



**HAL**  
open science

## Memory aware HLS and the implementation of ageing vectors

Gwenolé Corre, Eric Senn, Nathalie Julien, Eric Martin

► **To cite this version:**

Gwenolé Corre, Eric Senn, Nathalie Julien, Eric Martin. Memory aware HLS and the implementation of ageing vectors. EUROMICRO Symposium on Digital System Design, 2004, Rennes, France. pp.88-95. hal-00077366

**HAL Id: hal-00077366**

**<https://hal.science/hal-00077366>**

Submitted on 30 May 2006

**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.

# Memory aware HLS and the implementation of ageing vectors

Gwenolé Corre, Eric Senn, Nathalie Julien and Eric Martin  
LESTER / University of South Brittany  
BP92116, 56321 Lorient cedex, France  
Email: gwenole.corre@univ-ubs.fr

## Abstract

*We introduce a new approach to take into account the memory architecture and the memory mapping in behavioral synthesis. We formalize the memory mapping as a set of constraints for the synthesis, and defined a Memory Constraint Graph and an accessibility criterion to be used in the scheduling step. We present a new strategy for implementing signals (ageing vectors). We formalize the maturing process and explain how it may generate memory conflicts over several iterations of the algorithm. The final Compatibility Graph indicates the set of valid mappings for every signal. Several experiments are performed with our HLS tool GAUT. Our scheduling algorithm exhibits a relatively low complexity that permits to tackle complex designs in a reasonable time.*

## 1. Introduction

Behavioral synthesis, which is the process of generating automatically an RTL design from an algorithmic description, is an important research area in design automation. Many behavioral specifications, especially in digital signal and image processing, use arrays to represent, store and manipulate ever growing amounts of data. The ITRS roadmap indicates that, in 2011, 90 % of the SoC area will be dedicated to the memory [1]. Applications are indeed becoming more and more complex, and memory will take a more and more important place in future signal processing systems. This place is vital, strategic, for memory now appears as a terrific bottleneck in real-time systems. Indeed, performances are highly dependent on the memory architecture (hierarchy, number of banks) together with the way data are placed and transferred. To tackle the complexity of memory design, we consider as essential to take into account memory accesses directly during the behavioral synthesis,

assuming that a reasonable trade-off between the design time and the quality of the results is reached. In the context of HLS, several scheduling techniques actually include memory issues. Among them, most try to reduce the memory cost by estimating the needs in terms of number of registers for a given scheduling, but work only with scalars [2]. Some of them really schedule the memory accesses [3]. They include precise temporal models of those accesses, and try to improve performances without considering the possibility of simultaneous accesses which would ease the subsequent task of register and memory allocation. Works in [4] include the memory during HLS, but is dedicated to control intensive applications. In [5], a first scheduling (force directed) is performed on a Data Flow Graph (DFG); the memory accesses are then rescheduled after the selection and memory allocation to reduce the overall memory cost. The complexity of this scheduling algorithm, however, does not allow to target realistic applications in a reasonable time. In [4], memory accesses are represented as multi-cycle operations in a Control and Data Flow Graph (CDFG). Memory vertices are scheduled as operative vertices by considering conflicts among data accesses. This technique is used in some industrial HLS tools that include memory mapping in their design flow (Monet, Behavioral Compiler) [6]. Memory accesses are regarded as Input/Output. The I/O behavior and number of control step are managed in function of the scheduling mode [7]. In practice, the number of nodes in their input specifications must be limited, to obtain a realistic and satisfying architectural solution. This limitation is again mainly due to the complexity of the algorithms which are used for the scheduling.

In this paper, we propose a new and simple technique to take into account the memory mapping in the architectural synthesis. Indeed, our aim is to produce a simple algorithm to achieve the synthesis of even complex designs in a reasonable time. We focus on the definition of a memory mapping file that

is used in the synthesis process. We introduce an original scheduling in the synthesis flow, to obtain an optimized RTL design. This scheduling technique is described in section 2 with the formalism to resolve scheduling under memory constraint. The special case of ageing data is discussed in section 3. Experimental results are presented in section 4.

## 2. Memory integration

### 2.1. Memory aware synthesis

We introduce memory synthesis in the standard HLS design flow. A Signal Flow Graph (SFG) is first generated from the algorithmic specification. In the new approach, this SFG is parsed and a memory table is created (see Fig. 1). This memory table is then completed by the designer who can select the variable implementation (memory or register) and place the variable in the memory hierarchy (which bank). The resulting table is the memory mapping that will be used in the synthesis. In the standard flow, the processing unit is synthesized without any knowledge on the memory mapping. The memory architecture is designed afterward and a lot of optimization opportunities are definitely lost.

Name	Class	Implementation	Bank	Address	Initial Value
adapt	Variable	Register	-1	-1	0
deux_mu	Constant	Memory	0	0	0
h(0)	LoopBack	Memory	0	1	0
h(1)	LoopBack	Memory	0	2	0
h(2)	LoopBack	Memory	0	3	0
h(3)	LoopBack	Memory	0	4	0
x(0)	Constant	Memory	1	0	0
x(1)	Delay	Memory	1	1	0
x(2)	Delay	Memory	1	2	0
x(3)	Delay	Memory	1	3	0
y	Variable	Register	-1	-1	0
y0001	Variable	Register	-1	-1	0
y0002	Variable	Register	-1	-1	0
y0003	Variable	Register	-1	-1	0

Figure 1: Memory table

The memory mapping file contains information about every data structure in the algorithm (mainly arrays in DSP applications) and its allocation in memory (bank number and physical address). Scalars can also be defined. This memory table represents all data vertices extracted from a SFG. This data distribution can be static or dynamic. In the case of a static placement, the data stay at the same place during the whole execution. If the placement is dynamic, data can be transferred between different levels in the memory hierarchy. Thus, several data can share the same location in the circuit memory.

The memory mapping file explicitly describes the data transfers to occur during the algorithm execution. Direct Memory Address (DMA) directives will be added to the code to achieve these transfers. The definition of the memory architecture will be performed in the first step of the overall design flow. To achieve this task, advanced compilers such as Rice HPF compiler, Illinois Polaris or Stanford SUIF could be used [8]. Indeed, these compilers automatically perform data distribution across banks, determine which access goes to which bank, and then schedule to avoid bank conflicts. The Data Transfer and Storage Exploration (DTSE) method from IMEC and the associated tools (ATOMIUM, ADOPT) are also a good mean to determine a convenient data mapping [9].

### 2.2. Signal Flow Graph

The input of our HLS tool is an algorithmic description that specifies the circuit’s functionality at the behavioral level, disregarding any potential implementation solutions. This initial description is compiled in order to obtain an intermediate representation: the Signal Flow Graph (SFG), Fig. 2. A Signal Flow Graph is a directed polar graph  $SFG(V, E)$  where the set of vertices  $V = \{v_0, \dots, v_n\}$  represents the operations,  $v_0$  and  $v_n$  are respectively the source vertex and the sink vertex. The set of edges  $E = \{(v_i, v_j)\}$  represents the dependencies between the operations vertices. The Signal Flow Graph contains  $|V| = n + 1$  vertices. A vertex represents one of the following operations: arithmetic, logical, data or delay. The difference between a Signal Flow Graph and Data Flow Graph resides in the introduction of delay operators ( $z^{-k}$ ). These operators are necessary to express the use of data whose value was computed in a preceding iteration of the algorithm. An edge  $E_{i, j} = (v_i, v_j)$  represents a data dependence between operations  $v_i$  and  $v_j$  such as for any iteration of the SFG, operation  $v_i$  must start its execution before that of  $v_j$ . For the data dependencies, the execution of  $v_j$  can start only after the completion of operation  $v_i$ .

### 2.3. Memory Constraint Graph

As outlined in section ; all data vertices are extracted from the SFG to construct the memory table. The designer can choose the data to be placed in memory and defines a memory mapping. For every memory in the memory table, we construct a weighted Memory Constraint Graph (MCG). It represents conflicts and scheduling possibilities

between all nodes placed in this memory. The MCG is constructed from the SFG and the memory mapping file. We parse the SFG to find the data vertices. The memory table is constructed, where the designer adds mapping information. A Memory Constraint Graph is a cyclic directed polar graph  $MCG(V', E', W')$  where  $V' = \{v'0, \dots, v'n\}$  is the set of data vertices placed in memory. A memory Constraint Graph contains  $|V'| = n + 1$  vertices which represent the memory size, in term of memory elements. The set of edges  $E' = \{(v'i, v'j)\}$  represents possible consecutive memory accesses, and  $W'$  is a function that represents the access delay between two data nodes.  $W'$  has only two possible values:  $W_{seq}$  (sequential) for an adjacent memory access in memory, or  $W_{rand}$  (randomize) for a non adjacent memory access. Weight depends on the data placement defined in the memory file. There are as much sub-graphs as memory banks in memory. The memory table gives the number of banks and the address of every data in memory. This permits to construct the MGC. Figure 2 shows a MCG for the LMS filter with two simple port memory banks. The input samples  $x(i)$  are placed consecutively in one bank. The filter coefficients  $h(i)$  are placed consecutively in one another bank (dotted edges represent edges where  $W = W_{seq}$ ).

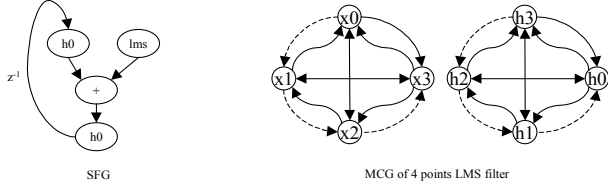


Figure 2: Signal Flow Graph and Memory constraint graph LMS

Memory constraint graphs are used during the scheduling process to determine the accessibility criterion and the time of every memory access.

## 2.4. Scheduling under memory constraint

The classical list scheduling algorithm relies on heuristics in which ready operations (operations to be scheduled) are listed by priority order. In our tool, an early scheduling is performed. In this scheduling, the priority function depends on the mobility criterion. This mobility is computed, for each cycle, as the difference, in number of cycles, between the current cycle and the operation deadline. Whenever two ready operations need to access the same resource (this is a so called resource conflict), the operation with the lower mobility has the highest priority and is scheduled. The

other is postponed. To perform a scheduling under memory constraint, we introduce fictive memory access operators and add an accessibility criterion based on the MCG. A memory has as much access operators as access ports. The memory is declared accessible if one of its fictive memory access operators is idle. Several operations can try to access the same memory in the same cycle; accessibility is used to determine which operations are really executable. The list of ready operations is still organized according to the mobility criterion, but all the operations that do not match the accessibility condition are removed from this list. To schedule an operation that involves an access to the memory, we check if the data is not in a busy memory bank. If a memory bank is not available, every operation that needs to access this memory will not be scheduled, no matter its priority level. The MCG is also used to compute the shortest memory access sequence when it's possible.

Our scheduling technique is illustrated in Fig. 3. The memory table is extracted from the SFG. The designer has defined two different memory mappings in memory table 1 and in memory table 2. Data a,b and c are placed at address in bank0. The constant cst is not stored in RAM. Our tool constructs two Memory Constraint Graphs  $MCG_1$  and  $MCG_2$ . For mapping 1, the sequential access sequence is  $a \rightarrow b \rightarrow c$ : it includes two dotted edges (with weight  $W_{seq}$ )  $a \rightarrow b$  and  $b \rightarrow c$ .  $MCG_2$  contains two different dotted edges:  $a \rightarrow c$  and  $c \rightarrow b$ . To deal with the memory bank access conflicts, We define a table of access for each port of a memory bank. In our example, the table has only one line for the single port memory bank0. The table of memory access has  $Data\_rate/Sequential\_access\_time$  elements. The value of each element of the table indicates if a fictive memory access operator is idle or not at the current  $C\_step$ . We use the MCG to produce a scheduling that permits to access the memory in burst mode. If two operations have the same mobility and request to the same memory bank, the operation that is scheduled is the operation that involves an access at an address consecutive with the preceding access. For example, operations `add2` and `mult1` have the same mobility. At `c_step cs_2`, they are both executable, the operands (stored in bank0) of `add2` and `mult1` are respectively data b and data c, the latest data accessed in bank0 is data a.  $MCG_1$  indicates that the access sequence  $a \rightarrow b$  is shorter than access  $a \rightarrow c$ . We schedule `add2` at `c_step cs_2` to favorize the sequential access. On a contrary, for mapping 2,  $MCG_2$  indicates that `mult1` must be scheduled before `add2`.

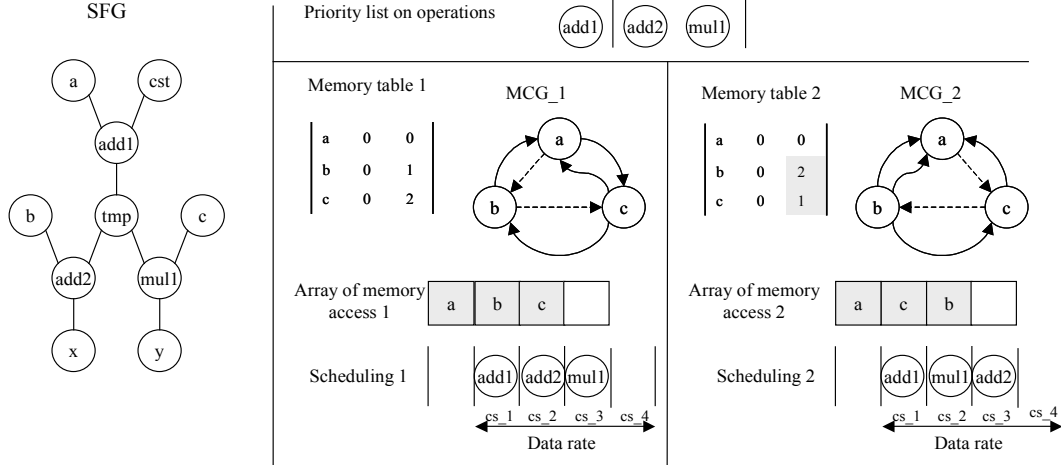


Figure 3: scheduling with MCG

### 3. Implementing ageing vectors

Signals are the input and output flows of the applications. A mono-dimensional signal  $x$  is a vector of size  $n$ , if  $n$  values of  $x$  are needed to compute the result. Every cycle, a new value for  $x$  ( $x[n+1]$ ) is sampled on the input, and the oldest value of  $x$  ( $x[0]$ ) is discarded. We called  $x$  an ageing, or maturing, vector or data. Ageing vectors are stored in RAM. A straightforward way to implement, in hardware, the maturing of a vector, is to write its new value always at the same address in memory, at the end of the vector in the case of a 1D signal for instance (that is how Monet works). Obviously, that involves to shift every other values of the signal in the memory to free the place for the new value. This shifting necessitates  $n$  reads and  $n$  writes in the memory, which is very time and power consuming. In GAUT, the new value is stored at the address of the oldest one in the vector. Only one write is needed. Obviously, the address generation is more difficult in this case, because the addresses of the samples called in the algorithm change from on cycle to the other. The Figure 4 illustrates this difficulty. In the following code a signal  $x$  is accessed; it includes  $N = 4$  elements.

```

ALGORITHM 1
x(0):=x input;
tmp := x(0);
for (i=1;i=N-1;1++) tmp=tmp+x(i);
for (i=N-1 ; i = 1 ; i--) x(i)=x(i-1);

```

The logical address of an element of  $x$  ( $x[0]$  for instance) changes from an iteration to the other. The logical address of  $x[0]$  is that of  $x3$  in iteration 3,  $x4$  in iteration 4,  $x5$  in iteration 5 etc. With GAUT, we make the distinction between physical and logical addresses. The logical address points on a memory

element that contains the physical address of the data. The physical address points on the memory element that contains the value of the data. Once determined, the physical address of a data never changes. In our example, for instance, the physical address of data  $x3$  from vector  $x$  will remains the same as long as  $x3$  is alive in the memory.

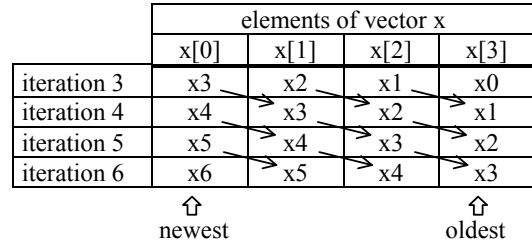


Figure 4: The maturing process

We have developed a new methodology to resolve the synthesis of our logical address generators. The advantage is a lower latency, since we avoid  $n$  reads and writes of the ageing vector, and a resulting lower power consumption. Indeed, the power consumption of a memory increases with the number of accesses.

This methodology is based on an oriented graph that traces the evolution of the logical addresses in a vector during the execution of one iteration of the algorithm: the *Logical Address Graph* (LAG). The LAG is a couple  $LAG = (V, E)$ ; it is defined for each ageing vector in the algorithm.  $V$  is the set of vertices  $V = \{v_0, v_1, \dots, v_{N-1}\}$  where vertex  $v_i$  is the  $i^{th}$  element of the vector. With  $x$  a vector of size  $N$ ;  $\text{card}(V) = N$ .  $E$  is the set of edges  $E = \{e_1, e_2, \dots, e_M\}$  where edge  $e = (v_i, v_j)$  links 2 elements  $v_i$  and  $v_j$  if the  $j^{th}$  element of the vector ( $x[j]$ ) is accessed immediately after the  $i^{th}$  element of the vector ( $x[i]$ ).  $E \subseteq V \times V$ . The weighting

function  $f$  is associated to the LAG;  $f : V \times V \rightarrow \mathbb{N}$ . For every edge  $e = (v_i, v_j)$ ,  $f$  gives the weight  $f_{ij}$  with  $f_{ij} = (j - i) \% N$ . % expresses the modulo. Figure 5 represents the LAG for the preceding example.

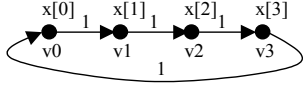


Figure 5: LAG for algorithm 1

The weight  $f_{ij}$  is used to calculate the logical address of the next access to vector  $x$ . Suppose that 0 is the logical address of  $x[0]$ . Then  $x[1]$  is the next access to  $x$  and its logical address is  $0 + f_{01} = 1$ . The logical address of  $x[2]$  is 2 and the logical address of  $x[3]$  is 3. The next data to be accessed is  $x[0]$ . Its address is still  $(3 + 1) \% 4 = 0$  in this iteration. However, to calculate the address of  $x[0]$  in the next iteration, we ought to take into account the ageing of vector  $x$ . In our example, the values in vector  $x$  are shifted so that the logical address of element  $x[i]$  at the iteration  $o+1$ , noted  $@x[i]^{o+1}$  is the logical address of element  $x[i]$  at the iteration  $o$  plus 1:  $@x[i]^{o+1} = @x[i]^o + 1$ . In general, we define the *ageing factor*  $k$  as the difference between the logical address of element  $x[i]$  at the iteration  $o+1$  and the logical address of element  $x[i]$  at the iteration  $o$ . In our example,  $k = 1$ .

$$k = @x[i]^{o+1} - @x[i]^o \quad (1)$$

Eventually, to calculate the logical address of  $x[0]$ , we add (modulo  $N$ ), to the logical address of  $x[3]$  in the preceding iteration, the weight  $f_{30}$  and the *ageing factor*  $k$  so that  $@x[0] = (@x[3] + f_{30} + k) \% N$ . More generally, if  $x[i]$  is the last element of  $x$  accessed in iteration  $o$  and  $x[j]$  is the first element of  $x$  accessed in iteration  $o+1$ , and with  $N$  the size of  $x$ :

$$@x[j]^{o+1} = (@x[i]^o + f_{ji} + k) \% N \quad (2)$$

Consider the algorithm below. This algorithm was synthesized with the following mapping:  $x[0]$  and  $x[1]$  are in a memory bank,  $x[2]$ ,  $x[3]$  and  $x[4]$  are in another bank. The relation with the logical addresses is determined for the first iteration. So  $@x[0]$  ( $= 0$ ) and  $@x[1]$  ( $= 1$ ) are in the first bank,  $@x[2]$  ( $= 2$ ),  $@x[3]$  ( $= 3$ ) and  $@x[4]$  ( $= 4$ ) are in the second.

**ALGORITHM 2**

```
x(0):=x input;
tmp=x(0);
tmp=tmp+ x(1);
tmp1=x(2);
tmp1=tmp1+x(3);
tmp1=tmp1+x(4);
for (i=N-1 ; i= 1 ; i--) x(i)=x(i-1);
```

The LAG for vector  $x$  is represented Figure 6. The chronogram of accesses is presented figure 7.

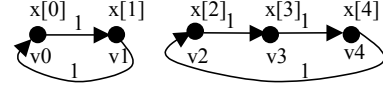


Figure 6: LAG for algorithm 2

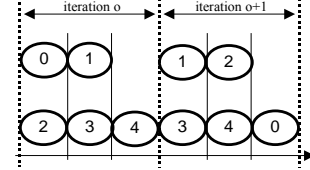


Figure 7: Chronogram of accesses

In iteration  $o$ , several concurrent accesses to the memory appear:  $@x[0] = 0$  with  $@x[2] = 2$ , and  $@x[1] = 1$  with  $@x[3] = 3$ . These parallel accesses do not generate conflict for involved data are in distinct memory banks. In the next iteration however, the logical addresses for  $x[1]$  and  $x[3]$  are respectively 2 and 4. A memory conflict is generated since these two logical addresses are mapped in the same memory bank. The concurrent accesses are computed for  $N$  successive iterations of the algorithm to obtain the concurrent accesses table (see table 1).

Table 1: Logical addresses evolution and concurrent accesses table

	x[0]	x[1]	x[2]	x[3]	x[4]	concurrent accesses table
it o	0	1	2	3	4	(0,2) / (1,3)
it o+1	1	2	3	4	0	(1,3) / (2,4)
it o+2	2	3	4	0	1	(2,4) / (3,0)
it o+3	3	4	0	1	2	(3,0) / (4,1)
it o+4	4	0	1	2	3	(4,1) / (0,2)
it o+5	0	1	2	3	4	(0,2) / (1,3)
logical addresses						

The *set of concurrent accesses* (SCA) is the set of all the concurrent accesses in the concurrent accesses table. In our example,  $SCA = \{(0, 2), (1, 3), (2, 4), (3, 0), (4, 1)\}$ . A *Concurrent Accesses Graph* (CAG) is constructed from this set of concurrent accesses. A CAG is a couple  $CAG = (L, A)$ .  $L$  is the set of vertices  $L = \{l_0, l_1, \dots, l_{N-1}\}$  where vertex  $l_i$  is the logical address of the  $i^{th}$  element of the vector in the first iteration. With  $x$  a vector of size  $N$ ;  $card(L) = N$ . A

is the set of edges  $A = \{a_1, a_2, \dots, a_M\}$  where edge  $a = (l_i, l_j)$  links 2 elements  $l_i$  and  $l_j$  if the couple  $(l_i, l_j)$  is included in the set of concurrent accesses.  $A \subseteq L \times L$ . Figure 8(a) gives the conflict graph for our example.

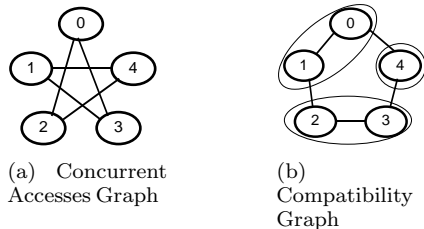


Figure 8: CAG and CG

In this case, the synthesis is not possible. GAUT issues a message to indicate that the data mapping is not valid. To help in determining a valid data mapping, a *Compatibility Graph* (CG) is constructed. The CG is orthogonal to the former Conflict Graph (Figure 8(b)). The minimum number of memory banks is easily computed from the Compatibility Graph. In our example, the minimum number of memory banks is 3. A possible mapping is to place  $x[0]$  and  $x[1]$  in a first bank,  $x[2]$  and  $x[3]$  in a second bank, and  $x[4]$  in a third bank. It is remarkable that these results actually depend on the scheduling, and therefore on the timing constraint provided to the tool. With a different timing constraint, the conflict and compatibility graphs change, as well as the set of valid data mappings.

Similar results are obtained when pipelined architectures are synthesized. The chronogram of accesses for algorithm 1 is presented on Figure 9. When the architecture is pipelined, this chronogram is modified as shown on Figure 10.

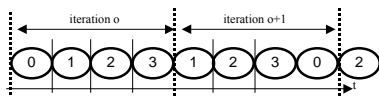


Figure 9: Non-pipelined architecture

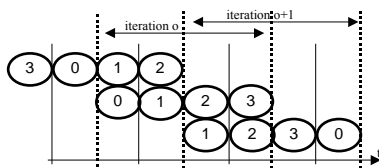


Figure 10: Pipelined architecture

The situation is similar to the situation with the algorithm 2: concurrent accesses appear and a concurrent accesses table is determined. The difference is that the conflicts arise between logical addresses that are calculated over several successive iterations (2 in this example).  $@x[2]^o$  is in concurrence with  $@x[0]^{o+1}$ , and  $x[3]^o$  is in concurrence with  $@x[1]^{o+1}$ . The set of concurrent accesses  $SCA = \{(0, 1), (1, 2), (2, 3), (3, 0)\}$ . The CAG and CG are computed from this SCA (Figure 11). The data mapping is verified. The minimum number of banks is 2, and the only valid mapping with 2 banks is to place  $x[0]$  and  $x[2]$  in a bank, and  $x[1]$  and  $x[3]$  in another bank.

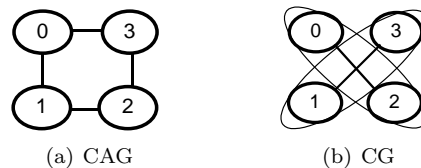


Figure 11: CAG and CG

## 4. Results

### 4.1. GAUT vs Monet

Several syntheses were performed, both with GAUT and the industrial behavioral synthesis tool Monet. We chose the elliptic and the Kalman filters which are the biggest applications in the HLSynth'92 benchmarks [10], and two classical digital algorithms: a FIR filter and an echo cancellation algorithm, the LMS. Table 2, indicates the synthesis time in seconds and the architecture's latency in number of cycles (the same real-time constraint was given to the tools, the clock cycle is 10ns). Required hardware resources are also indicated: the number of registers (Reg), of multiplexers (Mux), demultiplexers (Demux), of glue logic elements (which are tri-states in GAUT), and the number of RAM and ROM memories. The two last columns give the number of read and write in those memories. Single port SRAM were used to store data. Syntheses were executed on SUN Blade 2000 workstations.

Hardware resources are always lower in architectures synthesized with GAUT, although the same number of arithmetic operators is needed. The latency, which is the delay between the input of the first data and the first result on the output, is also lower with GAUT. A ROM is needed with GAUT for the FIR filter, since GAUT stores every static coefficient in ROM. Those coefficients are wired with Monet.

Table 2: GAUT vs Monet

		Synthtime	Lat (cycle)	Reg	Mux	Demux	Tri	Glue	RAM	ROM	Nb_R	Nb_W
elliptic	Monet	1s	20	19	16	15	–	27	–	–	–	–
	Gaut	1s	20	12	6	9	24	–	–	–	–	–
Kalman	Monet	1s	600	36	12	20	–	34	–	–	–	–
	Gaut	1s	60	14	11	10	29	–	–	–	–	–
FIR 16	Monet	2s	48	4	6	2	–	7	1	–	32	16
	Gaut	1.4s	19	4	2	1	1	–	1	1	32	1
LMS 32	Monet	6s	132	38	28	18	–	25	2	–	128	64
	Gaut	1.4s	100	19	3	3	23	–	2	–	128	33

Dynamic coefficients, whose value is changed during the execution of the algorithm, which is the case for an adaptative filtering like the LMS, are stored in RAM, together with signals (ageing vectors). The advantages of our approach appear clearly here: the latency is lower with GAUT since we avoid the  $n$  reads and writes of the ageing vector performed with Monet. As a result, the power consumption decreases. Indeed, the power consumption of a memory increases with the number of accesses.

The synthesis time, together with the reduction of hardware resources and memory accesses, exhibit the efficiency of our scheduling technique. In fact, the difference between the synthesis time with GAUT and with a behavioral synthesizer like Monet increases with the complexity of the application. We have measured the synthesis times for the FIR and the LMS filters, with an increasing complexity. Table 3 presents the results for the FIR for 64, 128, 256, 512, and 1024 points. Table 4 presents the results for the LMS filter for the same increasing complexities. It can be observed that, even if the difference between the synthesis time with GAUT and Monet is relatively small for small designs, it becomes enormous when the design’s complexity increases. Indeed, it becomes hours, then days or weeks for the FIR 1024 and the LMS 512 and 1024. In fact, every memory access is a node to be schedule in Monet, and the scheduling algorithm has a strong complexity. The difference in latency is comparatively stable: the latency with Monet varies from about 2 to 3 times the latency with GAUT.

## 4.2. Memory exploration

Once found the required number of operators, one can try different numbers of memory banks as well as several data mappings, and evaluate their impact on the final application’s performance and power consumption. In our example, for a data rate equals to 7.66s, we decide to allocate one multiplier and two ALU, and to share the memory in two banks.

Table 3: Synthesis of the FIR filter

FIR	Tool	cycles	Reads	Writes	Time
64	Monet	192	128	64	7s
	Gaut	67	128	1	1.7s
128	Monet	384	256	128	45s
	Gaut	131	256	1	2s
256	Monet	768	512	256	20mn
	Gaut	259	512	1	3.8s
512	Monet	1536	1024	512	7h17mn
	Gaut	512	1024	1	4.9s
1024	Monet	3072	2048	1024	... days
	Gaut	1027	2048	1	9s

Table 4: Synthesis of the LMS filter

LMS	Tool	cycles	Reads	Writes	Time
64	Monet	260	256	128	14s
	Gaut	196	256	65	1.9s
128	Monet	516	512	256	7mn30s
	Gaut	388	512	129	2.6s
256	Monet	1028	1024	512	3h30mn
	Gaut	772	1024	257	5.3
512	Monet	2052	2048	1027	... days
	Gaut	1540	2048	513	9.6
1024	Monet	4010	4096	2048	... weeks
	Gaut	3076	4096	1025	64

We then apply several memory mappings to the synthesis process, and observe the impact on the resulting circuit’s power consumption. The circuit is a FPGA Xilinx Virtex XC400. Its consumption is computed with the Xilinx tool: XPower. One memory unit is generated for each memory mapping. The memory unit power consumption Pmu is provided in Table 5, together with the overall consumption Ptot, and the processing unit consumption Ppu. In the map2.16 mapping, the first sixteen real and imaginary FFT samples are mapped in the first bank, the remaining sixteen samples are in the second bank. In the map2.8, map2.4, and map2.2 mappings, samples are mapped respectively 8 by 8, 4 by 4, and 2 by 2, in the first and second memory banks. In the map2.1



mapping, even samples are in the first bank, odd samples in the second.

Table 5: Power consumption

Power (mW)	P <sub>mu</sub>	P <sub>pu</sub>	P <sub>tot</sub>	$\Delta$ P <sub>tot</sub> %
map2_16 (2banks)	12.71	34.35	47.06	...
map2_8 (2banks)	12.83	21.85	34.68	-26
map2_4 (2banks)	12.9	17.77	30.67	-35
map2_2 (2banks)	12.99	21.63	34.62	-26
map2_1 (2banks)	12.87	19.26	32.13	-31

Every memory unit invariably contains 2 RAM, and 4 FSM to drive the write and read accesses between the busses and the RAM. As a result, and because the number of memory accesses is also constant, there are very few variations on the memory unit power consumption (less than 2.5%, only due to small changes in I/O signals commutations and some logic blocs). Variations of the processing unit consumption are much more important, for they represent from 58% to 64.1% of the overall power consumption. The lowest overall consumption is obtained with the map2.4 mapping (35% lower than the map2\_16), even if the memory unit consumption is slightly higher in this case (1.5%) than the lowest one.

## 5. Conclusion

In this paper, we present two recent improvements to our High-Level Synthesis tool GAUT. Our goal is to take into account the memory architecture and the memory mapping in the synthesis process. We first define the memory mapping constraint and include it in the synthesis design flow. We introduce the Memory Constraint Graph, and an accessibility criterion to enhance the scheduling algorithm. We show that a peculiar attention must be paid to signals, or ageing vectors. We formalize the maturing process and explain how it may generate memory conflicts over several iterations of the algorithm. We define the Logical Accesses Graph, and the Concurrent Accesses Table, which are used to construct the Concurrent Accesses Graph, and the Compatibility Graph. The Compatibility Graph indicates the minimum number of memory banks for the scheduling, and helps in finding a valid mapping for signals.

Several experiments were made, to explore the efficiency of our approach. The comparison with an industrial behavioral synthesis tool exhibits several advantages for GAUT. It appears firstly that GAUT uses less hardware resources, and reduces the count of memory accesses, which lead to a lower latency and a lower power consumption. Secondly, GAUT is able to

tackle complex designs, and to perform the synthesis in a reasonable time. Memory aware synthesis and GAUT appear very efficient for exploring the design space and for balancing optimizations between the processing unit and the memory unit. It permits to determine the best memory architecture, i.e. the best number of memory banks, as well as the best memory mapping, to meet the application constraints, and to finally reach a reasonable tradeoff between time, power consumption, and area.

In the future, the scheduling step will be enhanced with an anticipated read model for the data, which should allow to speedup the processing unit. The presented strategy for implementing ageing vectors will be reversed, in order to automatize the determination of the memory mapping for this type of data.

## References

- [1] ITRS homepage. [Online]. Available: <http://public.itrs.net/>
- [2] C. Gebotys, "Low energy memory and register allocation using network flow," in *Proc. Design Automation Conference DAC'97*, June 1997, pp. 435–440.
- [3] N. Passos, E. Sha, and L.-F. Chao, "Multi-dimensional interleaving for time-and-memory design optimization," in *Proc. IEEE Int. Conf. On Computer Design ICCD'95*, Oct. 1995, pp. 440–445.
- [4] P. Ellervee, "High-level synthesis of control and memory intensive applications," Ph.D. dissertation, Royal Institut of Technology, Jan. 2000.
- [5] J. Seo, T. Kim, and P. Panda, "An integrated algorithm for memory allocation and assignment in high-level synthesis," in *Proc. Design Automation Conference DAC'01*, June 2001, pp. 608–611.
- [6] H. Ly, D. Knapp, R. Miller, and D. McMillen, "Scheduling using behavioral templates," in *Proc. Design Automation Conference DAC'95*, June 1995, pp. 101–106.
- [7] D. Knapp, T. Lyand, *et al.*, "Behavioral synthesis methodology for HDL-based specification and validation," in *Proc. Design Automation Conference DAC'95*, June 1995.
- [8] P. Panda *et al.*, "Data and memory optimization techniques for embedded systems," *Transactions on Design Automation of Electronic Systems*, vol. 6, no. 2, pp. 149–206, 2001.
- [9] F. Catthoor, K. Danckaert, C. Kulkarni, and T. Omns, *Data Transfer and Storage (DTS) architecture issues and exploration in multimedia processors*. NewYork: Marcel Dekker Inc., 2000.
- [10] HLSynth'92 benchmark information. [Online]. Available: [http://www.cbl.ncsu.edu/CBL\\_Docs/hls92.html](http://www.cbl.ncsu.edu/CBL_Docs/hls92.html)