



## Imprecise Arithmetic for Low Power Image Processing

**Albicocco, Pietro; Cardarilli, Gian Carlo ; Nannarelli, Alberto; Petricca, Massimo; Re, Marco**

*Published in:*

2012 Conference Record of the Forty Sixth Asilomar Conference on Signals, Systems and Computers (ASILOMAR)

*Link to article, DOI:*

[10.1109/ACSSC.2012.6489164](https://doi.org/10.1109/ACSSC.2012.6489164)

*Publication date:*

2012

[Link back to DTU Orbit](#)

*Citation (APA):*

Albicocco, P., Cardarilli, G. C., Nannarelli, A., Petricca, M., & Re, M. (2012). Imprecise Arithmetic for Low Power Image Processing. In 2012 Conference Record of the Forty Sixth Asilomar Conference on Signals, Systems and Computers (ASILOMAR) (pp. 983-987 ). IEEE. (Asilomar Conference on Signals, Systems and Computers. Conference Record). DOI: 10.1109/ACSSC.2012.6489164

## DTU Library

Technical Information Center of Denmark

---

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

# Imprecise Arithmetic for Low Power Image Processing

Pietro Albicocco, Gian Carlo Cardarilli, Alberto Nannarelli\*, Massimo Petricca and Marco Re

Department of Electrical Engineering, University of Rome "Tor Vergata", Rome, Italy

\* Dept. of Informatics & Mathematical Modeling, Technical University, Denmark

**Abstract**—Sometimes reducing the precision of a numerical processor, by introducing errors, can lead to significant performance (delay, area and power dissipation) improvements without compromising the overall quality of the processing. In this work, we show how to perform the two basic operations, addition and multiplication, in an imprecise manner by simplifying the hardware implementation. With the proposed "sloppy" operations, we obtain a reduction in delay, area and power dissipation, and the error introduced is still acceptable for applications such as image processing.

## I. INTRODUCTION

There are several fields of application of computer arithmetic that can tolerate some imprecision. For example, in audio and image processing or in wireless communication, it might be desirable to get better performance (faster, smaller, less power-hungry systems) at expenses of some quality degradation.

Recently, a few papers have addressed this issue of designing imprecise hardware to save power [1], [2], [3], [4].

In this work, we introduce a systematic way of having imprecise arithmetic operations for the two most common operations: addition and multiplication. We liked the term "sloppy" introduced in [5], and we will use this term in the paper to refer to imprecise arithmetic operations.

## II. SLOPPY ADDITION

Ignoring the least significant bits of an addition, by implementing a truncated adder saves area and it is faster at expenses of a truncation error. Instead of completely ignoring the least significant bits, in the "sloppy" approach we do not propagate the carry in those bits.

Assuming that we are operating on positive integers, and defining position  $k$  as the bit of weight  $2^k$  in a  $n$ -bit word, we can ignore the carry up to position  $k$  when implementing the addition.

The bit-level algorithm to implement this sloppy adder is the following:

```

c = 0 // carry
if (i < k) then
  si = ai XOR bi
else
  si = ai XOR bi XOR c
  c = (ai AND bi) OR (ai AND c) OR (bi AND c)
end if

```

For example, the addition  $103 + 70$  ( $n = 8$ ,  $k = 4$ ) is

	sloppy	+	precise
A :	0110 0111	+	0110 0111
B :	0100 0110	+	0100 0110
c :	100- ----	=	0100 110-
-----			
S :	1010 0001		1010 1101

That is, the sloppy adder computes 161 (exact value is 173) introducing an error  $\epsilon = 12$ .

By looking at the bits of weight  $< 2^k$ , we notice that the XOR of two ones produces a zero sum bit ( $1 \oplus 1 = 0$ ). Because the carry is not computed (or propagated), in position  $k$  an error  $2^{k+1}$  is generated. The error can be halved to  $2^k$  by computing the OR of the two bits in place of the XOR. For the example above we have:

	sloppy (OR-ing)
A :	0110 0111
B :	0100 0110
c :	100- ----
-----	
S :	1010 0111

and the error is reduced from  $\epsilon = 12$  to  $\epsilon = 6$  (halved).

By simulating all possible combinations of the operands for the 8-bit addition ( $k = 4$ ), we found that by obtaining the sum by OR-ing the  $k$  least-significant bits the average error is  $\epsilon_{mean} = 3.75$ , while by XOR-ing, it is  $\epsilon_{mean} = 7.5$ .

We show in Fig. 1 the comparison of the hardware implementation of the sloppy adder used in the above example ( $n = 8$ ,  $k = 4$ ) and an error-free 8-bit carry-propagate adder (CPA). The data<sup>1</sup> on delay, area and power dissipation are reported in Table I.

In a rough evaluation, we considered lowering the supply voltage  $V_{DD}$  in the sloppy adder to match the delay of the error-free adder (1.0 ns). In our library, when  $V_{DD}$  is lowered from 1.0 V to 0.7 V the delay doubles. In the expression for the power dissipated by a circuit containing  $N$  gates

$$P_{1.0V} = V_{DD}^2 f \cdot \sum_{i=1}^N a_i C_i \Rightarrow 20 \mu W = (1.0 V)^2 \cdot \mathcal{K}$$

we assume that the switching capacitance  $a_i C_i$  does not change when scaling  $V_{DD}$ . Therefore,  $\mathcal{K} = 20$  is constant:

$$P_{0.7V} = (0.7)^2 \cdot 20 \simeq 10 \mu W$$

<sup>1</sup>The adders are synthesized with radix-4 carry-look-ahead iterative carry network [6].

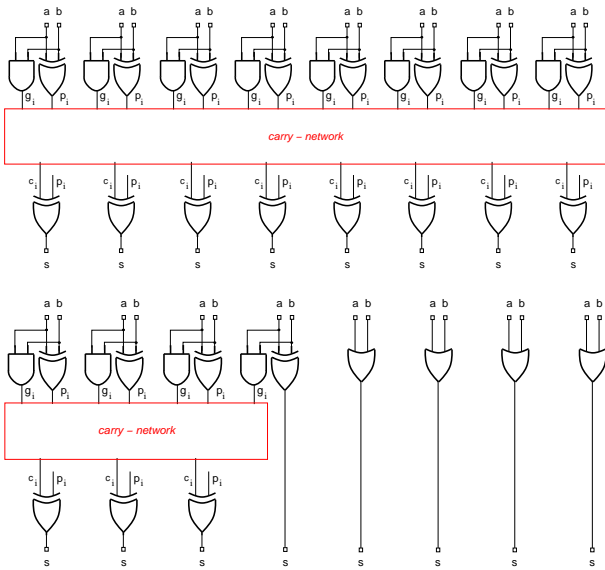


Fig. 1. Implementation of 8-bit error-free (top) and sloppy  $k = 4$  (bottom) adders.

	CPA 8-bit	sloppy	ratio
max. delay [ps]	999	495	2.00
Area [ $\mu\text{m}^2$ ]	191	112	1.70
Power [ $\mu\text{W}$ ]	42	20	2.10

TABLE I  
SYNTHESIS DATA OF ADDERS IN FIG. 1.

That is, with the sloppy adder the power is reduced to 1/4 at same adder speed.

The natural competitor of the sloppy adder is the truncated adder. We performed a comparison between our imprecise adder and the truncated adder by implementing a 16-bit adder with a Carry Look-Ahead (CLA) network to propagate the carry for different sloppy/truncated configurations. The output bits affected by errors are shown as  $\circ$  in Fig. 2 for truncation  $t = 8$  and sloppy bits  $k = 8$ .

Gate-level netlists are generated, by a C program, for each unit under test. The netlists are synthesized (unconstrained) to optimize buffering and cells' drive strength according to the actual fan-out.

In the comparison, we are interested in relating the introduced error to the power dissipation.

Fig. 3 shows the error introduced by the imprecise adders as a function of the number of imprecise/truncated bits (4, 8 and 12 bits). In Fig. 4 we show the power dissipation of each imprecise adder as function of the error. The sloppy adder turned out to dissipate lower power than the truncated adder for the same error level.

### III. SLOPPY MULTIPLICATION

Parallel multiplication  $p = x \cdot y$  can be divided into three steps:

- 1) generation of Partial Products (PPs);
- 2) carry-free reduction from  $n$  PPs to 2 operands;
- 3) carry-propagate two operands addition.

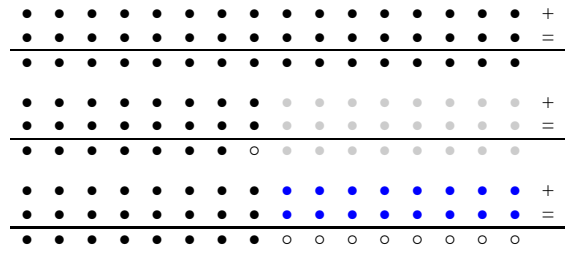


Fig. 2. Bit array and errors  $\circ$  for adders: precise 16-bit adder (top), truncated  $t = 8$  16-bit adder (middle), sloppy  $k = 8$  16-bit adder (bottom).

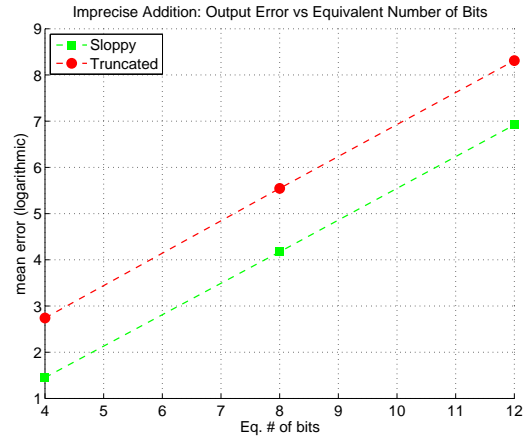


Fig. 3. Error of sloppy and truncated adders implementations for different number of imprecise bits.

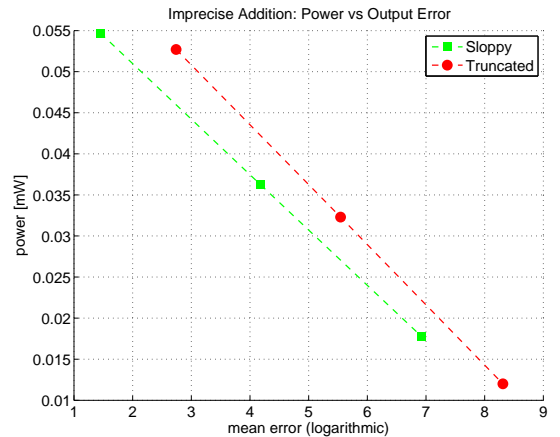


Fig. 4. Power dissipation of sloppy and truncated adders as function of the introduced error.

We use a sloppy approach for step 1 only, as step 2 is quite delay-efficient (no carry propagation) and step 3 has been addressed in the previous section.

We consider radix-4 multiplication because for  $n \times n$  bit operands the unit is smaller: only  $\frac{n}{2}$  PPs are generated. In radix-4 multiplication, the radix-4 digits of the multiplier  $y$  are recoded into signed-digit representation to avoid multiples of 3 and carry propagation as explained in [6]. The resulting architecture (for one digit) recoder plus PP generation (rec+PPgen)

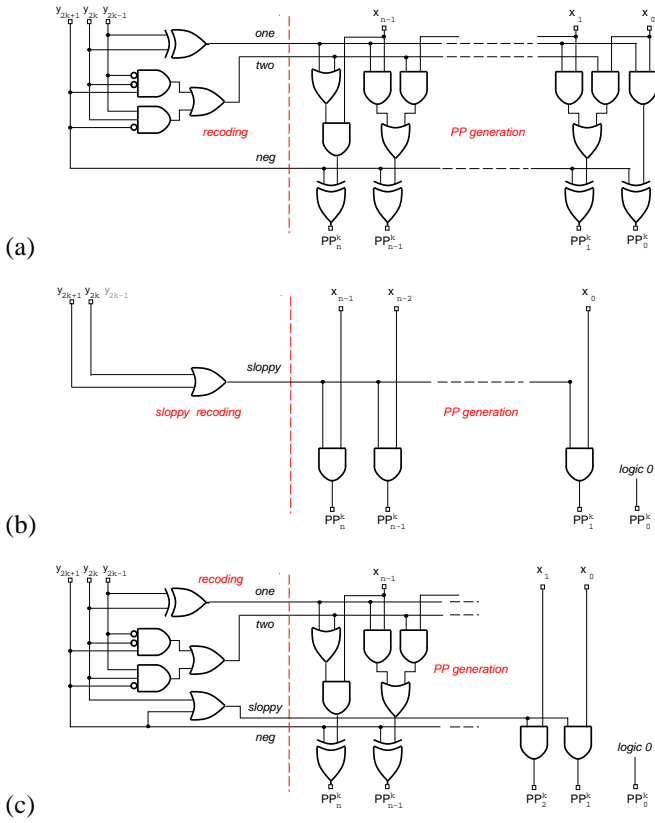


Fig. 5. Implementation of radix-4 rec+PPgen: (a) error-free, (b) whole row sloppy, and (c) sloppy columns in rows.

is sketched in Fig. 5(a).

Similarly to what was done for the addition, we have a sloppy rec+PPgen for the least-significant digits of  $y$ . The recoding is performed as shown in Table II.

The resulting hardware implementation is greatly simplified as shown in Fig. 5(b). Fig. 6(b) shows how the sloppy bits  $\circ$  are arranged in the array. As the average error for sloppy recoding is zero (Table II), for patterns in which two adjacent digits in  $y$  are  $1 = (01)_2$  and  $3 = (11)_2$  the errors on two different rows compensate in the internal columns of the array. This is shown in the example of Table III for  $(0111)_4 \times (0231)_4$ .

From Fig. 6(b) it is clear that the error due to sloppy rows can propagate well into the most-significant digits of the product. To limit this propagation, we opt for a hybrid row in which only the least-significant digits of the row are computed sloppy as shown in Fig. 5(c). With this scheme, called in the following *sloppy-columns*, the propagation of the error can be limited to

radix-4 digit		PP <sub>k</sub>		error
$y_{2k+1}$	$y_{2k}$	standard	sloppy	$\epsilon_k$
0	0	0	0	0
0	1	$x \cdot 4^k$	$2x \cdot 4^k$	$x \cdot 4^k$
1	0	$2x \cdot 4^k$	$2x \cdot 4^k$	0
1	1	$3x \cdot 4^k$	$2x \cdot 4^k$	$-x \cdot 4^k$

TABLE II  
SLOPPY RADIX-4 RECODING.

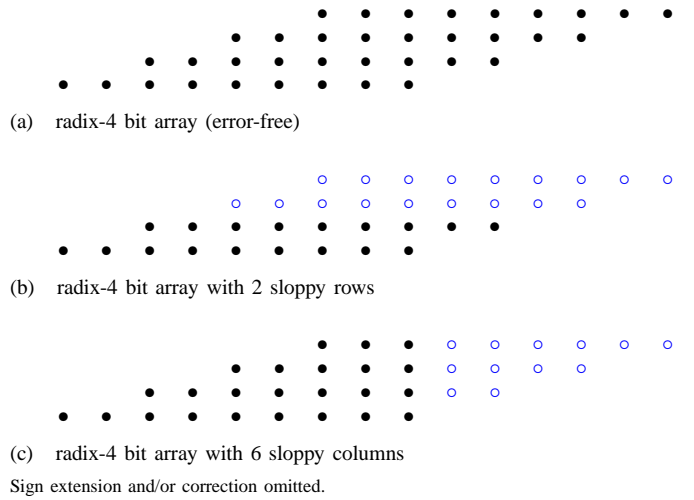


Fig. 6. Bit array and sloppy bits  $\circ$  for radix-4  $8 \times 8$  multipliers. Sign extension and/or correction omitted.

Example:  $(0111)_4 \times (0231)_4 = (0032301)_4$

			<b>0</b>	<b>2</b>	<b>2</b>	<b>2</b>	<i>sloppy</i>
		<b>0</b>	<b>2</b>	<b>2</b>	<b>2</b>		<i>sloppy</i>
	0	2	2	2			<i>regular</i>
0	0	0	0				<i>regular</i>
	0	0	3	1	3	0	2

Errors (boldface) are in radix-4 columns 0 and 3.

TABLE III  
EXAMPLE OF ERROR COMPENSATION IN INTERNAL COLUMNS.

a given column Fig. 6(c).

Again, a competitor of the sloppy scheme is the truncated one. To compare performance and error introduced, we implemented a  $12 \times 12$ -bit multiplier (two's complement) in the following schemes:

- 1) **r2-mult** a radix-2 standard multiplier;
- 2) **r4-mult** a radix-4 standard multiplier (with PPs generation as in Fig. 5(a));
- 3) **r2-trunc** a r2-mult with  $t$  truncated bits;
- 4) **r4-trunc** a r4-mult with  $t$  truncated bits;
- 5) **sloppy-rows** a radix-4 multiplier with PPs generation as in Fig. 5(b) for  $k$  multiplier radix-4 digits (rows).
- 6) **sloppy-cols** a radix-4 multiplier with PPs generation as in Fig. 5(c) for  $t$  radix-2 columns (bits).

As done for the adder, we report the mean error as function of the imprecise digits/bits in Fig. 7 and the power dissipation as function of the error in Fig. 8. The power dissipation of the precise multipliers is  $P_{r2} = 0.53 \text{ mW}$  for the radix-2 and  $P_{r4} = 0.47 \text{ mW}$  for the radix-4 multiplier. The power figures do not include the final carry-propagate adder.

Fig. 8 shows that among the truncated schemes, radix-4 is by far more power efficient than radix-2 because of the reduced number of PPs. Moreover, from Fig. 8 we derive that the sloppy row schemes with  $k = 1, 2, 3$  have very similar characteristics (error and power) to those of radix-4 truncated

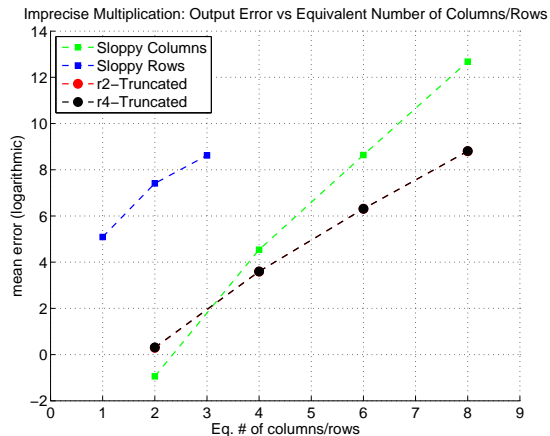


Fig. 7. Error of multiplier implementations for different imprecise schemes. Error curves for radix-2 and radix-4 truncated schemes overlap.

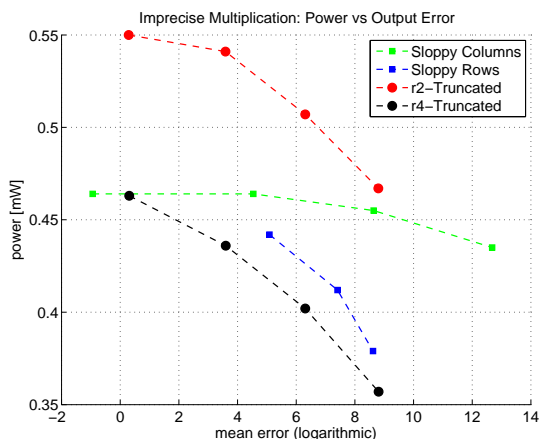


Fig. 8. Power dissipation for imprecise multipliers as function of the error.

to  $t = 4, 6, 8$  bits, respectively.

#### IV. APPLICATIONS IN IMAGE PROCESSING

To verify the figures found in the stand-alone characterization of the imprecise operators, we implement some common image processing algorithms in imprecise hardware and evaluate the performance.

As sample pictures, we used the two grayscale (each pixel is an unsigned 8-bit integer) images of Fig. 9 (upper part).

##### A. Image Filtering

We use the sloppy adder defined in Sec. II with  $k = 4$  sloppy bits to process two  $256 \times 256$  grayscale images of Fig. 9 (top) for the following bidimensional filters:

- 1) an averaging (low-pass) filter;
- 2) a sharpening filter;
- 3) an edge-detection unit.

These filters can be implemented in the the spatial domain by addition and shift operations.

The error is evaluated by taking the absolute value of the difference between the precise  $I^{ef}$  and sloppy  $I^{sl}$  value of

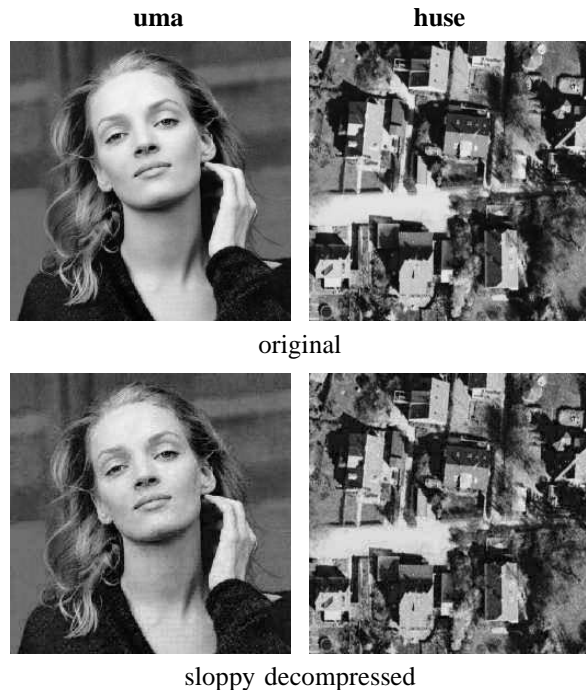


Fig. 9. Original pictures (top) and after decoding by **sloppy-row-2** IDCT (bottom).

	<i>smoothing</i>		<i>sharpening</i>		<i>edge det.</i>	
	$\epsilon_{max}$	$\bar{\epsilon}$	$\epsilon_{max}$	$\bar{\epsilon}$	$\epsilon_{max}$	$\bar{\epsilon}$
<b>uma</b>	26	7.2	60	18.9	64	9.0
<b>huse</b>	28	7.8	59	17.5	68	9.2

TABLE IV  
ERROR IN 2D-FILTERED IMAGES.

intensity (luminosity) per pixel  $(i, j)$ :

$$\epsilon_{i,j} = |I_{i,j}^{ef} - I_{i,j}^{sl}|$$

The maximum error  $\epsilon_{max}$  and the average error  $\bar{\epsilon} = \frac{\sum \epsilon_{i,j}}{N^2}$  are reported in Table IV for the different types of filtering. The results show that the degradation is independent of the image (**uma** is a portrait, while **huse** has greater detail). Depending on the filter mask, we can change the design of the sloppy adder to obtain larger savings. For example, for edge-detection, a sloppy adder with  $k = 6$  has an average error  $\bar{\epsilon} = 28$ , but visually, the degradation is not noticeable.

##### B. Inverse Discrete Cosine Transformation (IDCT)

Now we combine the imprecise multiplier schemes with an error-free adder in a multiply-add (and accumulate) unit (Fig. 10) which can be used for the trivial implementation of the Inverse Discrete Cosine Transform (IDCT), which is part of the JPEG decompression algorithm.

For the unit of Fig. 10 we opted for carry-save (error-free) accumulation to keep separate the imprecision due to the multiplier and to the adder. Based on the results of software simulations, we decided not to use a sloppy adder as the extra

