

The background of the entire page is a close-up photograph of a server rack. It features a grid of vertical server bays. Each bay has a top section with two indicator lights (one red, one green) and a bottom section with a label. The labels are partially visible and read 'SC' and 'iPSC'.

# iPSC™ Technical Report

No. 2  
Communication Utilities  
for the iPSC™

SC

iPSC

iPSC

iPSC

iPSC

iPSC

iPSC

iPSC

iP

# Communication Utilities For The iPSC™

Cleve Moler

David S. Scott

Intel Scientific Computers

August 1986

## ABSTRACT

A library of communication subroutines for the iPSC is described which provides facilities for global message broadcast, global sums and other commutative arithmetic operations, embedding of meshes into the hypercube, and very long messages. The descriptions correspond to the FORTRAN versions of the codes. Versions written in C are also available.

## 1. INTRODUCTION

The communication primitives in the basic iPSC operating system provide for sending a message from one single process to another. We have found it useful in applications programming to expand these facilities to provide some higher level constructs for interprocess communication. Included are:

- Global Broadcast -- send a message from one node to all other nodes.
- Global Operations -- commutative arithmetic operations, such as sum or maximize, with contributions from each node.
- Global Concatenation -- collect messages from each node.
- Meshes -- provide each node with left - right, or north - east - south - west, neighbors.
- Long messages -- handle messages longer than 16K bytes.
- Sparse Global Send -- Send and receive with selected subset of nodes.

The iPSC system send and receive calls use both a node number and a process id to specify the destination of a message. This allows multiprocessing of the nodes. However, the utilities described here assume only one participating process on each node. For simplicity it is assumed that each participating process has the same PID. For this reason the PID is not included in the calls. It is obtained internally by a call to MYPID().

---

iPSC™ is a trademark of Intel Corporation

Several of the subroutines allow operations involving set of nodes called subcubes. By a subcube of dimension  $d$  we mean the set of  $p = 2^d$  nodes whose addresses differ in their low order  $d$  bits. For example, the subcube of dimension 3 containing node 13 consists of nodes 8 through 15. Note that such subcubes correspond to neighboring sets of boards in the iPSC hardware; a 5 dimensional subcube is a full 32 node iPSC system (or cabinet), a 3 dimensional subcube is either the left or the right half of a horizontal row of boards, and so on.

An example problem which uses some of the subroutines is given after all of the subroutine descriptions.

## 2. GLOBAL BROADCAST

This section describes a procedure for one node (called the root) to distribute the same message to all the other nodes. It is also possible to limit the broadcast to a lower order subcube, or to initiate the broadcast from the host. It is assumed that each participating process has the same PID and that the host knows this common PID if it is the root. The receiving processes call GRECVW while the root calls GSENDW. If the host is the root, it just uses SENDMSG to send the message to any one of the receiving processes.

The calling sequences are:

```
SUBROUTINE GSENDW (CI, TYPE, BUF, LEN, DIM)
INTEGER CI, TYPE, LEN, DIM
BYTES BUF(*)
```

and

```
SUBROUTINE GRECVW (CI, TYPE, BUF, LEN, CNT, DIM)
INTEGER CI, TYPE, LEN, CNT, DIM
BYTES BUF(*)
```

The parameters are

- CI      The channel number to be used (already opened).
- TYPE   The message type (must be the same for all calls).
- BUF    The message buffer. The type of BUF is irrelevant.
- LEN    The length of the message in bytes for sending.  
        The length of the buffer in bytes for receiving.
- CNT    Returned as the length (in bytes) of the message received  
        NOTE: CNT > LEN is an error return condition since only  
        LEN bytes will be stored in BUF and the rest of the  
        message will be lost. In this case an error message  
        will be syslogged but the remnants of the message  
        will be sent on.
- DIM    The dimension of the subcube participating in the send  
        NOTE: only those nodes which will participate in the  
        send can call GRECVW. Any other node which calls it  
        will never return.

The implementation of GSENDW and GRECVW involves a subcube spanning tree, rooted at the root node. This is also known as "induction on the cube dimension." GSENDW sends the message from the root to all of its nearest neighbors. In GRECVW, if a node receives the message from a neighbor, it then sends it on to some, but not all, of the remaining neighbors. If it receives it from the host, it sends it on to all of its neighbors. The hypercube neighbors of a node are those nodes whose addresses differ from the

given address in exactly one bit. The list of neighbors can be ordered by increasing bit number. In [1] it was shown that each node should send the data to all neighbors on this list after the node they receive it from. In particular, it is not necessary for the receiving nodes to know the root of the broadcast.

It is likely that future versions of the iPSC node operating system will offer global broadcast as a primitive operation. Calling SEND or SENDW with NID = -1 will initiate the broadcast of a message to all the nodes in the cube. Other negative destination node numbers will initiate broadcasts to subcubes. The message will be received with an RECV or RECVW and the operating system will automatically manage the spanning tree.

### 3. GLOBAL COMMUTATIVE OPERATIONS

Sometimes all the nodes have some data that must be accumulated. The most common example is computing a distributed innerproduct. In these situations an inverse or dual of the global send is needed. All the participants call the subroutine GOP. Four different operations are included in the distributed version: addition, multiplication, maximum, and minimum. It is easy to modify the source code to include other operations. If the operation is applied to a vector, then it is applied component-wise to the elements. The implementation allows accumulation over a lower order subcube. If root is specified as the host, then the message is accumulated in the lowest numbered node and then sent on to the host (with the same TYPE and same PID). The calling sequence is:

```
SUBROUTINE GOP (CI, TYPE, X, N, OP, ROOT, DIM, WORK)
INTEGER CI, TYPE, N, ROOT, DIM
CHARACTER*1 OP
DOUBLE PRECISION X(N), WORK(N)
```

#### Input

CI Channel number (previously opened)

TYPE Message type. Must be the same for all participating processes. There must be no other messages of this type in the system.

X The input vector to be used in the operation.

N The length of the vectors.

OP '+' sum  
 '\*' product  
 'M' maximum  
 'm' minimum

ROOT Node number of root process. (-32768 for the host)

DIM The size of the subcube participating.

#### On output

X For the root process, X contains the desired result.  
For all other processes, X contains the partial result for their subtrees.

#### Workspace

WORK Used to receive other contributions.

#### 4. GLOBAL CONCATENATION

Closely related to the global commutative operation is the need to collect a distributed vector in one processor. The subroutine provided for this purpose is GCAT, which globally concatenates a vector. It is possible to concatenate over a lower order subcube.

The actual type of BUF is irrelevant, it is concatenated by byte. If the root is zero, then the contributions are in ascending order (ie. nodes 0, 1, 2, 3,...). If any other node is root, then the order obtained is the exclusive-or of the ascending order with the root. If the root is the host (-32768), the lowest numbered node acts as the root, but then sends the message on to the host (with the same type and pid). Many participating nodes receive contributions from several children. The resulting concatenation is sensitive to the order in which these contributions are received. To insure that the contributions are received in order it is necessary to use type numbers in addition to the specified number TYPE. In particular all the type numbers TYPE, TYPE+1, TYPE+2, ... TYPE+DIM-1 are used.

The calling sequence is:

```
SUBROUTINE GCAT (CI, TYPE, BUF, BUFLen, MSGLEN, CNT, ROOT, DIM)
INTEGER CI, TYPE, BUFLen, MSGLEN, CNT, ROOT, DIM
BYTES BUF(*)
```

Input:

CI	Channel number (previously opened).
TYPE	Message type. Must be the same for all participating processes. There must be no other messages of this type or of types TYPE+1, TYPE+2, ..., TYPE+DIM-1 in the system.
BUF	The input vector to be concatenated. The actual type of BUF is irrelevant.
BUFLen	The length of BUF (in bytes).
MSGLEN	The length of this contribution (in bytes).
ROOT	Node id of root process (which will get the full message).
DIM	The size of the subcube participating.

Output:

BUF	For the root process, BUF contains the desired result. For all other processes, BUF contains the partial result for their subtrees.
CNT	The number of meaningful bytes in BUF.

#### 5. EMBEDDING MESHES

Many applications require computations on a rectangular mesh. It is natural to subdivide the mesh into large rectangles and assign one processor to each rectangle. If rectangles only need to communicate with adjacent rectangles, it is desirable to map adjacent rectangles onto nearest neighbors in the hypercube. MESH computes an embedding of a rectangular mesh of nodes into the hypercube that preserves nearest neighbor connections. The subroutine takes the dimensions of the mesh, the boundary conditions for the mesh, and the node number of interest and computes the coordinates of the node and the node numbers of its nearest neighbors in each coordinate direction. A node number of -1 is used as NIL to indicate a

nonexistent neighbor for a boundary node. If the given node number is not used in the specified mesh, then all of the output parameters are set to -1.

The calling sequence is

```
SUBROUTINE MESH (ME, DIM, LINDIM, PERBC, COORDS, PRED, SUCC, WORK)
INTEGER ME, DIM, LINDIM(DIM), PERBC(DIM), COORDS(DIM),
      PRED(DIM), SUCC(DIM), WORK(DIM)
```

where

Input:

ME	The node number of interest.
DIM	The physical dimension of the mesh (usually 1, 2, or 3).
LINDIM	The linear dimension (number of processors) in each direction of the mesh.
PERBC	A flag indicating whether a particular direction has periodic boundary conditions. PERBC should be set as follows:  PERBC(I) = 0 for non periodic (closed) boundaries PERBC(I) = 1 for periodic boundaries

Output:

COORDS	The coordinates of ME in the mesh, starting with 0, so that the largest the I-th coordinate can be is LINDIM(I)-1. If ME is not part of the mesh all coordinates are returned as -1.
PRED	PRED(I) is the node number preceding ME in the I-th direction. -1 is used when there is no predecessor.
SUCC	SUCC(I) is the node number succeeding ME in the I-th direction. -1 is used when there is no successor.

Workspace

WORK	holds the result of rounding the elements of LINDIM up to the next power of two.
------	--

If messages must be sent diagonally in the mesh, then the diagonal neighbors can be identified by calls to mesh with the neighbor's node numbers.

Note that if the mesh is too large to fit in the hypercube, MESH will send an error message to the system log, but it will then return the correct (but perhaps nonexistent) neighboring node numbers.

## 6. LONG MESSAGES

The messages processed by the primitives SEND and RECV are limited to a maximum of 16K bytes. Three subroutines are provided to allow node-node and host-node transmission of long messages. A node-host utility is impossible since messages cannot be filtered by type on the host.

The three calling sequences are

For sending from a node:

```
SUBROUTINE LSENDW (CI, TYPE, BUF, LEN4, NODE, PID)
INTEGER CI, TYPE, NODE, PID
INTEGER*4 LEN4
BYTES BUF(LEN4)
```

For sending from the host:

```
SUBROUTINE LSENDMSG (CI, TYPE, BUF, LEN4, NODE, PID)
INTEGER CI, TYPE, NODE, PID
INTEGER*4 LEN4
BYTES BUF(LEN4)
```

For receiving (on a node):

```
SUBROUTINE LRECVW (CI, TYPE, BUF, LEN4, CNT4, NODE, PID)
INTEGER CI, TYPE, NODE, PID
INTEGER*4 LEN4, CNT4
BYTES BUF(LEN4)
```

Where

- CI     A channel number (previously opened).
- TYPE   The type of the message.
- BUF    The message buffer (which may actually be any type).
- LEN4   INTEGER\*4. The length of the message in bytes for sending.  
        The length of BUF in bytes for receiving.
- CNT4   INTEGER\*4. The length of the received message.
- NODE   The source node.
- PID    The source pid.

The subroutines send the message in 16K byte chunks and terminate with a final chunk which is smaller than 16K bytes. If the message is exactly a multiple of 16K bytes, then a zero length message is sent, however the receiving buffer must be large enough to receive one more byte (to prevent an illegal subscript when the receive is called).

## 7. SPARSE GLOBAL SEND

In some applications a specific node may have a message which must be sent to a specified subset of the other nodes. Two subroutines are provided for this purpose. SPGSNW initiates a sparse send and all receiving nodes must call SPGRCW. The sending node specifies the target nodes by a list of integers. The calling sequences are

For sending:

```
SUBROUTINE SPGSNW (CI, TYPE, BUF, LEN, LIST, NUMBER)
INTEGER CI, TYPE, LEN, LIST(*), NUMBER
BYTES BUF(*)
```

For receiving:

```
SUBROUTINE SPGRCW (CI, TYPE, BUF, LEN, CNT)
INTEGER CI, TYPE, LEN, CNT
BYTES BUF(*)
```

The parameters are

CI      An open channel.

TYPE    A legal message type (must be the same for all participating calls)

BUF     The buffer.

LEN     For sending: the length of the message in bytes.  
         For receiving: the length of BUF.  
         \*\*\*Note: SEE RESTRICTIONS BELOW\*\*\*

CNT     The length of the received message in bytes.

LIST    The list of node numbers which should receive the data.  
         The list is destroyed.

NUMBER      The number of node numbers in list.

Restrictions on use:

SPGSNW appends a list of integer targets to the actual message. Thus the sending buffer BUF must be long enough to hold these integers. It is sufficient if the length of the sending buffer is at least  $LEN + INTSIZE * NUMBER$  where INTSIZE is the number of bytes per integer.

SPGRCW receives this appended message and so the LEN of the receiving buffer must be at least  $CNT + INTSIZE * NUMBER$ .

Finally, to simplify the appending process and avoid possible alignment problems, the length of the message is implicitly rounded up to a multiple of INTSIZE. After the integers are stripped off this rounded length is returned as CNT. Therefore to avoid the accidental inclusion of garbage bytes, the sending LEN should also be a multiple of INTSIZE.

If the initiating node is called the **root**, and all the other participating nodes are called **targets**, then the sparse global send does not just send from the root to all of the targets. Instead, an attempt is made to construct an optimal fanout tree. SPGSNW calls SPLIT to choose a first target for the message. In addition to choosing this target, SPLIT divides the remaining targets between those to be processed by the first target and those to be processed by the root. The list of nodes for the target is appended to the message and sent. The rest of the nodes are iteratively processed by the root.

## 8. ERROR MESSAGES

Some utilities return error messages if called improperly. All such messages are syslogged, are in CAPITALS, and begin with the routine name. The messages are

GCAT: CALLED BY NON PARTICIPANT

GCAT: BUFFER LENGTH EXCEEDED

GOP: CALLED BY NON PARTICIPANT

GOP: LONG MESSAGE

Received too long a message

GOP: SHORT MESSAGE

Received too short a message

GRECVW: MESSAGE TRUNCATED

LRECVW: TRUNCATED MESSAGE

MESH: TOO LARGE

The mesh is too large to fit in the hypercube

SPGRCW: TRUNCATED MESSAGE

## 9. INTERNAL SUBROUTINES

The utilities library includes five subprograms that are used internally by the main procedures. They are documented here for completeness and because some of them may be useful in their own right. The modules are:

LOGICAL FUNCTION BETWEEN (I, J, K)  
INTEGER I, J, K

BETWEEN decides whether node J is between node I and node K, that is whether node J lies on a shortest path in the hypercube topology from node I to node K. This function is used by SPLIT.

INTEGER FUNCTION GRAY (I)  
INTEGER I

GRAY returns the Ith integer in the binary reflected gray code. (This code orders the node numbers so that consecutive numbers are adjacent nodes in a hypercube.)

INTEGER FUNCTION GINV (I)  
INTEGER I

GINV is the inverse of GRAY. It returns the location of I in the gray code ordering.

SUBROUTINE SPLIT (NUMBER, LIST, FIRST)  
INTEGER NUMBER, LIST(NUMBER), FIRST

SPLIT is used by SPGSNW (sparse global send). LIST is the list of NUMBER node numbers (**targets**) that must receive the message. Given the source node number r, the hypercube topology induces a partial ordering on the nodes where  $b < c$  is defined as b lies on a shortest path from r to c. SPLIT takes the list of targets, finds a minimal target node k, and then splits the remaining targets into the set of all nodes d for which  $k < d$  and the complement. FIRST is returned as the number of targets that are greater than k. These targets are put at the beginning of LIST, followed by initial target number k, followed by the

remaining targets.

```
INTEGER FUNCTION XOR (I,J)
INTEGER I, J
```

XOR returns the bitwise exclusive-or of integers I and J. Most FORTRAN compilers have a function that computes the exclusive-or but the name is different for different compilers. This function invokes the correct name. XOR is useful since XOR(I, 2\*\*K) is the K-th neighbor of node I in the hypercube.

## 10. EXAMPLE PROBLEM

Suppose a vector X(N) is distributed among the P processors of a hypercube. For convenience assume that  $N = P * M$ . Assume that this vector must be normalized and then the entire vector must be available to node 0. The following code fragment will accomplish these tasks.

```
INTEGER I, CI, GOTYPE, GSTYPE, GCTYPE, CUBEDIM, COPEN, MYNODE, CNT, ROOT
DOUBLE PRECISION X(M), Y(N), DOT, NORM, DUMMY
PARAMETER (DPSIZE = 8)
DATA GOTYPE /100/, GSTYPE /200/, GCTYPE /300/
C
  CI = COPEN(MYPID())
C
C Assume that X, M, and N have been
C   initialized
C
  ROOT = 0
C
C Compute the local dot product
C
  DOT = 0
  DO 1 I = 1, M
    DOT = DOT + X(I)*X(I)
  1 CONTINUE
C
C Sum local contributions in ROOT
C
  CALL GOP (CI, GOTYPE, DOT, 1, '+', ROOT, CUBEDIM(), DUMMY)
C
C Distribute norm to all the processors
C
  IF (MYNODE() .EQ. ROOT) THEN
    NORM = DSQRT (DOT)
    CALL GSENDW (CI, GSTYPE, NORM, DPSIZE, CUBEDIM())
  ELSE
    CALL GRECVW (CI, GSTYPE, NORM, DPSIZE, CNT, CUBEDIM())
  ENDIF
C
C Do the local normalization
C
  DO 2 I = 1, M
    X(I) = X(I)/NORM
  2 CONTINUE
C
C Copy X to the concatenation buffer
C
```

```

DO 3 I = 1,M
  Y(I) = X(I)
3 CONTINUE
C
C Concatenate the vector in ROOT
C
  CALL GCAT(CI, GCTYPE, Y, N*DPSIZE, M*DPSIZE, CNT, ROOT, CUBEDIM())

```

The array Y in Node 0 contains all of the vector X. Whether or not X is properly sorted depends on how the vector X was distributed on the nodes. If Node 0 contained X(I), I = 1, M and Node 1 contained X(I), I = M+1, 2\*M etc., then X would be properly sorted.

Note that the naive way in which the norm is computed in this example can lead to problems with numerical underflow or overflow. Safer and more complicated implementations are possible and ought to be used in practice.

## 11. INDEX OF SUBROUTINES

BETWEEN	hypercube ordering
GCAT	global concatenation
GINV	gray code inverse
GOP	global commutative operation
GRAY	gray code
GRECVW	global receive
GSENDW	global send from node
LRECVW	long message receive
LSENDMSG	long message send from host
LSENDW	long message send from node
MESH	nearest neighbor mesh embedding
SPGRCW	sparse global receive
SPGSNW	sparse global send
SPLIT	node list separator
XOR	bitwise exclusive-or

## 12. REFERENCES

- [1] J. E. Brandenburg and D. S. Scott  
 "Embeddings of Communication Trees and Grids into Hypercubes,"  
 iPSC Technical Report #1

### 13. FORTRAN LISTINGS

The following is a listing of the FORTRAN codes described in this report. There are three slightly different versions of the codes because three different FORTRAN compilers are of interest (f77 under Unix for use with the hypercube simulator, ftn286, and rmfort). The only distinctions between the versions are:

1. Whether integers are two bytes or four.  
f77 4 bytes  
ftn286 2 bytes  
rmfort 4 bytes
2. How byte addressable buffers are declared.  
f77 character\*1  
ftn286 logical\*1  
rmfort logical\*1
3. How exclusive-or computations are implemented.  
see the listing for XOR

The actual listings given here are the rmfort version. All three versions are distributed with the operating system (as well as the C version), both as source and as libraries.

```
LOGICAL FUNCTION BETWEEN(I, J, K)
INTEGER I,J,K
```

```
C
C BETWEEN decides whether node J lies on a shortest path from node I to
C node K. This function is used by SPGSNW and SPGRCW (sparse global send
C and receive).
```

```
C
C Calls: XOR
C
```

```
INTEGER XOR
BETWEEN = XOR(I,K) .EQ. (XOR(I,J) + XOR(J,K))
RETURN
END
```

```

SUBROUTINE GCAT (CI, TYPE, BUF, BUFLen, MSGLEN, CNT, ROOT, DIM)
INTEGER CI, TYPE, BUFLen, MSGLEN, CNT, ROOT, DIM
C
C byte indexed buffers are of type:
C
C   for ftn286 and RM FORTRAN:
C   LOGICAL*1 BUF (BUFLen)
C
C   for the simulator under f77:
C   CHARACTER*1 BUF (BUFLen)
C
C Global concatenation operation using spanning tree.  If the root
C is zero then the messages are concatenated in numerical order.
C
C All participating processes must have the same process id (PID).
C
C Input..
C
C   CI          channel number (previously opened).
C   TYPE        type number.  must be the same for all participating
C               processes.  There must be no other messages of this type
C               OR OF TYPES TYPE+1, TYPE+2, ..., TYPE+DIM-1
C               in the system.
C   BUF         the input vector to be concatenated (which may actually
C               be any type).
C   BUFLen     the length of BUF (in BYTES).
C   MSGLEN     the length of this contribution (in BYTES).
C   ROOT       Node id of root process (which will get the final message).
C               (if -32768, then the smallest node number in the active
C               subcube acts as root and then forwards the message to the host).
C   DIM        the size of the subcube participating.
C
C Output..
C
C   BUF        for the root process, BUF contains the desired result.
C               for all other processes, BUF contains the partial result
C               for their subtrees.
C   CNT        the number of meaningful bytes in BUF.
C
C Errors Conditions
C
C   If called by a nonparticipating node, an error message is put on
C   the syslog and then the subroutine exits.
C
C   If a result longer than BUFLen BYTES is attempted an error
C   message is syslogged and the subroutine terminates.
C
C Calls: MYNODE, MYPID, RECVW, SENDW, SYSLOG, XOR
C
C   INTEGER BIT, CURLen, DIFF, I, IGNORE, ME, MYNODE,
C   *       MYPID, P, PARENT, PID, TROOT, XOR
C
C   ME = MYNODE ()
C   P = 2**DIM

```

```

C
C Find temporary root (either the real root, or the lowest
C numbered node in the active subcube--found by zeroing the
C DIM lowest bits in mynode).
C
C TROOT = MAX0 ((ME/P) *P, ROOT)
C
C CURLEN = MSGLEN
C PID = MYPID ()
C DIFF = XOR (ME, TROOT)
C IF (DIFF .GE. P) THEN
C     CALL SYSLOG (MYPID (), 'GCAT: CALLED BY NON PARTICIPANT')
C     RETURN
C ENDIF
C
C Accumulate contributions from children, if any.
C
C BIT = 1
C
C DO 10 I = 1, DIM
C     IF (MOD (DIFF, 2*BIT) .NE. 0) GO TO 20
C     CALL RECVW (CI, TYPE+I-1, BUF (CURLen+1), BUFLen-CURLen,
C *         CNT, IGNORE, IGNORE)
C
C     IF (CNT .GT. BUFLen-CURLen) THEN
C         CALL SYSLOG (PID, 'GCAT: BUFFER LENGTH EXCEEDED')
C         RETURN
C     ENDIF
C
C     CURLen = CURLen + CNT
C     BIT = 2*BIT
C 10 CONTINUE
C
C If exit normally, then this is the root.
C
C CNT = CURLen
C IF (ROOT .LT. 0) CALL SENDW (CI, TYPE, BUF, CNT, -32768, PID)
C RETURN
C
C Pass result back to parent.
C
C 20 CONTINUE
C PARENT = XOR (ME, BIT)
C CALL SENDW (CI, TYPE+I-1, BUF, CURLen, PARENT, PID)
C CNT = CURLen
C RETURN
C END

```

```
INTEGER FUNCTION GINV(N)
INTEGER N
C
C   GINV = inverse of GRAY function
C
C   Answers the inevitable question, what task is being
C   carried out on node n?
C
   INTEGER T, XOR
   T = N
   GINV = N
1  T = T/2
   IF (T .EQ. 0) RETURN
   GINV = XOR(GINV, T)
   GO TO 1
END
```

```

SUBROUTINE GOP (CI, TYPE, X, N, OP, ROOT, DIM, WORK)
INTEGER CI, TYPE, N, ROOT, DIM
CHARACTER*1 OP
DOUBLE PRECISION X(N), WORK(N)

C
C Global vector commutative operation using spanning tree.
C
C All participating processes must have the same process id (PID).
C
C Input..
C
C   CI      channel number (previously opened).
C   TYPE    message type. Must be the same for all participating
C           processes. There must be no other messages of this type
C           in the system.
C   X       the input vector to be used in the operation.
C   N       the length of the vector.
C   OP      '+'  sum
C           '*'  product
C           'M'  maximum
C           'm'  minimum
C   ROOT    Node id of root process (which will get the final message).
C           (if -32768, then the smallest node number in the active
C           subcube acts as root and then forwards the message to the host)
C   DIM     the size of the subcube participating.
C
C Output..
C
C   X       for the root process, X contains the desired result.
C           for all other processes, X contains the partial result
C           for their subtrees.
C
C Workspace
C
C   WORK    used to receive other contributions.
C
C Errors Conditions
C
C   If called by a nonparticipating node, an error message is
C   syslogged and then the subroutine exits.
C
C   If a message longer than N elements is received, only the first N
C   elements are saved, an error message is syslogged,
C   and then the computation continues with the truncated results.
C
C   If a message shorter than N elements is received, then an error
C   message is syslogged and the computation continues.
C
C Calls: MYNODE, MYPID, RECVW, SENDW, SYSLOG, XOR
C
C   INTEGER BIT, BYTES, CNT, DIFF, DPSIZE, I, IGNORE, ME, MYNODE,
C   *   MYPID, P, PARENT, PID, TROOT, XOR
C   PARAMETER (DPSIZE = 8)

```

```

C
ME = MYNODE()
P = 2**DIM

C
C Find temporary root (either the real root, or the lowest
C numbered node in the active subcube--found by zeroing the
C DIM lowest bits in mynode).
C
TROOT = MAX0((ME/P)*P, ROOT)

C
PID = MYPID()
DIFF = XOR(ME, TROOT)
IF (DIFF .GE. P) THEN
    CALL SYSLOG(MYPID(), 'GOP: CALLED BY NON PARTICIPANT')
    RETURN
ENDIF

C
C Accumulate contributions from children, if any
C
BIT = P/2
BYTES = DPSIZE*N
5 IF (BIT .LE. DIFF) GO TO 20
    CALL RECVW(CI, TYPE, WORK, BYTES, CNT, IGNORE, PID)
    IF (CNT .GT. BYTES) CALL SYSLOG(PID, 'GOP: LONG MESSAGE')
    IF (CNT .LT. BYTES) CALL SYSLOG(PID, 'GOP: SHORT MESSAGE')
    DO 10 I = 1, N
        IF (OP .EQ. '+') X(I) = X(I) + WORK(I)
        IF (OP .EQ. '*') X(I) = X(I) * WORK(I)
        IF (OP .EQ. 'M') X(I) = DMAX1(X(I), WORK(I))
        IF (OP .EQ. 'm') X(I) = DMIN1(X(I), WORK(I))
10    CONTINUE
        BIT = BIT/2
    GO TO 5

C
C Pass result back to parent
C
20 CONTINUE
IF (BIT .NE. 0) THEN
    PARENT = XOR(ME, BIT)
    CALL SENDW(CI, TYPE, X, BYTES, PARENT, PID)
ELSE
    IF (ROOT .LT. 0) CALL SENDW(CI, TYPE, X, BYTES, -32768, PID)
ENDIF
RETURN
END

```

```
INTEGER FUNCTION GRAY(J)
INTEGER J, XOR
C
C GRAY(J) = J-th entry in gray code
C domain and range both 0..127
C
C This is the standard way to obtain a one-dimensional mesh
C or ring of tasks on the hypercube. If the j-th task is
C carried out on node number GRAY(J), then each task is
C connected directly to its predecessor and its successor
C in the hypercube topology.
C
GRAY = XOR(J, J/2)
RETURN
END
```

```

SUBROUTINE GRECVW(CI, TYPE, BUF, LEN, CNT, DIM)
INTEGER CI, TYPE, BUF(*), LEN, CNT, DIM
C
C Global send participant.  Receives message from unknown source and
C   sends it on to some neighbors.
C
C All participating processes must have the same process id (PID).
C
C Input..
C
C   CI          channel number (previously opened).
C   TYPE        message type.  Must be the same for all participating
C               processes.  There must be no other messages of this type
C               in the system.
C   LEN         the length of BUF in BYTES.
C   DIM         the dimension of the subcube participating in the send.
C
C Output..
C
C   BUF         the message (which may actually be any type).
C   CNT         the length (in BYTES) of the message received.
C
C Error Conditions
C
C   If a message longer than LEN bytes is received then only
C   LEN bytes will be stored in BUF and the rest of the
C   message will be lost.  In this case an error message
C   will be sent to syslog but the remnants of the message
C   will be sent on.
C
C   NOTE: only those nodes which will participate in the
C   send can call GRECVW.  Any other node which calls it
C   will never return.
C
C Calls: MYNODE, MYPID, RECVW, SENDW, SYSLOG, XOR
C
C   INTEGER BIT, I, LENOUT, ME, MYNODE, MYPID, NODE, P, PID,
C   *   PIDIN, XOR
C
C   P = 2**DIM
C   ME = MYNODE()
C   PID = MYPID()
C   CALL RECVW(CI, TYPE, BUF, LEN, CNT, NODE, PIDIN)
C
C   LENOUT = CNT
C   IF (CNT .GT. LEN) THEN
C     CALL SYSLOG(PID, 'GRECVW: MESSAGE TRUNCATED')
C     LENOUT = LEN
C   ENDIF
C
C   BIT = 2*XOR (ME, NODE)
C
C Check to see if received from host.
C

```

```
IF (NODE .LT. 0) BIT = 1
C
DO 10 I = 1, DIM
  IF (BIT .EQ. P) RETURN
  NODE = XOR (ME, BIT)
  CALL SENDW(CI, TYPE, BUF, LENOUT, NODE, PID)
  BIT = 2*BIT
10 CONTINUE
END
```

```

SUBROUTINE GSENDW(CI, TYPE, BUF, LEN, DIM)
INTEGER CI, TYPE, BUF(*), LEN, DIM
C
C Global send of data. Other participants call grecvw.
C
C All participating processes must have the same process id (PID).
C
C Input
C
C CI          channel number (previously opened).
C TYPE       message type. Must be the same for all participating
C            processes. There must be no other messages of this type
C            in the system.
C BUF        the message buffer (which may actually be any type).
C LEN        the length of the buffer in BYTES
C DIM        the dimension of the subcube
C
C Calls:     MYNODE, MYPID, SENDW, XOR
C
C           INTEGER BIT, I, ME, NODE, MYNODE, MYPID, PID, XOR
C
C           ME = MYNODE()
C           PID = MYPID()
C           BIT = 1
C
C           DO 10 I = 1, DIM
C             NODE = XOR(ME, BIT)
C             CALL SENDW(CI, TYPE, BUF, LEN, NODE, PID)
C             BIT = 2*BIT
10 CONTINUE
RETURN
END

```

```

SUBROUTINE LRECVW(CI, TYPE, BUF, LEN4, CNT4, NODE, PID)
INTEGER CI, TYPE, NODE, PID
INTEGER*4 LEN4, CNT4

C
C byte indexed buffers are of type:
C
C     for ftn286 and RM FORTRAN:
C     LOGICAL*1 BUF(*)
C
C     for the simulator under f77:
C     CHARACTER*1 BUF(*)
C
C send/rcv cannot process messages larger than 16K bytes. LSENDW and
C LRECVW can be used to send large messages between nodes. LSENDMSG and
C LRECVW can be used to send large messages from the host to a node.
C There is no subroutine to receive large messages on the host since
C the semantics of type is different for RECVMSG. The message
C is sent in 16K byte chunks and terminated with a smaller chunk. If
C the message is an even multiple of 16K bytes, this requires sending
C a zero byte message as a terminator and requires that LEN4 >= CNT4 + 1
C
C Input..
C
C   CI          channel number (previously opened).
C   TYPE        message type. There must be no other messages of this type
C               received by the target node.
C   BUF         the message buffer (which may actually be any type).
C   LEN4        an INTEGER*4, the length of BUF in BYTES.
C
C Output..
C
C   CNT4        an INTEGER*4, the length of the received message in BYTES.
C   NODE        the source node.
C   PID         the source pid.
C
C Calls: RECVW, SYSLOG
C
C
C     INTEGER*4 REMAIN, INDEX
C     INTEGER LEN2, MYPID, CNT2
C     PARAMETER (MAXLEN = 16384)
C
C     REMAIN = LEN4
C     INDEX = 1
C     CNT4 = 0
C     LEN2 = MAXLEN
C

```

```
10 CONTINUE
   IF (LEN4 - INDEX + 1 .LT. MAXLEN) LEN2 = LEN4 - INDEX + 1
   CALL RECVW(CI, TYPE, BUF(INDEX), MAXLEN, CNT2, NODE, PID)
   IF (CNT2 .GT. LEN2)
*    CALL SYSLOG(MYPID(), 'LRECVW: TRUNCATED MESSAGE')
   INDEX = INDEX + CNT2
   CNT4 = CNT4 + CNT2
   IF (CNT2 .EQ. MAXLEN) GO TO 10
c
RETURN
END
```

```

SUBROUTINE LSENDMSG(CI, TYPE, BUF, LEN4, NODE, PID)
INTEGER CI, TYPE, NODE, PID
INTEGER*4 LEN4
C
C byte indexed buffers are of type:
C
C     for ftn286 and RM FORTRAN:
LOGICAL*1 BUF (*)
C
C     for the simulator under f77:
CHARACTER*1 BUF (*)
C
C see LRECVW for details
C
C Input..
C
C   CI       a channel number (previously opened).
C   TYPE     message type. Must be the same as the corresponding
C           LRECVW call.
C   BUF     the message buffer (which may actually be any type).
C   LEN4    an INTEGER*4 which is the length of the message in BYTES.
C   NODE    the destination node.
C   PID     the destination pid.
C
C Calls SENDMSG
C
C   INTEGER*4 REMAIN, INDEX
C   INTEGER I, LEN2, MAXLEN, NUMMSG
C   PARAMETER (MAXLEN = 16384)
C
C   REMAIN = LEN4
C   INDEX = 1
C   LEN2 = MAXLEN
C   NUMMSG = LEN4/MAXLEN + 1
C
C   DO 10 I = 1, NUMMSG
C
C special settings for last message
C
C   IF (REMAIN .LT. MAXLEN) LEN2 = REMAIN
C   IF (REMAIN .EQ. 0) INDEX = 1
C   CALL SENDMSG(CI, TYPE, BUF (INDEX), LEN2, NODE, PID)
C   REMAIN = REMAIN - LEN2
C   INDEX = INDEX + LEN2
10 CONTINUE
C
C   RETURN
C   END

```

```

SUBROUTINE LSENDW(CI, TYPE, BUF, LEN4, NODE, PID)
INTEGER CI, TYPE, NODE, PID
INTEGER*4 LEN4
C
C byte indexed buffers are of type:
C
C     logical*1 for ftn286 and RM FORTRAN
LOGICAL*1 BUF(*)
C
C     character*1 for f77 and the simulator
CHARACTER*1 BUF(*)
C
C see LRECVW for details
C
C Input..
C
C CI          a channel number (previously opened).
C TYPE       message type. Must be the same as the corresponding
C           LRECVW call.
C BUF        the message buffer (which may actually be any type).
C LEN4       an INTEGER*4 which is the length of the message in BYTES.
C NODE       the destination node.
C PID        the destination pid.
C
C Calls:  sendw
C
C     INTEGER*4 REMAIN, INDEX
C     INTEGER I, LEN2, NUMMSG
C     PARAMETER (MAXLEN = 16384)
C
C     REMAIN = LEN4
C     INDEX = 1
C     LEN2 = MAXLEN
C     NUMMSG = LEN4/MAXLEN + 1
C
C     DO 10 I = 1, NUMMSG
C
C special parameters for last message
C
C     IF (REMAIN .LT. MAXLEN) LEN2 = REMAIN
C     IF (REMAIN .EQ. 0) INDEX = 1
C
C     CALL SENDW(CI, TYPE, BUF (INDEX), LEN2, NODE, PID)
C     REMAIN = REMAIN - LEN2
C     INDEX = INDEX + LEN2
10 CONTINUE
C
C     RETURN
C     END

```

```

SUBROUTINE MESH(ME, DIM, LINDIM, PERBC, COORDS, PRED, SUCC,
*   WORK)
INTEGER ME, DIM, LINDIM(DIM), PERBC(DIM), COORDS(DIM),
*   PRED(DIM), SUCC(DIM), WORK(DIM)
C
C MESH Embeds a rectangular mesh into the hypercube.
C The location of the current node in the mesh is computed
C as well as the neighboring node numbers.
C If mesh is called with too large a mesh, some of the
C neighboring numbers may not exist in the cube.
C
C Input..
C
C ME the node number of interest.
C DIM the physical dimension of the mesh (usually 1, 2, or 3).
C LINDIM the linear dimension (number of processors) in each
C direction of the mesh.
C PERBC a flag indicating whether a particular direction has
C periodic boundary conditions. perbc should be set as
C follows:
C
C PERBC(I) = 0 for non periodic (terminated) boundaries
C PERBC(I) = 1 for periodic boundaries
C
C Output..
C
C COORDS Gives the coordinates of me in the mesh,
C starting with 0 so that the largest the ith coordinate
C can be is LINDIM(I)-1. If me is not part
C of the mesh all coordinates are returned as -1.
C PRED PRED(I) is the node number preceding me in the ith
C direction. -1 is used when there is no predecessor.
C SUCC SUCC(I) is the node number succeeding me in the ith
C direction. -1 is used when there is no successor.
C
C Workspace
C
C WORK holds the result of rounding the elements of LINDIM
C up to the next power of two.
C
C Calls: CUBEDIM, GRAY, GINV, MYPID, SYSLOG
C
C INTEGER CUBEDIM, D, GRAY, GINV, I, J, MYPID, SIZE, SLICE, TEMP
C
C Round up LINDIM to next power of 2.
C
C SIZE = 1
C DO 1 I = 1, DIM
C D = 0
C IF (LINDIM(I) .GT. 1) D = ALOG (FLOAT (LINDIM(I) - 1)) /ALOG (2.0) +1
C WORK (I) = 2**D
C SIZE = SIZE*WORK (I)
1 CONTINUE

```



```

SUBROUTINE SPGRCW(CI, TYPE, BUF, LEN, CNT)
INTEGER CI, TYPE, BUF(*), LEN, CNT
C
C This subroutine participates in a sparse global send. It receives
C a message with a list of node numbers appended which must also
C receive the message. It calls SPGSNW to forward the message.
C The process is started by a call to SPGSNW.
C
C All participating processes must have the same process id (PID).
C
C Input..
C
C   CI          channel number (previously opened).
C   TYPE        message type. Must be the same for all participating
C               processes. There must be no other messages of this type
C               in the system.
C   LEN         the length of BUF in BYTES. LEN must be larger than the
C               expected message. In particular, the buffer must be able
C               to hold 2**CUBEDIM() additional integers.
C
C Output..
C
C   BUF         contains the message.
C   CNT         the length of the received message in BYTES.
C
C Error Conditions
C
C   If the received message is truncated then an error message is
C   syslogged and the subroutine returns.
C
C Calls: MYPID, RECVW, SPGSNW, SYSLOG
C
C   INTEGER MYPID, NODE, NUMBER, PID, PIDIN
C   PARAMETER (INTSIZE = 4)
C
C   PID = MYPID()
C   CALL RECVW(CI, TYPE, BUF, LEN, CNT, NODE, PIDIN)
C   IF (CNT .GT. LEN) THEN
C       CALL SYSLOG(PID, 'TRUNCATED MESSAGE IN SPGRCW')
C       RETURN
C   ENDIF
C
C   NUMBER = BUF (CNT/INTSIZE)
C   CNT = CNT - INTSIZE*(NUMBER + 1)
C   CALL SPGSNW(CI, TYPE, BUF, CNT, BUF (CNT/INTSIZE + 1), NUMBER)
C   RETURN
C
C   END

```

```

SUBROUTINE SPGSNW(CI, TYPE, BUF, LEN, LIST, NUMBER)
INTEGER CI, TYPE, BUF(*), LEN, LIST(*), NUMBER
C
C This subroutine sends a buffer to a subset of nodes specified by
C an index list. All the nodes which are to receive the buffer must
C call SPGRCW.
C
C All participating processes must have the same process id (PID).
C
C Input..
C
C   CI          channel number (previously opened).
C   TYPE        message type. Must be the same for all participating
C               processes. There must be no other messages of this type
C               in the system.
C   BUF         the message buffer (which may actually be any type).
C   LEN         the number of BYTES of data in BUF. SPGRCW can only
C               receive messages that are a multiple of INTSIZE.
C               Therefore LEN must be a multiple of INTSIZE.
C               Furthermore the LIST is appended to the message and so
C               BUF must have length at least LEN + INTSIZE*NUMBER.
C   LIST        the list of node numbers which should receive the data.
C   NUMBER      the number of node numbers in LIST.
C
C Output..
C
C   the list of node numbers is destroyed
C
C Calls: MYPID, SENDW, SPLIT
C
C   PARAMETER (INTSIZE = 4)
C   INTEGER BUFLen, FIRST, I, J, MYPID, NODE, PID
C
C   PID = MYPID()
C   BUFLen = LEN/INTSIZE
10 IF (NUMBER .EQ. 0) RETURN
   CALL SPLIT(NUMBER, LIST, FIRST)
   NODE = LIST(FIRST+1)
   DO 20 J = 1, FIRST
     BUF(BUFLen + J) = LIST(J)
20   CONTINUE
C
   BUF(BUFLen + FIRST + 1) = FIRST
   CALL SENDW(CI, TYPE, BUF, LEN + INTSIZE*(FIRST+1), NODE, PID)
   J = 1
   DO 30 I = FIRST+2, NUMBER
     LIST(J) = LIST(I)
     J = J+1
30   CONTINUE
C
   NUMBER = NUMBER - FIRST - 1
GO TO 10
C
END

```

```

SUBROUTINE SPLIT(NUMBER, LIST, FIRST)
  INTEGER NUMBER, LIST(*), FIRST
C
C SPLIT takes a set of node numbers, finds a child node (one with
C no other node between me and it), and then separates the remaining
C nodes into the descendents of the child (all nodes for which the
C child is between me and the node) and the rest. The descendents
C are put first, the child in the middle and the rest afterward.
C
C Input..
C
C   NUMBER    the number of node numbers.
C   LIST      the list of node numbers.
C
C Output..
C
C   LIST      sorted so that the descendents of the child are first,
C              the child in the middle, and the rest last.
C   FIRST     returned as the number of descendents of the child.
C
C Calls: MYNODE
C
C   INTEGER CHILD, HI, KEY, LO, ME, MYNODE, TEMP, TEST
C   LOGICAL BETWEEN
C   ME = MYNODE()
C
C find key
C
C   KEY = 1
C   TEST = 2
C   10 IF (TEST .GT. NUMBER) GO TO 20
C   IF (BETWEEN(ME, LIST(TEST), LIST(KEY))) KEY = TEST
C   TEST = TEST+1
C   GO TO 10
C
C swap key to front
C
C   20 CHILD = LIST(KEY)
C   LIST(KEY) = LIST(1)
C
C split rest of list. put descendents of key first.
C
C   LO = 2
C   HI = NUMBER
C
C move lo upward
C
C   30 IF (LO .GT. HI) GO TO 60
C   IF (.NOT. BETWEEN(ME, CHILD, LIST(LO))) GO TO 40
C   LO = LO + 1
C   GO TO 30
C
C move hi downwards
C

```

```
40 IF (LO .GE. HI) GO TO 60
    IF (BETWEEN(ME, CHILD, LIST(HI))) GO TO 50
    HI = HI - 1
    GO TO 40
c
c swap
c
50 TEMP = LIST(LO)
    LIST(LO) = LIST(HI)
    LIST(HI) = TEMP
    LO = LO + 1
    HI = HI - 1
    GO TO 30
c
c put child in the middle
c
60 FIRST = LO - 2
    LIST(1) = LIST(FIRST+1)
    LIST(FIRST+1) = CHILD
    RETURN
    END
```

```
      INTEGER FUNCTION XOR (M,N)
C
C  exclusive or
C
C  Pick one of the following:
C
C  UNIX 4.2, f77:
C    XOR = OR (M,N) -AND (M,N)
C
C  Intel FTN286:
C    XOR = M.NEQV.N
C
C  Ryan-McFarland Fortran
C    XOR = IEOR (M,N)
C
      RETURN
      END
```

Printed in USA/SC6054/0886/1k/WCP/VJC  
Order Number 280323-001



**INTEL SCIENTIFIC COMPUTERS**

15201 N.W. Greenbrier Parkway  
Beaverton, Oregon 97006  
(503) 629-7629