

Adventures in Supercomputing with R

Lecture 10: Distributed Matrix Computation (continued)

George Ostrouchov

Oak Ridge National Laboratory and University of Tennessee

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

qr_allreduce.r

```
library(cop, quiet = TRUE)

rank = comm.rank()
size = comm.size()

rows = 3
cols = 3
xb = matrix((1:(rows*cols*size))^2, ncol = cols) # a full matrix
xa = xb[(1:rows) + rank*rows, ] # split by row blocks

comm.print(xa, all.rank = TRUE)
comm.print(xb)

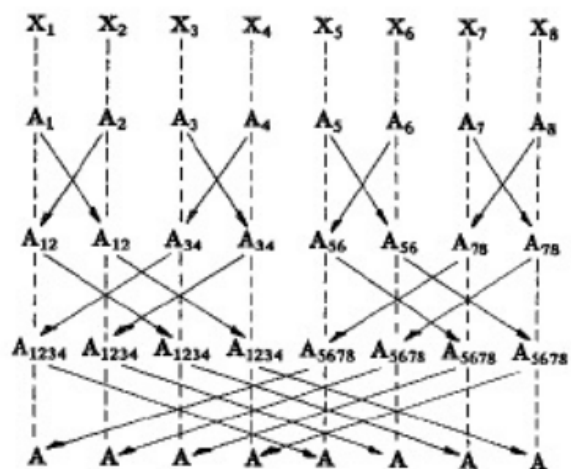
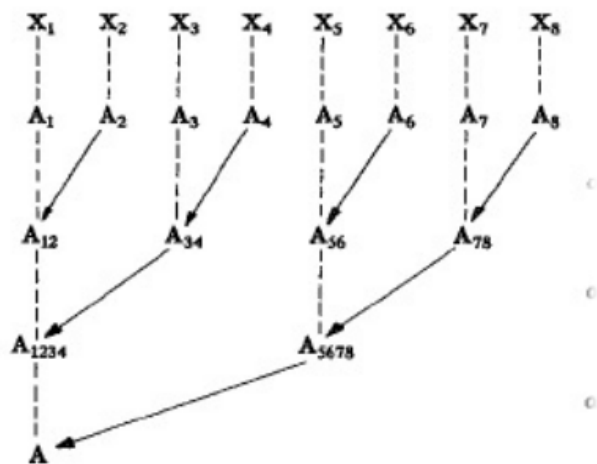
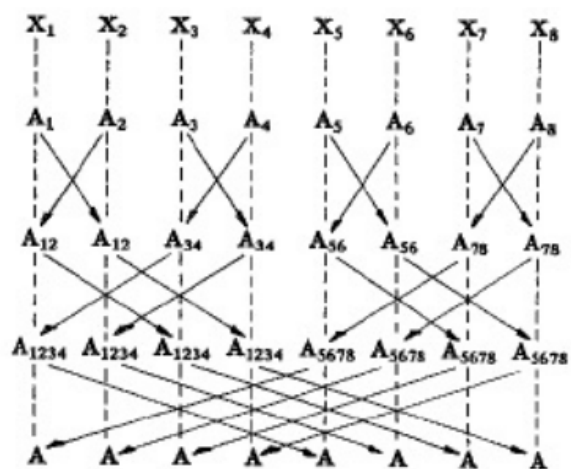
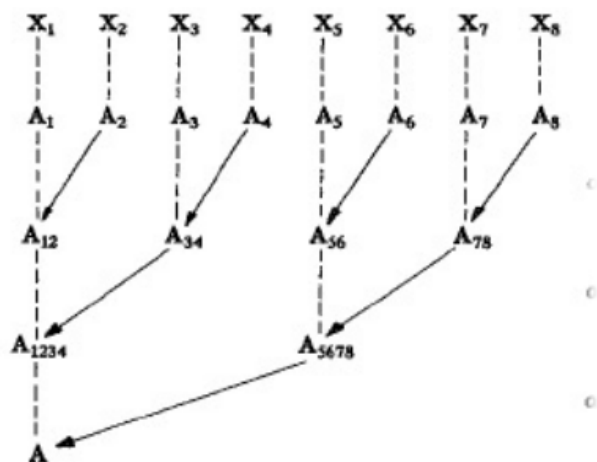
## compute usual QR from full matrix
rb = qr.R(qr(xb))
comm.print(rb)

## compute QR from gathered local QRs
rloc = qr.R(qr(xa)) # compute local QRs
rra = allgather(rloc) # gather them into a list
rra = do.call(rbind, rra) # rbind list elements
comm.print(rra) # print combined local QRs
ra = qr.R(qr(rra)) # QR the combined local QRs
comm.print(ra)

## use cop package to do it again via qr_allreduce
ra = qr_allreduce(xa)
comm.print(ra)
```

qr_allreduce.r

- cop version is faster because it does a reduce rather than a gather
- gather communication and memory grows and combined QR is serial
- reduce communication does not grow. reduce combines QRs in parallel.
reduce
- Memory use scales with local size, not global size.
- QR is unique up to multiplication by a diagonal matrix of $\{+1, -1\}$



mnist_read_mpi.r

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

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

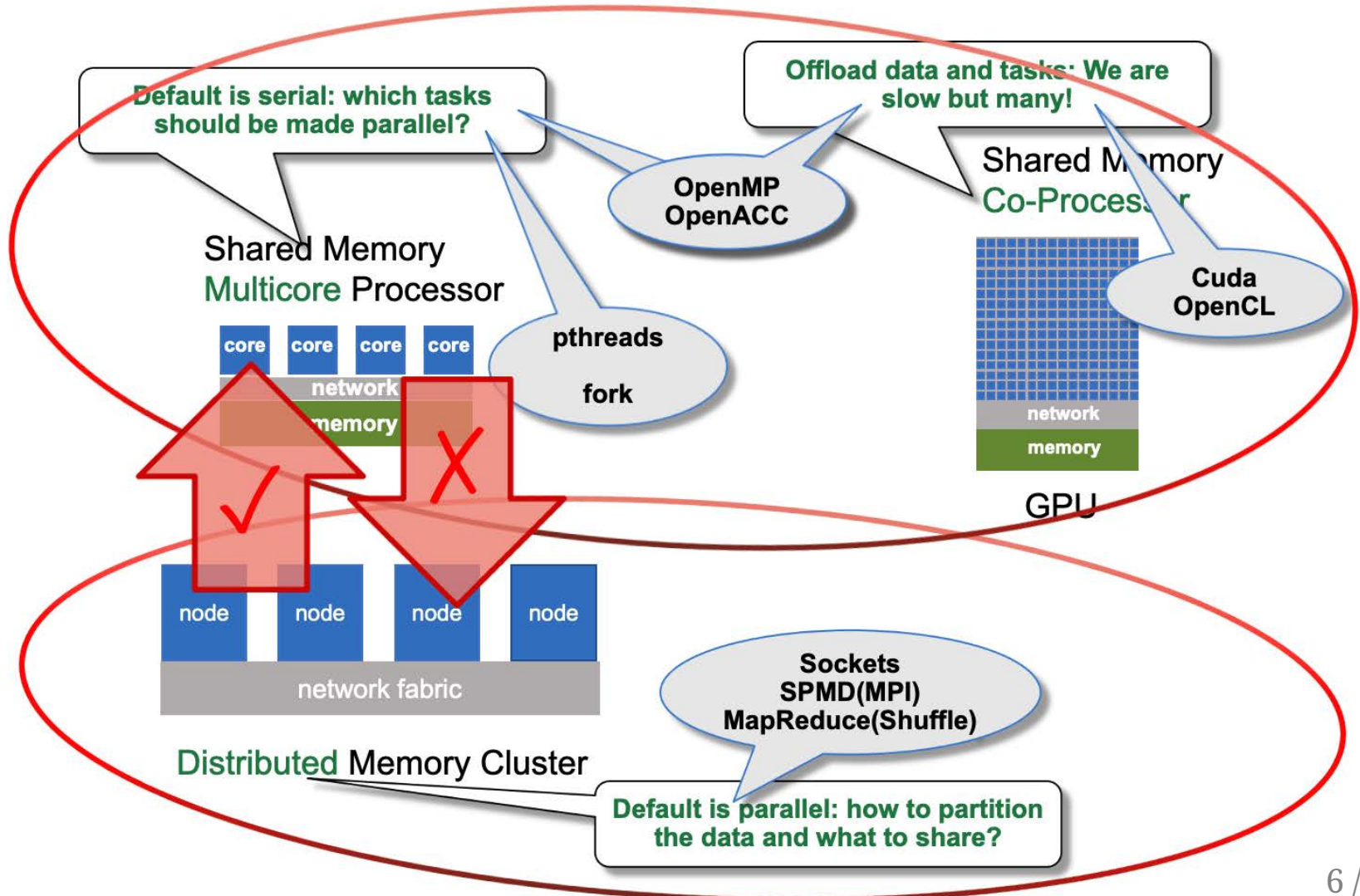
## get and broadcast dimensions to all processors
if ( comm.rank() == 0 ) {
  h5f = H5Fopen( filename, flags="H5F_ACC_RDONLY")
  h5d = H5Dopen( h5f, dataset1 )
  h5s = H5Dget_space( h5d )
  dims = H5Sget_simple_extent_dims( h5s )$size
  H5Dclose( h5d )
  H5Fclose( h5f )
} else dims = NA
dims = bcast( 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,
my_train_lab = as.character( h5read( filename, dataset2, index = list
H5close( )

dim( my_train ) = c( prod( dims[-length(dims)] ), length(my_ind) )
my_train = t( my_train ) # now it's rowblock
```

Distributed Programming Works in Shared Memory



Mapping Threads to Cores

Theory and Reality

- Operating system manages core affinity
- Operating system tasks can compete
- Core switching occurs frequently
- MPI is more rigid and more planned



pbdR-ScaLAPACK-PBLAS-BLACS-MPI



- MPI: Message Passing Interface - *de facto* standard for distributed communication in supercomputing
 - Used for data mostly via collective communication - high level
 - pbdMPI, kazaam, and cop R packages
- ScaLAPACK: Scalable LAPACK - Distributed version of LAPACK (uses PBLAS/BLAS but not LAPACK)
 - 2d Block-Cyclic data layout - mostly automated in pbdDMAT package
 - BLACS: Communication collectives for distributed matrix computation
 - PBLAS: BLAS - distributed BLAS (uses shared memory BLAS within blocks)
 - pbdDMAT and pbdML R packages - most matrix operations identical to serial through overloading operators and `ddmatrix` class

Processor Grid and Complex Block Cyclic Layout

RBigData/pbdDMAT package
on GitHub

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

Table: Processor Grid Shapes with 6 Processors

$$\begin{bmatrix}
 \begin{array}{cc} x_{11} & x_{12} \end{array} & \begin{array}{cc} x_{13} & x_{14} \end{array} & \begin{array}{cc} x_{15} & x_{16} \end{array} & \begin{array}{cc} x_{17} & x_{18} \end{array} & \begin{array}{cc} x_{19} & \end{array} \\
 \begin{array}{cc} x_{21} & x_{22} \end{array} & \begin{array}{cc} x_{23} & x_{24} \end{array} & \begin{array}{cc} x_{25} & x_{26} \end{array} & \begin{array}{cc} x_{27} & x_{28} \end{array} & \begin{array}{cc} x_{29} & \end{array} \\
 \hline
 \begin{array}{cc} x_{31} & x_{32} \end{array} & \begin{array}{cc} x_{33} & x_{34} \end{array} & \begin{array}{cc} x_{35} & x_{36} \end{array} & \begin{array}{cc} x_{37} & x_{38} \end{array} & \begin{array}{cc} x_{39} & \end{array} \\
 \begin{array}{cc} x_{41} & x_{42} \end{array} & \begin{array}{cc} x_{43} & x_{44} \end{array} & \begin{array}{cc} x_{45} & x_{46} \end{array} & \begin{array}{cc} x_{47} & x_{48} \end{array} & \begin{array}{cc} x_{49} & \end{array} \\
 \hline
 \begin{array}{cc} x_{51} & x_{52} \end{array} & \begin{array}{cc} x_{53} & x_{54} \end{array} & \begin{array}{cc} x_{55} & x_{56} \end{array} & \begin{array}{cc} x_{57} & x_{58} \end{array} & \begin{array}{cc} x_{59} & \end{array} \\
 \begin{array}{cc} x_{61} & x_{62} \end{array} & \begin{array}{cc} x_{63} & x_{64} \end{array} & \begin{array}{cc} x_{65} & x_{66} \end{array} & \begin{array}{cc} x_{67} & x_{68} \end{array} & \begin{array}{cc} x_{69} & \end{array} \\
 \hline
 \begin{array}{cc} x_{71} & x_{72} \end{array} & \begin{array}{cc} x_{73} & x_{74} \end{array} & \begin{array}{cc} x_{75} & x_{76} \end{array} & \begin{array}{cc} x_{77} & x_{78} \end{array} & \begin{array}{cc} x_{79} & \end{array} \\
 \begin{array}{cc} x_{81} & x_{82} \end{array} & \begin{array}{cc} x_{83} & x_{84} \end{array} & \begin{array}{cc} x_{85} & x_{86} \end{array} & \begin{array}{cc} x_{87} & x_{88} \end{array} & \begin{array}{cc} x_{89} & \end{array} \\
 \hline
 \begin{array}{cc} x_{91} & x_{92} \end{array} & \begin{array}{cc} x_{93} & x_{94} \end{array} & \begin{array}{cc} x_{95} & x_{96} \end{array} & \begin{array}{cc} x_{97} & x_{98} \end{array} & \begin{array}{cc} x_{99} & \end{array}
 \end{bmatrix}_{9 \times 9}$$

$$\begin{array}{c}
 \begin{bmatrix}
 \begin{array}{cc} x_{11} & x_{12} \end{array} & \begin{array}{cc} x_{17} & x_{18} \end{array} \\
 \begin{array}{cc} x_{21} & x_{22} \end{array} & \begin{array}{cc} x_{27} & x_{28} \end{array} \\
 \hline
 \begin{array}{cc} x_{51} & x_{52} \end{array} & \begin{array}{cc} x_{57} & x_{58} \end{array} \\
 \begin{array}{cc} x_{61} & x_{62} \end{array} & \begin{array}{cc} x_{67} & x_{68} \end{array} \\
 \hline
 \begin{array}{cc} x_{91} & x_{92} \end{array} & \begin{array}{cc} x_{97} & x_{98} \end{array}
 \end{bmatrix}_{5 \times 4}
 \end{array}
 \begin{array}{c}
 \begin{bmatrix}
 \begin{array}{cc} x_{13} & x_{14} \end{array} & \begin{array}{cc} x_{19} & \end{array} \\
 \begin{array}{cc} x_{23} & x_{24} \end{array} & \begin{array}{cc} x_{29} & \end{array} \\
 \hline
 \begin{array}{cc} x_{53} & x_{54} \end{array} & \begin{array}{cc} x_{59} & \end{array} \\
 \begin{array}{cc} x_{63} & x_{64} \end{array} & \begin{array}{cc} x_{69} & \end{array} \\
 \hline
 \begin{array}{cc} x_{93} & x_{94} \end{array} & \begin{array}{cc} x_{99} & \end{array}
 \end{bmatrix}_{5 \times 3}
 \end{array}
 \begin{array}{c}
 \begin{bmatrix}
 \begin{array}{cc} x_{15} & x_{16} \end{array} \\
 \begin{array}{cc} x_{25} & x_{26} \end{array} \\
 \hline
 \begin{array}{cc} x_{55} & x_{56} \end{array} \\
 \begin{array}{cc} x_{65} & x_{66} \end{array} \\
 \hline
 \begin{array}{cc} x_{95} & x_{96} \end{array}
 \end{bmatrix}_{5 \times 2}
 \end{array}$$

$$\begin{array}{c}
 \begin{bmatrix}
 \begin{array}{cc} x_{31} & x_{32} \end{array} & \begin{array}{cc} x_{37} & x_{38} \end{array} \\
 \begin{array}{cc} x_{41} & x_{42} \end{array} & \begin{array}{cc} x_{47} & x_{48} \end{array} \\
 \hline
 \begin{array}{cc} x_{71} & x_{72} \end{array} & \begin{array}{cc} x_{77} & x_{78} \end{array} \\
 \begin{array}{cc} x_{81} & x_{82} \end{array} & \begin{array}{cc} x_{87} & x_{88} \end{array}
 \end{bmatrix}_{4 \times 4}
 \end{array}
 \begin{array}{c}
 \begin{bmatrix}
 \begin{array}{cc} x_{33} & x_{34} \end{array} & \begin{array}{cc} x_{39} & \end{array} \\
 \begin{array}{cc} x_{43} & x_{44} \end{array} & \begin{array}{cc} x_{49} & \end{array} \\
 \hline
 \begin{array}{cc} x_{73} & x_{74} \end{array} & \begin{array}{cc} x_{79} & \end{array} \\
 \begin{array}{cc} x_{83} & x_{84} \end{array} & \begin{array}{cc} x_{89} & \end{array}
 \end{bmatrix}_{4 \times 3}
 \end{array}
 \begin{array}{c}
 \begin{bmatrix}
 \begin{array}{cc} x_{35} & x_{36} \end{array} \\
 \begin{array}{cc} x_{45} & x_{46} \end{array} \\
 \hline
 \begin{array}{cc} x_{75} & x_{76} \end{array} \\
 \begin{array}{cc} x_{85} & x_{86} \end{array}
 \end{bmatrix}_{4 \times 2}
 \end{array}$$

```

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

```

Exercise 8 ...

MPI version of crossvalidation for SVD basis functions

Added debugging strategies to Exercise 08 ...

More Debugging strategies

- Login with two terminal windows to cluster
 - Window1: submission, status, etc
 - Submit job
 - Check status, etc.
 - Window2: monitor with `top` on compute node
 - `check-pbs-jobs --check-all`
 - `ssh cn###`
 - `top`
- Multithreaded OpenBLAS show up as using more than 100% cpu
 - Example: with 4 threads, will show a little over 300%
- Fork with `mcLappy` will show as separate R process
- MPI rank will show as separate R process

Only 4 weeks remain!

Remaining topics to learn

- More distributed matrix computation
 - kazaam and pbdDMAT packages
- Distributed I/O
 - csv files, hdf5 files, and others

Scaling your research on IT4I

- **Your group:** ~5 minute presentation of your problem that needs HPC
- **All:** Discuss how to do it
- See Exercise 10

The Power of Speed

- < 1 minute: Not getting bored
- < 5 seconds: Can optimize parameters
- < 1 second: can provide automated optimization
- It's not just memory of a larger system
- Speed can bring optimization and search options that are otherwise difficult or even unthinkable