

Systolic solution of linear systems over GF(p) with partial pivoting

Bertrand HOCHET⁺, Patrice QUINTON⁺⁺ and Yves ROBERT⁺⁺⁺

⁺ CNRS, Laboratoire TIM3, BP 68, 38402 St Martin d'Hères Cedex, France
⁺⁺ CNRS, IRISA, 35042 Rennes Cedex, France
⁺⁺⁺ CNRS, Laboratoire TIM3 and IBM ECSEC, Via Giorgione 159, 00147 Roma, Italy

Abstract : We propose two systolic architectures for the Gaussian triangularization and the Gauss-Jordan diagonalization of large dense $n \times n$ matrices over GF(p), where p is a prime number. The solution of large dense linear systems over GF(p) is the major computational step in various algorithms issued from arithmetic number theory and computer algebra. The two proposed architectures implement the elimination with partial pivoting, although the operation of the array remains purely systolic. The last section is devoted to the design and layout of a CMOS 8 by 8 Gauss-Jordan diagonalization systolic chip over GF(2).

1. Introduction

The solution of large dense linear systems of algebraic equations in finite fields appears to be the computational kernel of many important algorithms. Let us mention three representative applications:

- the large integer factoring routines, where the final stage necessitates performing Gaussian elimination on a large dense matrix over GF(2) ([12] [13] [15]).
- the factorization of polynomials over GF(p), where p is a prime number, via Berlekamp's algorithm ([10]).
- linear prediction in the analysis of discrete signals ([11]).

In some of these applications, matrices of several thousands of rows and columns are required, and there are a large number of matrices to triangularize. Parallel processing therefore offers an interesting perspective for such computations. Indeed, Parkinson and Wunderlich report in [13] the implementation of Gaussian elimination over GF(2) on the ICL DAP (which has 4096 processors), and Poet [15] considers the use of special-purpose hardware to factor 140 digit numbers using the same algorithm.

In this paper, we introduce systolic architectures for the Gauss triangularization algorithm and the Gauss-Jordan diagonalization algorithm of general dense matrices over a finite field GF(p), where p is prime. Systolic arrays for Gaussian triangularization of dense real matrices are well known (see [1], [6] among others). However, partial pivoting has never been implemented in such arrays, due to the fact that the search for a pivot in a whole row or column would break down the locality and the regularity of the systolic design. Hence systolic arrays for Gaussian elimination over \mathbb{R} only apply to symmetric positive definite or diagonally dominant matrices. To triangularize general matrices, one must use an orthogonal

factorization via Givens rotations (see the systolic implementations of [1], [4], [6], [17] among others).

In a finite field GF(p), partial pivoting can not be avoided, since in average every other p element is zero. But it turns out that it can be solved much more efficiently than in \mathbb{R} , since any non-zero element in a given row or column can be chosen as the pivot. This simple consideration will be very important in our implementations.

The paper is organized as follows. First we concentrate on GF(2), and we defer the implementation over GF(p) for any prime p up to section 4. Section 2 is devoted to the design of a systolic array for Gaussian elimination over GF(2) of an n by $(n+1)$ matrix (A, b) .

When there are q linear systems $Ax_i = b_i$ to be solved, the previous array requires $n^2/2 + qn$ cells and $3n+q$ time-steps. There still remains q triangular systems to solve. Another strategy is to directly implement the Gauss-Jordan diagonalization algorithm with partial pivoting. We show in section 3 that the solution matrix $A^{-1}B$ (where $B = (b_1, b_2, \dots, b_q)$) can be computed within $4n+q$ time-steps on a systolic array of n^2 cells. In particular, the inverse of A can be computed within $5n$ time-steps.

We show in section 4 how to modify the operation of the elementary processors to deal with an implementation over GF(p), p prime number. Finally, we detail in section 5 the layout of a CMOS 8 by 8 Gauss-Jordan diagonalization systolic chip over GF(2) which has been designed using the system LUCIE [14] and which has been fabricated through the French Multi-Project Chip [2].

2. Gaussian elimination over GF(2)

Let A be a dense $n \times n$ matrix and b a n -vector. To solve the system

$$(1) Ax = b$$

we transform it into an equivalent triangular system

$$(2) Tx = b'$$

The transformation of the system (1) into the system (2) is done by triangularizing the matrix A , using Gaussian elimination with partial pivoting.

We briefly recall the operation of Gentleman and Kung's two-dimensional triangular array of orthogonally connected processors to triangularize a real dense $n \times (n+1)$ matrix (A, b) without pivoting. This array is depicted in figure 1. The total

This work has been supported by the Coordinated Research Program C3 of CNRS and Ministère de la Recherche et de la Technologie

number of cells in the array is $n[(n+1)/2+1]$. The array is composed of n rows, each row k including $n+2-k$ processors numbered from left to right $P_{k,1}, \dots, P_{k,n+2-k}$. The matrix (A,b) is fed into the array column by column. More specifically, column j of the matrix is input to processor $P_{1,j}$, one new element each time-step, beginning at time $t=j$. This input format is depicted in figure 1.

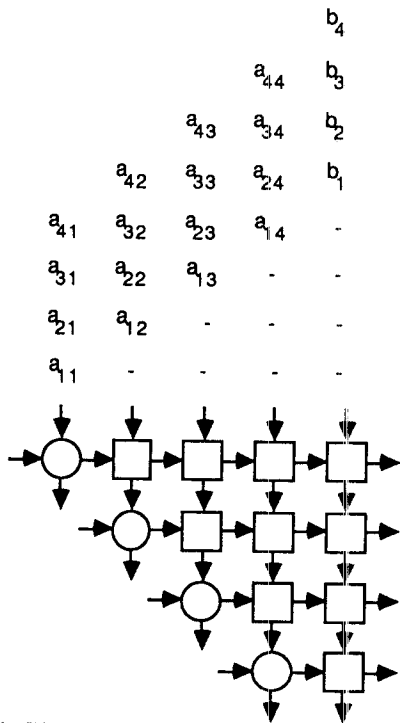


Figure 1 : The systolic array for Gaussian elimination ($n=4$)

The first row of (A,b) is stored in the upper row of the array, and, as any row numbered $i \geq 1$ is read by the array, it is combined with the first one in order to zero out element a_{i1} ; this combination corresponds to a transformation of the form

$$\begin{bmatrix} a_{11} & \dots & a_{1n} & b_1 \\ a_{i1} & \dots & a_{in} & b_i \end{bmatrix} := M_{i1} \begin{bmatrix} a_{11} & \dots & a_{1n} & b_1 \\ a_{i1} & \dots & a_{in} & b_i \end{bmatrix}, \quad M_{i1} = \begin{bmatrix} 1 & 0 \\ -a_{i1}/a_{11} & 1 \end{bmatrix}$$

The 2×2 matrix M_{i1} is computed by P_{11} at time $t = i$, and sent to the other cells of the first row; more precisely, cell $P_{1,k}$ ($k \geq 2$) performs the transformation

$$(a_{1k}, a_{ik})^t := M_{i1} \cdot (a_{1k}, a_{ik})^t$$

at time $t = i+k-1$.

Let

$$(A^{(k)}, b^{(k)}) = \begin{bmatrix} a_{11}^{(k)} & \dots & a_{1n}^{(k)} & t_1^{(k)} \\ a_{kk}^{(k)} & \dots & a_{kn}^{(k)} & t_k^{(k)} \\ \mathbf{0} & \dots & \dots & \dots \\ a_{nk}^{(k)} & \dots & a_{nn}^{(k)} & t_n^{(k)} \end{bmatrix}$$

denote the matrix obtained from (A,b) after the elimination of

the elements at positions (i,j) such that $i > j, j = 1, \dots, k-1$. Then $(a_{kk}^{(k)}, \dots, a_{kn}^{(k)}, b_k^{(k)})$ is stored in the k -th row of the array. When $(a_{ik}^{(k)}, \dots, a_{in}^{(k)}, b_i^{(k)})$, $i > k$, is read by this row of cells, it is combined with $(a_{kk}^{(k)}, \dots, a_{kn}^{(k)}, b_k^{(k)})$ using a 2×2 matrix M_{ik} , in order to set $a_{ik}^{(k)}$ to zero.

2.1. General description of the array

The key-idea to include partial pivoting in the algorithm is to generate matrices M_{ik} whose structure depends on the values of $a_{kk}^{(k)}$ and $a_{ik}^{(k)}$. For sake of simplicity, let us consider the case $k=1, i=2$, that is, when the first matrix M_{21} is generated by processor P_{11} . Element $a_{11} = a_{11}^{(1)}$ is stored in P_{11} , and $a_{12} = a_{12}^{(1)}$ is being input to P_{11} . Two cases must be considered:

- (a) if $a_{11} = 1$, row 1 will be chosen as the pivoting row, and the elimination of $a_{21}, a_{31}, \dots, a_{n1}$ will be done using usual Gaussian elimination. In our example:

- (a1) if $a_{21} = 0$, there is nothing to do, M_{21} is equal to I_2 , the identity matrix of order 2.

- (a2) if $a_{21} = 1$, add row 1 to row 2 to zero out a_{21} , that is

$$M_{21} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

Note that addition denotes here the addition in $GF(2)$, i.e. the exclusive or operation XOR.

- (b) if $a_{11} = 0$, we have to find out a pivoting row. Two cases occur:

- (b1) if $a_{21} = 0$, we do nothing, and choose $M_{21} = I_2$

- (b2) if $a_{21} = 1$, we exchange rows 1 and 2, and M_{21} is the permutation matrix

$$M_{21} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$

Row 2 is then stored in the first row of the array and acts as the pivoting row for phase 1, which consists in zeroing out all the elements but one of the first column of the matrix (A,b) .

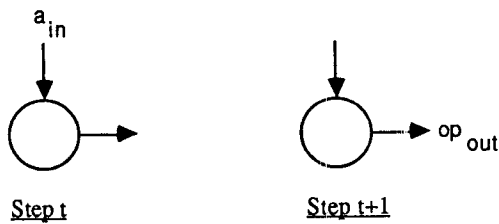
Note that in the case (b1), row 2 will not be modified during the first phase of the algorithm. Assume that a_{i1} , $i > 2$, is the first non-zero element that we find out: row i will be used as a pivoting row for phase 1, but since we already know that $a_{11}, a_{21}, \dots, a_{i-1,1}$ are all zero, we do not have to combine row i with row 1, 2, ..., $i-1$.

We can give an informal description of the algorithm as follows:

for $k := 1$ to $n-1$ do
{phase k of the algorithm }
begin

1. find out the first $j \geq k$ such that $a_{jk}^{(k)} = 1$
2. exchange row k and row j (if $j \neq k$), that is, choose row j as pivoting row

3. zero out the elements $a_{ik}^{(k)}$ such that $i > j$ and $a_{ik}^{(k)} = 1$ by adding row j to row i
 (here addition denotes the addition in GF(2), that is, the XOR operation)
 end

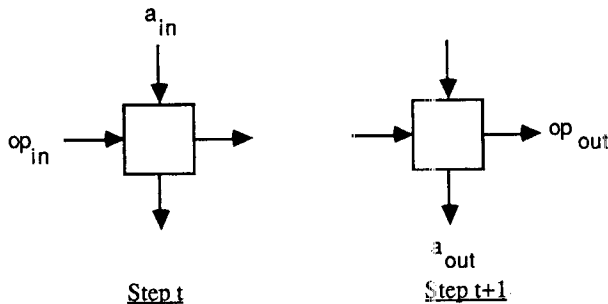


```

if init then {store a_in}
begin r := a_in ; init := false end
else {generate instruction }
case (a_in , r) of
begin
(0,0) : op_out := id {identity - still looking for pivot};
(0,1) : op_out := id {element already zero};
(1,0) : op_out := perm {exchange rows};
(1,1) : op_out := add {XOR operation}
end;

```

a) Operation of the circular cells



```

if init then {store a_in}
begin r := a_in ; init := false end
else {update a_in (and r if a permutation is required)}
begin
case op_in of
begin
id : a_out := a_in {r is not modified};
add : a_out := a_in XOR r {r is not modified};
perm : a_out := r ; r := a_in {exchange a_in and r}
end;
op_out := op_in
end

```

b) Operation of the square cells

Figure 2 : Gaussian elimination over GF(2) - Operation of the processors

For a given k , note that the steps 2 and 3 are not executed if we can not find a j such that $j \geq k$ and $a_{jk} = 1$: in this case, the matrix A is rank deficient.

2.2. Operation of the processors

There are two types of processors in the array, respectively represented as circular and square processors in the figure 1. In the k -th row of the array:
 - the processor P_{k1} is a circular processor, which generates the matrices M_{ik} , $i > k$.
 - all the other processors are square processors, which apply the transformation relative to M_{ik} .

The operation of these two types of processors is detailed in the figure 2. The operation of a given processor depends on whether it is the first time it operates. In the description of the figure 2, we assume for simplicity that there is a boolean called init (for "initialization") stored in every processor which is set to "true" at the beginning of the computation. This boolean controls the operation of the processors. In an actual implementation, the instruction to start the computation would be input to processor P_{11} together with the first coefficient a_{11} and systolically propagated through the whole array (it can be easily checked that processor P_{kj} operates for the first time at step $3k+j-3$).

Circular processors

The first time they operate, circular processors store their input in their internal register. Afterwards, they compute the matrices M_{ik} and send them rightwards to the square processors.

Rather than directly generating the matrices M_{ik} , the circular processors transmit rightwards one of the following three instructions:

- *id*, which stands for the identity matrix
- *perm*, which requires a permutation of rows
- *add*, which requires an addition of rows.

This allows the instructions generated by the circular processors to be encoded using only two bits.

Square processors

Again, square processors store their input in their internal register the first time they operate. Afterwards, they update their vertical input and possibly their internal register, according to the instruction they receive from the left. See figure 2 for a complete description.

2.3. Unloading the array

At the end, the matrix (T,b') is stored in the array as follows (assuming $n=4$):

t_{11}	t_{12}	t_{13}	t_{14}	b'_1
	t_{22}	t_{23}	t_{24}	b'_2
		t_{33}	t_{34}	b'_3
			t_{44}	b'_4

There remains to unload the array. According to the general philosophy of the model, this must be done systolically, by propagating special boolean control instructions. Moreover, we would like to pipeline the unloading of the array with the computation phase, so that n additional steps are not necessary. Various schemes are possible. We outline a scheme where each processor sends the content of its register to the right when it has finished to operate. Let $t_{i,n+1} = b'_i$ for convenience.

At time $t=4$ in our example, P_{11} operates for the last time. It remains idle for one step and then sends t_{11} rightwards. P_{12} finishes to work at time $t=5$. It transmits t_{11} to the right at $t=6$, and then sends t_{12} at time $t=7$. The process goes on, and t_{1i} is output by P_{15} at time $t=9+i$, for $1 \leq i \leq n-1$. Similarly, we start unloading the second row at time $t=8$. We see that the first rows of the array are unloaded, while the last rows still operate, thus achieving the pipelining that is sought.

Indeed, in the general case, the last computation occurs in $P_{n,n+1}$ at time $3n-1$, and $t_{i,i+k}$ ($1 \leq i \leq n$, $0 \leq k \leq n+1-i$) is output from the array at time $2n+i+k+1$. The total computation time is $3n+2$.

We point out that the unloading of the array can not be simplified as in the case of Gaussian elimination over \mathbb{R} without pivoting: in the real case, the coefficients $t_{i,i+k}$ are not modified after being stored in the array, hence they can be sent downwards immediately. Here on the contrary, the $t_{i,i+k}$ can be modified by any following permutation of rows.

2.5. General remarks

We have described here an array for matrix triangularization using Gaussian elimination. There still remains a triangular system to be solved. Assuming the matrix A is non-singular, we can use another systolic array, the triangular system solver of Kung and Leiserson [9], which requires $2n$ additional steps to solve the system $Tx = b'$. However, this would require the triangular matrix to be stored in the host and reordered by diagonals before being fed in Kung and Leiserson's array. Rather, it is more efficient to use the Jordan elimination scheme depicted in [5] [16] which can be implemented on the same array and requires only n additional steps to provide the solution x , leading to a very efficient scheme of only $n^2/2$ cells and $4n$ time steps to completely solve the system $Ax = b$.

We point out that checking for the non-singularity of the matrix can be done on the fly very easily, simply by testing the content of the registers of the circular processors at the end of the triangularization phase.

3. The Gauss-Jordan diagonalization algorithm with partial pivoting over GF(2)

The systolic array presented in section 2 allows general dense systems of equations $Ax = b$ to be solved: first triangularize the matrix (A,b) , then solve the triangular system.

Assume now that there are q systems $Ax_i = b_i$ to solve, $1 \leq i \leq q$, and let B be the nxq matrix $B=(b_1, b_2, \dots, b_q)$. To triangularize (A,B) , we can simply extend the triangularization array by adding $q-1$ columns on the right, leading to a total number of $n[(n+1)/2+q]$ cells. If $PA = LT$, where P is a permutation matrix, L is lower triangular with unit diagonal, and T is upper triangular, we obtain the matrix $(T, L^{-1}PB)$ within $3n+q$ steps. There remains q triangular systems to be solved. This can be done within $2n+q$ steps using an array of qn cells (simply concatenate q triangular systems solvers of [9]). Hence we need a total number of $n^2/2 + 2qn + o(n)$ cells and of $5n+2q$ time-steps to solve the q linear systems. However, this evaluation does not take into account the fact that the matrix

$(T, L^{-1}PB)$ has to be stored in the host and reordered by diagonals before entering the second systolic array.

We concentrate in this section on the design of a systolic array which directly computes the product $A^{-1}B$, where A is a dense (nonsingular) nxn matrix and B is a dense nxq matrix. The inverse of A is computed via the Gauss-Jordan diagonalization algorithm with partial pivoting. The performance of the new array is much better for large q : only n^2 cells and $4n+q$ time-steps are required by the architecture that we are going to describe. First we briefly recall the well-known Gauss-Jordan diagonalization algorithm.

3.1. The Gauss-Jordan diagonalization algorithm

Let C denote the n by $n+q$ matrix $C = (A,B)$. We want to reduce C to the matrix $(I, A^{-1}B)$. When no pivoting is needed, the usual way to proceed is to pre-multiply C by n elementary nxn matrices J_1, J_2, \dots, J_n in order to obtain after n steps

$$J_n \dots J_2 J_1 C = (I, A^{-1}B)$$

Therefore, define $C_k = J_k \cdot C_{k-1}$, starting with $C_0 = C$. We choose J_k so that the first k columns of C_k are those of I_n , the identity matrix of order n . Thus C_k has the following structure:

$$C_k = \left[\begin{array}{c|c} \text{the first } k \text{ columns of } I_n & C_k^\circ \end{array} \right]$$

where C_k° denotes the $nx(n+q-k)$ matrix built up with the last $n+q-k$ columns of C_k .

In particular, C_n° is the desired matrix $A^{-1}B$. The matrix J_k only differs from I_n by its k -th column, which we denote for convenience as

$$[c_{1k}^{(k)} \dots c_{nk}^{(k)}]^t$$

The coefficients of C_k° are denoted $(c_{ij}^{(k)})$, $1 \leq i \leq n$, $k+1 \leq j \leq n+q$. Therefore $c_{ik}^{(k)}$ refers to the k -th column of J_k , while $c_{ij}^{(k)}$ with $k < j$ refers to the matrix C_k° . The values of the $c_{ij}^{(k)}$, $1 \leq i \leq n$, $k \leq j \leq n+q$, $1 \leq k \leq n$, are computed recursively by the following algorithm, starting from $c_{ij}^{(0)} = c_{ij}$:

```
{ Gauss-Jordan diagonalization algorithm }
for k := 1 to n
begin
  { compute  $J_k$  }
   $c_{kk}^{(k)} := 1 / c_{kk}^{(k-1)}$ 
  for i := 1 to n, i  $\neq$  k
     $c_{ik}^{(k)} := -c_{ik}^{(k-1)} * c_{kk}^{(k)}$ 
  { compute  $C_k^\circ$  }
  for j := k+1 to n+q
  begin
    for i := 1 to n, i  $\neq$  k
       $c_{ij}^{(k)} := c_{ij}^{(k-1)} + c_{ik}^{(k)} * c_{kj}^{(k-1)}$ ;
     $c_{kj}^{(k)} := c_{kk}^{(k)} * c_{kj}^{(k-1)}$ ;
  end;
end;
```

The matrix J_k can be viewed as the (commutative) product of n elementary matrices. The first $n-1$ ones apply the transformations M_{jk} of section 2, and are chosen so as to

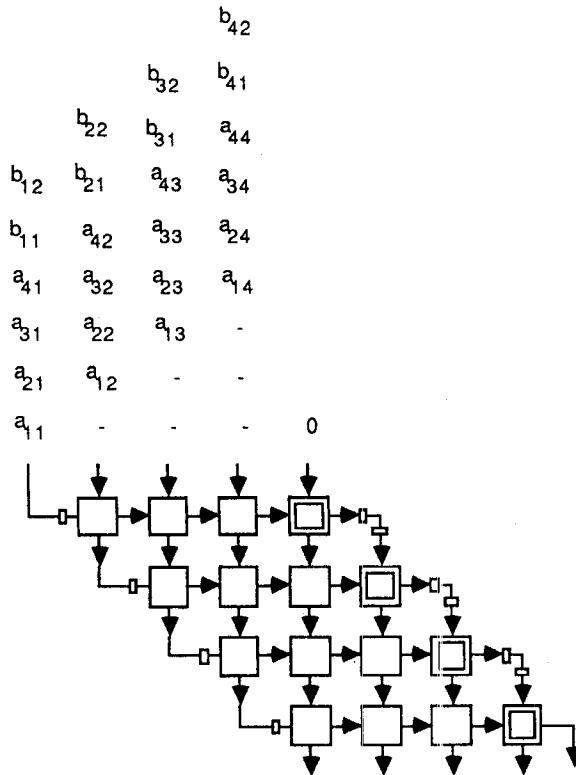


Figure 3 : The systolic array for computing $A^{-1}B$ ($n=4, q=2$)

annihilate the elements in position (i,k) , $i \neq k$. The n -th matrix only differs from I_n by its element in position (k,k) chosen so as to divide row k by the pivot $c_{kk}^{(k-1)}$.

Contrarily to the case of Gaussian elimination, every row of the matrix has to be combined with the pivoting row at each step. Provided that the organization of the array allows to do so, partial pivoting can be introduced just as before. During step k , column k is scanned out until a non-zero element in position (i,k) , say, is found. M_{ik} is still a permutation matrix, and row i is used as pivoting row and combined with all the rows of the matrix whose k -th element is non-zero.

We can give an informal description of the algorithm as follows:

```

for k := 1 to n do
{phase k of the algorithm }
begin
1. find out the first  $j \in \{k, k+1, \dots, n, 1, 2, \dots, k-1\}$ 
such that  $a_{jk}^{(k)} = 1$ 
2. if  $j \neq k$ , exchange row  $k$  and row  $j$ , that is, choose
row  $j$  as pivoting row
3. zero out the elements  $a_{ik}^{(k)}$  such that  $i \neq j$  and
 $a_{ik}^{(k)} = 1$  by adding row  $j$  to row  $i$ 
(here addition denotes the addition in  $GF(2)$ , that
is, the XOR operation)
end

```

For a given k , note that statements 2. and 3. are not executed if we can not find a non-zero element in column k , that

is, statement 1. does not succeed. In this case, the matrix A is singular. This leads to the systolic implementation which is now presented.

3.2. Systolic implementation

As we already mentioned, every row of the matrix must be combined with the pivoting row at each phase of the algorithm. That is the reason why we choose an implementation which is "dual" of that of the previous section. Rather than storing a coefficient of the matrix in the internal register of each cell, we store (an encoding of) the transformation matrix M_{ik} . Moreover, the matrix is input in a transposed fashion, so that each row can meet all the other rows while moving through the array.

Figure 3 depicts the Gauss-Jordan diagonalization array. It is a two-dimensional array of orthogonally connected processors with n rows. Each row k comprises n processors numbered from left to right $P_{k,1}, \dots, P_{k,n}$. There is a one-step delay cell at the left of each row, named P_{k0} for convenience. There are also two delay cells at the right of each row, as shown in figure 3. The operation of each processor is detailed in figure 4. There are two types of processors, represented as square cells (type 1) and double square cells (type 2). In the k -th row of the array, the rightmost processor $P_{k,n}$ is of type 2. All the other processors $P_{k,1}, \dots, P_{k,n-1}$ are of type 1. The cells operate as follows:

Square cells

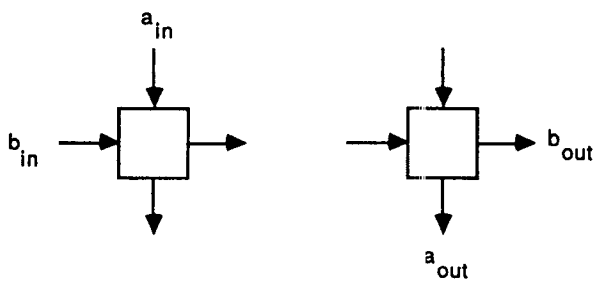
Square cells first initialize their internal register by storing the operation to be performed on all the subsequent coefficients of the two rows. Just as in the previous section, the operation is to be chosen in the set $\{id, perm, add\}$ according to the value of the first inputs. Once initialized, the cells update their inputs as indicated in figure 4.

Double square cells

Double-square cells first test for the non-singularity of the matrix: if the first input of cell $P_{k,n}$ is zero, then all the elements of column k after the first $k-1$ phases of the algorithm are zero, and the matrix is singular. To test this condition, a boolean variable, called *test*, is initialized to false and moves through all the double square cells. When it is output from $P_{n,n}$ at the end of the computation, its value is true if and only if the algorithm has failed, that is, if the matrix is singular. Once initialized, double square cells simply transmit their input (see figure 4).

The operation of a given processor in the array depends on whether it is the first data item it receives. As shown in figure 4, each processor has a control boolean *init* initialized to true. This boolean specifies the operation to be performed and the line along which the data is to be sent out. In an actual implementation, the instruction to start the computation would be input to the top-left corner of the array and systolically propagated to all the processors. It can be easily checked that processor P_{kj} operates for the first time at step $3k+j-2$, so that the start signal would move at full speed to the right and at one-third speed downwards.

The matrix A , followed by the matrix B , is fed into the array row by row. More specifically, row k of the matrix $C = (A,B)$ is input to processor $P_{1,k-1}$, one new element each time-step, beginning at time $t=k$. This input scheme is depicted in figure 3.



Step t

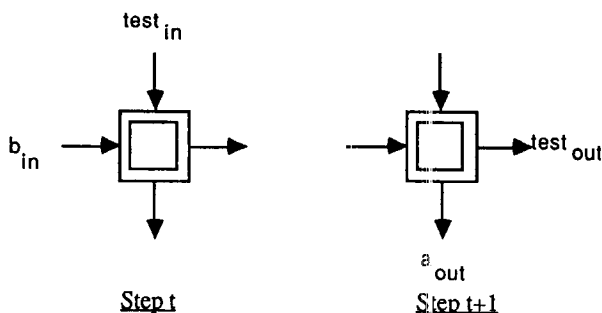
Step t+1

```

if init then {generate instruction and store it in internal register}
begin
  init := false;
  case (ain, bin) of
  begin
    (0,0) : op := id ; bout := bin {identity - still looking
              for pivot};
    (0,1) : op := id ; bout := bin {element already zero};
    (1,0) : op := perm ; bout := ain {exchange rows};
    (1,1) : op := add ; bout := bin {XOR operation}
  end
else {update ain (and bin if a permutation is required )}
  case op of
  begin
    id : aout := ain ; bout := bin ;
    add : aout := ain XOR bin ; bout := bin ;
    perm : aout := bin ; bout := ain ;
  end;

```

a) Operation of the square cells



Step t

Step t+1

```

if init then
{test for non singularity : bin=0 means that column is 0}
begin
  init := false ;
  testout := testin OR not (bin) { testout = 1 means that the
    matrix is already known to be singular}
end
else {transmit bin}
  aout := bin

```

b) Operation of the double square cells

Figure 4 : Gauss-Jordan diagonalization. Operation of the processors

The k-th row of the array is devoted to the computation

$$C_k^o = J_k C_{k-1}^o$$

Recall that J_k is the product of n elementary matrices, some of which possibly being permutation matrices. The leftmost delay cell P_{k0} transmits the input data arriving from the top to the right. Square cells of the row use the first data value arriving from the top to compute the instruction which they store in their internal register. Then they pass downwards all the following data after modification. Similarly the rightmost processor $P_{k,n}$ modifies and stores the first input data arriving from the left and passes downwards all the following data after modification. Thus after a row of $n+q$ input data flows through the whole array, its length is shortened by n to become a row of q output data. The n by $n+q$ matrix $C_0^o = C_0^o = C$ is input to the first row of the array, the n by $n+q-1$ matrix C_1^o is input to the second, ..., and the n by $q+1$ matrix C_{n-1}^o is input to the n -th row; finally the array outputs the n by q matrix $C_n^o = A^{-1}B$.

It is important to note that the rows of the matrix are "reordered" in some sense when moving through the array: the input of the processors $P_{k0}, P_{k1}, \dots, P_{k,n}$ in the k -th row of the array are respectively the rows $k, k+1, \dots, n, 1, \dots, k-1$ of C_{k-1}^o . The matrices $M_{ik}, i \neq k$, are computed and stored in the processors $P_{k,1}, \dots, P_{k,n-1}$. More precisely, M_{ik} is stored in $P_{k,(i-k) \bmod n}$. Note that M_{kk} is always the identity in $GF(2)$: either the matrix is singular, or there is already a 1 on the diagonal. After the processors are initialized, they perform the multiplication $C_k^o = J_k C_{k-1}^o$. The matrix C_k^o is output from the k -th row of the array in row order, the leftmost row of the matrix being now row $k+1$.

We can state the following: given a dense nonsingular $n \times n$ matrix A and a dense $n \times q$ matrix B , the orthogonal systolic array of n^2 processors can compute $A^{-1}B$ within $4n+q-2$ time-steps. In particular, we can compute the inverse A^{-1} of a dense $n \times n$ matrix within $5n-2$ steps.

3.3. A detailed example ($n=4, q=3$)

Let $n=4, q=3$, and A, B be the following matrices over $GF(2)$:

$$A = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix} \quad B = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

We want to describe the operation of the first row of the array, which is responsible for the first phase of the algorithm. For sake of clarity, assume that the columns of the matrix $C = (A, B)$ are input sequentially (rather than pipelined) to the array.

First $a_{11}=0$ is transmitted to the right by the delay cell P_{10} . When P_{11} receives $a_{21}=1$ and $a_{11}=0$, it generates the instruction *perm*, so that a_{21} replaces a_{11} and moves rightwards: row 2 will be chosen as the pivoting row for phase 1 of the algorithm. When P_{12} receives $a_{21}=1$ and $a_{31}=0$, it generates the instruction *id*. When P_{13} receives $a_{21}=1$ and $a_{41}=1$, it generates the instruction *add*. When $a_{21}=1$ is input to processor P_{14} , the control boolean *test* remains equal to false.

Once initialized, the processors P_{11} , P_{12} and P_{13} perform the instruction stored in their internal registers.

The operation of the first row of processors can be summarized as follows:

	P_{11}	P_{12}	P_{13}	P_{14}
Internal registers	perm	id	add	-
	↓	↓	↓	↓
Outputs	-	-	0	1
	-	1	1	1
	1	1	0	1
	0	1	1	1
	1	0	0	0
	0	0	1	-
	1	1	-	-
	0	-	-	-

In other words, the first row of the array outputs the matrix :

$$C_1^\circ = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

Therefore, during phase 2, row 3 is chosen as pivoting row; P_{21} , P_{22} and P_{23} generate respectively the instructions *perm*, *add* and *id*, so that :

$$C_2^\circ = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 \end{bmatrix}$$

During phase 3, row 3 is the pivoting row; P_{31} , P_{32} and P_{33} generate respectively the instructions *id*, *add*, and *id*, so that :

$$C_3^\circ = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 \end{bmatrix}$$

Finally, during phase 4, row 4 is the pivoting row; P_{41} , P_{42} and P_{43} generate respectively the instructions *add*, *id*, and *id*, so that at the end of the computation, the matrix $A^{-1}B$ is output from the array :

$$C_4^\circ = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

Since the first input to all the processors P_{14} , P_{24} , P_{34} and P_{44} was a 1, the final value of the boolean test is false. This guarantees that the algorithm has not failed, that is, that the matrix A is non-singular. The first column of $A^{-1}B$ is $x = (1, 1, 1, 1)^t$, the solution vector of section 2.3.

4. Extension to $GF(p)$, p prime number

The two previous arrays can be easily extended to deal with computations over a finite field $GF(p)$, where p is any prime number. The key idea is the same: the first non-zero element found out during the search is chosen as pivot.

The only modification in the operation of the processors is to replace the *add* instruction by the appropriate combination factor. Hence the processors now generate an instruction

$$op \in \{ id, perm, comb \}$$

Whenever $op=comb$, an appropriate combination factor is also generated, and transmitted or stored according to the program of the cell.

The arithmetic requirements are more complex in $GF(p)$ for general p than in $GF(2)$. A possible way to implement the inverse computation is a table look-up, whereas the multiplication can be done using the MMA Modular Multiplication Algorithm of [3], which consists of reducing each partial summation modulo p before going on: this prevents to perform a division by p at the end. Finally, the cost of an addition is twice that in \mathbb{R} (for the same number of bits) because $a+b \bmod p$ requires the computation of $a+b$ and $a+b-p$.

5. A 8 x 8 implementation over $GF(2)$

We briefly describe the design and layout of a systolic chip which implements the Gauss-Jordan diagonalization algorithm over $GF(2)$. The chip has been designed using the system LUCIE [14]. The technology is CMOS double metallisation with a 2 μm grid. The chip is currently being fabricated via the French Multi-Project Chip [2]. We have designed a 8 by 8 array, but of course the modularity of systolic arrays would permit to design larger arrays without any difficulty.

The only modification to the description of section 3.2. is the initialization of the array. As we already mentioned, there is a control boolean moving systolically through the whole array together with the coefficients of the matrix. The number of input pads is 13, distributed as follows: 8 pads for the rows of the matrix, 1 port for the initialization boolean, 2 pads for the alimentation and the ground, and 2 pads for the two non-overlapping clocks (the operation of the array is totally synchronous). We have 9 outputs pads, 8 for the rows of the solution matrix, and 1 for the boolean indicating whether the algorithm has failed or not.

More generally, the design of an $n \times n$ array will require $2n+6$ pads. As usual in two-dimensional systolic arrays implementations, the main limitation to the size of the array is due to the number of input/output pins rather than to area considerations.

The dimensions of a square cell are $380 \times 200 \mu m^2$, and those of a double square cell are $170 \times 200 \mu m^2$. The chip is $3100 \times 2800 \mu m^2$.

6. Conclusion

We have introduced two systolic arrays for solving dense linear systems over $GF(p)$, p prime number. These two arrays implement elimination algorithms with partial pivoting, although the operation of the array remains purely systolic.

Moreover, there is no feedback cycles in our arrays. The importance of designing arrays without feedback cycles is

emphasized in Kung and Lam [8]: acyclic implementations usually exhibit more favorable characteristics with respect to fault-tolerance, two-level pipelining, and problem decomposition (see Hwang and Cheng [7]) in general.

Using massive parallelism and pipelining, the systolic array concept allows a system implementor to design extremely efficient machines for specific computations. The suitability of the systolic model to chip design is well illustrated by the prototype array that we have designed within only a few weeks.

References

- [1] H.M. AHMED, J.M. DELOSME, M. MORF, "Highly concurrent computing structures for matrix arithmetic and signal processing," *IEEE Computer* 15, 1, 65-82 (1982)
- [2] C. ANTONINO, B. BOSCH, H. DELORI, A. GUYOT, J.F. PAILLOTIN, The French MPC 84-85, Internal Report, TIM3/INPG Grenoble, February 1986
- [3] G.R. BLAKLEY, "A computer algorithm for computing the product AB modulo M ," *IEEE Trans. Comput.* 32, 4 (1983), 497-500
- [4] A. BOJANCZYK, R.P. BRENT, H.T. KUNG, Numerically stable solution of dense systems of linear equations using mesh-connected processors, Technical Report, Carnegie Mellon University, 1981
- [5] M. COSNARD, Y. ROBERT, M. TCHUENTE, Matching parallel algorithms with architectures: a case study, IFIP Working Conference on Highly Parallel Computers, Nice, France, 24-26 March 1986
- [6] W.M. GENTLEMAN, H.T. KUNG, Matrix triangularisation by systolic arrays, Proc SPIE 298, Real-time Signal Processing IV, San Diego, California, 1981, 19-26
- [7] K. HWANG, Y.H. CHENG, "Partitioned matrix algorithm for VLSI arithmetic systems," *IEEE Trans. Computers* 31, 12 (1982), 1215-1224
- [8] H.T. KUNG, M.S. LAM, "Fault-tolerance and two-level pipelining in VLSI systolic arrays," *J. of Parallel & Distributed Computing* 1, 1 (1984), 32-63
- [9] H.T. KUNG, C.E. LEISERSON, "Systolic arrays for (VLSI)," in *Proc. of the Symposium on Sparse Matrices Computations*, I.S. Duff et al. eds, Knoxville, Tenn. (1978), 256-282
- [10] D.E. KNUTH, *The art of computer programming*, vol. 2, chap. 4.6.2., Addison WWesley (1969)
- [11] J. MAKHOUL, "Linear prediction: a tutorial review," *Proc. IEEE* 63, 4 (1975), 561-580
- [12] M.A. MORRISON, J. BRILLHART, "A method of factoring and the factorization of F_7 ," *Math. of Comput.* 29, 129 (1975), 183-205
- [13] D. PARKINSON, M. WUNDERLICH, "A compact algorithm for Gaussian elimination over $GF(2)$ implemented on highly parallel computers," *Parallel Computing* 1 (1984), 65-73
- [14] J.F. PAILLOTIN, Le système LUCIE, Internal Report, TIM3/INPG Grenoble, July 1985
- [15] R. POET, "The design of special purpose hardware to factor large integers," *Comput. Physics Communications* 37 (1985), 337-341
- [16] Y. ROBERT, M. TCHUENTE, "Résolution systolique de systèmes linéaires denses," *RAIRO Modélisation et Analyse Numérique* 19, 2 (1985), 315-326
- [17] R. SCHREIBER, P. KUEKES, "Systolic linear algebra machines in digital signal processing," in *VLSI and Modern Signal Processing*, S. Y. Kung et al. eds, Prentice Hall, Englewood Cliffs, NJ, 1985