

Low power Image Processing: Analog versus Digital comparison.

Jacques-Olivier Klein*, Lionel Lacassagne*, Sébastien Moutault*, Antoine Dupret*

Abstract — In this paper, a programmable analog retina is presented and compared with state of the art MPU for embedded imaging applications. The comparison is based on the energy requirement to implement the same image processing task. Results showed that analog processing requires lower power consumption than digital processing. In addition, the execution time is shorter since the size of the retina is reasonably large.

1 INTRODUCTION

Smart sensors, vision chips[3, 4, 5, 6] have potential to take an increasing part in navigation or surveillance systems: toys or industrial robots, car driving assistance... For this class of applications, one has to provide vision systems which feature high processing capabilities, low cost, compactness and reduced power consumption. In a previous paper[10] we introduced the architecture of the X-Cell, a universal analog computation cell. Compared to its digital counterpart, lower power consumption and reduced silicium area are expected. Such statement has to be proven with fairly quantitative study. Consequently, we propose a comparison between a vector of X-Cell dedicated to image processing called PARIS and a similar digital architecture comprising SIMD units: PowerPC G4 AltiVec. This comparison is performed using well-known algorithm, representative of image processing task: edge detector. We present a detailed implementation on both architectures and focus on the hot spots for an optimized implementation. Two benchmarks are provided, the first one is about the execution time only to estimated the efficiency of general purpose processor as a challenger to dedicated architectures, the second deals with the most embedded constraining criterion: power consumption.

2 PARIS ARCHITECTURE

In most vision chips, photodetectors form an array.

With our programmable approach, photodetectors are associated to memory elements, them also organized in array. These arrays are bordered on

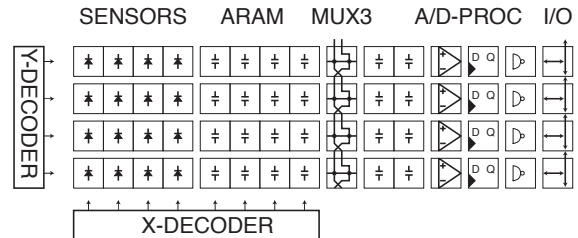


Figure 1: Array decoder architecture.

one of their side by a column of analog/digital processors (see Fig. 1). Operations are performed sequentially on columns while snapshot mode image acquisition is concurrently achieved. A decoder selects then the column reached by processors. Furthermore, each processor access to a set of rows by the way of a mux (MUX3). Finally, fully random addressing can be convenient for reading and writing images.

2.1 Architecture of rows

Each row of the retina is organized around two mixed analog-digital buses used to connect various functional units (see Figure 2). The functional units which can compose the row of a vision chip are: the rows of photosensors, the row of analog memory map, the set of analog registers, the Analog Processing Unit (X-Cell), the Boolean Processing Unit and few special registers. These last are notably required for I/O and global operators. In each processor, linear processings are handle by the analog processing unit. Boolean units associated to the condition register allows to achieve different operations according to locally stored values. Binary data stemming from a comparator are combined by the Boolean Processing Unit and can be written in a condition register. Mixed registers will then be modified wherever this condition is true. Such architecture paves the way to numerous linear, isotropic or not algorithms [8].

2.2 Generic functional units

Derived from [10], each functional unit is organized around one OTA, a set of capacitors associated to switches and of two buses: a global one, and a lo-

*Institut d'Electronique Fondamentale, Bât. 220, Univ. Paris-Sud, France, e-mail: Jacques-Olivier.Klein@ief.u-psud.fr, tel.: +33 1 69156572, fax: +33 1 69154000.

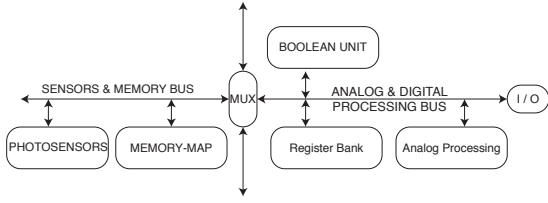


Figure 2: Architecture of rows

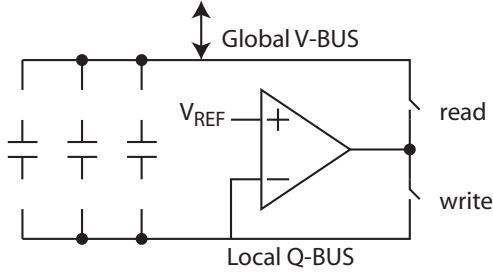


Figure 3: Generic Functional Unit

cal one (see Figure 3). The global bus, which is dedicated to inputs/outputs, is named $V-BUS$. It is intended to distribute a value represented by a voltage, therefore allowing to realize non-destructive copies. The voltage is forced by the output of one OTA or by the output of a digital cell. A voltage mode operating drastically reduce its sensitivity to parasitic capacitors. The local $Q-BUS$, is intended to realize charge transfers and balancing. The charge transfer is used to perform accumulations while division is based on charge balancing. The voltage of the $Q-BUS$ is set to V_{REF} by the output of one OTA thanks to a feedback. So, its parasitic capacitor keeps its charge and thus has little impact during the transfer of charges [10].

2.3 Operating with switched capacitors

All the functional units are based on switched capacitors structures. Four different operations are used. They are illustrated by an example on the scheme given figure 4. At the instant all the switches close, the charge of all the capacitors are modified:

1. The capacitor C_0 , is shorten, thus reset.
2. The capacitors C_1 and C_2 are also emptied of their charges, Q_1 and Q_2 , which flow by way of the $Q-BUS$ to capacitors C_3 and C_4 . It is a cumulative transfer of charges.

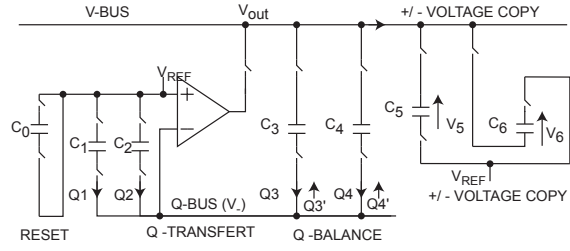


Figure 4: Operation of the switched capacitors

3. The total charge $Q_1 + Q_2 + Q_3 + Q_4$ divides up between two parallel capacitors, C_3 and C_4 , in proportion to their respective capacitance. It makes charge balancing.
4. The resulting voltage on capacitors C_3 and C_4 is copied onto capacitors C_5 and C_6 by way of the $V-BUS$. It makes a copy in voltage mode. The configuration of switches allows to do or not a change of sign by reversal of the target capacitor during the copy.

2.4 Analog processing unit

The analog processor is constituted by a set of capacitors associated to switches allowing various configurations. includes a set T of processing capacitors associated to registers-capacitors (cf. Fig. 5). To improve accuracy, each capacitor is an instance of a unitary capacitor C_u . Let define the *weight* of a set S of capacitors, the dimensionless quantity: $\frac{1}{C_u} \times \sum_{i \in S} C_i$, where C_i is the capacitance of the i th capacitor of S .

More general operation of the analog processor, multiplication-accumulation can be decomposed into three steps: *Load*, *Distribute*, *Accumulate*. For each of these 3 steps, a set of the implied capacitor (respectively L, B, A) is considered.

- During the first step (Load), the set $L \subset T$ (of weight l) is charged by one or more positive or negative copies. Each input voltage V_n is copied (positively or negatively) in one subset $L_n \subset L$ of capacitors (of weight l_n) so that $L = \bigcup_n L_n$ and $L_i \cap L_j = \emptyset$ for all $i \neq j$. After N loads, the charge Q_L , stored in set L , is $Q_L = \sum_n \pm l_n \times V_n \times C_u$
- During the second step (Balancing), the charge Q_L is distributed on the set $B \supset L$ (of weight b), so that each capacitor belonging to B has a voltage $V_B = \frac{1}{b} \times \sum_n \pm l_n \times V_n$
- Finally, the last step (Accumulation) consists in adding charges stored in a set of capacitors

$A \subset B$, of weight a , on a register-capacitor C_R of capacitance C_u . So: $V_{C_R}(t+1) = V_{C_R}(t) + \frac{a}{b} \sum_n \pm l_n \times V_n$

Hence, the realized operation is a set of multiplication/accumulation of coefficient $\frac{A}{B} \times L_n$. Obviously, if $B = 1$, and A is an integer lower than 8 step 2 can be omitted. As a consequence, the MAC instruction duration is 2 or 3 cycles. Table ?? describes a subset of the X-Cell instructions. AR a,d AAR represents any analog register and DR and DAC any digital register. LC stands for *Local Condition*.

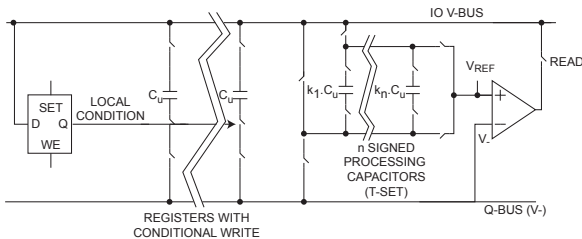


Figure 5: Architecture of the Analog Processing Unit.

Instruction	Description	Cyc.
MAC AR A	$AAC \leftarrow AAC + A \times AR$	2
MAC AR A/B	$AAC \leftarrow AAC + A/B \times AR$	3
ASTR AR	$AR \leftarrow LC?AAC : AR$	1
ARST	$AAC \leftarrow 0$	1
CMP	$DAC \leftarrow (ACC > 0)$	1
WHR	$LC \leftarrow DAC$	1
WHRN	$LC \leftarrow not DAC$	1
UWHR	$LC \leftarrow TRUE$	1
AND DR	$DAC \leftarrow DAC and DR$	1
OR DR	$DAC \leftarrow DAC or DR$	1
NAND DR	$DAC \leftarrow DAC and not DR$	1
SET	$DAC \leftarrow TRUE$	1
DRST	$DAC \leftarrow FALSE$	1
DSTR	$DR \leftarrow DAC$	1

Table 1: X-Cell Instruction subset

3 PHYSICAL IMPLEMENTATION

Two retinas prototypes were designed. Although the first, *PARIS I*, is based on a slightly different structure from the universal structure described here, its functioning is somewhat identical. It is consisted of 16×16 pixel array - each including a photosensors and 3 analog memory elements - associated to a minimal analog processor including only four capacitors: three for processing and one for register[8]. Its main characteristics are presented in the table 2.

Parameter	PARIS I	PARIS II
Resolution	16×16	256×256
Processor	16	256
Pixel Size	$50 \times 50 \mu m^2$	$25 \times 25 \mu m^2$
Max Frequency	10MHz	40MHz
Power cons.	30mW	800mW
MixtRegisters	2	3
Resolution Processing	7-bits	10-bits
Capacitors Boolean	2	4
Processor I/O	No 1 analog	Yes 1 analog 8 digital
Reduction Operator	No	1 global-OR 1 Mean Op

Table 2: Paris I and PARIS II parameters

This circuit has been successfully tested and operates properly [13]. It is currently being evaluated for applications in mobile robotics. The second circuit, *PARIS II*, was designed according to the principle described in this paper. It brings improvements with regard to *PARIS I*, notably on reading circuits of analog memory and photosensors [12]. Its main characteristics are presented in the table 2.

4 DERICHE BENCHMARK

In order to estimate the performance of the X-Cell architecture, we have decided to compare it to another SIMD vector architecture and to implement a de facto image processing algorithm like edge detection. The closest "software" architecture are the general purpose processor with multimedia SIMD extension (also called SWAR for SIMD Within A Register). The most embedded GPP are the PowerPC AltiVec and Intel Centrino. PowerPC has a more extensive SIMD ISA for image processing (crossbar capabilities, reductions and 8-bit multiplier) Centrino implements SSE2 but with only 16 multipliers, Pentium4 Prescott extends SSE2 instructions with reduction capabilities with SSE3, but can not be considered as an "embedded" processor. Note that an SoC version of the PowerPC G4 has been released by Motorola/Freescale Other embedded processors might be chosen for their SIMD: the ARM11 (SIMD in 32-bit registers: four 8-bit computations in parallel) or the latest Intel Xcale/PCA which includes a multimedia extension called Wireless MMX (64-bit registers for 8/16/32-bit integer and 32-bit FP).

Classical edges detector operators implemented

in artificial retinas FIR filters like Sobel, Prewitt or Roberts filters. Canny-Deriche filters have assert themselves for their robustness. These filters can be expressed as a non recursive filter like Canny's filter or a recursive filter like Deriche's one. Each have drawback and advantage : Deriche have a fixed complexity that does no depend on the smoother coefficient, but requires large memory to hold a complete image, Canny is more adapted to "data-flow" because the image must not be store in memory, only the current raw, but the filter size depends on the smoother coefficient.

X-Cell is well-adapted to Deriche filter: it has three planes to store 3 images, and the performances of the processor vector array are not limited by Deriche's filter structure, if the vector displacement is orthogonal to the filter. The Deriche's filter complexity has been reduced by a factor two by Garcia Lorca [16]. That is this filter that will be implemented.

The second order filter is:

$$y(n) = b_0x(n) + a_1y(n-1) + a_2y(n-2)$$

with:

$$\gamma = e^{-\alpha} \quad b_0 = (1-\gamma)^2 \quad a_1 = 2\gamma \quad a_2 = -\gamma^2$$

4.1 2D filter implementation

The Q8 fixed-radix code Deriche H & V smoothers are:

```
for(i=0; i<n; i++)
  for(j=0; j<n; j++)
    x0 = X[i][j]
    y1 = Y[i][j-1]
    y2 = Y[i][j-2]
    y0 = (b0.x0+a1.y1+a2.y2) >> 8
    Y[i][j] = y0
```

Deriche H

```
for(j=0; j<n; j++)
  for(i=0; i<n; i++)
    x0 = X[i][j]
    y1 = Y[i-1][j]
    y2 = Y[i-2][j]
    y0 = (b0.x0+a1.y1+a2.y2) >> 8
    Y[i][j] = y0
```

Deriche V

$$b_0=256 \times b_0 \quad a_1= 256 \times a_1 \quad a_2= 256 \times a_2$$

4.2 PowerPC AltiVec implementation

The three main problems to address for SIMD implementation are:

- cache impact
- recursive filter structure
- underflow

The horizontal filter does not generate cache miss whereas the vertical filter does. The solution is to permute the internal loop with the external loop of the filter to obtain an horizontal-like scan with a vertical filter. Such a permutation correspond to a cache blocking optimization [15].

```
for(i=0; i<n; i++)
  for(j=0; j<n; j++)
    x0 = X[i][j]
    y1 = Y[i-1][j]
    y2 = Y[i-2][j]
    y0 = (b0.x0+a1.y1+a2.y2) >> 8
    Y[i][j] = y0
```

Deriche VH

Between two iterations of the filter there is a loop-carried dependency. The solution proposed (figure 6) is to perform a block-transposition of pixel into a band, to process the band and then to perform a second block-transposition into the source image.

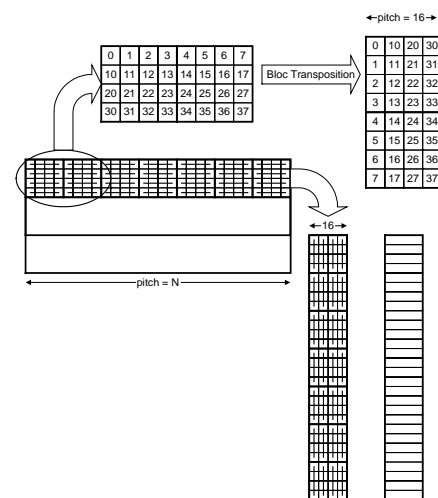


Figure 6: Deriche Band transposition

The last problem is about underflow: since the coef a_2 is negative, for a long set of zero input values, one can have $x_0 = 0$, $y_1 = 0$ but $y_2 \neq 0$, so an underflow can happen.

4.3 X-Cell implementation

The following pseudo-code sources describe the primitives used for edge detection. The program iterates on all this routines for each column.

```

ARST
MAC I1(i, 0) 0.375
MAC I2(i-1, 0) 0.875
MAC I2(i-2, 0) 0.25
ASTR I2(i, 0)

```

PARIS horizontal smoother

```

ARST
MAC I2(i, 0) -1
MAC I2(i, 1) 1
MAC I2(i+1, 0) 1
MAC I2(i+1, 1) 1
ASTR ARO
CMP

```

PARIS horizontal gradient

```

ARST
MACC ARO 1
WHRN
ASTR ARO
UWHR

```

Horizontal absolute value

```

ARST
MAC I2(i, 0) -1
MAC I2(i+1, 0) 1
MAC I2(i, 1) 1
MAC I2(i+1, 1) 1
ASTR AR1
CMP

```

vertical gradient

```

ARST
MACC AR1 1
WHRN
ASTR AR1
UWHR

```

Vertical absolute value

```

ARST
MACC ARO 1
MACC AR1 1
ASTR I2(i, 0)

```

Addition of the two previous results

Each one of this four filters requires, for each column, 1 reset, 3 MAC and 1 write-back instruction. These instructions are performed in eleven cycles, so the output image requires 2816 cycles and applying the four filters requires 11264 cycles. Gradient calculation costs twelve 2-cycle MAC and sixteen other instructions, that means 40 cycles for each column and 10240 for a entire image. All things considered, the algorithm is performed 21504 cycles, i.e. 0.54 ms at 40MHz.

5 RESULTS & ANALYSIS

To observe the impact of cache behavior we use the *cpp* (Cycle Per Pixel):

$$cpp = \frac{t \times F}{n^2}$$

n	128	256	512	1024
<i>cpp</i> Deriche H	2.95	2.85	3.31	3.87
<i>cpp</i> Deriche VH	4.86	4.88	5.24	6.19
<i>cpp</i> gradient	2.69	2.88	3.17	3.65
<i>cpp</i> total	10.5	10.61	11.72	13.71

Table 3: *cpp* for 128, 256, 512 and 1024 image size for PowerPC

n	128	256	512	1024
t(ms) Deriche H	0.048	0.187	0.868	4.058
t(ms) Deriche VH	0.080	0.320	1.374	6.491
t(ms) gradient	0.044	0.189	0.831	3.827
t(ms) total	0.172	0.696	3.073	14.376

Table 4: execution time (ms) for 128, 256, 512 and 1024 image size for PowerPC

The execution time on the Xcell does not suffer from cache misses: *cpp* is still constant: 11 cycles per Deriche filter, for a total of 44 for the four filters and 40 cycles for the gradient.

n	128	256	512	1024
time(ms) total	0.269	0.538	1.075	2.15

Table 5: execution time (ms) for 128, 256, 512 and 1024 image size for XCell

If we only compare the execution time, PowerPC and Xcell run at same speed (the G4 is even faster), for small images (128 and 256), when they fit the G4 cache. Such a comparison is biased since it does not take into account the required energy for these architecture.

The classical metric used to compare embedded processor is *Mips/Watt*. We do not believe that Mips or Mops is still an up-to-date metric since the latency of instructions may vary a lot, and so, counting the number of instructions could lead to erroneous conclusion except if you want your system to run the Dhrystone benchmark. Not very useful. We prefer the $t \times Watt$ (in Joule) which is the amount of energy required to apply an algorithm on an image. The idea is that if a processor is by far real-time for an application, its SoC version will use a *downclocked* version of the classic processor version, the energy remains constant but the power is smaller. For 256×256 images the classic G4 is 78 times faster than the realtime constraint (40 ms). Dividing its clock frequency by 10 will also reduce its power consumption by approximately 10, for a still realtime 5.1 ms execution.

$$E = t \times Watt$$

The technology used for the current XCell processor is $0.60 \mu\text{m}$. Switching from 0.60 to $0.25 \mu\text{m}$ will decrease the capacitor surface, that is the leaking capacitor, the required current and finally the consumption. A scale factor can be applied to estimate not a faster XCell but *smaller* XCell. The factor is $(0.60/0.25)^{1.5}$ the exponent is 1.5 and not 2 since it appears in the Literature that such a switch provides a factor that is smaller than the gain in surface, and closer to 1.5 than 2. For XCell we estimated the consumption of the micro-controller to 200 mW and 800mW for a 256 XCell vector. With such assumption, the result for the new criterion is:

n	128	256	512	1024
PowerPC (mJ)	1.72	6.96	30.73	143.76
XCell (mJ)	0.16	0.54	1.94	7.31
scaled XCell (mJ)	0.07	0.24	0.86	3.26
gain	10.7	12.9	15.9	19.7
scaled gain	23.9	29.0	35.6	44.1

Table 6: Comparison of required energy for PowerPC and XCell

With such criterion, the difference of performances for *extreme* embedded applications is more realistic from our point of view.

6 CONCLUSION

A programmable analog retina has been presented and compared with state of the art MPU for embedded imaging applications. The comparison is based on the energy requirement to implement the same image processing task. Each version has been independently optimized to fit the considered architecture. To complete the performance evaluation, an evaluation of 1GHz DSP C64x is planned. Right now, the validity of such an analog design has been demonstrated. Even when obsolete processes are used for the retina, results showed that analog processing requires lower power consumption than digital processing. In addition, the execution time is shorter since the size of the retina is reasonably large.

References

[1] Moutault S. Klein J.-O., Dupret A. "A universal switched capacitors operator for the automatic synthesis of analog computation circuits.", CAMP'03 - IEEE Sixth International Workshop on Computer Architecture for Machine Perception, New Orleans (USA), May 12-14, 2003, ISBN 0-7803-7971-3.

[2] B. Granado, L. Lacassagne, P. Garda, "Cab general purpose microprocessors simulate neural network in real time", IWAN (2), 1999, pp 21-29.

[3] C. A. Mead, "Analog VLSI and neural systems", Addison Wesley, 1989.

[4] A. Moini, "Vision Chips", Kluwer Academic, 1999.

[5] K. Kyuma, E. Lange, J. Ohta, A. Hermanns, B. Banish, and M. Oita, "Artificial Retinas-Fast, Versatile Image Processors", Nature, vol. 372, 1994.

[6] T. Bernard, B. Zavidovique, and F. Devos, "A Programmable Artificial Retina", IEEE Journal of Solid State Circuits, vol. 28, pp. 789-797, 1993.

[7] R. Etienne-Cummings, Z. Kalayijian and D. Cai, "A programmable focal-plane MIMD image processor chip", IEEE J. Solid State Circuit, vol. 36, N1, pp 64-73, January 2001.

[8] A. Dupret, J.-O. Klein, A. Nshare, "A programmable vision chip for CNN based algorithms", in Proc of 6th IEEE Int Workshop on Cellular Neural Networks and their Applications, pp.207-212, 23-25 may 2000, Catania, Italy

[9] A. Abo, "Design for reliability of low-voltage switched-capacitor circuits", PhD Thesis, University of California, Berkley, May 1999.

[10] S. Moutault, J.O. Klein, A. Dupret, "A universal switched capacitors operator for the automatic synthesis of analog computation circuits.", IEEE Sixth International Workshop on Computer Architecture for Machine Perception, May 11-14, 2003, New Orleans, LA.

[11] R. Carmona, S. Espejo, R. Dominguez-Castro, A. Rodriguez-Vazquez, T. Roska, T. Kozek, L.O. Chua, "A $0.5 \mu\text{m}$ CMOS CNN Analog random access memory chip for massive image processing", Proc of 5th IEEE Int Workshop on Cellular Neural Networks and their applications, pp 243-248, London, UK, April 1998.

[12] Nshare A., Klein J.O., Dupret A. "An improved ARAM for PARIS, an original vision chip", CNNA2002, Frankreset/Main, 2224 July 2002, World Publishing 2002, pp 347-354.

[13] A. Dupret, J.O. Klein, A. Nshare, "A DSP-like analog processing unit for smart image sensors" International Journal of Circuit Theory and Applications, 2002, vol. 30, pp.595-609.

- [14] R. Deriche. "Using Canny's criteria to derive a recursively implemented optimal edge detector". *The International Journal of Computer Vision*, 1(2):167-187, May 1987.
- [15] D. Demigny *et al* "Méthodes et architectures pour le traitement du signal en temps rel. *Traité Information, Commande et Communication*. Hermès. 2001.
- [16] F. Garcia Lorca "Filtres récursifs temps réel pour la détection de contours : optimisations algorithmiques et architecturales" PhD Thesis 1996.
- [17] L. Lacassagne, F. Lohier, PP. Garda, "Realtime execution of optimal edge detectors on RISC and DSP processors", ICASSP 1998.