



MICROCOPY RESOLUTION TEST CHART NATIONAL BUREAU OF STANDARDS 1963 A

•

# Presented at the 1<sup>St</sup> International Colloquium on Vector and Parallel Computing in Scientific Applications, Paris, March 1983

Contract N00014-82-K-0703

A SYSTOLIC ARCHITECTURE FOR SINGULAR VALUE DECOMPOSITION

Robert Schreiber Computer Science Department Stanford University Stanford, California 94305, USA

1 INTRODUCTION

**(1)** 

AUA13005

J Systelic arrays are highly parallel computing structures specific to particular computing tasks. They are well-suited for reliable and inexpensive implementation using many identical VLSI components. The designs consist of one and two-dimensional lattices of identical processing elements. Communication of data occurs only between neighboring cells. Control signals propagate through the array like data. These characteristics make it feasible to construct very large arrays.

Several modern methods in digital signal processing require real time solution of some of the basic problems of linear algebra <del>[13].</del> Fortunately systolic arrays have heen developed for many of these problems [[4,10,12]. But several gaps remain. Only partially satisfying results have been obtained for the eigenvalue and singular value decompositions, for example. Thus durance-t

 $\mathbf{A} = \mathbf{U} \ \mathbf{I} \ \mathbf{V}^{\mathsf{T}}$ 

accomodated.

Here we consider a systolic array for the singular value decomposition (SVD). An SVD of an m x n (m > n) matrix A is a factorization >

where U is m x n with orthonormal columns,  $\mathbb{C}$  = diag $(z_1, z_2, \dots, z_n)$  with  $z_1 \ge z_2 \ge \dots \ge z_n$ , and V is orthogonal. There are many important applications of the SVD [1,6,13].

There have been several earlier investigations of parallel SVD algorithms and arrays. First, Finn, Luk, and Pottle describe a systolic structure of  $n^2/2$  processors and two algorithms that use it. But the convergence of their algorithms has not been proved and may be slow [3]. Heller and Ipsen [8] describe an array for computing the singular values of a banded matrix with bandwidth w. They use O(w) processors and  $O(wn^2)$  time. Brent and Luk [2] describe an n/2 processor linear array that implements a one-sided orthogonalization method and converges reliably in O(n log n) time. Unfortunately the processors in this array are quite complex, and it is not clear that matrices with more than n columns can be efficiently

the author discossed

In this paper we discuss two top.cs. First, Is we showshow an architecture for computing the eigenvalues of a symmetric matrix can be modified to compute singular values and vectors. Second,

Ac using VLSI chips of the implementation using VLSI chips of these systolic eigenvalue and SVD arrays.

The SVD is often used to regularize ill-conditioned problems. In these there are p < n large singular values and n-p that are much smaller. What is needed is the pseudoinverse of the rank p matrix closest (with respect to the 2-norm) to A,

$$A_{(p)} = u_1 \sigma_1 v_1^T + \dots + u_p \sigma_p v_p^T$$

We have recently developed a new algorithm to com-

 $\mathbf{A}_{(\mathbf{p})}^{\mathsf{T}}$  that involves nothing but a sequence of matrix-matrix products, for which systolic arrays are well-known (see, e.g., [9].) An alternate form of the algorithm can be used to compute the related orthogonal projection matrix

 $P_{(p)} = v_1 v_1^T + \dots + v_p v_p^T$ 

2 AN SVD ARCHITECTURE

Let A be a given matrix. The singular values of A will be obtained in two phases:

> 1. A is reduced to an upper triangular matrix B with bandwidth k+1,

 $b_{ij} = 0$  if i > j or i < j - k.

and B = QAP where Q and P are orthogonal.

2. B is diagonalized by an iterative process equivalent to implicitly shifted QR iteration on B B.

With k=1 this is the standard method of Colub and Reinsch [7]. The reason for allowing k>l is an increase in the parallelism. In phase 1, kn processors are employed; the time is O(mn/k). In phase 2, 2k processors are used; the time per iteration is 6n+O(k).

### 2.1 Reduction to banded form

83

The reduction step uses a k x n trapezoidal array that has been described in detail previously [12]. Let the m x n matrix X be partitioned as



 $\overline{\mathbf{03}}$ 

16



where  $X_{11}$  is k x k. The array applies a sequence of Givens rotations to the rows of X to zero the first k columns below the main diagonal. If Q is the product of these rotations, then

$$QX = \begin{bmatrix} R_{11} & Y_{12} \\ 0 & Y_{22} \end{bmatrix}$$

where  $R_{11}$  is k x k upper triangular.  $R_{11}$ ,  $Y_{12}$ ,

 $Y_{22}$ , and the parameters of the rotations that make up Q all flow out from the array. The time required is m. (Here and below we give "times" in units of the time required for an individual cell in the array to carry out its computation.)

Now let

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$$

be the given matrix. Send A through the array to produce

 $Q_{1}^{T} = \begin{bmatrix} R_{11} & C_{12} \\ 0 & C_{22} \end{bmatrix}$ 

Next send  $[C_{12}^T, C_{22}^T]$  through to produce

$$P_1[C_{12}^T, C_{22}^T] = [L_{12}^T, \overline{A}_{22}^T]$$

(Although the input matrix has m columns, the array can handle this factorization in time m by making  $|\mathbf{m}/\mathbf{n}|$  passes over the data [12]. Now continue this process using  $\overline{A}_{22}$  in place of A. After  $J \Xi [n/k]$ such steps we have produced a k+l diagonal, upper triangular matrix B,



such that A = QBP where Q and P are orthogonal. The total time used is mJ = mn/k.

The transposition of data required can be done by a specialized switching device, a "systolic shifter," described earlier [12].

When singular vectors are to be computed, the rotations generated by the array may be applied to identity matrices of order m and n. This can be done by the array. These matrices accumulate the product of the rotations used, that is the orthogonal matrices Q and P above.

### 2.2 QR iteration

Now we consider QR iteration to get the singular values of B, hence those of A. We shall generate a sequence of matrices  $\{B^{(1)}\}$  having the same structure as B and converging to a diagonal matrix.  $B^{(0)} = B$  and  $B^{(i+1)} = P^{(i)}B^{(i)}O^{(i)}$  where  $P^{(i)}$  and Q<sup>(i)</sup> are orthogonal.

First we consider QR iteration on  $\mathbf{B}^T\mathbf{B}$  without shifts. This can be realized by the procedure

1. Find  $Q^{(1)}$  such that  $L^{(1)} = B^{(1)}Q^{(1)}$ 

is lower triangular,

2. Find P<sup>(1)</sup> such that

$$B^{(i+1)} = P^{(i)}L^{(i)}$$

is upper triangular.

Both steps of this procedure can be carried out by the Heller-Ipsen (HI) array [8]. This is a k x w rectangular array for QR factorization of w-diagonal matrices. In this array, plane rotations are generated at the left edge and move to the right, affecting a pair of matrix rows. Take w = k+1.

 $B^{(1)}$  enters the matrix at the bottom, each diagonal entering, one element at a time, into one of the processors. The array annihilates the elements of the upper triangle of  $B^{(1)}$ . This causes fill-in of k diagonals in the lower triangle. The result-

ing matrix  $L^{(1)}$  emerges from the top in the same diagonal-per-processor format. It immediately enters a second array. This array annihilates the lower triangle of  $L^{(1)}$  and the resulting upper triangular matrix  $B^{(i+1)}$  emerges from the top (Fig. 1) The time is 2n+4k per iteration: element  $a_n$  enters

the bottom array at time 2n, leaves at the upper left corner at time 2n+2k, and leaves the top array at time 2n+4k.

Unshifted QR converges slowly. The rate of convergence of  $b_{11}$  to  $\sigma_1$  is  $\sigma_2^2/\sigma_1^2$ . In some situations this may be adequate and the simplicity of the structure used is then a real advantage. It is also easy to pipeline the iterations. As B(i+1) comes out of the second array it can be sent directly into another pair of arrays to begin the (i+1)th iterations, etc. As many as n/kk iterations can be effectively pipelined; any more and the pipe length exceeds a, so that the pipe never gets full. If we choose  $k = O(n^{1/2})$  and pipeline  $n/4k = O(n^{1/2})$  iterations choes the number of processors in both arrays is  $O(n^{3/2})$  and the total time, assuming O(n) iterations of QR are required, is also  $O(n^{3/2})$ . These considerations also apply to the array implementation of the implicitly shifted QR algorithm that is discussed below, with one important proviso. When pipelin-ing the iterations, some strategy for choosing several shifts in advance must be used.

### 2.2.1 Implicitly shifted QR iteration

To obtain adequate convergence speed we need to incorporate shifts. Following Stewart [14], suppose that one QR iteration with shift  $\lambda$  is performed on B<sup>1</sup>B, and the orthogonal matrix so generated is Q. Then proceed as follows:

- Let Q<sub>0</sub> be any matrix whose first k columns are the same as those of Q;
- Using the same technique as in Section 2.1, reduce BQ to upper triangular k+1 diagonal form, yielding a matrix B'.

To use the trapezoidal array as described above to carry out step 2 would be inefficient. Rather we proceed as follows.  $Q_0^+$  is composed of plane rotations that zero the first k columns of  $(B^TB-\lambda)$ below the main diagonal. Applying  $Q_0$  to B causes fill-in in the k diagonals below the main diagonal, confined to rows 2, 3, ..., 2k. See Fig. 2 for the case k=2.



Figure 1



Structure of  $BQ_0$  , k=2

Let the first 2k rows of BQ be sent into a k x 2k+1 HI array. By a sequence of plane rotations applied to the rows, the array removes the "bulge" in the lower triangle, adding a bulge of the same shape in the first 3k columns of the upper triangle. This data flows directly into another k x 2k+1 HI array that removes the elements to the right of the  $k^{ch}$  superdiagonal and causes a new bulge to appear in the lower triangle, in columns k+1 through 3k-1 and extending to row 3k. (The second HI array is the mirror image of the first. Rotations are generated at its right edge and move left, affecting pairs of matrix columns. Let  $P_1$ and  $Q_1$  be the orthogonal matrices implicitly used by the two HI arrays. The matrix emerging from the second array is

$$\mathbf{B}_1 = \mathbf{P}_1 \mathbf{B} \mathbf{Q}_0 \mathbf{Q}_1$$

and it has the form





time



## Figure 3

"Chasing the bulge" with two k x 2k+1 Heller-Ipsen arrays

Now we do exactly the same thing to  $B_{21}$ , etc. This yields matrices

 $B_{j} = P_{j}B_{j,j-1}Q_{j}, j=2,...,J$ 

with

$$\mathbf{B}_{j} = \begin{bmatrix} \mathbf{B}_{j,j} \\ \mathbf{0} & \mathbf{B}_{j+1,j} \end{bmatrix}$$

and J = [(n-1)/k]. Finally  $B' = P_J \dots P_1 BQ_0 \dots Q_J$  is the matrix we require.

The time needs is 6n. It takes 2k steps for an HI array to start producing output. Thus, the second array starts its output at time 4k. The first element of  $B_{j+1,j}$ , which is the  $(k+1)^{st}$  element

of the main diagonal to come out of the second array, comes out at time 6k. By this time the first arrays inputs have become idle, so this element can immediately reenter. Therefore one step, from  $B_j$  to  $B_{j+1}$ , takes time 6k. There are  $\lceil (n-1)/k \rceil$  such steps, hence about 6n time for the whole process.

### 2.3 Complex matrices

In signal processing applications, complex matrices often arise. Here we discuss the algorithms to be used for QR iteration with complex matrices. Essentially we show that the plane rotations used can be of a special form:

(1)  $x^{1} = c^{+}x + \sigma y$  $y^{1} = -\sigma x + c y$ 

where x, y, and c are complex and  $\sigma$  is real. This saves 1/4 of the multiplications used by a fully complex plane rotation with complex s instead of  $\tau \sim -12$  are used instead of 16. We shall call these c,  $\sigma$  rotations.

It is possible to compute the SVD of a complex m x n matrix A =  $A_R$ +i $A_1$  using real arithmetic. One finds the SVD of the 2m x 2n real matrix

 $\begin{bmatrix} A_{R} & -A_{I} \\ A_{I} & A_{R} \end{bmatrix}$ 

Among the 2n singular values each singular value of A occurs twice, and the singular vectors are of the form  $[x_1^T, x_1^T]$  where  $x \approx x_n + ix_1$  is a singular vector of A. But the cost is much greater. In units where the cost of doing an m x n real SVD is one, the cost for the real 2m x 2n SVD is 8 while the complex m x n approach costs 3 (not 4, since the use of the c,o rotations saves 1/4 of the work).

We now show that the c,o rotations suffice. To start, we note that the banded matrix B produced by the reduction phase can always be chosen to have positive real elements on its main and  $k^{th}$  superdiagonals. Indeed the reduction B = QAP to k+1 diagonal, upper triangular form is not unique:

 $B = QD_1 (D_1^{-1}AD_2^{-1}) D_2P$ 

is also such a reduction for any unitary diagonal matrices  $D_1$  and  $D_2$ . These can always be chosen to give B the stated property. In fact, the trapexoidal array can do this automatically [12]. When it generates a rotation to zero some matrix element, the second element of the pair (x,y) for instance, it chooses the paramenters so that the result of the rotation is the pair

 $((|\mathbf{x}|^2 + |\mathbf{y}|^2)^{1/2}, 0)$ 

Furthermore, the elements to be zeroed are the real elements resulting from previous rotations. The rotations to do the zeroing can, for this reason, be taken to be  $c_{1,0}$  rotations.

Now we look at the second phase. Because of the structure of B, the main, k<sup>th</sup> super and k<sup>th</sup> subdiagonals of B<sup>T</sup>B are all real. The rotations that comprise Q<sub>0</sub> can be taken to be c,  $\sigma$  rotations since they zero real elements. And by keeping track of the locations of real elements one can show that in BQ<sub>0</sub> all elements of the outer diagonals are real. Again because the elements to be annihilated are real, c,  $\sigma$  rotations can be used to eliminate the bulge. A matrix with the same structure as BQ<sub>0</sub> results, and the proof therefore follows by induction.

### 2.4 An alternate scheme

Gene Golub has pointed out that the eigenvalues of the 2n x 2n matrix

-	0	8
C =	BT	0

are the singular values of A taken with positive and negative sign, and if  $(x^T, y^T)$  is an eigenvector of C then x is a left singular vector of B and y is a right singular vector of B [5]. Thus we may attempt to find the eigendecomposition of C. After a symmetric interchange of rows and columns corresponding to the permutation (n+1, 2, n+2, 2, ..., 2n, n), C is a symmetric 4k-1 diagonal matrix. A 2k-1 x 4k-1 HI array can implement one step of the QR method with shifts for this matrix in n + O(k)time [10]. In the complex case, both C and the permuted C have real outermost diagonals, so c,o rotations can be used. Thus, although twice as much hardware is used, the time per iteration is 1/6 as great as for the previous scheme.

#### 3 VLSI IMPLEMENTATION

Now we consider how to build the cells of the HI array. The fundamental unit we use in this construction is a multiply-add cell, whose function is this:



Outputs leave the cell one clock after inputs enter.

Although other primitive units (CORDIC blocks, for example) might be used, we feel that the multiply-add is a good basis for such an investigation. Currently, a floating point multiply-add is about what can be integrated on a single chip. It is almost universally useful. Indeed, the multiplyadd pair is often the inner loop in numerical linear algebraic computations. Even when larger cells and pieces of arrays can be integrated into single chips, designs based on the multiply-add primitive will be useful.

We shall discuss implementation of the HI array cells for complex data. The real case was discussed earlier [11] as were the cells of the trapezoidal array [12]. The complex RI array triangularizes a banded input matrix using  $c, \sigma$  rotations of the form (1). The rotations are applied to a pair (x,y) of matrix elements by an internal cell



after having been generated by a boundary cell



by

 $x' = (n^{2} + |x|^{2})^{1/2}$  c = x / x'  $\sigma = n / x'$ In the internal cell computation, 4 quantities are computed, each requiring 3 multiplies and 2 adds. Let x and x denote the real red forcing the real configuration.

are computed, each requiring 3 multiplies and 2 adds. Let  $z_R$  and  $z_r$  denote the real and imaginary parts of the complex quantity z. The computed values are

 $\begin{aligned} \mathbf{x}_{\mathbf{R}}^{*} &= \mathbf{c}_{\mathbf{R}}\mathbf{x}_{\mathbf{R}} + \mathbf{c}_{\mathbf{I}}\mathbf{x}_{\mathbf{I}} - \mathbf{\sigma}\mathbf{y}_{\mathbf{R}}, \\ \mathbf{x}_{\mathbf{T}}^{*} &= \mathbf{c}_{\mathbf{R}}\mathbf{x}_{\mathbf{T}} - \mathbf{c}_{\mathbf{T}}\mathbf{x}_{\mathbf{R}} - \mathbf{\sigma}\mathbf{y}_{\mathbf{T}}; \\ \mathbf{y}_{\mathbf{R}}^{*} &= \mathbf{c}_{\mathbf{R}}\mathbf{y}_{\mathbf{R}} - \mathbf{c}_{\mathbf{I}}\mathbf{y}_{\mathbf{I}} + \mathbf{\sigma}\mathbf{x}_{\mathbf{R}}, \\ \mathbf{y}_{\mathbf{I}}^{*} &= \mathbf{c}_{\mathbf{R}}\mathbf{y}_{\mathbf{I}} + \mathbf{c}_{\mathbf{T}}\mathbf{y}_{\mathbf{R}} + \mathbf{\sigma}\mathbf{x}_{\mathbf{I}}. \end{aligned}$ 

Using 4 multiply-add chips we can construct a compound cell that gives these results in the least possible time, 3 clocks. We assume that complex quantities are represented in "word serial" form, with the real part preceding the imaginary part on the same data path. A schedule using 4 chips that achieves the minimum latency is shown in Table 1.

Table	1.	Schedule	for	Internal	HI	Cell	

	Chij	2	Input/Output			<u> </u>	
time	C1	C2	c,	У	X	C3	C4
0			CR				
1	c <sub>R</sub> y <sub>R</sub>		CI CI	y <sub>R</sub>	×R	CR <sup>X</sup> R	
2	c <sub>I</sub> y <sub>R</sub>	-cIAI	σ	y <sub>I</sub>	×I	-c <sub>I</sub> x <sub>R</sub>	c <sup>I</sup> xI
3	σ×R	<sup>c</sup> R <sup>y</sup> I	CR			-σ у <sub>R</sub>	CRXI
4		σ×I	C I	y'R	×'R		-σ y <sub>I</sub>
5			σ	y'I	x'I		

The computation at the boundary cell is this: given inputs x and  $\mathfrak{n}^2$ , compute

 $n'^{2} = n^{2} + x_{R}^{2} + x_{I}^{2}$   $n' = \sqrt{n'^{2}}$   $c_{R} = x_{R} / n'$   $c_{I} = x_{I} / n'$   $\sigma = n / n'$ 

A second primitive, for divide and square root, is needed to implement the boundary cell. We assume that a chip for computing

$$(a,b) ===> a / b^{1/2}$$

is available. A compund cell using one multiplyadd and two of these square root chips can produce results at the rate required to keep up with the internal cell. A schedule is shown in Table 2. The overall array timing is now that of the "ideal" HI array in which everything happens in a single cycle (of length 3 chip clocks). The cells are used 1/2 of the time, but two independent problems can be solved simultaneously, making full use of the hardware.

Table	2.	Schedule	for H	I B	oundarv Cell	

<u>t18</u>	e	Chips			
	mult-add	sqrt #1	sgrt #2		
1	$\rho^2 + x_p^2$			02,x.	
2	+x <sup>2</sup> <sub>T</sub> (=p <sup>12</sup> )		1	x,	
3	-	$x_{p}[p^{+2}]^{-1/2}$	}	p _	
4		$x_{\tau} \{\rho^{2}\}^{-1/2}$	ļ	p'2,c	
5		P[P <sup>,2</sup> ] <sup>-1/2</sup>	[p <sup>2</sup> ] <sup>-1/2</sup>	- c, ^	
6	1			0.0'	

#### ACKNOWLEDGEMENT

This research was partially supported by the Office of Naval Research under Contract N00014-82-K-0703 and by ESL, Incorporated, Sunnyvale, Calif.

#### REFERENCES

- /1/ H. C. Andrews and C. L. Patterson : "Singular Value Decomposition and Digital Image Processing", <u>IEEE Trans. Acoustics, Speech, and Signal Processing ASSP-24</u> (1976), pp. 26-53.
- /2/ Richard P. Brent and Franklin T. Luk : "A Systolic Architecture for the Singular Value Decomposition", Report TR-CS-82-09, Department of Computer Science, The Australian National University, Camberra, 1982.
- /3/ Alan M. Finn, Franklin T. Luk, and Christopher Pottle: "Systolic Array Computation of the Singular Value Decomposition", <u>Real Time</u> <u>Signal Processing V, SPIE Vol. 341</u>, Bellingham Wash., Society of Photo-optical Instrumentation Engineers, 1982.
- /4/ W. M. Gentleman and H. T. Kung : "Matrix Triangularization by Systolic Array", <u>Real Time</u> <u>Signal Processing IV</u>, <u>SPIE Vol. 298</u>, <u>Bell-</u> ingham, Wash., Society of Photo-optical Instrumentation Engineers, 1981.

/5/ Gene Golub. Private communication.

- /6/ G. H. Golub and F. T. Luk : "Singular Value Decomposition: Applications and Computations", ARO Report 77-1, <u>Trans. of the 22<sup>nd</sup> Conf. of Army Mathematicians</u> (1977), pp. 577-605.
- /7/ G. H. Colub and C. Reinsch : "Singular Value Decomposition and Least Squares Solutions", <u>Numer. Math.</u> 14, (1970), pp. 403-420.

- /8/ Don E. Heller and Ilse C. F. Ipsen : "Systolic Networks for Orthogonal Decompositions, with Applications", <u>SIAM J. Scient. and Stat.</u> <u>Comput.</u>, to appear.
- /9/ H. T. Kung and Charles Leiserson : "Systolic Arrays for (VLSI)", in Carver Mead and Lynn Conway, <u>Introduction to VLSI Systems</u>, Reading, Mass., Addison-Wesley, 1980.
- /10/ Robert Schreiber : "Systolic Arrays for Eigenvalue Computation", <u>Real Time Signal Processing V, SPIE Vol. 341</u>, Bellingham, Wash., Society of Photo-optical Instrumentation Engineers, 1982
- /11/ Robert Schreiber : "Systolic Arrays for Eigenvalues", Proc. of the Inter-American Work-<u>shop in Numerical Analysis</u>, New York, Springer-Verlag, to appear.
- /12/ Robert Schreiber and Philip J. Kuekes: "Systolic Linear Algebra Machines in Digital Signal Processing", in Sun-Yuan Kung, ed., Proc. of the USC Workshop on VLSI and Modern Signal Processing, Englewood Cliffs, New Jersey, Prentice-Hall, to appear.
- /13/ J. M. Speiser and H. J. Whitehouse : "Architectures for Real-Time Matrix Operations", <u>Proc. Government Microcircuit Applications</u> <u>Conf.</u>, held in Houston, Tex., 1980.
- /14/ G. W. Stewart : <u>Introduction to Matrix Computations</u>, New York, Academic Press, 1973.

.

