

## 2.2. Data Distributions and Software Conventions

The ScaLAPACK software library provides routines that operate on three types of matrices: **in-core dense matrices**, **in-core narrow band matrices** and **out-of-core dense matrices**. On entry, these routines assume that the *data has been distributed on the processors according to a specific data decomposition scheme*. Conventional arrays are used to store locally the data when it resides in the processors memory. The data layout information as well as the local storage scheme for these different matrix operands is conveyed to the routines via a simple array of integers called an **array descriptor**. The first entry of this array identifies the type of the descriptor, i.e., the data distribution scheme it describes. This paragraph first presents the fundamental concepts of **process grids**, **communication contexts** and **array descriptors**. Then, for each of the three possible matrix operand mentioned above, the data distribution scheme and the corresponding descriptor array used by ScaLAPACK are discussed in detail. Finally, the software conventions common to all ScaLAPACK routines are presented.

So to realize data transmission (here Matrix elements) between processes it is necessary:

- to generate a grid of processes;
- to describe a communication medium, i.e. a communication context;
- to determine (initialize) a descriptor for each array.

### 2.2.1 Process Grid, Scoped Operations and Context

The  $P$  processes of an abstract parallel computer are often represented as a one-dimensional linear array of processes labeled  $0, 1, \dots, P$ . For reasons described below, it is often more convenient to map this one-dimensional array of processes into a two-dimensional rectangular grid, or **process grid** by using **row-major order** (the numbering of the processes increases sequentially across each row) or by using **column-major order** (the numbering of the processes proceeds down each column of the process grid).

This grid will have  $P_r$  process rows and  $P_c$  process columns, where  $P_r \times P_c = P$ . A process can now be referenced by its row and column coordinates,  $(p_r, p_c)$ , within the grid, where  $0 \leq p_r < P_r$ , and  $0 \leq p_c < P_c$ . An example of such a mapping is shown in the following figure, where  $P_r = 2$  and  $P_c = 4$ :

	0	1	2	3
0	0	1	2	3
1	4	5	6	7

An operation which involves more than just a sender and a receiver is called a scoped operation. All processes that participate in a scoped operation are said to be within the operation's scope. On a system using a linear array of processes, the only natural scope is all processes. Using a 2D grid, we have 3 natural scopes, as shown in the following table:

<i>Scope</i>	<i>Meaning</i>
Row	All processes in a process row participate.
Column	All processes in a process column participate.
All	All processes in the process grid participate.

These groupings of processes are of particular interest to the linear algebra programmer, since distributed data decompositions of a linear algebra matrix tend to follow this process mapping. Viewing the rows/columns of the process grid as essentially autonomous subsystems provides the programmer with additional levels of parallelism.

The BLACS routine `BLACS_GRIDINIT` performs the task of mapping the processes to the process grid. By default, `BLACS_GRIDINIT` assumes a row-major ordering, although a column-major ordering can also be specified. The companion routine `BLACS_GRIDMAP` is a more general form of `BLACS_GRIDINIT` and allows the user to define the mapping of the processes.

In ScaLAPACK, and thus the BLACS, each process grid is enclosed in a context. Similarly, a context is associated with every global matrix in ScaLAPACK. The use of a context provides the ability to have separate "universes" of message passing. This means that a process grid can safely communicate even if other process grids are also communicating. Thus, a context is a powerful mechanism for avoiding unintentional nondeterminism in message passing and provides support for the design of safe, modular software libraries. In MPI, this concept is referred to as a communicator.

A context partitions the communication space. A message sent from one context cannot be received in another context. The use of separate communication contexts by distinct libraries (or distinct library routine invocations) insulates communication internal to a specific library routine from external communication that may be going on within the user's program. In most respects, we can use the terms process grid and context interchangeably. So context allows us to do the following:

- create arbitrary groups of processes,
- create an indeterminate number of overlapping and/or disjoint process grids,
- isolate the process grid so that they do not interfere with each other.

In the BLACS, there are two grid creation routines (`BLACS_GRIDINIT` and `BLACS_GRIDMAP`) which create a process grid and its enclosing context. These routines return context handles, which are simple integers, assigned by the BLACS to identify the context. Subsequent BLACS routines will be passed these handles, which allow the BLACS to determine from which context/grid a routine is being called. The user should never actually manipulate these handles; they are opaque data objects which are only meaningful for the BLACS routines.

A defined context consumes resources. It is therefore advisable to release contexts when they are no longer needed. This is done via the routine `BLACS_GRIDEXIT`. When the entire BLACS system is shut down (via a call to `BLACS_EXIT`), all outstanding contexts are automatically freed.

We describe the BLACS routines that deal with grid/context creation, and processing before the grid/context has been defined

```
BLACS_PINFO (MYPNUM,NPROCS)
void Cblacs_pinfo(int*,int*);
```

MYPNUM (output) INTEGER

*An integer between 0 and (NPROCS-1) which uniquely identifies each process.*

NPROCS (output) INTEGER

*The number of processes available for BLACS use.*

These routine is used when some initial system information is required before the BLACS are set up. Primary purpose of calling this routine is to get the number of processors that the program will use (specified on parallel run command line) and the ID number of the processor that called the routine.

```
BLACS_GET(-1,0,ICONTEXT)
void Cblacs_get(int,int,int*);
```

ICONTEXT (output) INTEGER

BLACS\_GET is a general-purpose values for a set of BLACS internal parameters. With the arguments shown above, the BLACS\_GET call will get the context (argument ICONTEXT) for ScaLAPACK use. (Note: the first argument is ignored in this form of the call.)

```
BLACS_GRIDINIT(ICONTEXT,ORDER,NPROW,NPCOL)
void Cblacs_gridinit(int*,const char*,int,int);
```

ICONTEXT (input/output) INTEGER

*On input, an integer handle indicating the system context to be used in creating the BLACS context. The user may obtain a default system context via a call to BLACS\_GET. On output, the integer handle to the created BLACS context.*

ORDER (input) CHARACTER\*1

*Indicates how to map processes to BLACS grid. Choices are: 'R' - use row-major natural ordering; 'C' - use column-major natural ordering; ELSE - use row-major natural ordering.*

NPROW (input) INTEGER

*Indicates how many process rows the process grid should contain.*

NPCOL (input) INTEGER

*Indicates how many process columns the process grid should contain.*

All codes that use the ScaLAPACK (BLACS) must call this routine, or its companion routine BLACS\_GRIDMAP. These routines take the available processes, and assign, or map, them into a BLACS process grid. In other words, they establish how the BLACS coordinate system will map into the native machine's process numbering system. Each BLACS grid is contained in a context (its own message passing universe), so that it does not interfere with distributed operations which occur within other grids/contexts. These grid creation routines may be called repeatedly in order to define additional contexts/grids.

The creation of a grid requires input from all processes which are defined to be in it. It is therefore a globally-blocking (sometimes called synchronous) operation which means that processes belonging to more than one grid will have to agree on which grid formation will be serviced first. Note that these routines map already-existing processes to a grid: the processes are not created dynamically. On most parallel machines, the processes will be actual processors (hardware), and they are "created" when the user runs his executable. This routine creates a simple

NPROW×NPCOL process grid. This process grid will use the first NPROW×NPCOL processes, and assign them to the grid in a row or column-major natural ordering. If these process-to-grid mappings are unacceptable, more complex companion routine BLACS\_GRIDMAP must be called instead.

```
BLACS_GRIDMAP(ICONTEXT,USERMAP,LDUMAP,NPROW,NPCOL)
void Cblacs_gridmap(int*, int*, int, int, int);
```

ICONTEXT (input/output) INTEGER

*On input, an integer handle indicating the system context to be used in creating the BLACS context. The user may obtain a default system context via a call to BLACS\_GET. On output, the integer handle to the created BLACS context.*

USERMAP (input) INTEGER array, dimension (LDUMAP,NPCOL)

*Input array indicating the process-to-grid mapping.*

LDUMAP (input) INTEGER

*The leading dimension of the 2D array USERMAP.*

NPROW (input) INTEGER

*Indicates how many process rows the process grid should contain.*

NPCOL (input) INTEGER

*Indicates how many process columns the process grid should contain.*

These routines similar the BLACS\_GRIDINIT take the available processes, and assign, or map, them into a BLACS process grid. This routine allows the user to map processes to the process grid in an arbitrary manner. USERMAP(i,j) holds the process number of the process to be placed in {i,j} of the process grid. On most distributed systems, this process number will simply by a machine defined number between 0,...,NPROCS-1. BLACS\_GRIDMAP is not for the inexperienced user – BLACS\_GRIDINIT is much simpler. BLACS\_GRIDINIT simply performs a GRIDMAP where the first NPROW×NPCOL processes are mapped into the current grid in a row- or column-major natural ordering. BLACS\_GRIDMAP allows the experienced user to take advantage of the processors' actual network (i.e. he can map nodes that are physically connected to be neighbors in the BLACS grid, etc.). BLACS\_GRIDMAP also opens the way for multigridding: the user can separate his nodes into arbitrary grids, join them together at some later date, and then re-split them into new grids. BLACS\_GRIDMAP also provides the ability to make arbitrary grids or subgrids (e.g., a "nearest neighbor" grid), which can greatly facilitate operations among groups of processes which do not fall on a row or column of the main process grid.

The following routines destroy grids and free resources.

```
BLACS_FREEBUFF(ICONTEXT,WAIT)
void Cblacs_freebuf(int,int);
```

ICONTEXT (input) INTEGER

*Integer handle indicating the BLACS context.*

WAIT (input) INTEGER

*Whether to wait on non-blocking operations:*

*IF (WAIT .EQ. 0) THEN*

*Do not wait on operations, free only unused buffers.*

*ELSE*

*If necessary, wait in order to free all buffers.*

*END IF*

The BLACS have at least one internal buffer that is used for packing messages

(the number of internal buffers varies, depending on which BLACS you are using). On systems where memory is tight, keeping this buffer(s) around may become expensive. Calling this routine will release the BLACS buffer(s). However, the next call to a communication routine which requires packing will cause the buffer to be reallocated. The parameter WAIT determines whether the BLACS should wait for any non-blocking operations to complete or not. If WAIT=0, the BLACS will free any buffers that can be freed without waiting. If WAIT is not 0, the BLACS will free all internal buffers, possibly causing the call to block while the BLACS wait for internal non-blocking operations complete.

```
BLACS_GRIDEXIT(ICONTXT)
void Cblacs_gridexit(int);
    ICONTXT (input) INTEGER
        Integer handle indicating the BLACS context to be freed.
```

Contexts consume resources, and therefore the user should release them when they are no longer needed. BLACS\_GRIDEXIT frees a context. After the freeing of a context, the context no longer exists, and its handle may be re-issued by the BLACS if a new context is defined.

```
BLACS_EXIT(CONTINUE)
void Cblacs_exit(int);
    CONTINUE (input) INTEGER.
        If CONTINUE is non-zero, it is assume the user will continue using the machine after the BLACS are done. Otherwise, it is assumed that no message passing will be done after the BLACS_EXIT call.
```

This routine should be called when a process has finished all use of the BLACS. It frees all BLACS contexts and releases all memory the BLACS have allocated. CONTINUE indicates whether the user will be using the underlying communication platform after the BLACS are finished.

The following routines return information involving the process grid. Also included here is the barrier routine.

```
BLACS_GRIDINFO(ICONTXT,NPROW,NPCOL,MYPROW,MYPCOL)
void Cblacs_gridinfo(int,int*,int*,int*,int*);
    ICONTXT (input) INTEGER
        Integer handle indicating the BLACS context to be queried.
    NPROW (output) INTEGER
        On output, the number of process rows in ICONTXT's process grid.
    NPCOL (output) INTEGER
        On output, the number of process columns in ICONTXT's process grid.
    MYPROW (output) INTEGER
        On output, the calling process's row coordinate in the process grid.
    MYPCOL (output) INTEGER
        On output, the calling process's column coordinate in the process grid.
```

Returns information about the process grid contained in the context whose handle is ICONTXT. If the context handle is invalid, all quantities are returned as -1.

```
INTEGER FUNCTION BLACS_PNUM(ICONTXT,PROW,PCOL)
```

```
int Cblacs_pnum(int, int, int);
```

```
ICONTEXT (input) INTEGER
```

*Handle indicating the BLACS context to be queried.*

```
PROW (input) INTEGER
```

*The row coordinate of the process whose system process number is to be determined.*

```
PCOL (input) INTEGER
```

*The column coordinate of the process whose system process number is to be determined.*

This function returns the system process number of the process at PROW, PCOL in the process grid.

```
BLACS_PCOORD(ICONTEXT, PNUM, PROW, PCOL)
```

```
void Cblacs_pcoord(int, int, int*, int*);
```

```
ICONTEXT (input) INTEGER
```

*Integer handle indicating the BLACS context.*

```
PNUM (input) INTEGER
```

*The process number whose coordinates is to be determined. This is the process number of the underlying machine.*

```
PROW (output) INTEGER
```

*On output, the row coordinate of process PNUM in the BLACS grid.*

```
PCOL (output) INTEGER
```

*On output, the column coordinate of process PNUM in the BLACS grid.*

Given the system process number, returns the row and column coordinates in the BLACS' process grid.

A call to the ScaLAPACK TOOLS routine SL\_INIT initializes the process grid. This routine initializes a PrxPc (denoted NPROW×NPCOL in the source code) process grid by using a row-major ordering of the processes, and obtains a default system context. The user can then query the process grid to identify each process's coordinates (MYROW,MYCOL) ) via a call to BLACS\_GRIDINFO.

A typical code fragment to accomplish this task would be the following:

```
CALL SL_INIT(ICONTEXT, NPROW, NPCOL)
```

```
CALL BLACS_GRIDINFO(ICONTEXT, NPROW, NPCOL, MYROW, MYCOL)
```

where details of the calling sequence for SL\_INIT are provided below.

```
SL_INIT(ICONTEXT, NPROW, NPCOL)
```

```
ICONTEXT (global output) INTEGER
```

*ICONTEXT specifies the BLACS context identifying the created process grid.*

```
NPROW (global input) INTEGER
```

*NPROW specifies the number of process rows in the process grid to be created.*

```
NPCOL (global input) INTEGER
```

*NPCOL specifies the number of process columns in the process grid to be created.*

**Example 2.2.1.** Using these routines we can elaborate the program that realizes the steps 2), 3) and 8) by means of ScaLAPACK system.

The program of this type in C++ is represented bellow.

```
#include <stdio.h>
#include <stdio.h>
#include <string.h>
#include <stdlib.h>
#include <cmath>
```

```

#include "mpi.h"
extern "C" {
// Cblacs declarations
    void Cblacs_pinfo(int*, int*);
    void Cblacs_get(int, int, int*);
    void Cblacs_gridinit(int*, const char*, int, int);
    void Cblacs_gridinfo(int, int*, int*, int*, int*);
    void Cblacs_pcoord(int, int, int*, int*);
    void Cblacs_gridexit(int);
    void Cblacs_barrier(int, const char*);
    int Cblacs_pnum(int, int, int);
}

int main(int argc, char **argv)
{
int ICTXT,MYID,NPROCS,NPROW,NPCOL;
int P,Q,MYROW,MYCOL;
int BLACS_PNUM;
int rank, namelen;
char processor_name[MPI_MAX_PROCESSOR_NAME];
MPI_Status status;
MPI_Init(&argc,&argv);
MPI_Get_processor_name(processor_name,&namelen);
MPI_Comm_rank(MPI_COMM_WORLD,&rank);
// Initializing the BLACS library (step 2))
Cblacs_pinfo(&MYID,&NPROCS);
Cblacs_get(-1, 0, &ICTXT);
// Calculations to make the most square-like proc grid possible
NPROW=(int)sqrt(NPROCS);
NPCOL=NPROCS/NPROW;
// Creating and using the processor grid (step 3))
Cblacs_gridinit(&ICTXT, "Row-major", NPROW,NPCOL);
Cblacs_gridinfo(ICTXT,&NPROW,&NPCOL,&MYROW,&MYCOL);
sleep(MYID);
    if ((MYROW==0)&(MYCOL==0))
printf("===== RESULT OF THE PROGRAM %s \n",argv[0]);
Cblacs_pcoord(ICTXT, rank, &P,&Q);
Cblacs_barrier(ICTXT, "All");
printf("Hello, I am the process number %d on the compute hostname %s , and have a
grid position (%d,%d)\n",rank, processor_name, P,Q);
Cblacs_barrier(ICTXT, "All");
if (rank ==0) printf("=== \n");
Cblacs_barrier(ICTXT, "All");
sleep(2*MYID);
printf("A grid position (%d,%d) is the process number %d on the compute hostname %s
\n",P,Q,Cblacs_pnum(ICTXT, P, Q), processor_name);
// Release the proces gird and blacs library (step 8))
Cblacs_gridexit(ICTXT);
MPI_Finalize();
return 0;
}

```

## Result of compilation and executing

```

[Hancu_B_S@hpc ScaLAPACK_Exemple_Curs_Online]$ ./mpiCC_ScL -o Example2.2.1.exe
Example2.2.1.cpp
[Hancu_B_S@hpc ScaLAPACK_Exemple_Curs_Online]$ /opt/openmpi/bin/mpirun -n 12 -host
compute-0-0,compute-0-1,compute-0-10 Example2.2.1.exe
===== RESULT OF THE PROGRAM Example2.2.1.exe
Hello, I am the process number 3 on the compute hostname compute-0-0.local , and have a grid
position (0,3)

```

```

Hello, I am the process number 0 on the compute hostname compute-0-0.local , and have a grid
position (0,0)
Hello, I am the process number 7 on the compute hostname compute-0-1.local , and have a grid
position (1,3)
Hello, I am the process number 8 on the compute hostname compute-0-10.local , and have a grid
position (2,0)
Hello, I am the process number 6 on the compute hostname compute-0-0.local , and have a grid
position (1,2)
Hello, I am the process number 4 on the compute hostname compute-0-1.local , and have a grid
position (1,0)
Hello, I am the process number 2 on the compute hostname compute-0-10.local , and have a grid
position (0,2)
Hello, I am the process number 9 on the compute hostname compute-0-0.local , and have a grid
position (2,1)
Hello, I am the process number 10 on the compute hostname compute-0-1.local , and have a grid
position (2,2)
Hello, I am the process number 11 on the compute hostname compute-0-10.local , and have a grid
position (2,3)
Hello, I am the process number 5 on the compute hostname compute-0-10.local , and have a grid
position (1,1)
Hello, I am the process number 1 on the compute hostname compute-0-1.local , and have a grid
position (0,1)
===
A grid position (0,0) is the process number 0 on the compute hostname compute-0-0.local
A grid position (0,1) is the process number 1 on the compute hostname compute-0-1.local
A grid position (0,2) is the process number 2 on the compute hostname compute-0-10.local
A grid position (0,3) is the process number 3 on the compute hostname compute-0-0.local
A grid position (1,0) is the process number 4 on the compute hostname compute-0-1.local
A grid position (1,1) is the process number 5 on the compute hostname compute-0-10.local
A grid position (1,2) is the process number 6 on the compute hostname compute-0-0.local
A grid position (1,3) is the process number 7 on the compute hostname compute-0-1.local
A grid position (2,0) is the process number 8 on the compute hostname compute-0-10.local
A grid position (2,1) is the process number 9 on the compute hostname compute-0-0.local
A grid position (2,2) is the process number 10 on the compute hostname compute-0-1.local
A grid position (2,3) is the process number 11 on the compute hostname compute-0-10.local
[Hancu_B_S@hpc ScaLAPACK_Exemple_Curs_Online]$

```