

# Computation of the Singular Value Decomposition Using Mesh-Connected Processors

RICHARD P. BRENT,\* FRANKLIN T. LUK,†  
and CHARLES VAN LOAN‡

*Abstract*—A cyclic Jacobi method for computing the singular value decomposition of an  $m \times n$  matrix ( $m \geq n$ ) using systolic arrays is proposed. The algorithm requires  $O(n^2)$  processors and  $O(m + n \log n)$  units of time.

*Key words and phrases*: Systolic arrays, singular value decomposition, cyclic Jacobi method, real-time computation, VLSI.

## 1. INTRODUCTION

A singular value decomposition (SVD) of a matrix  $A \in \mathbf{C}^{m \times n}$  is given by

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

where  $U \in \mathbf{C}^{m \times m}$  and  $V \in \mathbf{C}^{n \times n}$  are unitary matrices and  $\Sigma \in \mathbf{R}^{m \times n}$  is a real non-negative “diagonal” matrix. Since  $A^H = V\Sigma^T U^H$ , we may assume  $m \geq n$  without loss of generality. Let

$$\begin{aligned} U &= [u_1, \dots, u_m], & \Sigma &= \text{diag}(\sigma_1, \dots, \sigma_n) \quad \text{and} \\ V &= [v_1, \dots, v_n]. \end{aligned} \quad (1.2)$$

We refer to  $\sigma_i$  as the  $i$ -th singular values of  $A$ ,  $u_i$  and  $v_i$  as the corresponding left and right singular vectors. The singular values may be arranged in any

---

\*Centre for Mathematical Analysis, Australian National University, Canberra, A.C.T. 2601, Australia.

†School of Electrical Engineering, Cornell University, Ithaca, New York 14853.

‡Department of Computer Science, Cornell University, Ithaca, New York 14853.

order, for if  $P \in \mathbf{R}^{m \times m}$  and  $Q \in \mathbf{R}^{n \times n}$  are permutation matrices such that  $P\Sigma Q$  remains "diagonal," then

$$A = (UP^T)(P\Sigma Q)(Q^TV^H)$$

is also an SVD. It is customary to choose  $P$  and  $Q$  so that the singular values are arranged in non-increasing order:

$$\sigma_1 \geq \cdots \geq \sigma_r > 0, \quad \sigma_{r+1} = \cdots = \sigma_n = 0, \quad (1.3)$$

with  $r = \text{rank}(A)$ . If  $A$  is real, then the unitary matrices  $U$  and  $V$  are real and hence orthogonal. To simplify the presentation, we shall assume that the matrix  $A$  is real. The algorithms presented in this paper can be readily extended to handle the complex case.

The standard method for computing (1.1) is the Golub-Kahan-Reinsch SVD algorithm ([9] and [11]) that is implemented in both EISPACK [7] and LINPACK [4]. It requires time  $O(mn^2)$ . However, the advent of massively parallel computer architectures has aroused much interest in parallel SVD procedures, e.g., [1], [5], [13], [16], [18], [21], and [22]. Such architectures may turn out to be indispensable in settings where real-time computation of the SVD is required (cf. [24] and [25]). Speiser and Whitehouse [24] survey parallel processing architectures and conclude that systolic architectures offer the best combination of characteristics for utilizing VLSI/VHSIC technology to do real-time signal processing (see also [15] and [25]). On a linear systolic array, the most efficient SVD algorithm is the Jacobi-like algorithm given by Brent and Luk [1]; it requires  $O(mn)$  time and  $O(n)$  processors. This is not surprising because Jacobi-type procedures are very amenable to parallel computations (cf. [1], [18], [20], and [21]).

In this paper we present a modification of the two-sided Jacobi SVD method for square matrices detailed in Forsythe and Henrici [6] and show how it can be implemented on a quadratic systolic array. The array is very similar to the one proposed in Brent and Luk [1] that uses  $O(n^2)$  processors to solve  $n$ -by- $n$  symmetric eigenproblems in  $O(n \log n)$  time. However, to handle  $m$ -by- $n$  SVD problems, our algorithm requires a "pre-processing" step in which the QR-factorization  $A = QR$  is computed. This can be done in  $O(m)$  time using the systolic array described in [8]. Our modified Brent-Luk array then computes the SVD of  $R$  in  $O(n \log n)$  time.  $O(n^2)$  processors are required for the overall algorithm. For an extension of this work, see the paper [2] by the authors on computing the generalized SVD of two given matrices.

The paper is organized as follows. In Section 2, basic concepts and calculations are discussed by reviewing the Jacobi algorithm for the symmetric eigenproblem. Our modified Forsythe-Henrici scheme is then presented in Sec-

tions 3 and 4 where we discuss two-by-two SVD problems and “parallel orderings,” respectively. In Section 5, we describe a systolic array for implementing the scheme. Section 6 discusses how this array can be used to solve rectangular SVD problems. We also describe in that final section a block Jacobi method that might be of interest when our array is too small to handle a given SVD problem.

## 2. BACKGROUND

The classical method of Jacobi uses a sequence of plane rotations to diagonalize a symmetric matrix  $A \in \mathbf{R}^{n \times n}$ . We denote a Jacobi rotation of an angle  $\theta$  in the  $(p, q)$  plane by  $J(p, q, \theta) \equiv J$ , where  $p < q$ . The matrix  $J$  is the same as the identity matrix except for four strategic elements:

$$\begin{aligned} J_{pp} &= c, & J_{pq} &= s, \\ J_{qp} &= -s, & J_{qq} &= c, \end{aligned}$$

where  $c = \cos(\theta)$  and  $s = \sin(\theta)$ . Letting  $B = J^T A J$ , we get

$$\begin{bmatrix} b_{pp} & b_{pq} \\ b_{qp} & b_{qq} \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}^T \begin{bmatrix} a_{pp} & a_{pq} \\ a_{qp} & a_{qq} \end{bmatrix} \begin{bmatrix} c & s \\ -s & c \end{bmatrix}.$$

If we choose the cosine-sine pair  $(c, s)$  such that

$$b_{pq} = b_{qp} = a_{pq}(c^2 - s^2) + (a_{pp} - a_{qq})cs = 0, \quad (2.1)$$

then  $B$  becomes “more diagonal” than  $A$  in the sense that

$$\text{off}(B) = \text{off}(A) - 2a_{pq}^2,$$

where

$$\text{off}(C) \equiv \sum_{i \neq j} c_{ij}^2 \quad \text{for } C = (c_{ij}) \in \mathbf{R}^{m \times n}.$$

Assuming that  $a_{pq} \neq 0$ , we have from equation (2.1) that

$$\rho \equiv \frac{a_{qq} - a_{pp}}{2a_{pq}} = \frac{c^2 - s^2}{2cs} = \text{ctn}(2\theta), \quad (2.2)$$

and that  $t \equiv \tan(\theta)$  satisfies

$$t^2 + 2\rho t - 1 = 0.$$

By solving this quadratic equation and using a little trigonometry, we find two possible solutions to (2.1):

$$t = -\text{sign}(\rho) [|\rho| + \sqrt{1 + \rho^2}], \quad c = \frac{1}{\sqrt{1 + t^2}}, \quad s = ct, \quad (2.3)$$

and

$$t = \frac{\text{sign}(\rho)}{|\rho| + \sqrt{1 + \rho^2}}, \quad c = \frac{1}{\sqrt{1 + t^2}}, \quad s = ct. \quad (2.4)$$

The angle  $\theta$  associated with (2.4) is the smaller of the two possible rotation angles and satisfies  $0 \leq |\theta| \leq \pi/4$ .

By systematically "zeroing" off-diagonal entries,  $A$  can be effectively diagonalized. To be specific, consider the iteration:

#### *Algorithm Jacobi*

```

do until  $\text{off}(A) < \zeta$ 
  for  $p = 1, \dots, n - 1$ 
    for  $q = p + 1, \dots, n$ 
      begin
        Determine  $c$  and  $s$  via (2.2) and (2.4);
         $A := J(p, q, \theta)^T A J(p, q, \theta)$ 
      end.

```

The parameter  $\zeta$  is some small machine-dependent number. Each pass through the "until" loop is called a "sweep." In this scheme a sweep consists of zeroing the off-diagonal elements according to the "row ordering" [6]. This ordering is amply illustrated by the  $n = 4$  case:

$$(p, q) = (1, 2), (1, 3), (1, 4), (2, 3), (2, 4), (3, 4).$$

It is well known that Algorithm Jacobi always converges [6] and that the asymptotic convergence rate is quadratic [26]. The rigorous proof of these results requires that (2.4) rather than (2.3) be used. Brent and Luk [1] conjecture that  $O(\log n)$  sweeps are required by Algorithm Jacobi for typical values

of  $\zeta$ , e.g.,  $\zeta = 10^{-12} \text{off}(A_0)$ , where  $A_0$  denotes the original matrix. For most problems this usually means about six to ten sweeps (cf. [1] and [19]).

The practical computation of  $c$  and  $s$  requires that we guard against overflow. If  $\epsilon$  denotes the machine precision, then this can usually be accomplished as follows:

*Algorithm CS*

$$\mu_1 := a_{qq} - a_{pp};$$

$$\mu_2 := 2a_{pq};$$

if  $|\mu_2| \leq \epsilon |\mu_1|$  then

begin

$$c := 1;$$

$$s := 0$$

end

else

begin

$$\rho := \frac{\mu_1}{\mu_2};$$

$$t := \frac{\text{sign}(\rho)}{|\rho| + \sqrt{1 + \rho^2}};$$

$$c := \frac{1}{\sqrt{1 + t^2}};$$

$$s := ct$$

end.

Jacobi methods for the symmetric eigenproblem are of interest because they lend themselves to parallel computation (see, e.g., [1] and [20]). Of particular interest to us is the "parallel ordering" described in [1] and illustrated in the  $n = 8$  case by

$$\begin{aligned} (p, q) = & (1, 2), (3, 4), (5, 6), (7, 8), \\ & (1, 4), (2, 6), (3, 8), (5, 7), \\ & (1, 6), (4, 8), (2, 7), (3, 5), \\ & (1, 8), (6, 7), (4, 5), (2, 3), \\ & (1, 7), (8, 5), (6, 3), (4, 2), \\ & (1, 5), (7, 3), (8, 2), (6, 4), \\ & (1, 3), (5, 2), (7, 4), (8, 6). \end{aligned}$$

Note that the rotation pairs associated with each "row" of the above can be calculated concurrently. Brent and Luk [1] have developed a systolic array that exploits this concurrency to such an extent that a sweep with the parallel ordering can be executed in  $O(n)$  time. Our aim is to extend their work to the SVD problem.

At first glance this seems unnecessary, since software (or hardware) for the symmetric eigenvalue problem can in principle be used to solve the SVD problem. For example, we may compute the eigenvalue decomposition

$$V^T(A^T A)V = \text{diag}(\sigma_1^2, \dots, \sigma_n^2),$$

where  $V = [v_1, \dots, v_n]$  is orthogonal and the  $\sigma_i$  satisfies

$$\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_n = 0,$$

with  $r = \text{rank}(A)$ . We next calculate the vectors

$$u_i = \frac{Av_i}{\sigma_i} \quad (i = 1, \dots, r), \quad (2.5)$$

and then determine the vectors  $u_{r+1}, \dots, u_m$  so that the matrix  $U = [u_1, \dots, u_m]$  is orthogonal. The factorization  $U^T A V = \text{diag}(\sigma_1, \dots, \sigma_n)$  gives the SVD of  $A$ . Thus, one can theoretically compute an SVD of  $A$  via an eigenvalue decomposition of  $A^T A$ . Unfortunately, well-known numerical difficulties are associated with the explicit formation of  $A^T A$ .

A way around this difficulty is to apply the Jacobi method implicitly. This is the gist of the "one-sided" Hestenes approach in which the matrix  $V$  is determined so that the columns of  $AV$  are mutually orthogonal [14]. Implementations are discussed in Luk [18] and in Brent and Luk [1]. In the latter reference, a systolic array is developed that is tailored to the method. However, inner products of  $m$ -vectors are required for each  $(c, s)$  computation. Because of this, the speed of their parallel algorithm is  $O(mn \log n)$  for a linear array of processors, and  $O(n \log m \log n)$  for a two-dimensional array of  $O(mn)$  processors with some special interconnection patterns for inner-product computations. Another drawback of the one-sided Jacobi method is that it does not directly generate the vectors  $u_{r+1}, \dots, u_m$ . This is an inconvenience in the systolic array setting since one would need a special architecture to carry out the matrix-vector multiplications in (2.5). It should be pointed out, however, that these vectors are not necessary for solving the least squares problem  $\|Ax - b\|_2 = \min$ .

Yet another approach to the SVD problem is to compute the eigenvalue decomposition of the  $(m + n) \times (m + n)$  symmetric matrix

$$C = \begin{bmatrix} O & A \\ A^T & 0 \end{bmatrix}.$$

Note that if

$$\begin{bmatrix} O & A \\ A^T & 0 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \sigma \begin{bmatrix} u \\ v \end{bmatrix},$$

then  $A^T A v = \sigma^2 v$  and  $A A^T u = \sigma^2 u$ . Thus, the eigenvectors of  $C$  are “made up” of the singular vectors of  $A$ . It can also be shown that the spectrum of  $C$  is given by

$$\lambda(C) = \{ \pm\sigma_1, \dots, \pm\sigma_n, 0, \dots, 0 \}.$$

The disadvantages of this approach to the SVD are that  $C$  has expanded dimension and that recovering the singular vectors may be a difficult numerical task (Golub and Kahan [9]). In addition, the case  $\text{rank}(A) < n$  requires extra work to generate  $v_{r+1}, \dots, v_n$ .

To summarize, it is preferable from several different points of view *not* to approach the SVD problem as a symmetric eigenvalue problem.

### 3. THE TWO-BY-TWO SVD PROBLEM

Forsythe and Henrici [6] extend the Jacobi eigenvalue algorithm to the computation of an SVD of a square matrix  $A \in \mathbf{R}^{n \times n}$ . They effectively diagonalize  $A$  via a sequence of two-by-two SVDs. In this section, we discuss how to compute cosine-sine pairs  $(c_1, s_1)$  and  $(c_2, s_2)$  such that

$$\begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix}^T \begin{bmatrix} w & x \\ y & z \end{bmatrix} \begin{bmatrix} c_2 & s_2 \\ -s_2 & c_2 \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}, \quad (3.1)$$

where

$$A = \begin{bmatrix} w & x \\ y & z \end{bmatrix}$$

is given. It is always possible to choose the above rotations such that  $|d_1| \geq |d_2|$ ; however, since  $\det(A) = d_1 d_2$ , it may not be possible to achieve  $d_1 \geq d_2 \geq 0$ . Consequently, we refer to (3.1) as an "unnormalized" SVD. To generate a genuine, "normalized" SVD, it may be necessary to use reflections as well as rotations. We shall return to this detail at the end of this section.

The following approach is suggested in [6] for computing (3.1). Let  $\theta_1$  and  $\theta_2$  be the angles generating the cosine-sine pairs  $(c_1, s_1)$  and  $(c_2, s_2)$ , respectively. Then  $\theta_1$  and  $\theta_2$  are solutions to the equations:

$$\tan(\theta_1 + \theta_2) = \frac{y + x}{z - w}, \quad \tan(-\theta_1 + \theta_2) = \frac{y - x}{z + w}.$$

We may use the following procedure to determine the cosine-sine pairs:

*Algorithm FHSVD*

$\mu_1 := z - w;$

$\mu_2 := y + x;$

$\left\{ \text{Find } \chi_1 = \cos\left(\frac{\theta_1 + \theta_2}{2}\right) \text{ and } \sigma_1 = \sin\left(\frac{\theta_1 + \theta_2}{2}\right) \right\}$

if  $|\mu_2| \leq \epsilon |\mu_1|$  then

begin

$\chi_1 := 1;$

$\sigma_1 := 0$

end

else

begin

$\rho_1 := \frac{\mu_1}{\mu_2};$

$\tau_1 := \frac{\text{sign}(\rho_1)}{|\rho_1| + \sqrt{1 + \rho_1^2}};$

$\chi_1 := \frac{1}{\sqrt{1 + \tau_1^2}};$

$\sigma_1 := \chi_1 \tau_1$

end;

$$\mu_1 := z + w;$$

$$\mu_2 := y - x;$$

$$\left\{ \text{Find } \chi_2 = \cos\left(\frac{-\theta_1 + \theta_2}{2}\right) \text{ and } \sigma_2 = \sin\left(\frac{-\theta_1 + \theta_2}{2}\right) \right\}$$

if  $|\mu_2| \leq \epsilon |\mu_1|$  then

begin

$$\chi_2 := 1;$$

$$\sigma_2 := 0$$

end

else

begin

$$\rho_2 := \frac{\mu_1}{\mu_2};$$

$$\tau_2 := \frac{\text{sign}(\rho_2)}{|\rho_2| + \sqrt{1 + \rho_2^2}};$$

$$\chi_2 := \frac{1}{\sqrt{1 + \tau_2^2}};$$

$$\sigma_2 := \chi_2 \tau_2$$

end;

{ Find  $(c_1, s_1)$  and  $(c_2, s_2)$  }

$$c_1 := \chi_1 \chi_2 + \sigma_1 \sigma_2;$$

$$s_1 := \sigma_1 \chi_2 - \chi_1 \sigma_2;$$

$$c_2 := \chi_1 \chi_2 - \sigma_1 \sigma_2;$$

$$s_2 := \sigma_1 \chi_2 + \chi_1 \sigma_2.$$

An alternative method for computing (3.1) is to first symmetrize  $A$ :

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix}^T \begin{bmatrix} w & x \\ y & z \end{bmatrix} = \begin{bmatrix} p & q \\ q & r \end{bmatrix}, \quad (3.2)$$

and then diagonalize the result:

$$\begin{bmatrix} c_2 & s_2 \\ -s_2 & c_2 \end{bmatrix}^T \begin{bmatrix} p & q \\ q & r \end{bmatrix} \begin{bmatrix} c_2 & s_2 \\ -s_2 & c_2 \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}.$$

Equation (3.1) holds by setting

$$\begin{bmatrix} c_2 & s_2 \\ -s_2 & c_2 \end{bmatrix}^T \begin{bmatrix} c & s \\ -s & c \end{bmatrix}^T = \begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix}^T, \quad (3.3)$$

that is,

$$c_1 = c_2c - s_2s,$$

$$s_1 = s_2c + c_2s.$$

The formulas for  $c$  and  $s$  are easily derived. From (3.2) we have

$$cx - sz = sw + cy.$$

Thus, if  $x \neq y$ , then we have

$$\rho \equiv \text{ctn}(\theta) = \frac{w + z}{x - y}, \quad s = \frac{\text{sign}(\rho)}{\sqrt{1 + \rho^2}}, \quad c = s\rho.$$

For reasons that will be explained in Section 6, it is necessary that our two-by-two SVD algorithm handles the special cases  $x = z = 0$  and  $y = z = 0$  as follows:

$$\begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix}^T \begin{bmatrix} w & 0 \\ y & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & 0 \end{bmatrix}, \quad (3.4)$$

and

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^T \begin{bmatrix} w & x \\ 0 & 0 \end{bmatrix} \begin{bmatrix} c_2 & s_2 \\ -s_2 & c_2 \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & 0 \end{bmatrix}. \quad (3.5)$$

In other words,  $\theta_2$  should be zero if  $x = z = 0$ , and  $\theta_1$  should be zero if  $y = z = 0$ . Our scheme as described already produces (3.4). However, (3.5) is only achieved if  $|w| > |x|$ . To rectify this, we apply our algorithm to the transpose problem and obtain

$$\begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix}^T \begin{bmatrix} w & 0 \\ x & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & 0 \end{bmatrix}.$$

We then set  $c_2 = c_1$ ,  $s_2 = s_1$ ,  $c_1 = 1$ , and  $s_1 = 0$ . Let us present the overall scheme:

*Algorithm USVD*

```

flag := 0;
if y = 0 and z = 0 then
  begin
    y := x;    x := 0;    flag := 1
  end;

 $\mu_1 := w + z;$ 
 $\mu_2 := x - y;$ 

if  $|\mu_2| \leq \epsilon |\mu_1|$  then
  begin
    c := 1;    s := 0
  end
else
  begin
     $\rho := \frac{\mu_1}{\mu_2};$ 

     $s := \frac{\text{sign}(\rho)}{\sqrt{1 + \rho^2}};$ 
    c := s  $\rho$ 
  end;

 $\mu_1 := s(x + y) + c(z - w);$     { = r - p }
 $\mu_2 := 2(cx - sz);$             { = 2q }

if  $|\mu_2| \leq \epsilon |\mu_1|$  then
  begin
    c2 := 1;    s2 := 0
  end
else
  begin
     $\rho_2 := \frac{\mu_1}{\mu_2};$ 

     $t_2 := \frac{\text{sign}(\rho_2)}{|\rho_2| + \sqrt{1 + \rho_2^2}};$ 

```

$$c_2 := \frac{1}{\sqrt{1 + t_2^2}};$$

$$s_2 := c_2 t_2$$

end;

$$c_1 := c_2 c - s_2 s;$$

$$s_1 := s_2 c + c_2 s;$$

$$d_1 := c_1(wc_2 - xs_2) - s_1(yc_2 - zs_2);$$

$$d_2 := s_1(ws_2 + xc_2) + c_1(ys_2 + zc_2);$$

if flag = 1 then  
 begin  
 $c_2 := c_1; \quad s_2 := s_1;$   
 $c_1 := 1; \quad s_1 := 0$   
 end.

If a "normalized" SVD is required, then further computations must be performed. The diagonal elements must be sorted by modulus and then (if necessary) premultiplied by -1, as shown in the following algorithm.

*Algorithm NSVD*

Apply Algorithm USVD;

if  $|d_2| > |d_1|$  then  
 begin  
 $\tau := c_1; \quad c_1 := -s_1; \quad s_1 := \tau;$   
 $\tau := c_2; \quad c_2 := -s_2; \quad s_2 := \tau;$   
 $\tau := d_1; \quad d_1 := d_2; \quad d_2 := \tau$   
 end;  
 $\kappa := 1;$   
 if  $d_1 < 0$  then  
 begin  
 $d_1 := -d_1;$   
 $c_1 := -c_1;$   
 $s_1 := -s_1;$   
 $\kappa := -\kappa$   
 end;  
 if  $d_2 < 0$  then  
 begin  
 $d_2 := -d_2;$   
 $\kappa := -\kappa$   
 end.

The “normalized” SVD is given by

$$\begin{bmatrix} c_1 & s_1\kappa \\ -s_1 & c_1\kappa \end{bmatrix}^T \begin{bmatrix} w & x \\ y & z \end{bmatrix} \begin{bmatrix} c_2 & s_2 \\ -s_2 & c_2 \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}. \quad (3.6)$$

Note that if  $\kappa = -1$  then the left transformation is a “reflection,” i.e., a two-by-two orthogonal matrix of the form

$$\begin{bmatrix} \cos(\theta_1) & -\sin(\theta_1) \\ -\sin(\theta_1) & -\cos(\theta_1) \end{bmatrix}.$$

We denote the corresponding  $n$ -by- $n$  orthogonal transformation by  $\hat{J}(p, q, \theta_1, \kappa) \equiv \hat{J}$ , where  $p < q$ . The matrix  $\hat{J}$  is the same as the identity except for the four elements:

$$\begin{aligned} \hat{J}_{pp} &= \cos(\theta_1), & \hat{J}_{pq} &= \kappa \sin(\theta_1), \\ \hat{J}_{qp} &= -\sin(\theta_1), & \hat{J}_{qq} &= \kappa \cos(\theta_1), \end{aligned}$$

and it is a reflection whenever  $\kappa = -1$ .

#### 4. SOME TWO-SIDED JACOBI SVD PROCEDURES

By solving an appropriate sequence of two-by-two SVD problems, we may compute the SVD of a general  $n \times n$  matrix  $A$ . Analogous to Algorithm Jacobi for the symmetric eigenvalue problem, we propose the following routine:

##### *Algorithm JSVD*

```

do until  $off(A) < \zeta$ 
  for each  $(p, q)$  in accordance with the chosen ordering
    begin
      Set  $w = a_{pp}$ ,  $x = a_{pq}$ ,
           $y = a_{qp}$ ,  $z = a_{qq}$ ;

      Compute  $(c_1, s_1)$ ,  $(c_2, s_2)$  and  $\kappa \in \{-1, 1\}$  so that
        the  $(p, q)$  and  $(q, p)$  entries
        of  $\hat{J}(p, q, \theta_1, \kappa)^T A \hat{J}(p, q, \theta_2, \kappa)$  are zero;

       $A := \hat{J}(p, q, \theta_1, \kappa)^T A \hat{J}(p, q, \theta_2, \kappa)$ 
    end.

```

We have experimentally compared the three-angle formulas FHSVD, USVD, and NSVD. The convergence of Algorithm JSVD in conjunction with any one of these formulas and with any ordering has not been rigorously proved. However, there are some important theoretical observations that can be made.

Let  $\theta_1$  and  $\theta_2$  be the angles generating the cosine-sine pairs  $(c_1, s_1)$  and  $(c_2, s_2)$ , respectively, and suppose that  $-b \leq \theta_1, \theta_2 \leq b$ . Although Forsythe and Henrici [6] prove that JSVD with the "row ordering" always converges if  $b < \pi/2$ , they also demonstrate that this condition may fail to hold. As a remedy, they suggest an under- or over-rotation variant of FHSVD and prove its convergence. However, we do not favor this variant, believing that it will eliminate the empirically observed quadratic convergence of Algorithm JSVD. One can show that the bound  $b$  equals (i)  $\pi/2$  for FHSVD, (ii)  $3\pi/4$  for USVD, and (iii)  $5\pi/4$  for NSVD even if no reflection is involved. We conjecture that the smaller the rotation angles are the faster the procedure JSVD will converge. Our experimental results agree with the conjecture (see Table 1).

We also point out that the computed diagonal matrix in JSVD may not have sorted diagonal entries even if the "normalized" SVD of each two-by-two submatrix is calculated. To obtain the "normalized" SVD of  $A$  we must sort the diagonal entries and permute the columns of  $U$  and  $V$  accordingly.

It is quite obvious that the "row ordering" is not amenable to parallel processing. A new "parallel ordering" has been introduced in [1] for doing  $\lfloor n/2 \rfloor$  rotations simultaneously. For the symmetric eigenvalue problem, the latter ordering has been found to be superior to the former both empirically [1] and theoretically (see [1] and [12]), even on a single-processor machine. Our consideration of systolic array implementation necessitates the use of the "parallel ordering" in this paper.

**Table 1. Average and maximum number (in parentheses) of sweeps using the "parallel ordering"**

$n$	Trials	FHSVD	USVD	NSVD
4	1000	2.97(3.67)	2.97(4.00)	2.97(4.33)
6	1000	3.73(4.80)	3.76(4.87)	4.17(5.40)
8	1000	4.19(5.07)	4.21(5.14)	4.76(5.96)
10	1000	4.51(5.38)	4.55(5.44)	5.18(6.40)
20	100	5.50(6.13)	5.54(6.01)	6.44(7.44)
30	100	6.03(6.52)	6.09(6.80)	7.30(7.98)
40	100	6.36(6.98)	6.40(6.98)	7.94(8.50)
50	100	6.66(7.11)	6.72(7.34)	8.51(9.00)

We have applied each of the three-angle formulas with the “parallel ordering” to random  $n \times n$  matrices  $A$ , whose elements were uniformly and independently distributed in  $[-1, 1]$ . The stopping criterion was that  $\text{off}(A)$  was reduced to  $10^{-12}$  times its original value. Table 1 gives our simulation results. We conclude that our “unnormalized” SVD method USVD performs just as well as the more complicated Forsythe-Henrici scheme FHSVD and marginally faster than the “normalized” SVD method NSVD.

We have also compared the “row” and “parallel” orderings for the method USVD. The test data and convergence criterion are as described in the last paragraph. Our simulation results are presented in Table 2. The convergence rates for the two orderings appear to be about the same.

**Table 2. Average and maximum number (in parentheses) of sweeps using the USVD method**

$n$	Trials	Row ordering	Parallel ordering
4	1000	3.22(4.17)	2.97(4.00)
6	1000	4.00(5.13)	3.76(4.87)
8	1000	4.52(5.50)	4.21(5.14)
10	1000	4.83(6.02)	4.55(5.44)
20	100	5.81(6.59)	5.54(6.01)
30	100	6.26(6.79)	6.09(6.80)
40	100	6.60(7.00)	6.40(6.98)
50	100	6.78(7.62)	6.72(7.34)
80	30	7.31(7.81)	7.30(7.79)
100	10	7.46(7.79)	7.56(8.00)
120	5	7.69(7.90)	7.73(7.98)
150	3	7.76(7.80)	7.73(8.03)
170	2	7.86(7.92)	8.02(8.02)
200	1	7.90(7.90)	8.10(8.10)
230	1	8.35(8.35)	8.43(8.43)

In the remainder of the paper, we will focus our attention exclusively on the USVD method for solving two-by-two SVD problems.

## 5. A SYSTOLIC ARRAY

We now describe a systolic array for implementing the Jacobi SVD method JSVD for an  $n \times n$  real matrix  $A$ . Our array is very similar to the one detailed in [1]. For pedagogic purposes, we first idealize the array so that it has the ability to broadcast the rotation parameters in constant time.

Assume that  $n$  is even and that we have a square array of  $n/2$  by  $n/2$  processors, each containing a  $2 \times 2$  submatrix of  $A$ . Initially, processor  $P_{ij}$  contains

$$\begin{bmatrix} a_{2i-1,2j-1} & a_{2i-1,2j} \\ a_{2i,2j-1} & a_{2i,2j} \end{bmatrix}$$

( $i, j = 1, \dots, n/2$ ). Processor  $P_{ij}$  is connected to its "diagonally" nearest neighbors  $P_{i\pm 1, j\pm 1}$  ( $1 < i, j < n/2$ ). For the boundary processors and for a complete picture, see Figures 1 and 2. In general,  $P_{ij}$  contains four real numbers

$$\begin{bmatrix} \alpha_{ij} & \beta_{ij} \\ \gamma_{ij} & \delta_{ij} \end{bmatrix}.$$

The diagonal processors  $P_{ii}$  ( $i = 1, \dots, n/2$ ) act differently from the off-diagonal processors  $P_{ij}$  ( $i \neq j, 1 \leq i, j \leq n/2$ ). At each time step, the diagonal processors  $P_{ii}$  compute rotation pairs  $(c_i^L, s_i^L)$  and  $(c_i^R, s_i^R)$  to annihilate their off-diagonal elements  $\beta_{ii}$  and  $\gamma_{ii}$  (cf. Algorithm USVD):

$$\begin{bmatrix} c_i^L & s_i^L \\ -s_i^L & c_i^L \end{bmatrix}^T \begin{bmatrix} \alpha_{ii} & \beta_{ii} \\ \gamma_{ii} & \delta_{ii} \end{bmatrix} \begin{bmatrix} c_i^R & s_i^R \\ -s_i^R & c_i^R \end{bmatrix} = \begin{bmatrix} \alpha'_{ii} & 0 \\ 0 & \delta'_{ii} \end{bmatrix}.$$

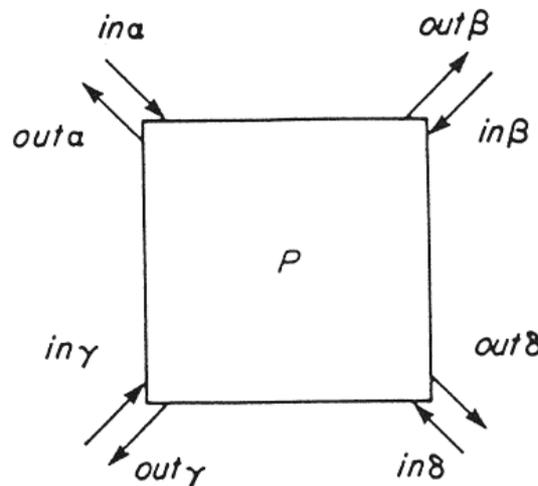


Figure 1. Diagonal input and output lines for processors.

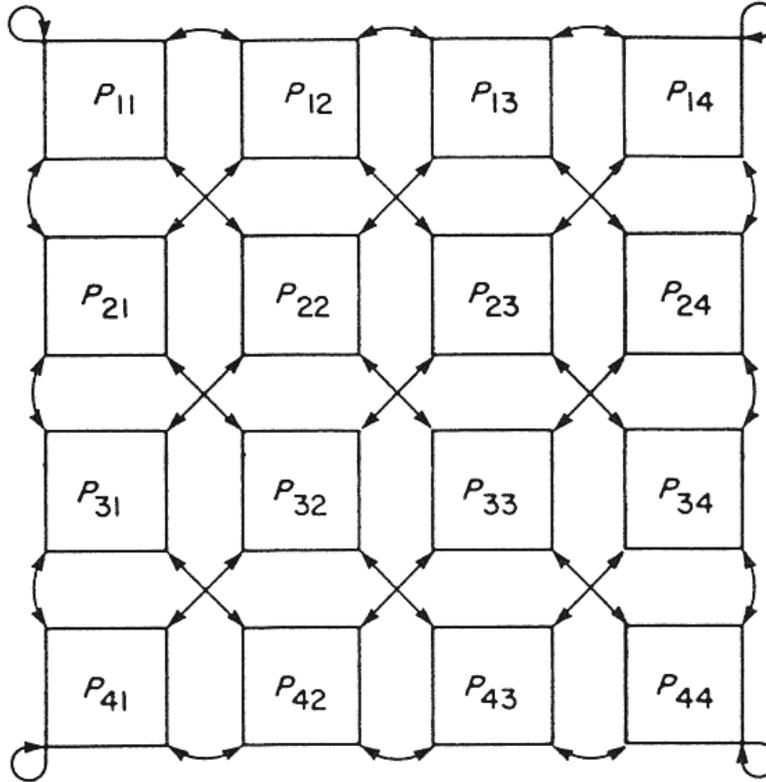


Figure 2. "Diagonal" connections,  $n = 8$ . (here  $\leftrightarrow$  stands for  $\rightleftarrows$ )

To complete the rotations that annihilate  $\beta_{ii}$  and  $\gamma_{ii}$  ( $i = 1, \dots, n/2$ ), the off-diagonal processors  $P_{ij}$  ( $i \neq j$ ) must perform the transformations

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

where

$$\begin{bmatrix} \alpha'_{ij} & \beta'_{ij} \\ \gamma'_{ij} & \delta'_{ij} \end{bmatrix} = \begin{bmatrix} c_i^L & s_i^L \\ -s_i^L & c_i^L \end{bmatrix}^T \begin{bmatrix} \alpha_{ij} & \beta_{ij} \\ \gamma_{ij} & \delta_{ij} \end{bmatrix} \begin{bmatrix} c_j^R & s_j^R \\ -s_j^R & c_j^R \end{bmatrix}.$$

We assume that the diagonal processor  $P_{ii}$  broadcasts the rotation parameters  $(c_i^L, s_i^L)$  to other processors on the  $i$ -th row, and  $(c_i^R, s_i^R)$  to other processors on the  $i$ -th column in constant time, so that the off-diagonal processor  $P_{kj}$  has access to the parameters  $(c_k^L, s_k^L)$  and  $(c_j^R, s_j^R)$  when required.

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. We

have assumed diagonal connections for convenience. They can easily be simulated if only horizontal and vertical connections are available (cf. [1]). The diagonal outputs and inputs are illustrated in Figure 1, where the subscripts  $(i, j)$  are omitted. They are connected in the obvious way, as shown in Figure 2. For example,

$$\text{out } \beta_{ij} \text{ is connected to } \begin{cases} \text{in } \gamma_{i-1, j+1} & \text{if } i > 1, \quad j < n/2 \\ \text{in } \alpha_{i, j+1} & \text{if } i = 1, \quad j < n/2 \\ \text{in } \delta_{i-1, j} & \text{if } i > 1, \quad j = n/2 \\ \text{in } \beta_{i, j} & \text{if } i = 1, \quad j = n/2 \end{cases}$$

Formally, we can specify the diagonal interchanges performed by processor  $P_{ij}$  as follows.

#### *Algorithm Interchange*

$$\text{if } i = 1 \text{ and } j = 1 \text{ then } \begin{bmatrix} \text{out } \alpha \leftarrow \alpha; & \text{out } \beta \leftarrow \beta; \\ \text{out } \gamma \leftarrow \gamma; & \text{out } \delta \leftarrow \delta; \end{bmatrix}$$

$$\text{else if } i = 1 \text{ then } \begin{bmatrix} \text{out } \alpha \leftarrow \beta; & \text{out } \beta \leftarrow \alpha; \\ \text{out } \gamma \leftarrow \delta; & \text{out } \delta \leftarrow \gamma; \end{bmatrix}$$

$$\text{else if } j = 1 \text{ then } \begin{bmatrix} \text{out } \alpha \leftarrow \gamma; & \text{out } \beta \leftarrow \delta; \\ \text{out } \gamma \leftarrow \alpha; & \text{out } \delta \leftarrow \beta; \end{bmatrix}$$

$$\text{else } \begin{bmatrix} \text{out } \alpha \leftarrow \delta; & \text{out } \beta \leftarrow \gamma; \\ \text{out } \gamma \leftarrow \beta; & \text{out } \delta \leftarrow \alpha; \end{bmatrix}$$

{ wait for outputs to propagate to inputs of adjacent processors }

$$\begin{bmatrix} \text{in } \alpha \leftarrow \alpha; & \text{in } \beta \leftarrow \beta; \\ \text{in } \gamma \leftarrow \gamma; & \text{in } \delta \leftarrow \delta. \end{bmatrix}$$

It is clear that a diagonal processor  $P_{ii}$  might omit rotations if its off-diagonal elements  $\beta_{ii}$  and  $\gamma_{ii}$  were sufficiently small. All that is required is to broadcast  $(c_i^L, s_i^L) = (1, 0)$  and  $(c_i^R, s_i^R) = (1, 0)$  along processor row and column  $i$ , respectively. As is well known (cf. [1]), a suitable threshold strategy guaran-

tees convergence, although we do not know any example for which our algorithm fails to give convergence even without a threshold strategy.

If singular vectors are required, the matrices  $U$  and  $V$  of singular vectors can be accumulated at the same time that  $A$  is being diagonalized. To compute  $U$ , each systolic processor  $P_{ij}$  ( $1 \leq i, j \leq n/2$ ) needs four additional memory cells

$$\begin{bmatrix} \mu_{ij} & \nu_{ij} \\ \sigma_{ij} & \tau_{ij} \end{bmatrix}.$$

During each step it updates

$$\begin{bmatrix} \mu_{ij} & \nu_{ij} \\ \sigma_{ij} & \tau_{ij} \end{bmatrix} \leftarrow \begin{bmatrix} \mu_{ij} & \nu_{ij} \\ \sigma_{ij} & \tau_{ij} \end{bmatrix} \begin{bmatrix} c_j^L & s_j^L \\ -s_j^L & c_j^L \end{bmatrix}.$$

Each processor transmits its

$$\begin{bmatrix} \mu & \nu \\ \sigma & \tau \end{bmatrix}$$

values to its adjacent processors in the same way as the

$$\begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix}$$

values (see Algorithm Interchange). Initially we set  $\mu_{ij} = \nu_{ij} = \sigma_{ij} = \tau_{ij} = 0$  ( $i \neq j$ ), and  $\mu_{ii} = \tau_{ii} = 1$ ,  $\sigma_{ii} = \nu_{ii} = 0$ . After a sufficiently large and integral number of sweeps, we have  $U$  defined by

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

The matrix  $V$  can be computed in an identical manner, but the rotation

$$\begin{bmatrix} c_j^L & s_j^L \\ -s_j^L & c_j^L \end{bmatrix} \text{ is replaced by } \begin{bmatrix} c_j^R & s_j^R \\ -s_j^R & c_j^R \end{bmatrix}.$$

If the dimension  $n$  is odd, we either modify the algorithm for processors  $P_{1i}$  and  $P_{i1}$  ( $i = 1, \dots, \lfloor n/2 \rfloor$ ) in a manner analogous to that used in [1], or simply border  $A$  by a zero row and column. The details of handling problems where  $n$  is larger or smaller than the effective dimension of the array are presented in the next section.

Instead of broadcasting the parameter pairs  $(c_i^L, s_i^L)$  and  $(c_i^R, s_i^R)$ , it may be preferable to broadcast only the corresponding tangent values  $t_i^L$  and  $t_i^R$  and let each off-diagonal processor  $P_{ij}$  compute  $(c_i^L, s_i^L)$  and  $(c_j^R, s_j^R)$  from  $t_i^L$  and  $t_j^R$  (cf. (2.4)). Communication costs are thus reduced at the expense of requiring off-diagonal processors to compute two square roots per time step. This is not significant since the diagonal processors must compute three square roots per step in any case. In what follows a "rotation parameter" may mean either  $t_i$  or the pair  $(c_i, s_i)$  (superscripts omitted).

Let us show how we may avoid the idealization that the rotation parameters are broadcast along processor rows and columns in constant time. We can retain time  $O(n)$  per sweep for our algorithm by merely transmitting rotation parameters at constant speed between adjacent processors.

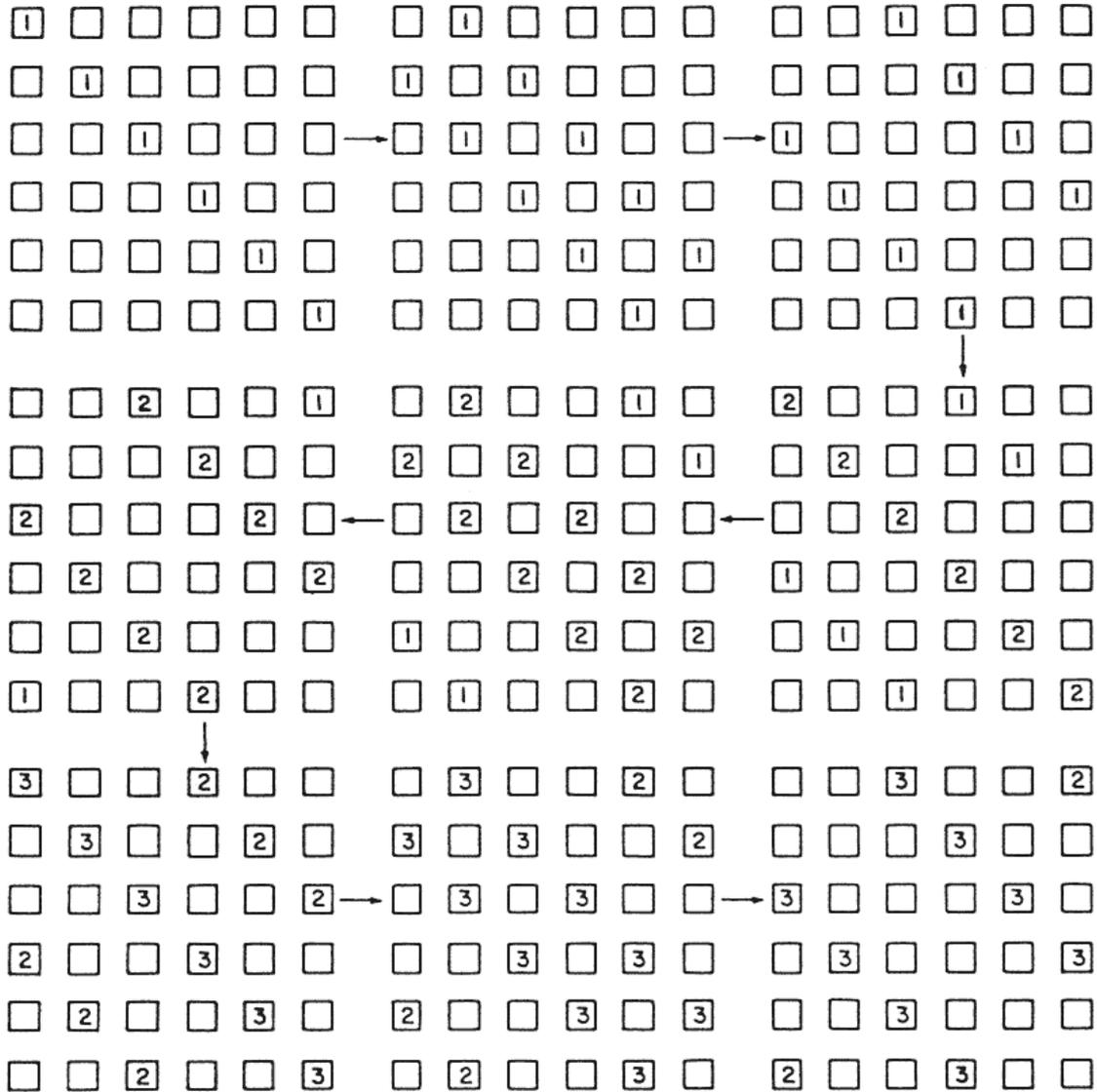
Let  $\Delta_{ij} = |i - j|$  denote the distance of processor  $P_{ij}$  from the diagonal. The operation of processor  $P_{ij}$  will be delayed by  $\Delta_{ij}$  time units relative to the operation of the diagonal processors, in order to allow time for the rotation parameters to be propagated at unit speed along each row and column of the processor array.

A processor cannot commence a rotation until data from earlier rotations are available on all its input lines. Thus, processor  $P_{ij}$  needs data from its four neighbors  $P_{i\pm 1, j\pm 1}$  ( $1 < i, j < n/2$ ). For the other cases, see Figure 2. Since

$$|\Delta_{ij} - \Delta_{i\pm 1, j\pm 1}| \leq 2,$$

it is sufficient for processor  $P_{ij}$  to be idle for two time steps while waiting for the processors  $P_{i\pm 1, j\pm 1}$  to complete their (possibly delayed) steps. Thus, the price paid to avoid broadcasting is that each processor is active for only one-third of the total computation. This is a special case of the general technique of Leiserson and Saxe [17], and is illustrated in Figure 3. A similar inefficiency occurs with many other systolic algorithms, see, e.g., [1] and [15]. The fraction can be increased to almost unity if the rotation parameters are propagated at greater than unit speed, or if multiple problems are interleaved.

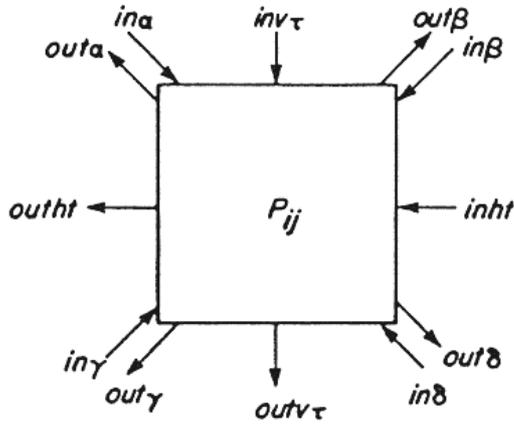
A typical processor  $P_{ij}$  ( $1 < j \leq i < n/2$ ) has input and output lines as shown in Figure 4 (with subscripts  $(i, j)$  or  $(i, i)$  omitted). Figure 4 differs from Figure 1 in that it shows the horizontal and vertical lines *inht*, *outht*, *invt*, *outvt* for transmitting the rotation parameters. Processors interconnect as shown in Figure 5.



**Figure 3.** An example ( $n = 12$ ) showing how the computations are staggered to avoid broadcasting. The value inside each box indicates the iteration number.

Suppose that the matrix  $A$  is available in the systolic array at time  $T = 0$ . Then the operation of processor  $P_{ij}$  proceeds as described in Algorithm Processor. We also assume that each time step has non-overlapping read and write phases; the result of a write at step  $T$  should be available at the read phase of step  $T + 1$ ,  $T + 2$ , and  $T + 3$  in a neighboring processor, but should not interfere with a read at step  $T$  in a neighboring processor. The first time steps at which data are available on various processors' input lines are indicated in Figure 5.

Subdiagonal ( $1 < j < i < n/2$ )



Diagonal ( $1 < i < n/2$ )

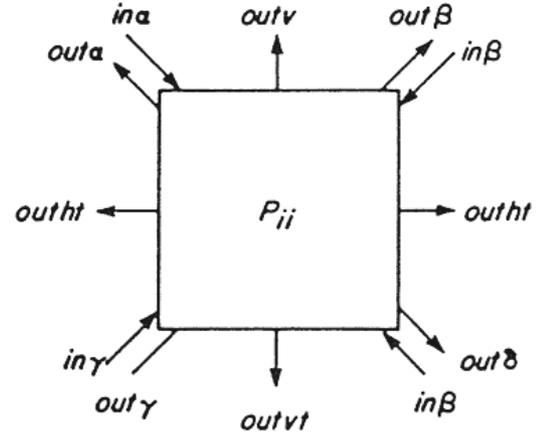


Figure 4. Input and output lines for typical subdiagonal and diagonal processors.

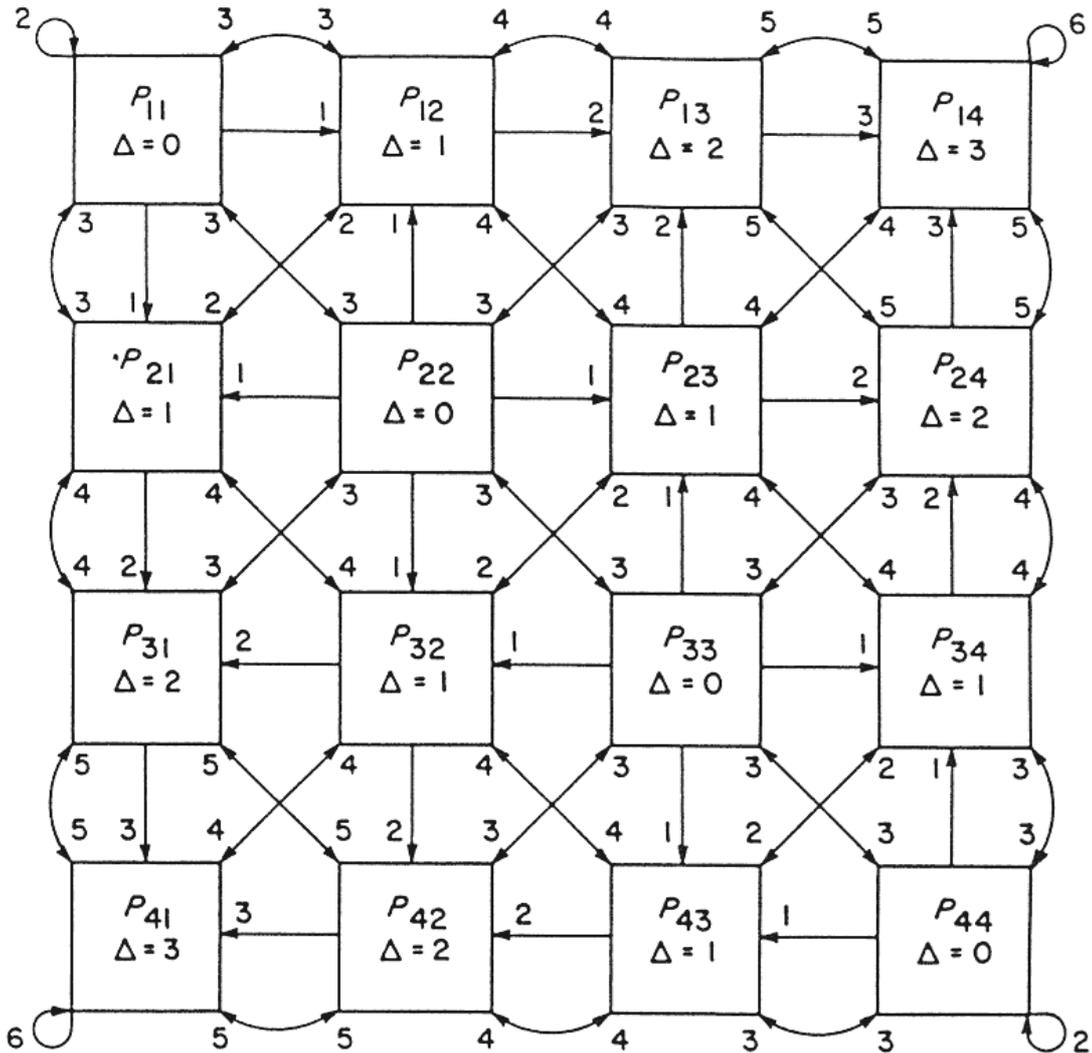


Figure 5. Interprocessor connections ( $n = 8$ ). (The first times at which inputs are available are indicated.)

*Algorithm Processor*

if  $(T \geq \Delta)$  and  $(T - \Delta \equiv 0 \pmod{3})$  then

begin

$$\text{if } T \neq \Delta \text{ then } \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \leftarrow \begin{bmatrix} in \alpha & in \beta \\ in \gamma & in \delta \end{bmatrix};$$

if  $\Delta = 0$  then { diagonal processor }

$$\text{Use (3.5) to determine } t^L, t^R \text{ and } \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \leftarrow \begin{bmatrix} \alpha' & 0 \\ 0 & \delta' \end{bmatrix}$$

else { off-diagonal processor }

begin

$$t^L \leftarrow inht; \quad t^R \leftarrow invt;$$

Use (2.4) to recover  $(c^L, s^L)$  and  $(c^R, s^R)$ ;

$$\begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \leftarrow \begin{bmatrix} c^L & s^L \\ -s^L & c^L \end{bmatrix}^T \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \begin{bmatrix} c^R & s^R \\ -s^R & c^R \end{bmatrix}$$

end;

$$outht \leftarrow t^L; \quad outvt \leftarrow t^R;$$

if  $i > j$  then set *out*  $\beta$  as in Algorithm Interchange;

if  $i < j$  then set *out*  $\gamma$  as in Algorithm Interchange

end

else if  $(T \geq \Delta)$  and  $(T - \Delta \equiv 1 \pmod{3})$  then

begin

if  $(i = 1)$  or  $(j = 1)$  then set *out*  $\alpha$  as in Algorithm Interchange;

if  $(i = n/2)$  or  $(j = n/2)$  then set *out*  $\delta$  as in Algorithm Interchange

end

else if  $(T \geq \Delta)$  and  $(T - \Delta \equiv 2 \pmod{3})$  then

begin

if  $(i > 1)$  and  $(j > 1)$  then set *out*  $\alpha$  as in Algorithm Interchange;

if  $i \leq j$  then set *out*  $\beta$  as in Algorithm Interchange;

if  $i \geq j$  then set *out*  $\gamma$  as in Algorithm Interchange;

if  $(i < n/2)$  and  $(j < n/2)$  then set *out*  $\delta$  as in Algorithm Interchange

end

else

Do nothing this time step.

Algorithm Processor does not compute singular vectors, but may easily be modified to do so. We have also omitted a termination criterion. The simplest is to perform a fixed number  $S$  (say conservatively, 10) sweeps; then processor  $P_{ij}$  halts when  $T = 3S(n - 1) + \Delta_{ij} + 3$ , since a sweep takes  $3(n - 1)$  time steps. A more sophisticated criterion is to stop when no nontrivial rotations were performed during the previous sweep. This requires communication along the diagonal, that can be done in  $n/2$  time steps.

## 6. HANDLING RECTANGULAR, UNDERSIZED, AND OVERSIZED PROBLEMS

In the previous two sections we assume that the matrix  $A$  is square. If  $A$  has more rows than columns, then one of two strategies can be adopted.

One approach is to apply our "square" algorithm to the matrix

$$\bar{A} = [A, 0] = (\bar{a}_{ij}) \in \mathbf{R}^{m \times m}.$$

Note that in the SVD problem:

$$\begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix}^T \begin{bmatrix} \bar{a}_{pp} & \bar{a}_{pq} \\ \bar{a}_{qp} & \bar{a}_{qq} \end{bmatrix} \begin{bmatrix} c_2 & s_2 \\ -s_2 & c_2 \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix},$$

we will have

$$\begin{bmatrix} c_2 & s_2 \\ -s_2 & c_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

whenever  $q > n$ , because the symmetrization step (3.2) will diagonalize the two-by-two matrix (see equation (3.4)). Thus, the sparsity structure of  $\bar{A}$  will be preserved during the iteration, and we will emerge with the factorization:

$$U^T [A \quad 0] \begin{bmatrix} V & 0 \\ 0 & I \end{bmatrix} = \text{diag}(\sigma_1, \dots, \sigma_n, 0, \dots, 0),$$

from which we get

$$U^T A V = \text{diag}(\sigma_1, \dots, \sigma_n)$$

as the desired SVD. A drawback with this approach is the expanded dimension of  $\bar{A}$ , especially painful if  $m \gg n$ . On the other hand, no additional

hardware is required so long as the array of Section 5 can handle  $m$ -by- $m$  problems.

In contrast, the second approach to handling rectangular SVD problems is more efficient but entails the interfacing of an SVD array with a QR array. The idea is to compute first the QR factorization of  $A$ :

$$A = Q \begin{bmatrix} R \\ 0 \end{bmatrix},$$

where  $R \in \mathbf{R}^{n \times n}$  is upper triangular, and then compute the square SVD:

$$W^T R V = \Sigma = \text{diag}(\sigma_1, \dots, \sigma_n).$$

Defining

$$U = Q \begin{bmatrix} W & 0 \\ 0 & I \end{bmatrix},$$

we get

$$U^T A V = \text{diag}(\sigma_1, \dots, \sigma_n).$$

Actually, in the conventional one-processor setting, it is advisable to compute the QR factorization first anyway, especially if  $m \gg n$ . See Chan [3].

We have not looked into the details of the interface between the QR and the SVD arrays. It would be preferable to have a single array with enough generality to carry out both phases of the computation. For example, perhaps the QR factorization could be obtained as  $A$  is "loaded" into the combination QR-SVD array.

We conclude with some remarks about the handling of SVD problems whose dimension differs from the effective dimension of the systolic array of Section 5. To fix the discussion, suppose that  $A$  is an  $n$ -by- $n$  matrix whose SVD we want and that our array can handle SVD problems with maximum dimension  $N$ .

If  $n < N$ , then it is natural to have the array compute the SVD of

$$\bar{A} = \begin{bmatrix} A & 0 \\ 0 & 0 \end{bmatrix},$$

i.e.,

$$\bar{U}^T \bar{A} \bar{V} = \text{diag}(\sigma_1, \dots, \sigma_n, 0, \dots, 0).$$

Our method of computing this SVD (based on JSVD) ensures that

$$\bar{U} = \begin{bmatrix} U & 0 \\ 0 & I \end{bmatrix} \quad \text{and} \quad \bar{V} = \begin{bmatrix} V & 0 \\ 0 & I \end{bmatrix},$$

whence  $U^T A V = \text{diag}(\sigma_1, \dots, \sigma_n)$ . This follows because of (3.4) and (3.5). (Whenever a two-by-two SVD involves an index greater than  $n$ , the special form of the orthogonal update guarantees that  $A$ 's block structure is preserved.) Let us point out that an SVD procedure need not produce a  $\bar{U}$  and a  $\bar{V}$  matrix with the above block structure in the case  $\text{rank}(\bar{A}) = \text{rank}(A) < n$ . For example, if  $N = 3$  and  $A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ , then one of the infinitely many SVD's of  $\bar{A}$  is

$$\begin{bmatrix} p & p^2 & p^2 \\ p & -p^2 & -p^2 \\ 0 & -p & p \end{bmatrix}^T \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} p & p^2 & p^2 \\ p & -p^2 & -p^2 \\ 0 & -p & p \end{bmatrix} = \begin{bmatrix} \sqrt{2} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

where  $p = 1/\sqrt{2}$ . Thus, further computations are necessary before the SVD of  $A$  can be obtained.

Lastly, we discuss how oversized SVD problems might be handled. Partition the matrix  $A \in \mathbf{R}^{n \times n}$  so that

$$A = \begin{bmatrix} A_{11} & \cdots & A_{1k} \\ \vdots & \ddots & \vdots \\ A_{k1} & \cdots & A_{kk} \end{bmatrix},$$

where each  $A_{ij}$  is  $N/2$ -by- $N/2$ . (Assume that  $N$ , the dimension of the systolic array, is even so that  $n = kN/2$ .) One way to compute the SVD of  $A$  might be a block Jacobi scheme. See also Schreiber [23]. In this scheme we repeatedly pick  $(p, q)$  satisfying  $1 \leq p < q \leq k$  and use the array of Section 5 to compute the  $N$ -by- $N$  SVD:

$$\begin{bmatrix} U_{pp} & U_{pq} \\ U_{qp} & U_{qq} \end{bmatrix}^T \begin{bmatrix} A_{pp} & A_{pq} \\ A_{qp} & A_{qq} \end{bmatrix} \begin{bmatrix} V_{pp} & V_{pq} \\ V_{qp} & V_{qq} \end{bmatrix} = \begin{bmatrix} D_p & 0 \\ 0 & D_q \end{bmatrix}.$$

Next, we construct the  $n$ -by- $n$  orthogonal matrix  $U$  so that it is equal to the identity matrix except for the four strategic blocks in the  $(p, p)$ ,  $(p, q)$ ,  $(q, p)$  and  $(q, q)$  positions. Those blocks assume the values as given by (6.1). The  $n$ -by- $n$  orthogonal matrix  $V$  is constructed in an identical manner. Then the matrix  $B = U^T A V$  will have the property that

$$\text{off}(B) = \text{off}(A) - 2\|A_{pq}\|_F^2 - \text{off}(A_{pp}) - \text{off}(A_{qq}).$$

The indices  $(p, q)$  can be chosen according to the “row” or the “parallel” ordering. In the latter case, we could exploit a block systolic array (see Figure 6).

The diagonal arrays are SVD arrays. The off-diagonal arrays are matrix-multiply arrays. In this scheme, the blocks move around this array of arrays in exactly the same fashion as the  $a_{ij}$  do in the array of Section 5. Obviously,

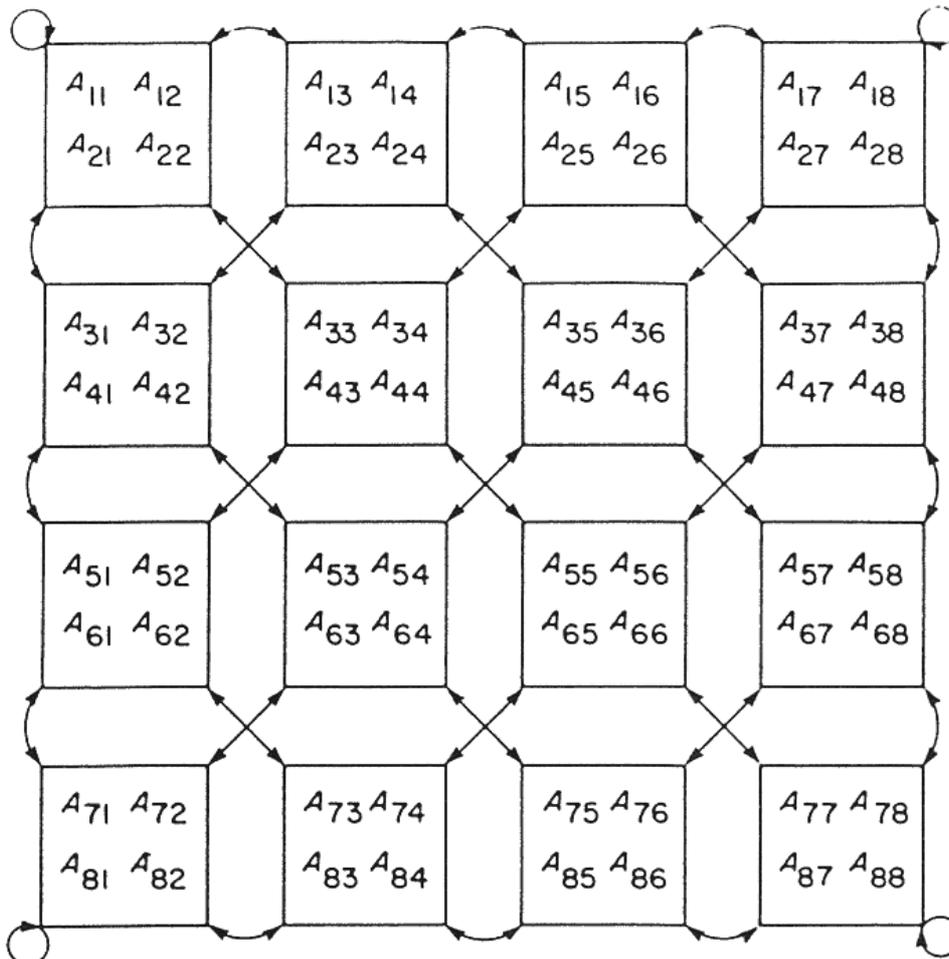


Figure 6. Interarray connections ( $k = 8$ ). (The initial state is illustrated.)

the success of this technique will depend upon the nature of the interconnections, and a host of other unexamined issues.

## 7. REFERENCES

- [1] R. P. Brent and F. T. Luk, The solution of singular-value and symmetric eigenvalue problems on multiprocessor arrays, *SIAM J. Sci. Statist. Comput.* 6 (1985), pp. 69-84.
- [2] R. P. Brent, F. T. Luk, and C. Van Loan, Computation of the generalized singular value decomposition using mesh-connected processors, *Proc. SPIE*, Vol. 431, *Real Time Signal Processing VI*, 1983, pp. 66-71.
- [3] T. F. Chan, An improved algorithm for computing the singular value decomposition, *ACM Trans. Math. Softw.* 8, 1982, pp. 72-83.
- [4] J. J. Dongarra, C. B. Moler, J. R. Bunch, and G. W. Stewart, *LINPAK Users' Guide*, SIAM, Philadelphia, 1979.
- [5] A. M. Finn, F. T. Luk, and C. Pottle, Systolic array computation of the singular value decomposition, *Proc. SPIE*, Vol. 341, *Real-Time Signal Processing V*, 1982, pp. 34-43.
- [6] 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, pp. 1-23.
- [7] B. S. Garbow, J. M. Boyle, J. J. Dongarra, and C. B. Moler, *Matrix Eigensystem Routines—EISPACK Guide Extension*, Springer-Verlag, Berlin, 1977.
- [8] W. M. Gentleman and H. T. Kung, Matrix triangularization by systolic arrays, *Proc. SPIE*, Vol. 298, *Real-Time Signal Processing IV*, 1981, pp. 19-26.
- [9] G. H. Golub and W. Kahan, Calculating the singular values and pseudo-inverse of a matrix, *J. SIAM Ser. B: Numer. Anal.* 2, 1965, pp. 205-224.
- [10] G. H. Golub and F. T. Luk, Singular value decomposition: applications and computations, *Trans. 22nd Conf. Army Mathematicians*, ARO Report 77-1, 1977, pp. 577-605.
- [11] G. H. Golub and C. Reinsch, Singular value decomposition and least squares solutions, in [27], pp. 134-151.
- [12] E. R. Hansen, On cyclic Jacobi methods, *J. Soc. Indus. Appl. Math* 11, 1963, pp. 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, Cambridge, Massachusetts, 1982, pp. 113-122.
- [14] M. R. Hestenes, Inversion of matrices by biorthogonalization and related results, *J. Soc. Indust. Appl. Math* 6, 1958, pp. 51-90.
- [15] H. T. Kung, Why systolic architectures?, *IEEE Computer* 15, 1982, pp. 37-46.
- [16] 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, November 1982, pp. 89-98.
- [17] C. E. Leiserson and J. B. Saxe, Optimizing synchronous systems, *J. VLSI Computer Systems* 1, 1983, pp. 41-67.
- [18] F. T. Luk, Computing the singular-value decomposition of the ILLIAC IV, *ACM Trans. Math. Softw.* 6, 1980, pp. 524-539.
- [19] H. Rutishauser, The Jacobi method for real symmetric matrices, in [27], pp. 202-211.
- [20] A. H. Sameh, On Jacobi and Jacobi-like algorithms for a parallel computer, *Math. Comput.* 25, 1971, pp. 579-590.
- [21] A. H. Sameh, Solving the linear least squares problem on a linear array of processors, *Proc. Purdue Workshop on Algorithmically-specialized Computer Organizations*, 1982.
- [22] R. Schreiber, A systolic architecture for singular value decomposition, *Proc. 1st Internat. Coll. on Vector and Parallel Computing in Scientific Applications*, Paris, France, March 1983.
- [23] R. Schreiber, On the systolic arrays of Brent, Luk, and Van Loan, *Proc. SPIE*, Vol. 431, *Real Time Signal Processing VI*, 1983, pp. 72-76.

- [24] J. M. Speiser and H. J. Whitehouse, Architectures for real-time matrix operations, *Proc. 1980 Government Microcircuits Applications Conf.*, Houston, Texas, November 1980.
- [25] H. J. Whitehouse, J. M. Speiser, and K. Bromley, Signal processing applications of systolic array technology, *Proc. USC Workshop on VLSI and Modern Signal Processing*, Los Angeles, California, November 1982, pp. 5-10.
- [26] J. H. Wilkinson, Note on the quadratic convergence of the cyclic Jacobi process, *Numer. Math* 4, 1962, pp. 296-300.
- [27] J. H. Wilkinson and C. Reinsch, eds., *Handbook for Automatic Computation, Vol. 2 (Linear Algebra)*, Springer-Verlag, New York, 1971.