

Adventures in Supercomputing with R

Lecture 11: Distributed, I/O, ScaLAPACK

George Ostrouchov

Oak Ridge National Laboratory and University of Tennessee

2022/02/20 (updated: 2022-04-27)

Exercise 8 discussion ...

Knit RBigData/KPMS-IY4I-EX/exercise/Exercise08.Rmd and read notes.

See code:

RBigData/KPMS-IY4I-EX/mpi/mnist_svd_cv_mpi.R

and

RBigData/KPMS-IY4I-EX/mpi/mnist_svd_cv_mpi.sh.

MPI Cooperation

Helper functions coordinate by knowing `comm.rank()` and `comm.size()`

```
pbdIO::comm.chunk(15, form = "vector", all.rank = TRUE, p = 4)
```

```
## [[1]]  
## [1] 1 2 3 4  
##  
## [[2]]  
## [1] 5 6 7 8  
##  
## [[3]]  
## [1] 9 10 11 12  
##  
## [[4]]  
## [1] 13 14 15
```

* When running with `mpi run`, omit `p` parameter. Here it is used because my notes run serial R.

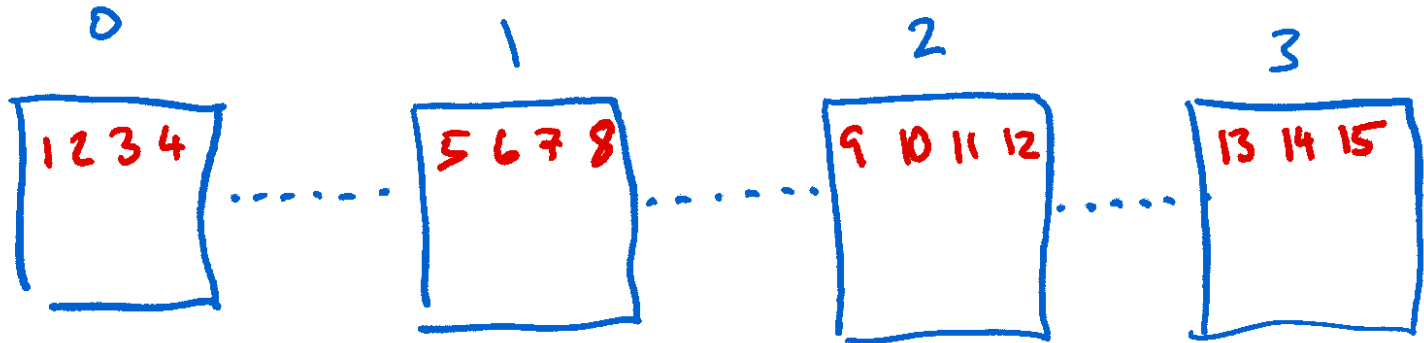
** Omitting `all.rank` delivers only the correct vector to each rank.

*** `pbdIO::` is omitted when `library(pbdIO)` was done earlier.

```
mpirun -np 4 Rscript your-code.R
```

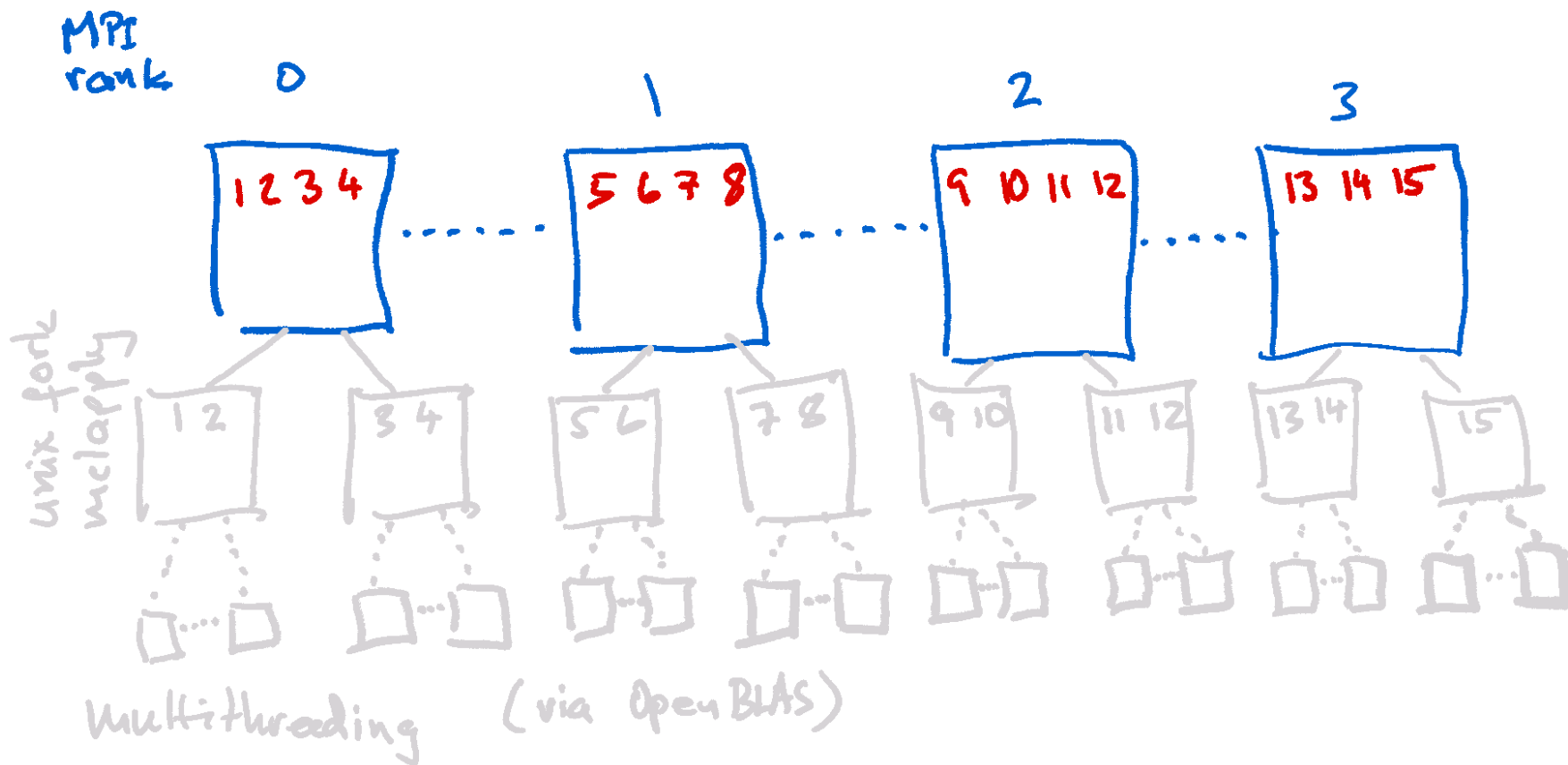
```
comm.chunk(15, form = "vector")
```

MPI
rank



```
mpirun -np 4 Rscript your-code.R`
```

```
comm.chunk(15, form = "vector")
```



Modes of parallel computing

- Shared memory: multithreading (C, C++, Fortran, etc.)
 - OpenMP Compiler directives: pragma
 - Memory efficient: One copy of memory and program
 - Available via libraries: OpenBLAS
 - Program your own via C, C++, Rcpp, etc.
- Shared memory: unix fork
 - Operating system mediated
 - Memory efficient: copy-on-write
- Distributed memory: MPI
 - Can use memory of multiple nodes
 - Replicates memory unless managed by distributing objects
 - Can operate in shared memory ("push MPI into nodes")

Parallel I/O

- A large object that fits on a node
 - Stay below 1/3 available memory
 - Gain speed by distributing
- Distributing data:
 - Each node (or rank) reads own data
 - Need parallel file system (Lustre or GPFS) for speedup
 - Karolina SCRATCH filesystem: Lustre
 - Accessible via the Infiniband network
 - <https://docs.it4i.cz/barbora/storage/#understanding-the-lustre-filesystems>
 - Number of file stripes

Parallel I/O with skinny matrix methods

- rhdf5 and kazaam packages
- Reader `KPMS-IT4I-EX/mpi/mnist_read_mpi.R`
- SVD example `KPMS-IT4I-EX/mpi/mnist_kazaam.R`

mnist_read_mpi.R

```
suppressMessages(library(rhdf5))
suppressMessages(library(pbdIO))

filename = "/scratch/project/dd-21-42/data/mnist/train.hdf5"
dataset1  = "image"
dataset2  = "label"

## get and broadcast dimensions to all processors
if (comm.rank() == 0) {
  h5f = H5Fopen(filename, flags="H5F_ACC_RDONLY")
  h5d = H5Dopen(h5f, dataset1)
  h5s = H5Dget_space(h5d)
  dims = H5Sget_simple_extent_dims(h5s)$size
  H5Dclose(h5d)
  H5Fclose(h5f)
} else dims = NA
dims = bcast(dims)

## get my local indices for contiguous data read
nlast = dims[length(dims)] # last dim moves slowest
my_ind = comm.chunk(nlast, form = "vector")

## parallel read of local data
my_train = as.double(h5read(filename, dataset1,
                             index = list(NULL, NULL, my_ind)))
my_train_lab = as.character(h5read(filename, dataset2,
                                    index = list(my_ind)))

H5close()
```

mnist_kazaam.R

```
suppressMessages(library(kazaam))
source("mnist_read_mpi.R") # reads blocks of rows

## create shaq class distributed matrix from local pieces
sq_train = shaq(my_train)

## svd (shaq class: tall-skinny matrix)
train_svd = svd(sq_train, nu = 0, nv = 5)
comm.cat("kazaam top(5) singular values:", train_svd$d[1:5], "\n")

cp_train = allreduce(crossprod(my_train))
train_svd2 = eigen(cp_train)
comm.cat("direct top(5) singular values:", sqrt(train_svd2$values[1:5]), "\n")

finalize()
```

ScaLAPACK Processor Grid and Block Cyclic Layout

RBigData/pbdDMAT package
on GitHub

$$\begin{array}{c}
 [0 \ 1 \ 2 \ 3 \ 4 \ 5] \\
 \\
 \begin{array}{ccc}
 \begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{bmatrix} & \begin{bmatrix} 0 & 1 \\ 2 & 3 \\ 4 & 5 \end{bmatrix} & \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix} \\
 \text{(a) } 1 \times 6 & \text{(b) } 2 \times 3 & \text{(c) } 3 \times 2 & \text{(d) } 6 \times 1
 \end{array}
 \end{array}$$

Table: Processor Grid Shapes with 6 Processors

x ₁₁	x ₁₂	x ₁₃	x ₁₄	x ₁₅	x ₁₆	x ₁₇	x ₁₈	x ₁₉
x ₂₁	x ₂₂	x ₂₃	x ₂₄	x ₂₅	x ₂₆	x ₂₇	x ₂₈	x ₂₉
x ₃₁	x ₃₂	x ₃₃	x ₃₄	x ₃₅	x ₃₆	x ₃₇	x ₃₈	x ₃₉
x ₄₁	x ₄₂	x ₄₃	x ₄₄	x ₄₅	x ₄₆	x ₄₇	x ₄₈	x ₄₉
x ₅₁	x ₅₂	x ₅₃	x ₅₄	x ₅₅	x ₅₆	x ₅₇	x ₅₈	x ₅₉
x ₆₁	x ₆₂	x ₆₃	x ₆₄	x ₆₅	x ₆₆	x ₆₇	x ₆₈	x ₆₉
x ₇₁	x ₇₂	x ₇₃	x ₇₄	x ₇₅	x ₇₆	x ₇₇	x ₇₈	x ₇₉
x ₈₁	x ₈₂	x ₈₃	x ₈₄	x ₈₅	x ₈₆	x ₈₇	x ₈₈	x ₈₉
x ₉₁	x ₉₂	x ₉₃	x ₉₄	x ₉₅	x ₉₆	x ₉₇	x ₉₈	x ₉₉

9 × 9

x ₁₁	x ₁₂	x ₁₇	x ₁₈	x ₁₃	x ₁₄	x ₁₉	x ₁₅	x ₁₆
x ₂₁	x ₂₂	x ₂₇	x ₂₈	x ₂₃	x ₂₄	x ₂₉	x ₂₅	x ₂₆
x ₅₁	x ₅₂	x ₅₇	x ₅₈	x ₅₃	x ₅₄	x ₅₉	x ₅₅	x ₅₆
x ₆₁	x ₆₂	x ₆₇	x ₆₈	x ₆₃	x ₆₄	x ₆₉	x ₆₅	x ₆₆
x ₉₁	x ₉₂	x ₉₇	x ₉₈	x ₉₃	x ₉₄	x ₉₉	x ₉₅	x ₉₆
x ₃₁	x ₃₂	x ₃₇	x ₃₈	x ₃₃	x ₃₄	x ₃₉	x ₃₅	x ₃₆
x ₄₁	x ₄₂	x ₄₇	x ₄₈	x ₄₃	x ₄₄	x ₄₉	x ₄₅	x ₄₆
x ₇₁	x ₇₂	x ₇₇	x ₇₈	x ₇₃	x ₇₄	x ₇₉	x ₇₅	x ₇₆
x ₈₁	x ₈₂	x ₈₇	x ₈₈	x ₈₃	x ₈₄	x ₈₉	x ₈₅	x ₈₆

5 × 4 5 × 3 5 × 2

4 × 4 4 × 3 4 × 2

```

> x <- as.rowblock( x )
> x <- as.blockcyclic( x )
> x <- redistribute( x, bldim=c(8, 8), ICTXT = 0 )

```

ScaLAPACK Complex Communication

- Several communicators (global, row wise, column wise)

Converters for Bolck-Cyclic Layout

`as.matrix(X)`



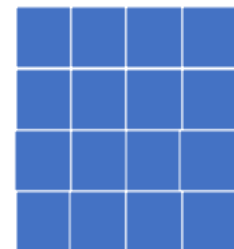
`as.rowblock(X)`



`as.colblock(X)`



`as.blockcyclic(X)`



ddmatrix Examples

KPMS-IT4I-EX/dmat