

THE SOLUTION OF SINGULAR-VALUE AND EIGENVALUE
PROBLEMS ON SYSTOLIC ARRAYS

Richard P. Brent and Franklin T. Luk

0. SUMMARY

Parallel algorithms are presented for computing a singular-value decomposition of an $m \times n$ matrix ($m \geq n$) and an eigenvalue decomposition of an $n \times n$ symmetric matrix. A linear array of $O(n)$ processors is proposed for the singular-value problem and the associated algorithm requires time $O(mnS)$, where S is the number of Jacobi sweeps (typically $S \leq 10$). A square array of $O(n^2)$ processors with nearest-neighbor communication is proposed for the eigenvalue problem; the associated algorithm requires time $O(nS)$.

1. INTRODUCTION

A singular-value decomposition (SVD) of a real $m \times n$ ($m \geq n$) matrix A is its factorization into the product of three matrices:

$$(1.1) \quad A = U \Sigma V^T,$$

where U is an $m \times n$ matrix with orthonormal columns, Σ is an $n \times n$ nonnegative diagonal matrix and the $n \times n$ matrix V is orthogonal. This decomposition has many important scientific and engineering applications (cf. [1,11,27,28]).

If the matrix A is square (i.e., $m=n$) and symmetric, we may adjust the sign of the elements of Σ so that $U = V$. We then obtain an eigenvalue decomposition:

$$(1.2) \quad A = U D U^T,$$

where U is orthogonal and D diagonal. The advent of massively parallel computer architectures has aroused much interest in parallel singular-value and eigenvalue procedures, e.g. [2,4,5,6,7,9,13,14,16,18,20,22,23,24,25]. Such architectures may turn out to be indispensable in settings where real-time computation of the decompositions is required [27,28]. Speiser and Whitehouse [27] survey parallel processing architectures and conclude that systolic arrays offer the best combination of characteristics for utilizing VLSI/VHSIC technology to do real-time signal processing. (See also [17,28].)

In this paper we present an array of $O(n)$ linearly-connected processors which computes an SVD in time $O(mnS)$. Here S is a slowly growing function of n which is conjectured to be $O(\log n)$; for practical purposes S may be regarded as a constant (see [21]). Our array implements a one-sided orthogonalization method due to Hestenes [15]. His method is essentially the serial Jacobi procedure for finding an eigenvalue decomposition of the matrix $A^T A$, and has been used by Luk [20] on the ILLIAC IV computer. We also describe how one may implement a Jacobi method on a two-dimensional array of processors to compute an eigenvalue decomposition of a symmetric matrix. Our array requires $O(n^2)$ processors and $O(nS)$ units of time. Assuming that $S = O(\log n)$, this time requirement is within a factor $O(\log n)$ of that necessary for the solution of n linear equations in n unknowns on a systolic array [2,3,17].

Results similar to ours have been reported in the literature. For computing the SVD, Sameh [23] describes an implementation of Hestenes' method on a ring of $O(n)$ processors. Suppose n is even (the result is similar for an odd n). At each orthogonalization step $\frac{n}{2}$ column rotations are performed. Sameh's permutation scheme requires $3n - 2$ steps to ensure the execution of every possible pairwise rotation at least once; our permutation scheme (presented in Section 3) requires only $n - 1$ steps.

Parallel Jacobi methods for computing eigenvalues are given in [7,16,22]. However, the procedure of Sameh [22] may be unsuitable for multiprocessor arrays. For simplicity, assume again that n is even, so $\frac{n}{2}$ off-diagonal elements can be set to zero at each elimination step. Let us compare the number of permutations necessary for the annihilation of each off-diagonal element at least once. Our procedure (see Sections 3 and 6) requires $n - 1$ permutations, which is optimal; that of Chen and Irani [7] requires n permutations. The scheme of Kuck and Sameh [16] is worse. Their basic scheme appears to cycle every $2n - 2$ steps and to miss some off-diagonal elements. A modification ("the second row and column are shifted to the n -th position after every $(n - 1)$ orthogonal transformations") can be made to overcome this problem, but the modified scheme requires $(n - 1)^2$ permutations [7].

Let us generalize the notion of a "sweep" and use it to denote a minimum-length sequence of rotations that eliminates each off-diagonal element at least once [7]. It is probably fair to assume

that the Jacobi procedures in [7, 16] and in this paper require an equal number (say S) of sweeps for convergence. For the algorithms presented in this paper a sweep always consists of $n(n - 1)/2$ rotations (the minimal number possible), but this is not the case for the Chen and Irani or Kuck and Sameh algorithms mentioned above. The architecture proposed in [7] is a linear array of $O(n)$ processors; the associated Jacobi method requires time $O(n^2S)$. The architecture described in [16] is a square array of $O(n)$ processors, with boundary wraparounds and a broadcast unit. The associated algorithm requires time $O(n^3S)$. In comparison, our procedure requires $O(n^2)$ processors and $O(nS)$ units of time.

The principal results of this paper were first reported in [4,5]. A related SVD algorithm is presented by the authors and Van Loan. It requires $O(n^2)$ processors and $O(nS)$ time to compute the singular values of an $n \times n$ matrix. For a generalization of this result, see [6].

This paper is organized as follows. Sections 2-4 are devoted to the singular-value problem and Sections 5-7 to the eigenvalue problem. Hestenes' method is reviewed in Section 2. The new ordering is described in Section 3 and the corresponding SVD algorithm in Section 4. The serial Jacobi method is outlined in Section 5. Details are filled in and some variations and extensions of the basic algorithm are mentioned in Section 7.

The SVD algorithm described in Sections 3-4 below is being implemented on an experimental 64-processor systolic array by Speiser at the Naval Ocean Systems Center (San Diego).

2. HESTENES' METHOD

We wish to compute an SVD of an $m \times n$ matrix A , where $m \geq n$. An idea is to generate an orthogonal matrix V such that the transformed matrix $AV = W$ has orthogonal columns. Normalizing the euclidean length of each nonnull column to unity, we get the relation

$$(2.1) \quad W = \tilde{U} \Sigma,$$

where \tilde{U} is a matrix whose nonnull columns form an orthonormal set of vectors and Σ is a nonnegative diagonal matrix. An SVD of A is given by

$$(2.1') \quad A = \tilde{U} \Sigma V^T.$$

As a null column of \tilde{U} is always associated with a zero diagonal element of Σ , there is no essential difference between (1.1) and (2.1').

Hestenes [15] uses plane rotations to construct V . He generates a sequence of matrices $\{A_k\}$ using the relation

$$A_{k+1} = A_k Q_k,$$

where $A_1 = A$ and Q_k is a plane rotation. Let $A_k \equiv (a_{rs}^{(k)}, \dots, a_{rs}^{(k)})$ and $Q_k \equiv (q_{rs}^{(k)})$, and suppose Q_k represents a rotation in the (i, j) plane, with $i < j$, i.e.

$$(2.2) \quad \begin{aligned} q_{ii}^{(k)} &= \cos \theta, & q_{ij}^{(k)} &= \sin \theta, \\ q_{ji}^{(k)} &= -\sin \theta, & q_{jj}^{(k)} &= \cos \theta. \end{aligned}$$

We note that postmultiplication by Q_k affects only $\tilde{a}_i^{(k)}$ and $\tilde{a}_j^{(k)}$, and that

$$(2.3) \quad (\tilde{a}_i^{(k+1)}, \tilde{a}_j^{(k+1)}) = (\tilde{a}_i^{(k)}, \tilde{a}_j^{(k)}) \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}.$$

The rotation angle θ is chosen so that the two new columns are orthogonal.

Adopting the formulas of Rutishauser [21], we let

$$(2.4) \quad \alpha \equiv \|\tilde{a}_i^{(k)}\|_2^2, \quad \beta \equiv \|\tilde{a}_j^{(k)}\|_2^2, \quad \gamma \equiv \tilde{a}_i^{(k)T} \tilde{a}_j^{(k)}.$$

We set $\theta = 0$ if $\gamma = 0$; otherwise we compute

$$(2.5) \quad \begin{aligned} \xi &= \frac{\beta - \alpha}{2\gamma}, \\ t &= \frac{\text{sign}(\xi)}{|\xi| + \sqrt{1 + \xi^2}}, \\ \cos \theta &= \frac{1}{\sqrt{1 + t^2}}, \end{aligned}$$

and $\sin\theta = t \cos\theta$.

The rotation angle always satisfies

$$(2.6) \quad |\theta| \leq \frac{\pi}{4}.$$

However, there remains the problem of choosing (i, j) , which is usually done according to some fixed cycle. An objective is to go through all column pairs exactly once in any sequence (a sweep) of $n(n-1)/2$ rotations. A simple sweep consists of a cyclic-by-rows ordering:

$$(2.7) \quad (1,2), (1,3), \dots, (1,n), (2,3), \dots, (2,n), (3,4), \dots, (n-1,n).$$

Forsythe and Henrici [10] prove that, subject to (2.6), the cyclic-by-rows Jacobi method always converges. Convergence of the cyclic-by-rows Hestenes' method thus follows.

Unfortunately, the cyclic-by-rows scheme is apparently not amenable to parallel processing. In Section 3 we present an ordering that enables us to do $\lfloor \frac{n}{2} \rfloor$ rotations simultaneously. The (theoretical) price we pay is the loss of guaranteed convergence. Hansen [12] discusses the convergence properties associated with various orderings for the serial Jacobi method. He defines a certain "preference factor" for comparing different ordering schemes. Our new ordering is in fact quite desirable, for it asymptotically optimizes the preference factor as $n \rightarrow \infty$. Thus, although the convergence proof of [10] does not apply, we expect convergence in practice to be faster than for the cyclic-by-rows ordering. Simulation results support this conclusion.

To enforce convergence, we may choose a threshold approach [30, pp.277-278]. That is, we associate with each sweep a threshold value, and when making the transformations of that sweep, we omit any rotation based on a normalized inner product

$$\frac{\tilde{a}_{\sim i}^{(k)T} \tilde{a}_{\sim j}^{(k)}}{\|\tilde{a}_{\sim i}^{(k)}\|_2 \|\tilde{a}_{\sim j}^{(k)}\|_2},$$

which is below the threshold value. Although such a strategy guarantees convergence, we do not know any example for which our new ordering fails to give convergence even without using thresholds. Our method, like the cyclic-by-rows method, is ultimately quadratically convergent [29].

The plane rotations are accumulated if the matrix V is desired. We compute

$$V_{k+1} = V_k O_k,$$

with $V_1 = I$. Denoting the r -th Column of V_k (respectively V_{k+1}) by $\tilde{v}_{\sim r}^{(k)}$ (respectively $\tilde{v}_{\sim r}^{(k+1)}$), we may update both A_k and V_k simultaneously:

$$(2.8) \quad \begin{pmatrix} \tilde{a}_{\sim i}^{(k+1)} & \tilde{a}_{\sim j}^{(k+1)} \\ \tilde{v}_{\sim i}^{(k+1)} & \tilde{v}_{\sim j}^{(k+1)} \end{pmatrix} = \begin{pmatrix} \tilde{a}_{\sim i}^{(k)} & \tilde{a}_{\sim j}^{(k)} \\ \tilde{v}_{\sim i}^{(k)} & \tilde{v}_{\sim j}^{(k)} \end{pmatrix} \begin{pmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{pmatrix}.$$

3. GENERATION OF ALL PAIRS

In this section we show how $O(n)$ linearly-connected processors can generate all pairs (i,j) , $1 \leq i < j \leq n$, in $O(n)$ steps. The application to the computation of the SVD and of the symmetric eigenvalue decomposition is described in Section 4 and in Sections 6-7, respectively.

First, suppose n is even. We use $n/2$ processors $P_1, \dots, P_{n/2}$, where P_k and P_{k+1} communicate ($k = 1, 2, \dots, n/2 - 1$). Each processor P_k has registers L_k and R_k , output lines $outL_k$ and $outR_k$, and input lines inL_k and inR_k , except that $outL_1$, inL_1 , $outR_{n/2}$ and $inR_{n/2}$ are omitted. The output $outR_k$ is connected to the input inL_{k+1} as shown in Figure 1.

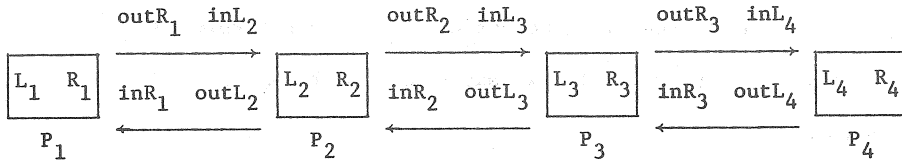


Figure 1: Inter-processor connections for $n = 8$

Initially $L_k = 2k - 1$ and $R_k = 2k$. At each time step processor P_k executes the following program:

```

if  $L_k < R_k$  then process  $(L_k, R_k)$  else process  $(R_k, L_k)$ ;
if  $k=1$  then  $outR_k := R_k$ 
    else if  $k < n/2$  then  $outR_k := L_k$ ;
if  $k > 1$  then  $outL_k := R_k$ ;
{wait for outputs to propagate to inputs of adjacent processors}
if  $k < n/2$  then  $R_k := inR_k$  else  $R_k := L_k$ ;
if  $k > 1$  then  $L_k := inL_k$ ;

```

Here "process (i,j) " means perform whatever operations are required on the pair (i,j) , $1 \leq i < j \leq n$. The operation of the systolic array is illustrated in Figure 2.

We see that the index 1 stays in the register L_1 of processor P_1 . Indices 2, ..., n travel through a cycle of length $n-1$ consisting of the registers $L_2, L_3, \dots, L_{n/2}, R_{n/2}, R_{n/2-1}, \dots, R_1$. During any $n-1$ consecutive steps a pair (i,j) or (j,i) can appear in a register pair (L_k, R_k) at most once. A parity argument shows that (i,j) and (j,i) can not both occur (see Figure 2). Since there are $n/2$ register pairs at each of $n-1$ time steps, each possible pair (i,j) , $1 \leq i < j \leq n$, is processed exactly once during a cycle of $n-1$ consecutive steps.

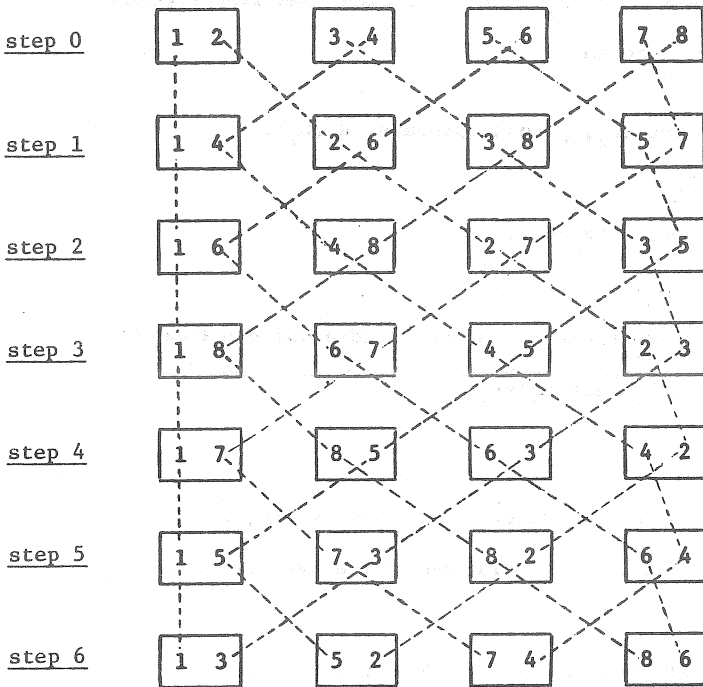


Figure 2: Full cycle of the systolic array for $n = 8$

If n is odd, we use $\lceil n/2 \rceil$ processors but initialize $L_k = 2k - 2$, $R_k = 2k - 1$ for $k=1, \dots, \lceil \frac{n}{2} \rceil$ and omit any "process" calls from processor P_1 .

It is interesting to note that similar permutations are "well known" for use in chess and bridge tournaments, but have apparently not been applied to parallel computation.

4. A ONE-DIMENSIONAL SYSTOLIC ARRAY FOR SVD COMPUTATION

Assume that n is even (else we can add a zero column to A or modify the algorithm as described at the end of Section 3). We use $n/2$ processors $P_1, \dots, P_{n/2}$, as described in Section 3, except that L_k and R_k are now local memories large enough to store a column of A (i.e., L_k and R_k each has at least m floating-point

words). Shift registers or other sequential access memories are sufficient as we do not need random access to the elements of each row.

Suppose processor P_k contains column \tilde{a}_i^C in L_k and column \tilde{a}_j^C in R_k . It is clear that P_k can implement the column orthogonalization scheme in time $O(m)$ by making one pass through \tilde{a}_i^C and \tilde{a}_j^C to compute the inner products (2.4), and another pass to perform the transformations (2.3) or (2.8). Adjacent processors can then exchange columns in the same way that the processors of Section 3 exchange indices. This takes time $O(m)$ if the bandwidth between adjacent processors is one floating-point word. (Alternatively, exchanges can be combined with the transformations (2.3) or (2.8).)

Consequently, we see that $n/2$ processors can perform a full sweep of the Hestenes method in $n - 1$ steps of time $O(m)$ each, i.e., in total time $O(mn)$. Initialization requires that the $(2k-1)$ -th and $2k$ -th columns of A be stored in the local memory of processor P_k for $k = 1, \dots, n/2$; clearly this can also be performed in time $O(mn)$.

The process is iterative. Suppose S sweeps are required to orthogonalize the columns to full machine accuracy. We then have a systolic array of $n/2$ processors which computes the SVD in time $O(mnS)$. By comparison, the serial Hestenes algorithm takes time $O(mn^2S)$. Our simulation results suggest that S is $O(\log n)$, although for practical purposes we can regard S as a constant in the range six to ten [21].

After an integral number of sweeps the columns of the matrix $W \equiv AV$ (see (2.1)) will be stored in the systolic array (two per processor). If V is required, it can be accumulated at the same time that W is accumulated, at the expense of increasing each processor's local memory (but the computation time remains $O(mnS)$): see (2.8).

5. SERIAL JACOBI METHOD

We now consider the related problem of diagonalizing a given $n \times n$ symmetric matrix $A = A_1$. The serial Jacobi method generates a sequence of symmetric matrices $\{A_k\}$ via the relation

$$A_{k+1} = Q_k^T A_k Q_k ,$$

where Q_k is a plane rotation. Let $A_k \equiv (a_{rs}^{(k)})$ and suppose Q_k represents a rotation through angle θ in the (i,j) plane, with $i < j$ (see (2.2)). We choose the rotation angle to annihilate the (i,j) element of A_k . If $a_{ij}^{(k)} = 0$, we do not rotate, i.e., $\theta = 0$. Otherwise we use the formulas in [21] to compute $\sin\theta$ and $\cos\theta$:

$$\begin{aligned} \xi &= \frac{a_{jj}^{(k)} - a_{ii}^{(k)}}{2a_{ij}^{(k)}} , \\ t &= \frac{\text{sign}(\xi)}{|\xi| + \sqrt{1 + \xi^2}} = \tan\theta , \\ \cos\theta &= \frac{1}{\sqrt{1 + t^2}} , \text{ and} \\ \sin\theta &= t \cos\theta . \end{aligned} \tag{5.1}$$

Note that the rotation angle θ may be chosen to satisfy

$$|\theta| \leq \frac{\pi}{4} .$$

The new matrix A_{k+1} differs from A_k only in rows and columns i and j . The modified values are defined by

$$\begin{aligned}
 a_{ii}^{(k+1)} &= a_{ii}^{(k)} - t a_{ij}^{(k)}, \\
 a_{jj}^{(k+1)} &= a_{jj}^{(k)} + t a_{ij}^{(k)}, \\
 (5.2) \quad a_{ij}^{(k+1)} &= a_{ji}^{(k+1)} = 0, \\
 a_{iq}^{(k+1)} &= a_{qi}^{(k+1)} = \cos\theta a_{iq}^{(k)} - \sin\theta a_{jq}^{(k)} \\
 a_{jq}^{(k+1)} &= a_{qj}^{(k+1)} = \sin\theta a_{iq}^{(k)} + \cos\theta a_{jq}^{(k)} \left. \vphantom{a_{iq}^{(k+1)}} \right\} (q \neq i, j).
 \end{aligned}$$

Again we choose (i, j) in accordance to the new ordering introduced in Section 3. The comments that were made in Section 2 concerning various aspects (convergence proof, convergence rate, threshold approach, etc.) of the Hestenes method apply equally well here to the Jacobi procedure.

6. AN IDEALIZED SYSTOLIC ARCHITECTURE

In this section we describe an idealized systolic architecture for implementing the Jacobi method to compute an eigenvalue decomposition of A . The architecture is idealized in that it assumes the ability to broadcast row and column rotation parameters in constant time. In Section 7 we mention how to avoid this assumption.

Assume that the order n is even and that we have a square array of $n/2$ by $n/2$ processors, each processor containing an 2×2 submatrix of $A \equiv (a_{ij})$. Initially processor P_{ij} contains $\begin{pmatrix} a_{2i-1, 2j-1} & a_{2i-1, 2j} \\ a_{2i, 2j-1} & a_{2i, 2j} \end{pmatrix}$ for $i, j = 1, \dots, n/2$, and P_{ij} is connected to its nearest neighbors $P_{i\pm 1, j}$ and $P_{i, j\pm 1}$ (see Figure 3). In general P_{ij} contains four real numbers $\begin{pmatrix} \alpha_{ij} & \beta_{ij} \\ \gamma_{ij} & \delta_{ij} \end{pmatrix}$,

where $\alpha_{ij} = \alpha_{ji}$, $\delta_{ij} = \delta_{ji}$ and $\beta_{ij} = \gamma_{ji}$ by symmetry.

The diagonal processors P_{ii} ($i = 1, \dots, n/2$) act differently from the off-diagonal processors P_{ij} ($i \neq j$, $1 \leq i, j \leq n/2$). Each time step the diagonal processors P_{ii} compute rotations $\begin{pmatrix} c_i & s_i \\ -s_i & c_i \end{pmatrix}$ to annihilate their off-diagonal elements β_{ii} and γ_{ii} , (actually $\beta_{ii} = \gamma_{ii}$), i.e., so that $c_i^2 + s_i^2 = 1$ and

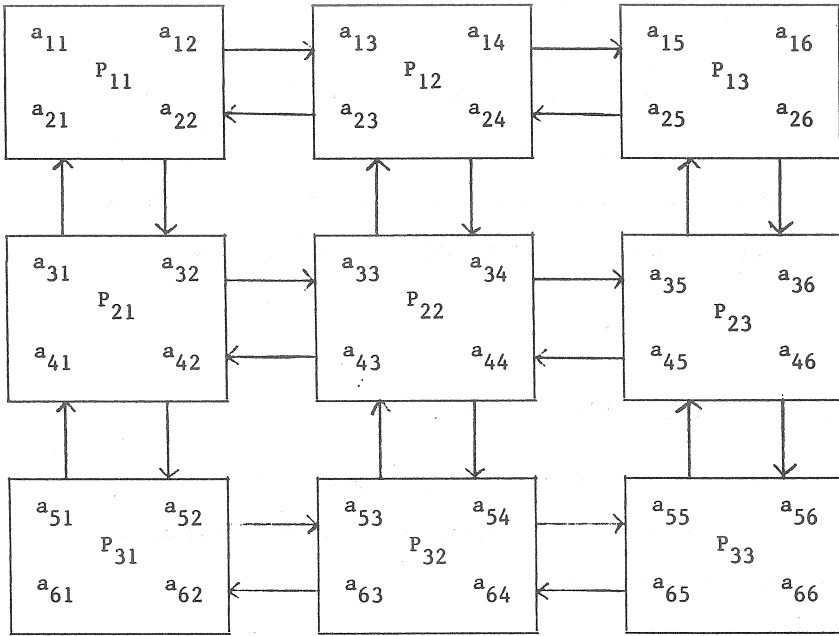


Figure 3: Initial configuration (idealized, $n = 6$)

$$\begin{pmatrix} c_i & -s_i \\ s_i & c_i \end{pmatrix} \begin{pmatrix} \alpha_{ii} & \beta_{ii} \\ \gamma_{ii} & \delta_{ii} \end{pmatrix} \begin{pmatrix} c_i & s_i \\ -s_i & c_i \end{pmatrix} = \begin{pmatrix} \alpha'_{ii} & 0 \\ 0 & \delta'_{ii} \end{pmatrix} \text{ is diagonal. From (5.1) and (5.2)}$$

with a change of notation we find that

$$(6.1) \quad \begin{pmatrix} c_i \\ s_i \end{pmatrix} = \frac{1}{\sqrt{1+t_i^2}} \begin{pmatrix} 1 \\ t_i \end{pmatrix}$$

and

$$\begin{pmatrix} \alpha'_{ii} \\ \delta'_{ii} \end{pmatrix} = \begin{pmatrix} \alpha_{ii} \\ \delta_{ii} \end{pmatrix} + \tau_i \beta_{ii} \begin{pmatrix} -1 \\ 1 \end{pmatrix},$$

where

$$(6.2) \quad \tau_i = \begin{cases} 0 & \text{if } \beta_{ii} = 0 \\ \frac{\text{sign}(\xi_i)}{|\xi_i| + \sqrt{1 + \xi_i^2}} & \text{if } \beta_{ii} \neq 0, \end{cases}$$

and

$$\xi_i = \frac{\delta_{ii} - \alpha_{ii}}{2\beta_{ii}}$$

To complete the rotations which annihilate β_{ii} and γ_{ii} ,

$i = 1, \dots, n/2$, the off-diagonal processors P_{ij} ($i \neq j$) must perform

the transformations $\begin{pmatrix} \alpha_{ij} & \beta_{ij} \\ \gamma_{ij} & \delta_{ij} \end{pmatrix} \leftarrow \begin{pmatrix} \alpha'_{ij} & \beta'_{ij} \\ \gamma'_{ij} & \delta'_{ij} \end{pmatrix}$, where

$$\begin{pmatrix} \alpha'_{ij} & \beta'_{ij} \\ \gamma'_{ij} & \delta'_{ij} \end{pmatrix} = \begin{pmatrix} c_i & -s_i \\ s_i & c_i \end{pmatrix} \begin{pmatrix} \alpha_{ij} & \beta_{ij} \\ \gamma_{ij} & \delta_{ij} \end{pmatrix} \begin{pmatrix} c_j & s_j \\ -s_j & c_j \end{pmatrix}. \quad \text{We assume that the diagonal}$$

processor P_{ii} broadcasts the rotation parameters c_i and s_i to processors P_{ij} and P_{ji} ($j = 1, \dots, n/2$) in constant time, so that the off-diagonal processor P_{ij} has access to the parameters c_i, s_i, c_j and s_j when required. (This assumption is removed in Section 8.)

To complete a step, columns (and corresponding rows) are interchanged between adjacent processors so that a new set of n off-diagonal elements is ready to be annihilated by the diagonal processors during the next time step. This is done in two sub-steps. First, adjacent columns are exchanged as in the SVD algorithm described in Sections 3-4 and as illustrated in Figure 2.

Next, the same permutation is applied to rows, so as to maintain symmetry.

Formally, we can specify the operations performed by a processor P_{ij} with outputs $\text{outh}\alpha_{ij}, \dots, \text{outh}\delta_{ij}, \text{outv}\alpha_{ij}, \dots, \text{outv}\delta_{ij}$, and inputs $\text{inh}\alpha_{ij}, \dots, \text{inv}\delta_{ij}$ by Program 1. Note that outputs of one processor are connected to inputs of adjacent processors in the obvious way, e.g. $\text{outh}\beta_{ij}$ is connected to $\text{inh}\alpha_{i,j+1}$


```

{subscripts (i,j) omitted if no ambiguity results}
{column interchanges}
if i = 1 then [outh $\beta$   $\leftarrow$   $\beta$ ; outh $\delta$   $\leftarrow$   $\delta$ ]
    else if i < n/2 then [outh $\beta$   $\leftarrow$   $\alpha$ ; outh $\delta$   $\leftarrow$   $\gamma$ ];
if i > 1 then [outh $\alpha$   $\leftarrow$   $\beta$ ; outh $\gamma$   $\leftarrow$   $\delta$ ];
{wait for outputs to propagate to inputs of adjacent processors}
if i < n/2 then [ $\beta$   $\leftarrow$  inh $\beta$ ;  $\delta$   $\leftarrow$  inh $\delta$ ]
    else [ $\beta$   $\leftarrow$   $\alpha$ ;  $\delta$   $\leftarrow$   $\gamma$ ];
if i > 1 then [ $\alpha$   $\leftarrow$  inh $\alpha$ ;  $\gamma$   $\leftarrow$  inh $\gamma$ ];
{row interchanges}
if j = 1 then [outv $\gamma$   $\leftarrow$   $\gamma$ ; outv $\delta$   $\leftarrow$   $\delta$ ]
    else if j < n/2 then [outv $\gamma$   $\leftarrow$   $\alpha$ ; outv $\delta$   $\leftarrow$   $\beta$ ];
if j > 1 then [outv $\alpha$   $\leftarrow$   $\gamma$ ; outv $\beta$   $\leftarrow$   $\delta$ ];
{wait for outputs to propagate to inputs of adjacent processors}
if j < n/2 then [ $\gamma$   $\leftarrow$  inv $\gamma$ ;  $\delta$   $\leftarrow$  inv $\delta$ ]
    else [ $\gamma$   $\leftarrow$   $\alpha$ ;  $\delta$   $\leftarrow$   $\beta$ ];
if j > 1 then [ $\alpha$   $\leftarrow$  inv $\alpha$ ;  $\beta$   $\leftarrow$  inv $\beta$ ];

```

Program 1: Column and row interchanges for idealized processor P_{ij}

($1 \leq i \leq n/2$, $1 \leq j < n/2$): see Figure 4. Note that, in Figure 4 and elsewhere, we have omitted subscripts (i,j) if no ambiguity arises, e.g. inv α is used instead of inv α_{ij} .

The only difference between the data flow here and that in Section 4 is that here rows are permuted as well as columns in order to maintain the symmetry of A and move the elements to be annihilated during the next time step into the diagonal processors. Hence, from Section 3 it is clear that a complete sweep is performed every $n - 1$ steps, because each off-diagonal element of A is moved into one of the diagonal processors in exactly one of the steps. Each sweep takes time $O(n)$ so, assuming that $O(\log n)$ sweeps are

required for convergence, the total time required to diagonalize A is $O(n \log n)$.

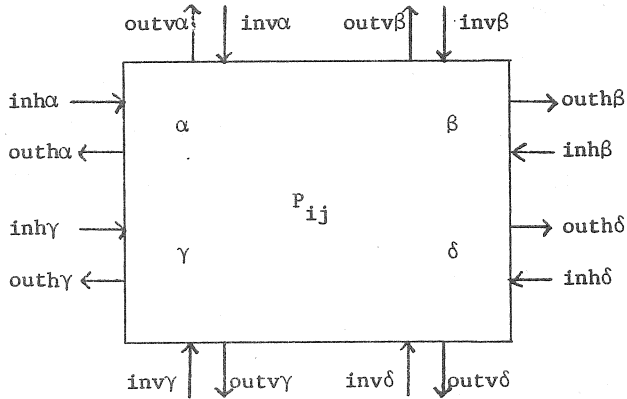


Figure 4: Input and output lines for idealized processor P_{ij} with
nearest-neighbour connections

7. FURTHER DETAILS

Several assumptions were made in Section 6 to simplify the exposition. In this section we show how to remove these assumptions.

7.1 Threshold strategy

It is clear that a diagonal processor P_{ii} might omit rotations if its off-diagonal elements $\beta_{ii} = \gamma_{ii}$ were sufficiently small. All that is required is to broadcast $\begin{pmatrix} c_i \\ s_i \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ along processor row and column i .

As discussed in Section 2, a suitable threshold strategy guarantees convergence, although we do not know any example for which our ordering fails to give convergence even without a threshold strategy.

7.2 Computation of eigenvectors

If eigenvectors are required, the matrix U of eigenvectors can be accumulated at the same time as A is being diagonalized. Each systolic processor P_{ij} ($1 \leq i, j \leq n/2$) needs four additional memory cells

$\begin{pmatrix} \mu_{ij} & v_{ij} \\ \sigma_{ij} & \tau_{ij} \end{pmatrix}$, and during each step sets

$$\begin{pmatrix} \mu_{ij} & v_{ij} \\ \sigma_{ij} & \tau_{ij} \end{pmatrix} \leftarrow \begin{pmatrix} \mu_{ij} & v_{ij} \\ \sigma_{ij} & \tau_{ij} \end{pmatrix} \begin{pmatrix} c_j & s_j \\ -s_j & c_j \end{pmatrix}$$

Each processor transmits its $\begin{pmatrix} \mu & v \\ \sigma & \tau \end{pmatrix}$ values to adjacent processors in the

same way as its $\begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix}$ values (see Program 1). Initially

$\mu_{ij} = v_{ij} = \sigma_{ij} = \tau_{ij} = 0$ if $i \neq j$, and $\mu_{ii} = \tau_{ii} = 1$, $\sigma_{ii} = v_{ii} = 0$.

After a sufficiently large (integral) number of sweeps, we have U defined to working accuracy by

$$\begin{pmatrix} u_{2i-1,2j-1} & u_{2i-1,2j} \\ u_{2i,2j-1} & u_{2i,2j} \end{pmatrix} = \begin{pmatrix} \mu_{ij} & v_{ij} \\ \sigma_{ij} & \tau_{ij} \end{pmatrix}.$$

7.3 Diagonal connections

In Program 1 we assumed that only horizontal and vertical nearest-neighbour connections were available. Except at the boundaries, diagonal connections are more convenient. This is illustrated in Figures 5 and 6 (with subscripts (i,j) omitted).

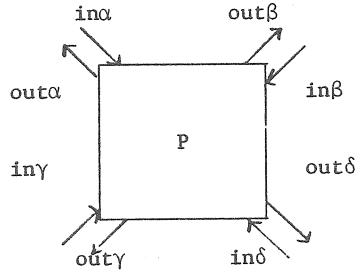


Figure 5: Diagonal input and output lines for processor P_{ij}

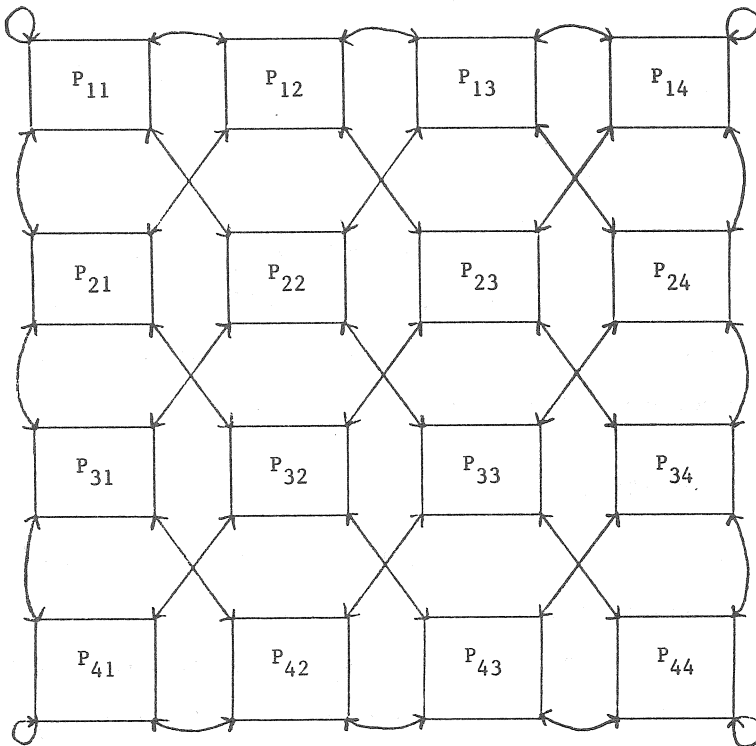


Figure 6: "Diagonal" connections, $n = 8$

(here and below \longleftrightarrow stands for \rightleftarrows)

Diagonal outputs and inputs are connected in the obvious way, as shown in Figure 6. Horizontal and vertical connections (not shown) are still required for the transmission of rotation parameters.

7.4 Taking full advantage of symmetry

Because A is symmetric and our transformations preserve symmetry, only a triangular array of $\frac{1}{2} \cdot \frac{n}{2} \cdot \frac{n}{2} + 1 = n(n+2)/8$ systolic processors is necessary for the eigenvalue computation. In the description above, simply replace any reference to a below-diagonal element a_{ij} (or processor P_{ij}) with $i > j$ by a reference to the corresponding above-diagonal element a_{ji} (or processor P_{ji}). Note, however, that this idea complicates the programs, and cannot be used if eigenvectors as well as eigenvalues are to be computed.

7.5 Odd n

So far we assumed n to be even. For odd n we can modify the program for processors P_{1i} and P_{i1} ($i = 1, \dots, \lceil \frac{n}{2} \rceil$) in a manner analogous to that used in Section 3, or simply border A by a zero row and column. For simplicity we continue to assume that n is even.

7.6 Rotation parameters

In Section 6 we assumed that the diagonal processor P_{ii} would compute c_i and s_i according to (6.1), and then broadcast both c_i and s_i along processor row and column i . It may be preferable to broadcast only t_i (given by (6.2)) and let each off-diagonal processor P_{ij} compute c_i , s_i , c_j and s_j from t_i and t_j . Thus communication costs are reduced at the expense of requiring off-diagonal processors to compute two square roots per time step (but this may not be significant since the diagonal processors must compute one or two square roots per step in any case). In what follows a "rotation parameter" may mean either t_i or the pair (c_i, s_i) .

7.7 Avoiding broadcast of rotation parameters

The most serious assumption of Section 6 is that rotation parameters computed by diagonal processors can be broadcast along rows and columns in constant time. However, it is possible to avoid this assumption, using a special case of the general technique of Leiserson and Saxe [19]. For the details, see [5]. The conclusion is that we only need to transmit rotation parameters at constant speed between adjacent processors.

7.8 Solving large problems on small systolic arrays

We have assumed that an array of $\lceil \frac{n}{2} \rceil$ by $\lceil \frac{n}{2} \rceil$ systolic processors is available. In practice the systolic array would have a fixed number of processors, and a large problem might have to be decomposed in some manner in order to fit on the available hardware. This is an interesting problem of some practical significance, but space limitations prevent us from discussing it here. For some ideas (which might be improved) on how to solve it, see [26].

8. CONCLUSION

We have presented a linear array of $\lceil \frac{n}{2} \rceil$ processors, each able to perform floating-point operations (including square roots) and with $O(m)$ local storage, for computing the SVD of a real $m \times n$ matrix in time $O(mn \log n)$, with a small constant. We have also described how a square array of $\lceil \frac{n}{2} \rceil$ by $\lceil \frac{n}{2} \rceil$ processors, each with similar arithmetical capabilities but with only $O(1)$ local storage, and having connections to nearest horizontal and vertical (and preferably also diagonal) neighbors, can compute the eigenvalues and eigenvectors of a real symmetric matrix in time $O(n \log n)$. The constant is sufficiently small that the method is competitive with the usual $O(n^3)$ serial algorithms, e.g., tridiagonalization followed by the QR iteration, for quite small n . The speedup should be significant for real-time computations with moderate or large n . For further results along these lines, see [6].

Acknowledgement

A revised and expanded version of this paper is to appear in SIAM Journal on Scientific and Statistical Computing. The work of the second author was supported in part by the U.S. Army Office under grant DAAG 29-79-CO124 and the National Science Foundation under grant MCS-8213718, and in part by the Mathematical Sciences Research Centre and the Centre for Mathematical Analysis, ANU.

REFERENCES

- [1] H.C. Andrews and C.L. Patterson, "Singular value decomposition and digital image processing", *IEEE Trans. Acoustics, Speech and Signal Processing ASSP-24* (1976), 26-53.
- [2] A. Bojańczyk, R.P. Brent and H.T. Kung, "Numerically stable solution of dense systems of linear equations using mesh-connected processors", *SIAM J. Sci. Statist. Comput.* 5 (1984), to appear.
- [3] R.P. Brent and F.T. Luk, "Computing the Cholesky factorization using a systolic architecture", *Proc. 6-th Australian Computer Science Conference* (1983), 295-302.
- [4] R.P. Brent and F.T. Luk, "A systolic architecture for the singular value decomposition", Tech. Report TR-CS-82-09, Dept. of Computer Science, Aust. Nat. Univ., August, 1982.
- [5] R.P. Brent and F.T. Luk, "A systolic architecture for almost linear-time solution of the symmetric eigenvalue problem", Tech. Report TR-CS-82-10, Dept. of Computer Science, Aust. Nat. Univ., 1982.
- [6] R.P. Brent, F.T. Luk and C. Van Loan, "Computation of the generalized singular value decomposition using mesh-connected processors", *Proceedings SPIE Volume 431, Real Time Signal Processing VI*, Society of Photo-Optical Instrumentation Engineers, Bellingham, Washington, 1983, 66-71. (Also available as Report CMA-R31-83, CMA, ANU, Aug. 1983.)
- [7] K-W. Chen and K.B. Irani, "A Jacobi algorithm and its implementation on parallel computers", *Proc. 18-th Annual Allerton Conference on Communication, Control and Computing* (1980), 564-573.
- [8] P.J. Eberlein and J. Boothroyd, "Solution to the eigenproblem by a norm reducing Jacobi type method", in [31], 327-338.
- [9] A.M. Finn, F.T. Luk and C. Pottle, "Systolic array computation of the singular value decomposition", *Proc. SPIE Symp. East 1982, Vol. 341, Real Time Signal Processing V* (1982), 35-43.

- [10] G.E. Forsythe and P. Henrici, "The cyclic Jacobi method for computing the principal values of a complex matrix", *Trans. Amer. Math. Soc.* 94 (1960), 1-23.
- [11] G.H. Golub and F.T. Luk, "Singular value decomposition: applications and computations", ARO Report 77-1, *Trans. of 22nd Conf. of Army Mathematicians* (1977), 577-605.
- [12] E.R. Hansen, "On cyclic Jacobi methods", *J. Soc. Indust. Appl. Math.* 11 (1963), 448-459.
- [13] D.E. Heller and I.C.F. Ipsen, "Systolic networks for orthogonal equivalence transformations and their applications", *Proc. 1982 Conf. on Advanced Research in VLSI*, MIT (1982), 113-122.
- [14] D.E. Heller and I.C.F. Ipsen, "Systolic networks for orthogonal decompositions", *SIAM J. Sci. Statist. Comput.* 4 (1983), 261-269.
- [15] M.R. Hestenes, "Inversion of matrices by biorthogonalization and related results", *J. Soc. Indust. Appl. Math.* 6 (1958), 51-90.
- [16] D.J. Kuck and A.H. Sameh, "Parallel computation of eigenvalues of real matrices", *Information Processing 1971*, North-Holland, Amsterdam, (1972), 1266-1272.
- [17] H.T. Kung, "Why systolic architectures", *IEEE Computer* 15, 1 (1982), 37-46.
- [18] S.Y. Kung and R.J. Gal-Ezer, "Linear or square array for eigenvalue and singular value decompositions?", *Proc. USC Workshop on VLSI and Modern Signal Processing*, Los Angeles, California (Nov. 1982), 89-98.
- [19] C.E. Leiserson and J.B. Saxe, "Optimizing synchronous systems", *J. VLSI and Computer Systems* 1 (1983), 41-67.
- [20] F.T. Luk, "Computing the singular-value decomposition on the ILLIAC IV", *ACM Trans. Math. Softw.* 6 (1980), 524-539.

- [21] H. Rutishauser, "The Jacobi method for real symmetric matrices", in [31], 202-211.
- [22] A.H. Sameh, "On Jacobi and Jacobi-like algorithms for a parallel computer", *Math. Comput.* 25 (1971), 579-590.
- [23] A.H. Sameh, "Solving the linear least squares problem on a linear array of processors," *Proc. Purdue Workshop on Algorithmically-specialized Computer Organizations* (1982).
- [24] R. Schreiber, "Systolic arrays for eigenvalue computation", *Proc. SPIE Symp. East 1982, Vol. 341, Real-Time Signal Processing* (1982).
- [25] R. Schreiber, "A systolic architecture for singular value decomposition", *Proc. 1st Intern. Coll. on Vector and Parallel Computing in Scientific Applications*, Paris, France (1983).
- [26] R. Schreiber, "On the systolic arrays of Brent, Luk and Van Loan", *Proceedings SPIE Vol. 431, Real-Time Signal Processing VI*, Society of Photo-Optical Instrumentation Engineers, Bellingham, Washington, 1983, 72-78.
- [27] J.M. Speiser and H.J. Whitehouse, "Architecture for real-time matrix operations", *Proc. 1980 Government Microcircuits Applications Conf.*, Houston, Texas (Nov. 1980).
- [28] H.J. Whitehouse, J.M. Speiser and K. Bronley, "Signal processing applications of systolic array technology", *Proc. USC Workshop on VLSI and Modern Signal Processing*, Los Angeles, California (Nov. 1982), 5-10.
- [29] J.H. Wilkinson, "Note on the quadratic convergence of the cyclic Jacobi process", *Numer. Math.* 4 (1962), 296-300.
- [30] J.H. Wilkinson, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.

- [31] J.H. Wilkinson and C. Reinsch (editors), *Handbook for Automatic Computation, Vol. 2 (Linear Algebra)*, Springer-Verlag, Berlin, 1971.

Richard P. Brent
Centre for Mathematical Analysis
Australian National University
GPO Box 4, Canberra, ACT 2601
AUSTRALIA

Franklin T. Luk
Department of Computer Science
Cornell University
ITHACA NY 14853
U.S.A.