
**EFFICIENT PORTABLE PARALLEL
MATRIX COMPUTATIONS**

James Walter Juszczak

Department of Computer Sciences
University of Texas at Austin
Austin, TX 78712-1188

TR-89-38

December 1989

Abstract

In this thesis we exercise a method of developing parallel algorithms for matrix computations that facilitates efficient and portable implementations. The method includes defining a set of communication primitives, selecting a storage scheme, embedding a logical communications topology in the physical architecture, and synchronizing data flow and computations to reduce the overhead of communications. Several algorithms are implemented using the column wrapped storage scheme and communication primitives which are independent of the underlying parallel architecture. Theoretical and experimental results are presented.

Table of Contents

Abstract	ii
Table of Contents	iii
List of Tables	iv
List of Figures	v
1. Introduction	1
1.1 Interprocessor Communications	2
1.2 Static Load Balancing	3
1.3 Architectures - ring, mesh, hypercube	4
1.3.1 SYMULT S2010	5
1.4 Notation	7
1.5 Overview	8
2. Communications	10
2.1 The Primitives	10
2.1.1 Node-to-Neighbor	11
2.1.2 One-Way-Shift	11
2.1.3 Broadcast	11
2.1.4 Total Exchange	13
2.1.5 Data Transpose	13
2.1.6 Vector Sum	14

2.1.7	Inner Product	14
2.1.8	Global Compare	14
3.	Systems of Linear Equations	15
3.1	Gaussian Elimination with Partial Pivoting	16
3.1.1	Sequential Algorithm	16
3.1.2	Parallel Algorithm	19
3.1.3	Implementation and Numerical Experiments	21
3.2	Cholesky Decomposition	23
3.2.1	Sequential Algorithm	24
3.2.2	Parallel Algorithm	26
3.2.3	Implementation and Numerical Experiments	28
3.3	Triangular Solve	30
3.3.1	Sequential Algorithm	31
3.3.2	Parallel Algorithm	32
3.3.3	Implementation and Numerical Experiments	38
4.	Q-R Factorization	42
4.1	Householder Orthogonalization	43
4.1.1	Sequential Algorithm	43
4.1.2	Parallel Algorithm	45
4.1.3	Implementation and Numerical Experiments	47
4.2	Modified Gram-Schmidt Method	49
4.2.1	Sequential Algorithm	50
4.2.2	Parallel Algorithm	51
4.2.3	Implementation and Numerical Experiments	52

5. Householder Reduction to Hessenberg Form	56
5.1 Sequential Algorithm	56
5.2 Parallel Algorithm	58
5.3 Implementation and Numerical Experiments	61
6. Conclusion	65
6.1 Summary	65
6.2 Other Applications and Future Work	66
A. Parallel Matrix Subroutines in C	68
A.1 PGEFA	68
A.2 PPOFA	71
A.3 PTRSL	74
A.4 PQRDC	77
A.5 PMGS	80
A.6 PGEHR	83
BIBLIOGRAPHY	86

List of Tables

3.1	Column index below which delay is zero	37
6.1	Algorithm Scorecard	66

List of Figures

1.1	Column Wrapped Storage Scheme for $p = 3$ and $n = 6$	3
1.2	Ring of Processors	4
1.3	Two Dimensional Mesh Network with Embedded Ring	5
1.4	Hypercubes with Dimensions 1, 2, and 3	6
3.1	Gaussian Elimination with Partial Pivoting	18
3.2	Parallel Gaussian Elimination with Partial Pivoting	20
3.3	PGEFA	22
3.4	Observed and expected timings for $p = 8$ during Gaussian Elimination with Partial Pivoting	22
3.5	Ratio of Observed to Expected run times for Gaussian elimination with 8 processors	23
3.6	Observed efficiencies attained during Gaussian Elimination with Par- tial Pivoting	24
3.7	Cholesky Decomposition	25
3.8	Parallel Cholesky Decomposition	27
3.9	PPOFA	28
3.10	Observed timings attained during Cholesky Decomposition	29
3.11	Observed and expected timings for $p = 8$ during Cholesky Decompo- sition	29
3.12	Observed efficiencies attained during Cholesky Decomposition	30
3.13	Back-Substitution	32
3.14	Distributed Back-Substitution	32
3.15	Parallel Back-Substitution	34

3.16	Data Flow During Parallel Back-Substitution	35
3.17	PTRSL	39
3.18	Observed timings attained during backward substitution	40
3.19	Observed and expected timings for $p = 4$ during backward substitution	40
3.20	Observed efficiencies attained during backward substitution	41
4.1	Householder Orthogonalization	44
4.2	Parallel Householder Orthogonalization	45
4.3	PQRDC	47
4.4	Observed and expected timings for $p = 16$ during Householder Or- thogonalization	48
4.5	Observed efficiencies attained during Householder Orthogonalization .	48
4.6	Modified Gram-Schmidt	51
4.7	Parallel Modified Gram-Schmidt	52
4.8	PMGS	53
4.9	Observed timings attained during Modified Gram-Schmidt	54
4.10	Observed and expected timings for $p = 8$ during Modified Gram-Schmidt	54
4.11	Observed efficiency attained during Modified Gram-Schmidt	55
5.1	Householder Reduction to Hessenberg Form	58
5.2	Parallel Householder Reduction to Hessenberg Form	60
5.3	PGEHR	62
5.4	Observed timings attained during Reduction to Upper Hessenberg Form	63
5.5	Observed and expected timings for $p = 16$ during Householder Reduc- tion to Hessenberg Form	63
5.6	Observed efficiencies attained during Reduction to Upper Hessenberg Form	64

Chapter 1

Introduction

The objective of this thesis is to present, analyze and test a method of developing parallel matrix algorithms. The method involves defining a set of communication primitives, selecting a storage scheme, embedding a logical communications topology in the physical architecture, and synchronizing data flow and computations so as to reduce the overhead of communications. A consequence of this work is the development of portable and efficient code to perform a variety of matrix computations in parallel. It is envisioned that ultimately a package of parallel linear algebra routines (PLAPACK) will emerge from this effort.

The recent focus on parallel computing has been in response to the need for higher performance computers. Problems have been encountered that require intensive computations, which cannot be completed in a reasonable amount of time on the current generation of computers. The options are to construct faster hardware or use off the shelf technology to build arrays of processors and perform computation in parallel. We will explore how to exploit the latter of these options in this thesis.

We begin by defining some terminology which we will use to describe the implementation and evaluation of parallel algorithms. Since the objective of parallel computation is to execute algorithms in less time, we will define a measurement by which we can compare algorithms. *Speedup* is the ratio of the time taken by a computer to execute an equivalent serial algorithm and the time taken by the same computer to execute the parallel algorithm using p processors ¹. *Efficiency* is the speedup divided by the number of processors. The efficiency provides a measure of the performance cost and will indicate whether an algorithm fully utilizes the

¹Speedup is sometimes defined as the ratio of the times taken by the fastest serial algorithm over the parallel algorithm.

available processors. Our goal is to approach linear speedup (i.e., speedup = p) and 100 per cent efficiency.

1.1 Interprocessor Communications

In this paper we will concentrate on MIMD (multiple data stream multiple instruction stream) machines with distributed memory. Each processor will operate in an asynchronous manner executing instructions and operating on data stored in its local memory. Each processor will have access to only its local memory and all communication and synchronization will be done by message passing.

Messages are passed through explicit calls to the communication library. This library contains several routines by which processors can send or receive messages. These primitives have been postulated as sufficient to implement a wide range of matrix algorithms [4, 14]. These are:

- node-to-neighbor
- broadcast
- total exchange
- data transpose
- one-way-shift
- vector sum
- inner product
- global compare

By restricting communications to this library of primitives it is possible to execute the calling routines on any machine on which these primitives have been implemented. It is through this communications library that we attain portability of the code. Moreover, by optimizing their implementation on a specific architecture we are able to take advantage of the particular network and achieve efficiency as well as portability. The communications primitives will be discussed in more detail in Chapter 2.

\mathbf{P}_0	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_0	\mathbf{P}_1	\mathbf{P}_2
×	×	×	×	×	×
×	×	×	×	×	×
×	×	×	×	×	×
×	×	×	×	×	×
×	×	×	×	×	×
×	×	×	×	×	×

Figure 1.1: Column Wrapped Storage Scheme for $p = 3$ and $n = 6$

1.2 Static Load Balancing

Having decided on a set of communication primitives, we must concern ourselves with the distribution of the problem among the processors in order to evenly balance the work load, maximizing parallel activity while minimizing communications. Assuming a homogeneous system, we can attain higher efficiencies by assigning equal amounts of work to each processor.

For each of the algorithms in this paper we employ the *column wrapped* storage scheme, depicted in Figure 1.1, to store the matrix in the distributed memories [3, 5, 6, 7]. Given an indexed ordering of p processors $\mathbf{P}_0, \dots, \mathbf{P}_{p-1}$, this scheme assigns the i th column ($i = 1, 2, \dots, n$) to processor $\mathbf{P}_{(i-1) \bmod p}$. Therefore, the first column is assigned to \mathbf{P}_0 , the p th column to \mathbf{P}_{p-1} , the $(p+1)$ th column to \mathbf{P}_0 and so on. If the number of processors divides the number of columns, n , then each processor receives an equal share of the matrix and presumably an equal share of the work. Otherwise $(n \bmod p)$ processors will have $\lceil n/p \rceil$ columns and the rest will have $\lfloor n/p \rfloor$ columns. For $n \gg p$ the imbalance is insignificant. One advantage of the column wrapped storage scheme is that it is simple; since contiguous blocks (columns) of the matrix have been extracted, accessing elements will be straightforward involving little indexing overhead. Note, that for the column wrapped storage scheme, p cannot

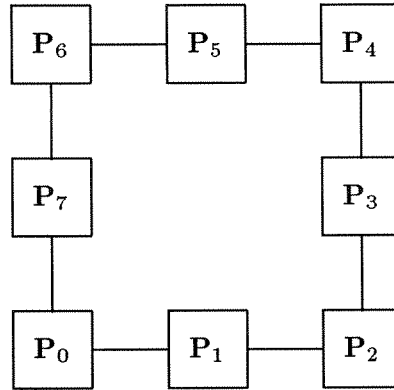


Figure 1.2: Ring of Processors

exceed n .

1.3 Architectures - ring, mesh, hypercube

In this section we will discuss three common distributed memory architectures: the ring, the mesh and the hypercube. In each of these models, it is possible to embed a logical ring of processors in the physical architecture so that neighbors in the ring are also neighbors in the underlying connection network. We will see that when an embedded ring is combined with the column wrapped storage scheme, high efficiencies can be obtained. This is advisable when considering some matrix algorithms since computations can be viewed as progressing across the matrix column by column.

The ring multiprocessor consists of p processors, $\mathbf{P}_0, \dots, \mathbf{P}_{p-1}$, connected in a ring as shown in Figure 1.2, where nodes P_i and P_j are neighbors if $(i + 1) \bmod p = j$ (or $(j + 1) \bmod p = i$).

It is assumed that each processor can simultaneously send to both neighbors or simultaneously send to one neighbor and receive from the other neighbor. In this thesis we assume that a ring can be embedded in the underlying physical architecture.

In a mesh network the processors are arranged in a d -dimensional lattice.

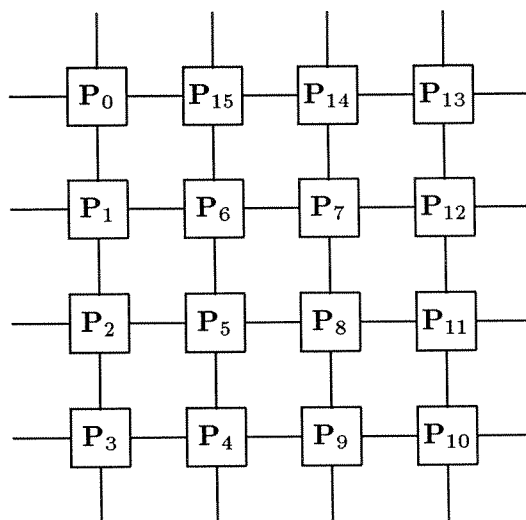


Figure 1.3: Two Dimensional Mesh Network with Embedded Ring

An interior processor has two neighbors in each dimension, giving it $2d$ neighbors. If the mesh has wrap-around connections then every processor has $2d$ neighbors. A processor can only communicate directly with these neighbors. Figure 1.3 shows a two dimensional mesh in which a ring has been embedded. Processors on the edges of the grid have wrap around connections to processors on the opposite edge.

A hypercube is a loosely coupled multiprocessor based on a binary n -cube network. An n dimensional hypercube has 2^n processors and each processor has n neighbors. The processors can be numbered P_0, \dots, P_{2^n-1} , so that neighboring processors differ in exactly one place of the binary representation of their index. Figure 1.4 shows 1, 2 and 3 dimensional hypercubes. The following sequence of processor indices describes an embedded ring in a 3 dimensional hypercube, 000, 010, 011, 001, 101, 111, 110, 100 .

By restricting ourselves to a ring, the simplest and least efficient architecture, we examine the worst case scenario.

1.3.1 SYMULT S2010

All algorithms discussed in this thesis were implemented and tested on the SYMULT Series 2010 parallel processor belonging to the Computer Sciences Department of The University of Texas at Austin. The S2010 is a multiple data multiple

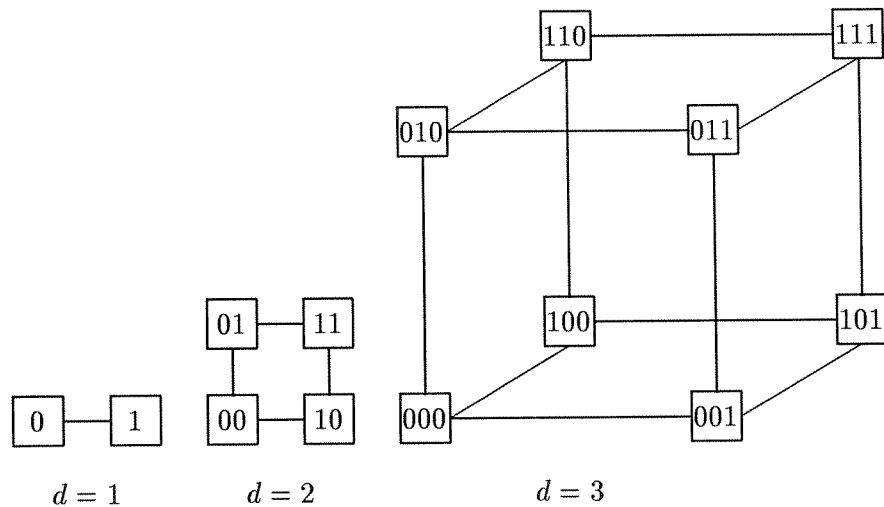


Figure 1.4: Hypercubes with Dimensions 1, 2, and 3

instruction parallel processing system. This machine has twenty-four Motorola 68020 microprocessors connected via automatic message routing device (AMRD) circuits in a 6×4 array topology. Each AMRD has five channels, four for communications to other AMRDs and one for access to its own local memory. Each AMRD can send up to 20 MBytes per second on any two channels simultaneously. When allocating a set of nodes on which to run an application a ring of an even number of nodes can be explicitly requested but owing to the AMRD circuits this is not necessary to obtain good results. Eight of the processors have 3 Mbytes of local memory and the rest each have 1 Mbyte. Applications utilizing more than 8 nodes were restricted in size because of the limited memory on 16 nodes.

All algorithms were implemented using double precision computer arithmetic in the C programming language on a SUN III (UNIX host), under a runtime system called the *Cosmic Environment*, which interacts with a node operating system running and on the Symult S2010 multicomputer called the *Reactive Kernel*. The host and node programs make procedure or function calls to routines from the set of communication primitives discussed in Chapter 2, and a set of BLAS routines written in C and optimized for the Motorola 68020 microprocessor ².

²courtesy of Jim Meyering at the University of Texas at Austin

When comparing the expected and observed timings of the algorithms, it is necessary first to measure the *flop rate* or speed at which a node of the multicomputer performs a computation. We will abide by convention and define a *flop* as the operation

$$y \leftarrow \alpha x + y \quad \alpha \in \mathbf{R}; x, y \in \mathbf{R}^n$$

in which two floating point operations are performed [2]. The flop rate of a single node of the SYMULT was timed through repeated calls to the BLAS routine, *saxpy*, on vectors of length 100, and a flop rate of 71 Kflops was attained. The time for a single flop is therefore, $\gamma = .000014$ seconds.

The cost of a communication between two processors is expressed as $\alpha + m\beta$ where α is the communication startup time, β is the per item (double precision floating point word) transmission time and m is the number of items being transmitted. The startup and transfer times were measured on the SYMULT S2010 and found to be: $\alpha \approx 450\mu\text{sec}$ and $\beta \approx 10\mu\text{sec/item}$ [15]. These measurements are with respect to the communication primitives (i.e., they include buffer manipulations in addition to those performed by the Cosmic Environment [10]).

1.4 Notation

The routines which perform the actual matrix computations will be discussed in their respective chapters. In doing so, we present the algorithms in a systematic way, using consistent notation discussed in this section.

The naming conventions and calling sequences of the LINPACK library developed at the Argonne National Laboratory were adopted where possible to conform to the LINPACK template.

When discussing the algorithms and the theory supporting them, the following notational conventions were used. The set of real numbers is denoted as \mathbf{R} . \mathbf{R}^n is the set of all n -dimensional vectors with real elements and $\mathbf{R}^{m \times n}$ is the set of all $m \times n$ real matrices. Following Stewart [11], lower case Greek letters are used to represent scalars and lowercase Latin letters represent vectors. Uppercase Latin letters denote matrices. Elements of a vector or matrix are denoted with the Greek letter that best corresponds to the Latin letter identifying that vector or matrix. For example, α_{ij} is the element in the i th row and j th column of the matrix A and ξ_i is

the i th component of the vector x . The letters n, m, i, j, k are reserved to represent integer dimensions and indices into vectors or matrices. The index i will hereafter be reserved to index processors in the ring. The superscript T is used to denote the transpose of a vector or matrix (e.g., $y = (v_1, v_2, \dots, v_n)^T$ denotes a column vector).

Superscripts appearing in parentheses indicate the iteration number during which the variable exists. $A^{(k)} = (a_1^{(k)} \dots a_n^{(k)})$ is a column partitioning of $A^{(k)}$. In a ring of p nodes, $\mathbf{P}_0, \dots, \mathbf{P}_{p-1}$, node \mathbf{P}_i has column k if $(k - 1) \bmod p = i$. We denote this as $k \in \mathbf{P}_i$ and this processor is also referred to as $\mathbf{P}(k)$. Therefore the symbols \mathbf{P}_i and $\mathbf{P}(k)$ serve double duty, representing either a processor or a set of indices. We hope that each instance of their use is made clear by the context in which they appear.

1.5 Overview

The intent of this thesis is to demonstrate to the reader through a few representative examples that an effective method of parallelizing matrix algorithms is being presented.

In Chapter 2, we examine in detail the problem of exchanging data in a ring of processors. We will describe algorithms to perform a set of data exchange operations considered sufficient to perform matrix computations. For those operations used in this thesis, we will derive the time complexity of the algorithm as it has been implemented for a ring of processors.

Chapter 3 is concerned with the solution of systems of linear equations. The three algorithms presented in this chapter are : Gaussian elimination with partial pivoting, Cholesky decomposition and triangular solve. The first two algorithms factor a matrix into the product of simpler matrices thereby simplifying subsequent calculations involving the matrix. Both of these algorithms are parallelized using only the broadcast primitive. The triangular solve algorithm receives the factored matrix and the right hand side of the equation, $Ax = b$, and solves for the vector x . The parallel version of this algorithm employs the node-to-neighbor communication primitive.

Chapter 4 is concerned with the QR-factorization of matrices which can be applied to solving the linear least squares problem. In this chapter, the QR decom-

position by Householder transformations and the Modified Gram Schmidt algorithms are both implemented in parallel using only the broadcast primitive.

In Chapter 5, we discuss the Householder reduction to upper Hessenberg form. The parallel algorithm requires the broadcast vector sum and the total exchange primitives.

For each algorithm, we will review the relevant theory in order to present and analyze the sequential algorithm and then discuss the development and complexity of the parallel solution. Issues pertaining to the implementation of these algorithms and the results of the experiments on the SYMULT S2010 will then be presented and compared to the theoretical expectations. Code for these routines can be found in the appendix.

In Chapter 6, we summarize the theoretical and experimental findings, discuss the potential and limitations of our approach and discuss future work and other applications.

Chapter 2

Communications

Many sequential algorithms cannot be decomposed into totally independent tasks, and distributed to separate processors for isolated computation. Instead these parallel tasks must communicate to share data and synchronize processing. It follows that improved interprocessor communications would benefit all but the most perfectly parallelizable algorithms. Typically the time to perform a communication far exceeds the time to perform a floating point operation and often the communication time of an algorithm dominates the computation time until problem sizes become very large.

In this chapter, we will examine in detail the problem of exchanging data in a ring of processors. We will describe and analyze algorithms to perform a set of data exchange operations considered sufficient to perform matrix computations. The theoretical analysis of implementing a set of primitives similar to ours on different parallel architectures is studied in [9]. By restricting communications to only these primitives the calling routines become independent of the underlying architecture. In addition, the programmer is provided with a conceptual tool with which parallel algorithms may be more easily developed.

2.1 The Primitives

In this section, we present the communication primitives some of which we will use to implement the matrix algorithms in the subsequent chapters. We will derive the time complexity of the various communication operations as they have been implemented for the ring topology. In the next sections, the term *node* will often be used to refer to a processor. In particular, the *source* node will refer to the processor that originates the communication and the *destination* node(s) to the processor(s) to which the message is directed.

We assume that sending a packet of m data items (e.g. double precision floating point numbers) between neighboring nodes in the ring requires time $\alpha + m\beta$, where α represents the communications startup time, and β is the per item transmission time. It is also assumed that each processor can simultaneously send to both neighbors or simultaneously send to one neighbor and receive from the other neighbor.

2.1.1 Node-to-Neighbor

Moving data from one processor to another represents the simplest data transfer operation. Here one node sends a packet of m items to a neighboring node. The time complexity for this operation is

$$\alpha + m\beta.$$

This operation is used in solving triangular systems as seen in Section 3.3 and it proves useful when performing timings and handling exceptions.

2.1.2 One-Way-Shift

This operation involves every node sending data to its right neighbor or every node sending data to its left neighbor. Each node sends a packet of m items to its neighbor. The time complexity for this operation is

$$\alpha + m\beta.$$

This operation has the same time complexity as the node-to-neighbor communication since all nodes communicate in parallel.

2.1.3 Broadcast

This data exchange operation involves transferring m data items from one processor, the source node, to all other processors. It occurs frequently in parallel numerical algorithms such as in Gaussian elimination, the QR decomposition and in the reduction of a matrix to upper Hessenberg form.

We will implement this operation in two different ways, which we will call the broadcast and the pipelined broadcast. The broadcast sends the entire

message in one direction around the ring. The pipelined broadcast takes advantage of parallelism in the communication by breaking the message into ν packets and transmitting several packets simultaneously in both directions around the ring.

In the *broadcast* operation the source node initiates the communication by sending the m data items to its right neighbor. All other nodes call the broadcast receive routine which instructs them to receive data from their left neighbor and send the data on to their right neighbor (if it is not the source node). This requires $(p-1)$ sends of the m data items. The time complexity for the broadcast operation is

$$(p-1)(\alpha + m\beta).$$

This implementation may seem rather naive. However, in inherently sequential situations, where it is best for the node adjacent to the source to receive all the information first, this implementation performs very well. This situation arises in the Gaussian elimination, Cholesky and QR decomposition, and modified Gram-Schmidt algorithms presented in Chapters 3 and 4.

In the *pipelined broadcast* the source node breaks the message into ν packets and sends them in both directions around the ring. The maximum distance that a packet will travel is $\lfloor p/2 \rfloor$. Without loss of generality let the source node be \mathbf{P}_0 . Packets sent to the right neighbor will follow the right path, $\mathbf{P}_0, \mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_{\lfloor p/2 \rfloor}$, and packets sent to the left neighbor will traverse the left path, $\mathbf{P}_0, \mathbf{P}_{p-1}, \dots, \mathbf{P}_{(\lfloor p/2 \rfloor + 1)}$. The right path is at least as long as the left path and since data transfers are occurring in parallel we need only consider the right path when considering complexity.

In step 1, \mathbf{P}_0 sends packet #1 to \mathbf{P}_1 . Next, in step 2, \mathbf{P}_0 sends packet #2 to \mathbf{P}_1 , and \mathbf{P}_1 sends packet #1 to \mathbf{P}_2 . After $\lfloor p/2 \rfloor$ steps, packet #1 has reached $\mathbf{P}_{\lfloor p/2 \rfloor}$ and the *pipe* is filled. After $(\nu - 1)$ more steps, the last packet reaches $\mathbf{P}_{\lfloor p/2 \rfloor}$ and the broadcast is complete. The pipelined broadcast then takes $\lfloor p/2 \rfloor + \nu - 1$ steps, where each step takes time $\alpha + m\beta/\nu$. The time complexity for this operation is

$$(\lfloor p/2 \rfloor + \nu - 1)(\alpha + m\beta/\nu).$$

This time is minimized, if the optimal packet size of

$$\sqrt{\frac{m\beta(\lfloor p/2 \rfloor - 1)}{\alpha}}$$

is chosen, making the minimum time complexity,

$$\left(\sqrt{m\beta} + \sqrt{(\lfloor p/2 \rfloor - 1)\alpha}\right)^2.$$

Note: we have not yet implemented the pipelined broadcast.

2.1.4 Total Exchange

The total exchange can be considered a multi-broadcast or $(p-1)$ consecutive one-way-shifts. Every processor sends a block of data of the same size, m , to every other processor. This communication is used in the Householder reduction to upper Hessenberg form, presented in Chapter 5.

Each node begins the communication by sending its m data items to its right neighbor. Once these messages arrive each node stores the received message then sends it on to its right neighbor. This cycle is repeated $(p-1)$ times, at which point all nodes possess all p messages. Since all nodes can send simultaneously the time complexity for this operation is

$$(p-1)(\alpha + m\beta).$$

2.1.5 Data Transpose

This operation is the simultaneous scattering of data packets from each processor to every other processor. As its name suggests, this operation acts much like the transpose operation for a matrix.

Each node \mathbf{P}_i has p packets, each of size h , denoted by x_{ij} , $0 \leq j < p$. Let $m = (p-1)h$ be the total amount of data each processor must send. After the data transpose, each node \mathbf{P}_i has packets x_{ji} , $0 \leq j < p$.

Every node begins the communication by sending its m data items or $(p-1)$ packets to its right neighbor. When each node receives the message of this size from its left neighbor, the node removes one packet from the message and sends the remainder, now of size $(p-2)h$ on to its right neighbor. This cycle is repeated $(p-1)$ times. During the k th cycle the message size is $m_k = (p-k)h$. Since all nodes can send simultaneously the time complexity for this operation is,

$$\sum_{k=1}^{p-1} (\alpha + (p-k)h) = (p-1)\alpha + \frac{pm}{2}\beta.$$

Note: we have not yet implemented the data transpose operation.

2.1.6 Vector Sum

In this operation each node \mathbf{P}_i owns a vector of length $n = hp$, where for simplicity we now assume h is an integer. Each processor divides its vector into p equal sections, indexed by i where $0 \leq i < p$. The sum of section i from all nodes will reside on \mathbf{P}_i on completion of the distributed vector sum.

All processors begin the operation by sending to their right neighbor the section of the vector whose total will eventually reside on their left neighbor. Then each processor receives the section from its left neighbor, adds to this section the corresponding section of its own vector and sends the result on to its right neighbor. This cycle is repeated $(p - 1)$ times. Since all nodes can send simultaneously the time complexity for this operation is,

$$(p - 1)(\alpha + h\beta + h\gamma).$$

The Householder reduction to Hessenberg form, seen in Chapter 5, requires this communication.

2.1.7 Inner Product

Two other primitives that we plan to include in this set are the inner product and global compare operations. We briefly describe them in this and the next section.

Let x and y be vectors of length $n = hp$, and assume that they are distributed among the processors so that each node has parts of x and y of length h . The inner product computes $x^T y$ leaving the result on each node.

2.1.8 Global Compare

Assume that $\xi_i \in \mathbf{P}_i$ for $i = 0, \dots, p-1$, then the global compare operation finds the largest ξ_i and distributes it to all nodes along with the index of the processor that originally owned it.

Chapter 3

Systems of Linear Equations

In this chapter, we will consider algorithms concerned with the problem of solving dense linear systems of equations, $Ax = b$, where $A \in \mathbf{R}^{n \times n}$ and $x, b \in \mathbf{R}^n$. Sequential algorithms which solve this problem are presented in [2, 11] and efficient implementations of these algorithms can be found in [1].

In Sections 3.1 and 3.2 we examine algorithms to factor a matrix into the product of a lower and upper triangular matrix or LU decomposition. Such decompositions are based on the existence of a unique LDU decomposition of a square general matrix A . In a LDU decomposition, L is unit lower triangular, D is diagonal and U is unit upper triangular. The LU decompositions differ in how the diagonal matrix is handled in the LDU decomposition. For dense matrices, $O(n^3)$ operations are required to perform these decompositions.

In Section 3.3 we consider the solution of the resulting triangular systems. Once the LU -decomposition has been obtained the problem, $LUx = b$, reduces to solving two triangular systems,

$$Ly = b$$

and

$$Ux = y.$$

The first system is solved by forward elimination and then the second can be solved by backward substitution. Each of these algorithms requires $O(n^2)$ operations and contribute little to the overall complexity compared to the decomposition. However, in practice it is common for the system $Ax = b_i$ to be solved for many different right hand side vectors, b_i . In this case, the matrix A is decomposed once and the repeated triangular solves become a significant part of the complexity. We implement and test a parallel back-substitution algorithm to solve an upper triangular system. The algorithm for a lower triangular system is analogous.

3.1 Gaussian Elimination with Partial Pivoting

The first of these algorithms is Gaussian elimination with partial pivoting for factoring a permuted general matrix PA into the product of a unit lower triangular matrix and an upper triangular matrix. In terms of the LDU decomposition, $PA = LDU = LU'$ where $U' = DU$ is upper triangular.

DEFINITION 3.1 *An elementary lower triangular matrix of order n and index k is of the form,*

$$M = I_n - me_k^T$$

where $m = (0, \dots, 0, \mu_{k+1}, \mu_{k+2}, \dots, \mu_n)^T$.

Elementary transformations are elementary lower triangular matrices that are used to introduce zero components in a vector and can be exploited to perform the reduction by Gaussian elimination. The following theorem states this property more precisely.

THEOREM 3.1 *Given a vector $x = (\xi_1, \dots, \xi_k, \dots, \xi_n)^T$ and $\xi_k \neq 0$, there exists an elementary transformation, M such that $Mx = (\xi_1, \dots, \xi_k, 0, \dots, 0)^T$. If $\xi_k = 0$, then M does not exist unless ξ_{k+1}, \dots, ξ_n are also zero.*

The following theorem states the existence of this factorization.

THEOREM 3.2 *Let A be an $n \times n$ matrix. Then there are elementary permutations P_i ($i = 1, 2, \dots, n - 1$), and elementary lower triangular matrices M_i of index i ($i = 1, 2, \dots, n - 1$), such that*

$$A_n = M_{n-1}M_{n-2} \dots M_1P_{n-1}P_{n-2} \dots P_1A$$

is upper triangular.

3.1.1 Sequential Algorithm

In Gaussian elimination, a matrix A is reduced to upper triangular form by premultiplying A by a sequence of elementary transformations, M_{n-1}, \dots, M_1 , so as to introduce zeros below the diagonal of the product matrix. The process begins by producing the elementary transformation M_1 and premultiplying $A^{(1)} = A$ to get $A^{(2)} = M_1A^{(1)}$. The matrix $A^{(2)}$ has zeros below the diagonal in its first column. In

the k th step, ($k = 1, \dots, n-1$), of the algorithm the matrix $A^{(k)}$ has zeros below the diagonal in the first $(k-1)$ columns. The matrix, $M_k = I_n - m_k e_k^T$ is determined by the vector,

$$m_k = (0, \dots, 0, \mu_{k+1,k}, \dots, \mu_{nk}),$$

where

$$\mu_{ik} = \alpha_{ik}^{(k)} / \alpha_{kk}^{(k)}, \quad (i = k+1, \dots, n).$$

Here the elements μ_{ik} are called the *multipliers* and α_{kk} , the k th *pivot* element.

Notice that because of the structure of M_k , premultiplying $A^{(k)}$ by M_k does not change the first $(k-1)$ columns or rows of $A^{(k)}$. In fact $A^{(k+1)}$ is identical to $A^{(k)}$ in the first k rows and $k-1$ columns. The effect of the k th step of the algorithm is to annihilate the subdiagonal elements in the k th column and to perform a rank one update on $A_{k+1,k+1}^{(k)}$, the $(n-k) \times (n-k)$ trailing principal submatrix of $A^{(k)}$. In the implementation the subdiagonal zero elements can be overwritten by the multipliers which were used to create the zeros. Premultiplication by M_k then alters the matrix $A^{(k)}$ to produce the matrix $A^{(k+1)}$ as follows,

$$A^{(k+1)} = M_k A^{(k)} = \begin{pmatrix} I_{k-1} & 0 \\ 0 & M'_k \end{pmatrix} \begin{pmatrix} A_{11}^{(k)} & A_{12}^{(k)} \\ 0 & A_{kk}^{(k)} \end{pmatrix}.$$

Here M'_k is an elementary transformation of index 1, and

$$A_{kk}^{(k)} = \begin{pmatrix} \alpha_{kk}^{(k)} & (a_{k,k+1}^{(k)})^T \\ a_{k+1,k}^{(k)} & A_{k+1,k+1}^{(k)} \end{pmatrix}.$$

This results in $A^{(k+1)}$ having the following form,

$$A^{(k+1)} = \begin{pmatrix} A_{11}^{(k)} & A_{12}^{(k)} \\ 0 & \begin{pmatrix} \alpha_{kk}^{(k)} & (a_{k,k+1}^{(k)})^T \\ 0 & A_{k+1,k+1}^{(k+1)} \end{pmatrix} \end{pmatrix} = \begin{pmatrix} A_{11}^{(k+1)} & A_{12}^{(k+1)} \\ 0 & A_{k+1,k+1}^{(k+1)} \end{pmatrix},$$

where the multipliers are given by

$$\mu_{k+1,k}^{(k+1)} = a_{k+1,k}^{(k)} / \alpha_{kk}^{(k)}$$

and

$$A_{k+1,k+1}^{(k+1)} = A_{k+1,k+1}^{(k)} - \mu_{k+1,k}^{(k+1)} (a_{k,k+1}^{(k)})^T.$$

If the k th pivot element is zero, the process cannot proceed unless all the sub-diagonal elements of the k th column of $A^{(k)}$ are also zero. Furthermore, the

Algorithm 3.1 *The following algorithm uses Gaussian elimination with partial pivoting to overwrite A with the LU decomposition of PA . The pivot row indices are stored in a separate vector of length $(n - 1)$.*

```

 $A_1 = A$ 
for  $k = 1, \dots, n - 1$ 
    find pivot to determine  $P_k$ 
    compute  $M_k = I_n - m_k e_k^T$ 
     $A \leftarrow A^{(k+1)} = M_k P_k A^{(k)}$ 

```

Figure 3.1: Gaussian Elimination with Partial Pivoting

algorithm becomes unstable if the pivot element is close to zero, since this can make the multipliers large and cause subtractive cancellation when updating elements in the submatrix [2, 11]. However, the algorithm can be made more stable by *partial pivoting* or swapping rows of the submatrix so that the largest element on or below the diagonal in the k th column is moved to the diagonal position. Performing a pivot with the i th row during the k th step of the algorithm ($i \geq k$) is equivalent to premultiplying the matrix $A^{(k)}$ by a permutation matrix P_k , which is of the form

$$P_k = \begin{pmatrix} I_{k-1} & 0 \\ 0 & P_{i1} \end{pmatrix},$$

where P_{i1} differs from I_{n-k+1} only in that the first and i th rows are interchanged. With pivoting, the k th step amounts to premultiplying the matrix $A^{(k)}$ by the permutation matrix P_k and the elementary triangular matrix M_k so that,

$$A^{(k+1)} = M_k P_k A^{(k)}$$

has zeros below the diagonal in the first k columns.

The result of this process is $M_{n-1} P_{n-1} \dots M_1 P_1 A = U$, where U is an upper triangular matrix. Letting $P = P_{n-1} \dots P_1$ and noting that $P_i^{-1} = P_i$, we have, $PA = P(M_{n-1} P_{n-1} \dots M_1 P_1)^{-1} U = M_1^{-1} \dots M_{n-1}^{-1} U$ or $PA = LU$, where $L = M_1^{-1} \dots M_{n-1}^{-1}$ is a unit lower triangular matrix. The elements of L below the diagonal are the multipliers μ_{ij} produced when determining M_j to zero out the j th column of $A^{(j)}$.

The time to find the pivot and compute the multipliers is insignificant for large n , compared to the last step of the loop (see Figure 3.1), in which the submatrix

is updated, requiring $(n - k)^2$ flops. If we sum this expression over all iterations we get,

$$T_1(n) \approx \sum_{k=1}^{n-1} (n - k)^2 \gamma \approx \frac{n^3}{3} \gamma. \quad (3.1)$$

Ignoring low order terms, this algorithm requires approximately $n^3/3$ flops.

3.1.2 Parallel Algorithm

We begin our discussion of the parallel algorithm to perform Gaussian elimination with pivoting by considering the distribution of the problem in a *column wrapped* fashion among a ring of p processors, $\mathbf{P}_0, \dots, \mathbf{P}_{p-1}$. As described in Section 1.2, the first column is assigned to \mathbf{P}_0 , the p th column to \mathbf{P}_{p-1} , the $(p + 1)$ th column to \mathbf{P}_0 and so on. If the number of processors divides the number of columns, n , then each processor receives an equal share of the matrix. Otherwise $(n \bmod p)$ processors will have $\lceil n/p \rceil$ columns and the rest will have $\lfloor n/p \rfloor$ columns.

Given that each processor has approximately (n/p) columns, we consider the sequential algorithm, which loops through a set of computations for each of the first $(n - 1)$ columns of the matrix. Within each iteration, the sequence of computations is: find the pivot row, determine the multipliers, and update the submatrix. The pivot row and multipliers can all be determined by processor $\mathbf{P}(k)$, the one owning the k th column. The submatrix, however, is distributed among all processors, so all processors are required to assist in applying the update.

First, note the need for synchronization and communication. There is a required order to the computations that is satisfied by completely applying the k th update to column j before applying the $(k + 1)$ th update to that same column. The k th iteration begins with finding the k th pivot row and then computing the multipliers before the $(n - k) \times (n - k)$ submatrix is updated. The columns of the submatrix can be updated in any order and hence the update can be performed in parallel by all processors once they have received the multipliers and pivot row index from $\mathbf{P}(k)$. This requires a communication from one processor to all others, each iteration. The *broadcast* primitive can be used to disperse this information as well as to synchronize the computations among all processors so that the order of computations is preserved.

Algorithm 3.2 *This algorithm uses the broadcast primitive alone to reduce a general matrix $A \in \mathbf{R}^{n \times n}$ to triangular form using elementary transformations with partial pivoting overwriting PA with LU . Pseudo code driving node \mathbf{P}_i in a ring of processors is given by,*

```

for  $k = 1, \dots, n - 1$ 
  if  $k \in \mathbf{P}_i$ 
    find pivot
    compute  $M_k = I_n - m_k e_k^T$  [[ $(n - k)\gamma$ ]
    broadcast  $m_k$  and pivot index [ $2(\alpha + (n - k + 1)\beta)$ ]
  else
    receive  $m_k$  and pivot index
    update  $a_j \leftarrow a_j - \alpha_{kj} m_k, j \in \mathbf{P}_i > k$  [[ $\lceil \frac{n-k}{p} \rceil (n - k) \gamma$ ]]

```

Figure 3.2: Parallel Gaussian Elimination with Partial Pivoting

One more observation before we specify the algorithm for any processor. We recall from Section 2.1.3 that the non-pipelined broadcast primitive involved the source node sending the complete message to its right neighbor in the ring. This right neighbor then receives the message and passes it on to its own right neighbor and so on, until the message reaches the left neighbor of the source node. Returning to the parallel algorithm for Gaussian elimination, during the k th iteration processor $\mathbf{P}(k)$ computes the pivot row and multipliers, then broadcasts this data to the other processors. $\mathbf{P}(k + 1)$ is the first processor to receive the data and will continue the broadcast operation by sending the data on to $\mathbf{P}(k + 2)$ at which point $\mathbf{P}(k + 1)$ is free to begin updating its part of the submatrix. The effect of implementing the broadcast in this manner is an overlap of the remainder of the broadcast with computation. As soon as $\mathbf{P}(k + 1)$ completes updating its part of the submatrix it may determine the $(k + 1)$ th pivot row, compute the multipliers, and begin the next broadcast. As the algorithm nears completion, (i.e, $k > (n - p)$) the data need not be broadcast to all other processors since the submatrix is contained in fewer than p processors' memories. However, this contributes little to the time complexity and can be overlooked in favor of the simplicity of the algorithm.

In Figure 3.2 the expressions in brackets indicate the effective contribution to the time complexity by that step in each iteration. Notice that, as discussed above, the effective time to perform the broadcast is equivalent to the time to perform two

node-to-neighbor communications with a $(n - k + 1)$ size vector. Here we assume that the sending processor is not free to compute until the message is received. The time complexity for a single iteration is then,

$$(n - k)\gamma + 2(\alpha + (n - k + 1)\beta) + \left(\left\lceil \frac{n - k}{p} \right\rceil (n - k)\gamma\right).$$

Noting that $\lceil x \rceil \leq (x + 1)$, and summing over the $(n - 1)$ iterations of the loop we arrive at an upper bound for the total time complexity.

$$T_p(n) \leq \frac{\gamma}{p} \sum_{k=1}^{n-1} (n - k)^2 + 2\gamma \sum_{k=1}^{n-1} (n - k) + 2(n - 1)\alpha + 2\beta \sum_{k=1}^{n-1} (n - k + 1)$$

Performing this sum and ignoring low order terms we get

$$T_p(n) \approx \left(\frac{n^3}{3p} + n^2\right) \gamma + 2n\alpha + n^2\beta. \quad (3.2)$$

This compares favorably with the sequential time complexity, given by (3.1). For a fixed value of p , as the problem size increases we have

$$\lim_{n \rightarrow \infty} \frac{T_1(n)}{T_p(n)} = p,$$

so the speedup approaches p and efficiency nears 100%. In (3.2) we retain the term, $n^2\gamma$, because of its increasing significance as p approaches n . This term results from the cost of computing M_k sequentially each iteration and as p approaches n it contributes, along with communication cost, to the inefficiency.

3.1.3 Implementation and Numerical Experiments

The routine PGEFA (**p**arallel **g**eneral matrix **f**actorization) drives each node in a ring of processors to factor a square general matrix by Gaussian elimination. The calling sequence and argument descriptions are given in Figure 3.3.

Figure 3.4 displays the expected and observed timings that were obtained by PGEFA with a ring of eight processors on the Symult S2010 for $n = 10, 20, \dots, 490$. Figure 3.5 plots the ratio of these observed results over the expected times. Discrepancies between these times are explained by the combined effects of several factors. For smaller problems, the effect of having dropped negative low order terms when approximating the time complexity tends to make the expected time greater than the

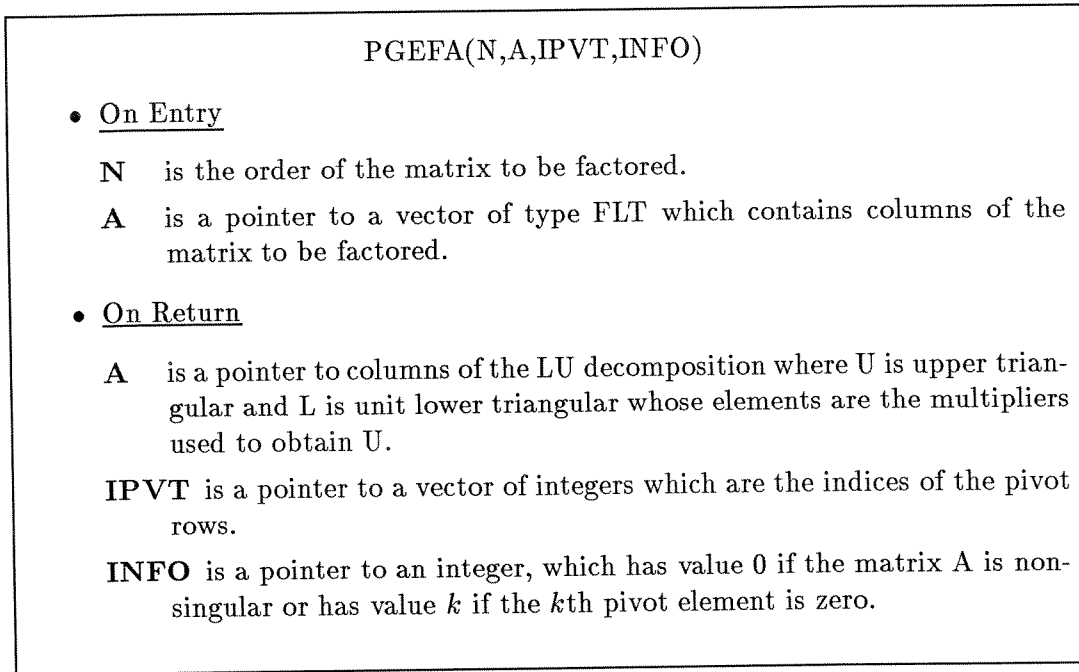
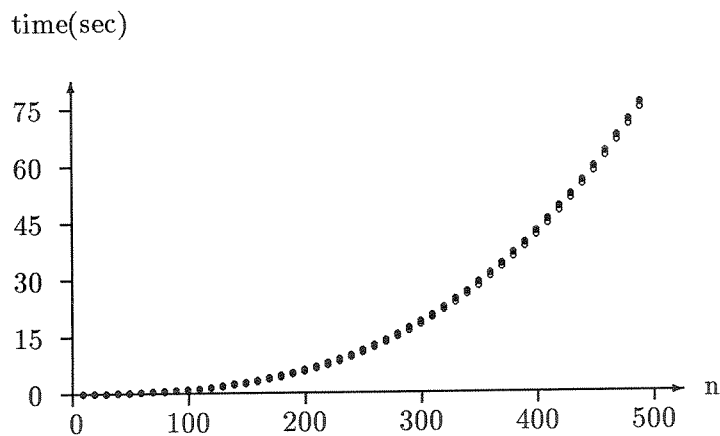


Figure 3.3: PGEFA

Figure 3.4: Observed and expected timings for $p = 8$ during Gaussian Elimination with Partial Pivoting (key: ●: observed, ○: expected)

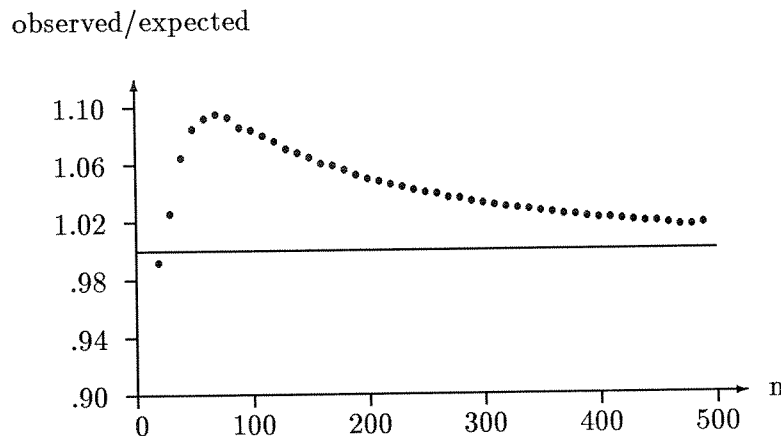


Figure 3.5: Ratio of Observed to Expected run times for Gaussian elimination with 8 processors

observed time. Also, for small n , γ is underestimated owing to the overhead in calling the BLAS routines. This tends to make the expected time less than the observed time. As the problem size increases this underestimation becomes less pronounced. In addition, the value used for α is an overestimation when n is small, because of a message buffer size of 256 bytes [15]. This tends to make the expected time greater than the observed time for message sizes less than 30 double precision floating point numbers. The improvement after $n > 70$ is caused by the dominating influence of the $(n^3/3)\gamma/p$ term, which eventually overshadows all other terms contributing to the time complexity. This agreement is better for fewer than eight processors and slightly worse for 16 or 24 processors.

Figure 3.6 plots the efficiencies attained by PGEFA for various ring and problem sizes. Reasonable efficiencies ($> 50\%$) are obtained once $n/p > 10$.

3.2 Cholesky Decomposition

We will now consider the decomposition of a symmetric positive definite matrix A into the product of a lower triangular matrix and its transpose. In terms of the LDU decomposition, $A = LDU = LDL^T = LD^{1/2}D^{1/2}L^T$. Letting $L' = LD^{1/2}$, we have $A = L'(L')^T$ where L' is a lower triangular matrix. This decomposition is

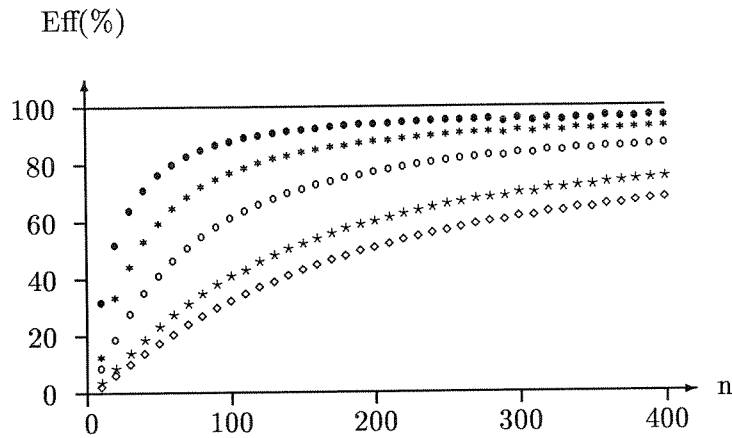


Figure 3.6: Observed efficiencies attained during Gaussian Elimination with Partial Pivoting (key: \bullet : $p = 2$, $*$: $p = 4$, \circ : $p = 8$, \star : $p = 16$, \diamond : $p = 24$)

known as the Cholesky decomposition and L' as the Cholesky triangle.

Obviously, this decomposition does not exist for every matrix; A must be symmetric and the elements of D must be nonnegative limiting A to the class of matrices which are symmetric positive semi-definite. The algorithms presented in Figures 3.7 and 3.8 require that A be positive definite.

DEFINITION 3.2 *A symmetric matrix, $A \in \mathbf{R}^{n \times n}$, is positive definite if and only if for every nonzero $x \in \mathbf{R}^n$,*

$$x^T A x > 0.$$

THEOREM 3.3 *Given that $A \in \mathbf{R}^{n \times n}$ is symmetric and positive definite, there is a unique lower triangular matrix L with positive diagonal elements such that $A = LL^T$.*

3.2.1 Sequential Algorithm

The proof of the existence of the Cholesky decomposition for positive definite matrices is constructive in that it suggests an algorithm with which we can compute the Cholesky triangle. However, we can also arrive at an algorithm to

Algorithm 3.3 *The following algorithm computes the Cholesky triangle by column and overwrites the lower half of A with L .*

```

for  $k = 1, \dots, n$ 
   $\lambda_{kk} = \alpha_{kk} \leftarrow (\alpha_{kk})^{1/2}$ 
   $\lambda_{ik} = \alpha_{ik} \leftarrow \alpha_{ik} / \lambda_{kk}$            ( $i > k$ )
  for  $j = k + 1, \dots, n$ 
    for  $i = j, \dots, n$ 
       $\alpha_{ij} \leftarrow (\alpha_{ij} - \lambda_{ik}\lambda_{jk})$ 

```

Figure 3.7: Cholesky Decomposition

compute L by considering entries of $A = LL^T$,

$$\alpha_{ij} = \sum_{k=1}^{\min(i,j)} \lambda_{ik}\lambda_{jk},$$

since the (k, j) th element of L^T is λ_{jk} . If we consider only elements in the lower triangle, where $i \geq j$ we have,

$$\alpha_{ij} = \sum_{k=1}^j \lambda_{ik}\lambda_{jk} = \left(\sum_{k=1}^{j-1} \lambda_{ik}\lambda_{jk} \right) + \lambda_{ij}\lambda_{jj}$$

which after rearranging terms yields the following equations,

$$\lambda_{ij} = \left(\alpha_{ij} - \sum_{k=1}^{j-1} \lambda_{ik}\lambda_{jk} \right) / \lambda_{jj} \quad (i > j)$$

and

$$\lambda_{jj} = \left(\alpha_{jj} - \sum_{k=1}^{j-1} \lambda_{jk}^2 \right)^{1/2}.$$

We now vary the order in which these computations are performed so that as the k th column of L is formed, the associated update of the lower half of the $(n - k) \times (n - k)$ trailing submatrix of A is applied before the $(k + 1)$ th column is formed [8], much like Gaussian elimination. The result is given by Algorithm 3.7.

Most of the work is done in the loop indexed by j , updating the lower *half* of the $(n - k) \times (n - k)$ submatrix. There are $((n - k)^2 + (n - k))/2$ elements in this

portion of the submatrix. Summing this over all iterations gives,

$$T_1(n) \approx \frac{1}{2} \sum_{k=1}^n ((n-k)^2 + (n-k))\gamma \approx \frac{n^3}{6}. \quad (3.3)$$

Ignoring the lower order terms, this algorithm requires approximately $n^3/6$ flops for large n . By taking advantage of symmetry, we are able to approximately halve the number of computations required by Gaussian elimination to factor a matrix into a product of triangular matrices.

3.2.2 Parallel Algorithm

Again we distribute the matrix A to the ring of processors in a column wrapped fashion. The parallelization of the algorithm proceeds almost identically to that for Gaussian elimination. Given that each processor has approximately (n/p) columns, we consider the sequential algorithm which loops through a set of computations to form each of the n columns of the matrix L . Within each iteration, the sequence of computations is to first determine λ_{kk} , then scale the k th column, and then update the lower half of the submatrix. In the k th iteration, the k th column of L can be completely determined by $\mathbf{P}(k)$, the processor that owns the k th column of A . The submatrix however, is distributed among all processors, so all processors are required to assist in applying the update. The (i, j) th element is updated as follows,

$$\alpha_{ij} \leftarrow (\alpha_{ij} - \lambda_{ik}\lambda_{jk}).$$

For a processor to update the j th column of the submatrix it must have the newly formed k th column of L . This requires a broadcast by $\mathbf{P}(k)$ of the k th column to all other processors. The updating of the submatrix can be performed in parallel by all processors owning columns of that submatrix once they have received the vector $(\lambda_{kk}, \dots, \lambda_{nk})^T$ from $\mathbf{P}(k)$.

Again, in Figure 3.8, the expressions in brackets indicate the effective contribution to the time complexity of that step in each iteration. The scaling of the k th column takes $(n-k)\gamma$ time and, as was the case in Gaussian elimination, the effective time to perform the broadcast is equivalent to the time to perform two node-to-neighbor communications with a $(n-k+1)$ size vector. Continuing with this reasoning, the $(k+1)$ th iteration may begin as soon as $\mathbf{P}(k+1)$ completes

Algorithm 3.4 *This algorithm drives processor \mathbf{P}_i in a ring of processors to factor a symmetric positive definite matrix by the Cholesky algorithm. The matrix is distributed among the processors in a column wrapped fashion.*

```

for  $k = 1, \dots, n$ 
  if  $k \in \mathbf{P}_i$ 
     $\lambda_{kk} = \alpha_{kk} \leftarrow (\alpha_{kk})^{1/2}$ 
     $\lambda_{ik} = \alpha_{ik} \leftarrow \alpha_{ik} / \lambda_{kk}$       ( $i > k$ )     $[(n - k)\gamma]$ 
    broadcast  $(\lambda_{kk}, \dots, \lambda_{nk})^T$        $[2(\alpha + (n - k + 1)\beta)]$ 
  else
    receive  $(\lambda_{kk}, \dots, \lambda_{nk})^T$ 
  for  $j \in \mathbf{P}_i$  and  $j > k$ 
     $\alpha_{ij} \leftarrow (\alpha_{ij} - \lambda_{ik}\lambda_{jk})$       ( $i \geq j$ )     $[\frac{1}{2p}(n - k)(n - k + p)\gamma]$ 

```

Figure 3.8: Parallel Cholesky Decomposition

updating the lower half of columns that it owns and that belong to the submatrix (i.e., columns $j \ni j \in \mathbf{P}(k + 1) \wedge j > k$). Processor $\mathbf{P}(k + 1)$ owns $[(n - k)/p]$ columns of this submatrix, and these columns have indices $j = k + 1 + ip$, for $i = 0, 1, \dots, [(n - k)/p] - 1$. In column j , there are $(n - j + 1)$ elements to update. The time for $\mathbf{P}(k + 1)$ to update its share of the submatrix is then,

$$\sum_{i=0}^{[(\frac{n-k}{p})]-1} (n - k - ip)\gamma,$$

and using the approximation $[(n - k)/p] \approx (n - k)/p + 1$ yields the third expression in brackets,

$$\sum_{i=0}^{\frac{n-k}{p}} (n - k - ip)\gamma = \frac{1}{2p}(n - k)(n - k + p)\gamma.$$

Summing these three expressions over all iterations gives the total time complexity of the algorithm,

$$\begin{aligned} T_p(n) &\approx \sum_{k=1}^n \left[\frac{1}{2p}(n - k)(n - k + p)\gamma + (n - k)\gamma \right. \\ &\quad \left. + 2(\alpha + (n - k + 1)\beta) \right] \\ &\approx \left(\frac{n^3}{6p} + \frac{3n^2}{4} \right) \gamma + 2n\alpha + n^2\beta. \end{aligned} \tag{3.4}$$

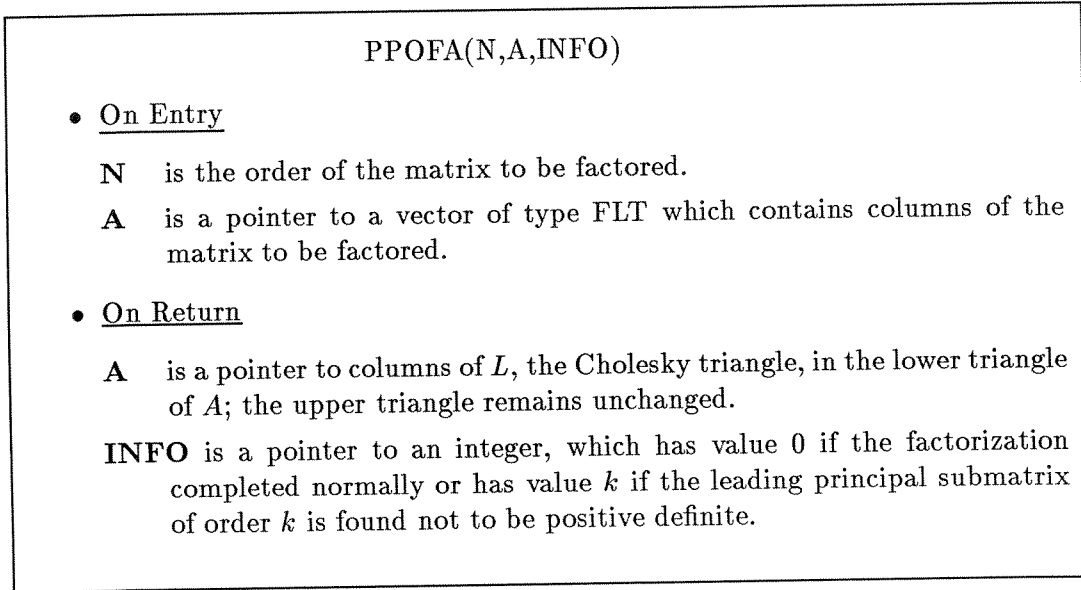


Figure 3.9: PPOFA

This algorithm also approaches 100% efficiency for large n and fixed p . Comparing equations (3.3) and (3.4) we see that,

$$\lim_{n \rightarrow \infty} \frac{T_1(n)}{T_p(n)} = p.$$

However, as p approaches n the two terms $n^2\beta$ and $(3n^2/4)\gamma$ gain significance and efficiency drops. The term $(3n^2/4)\gamma$ reflects the cost of the sequential portion of each iteration, where one processor determines λ_{kk} and scales the k th column.

3.2.3 Implementation and Numerical Experiments

The routine PPOFA (**p**arallel **p**ositive definite matrix **f**actorization) drives each node in a ring of processors to factor a symmetric positive definite matrix by the Cholesky algorithm. The calling sequence and argument descriptions are given in Figure 3.9.

Figure 3.10 plots the observed timings obtained by PPOFA with various ring sizes on the Symult S2010 for problems ranging up to $n = 490$. Figure 3.11 plots the expected and observed timings obtained by PPOFA with a ring of eight processors on the Symult S2010 for problem sizes up to $n = 490$. The expected and observed timings agree within 12%. Again, discrepancies between these times are

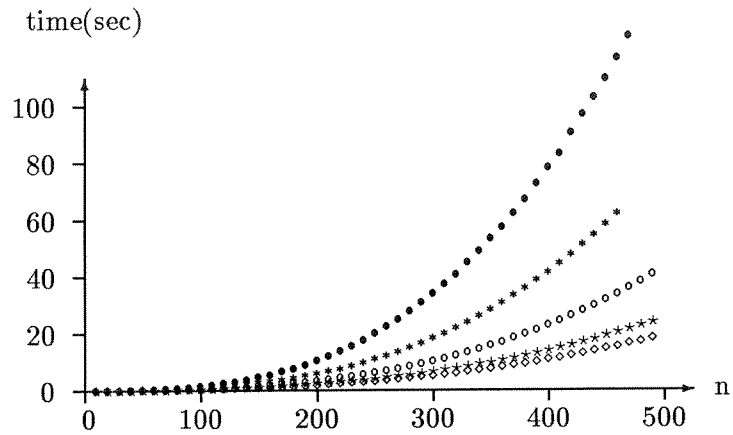


Figure 3.10: Observed timings attained during Cholesky Decomposition
(key:: \bullet : $p = 2$, $*$: $p = 4$, \circ : $p = 8$, \star : $p = 16$, \diamond : $p = 24$)

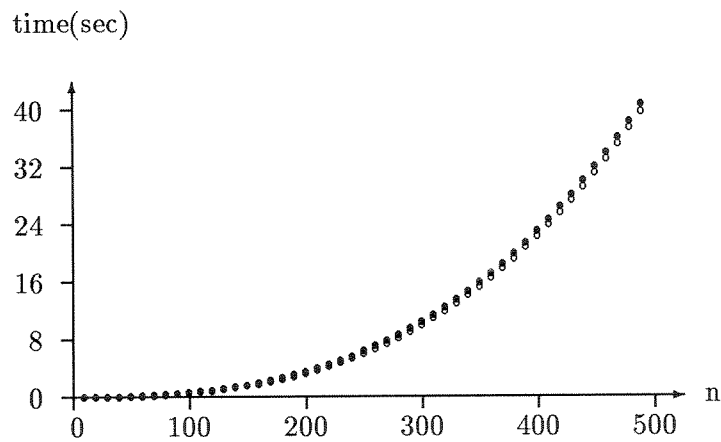


Figure 3.11: Observed and expected timings for $p = 8$ during Cholesky Decomposition
(key:: \bullet : observed, \circ : expected)

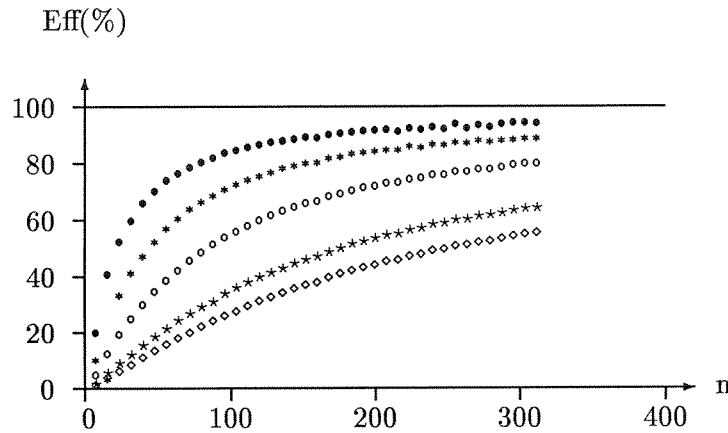


Figure 3.12: Observed efficiencies attained during Cholesky Decomposition
(key: ●: $p = 2$, *: $p = 4$, ○: $p = 8$, ☆: $p = 16$, ◇: $p = 24$)

explained by the combined effects of several factors. This agreement is better for fewer than eight processors and slightly worse for 16 and 24 processors. Within the range of tests performed, it is as low as .1% when $p = 2$ and $n = 400$, and as high as 20% when $p = 24$ and $n = 130$.

Figure 3.12 plots the observed efficiency attained by PPOFA for various ring sizes and problem sizes. Reasonable efficiencies can be obtained once $n/p > 15$. The Cholesky algorithm is initially less efficient than that for Gaussian elimination. The submatrix update is the portion of each iteration that is performed in parallel. Since in Cholesky's algorithm there is less to update, the ratio of useful computation to communication is smaller and efficiency suffers.

3.3 Triangular Solve

In this section we will discuss only the algorithm for backward substitution, which solves the following system for x ,

$$Ux = b,$$

where $U \in \mathbf{R}^{n \times n}$ is upper triangular and $x, b \in \mathbf{R}^n$. The algorithm for a lower triangular system is analogous.

3.3.1 Sequential Algorithm

The algorithm to solve an upper triangular system of linear equations is an $O(n^2)$ operation. It is best introduced by first considering a 2×2 system. In this case we have,

$$\begin{pmatrix} v_{11} & v_{12} \\ 0 & v_{22} \end{pmatrix} \begin{pmatrix} \xi_1 \\ \xi_2 \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}.$$

The unknowns can be found provided $v_{11}v_{22} \neq 0$, by first solving the equation,

$$\xi_2 = \beta_2/v_{22}$$

for ξ_2 and then substituting this value into

$$v_{11}\xi_1 + v_{12}\xi_2 = \beta_1,$$

and reducing this to a single equation with a single unknown,

$$\xi_1 = (\beta_1 - v_{12}\xi_2)/v_{11}.$$

Larger systems of order n can be solved in a similar manner, by solving for the last element of x and then updating the first $(n-1)$ elements of b , thereby reducing the remaining system to an upper triangular system of size $(n-1) \times (n-1)$. The process is repeated n times; each time a new element of x is found, the problem size is reduced to a smaller triangular system. The j th step of the algorithm is outlined below for $j = n, \dots, 1$ and the algorithm is listed in Figure 3.3.1. We enter the j th iteration with the $j \times j$ system $Ux = b$, which can be partitioned as

$$\begin{pmatrix} U_{11} & \tilde{u}_j \\ 0 & v_{jj} \end{pmatrix} \begin{pmatrix} \tilde{x} \\ \xi_j \end{pmatrix} = \begin{pmatrix} \tilde{b} \\ \beta_j \end{pmatrix}$$

where,

$$\tilde{b}_j = (v_{1j}, \dots, v_{j-1,j})^T, \quad \tilde{x} = (\xi_1, \dots, \xi_{j-1})^T \quad \text{and} \quad \tilde{b} = (\beta_1, \dots, \beta_{j-1})^T.$$

If $v_{jj} \neq 0$ we can solve for

$$\xi_j = \beta_j/v_{jj}$$

and then reduce the problem to a $(j-1) \times (j-1)$ upper triangular system by updating the first $(j-1)$ elements of b with

$$\tilde{b} \leftarrow \tilde{b} - \xi_j \tilde{u}_j$$

Algorithm 3.5 *The following algorithm solves the $n \times n$ upper triangular system of linear equations, $Ux = b$, provided $U = (u_1, \dots, u_n)$ is non-singular.*

```

for  $j = n, \dots, 1$ 
     $\xi_j \leftarrow \beta_j / v_{jj}$ 
     $\tilde{b} \leftarrow \tilde{b} - \xi_j \tilde{u}_j$ 
     $\tilde{b}, \tilde{u}_j \in \mathbf{R}^{j-1}$ 

```

Figure 3.13: Back-Substitution

Algorithm 3.6 *This algorithm drives processor \mathbf{P}_i in a ring of processors to solve the $n \times n$ upper triangular system of linear equations, $Ux = b$, provided $U = (u_1, \dots, u_n)$ is non-singular. The matrix is distributed among the processors in a column wrapped fashion.*

```

for  $j = n, \dots, 1$ 
    if  $j \in \mathbf{P}_i$ 
        if  $(j \neq n)$  receive( $b$ )
         $\xi_j \leftarrow \beta_j / v_{jj}$ 
         $b \leftarrow b - \xi_j u_j$ 
        if  $(j \neq 1)$  send left( $b$ )

```

Figure 3.14: Distributed Back-Substitution

and then letting $U \leftarrow U_{11}$, $x \leftarrow \tilde{x}$, and $b \leftarrow \tilde{b}$. Clearly, the solution exists and is unique if and only if the diagonal elements of U are non-zero.

In each pass through the loop, there is a divide operation and a vector update (saxpy) operation of order $(j - 1)$, totaling to approximately j flops. If we sum these over all n iterations we get a time complexity of

$$T_1(n) \approx \frac{n^2}{2} \gamma. \quad (3.5)$$

3.3.2 Parallel Algorithm

The matrix U is distributed to the ring of processors in a column wrapped fashion and the right hand side vector, b , is given to $\mathbf{P}(n)$, the processor that holds the n th column of U .

First, we consider a distributed version of the sequential algorithm to motivate the parallel algorithm (see Figure 3.14). In this algorithm we let $U = (u_1, u_2, \dots, u_n)$

be a column partitioning of U . The vector b starts at $\mathbf{P}(n)$ and is passed around the ring from $\mathbf{P}(n)$ to $\mathbf{P}(n-1)$ to \dots to $\mathbf{P}(1)$. At processor $\mathbf{P}(j)$, b is used to determine the j th component of x and is then updated before being passed on to the next processor. In this distributed algorithm, only one processor computes at any given time, hence there is no parallel activity. However, a slight modification allows processors to compute simultaneously. The resulting parallel algorithm was developed, analyzed and tested by Li and Coleman [5, 6].

In each loop iteration of the distributed sequential algorithm, the majority of work is done in the update of the vector b . The object then is to have all processors performing these updates simultaneously. However, each update of b requires that a new element of x be determined and these elements must be determined in order. For processor $\mathbf{P}(j)$ to compute ξ_j , it must have access to $\beta_j, \xi_{j+1}, \dots, \xi_n$, and $v_{j,j+1}, \dots, v_{jn}$. These data, however, are distributed among several processors.

Access to these values is provided by the auxiliary vectors s (sum) and t (partial sum). Initially, on processor $\mathbf{P}(n)$, s contains the last $(p-1)$ elements of b ,

$$\sigma_k = \beta_{n-k+1} \quad \text{for } k = 1, \dots, p-1$$

and t contains the first $(n-p+1)$ elements of b ,

$$\tau_k = \beta_k \quad \text{for } k = 1, \dots, n-p+1.$$

These vectors are initialized to zero on all other processors. As the algorithm proceeds, for any processor, τ_k will contain a partial sum of those terms to be divided by v_{kk} to determine ξ_k , that are resident on that particular processor, (including β_k on $\mathbf{P}(n)$). The vector s is used as a communications buffer in which the partial sums of all other contributing processors are gathered for determining the next $(p-1)$ elements of the vector x .

Processor $\mathbf{P}(j)$ loops through the following sequence of steps:

- receive the vector s from $\mathbf{P}(j+1)$, ($j \neq n$), which contains partial sums used to determine $\xi_j, \dots, \xi_{j+p-2}$,
- determine ξ_j ,
- shift s so that the next element of t is entered into the last position of s and add to each element of s both partial sums from t and the new term containing ξ_j ,

Algorithm 3.7 *This algorithm drives processor \mathbf{P}_i in a ring of processors to solve the upper triangular system of equations $Ax = b$ by back substitution. The matrix is distributed among the processors in a column wrapped fashion. On completion the vector x is distributed among the processors such that $\xi_j \in \mathbf{P}(j)$.*

```

for  $j = n, \dots, 1$ 
  if  $j \in \mathbf{P}_i$ 
    if  $(j \neq n)$  receive( $s$ )
     $\xi_j \leftarrow (\sigma_1 + \tau_j)/v_{jj}$ 
    for  $k = 1, \dots, p - 2$ 
       $\sigma_k = \sigma_{k+1} - v_{j-k-1,j}\xi_j + \tau_{j-k-1}$ 
     $\sigma_{p-1} = -v_{j-p+1,j}\xi_j + \tau_{j-p+1}$ 
    if  $(j \neq 1)$  send left( $s$ )
    for  $k = 1, \dots, j - p + 1$ 
       $\tau_k = \tau_k - v_{kj}\xi_j$ 

```

Figure 3.15: Parallel Back-Substitution

- send s to $\mathbf{P}(j - 1)$, ($j \neq 1$),
- use ξ_j to update the remaining elements of t .

This algorithm is presented in Figure 3.15 in greater detail.

Consider the following example where $p = 4$ and $n = 16$. Figure 3.16 shows the assignment of processors to columns as well as the data flow between processors via the vector s . Each column represents the elements of the vectors t and s on the processor denoted above the column during the iteration indicated below the column. In the j th column ($j = 16, \dots, 1$), the “x” identifies the element of s used to determine ξ_j . The three elements labeled “s” identify the elements of the buffer s which are updated and sent to $\mathbf{P}(j - 1)$. The position of each letter in the triangle corresponds to the position of the element of U which is used to update that particular element of s , t or x in that iteration.

Let us examine how the data becomes available to \mathbf{P}_1 so that it can compute ξ_6 . The equation (3.6) shows the data elements needed to determine ξ_6 . The processors on which these elements reside originally are indicated in subscripts, and the elements of t and s in which these terms are collected are highlighted in Fig-

buffer s around the ring of processors as the algorithm loops through each cycle. Each cycle then corresponds to the passage of s once around the ring of processors. In the first cycle, the elements $\xi_n, \dots, \xi_{n-p+1}$ are determined, and in the second cycle, the elements $\xi_{n-p}, \dots, \xi_{n-2p+1}$ are determined and so on until ξ_p, \dots, ξ_1 are determined in the m th cycle. Each cycle (except the first) begins with $\mathbf{P}_{p-1} = \mathbf{P}(n)$ receiving the vector s , determining the current element of x and then updating and sending the vector s to \mathbf{P}_{p-2} before updating the vector t . Each cycle (except the last) ends when \mathbf{P}_0 sends s to \mathbf{P}_{p-1} for it to begin the next cycle of iterations. When s arrives at \mathbf{P}_{p-1} , this processor may begin using and updating the vector immediately or it may have to finish updating the vector t from its previous iteration before initiating the next cycle of iterations.

The buffer, s , must pass completely around the ring exactly $(m - 1)$ times and then in the m th pass it travels all but the last link of the ring. The time complexity of the algorithm is the time for s to circle the ring $(m - 1)$ times plus time for the last partial circuit plus any delays between cycles. During a complete pass around the ring each processor performs p floating point operations and then sends the buffer to its neighbor. This takes time equal to $p\gamma + \alpha + p\beta$, making the time for a complete cycle, $p^2\gamma + p(\alpha + p\beta)$. The time for s to make $(m - 1)$ complete cycles is then

$$\begin{aligned} T_{complete} &= (m - 1)(p^2\gamma + p(\alpha + p\beta)) \\ &= (n - p)(\alpha + p\beta + p\gamma). \end{aligned}$$

The time to complete the last partial circuit must account for the shrinking of the vector s in each iteration and is given by,

$$\begin{aligned} T_{partial} &= \left(\sum_{j=0}^{p-1} (p - j)\gamma \right) + (p - 1)(\alpha + p\beta) \\ &= \frac{p(p + 1)}{2}\gamma + (p - 1)(\alpha + p\beta). \end{aligned}$$

Li and Coleman show that only $\mathbf{P}(n) = \mathbf{P}_{p-1}$ can delay s as it is passed around the ring of processors. It suffices to consider only this processor when determining the total time that s is delayed. In the i th cycle ($i = 1, \dots, m$), $\mathbf{P}(n)$ performs $(n - ip)$ flops to update the remainder of t . The time for s to circle the ring and return to the sender is $p(\alpha + p\beta) + p(p - 1)\gamma$. Whenever the time for $\mathbf{P}(n)$ to

p	$n_{delay}(p)$
2	79
4	180
8	449
16	1491
24	2680

Table 3.1: Column index below which delay is zero

update t exceeds the time for s to return to $\mathbf{P}(n)$, the vector s is delayed. Therefore the start of the $(i + 1)$ th cycle is delayed whenever

$$(n - ip)\gamma > p(\alpha + p\beta) + p(p - 1)\gamma,$$

or when

$$i < k \quad \text{where } k = \left(\frac{n}{p} - \frac{\alpha + p\beta}{\gamma} - (p - 1) \right),$$

The column index below which the delay is zero is a function of p , (see Table 3.1) and is found by letting $i = k - 1$,

$$n_{delay}(p) = n - p(k - 1) = \frac{p(\alpha + p\beta)}{\gamma} + p^2.$$

The entries in Table 3.1 can also be interpreted as being the problem size above which communications is almost completely overlapped by computation.

If the problem size is less than or equal to $n_{delay}(p)$, then no delay occurs and the time complexity is,

$$\begin{aligned} T_p(n) &= T_{complete} + T_{partial} & n \leq n_{delay}(p) \\ &= (n - 1)(\alpha + p\beta) + \left(n - \frac{p - 1}{2} \right) p\gamma. \end{aligned} \quad (3.7)$$

If the problem size is greater than $n_{delay}(p)$ then a delay occurs and the delay is $d_i = (n - ip)\gamma - p(\alpha + p\beta) - p(p - 1)\gamma$. Using the approximation that $[k] - 1 \approx k$ the total time that s is delayed is,

$$\begin{aligned} T_{delay} &\approx \sum_{i=1}^k d_i \\ &= \left(\sum_{i=1}^k (n - ip)\gamma \right) - k(p(\alpha + p\beta) + p(p - 1)\gamma) \\ &= \frac{1}{2}p\gamma k(k - 1). \end{aligned}$$

For problems larger than $n_{delay}(p)$ the time complexity is,

$$\begin{aligned} T_p &= T_{complete} + T_{partial} + T_{delay} && n > n_{delay}(p) \\ &\approx \left(\frac{pk^2}{2} + np - \frac{p^2}{2} \right) \gamma + (n-1)(\alpha + p\beta). \end{aligned}$$

Noting that $k \approx n/p$ for large n and ignoring low order terms we get

$$T_p(n) \approx \left(\frac{n^2}{2p} \right) \gamma + (n-1)(\alpha + p\beta) \quad n > n_{delay}(p). \quad (3.8)$$

For a fixed value of p , as the problem size becomes very large we have

$$\lim_{n \rightarrow \infty} \frac{T_1(n)}{T_p(n)} = p,$$

so the speedup approaches p and efficiency nears 100%, but as indicated in Table 3.1 these values are prohibitively large for even a moderate number of processors.

3.3.3 Implementation and Numerical Experiments

The routine PTRSL (**p**arallel **t**riangular **s**olve) drives each node in a ring of processors to solve the upper triangular system of equations $Ax = b$ by back substitution. The calling sequence and argument descriptions are given in Figure 3.17.

The observed execution times for two, four and sixteen processors to solve problems ranging from $n = 10$ to $n = 500$ are plotted in Figure 3.18. The transition from linear to quadratic behavior is apparent for two and four processors. This transition occurs when the problem size becomes large enough for T_{delay} to become non-zero. In the case of sixteen processors this requires n to be greater than 1400 and hence the timings remain linear in the range of problem sizes tested (see Table 3.1).

Figure 3.19 exhibits the agreement between observed and theoretical run times. Again, discrepancies between these times are explained by the combined effects of several factors. Figure 3.20 plots the observed efficiencies attained by the parallel algorithm for various ring and problem sizes. Reasonable efficiencies are attained for two and four processors once the problem size has grown large enough to mask the expense of sending the buffer around the ring with the useful computation of updating the t vector.

PTRSL(N,A,X,T,S,INFO)

- On Entry

N is the order of the matrix to be factored.

A is a pointer to a vector of type FLT which contains columns of the upper triangular matrix.

T is a pointer to a vector of type FLT which contains the first $(n - p + 1)$ elements of b on $\mathbf{P}(n)$ and 0 on all other processors.

S is a pointer to a vector of type FLT which contains the last $(p - 1)$ elements of b on $\mathbf{P}(n)$ and 0 on all other processors.

- On Return

A is a pointer to a vector of type FLT which contains columns of the unchanged upper triangular matrix.

X is a pointer to a vector of type FLT which contains the solution. The full solution vector is distributed among the processors such that $\xi_j \in \mathbf{P}(j)$.

INFO is a pointer to an integer, which has value 0 if the matrix A is non-singular or has value k if the k th diagonal element of A is zero.

Figure 3.17: PTRSL

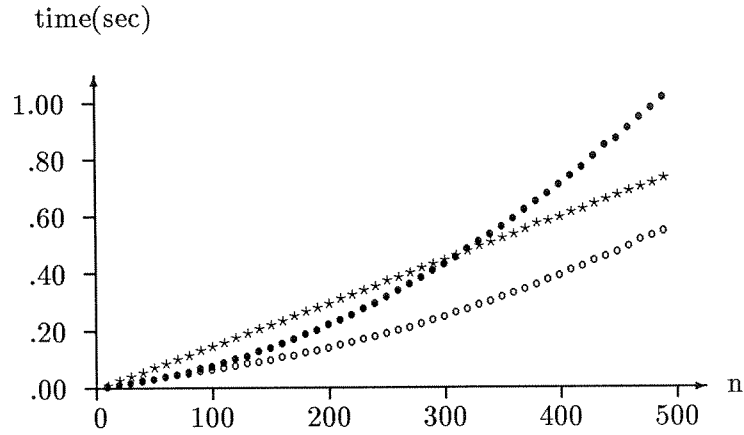


Figure 3.18: Observed timings attained during backward substitution
(key:: ●: $p = 2$, ○: $p = 4$, *: $p = 16$)

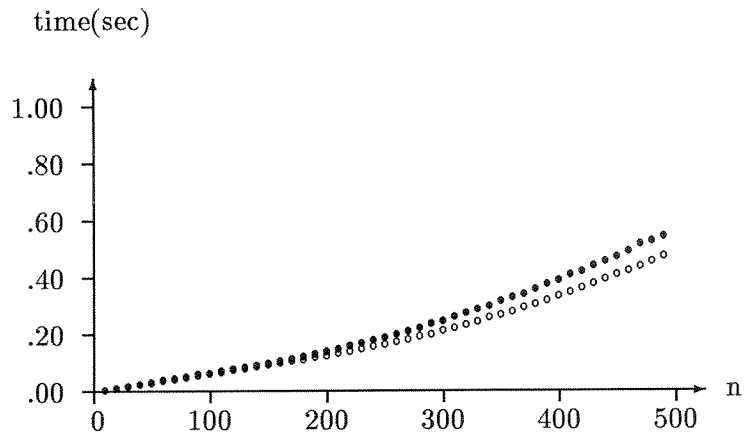


Figure 3.19: Observed and expected timings for $p = 4$ during backward substitution
(key:: ●: observed, ○: expected)

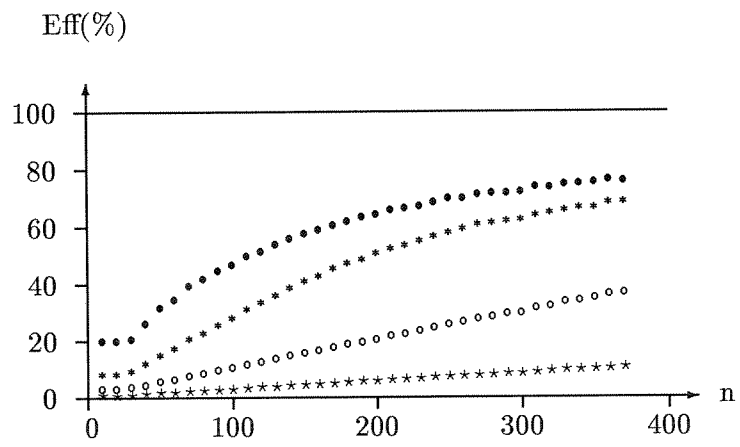


Figure 3.20: Observed efficiencies attained during backward substitution
(key:: ●: $p = 2$, *: $p = 4$, ○: $p = 8$, ☆: $p = 16$)

Chapter 4

Q-R Factorization

One application of the algorithms in this chapter is concerned with solving the least squares problem of minimizing $\|Ax - b\|_2$, where $A \in \mathbf{R}^{m \times n}$ for $m > n$ and $b \in \mathbf{R}^m$. We will examine the Householder and Gram-Schmidt methods for transforming this problem into an equivalent yet easier problem to solve. Householder orthogonalization factors A into the product of a unitary matrix Q and an upper trapezoidal matrix R . The Gram-Schmidt method factors A into the product of a matrix $Q \in \mathbf{R}^{m \times n}$ with orthonormal columns and $R \in \mathbf{R}^{n \times n}$, an upper triangular matrix with positive diagonal elements.

Once computed we can use this QR factorization of A to transform the least squares problem as follows,

$$Q^T A = R = \begin{bmatrix} R_1 & \\ & 0 \end{bmatrix} \begin{matrix} n \\ m - n \end{matrix}$$

where R_1 is upper triangular and

$$Q^T b = \begin{bmatrix} c \\ d \end{bmatrix} \begin{matrix} n \\ m - n \end{matrix}.$$

The least squares problem is then to minimize,

$$\begin{aligned} \|Ax - b\|_2^2 &= \|Q^T Ax - Q^T b\|_2^2 \\ &= \|R_1 x - c\|_2^2 + \|d\|_2^2 \end{aligned}$$

for any $x \in \mathbf{R}^n$. If $\text{rank}(A) = \text{rank}(R_1) = n$ then a unique x_{LS} which minimizes $\|Ax - b\|_2$ exists and is defined as the solution to the upper triangular system,

$$R_1 x_{LS} = c.$$

4.1 Householder Orthogonalization

A Householder transformation or elementary reflector has the form

$$H = I - \beta vv^T,$$

where $\beta = 2/v^T v$. Householder matrices are symmetric, orthogonal and involutory, and can be used to introduce a contiguous block of zeros into a vector. In fact if $x \in \mathbf{R}^n$, the Householder transformation,

$$H = I - 2 \frac{vv^T}{v^T v}$$

where

$$v = x \pm \|x\|_2 e_1,$$

has the property that

$$Hx = \mp \|x\|_2 e_1.$$

The vector v is chosen to be

$$v = x + \text{sign}(\xi_1) \|x\|_2 e_1,$$

to avoid introducing large relative error in the factor $\beta = 2/v^T v$ when x is close to a multiple of e_1 . This guarantees nearly perfect orthogonality in the computed H [2].

4.1.1 Sequential Algorithm

This algorithm (see Figure 4.1) factors $A \in \mathbf{R}^{m \times n}$ into the product QR where $Q \in \mathbf{R}^{m \times m}$ is unitary and $R \in \mathbf{R}^{m \times n}$ is upper trapezoidal, by applying Householder transformations to introduce zeros below the diagonal of A .

The algorithm begins with $A = A^{(1)} = (a_1^{(1)}, a_2^{(1)}, \dots, a_n^{(1)})$ and proceeds by first determining a Householder matrix, H_1 , such that

$$H_1 a_1^{(1)} = -\text{sign}(\alpha_{11}^{(1)}) \|a_1^{(1)}\|_2 e_1$$

and applying it to form

$$A^{(2)} = H_1 A^{(1)}.$$

Algorithm 4.2 *This algorithm drives processor \mathbf{P}_i , a node in a ring of processors to form the orthogonal QR decomposition of a general rectangular matrix, A , by applying Householder transformations. The matrix is distributed among the processors in a column wrapped fashion.*

```

for  $k = 1, \dots, n$ 
  if  $k \in \mathbf{P}_i$ 
    determine  $\tilde{H}_k : \tilde{H}_k(\alpha_{kk}, \dots, \alpha_{mk})^T = \rho_{kk} e_1$   $[(m - k + 1)\gamma]$ 
    broadcast  $v$  and  $\beta$   $[2(\alpha + (m - k)\beta)]$ 
  else
    receive  $v$  and  $\beta$ 
    update  $d_j \leftarrow d_j - \beta(v^T d_j)v, j \in \mathbf{P}_i > k$   $[(2\lceil \frac{n-k}{p} \rceil)(m - k + 1)\gamma]$ 

```

Figure 4.2: Parallel Householder Orthogonalization

of updating this submatrix is approximately $2(n - k)(m - k)\gamma$. The complexity of this algorithm is then,

$$T_1(m, n) \approx \sum_{k=1}^n 2(n - k + 1)(m - k)\gamma \approx n^2 m - \frac{n^3}{3}. \quad (4.1)$$

4.1.2 Parallel Algorithm

The parallelization of the Householder orthogonalization algorithm proceeds almost identically to that for Gaussian elimination. Again we distribute the matrix A to the ring of processors in a column wrapped fashion.

Given that each processor has approximately (n/p) columns, we consider the sequential algorithm which loops through a set of computations for each of the n columns of the matrix. Within each iteration, the sequence of computations is to determine \tilde{H}_k and then apply it from the left side to the $(m - k) \times (n - k)$ submatrix D_k . The Householder transformation, \tilde{H}_k , can be completely determined by processor $\mathbf{P}(k)$. The submatrix, D_k , however, is distributed among all processors, so all processors are required to assist in applying the update. Once \tilde{H}_k is computed, $\mathbf{P}(k)$ can broadcast v and β to the other processors and all processors can simultaneously update their portions of D_k .

The expressions in brackets, in Figure 4.2, indicate the effective contribution to the time complexity of that step in each iteration. As in the sequential

algorithm, computing the Householder transformation requires $(m - k + 1)\gamma$ time, and as was the case in both Gaussian elimination, and the Cholesky decomposition the effective time to perform the broadcast is equivalent to the time to perform two node-to-neighbor communications. The message here is a vector of $(m - k)$ double precision floating point numbers so the time to communicate in each iteration is $2(\alpha + (m - k)\beta)$. The $(k + 1)$ th iteration begins as $\mathbf{P}(k + 1)$ completes updating the columns of D_k that it owns. Processor $\mathbf{P}(k + 1)$ owns $\lceil (n - k)/p \rceil$ columns of this submatrix. The update of each column requires a dot product and saxpy operation each of order $(m - k)$.

Using the approximation $\lceil (n - k)/p \rceil \approx (n - k)/p + 1$, the time for $\mathbf{P}(k + 1)$ to update its share of the submatrix is then,

$$2 \left(\frac{n - k}{p} + 1 \right) (m - k)\gamma.$$

Summing these three expressions over all iterations gives the total time complexity of the algorithm,

$$\begin{aligned} T_p(m, n) &\approx \sum_{k=1}^n \left[(m - k + 1)\gamma + 2 \left(\frac{n - k}{p} + 1 \right) (m - k)\gamma \right. \\ &\quad \left. + 2(\alpha + (m - k)\beta) \right] \\ &\approx \left(n^2 m - \frac{n^3}{3} \right) \frac{\gamma}{p} + 2n \left(m - \frac{n}{2} \right) \gamma \\ &\quad + 2n \left(\alpha + \left(m - \frac{n}{2} \right) \beta \right). \end{aligned} \tag{4.2}$$

Again, the complexity consists of a dominating term which compares favorably with the sequential time complexity, given by (4.1). For a fixed value of p , as the problem size increases we have

$$\lim_{n \rightarrow \infty} \frac{T_1(n)}{T_p(n)} = p,$$

so the speedup approaches p and efficiency nears 100%. The complexity also contains a term resulting from the sequential computation and a term reflecting communication costs. Both of these terms become significant and reduce efficiency as p approaches n .

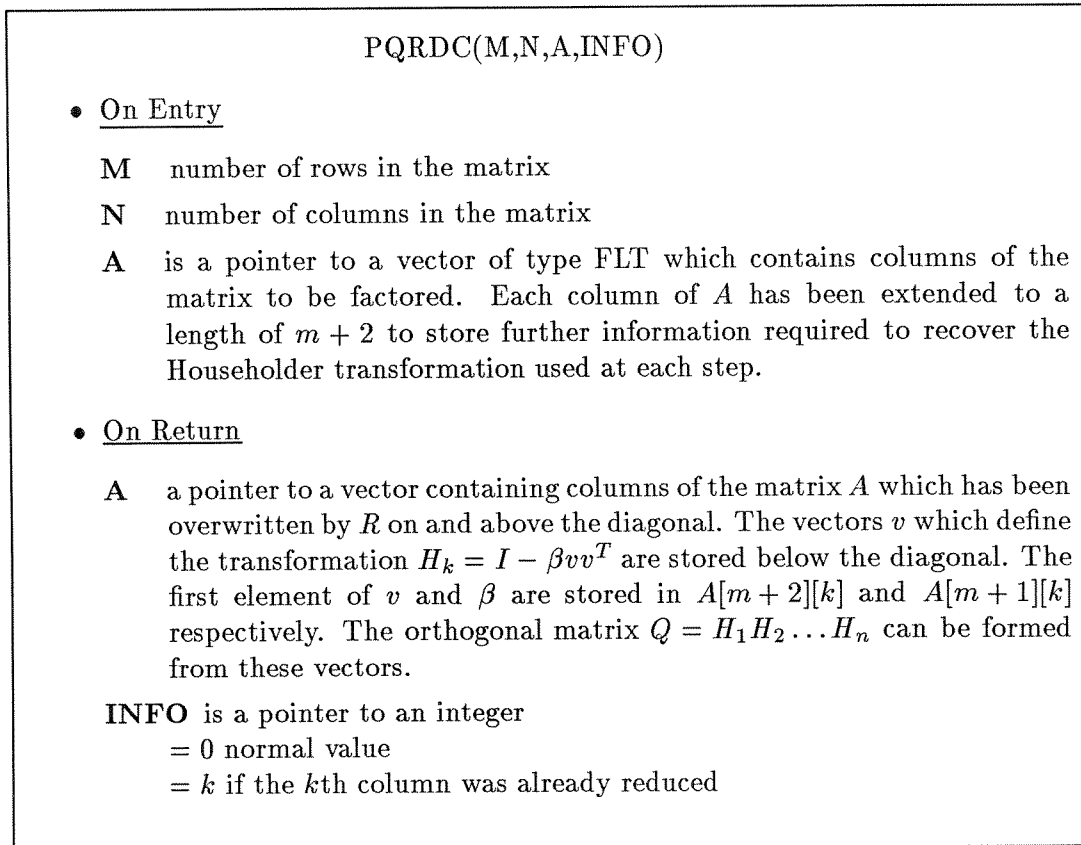


Figure 4.3: PQRDC

4.1.3 Implementation and Numerical Experiments

The routine PQRDC (**parallel Q-R decomposition**) drives each node in a ring of processors to form the orthogonal QR decomposition of a general rectangular matrix, A , by applying Householder transformations. The matrix is distributed among the processors in a column wrapped fashion. The calling sequence and argument descriptions are given in Figure 4.3.

Figure 4.4 plots the expected and observed timings obtained by the parallel algorithm with a ring of sixteen processors on the Symult S2010 for $n = 10, 20, \dots, 490$. For these timings, m was set to $n + 2$. Again, discrepancies between these times are explained by the combined effects of several factors, this time including the overestimation of the cost of an inner product.

Figure 4.5 plots the observed efficiency attained by the parallel algorithm for various ring sizes and problem sizes. Reasonable efficiencies can be obtained