



Rhpc: An R package for High-Performance Computing

Eiji NAKAMA [†] Junji NAKANO [‡]

[†]COM-ONE Ltd., Japan

[‡]The Institute of Statistical Mathematics , Japan

Workshop on Distributed Computing in R
Jan 26-27, 2015
HP Labs, Palo Alto, CA, USA



Outline

- 1 Introduction
- 2 Rhpc
- 3 NUMA effect
- 4 Concluding remarks

Outline

- 1 Introduction
- 2 Rhpc
- 3 NUMA effect
- 4 Concluding remarks

Introduction

- The Institute of Statistical Mathematics possesses three supercomputers. Their types are all different: a shared memory system, a distributed memory system and a cloud system.
- We hope to use R on them efficiently. For this purpose, we have developed a package `Rhpc`, which stands for “R for High Performance Computing”.
- `Rhpc` follows the way of `snow` package, and is tuned for just using MPI.
- We also noticed that NUMA (Non-Uniform Memory Access) support is important for present parallel computing on R.

Outline

- 1 Introduction
- 2 Rhpc
- 3 NUMA effect
- 4 Concluding remarks

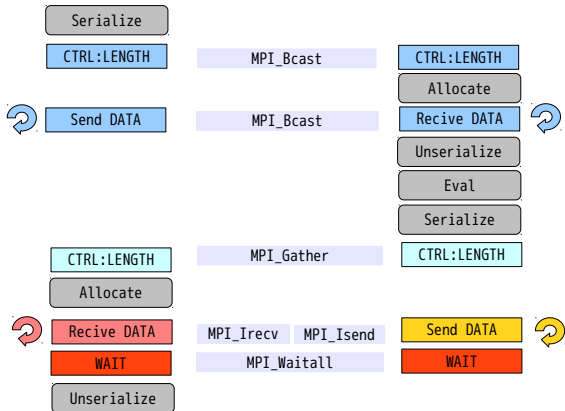
Rhpc functions (1)

- Rhpc supports long vector, so it can handle “big data” .
- Worker process is written by Embedding R.
- Functions to initialize and finalize MPI, and get MPI communication handle
 - Rhpc_initialize
 - Rhpc_getHandle (~ snow::makeMPICluster)
 - Rhpc_finalize (~ snow::stopCluster)
- Functions to call Rhpc_worker
 - Rhpc_worker_call (~ snow::clusterCall)
 - Rhpc_Export (~ snow::clusterExport)
 - Rhpc_EvalQ (~ snow::clusterEvalQ)

Rhpc functions (2)

- Apply type functions
 - Rhpc_lapply (~ snow::clusterApply)
 - Rhpc_lapplyLB (~ snow::clusterApplyLB)
- Function to set up random number generators
 - Rhpc_setupRNG (~ snow::clusterSetupRNGstream)
- Miscellaneous functions
 - Rhpc_worker_noback
(Experimental I/F for calling foreign SPMD program)
 - Rhpc_serialize, Rhpc_unserialize
(C functions for serialization and unserialization)

Rhpc_worker_call (2)

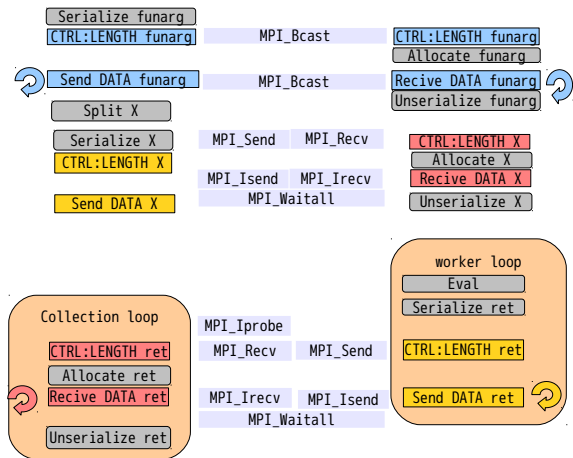


Rhpc_lapply (1)

Rhpc_lapply(cl, x, FUN, ...)

- cl
pointer to communicator
- x
vector or list.
Divided into smaller vectors according to the number of workers, and distributed to workers when the function is first executed.
One-sided communication is used asynchronously.
- FUN
Function name or string (string expresses function name)
Distributed by **collective communication** at first, then they **are not sent again**
- ... (argument)
Distributed by **collective communication** at first, then they **are not sent again.**

Rhpc_lapply (2)

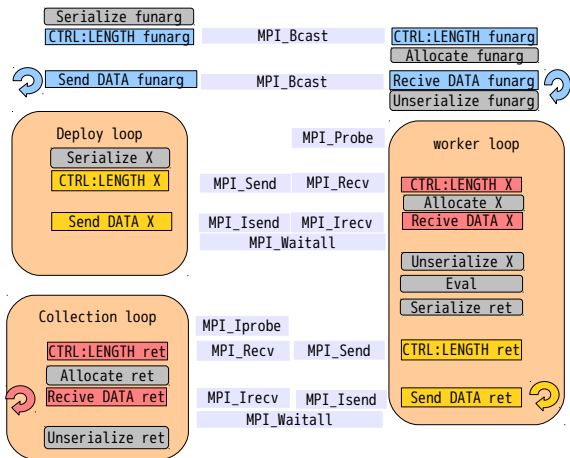


Rhpc_lapplyLB (1)

Rhpc_lapplyLB(cl, x, FUN, ...)

- cl
pointer to communicator
- x
vector or list.
Distributed to workers when the function is executed. **One-sided communication is used asynchronously.**
- FUN
Function name or string (string expresses function name)
Distributed by **collective communication** at first, then they **are not sent again.**
- ... (argument)
Distributed by **collective communication** at first, then they **are not sent again.**

Rhpc_lapplyLB (2)



Huge data example: snow (SOCK)

As snow(SOCK) utilizes pipe to serialize and unserialize data, it can handle huge data. However, **socket mechanism has difficulty to control processes on many nodes by many users in supercomputer environment**

use snow(SOCK)::clusterExport in one worker

```
> library(snow)
> cl<-makeCluster(1,type="SOCK")
> set.seed(123)
> N<-17e3
> M<-matrix(runif(N^2),N,N)
> sum(M)
[1] 144501466
> system.time(clusterExport(cl,"M"))
  user  system elapsed
 4.213   1.580   8.228
> f<-function()sum(M)
> clusterCall(cl,f)
[[1]]
[1] 144501466
```

Huge data example: snow(MPI)

At present, Rmpi cannot handle data more than 2GB, because the argument of MPI_send etc. should be int size.

use snow(MPI)::clusterExport in one worker

```
> library(snow)
> cl<-makeCluster(1,type="MPI")
> set.seed(123)
> N<-17e3
> M<-matrix(runif(N^2),N,N)
> sum(M)
[1] 144501466
> system.time(clusterExport(cl,"M"))
Error in mpi.send(x = serialize(obj, NULL), type = 4, dest = dest, tag = tag, :
  long vectors not supported yet: memory.c:3100
Calls: system.time ... sendData.MPIinode -> mpi.send.Robj -> mpi.send -> .Call
```

Huge data example: Rhpc

Rhpc can handle huge data, because it divides huge data into appropriate size and uses MPI repeatedly.

use Rhpc::Rhpc_Export in one worker

```
> library(Rhpc)
> Rhpc_initialize()
> cl<-Rhpc_getHandle()
> set.seed(123)
> N<-17e3
> M<-matrix(runif(N^2),N,N)
> sum(M)
[1] 144501466
> system.time(Rhpc_Export(cl,"M"))
  user system elapsed
 9.241  1.700 10.972
> f<-function()sum(M)
> Rhpc_worker_call(cl,f)
[[1]]
[1] 144501466
> Rhpc_finalize()
```


Many workers example (1): snow(MPI)

As clusterCall of snow starts workers sequentially, it becomes slow when the number of workers increases.

use snow(MPI)::clusterExport in 63 workers

```
> library(Rmpi)
> library(snow)
> cl<-makeMPIcluster()
> set.seed(123)
> N<-4e3
> length(cl)
[1] 63
> M<-matrix(runif(N^2),N,N)
> system.time(clusterExport(cl,"M"))
  user  system elapsed
26.715  10.903  37.761
> f<-function()sum(M)
> all.equal(rep(sum(M),length(cl)),unlist(clusterCall(cl,f)))
[1] TRUE
> stopCluster(cl)
```

Many workers example (1): Rhpc

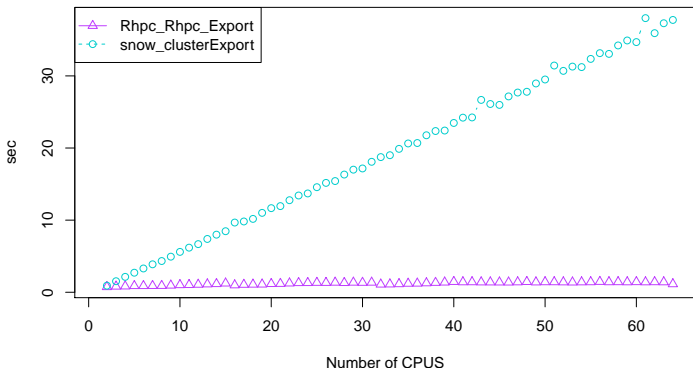
As Rhpc uses collective communication by MPI_Bcast, data transportation to workers is still fast even when the number of workers increases.

use `Rhpc::Rhpc_Export` in **63** workers

```
> library(Rhpc)
Loading required package: rlecuyer
> Rhpc_initialize()
> cl<-Rhpc_getHandle()
Detected communication size 64
> set.seed(123)
> N<-4e3
> Rhpc_numberOfWorkers(cl)
[1] 63
> M<-matrix(runif(N^2),N,N)
> system.time(Rhpc_Export(cl,"M"))
  user  system elapsed
1.012  0.116  1.139
> f<-function()sum(M)
> all.equal(rep(sum(M),Rhpc_numberOfWorkers(cl)),unlist(Rhpc_worker_call(cl,f)))
[1] TRUE
> Rhpc_finalize()
```

Many workers example (1): Rhpc and snow(MPI)

Export performance



Many workers example (2a): snow(MPI)

As the main parts of snow and Rmpi are written in R language, they are rather slow.

use snow(MPI)::clusterApply in 63 workers

```
> library(Rmpi)
> library(snow)
> cl<-makeMPIcluster()
> system.time(ans<-clusterApply(cl,1:10000,sqrt))
  user system elapsed
 1.423   0.005   1.429
> all.equal(sqrt(1:10000),unlist(ans))
[1] TRUE
> stopCluster(cl)
```

Many workers example (2b): snow(MPI)

As the load balancing function is written in R language, it becomes slow according to the number of parallel workers.

use snow(MPI)::clusterApplyLB in 63 workers

```
> library(Rmpi)
> library(snow)
> cl<-makeMPIcluster()
> system.time(ans<-clusterApplyLB(cl,1:10000,sqrt))
  user  system elapsed
4.395   0.003   4.413
> all.equal(sqrt(1:10000),unlist(ans))
[1] TRUE
> stopCluster(cl)
```

Many workers example (2a): Rhpc

As the main part of Rhpc is written in C language, it is efficient.

use Rhpc::Rhpc_lapply in 63 workers

```
> library(Rhpc)
Loading required package: rlecuyer
> Rhpc_initialize()
> cl<-Rhpc_getHandle()
Detected communication size 64
> system.time(ans<-Rhpc_lapply(cl, 1:10000, sqrt))
  user  system elapsed
0.045  0.001  0.046
> all.equal(sqrt(1:10000),unlist(ans))
[1] TRUE
> Rhpc_finalize()
```

Many workers example (2b): Rhpc

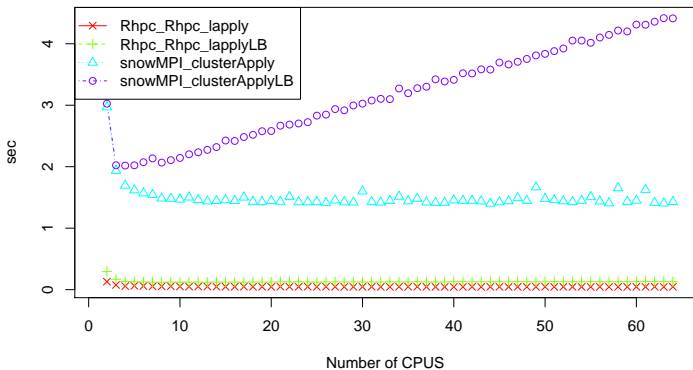
`Rhpc_lapplyLB` sends the argument `x` to an available worker. So it produces a little delay.

use `Rhpc::Rhpc_lapplyLB` in 63 workers

```
> library(Rhpc)
Loading required package: rlecuyer
> Rhpc_initialize()
> cl<-Rhpc_getHandle()
Detected communication size 64
> system.time(ans<-Rhpc_lapplyLB(cl, 1:10000, sqrt))
  user  system elapsed 
0.125  0.001  0.127 
> all.equal(sqrt(1:10000),unlist(ans))
[1] TRUE
> Rhpc_finalize()
```

Many workers example (2): Rhpc and snow(MPI)

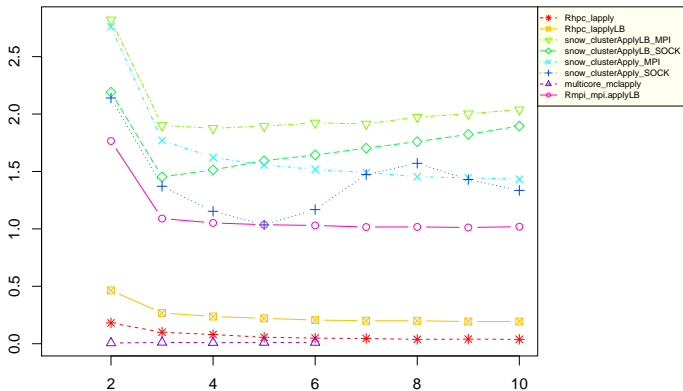
SQRT-Apply performance on Super-Computer



Several parallel apply functions

At present, Rhpc is a little slower than multicore on one CPU.

SQRT on Apply performance



Interface to foreign programs using MPI

C and/or Fortran programs using MPI adopt SPMD programming, in which Master(rank0) communicates with Workers(their ranks are more than 1) using communicators. Therefore, when Rhpc, which uses MPI, call such programs, exchanging MPI communicators is required. Rhpc provides Rhpc_worker_noback function to support it.

Rhpc sets global variable `_options_` (see options function):

- `Rhpc.mpi.f.comm`
Comuunicator for Fortran (R type: integer)
- `Rhpc.mpi.c.comm`
Communicator for C (R type: external pointer)
- `Rhpc.mpi.procs`
Communication size in MPI
- `Rhpc.mpi.rank`
Rank in MPI

Using '.Fortran', '.C' and '.Call' in R

Foreign programs (Fortran or C) using MPI are called in R

```

1 mpipif<-function(n)
2 {
3   ## Exported functions get values by getOption()
4   ## when they run on workers
5   out<-Fortran("mpipif",
6               comm=getOption("Rhpc.mpi.f.comm"),
7               n=as.integer(n),
8               outpi=as.double(0))
9   out$outpi
10 }
```

```

1 mpipic<-function(n)
2 {
3   ## Exported functions get values by getOption()
4   ## when they run on workers
5   out<-C("mpipic",
6          comm=getOption("Rhpc.mpi.f.comm"),
7          n=as.integer(n),
8          outpi=as.double(0))
9   out$outpi
10 }
```

```

1 mpipicall<-function(n)
2 {
3   ## Exported functions get values by getOption()
4   ## when they run on workers
5   out<-Call("mpipicall",
6            comm=getOption("Rhpc.mpi.c.comm"),
7            n=as.integer(n))
8   out
9 }
```

Note that '.C' cannot receive external pointer and integer communicator for Fortran should be used. Thus '.Call' is preferable for C programs.

Old serialization

At first, we realized serialization by calling R interpreter from C program as below. It was slow.

one-time call serialize

```
SXP boo = LCONS(install("serialize"),
                CONS(args,
                    CONS(R_NilValue,
                        CONS(ScalarLogical(FALSE),
                            CONS(R_NilValue,
                                CONS(R_NilValue,
                                    CONS(R_NilValue,
                                        R_NilValue))))));
boo = LCONS(install(".Internal"), CONS(boo, R_NilValue));
SXP ret = R_tryEval(boo, R_GlobalEnv, &errorOccurred);
```


New serialization

We extract code for serialization from R source code. It is faster about 5 times than the old serialization program. This idea is used to produce RApiSerialize package by Dirk Eddelbuettel.

new call serialize

```
out=Rhpc_serialize(args);
```



Outline

- 1 Introduction
- 2 Rhpc
- 3 NUMA effect**
- 4 Concluding remarks

Package RhpcBLASctl

- Functions to control the number of threads on BLAS, MKL, ACML and GotoBLAS(or OpenBLAS)
 - `blas_get_num_procs()`
 - `blas_set_num_threads(threads)`
- Functions to control the number of threads on OpenMP
 - `omp_get_num_procs()`
 - `omp_get_max_threads()`
 - `omp_set_num_threads(threads)`
- Functions to get the number of physical and logical cores
 - `get_num_cores()` # physical cores
 - `get_num_procs()` # logical cores



Simple example of RhpcBLASctl

Number of threads for parallel BLAS can be controlled to reserve cores.

parallel BLAS control

```
> library(RhpcBLASctl)
> A<-matrix(runif(3e3*3e3),3e3,3e3)
> blas_set_num_threads(1)
> blas_get_num_procs()
[1] 1
> system.time(A%*%A)
  user system elapsed
 3.992  0.012  4.003
> get_num_cores()
[1] 4
> blas_set_num_threads(get_num_cores())
> system.time(A%*%A)
  user system elapsed
 4.332  0.048  1.234
```


Example to use MPI and parallel BLAS: JOB script

gemm_s.sh: JOB script

```
#!/bin/bash
#PBS -q ice-s
#PBS -l select=1:ncpus=24:mem=120gb
#PBS -N gemm_s
#PBS -o gemm_s.out
#PBS -e gemm_s.err
#PBS -V
cd ${PBS_O_WORKDIR}
OMP_NUM_THREADS=1
OMP_DYNAMIC=FALSE
MKL_DYNAMIC=FALSE
export MKL_DYNAMIC
export OMP_DYNAMIC
export OMP_NUM_THREADS
mpirun --report-bindings --bind-to-socket --bysocket -np 3 \
/home/nakama/R/x86_64-unknown-linux-gnu-library/3.1/Rhpc/Rhpc \
CMD BATCH --no-save gemm.R gemm_s.Rout
```

Note that

- `MKL_DYNAMIC=FALSE`
specifies the number of threads manually
- `mpirun --bind-to-socket` option
specifies binding processes to processor sockets

Example to use MPI and parallel BLAS: Result

gemm_s.Rout:

```
> library(RhpcBLASctl)
> library(Rhpc)
> Rhpc_initialize()
reload mpi library /home/nakama/lib64/libmpi.so.1
rank 0/ 3(0) : r5i2n7 : 6750
> cl<-Rhpc_getHandle()
Detected communication size 3
> dummy<-Rhpc_EvalQ(cl,library(RhpcBLASctl))
> dummy<-Rhpc_worker_call(cl,blas_set_num_threads,8)
> unlist(Rhpc_worker_call(cl,blas_get_num_procs))
[1] 8 8
> f<-function(M) {
+ res<-sapply(1:5,function(x){x;system.time(M%*%M)
+ [[3]]})
+ sort(res)
+ }
> N<-3e3
> M<-matrix(runif(N*N),N,N)
> system.time(ans<-Rhpc_worker_call(cl,f,M))
user system elapsed
0.828 2.008 2.838
> (unlist(ans))
[1] 0.375 0.376 0.376 0.376 0.986 0.386 0.387 0.387
0.387 0.995
> Rhpc_finalize()
```

- We see elapsed times to check whether MKL works by using 8 threads or not
- If we use `-bind-to-core` option of `mpirun`, just 1 thread is available

