tan^{-1}\left(\frac{2^{-n}}{1 - 2^{-n}}\right) - \sum_{k=n+1}^{\infty} \tan^{-1}\left(\frac{2^{-k}}{1 + 2^{-k}}\right) \\ D_n &= \sum_{k=n+1}^{\infty} \tan^{-1}\left(\frac{2^{-k}}{1 + 2^{-k}}\right) \end{aligned}$$

For  $n = 1$ , we easily obtain:

$$\begin{aligned} 2\bar{B}_1 &\leq -2^{-1} - 2^{-3} \leq -2^{-1} \leq 2\bar{A}_1 \\ 2A_1 &\leq 2^{-2} \leq 2^{-1} \leq 2B_1 \\ 2C_1 &\leq 2^{-1} + 2^{-2} \leq 2^{-1} + 2^{-2} + 2^{-4} \leq 2^n C_n \end{aligned}$$

Using Taylor expansions, we get the following bounds:

$$\begin{aligned} 2^n \bar{A}_n &\geq -\frac{2}{3} 2^{-n} - \frac{2}{3} 2^{-2n} \geq -\frac{1}{4} \\ 2^n \bar{B}_n &\leq -1 + \frac{1}{2} 2^{-n} + \frac{4}{21} 2^{-2n} \leq -\frac{1}{2} \\ 2^n B_n &\geq 1 - \frac{1}{6} 2^{-n} \geq \frac{1}{2} \\ 2^n A_n &\leq \frac{2}{21} 2^{-n} \leq \frac{1}{4} \\ 2^n D_n &\geq 1 - \frac{1}{3} 2^{-n} \geq \frac{1}{2} + \frac{1}{4} + \frac{1}{16} \\ 2^n C_n &\leq \frac{4}{3} 2^{-n} + \frac{2}{3} 2^{-2n} \leq \frac{1}{2} \end{aligned}$$

These relations and the relations obtained for  $n = 1$  give:

$$\begin{aligned} 2^n \bar{B}_n &\leq -\frac{1}{2} - \frac{1}{8} \leq -\frac{1}{2} \leq 2^n \bar{A}_n \\ 2^n A_n &\leq \frac{1}{4} \leq \frac{1}{4} + \frac{1}{8} \leq 2^n B_n \\ 2^n C_n &\leq \frac{1}{2} + \frac{1}{4} \leq \frac{1}{2} + \frac{1}{4} + \frac{1}{16} \leq 2^n D_n \end{aligned}$$

From this, we deduce that  $\bar{B}_n < \bar{A}_n$ ,  $A_n < B_n$ ,  $C_n < D_n$ , and that the following parameters fulfill the requirements presented above:

$$\begin{aligned} \bar{A} &= \frac{1}{2} & A &= \frac{1}{4} & C &= \frac{1}{2} + \frac{1}{4} \\ p_1 &= 3 & p_2 &= 4 \end{aligned}$$