

# Programming with Big Data in R

Drew Schmidt and George Ostrouchov

useR! 2014



Wei-Chen Chen<sup>1</sup>  
George Ostrouchov<sup>2,3</sup>  
Pragneshkumar Patel<sup>3</sup>  
Drew Schmidt<sup>3</sup>



## Support

This work used resources of [National Institute for Computational Sciences](#) at the University of Tennessee, Knoxville, which is supported by the Office of Cyberinfrastructure of the U.S. National Science Foundation under Award No. ARRA-NSF-OCI-0906324 for NICS-RDAV center. This work also used resources of the [Oak Ridge Leadership Computing Facility](#) at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.

<sup>1</sup>Department of Ecology and Evolutionary Biology  
University of Tennessee, Knoxville TN, USA

<sup>2</sup>Computer Science and Mathematics Division  
Oak Ridge National Laboratory, Oak Ridge TN, USA

<sup>3</sup>Joint Institute for Computational Sciences  
University of Tennessee, Knoxville TN, USA



# About This Presentation

## Downloads

This presentation is available at: <http://r-pbd.org/tutorial>

# About This Presentation

## Installation Instructions

Installation instructions for setting up a **pbdR** environment are available:

<http://r-pbd.org/install.html>

This includes instructions for installing R, MPI, and **pbdR**.

# Contents

- 1 Introduction
- 2 Profiling and Benchmarking
- 3 The pbdR Project
- 4 Introduction to pbdMPI
- 5 The Generalized Block Distribution
- 6 Basic Statistics Examples
- 7 Data Input
- 8 Introduction to pbdDMAT and the ddmatrix Structure
- 9 Examples Using pbdDMAT
- 10 MPI Profiling
- 11 Wrapup

# Contents

## 1 Introduction

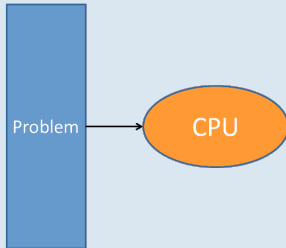
- A Concise Introduction to Parallelism
- A Quick Overview of Parallel Hardware
- A Quick Overview of Parallel Software
- Summary

## 1 Introduction

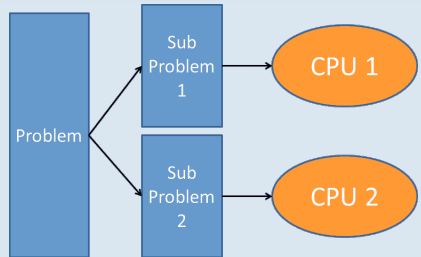
- A Concise Introduction to Parallelism
- A Quick Overview of Parallel Hardware
- A Quick Overview of Parallel Software
- Summary

## Parallelism

### Serial Programming



### Parallel Programming





## Difficulty in Parallelism

- ① *Implicit parallelism*: Parallel details hidden from user
- ② *Explicit parallelism*: Some assembly required. . .
- ③ *Embarrassingly Parallel*: Also called *loosely coupled*. Obvious how to make parallel; lots of independence in computations.
- ④ *Tightly Coupled*: Opposite of embarrassingly parallel; lots of dependence in computations.

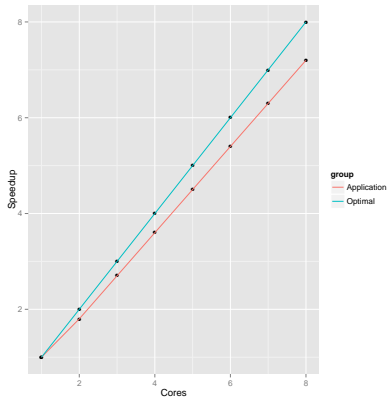
## Speedup

- *Wallclock Time*: Time of the clock on the wall from start to finish
- *Speedup*: unitless measure of improvement; more is better.

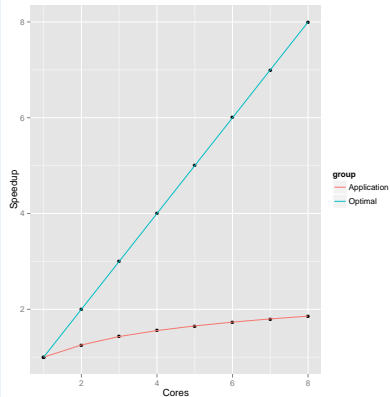
$$S_{n_1, n_2} = \frac{\text{Run time for } n_1 \text{ cores}}{\text{Run time for } n_2 \text{ cores}}$$

- $n_1$  is often taken to be 1
- In this case, comparing parallel algorithm to serial algorithm

## Good Speedup



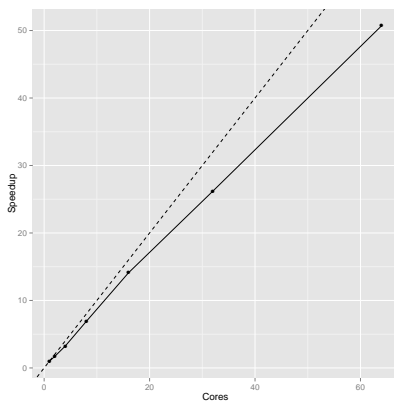
## Bad Speedup



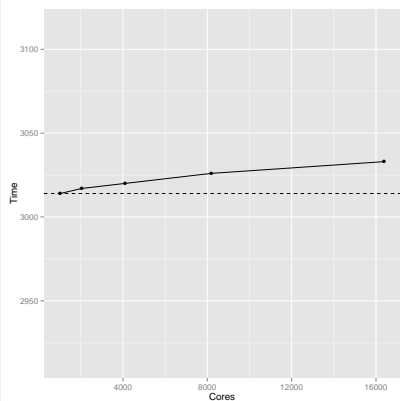
## Scalability and Benchmarking

- 1 *Strong*: Fixed **total** problem size.  
Less work per core as more cores are added.
- 2 *Weak*: Fixed **local** (per core) problem size.  
Same work per core as more cores are added.

## Good Strong Scaling



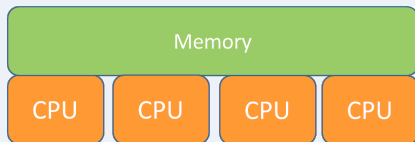
## Good Weak Scaling



# Shared and Distributed Memory Machines

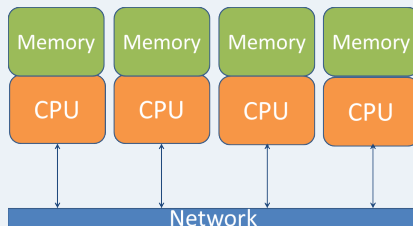
## Shared Memory

Direct access to read/change memory (one node)



## Distributed

No direct access to read/change memory (many nodes); requires communication



# Shared and Distributed Memory Machines

## Shared Memory Machines

Thousands of cores



*Nautilus*, University of Tennessee  
1024 cores  
4 TB RAM

## Distributed Memory Machines

Hundreds of thousands of cores



*Kraken*, University of Tennessee  
112,896 cores  
147 TB RAM

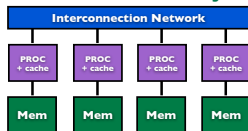
## 1 Introduction

- A Concise Introduction to Parallelism
- A Quick Overview of Parallel Hardware
- A Quick Overview of Parallel Software
- Summary

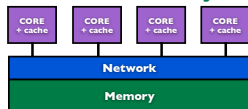


# Three Basic Flavors of Hardware

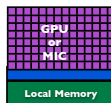
## Distributed Memory



## Shared Memory



## Co-Processor

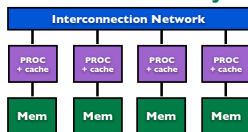


GPU: Graphical Processing Unit

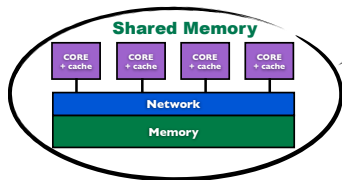
MIC: Many Integrated Core

# Your Laptop or Desktop

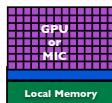
## Distributed Memory



## Shared Memory



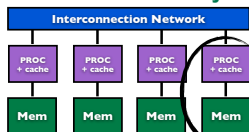
## Co-Processor



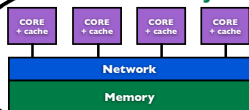
GPU: Graphical Processing Unit  
MIC: Many Integrated Core

# A Server or Cluster

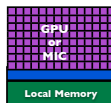
## Distributed Memory



## Shared Memory



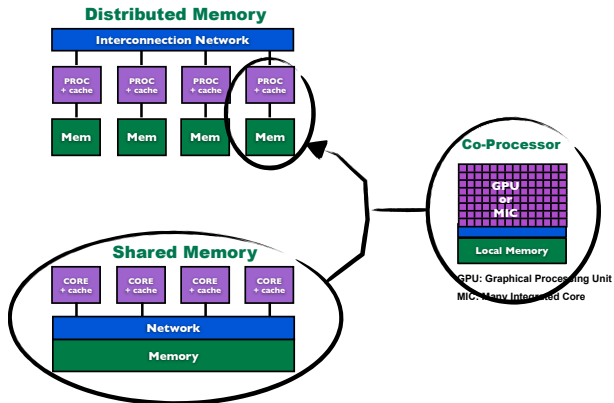
## Co-Processor



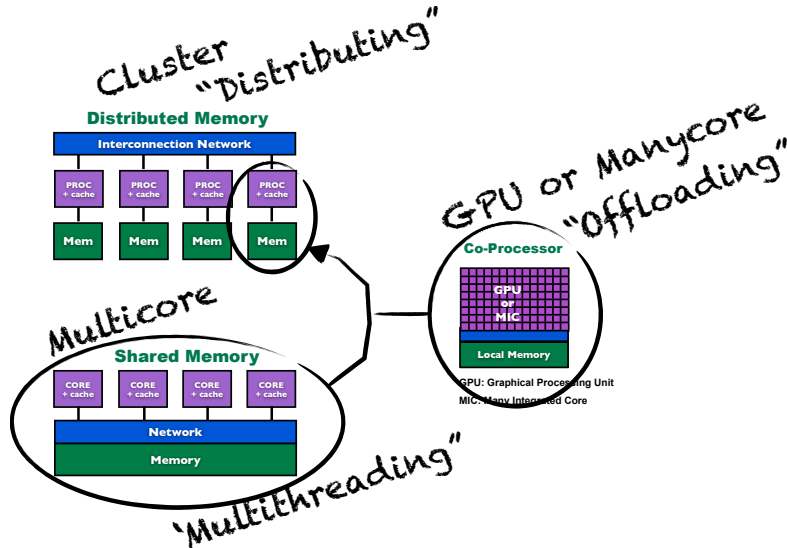
GPU: Graphical Processing Unit

MIC: Many Integrated Core

# Server to Supercomputer



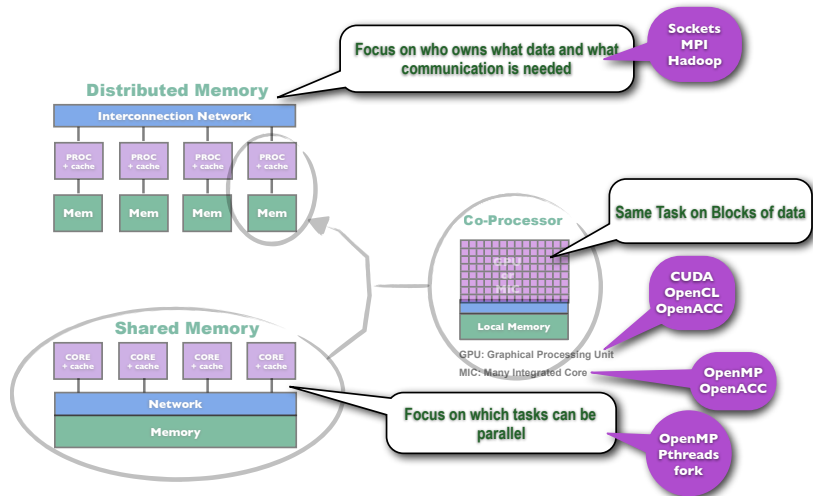
# Knowing the Right Words



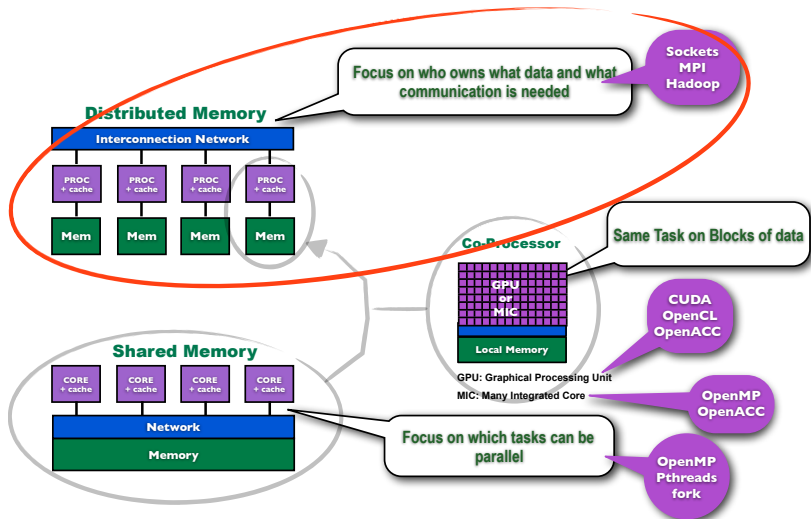
## 1 Introduction

- A Concise Introduction to Parallelism
- A Quick Overview of Parallel Hardware
- A Quick Overview of Parallel Software
- Summary

# "Native" Programming Models and Tools

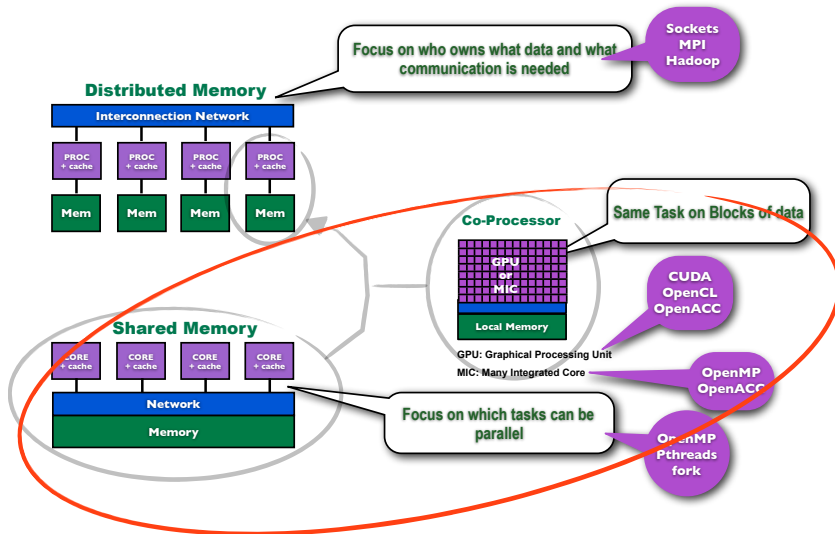


# 30+ Years of Parallel Computing Research

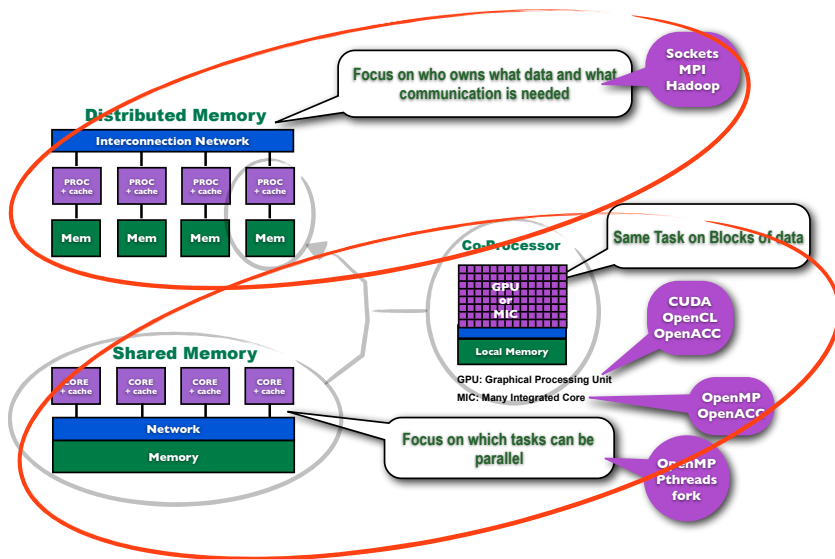




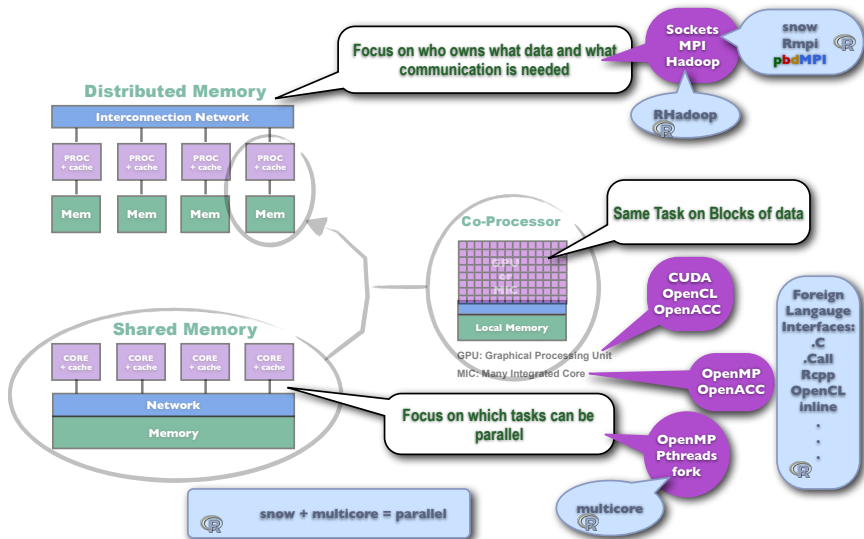
# Last 10 years of Advances



# Putting It All Together Challenge



# R Interfaces to Native Tools



## 1 Introduction

- A Concise Introduction to Parallelism
- A Quick Overview of Parallel Hardware
- A Quick Overview of Parallel Software
- Summary

## Summary

- Three flavors of hardware
  - Distributed is stable
  - Multicore and co-processor are evolving
  - Two memory models
  - Distributed works in multicore
- Parallelism hierarchy
- Medium to big machines have all three

# Contents

## 2 Profiling and Benchmarking

- Why Profile?
- Profiling R Code
- Advanced R Profiling
- Summary

## 2 Profiling and Benchmarking

- Why Profile?
- Profiling R Code
- Advanced R Profiling
- Summary

## Performance and Accuracy



*Sometimes  $\pi = 3.14$  is (a) infinitely faster than the “correct” answer and (b) the difference between the “correct” and the “wrong” answer is meaningless. . . . The thing is, some specious value of “correctness” is often irrelevant because it doesn’t matter. While performance almost always matters. And I absolutely detest the fact that people so often dismiss performance concerns so readily.*

— Linus Torvalds, August 8, 2008



## Why Profile?

- Because performance matters.
- Bad practices scale up!
- Your bottlenecks may surprise you.
- Because R is dumb.
- R users claim to be data people... so act like it!

# Compilers often correct bad behavior...

## A Really Dumb Loop

```
int main(){
    int x, i;
    for (i=0; i<10; i++)
        x = 1;
    return 0;
}
```

## clang -O3 example.c

```
main:
    .cfi_startproc
# BB#0:
    xorl    %eax,
        %eax
    ret
```

## clang example.c

```
main:
    .cfi_startproc
# BB#0:
    movl    $0, -4(%rsp)
    movl    $0, -12(%rsp)
.LBB0_1:
    cmpl    $10, -12(%rsp)
    jge     .LBB0_4
# BB#2:
    movl    $1, -8(%rsp)
# BB#3:
    movl    -12(%rsp), %eax
    addl    $1, %eax
    movl    %eax, -12(%rsp)
    jmp     .LBB0_1
.LBB0_4:
    movl    $0, %eax
    ret
```

# R will not!

## Dumb Loop

```
1 for (i in 1:n){  
2   tA <- t(A)  
3   Y <- tA %*% Q  
4   Q <- qr.Q(qr(Y))  
5   Y <- A %*% Q  
6   Q <- qr.Q(qr(Y))  
7 }  
8  
9 Q
```

## Better Loop

```
1   tA <- t(A)  
2  
3 for (i in 1:n){  
4   Y <- tA %*% Q  
5   Q <- qr.Q(qr(Y))  
6   Y <- A %*% Q  
7   Q <- qr.Q(qr(Y))  
8 }  
9  
10 Q
```

# Example from a Real R Package

## Exerpt from Original function

```
1 while(i<=N){  
2   for(j in 1:i){  
3     d.k <- as.matrix(x)[l==j,l==j]  
4     ...
```

## Exerpt from Modified function

```
1 x.mat <- as.matrix(x)  
2  
3 while(i<=N){  
4   for(j in 1:i){  
5     d.k <- x.mat[l==j,l==j]  
6     ...
```

By changing just 1 line of code, performance of the main method improved by **over 350%**!

## Some Thoughts

- R is slow.
- Bad programmers are slower.
- R isn't very clever (compared to a compiler).
- The Bytecode compiler helps, but not nearly as much as a compiler.

## 2 Profiling and Benchmarking

- Why Profile?
- Profiling R Code
- Advanced R Profiling
- Summary

## Timings

Getting simple timings as a basic measure of performance is easy, and valuable.

- `system.time()` — timing blocks of code.
- `Rprof()` — timing execution of R functions.
- `Rprofmem()` — reporting memory allocation in R .
- `tracemem()` — detect when a copy of an R object is created.
- The **rbenchmark** package — Benchmark comparisons.

## 2 Profiling and Benchmarking

- Why Profile?
- Profiling R Code
- Advanced R Profiling
- Summary



## Other Profiling Tools

- perf (Linux)
- PAPI
- MPI profiling: fpmapi, mpiP, TAU

# Profiling MPI Codes with pbdPROF

## 1. Rebuild pbdR packages

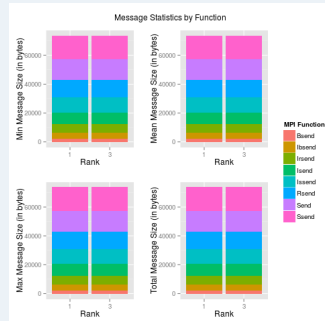
```
R CMD INSTALL pbdMPI_0.2-1.tar.gz \
  --configure-args= \
  "--enable-pbdPROF"
```

## 2. Run code

```
mpirun -np 64 Rscript my_script.R
```

## 3. Analyze results

```
1 library(pbdPROF)
2 prof <- read.prof( "output.mpiP")
3 plot(prof, plot.type="messages2")
```



## Profiling with pbdPAPI

- Performance Application Programming Interface
- High and low level interfaces
- Linux only :(



Function	Description of Measurement
<code>system.flips()</code>	Time, floating point instructions, and Mflips
<code>system.flops()</code>	Time, floating point operations, and Mflops
<code>system.cache()</code>	Cache misses, hits, accesses, and reads
<code>system.epc()</code>	Events per cycle
<code>system.idle()</code>	Idle cycles
<code>system.cpuormem()</code>	CPU or RAM bound*
<code>system.utilization()</code>	CPU utilization*

## 2 Profiling and Benchmarking

- Why Profile?
- Profiling R Code
- Advanced R Profiling
- Summary

## Summary

- *Profile, profile, profile.*
- Use `system.time()` to get a general sense of a method.
- Use **rbenchmark**'s `benchmark()` to compare 2 methods.
- Use `Rprof()` for more detailed profiling.
- Other tools exist for more hardcore applications (**pbdPAPI** and **pbdPROF**).

# Contents

- 3 The pbdR Project
  - The pbdR Project
  - pbdR Connects R to HPC Libraries
  - Using pbdR
  - Summary

- ### 3 The pbdR Project
- The pbdR Project
  - pbdR Connects R to HPC Libraries
  - Using pbdR
  - Summary

## Programming with Big Data in R (pbdR)

Striving for *Productivity, Portability, Performance*



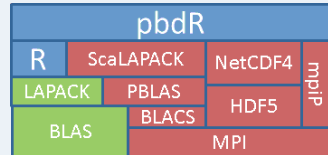
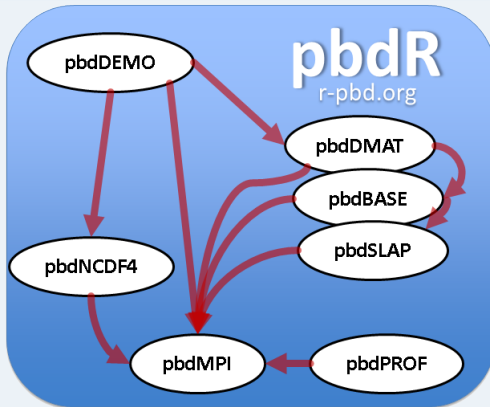
- *Free*<sup>a</sup> R packages.
- Bridging high-performance compiled code with high-productivity of R
- Scalable, big data analytics.
- Offers implicit and explicit parallelism.
- Methods have syntax *identical* to R.

---

<sup>a</sup>MPL, BSD, and GPL licensed



## pbdR Packages



## pbdR Motivation

Why HPC libraries (MPI, ScaLAPACK, PETSc, ...)?

- The HPC community has been at this for decades.
- *They're tested. They work. They're fast.*
- You're not going to beat Jack Dongarra at dense linear algebra.

## Simple Interface for MPI Operations with pbdMPI

### Rmpi

```
1 # int
2 mpi.allreduce(x, type=1)
3 # double
4 mpi.allreduce(x, type=2)
```

### pbdMPI

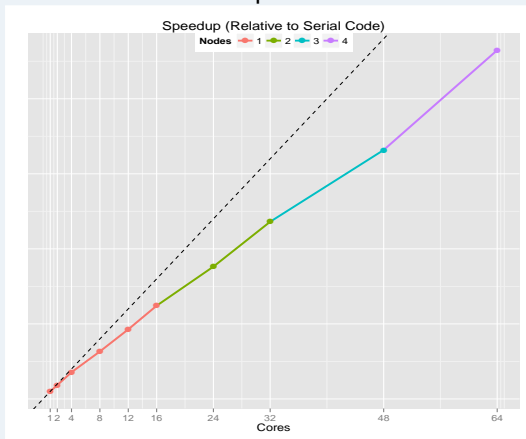
```
1 allreduce(x)
```

## Types in R

```
1 > is.integer(1)
2 [1] FALSE
3 > is.integer(2)
4 [1] FALSE
5 > is.integer(1:2)
6 [1] TRUE
```

Distributed Matrices and Statistics with **pbdDMAT**

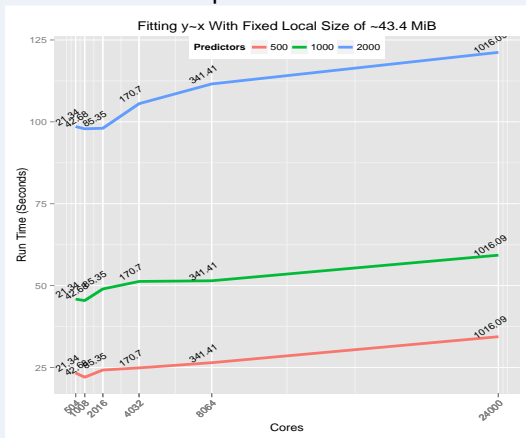
## Matrix Exponentiation



```
1 library(pbdDMAT)
2
3 dx <- ddmatrix("rnorm", 5000, 5000)
4 expm(dx)
```

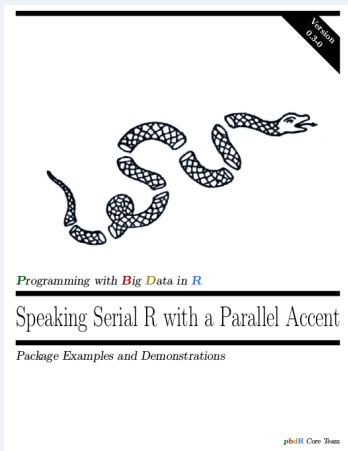
## Distributed Matrices and Statistics with pbdDMAT

## Least Squares Benchmark



```
x <- ddmatrix("rnorm", nrow=m, ncol=n)
y <- ddmatrix("rnorm", nrow=m, ncol=1)
mdl <- lm.fit(x=x, y=y)
```

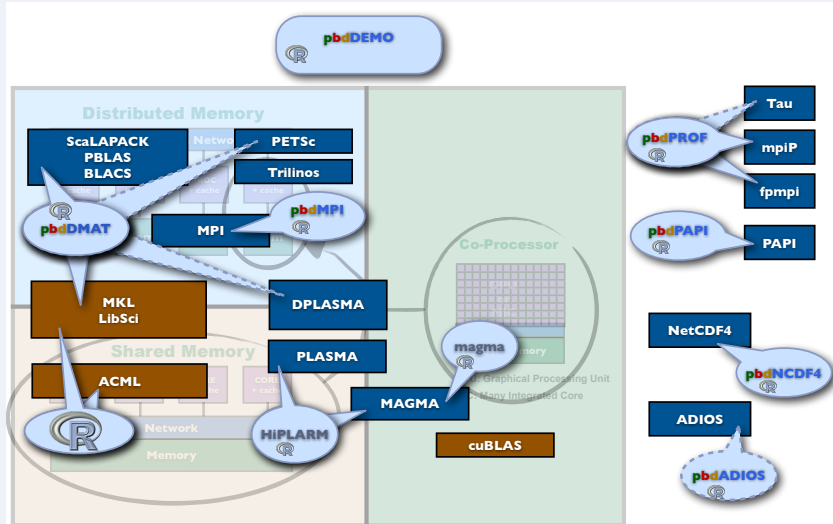
## Getting Started with HPC for R Users: **pbdDEMO**



- 140 page, textbook-style vignette.
- Over 30 demos, utilizing all\* packages.
- Not just a hello world!
- Demos include:
  - PCA
  - Regression
  - Parallel data input
  - Model-based clustering
  - Simple Monte Carlo simulation
  - Bayesian MCMC

- ### 3 The pbdR Project
- The pbdR Project
  - pbdR Connects R to HPC Libraries
  - Using pbdR
  - Summary

# R and pbdR Interfaces to HPC Libraries





- ### 3 The pbdR Project
- The pbdR Project
  - pbdR Connects R to HPC Libraries
  - Using pbdR
  - Summary

## pbdR Paradigms

**pbdR** programs are R programs!

Differences:

- Batch execution (non-interactive).
- Parallel code utilizes Single Program/Multiple Data (SPMD) style
- Emphasizes data parallelism.

## Batch Execution

- Running a serial R program in batch:

```
1 Rscript my_script.r
```

or

```
1 R CMD BATCH my_script.r
```

- Running a parallel (with MPI) R program in batch:

```
1 mpirun -np 2 Rscript my_par_script.r
```

## Single Program/Multiple Data (SPMD)

- SPMD is a programming *paradigm*.
- Not to be confused with SIMD.

### Paradigms

#### **Programming models**

OOP, Functional, SPMD, ...

### SIMD

#### **Hardware instructions**

MMX, SSE, ...

## Single Program/Multiple Data (SPMD)

SPMD is arguably the simplest extension of serial programming.

- Only one program is written, executed in batch on all processors.
- Different processors are autonomous; there is no manager.
- Dominant programming model for large machines for 30 years.

## Summary

- **pbdR** connects R to scalable HPC libraries.
- The **pbdDEMO** package offers numerous examples and explanations for getting started with distributed R programming.
- **pbdR** programs are R programs.

# Contents

## 4 Introduction to pbdMPI

- Managing a Communicator
- Reduce, Gather, Broadcast, and Barrier
- Other pbdMPI Tools
- Summary

## 4 Introduction to pbdMPI

- Managing a Communicator
- Reduce, Gather, Broadcast, and Barrier
- Other pbdMPI Tools
- Summary



## Message Passing Interface (MPI)

- *MPI*: Standard for managing communications (data and instructions) between different nodes/computers.
- *Implementations*: OpenMPI, MPICH2, Cray MPT, ...
- Enables parallelism (via communication) on distributed machines.
- *Communicator*: manages communications between processors.

## MPI Operations (1 of 2)

- **Managing a Communicator:** Create and destroy communicators.  
`init()` — initialize communicator  
`finalize()` — shut down communicator(s)
- **Rank query:** determine the processor's position in the communicator.  
`comm.rank()` — “who am I?”  
`comm.size()` — “how many of us are there?”
- **Printing:** Printing output from various ranks.  
`comm.print(x)`  
`comm.cat(x)`  
**WARNING:** only use these functions on *results*, never on yet-to-be-computed things.

## Quick Example 1

### Rank Query: 1\_rank.r

```
1 library(pbdMPI, quietly = TRUE)
2 init()
3
4 my.rank <- comm.rank()
5 comm.print(my.rank, all.rank=TRUE)
6
7 finalize()
```

Execute this script via:

```
1 mpirun -np 2 Rscript 1_rank.r
```

Sample Output:

```
1 COMM.RANK = 0
2 [1] 0
3 COMM.RANK = 1
4 [1] 1
```

## Quick Example 2

### Hello World: 2\_hello.r

```
1 library(pbdMPI, quietly=TRUE)
2 init()
3
4 comm.print("Hello, world")
5
6 comm.print("Hello again", all.rank=TRUE, quietly=TRUE)
7
8 finalize()
```

Execute this script via:

```
1 mpirun -np 2 Rscript 2_hello.r
```

Sample Output:

```
1 COMM.RANK = 0
2 [1] "Hello, world"
3 [1] "Hello again"
4 [1] "Hello again"
```

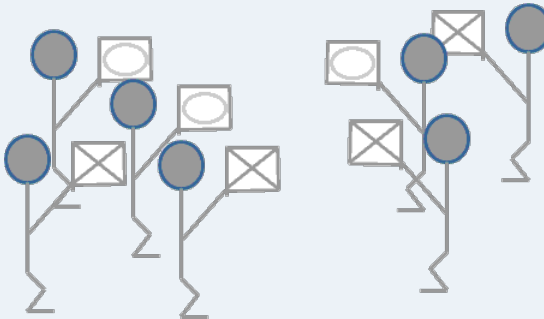
## 4 Introduction to pbdMPI

- Managing a Communicator
- Reduce, Gather, Broadcast, and Barrier
- Other pbdMPI Tools
- Summary

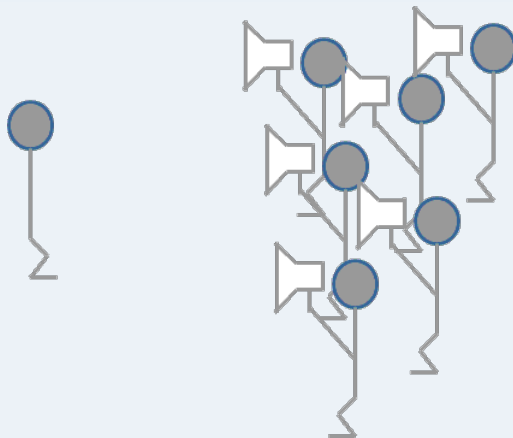
## MPI Operations

- 1 Reduce
- 2 Gather
- 3 Broadcast
- 4 Barrier

## Reductions — Combine results into single result

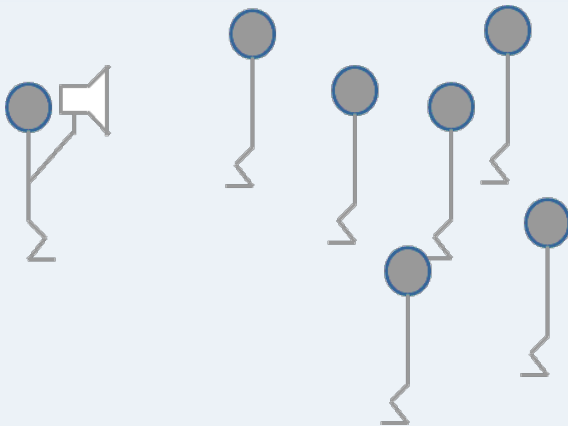


## Gather — Many-to-one

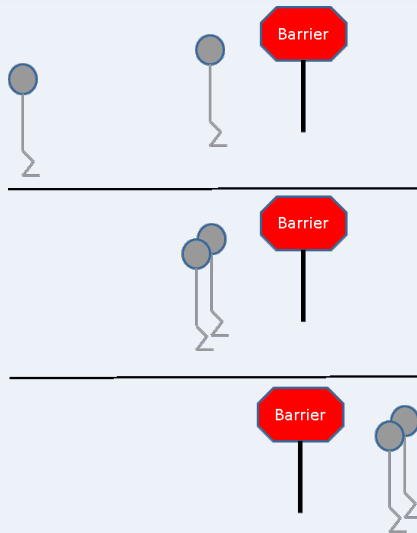




## Broadcast — One-to-many



## Barrier — Synchronization



## MPI Operations (2 of 2)

- **Reduction:** each processor has a number  $x$ ; add all of them up, find the largest/smallest, ....  
`reduce(x, op='sum')` — reduce to one  
`allreduce(x, op='sum')` — reduce to all
- **Gather:** each processor has a number; create a new object on some processor containing all of those numbers.  
`gather(x)` — gather to one  
`allgather(x)` — gather to all
- **Broadcast:** one processor has a number  $x$  that every other processor should also have.  
`bcast(x)`
- **Barrier:** “computation wall”; no processor can proceed until *all* processors can proceed.  
`barrier()`

## Quick Example 3

### Reduce and Gather: 3\_gt.r

```
1 library(pbdMPI, quietly=TRUE)
2 init()
3
4 comm.set.seed(diff=TRUE)
5
6 n <- sample(1:10, size=1)
7
8 gt <- gather(n)
9 comm.print(unlist(gt))
10
11 sm <- allreduce(n, op='sum')
12 comm.print(sm, all.rank=T)
13
14 finalize()
```

Execute this script via:

```
1 mpirun -np 2 Rscript 3_gt.r
```

Sample Output:

```
1 COMM.RANK = 0
2 [1] 2 8
3 COMM.RANK = 0
4 [1] 10
5 COMM.RANK = 1
6 [1] 10
```

## Quick Example 4

## Broadcast: 4\_bcast.r

```
1 library(pbdMPI, quietly=T)
2 init()
3
4 if (comm.rank()==0){
5   x <- matrix(1:4, nrow=2)
6 } else {
7   x <- NULL
8 }
9
10 y <- bcast(x, rank.source=0)
11
12 comm.print(y, rank=1)
13
14 finalize()
```

Execute this script via:

```
1 mpirun -np 2 Rscript 4_bcast.r
```

Sample Output:

```
1 COMM.RANK = 1
2   [,1] [,2]
3 [1,]   1   3
4 [2,]   2   4
```

## 4 Introduction to pbdMPI

- Managing a Communicator
- Reduce, Gather, Broadcast, and Barrier
- Other pbdMPI Tools
- Summary

## Random Seeds

**pbdMPI** offers a simple interface for managing random seeds:

- `comm.set.seed(seed=1234, diff=FALSE)` — All processors use the same seed.
- `comm.set.seed(seed=1234, diff=FALSE)` — All processors use the same seed.

## Other Helper Tools

**pbdMPI** Also contains useful tools for Manager/Worker and task parallelism codes:

- **Task Subsetting:** Distributing a list of jobs/tasks  
`get.jid(n)`
- **\*ply:** Functions in the \*ply family.  
`pbdApply(X, MARGIN, FUN, ...)` — analogue of `apply()`  
`pbdLapply(X, FUN, ...)` — analogue of `lapply()`  
`pbdSapply(X, FUN, ...)` — analogue of `sapply()`



## 4 Introduction to pbdMPI

- Managing a Communicator
- Reduce, Gather, Broadcast, and Barrier
- Other pbdMPI Tools
- Summary

## Summary

- Start by loading the package:

```
1 library(pbdMPI, quiet = TRUE)
```

- Always initialize before starting and finalize when finished:

```
1 init()  
2  
3 # ...  
4  
5 finalize()
```

# Contents

- 5 The Generalized Block Distribution
  - GBD: a Way to Distribute Your Data
  - Example GBD Distributions
  - Summary

## 5 The Generalized Block Distribution

- GBD: a Way to Distribute Your Data
- Example GBD Distributions
- Summary

## Distributing Data

**Problem:** How to distribute the data

$$X = \begin{bmatrix} x_{1,1} & x_{1,2} & x_{1,3} \\ x_{2,1} & x_{2,2} & x_{2,3} \\ x_{3,1} & x_{3,2} & x_{3,3} \\ x_{4,1} & x_{4,2} & x_{4,3} \\ x_{5,1} & x_{5,2} & x_{5,3} \\ x_{6,1} & x_{6,2} & x_{6,3} \\ x_{7,1} & x_{7,2} & x_{7,3} \\ x_{8,1} & x_{8,2} & x_{8,3} \\ x_{9,1} & x_{9,2} & x_{9,3} \\ x_{10,1} & x_{10,2} & x_{10,3} \end{bmatrix}_{10 \times 3}$$

?

## Distributing a Matrix Across 4 Processors: Block Distribution

	Data	Processors
$X =$	$x_{1,1}$ $x_{1,2}$ $x_{1,3}$	0
	$x_{2,1}$ $x_{2,2}$ $x_{2,3}$	1
	$x_{3,1}$ $x_{3,2}$ $x_{3,3}$	2
	$x_{4,1}$ $x_{4,2}$ $x_{4,3}$	3
	$x_{5,1}$ $x_{5,2}$ $x_{5,3}$	
	$x_{6,1}$ $x_{6,2}$ $x_{6,3}$	
	$x_{7,1}$ $x_{7,2}$ $x_{7,3}$	
	$x_{8,1}$ $x_{8,2}$ $x_{8,3}$	
	$x_{9,1}$ $x_{9,2}$ $x_{9,3}$	
	$x_{10,1}$ $x_{10,2}$ $x_{10,3}$	

$10 \times 3$

## Distributing a Matrix Across 4 Processors: Local Load Balance

	Data			Processors
$X =$	$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	0
	$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	1
	$x_{3,1}$	$x_{3,2}$	$x_{3,3}$	2
	$x_{4,1}$	$x_{4,2}$	$x_{4,3}$	3
	$x_{5,1}$	$x_{5,2}$	$x_{5,3}$	
	$x_{6,1}$	$x_{6,2}$	$x_{6,3}$	
	$x_{7,1}$	$x_{7,2}$	$x_{7,3}$	
	$x_{8,1}$	$x_{8,2}$	$x_{8,3}$	
	$x_{9,1}$	$x_{9,2}$	$x_{9,3}$	
	$x_{10,1}$	$x_{10,2}$	$x_{10,3}$	

 $10 \times 3$

## The GBD Data Structure

Throughout the examples, we will make use of the Generalized Block Distribution, or GBD distributed matrix structure.

- 1 GBD is *distributed*. No processor owns all the data.
- 2 GBD is *non-overlapping*. Rows uniquely assigned to processors.
- 3 GBD is *row-contiguous*. If a processor owns one element of a row, it owns the entire row.
- 4 GBD is globally *row-major*, locally *column-major*.
- 5 GBD is often *locally balanced*, where each processor owns (almost) the same amount of data. But this is not required.
- 6 The last row of the local storage of a processor is adjacent (by global row) to the first row of the local storage of next processor (by communicator number) that owns data.
- 7 GBD is (relatively) easy to understand, but can lead to bottlenecks if you have many more columns than rows.

$x_{1,1}$	$x_{1,2}$	$x_{1,3}$
$x_{2,1}$	$x_{2,2}$	$x_{2,3}$
$x_{3,1}$	$x_{3,2}$	$x_{3,3}$
$x_{4,1}$	$x_{4,2}$	$x_{4,3}$
$x_{5,1}$	$x_{5,2}$	$x_{5,3}$
$x_{6,1}$	$x_{6,2}$	$x_{6,3}$
$x_{7,1}$	$x_{7,2}$	$x_{7,3}$
$x_{8,1}$	$x_{8,2}$	$x_{8,3}$
$x_{9,1}$	$x_{9,2}$	$x_{9,3}$
$x_{10,1}$	$x_{10,2}$	$x_{10,3}$



## 5 The Generalized Block Distribution

- GBD: a Way to Distribute Your Data
- Example GBD Distributions
- Summary

## Understanding GBD: Global Matrix

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

Processors = 0 1 2 3 4 5

## Understanding GBD: Load Balanced GBD

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

Processors = 0 1 2 3 4 5

## Understanding GBD: Local View

$$\begin{bmatrix}
 x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\
 x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29}
 \end{bmatrix}_{2 \times 9}$$

$$\begin{bmatrix}
 x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\
 x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49}
 \end{bmatrix}_{2 \times 9}$$

$$\begin{bmatrix}
 x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\
 x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69}
 \end{bmatrix}_{2 \times 9}$$

$$\begin{bmatrix}
 x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79}
 \end{bmatrix}_{1 \times 9}$$

$$\begin{bmatrix}
 x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89}
 \end{bmatrix}_{1 \times 9}$$

$$\begin{bmatrix}
 x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99}
 \end{bmatrix}_{1 \times 9}$$

Processors = 0 1 2 3 4 5

## Understanding GBD: Non-Balanced GBD

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

Processors = 0 1 2 3 4 5

# Understanding GBD: Local View

[										]	0×9
[	X <sub>11</sub>	X <sub>12</sub>	X <sub>13</sub>	X <sub>14</sub>	X <sub>15</sub>	X <sub>16</sub>	X <sub>17</sub>	X <sub>18</sub>	X <sub>19</sub>	]	4×9
	X <sub>21</sub>	X <sub>22</sub>	X <sub>23</sub>	X <sub>24</sub>	X <sub>25</sub>	X <sub>26</sub>	X <sub>27</sub>	X <sub>28</sub>	X <sub>29</sub>		
	X <sub>31</sub>	X <sub>32</sub>	X <sub>33</sub>	X <sub>34</sub>	X <sub>35</sub>	X <sub>36</sub>	X <sub>37</sub>	X <sub>38</sub>	X <sub>39</sub>		
	X <sub>41</sub>	X <sub>42</sub>	X <sub>43</sub>	X <sub>44</sub>	X <sub>45</sub>	X <sub>46</sub>	X <sub>47</sub>	X <sub>48</sub>	X <sub>49</sub>		
[	X <sub>51</sub>	X <sub>52</sub>	X <sub>53</sub>	X <sub>54</sub>	X <sub>55</sub>	X <sub>56</sub>	X <sub>57</sub>	X <sub>58</sub>	X <sub>59</sub>	]	2×9
	X <sub>61</sub>	X <sub>62</sub>	X <sub>63</sub>	X <sub>64</sub>	X <sub>65</sub>	X <sub>66</sub>	X <sub>67</sub>	X <sub>68</sub>	X <sub>69</sub>		
[	X <sub>71</sub>	X <sub>72</sub>	X <sub>73</sub>	X <sub>74</sub>	X <sub>75</sub>	X <sub>76</sub>	X <sub>77</sub>	X <sub>78</sub>	X <sub>79</sub>	]	1×9
[										]	0×9
[	X <sub>81</sub>	X <sub>82</sub>	X <sub>83</sub>	X <sub>84</sub>	X <sub>85</sub>	X <sub>86</sub>	X <sub>87</sub>	X <sub>88</sub>	X <sub>89</sub>	]	2×9
	X <sub>91</sub>	X <sub>92</sub>	X <sub>93</sub>	X <sub>94</sub>	X <sub>95</sub>	X <sub>96</sub>	X <sub>97</sub>	X <sub>98</sub>	X <sub>99</sub>		

Processors = 0 1 2 3 4 5

- 5 The Generalized Block Distribution
  - GBD: a Way to Distribute Your Data
  - Example GBD Distributions
  - Summary

## Summary

- Need to distribute your data? Try splitting by row.
- May not work well if your data is square (or longer than tall).



# Contents

## 6 Basic Statistics Examples

- pbdMPI Example: Monte Carlo Simulation
- pbdMPI Example: Sample Covariance
- pbdMPI Example: Linear Regression
- Summary

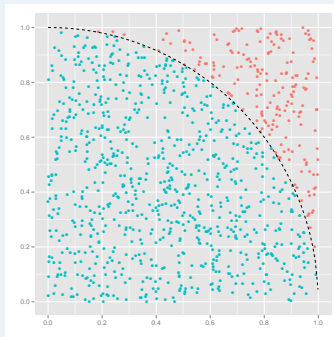
## 6 Basic Statistics Examples

- pbdMPI Example: Monte Carlo Simulation
- pbdMPI Example: Sample Covariance
- pbdMPI Example: Linear Regression
- Summary

## Example 1: Monte Carlo Simulation

Sample  $N$  uniform observations  $(x_i, y_i)$  in the unit square  $[0, 1] \times [0, 1]$ .  
Then

$$\pi \approx 4 \left( \frac{\# \text{ Inside Circle}}{\# \text{ Total}} \right) = 4 \left( \frac{\# \text{ Blue}}{\# \text{ Blue} + \# \text{ Red}} \right)$$



## Example 1: Monte Carlo Simulation GBD Algorithm

- 1 Let  $n$  be big-ish; we'll take  $n = 50,000$ .
- 2 Generate an  $n \times 2$  matrix  $x$  of standard uniform observations.
- 3 Count the number of rows satisfying  $x^2 + y^2 \leq 1$
- 4 Ask everyone else what their answer is; sum it all up.
- 5 Take this new answer, multiply by 4 and divide by  $n$
- 6 If my rank is 0, print the result.

## Example 1: Monte Carlo Simulation Code

### Serial Code

```
1 N <- 50000
2 X <- matrix(runif(N * 2), ncol=2)
3 r <- sum(rowSums(X^2) <= 1)
4 PI <- 4*r/N
5 print(PI)
```

### Parallel Code

```
1 library(pbdMPI, quiet = TRUE)
2 init()
3 comm.set.seed(seed=1234567, diff=TRUE)
4
5 N.gbd <- 50000 / comm.size()
6 X.gbd <- matrix(runif(N.gbd * 2), ncol = 2)
7 r.gbd <- sum(rowSums(X.gbd^2) <= 1)
8 r <- allreduce(r.gbd)
9 PI <- 4*r/(N.gbd * comm.size())
10 comm.print(PI)
11
12 finalize()
```

## Note

For the remainder, we will exclude loading, init, and finalize calls.

## 6 Basic Statistics Examples

- pbdMPI Example: Monte Carlo Simulation
- **pbdMPI Example: Sample Covariance**
- pbdMPI Example: Linear Regression
- Summary

## Example 2: Sample Covariance

$$\text{cov}(x_{n \times p}) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu_x)(x_i - \mu_x)^T$$



## Example 2: Sample Covariance GBD Algorithm

- 1 Determine the total number of rows  $N$ .
- 2 Compute the vector of column means of the full matrix.
- 3 Subtract each column's mean from that column's entries in each local matrix.
- 4 Compute the crossproduct locally and reduce.
- 5 Divide by  $N - 1$ .

## Example 2: Sample Covariance Code

### Serial Code

```
1 N <- nrow(X)
2 mu <- colSums(X) / N
3
4 X <- sweep(X, STATS=mu, MARGIN=2)
5 Cov.X <- crossprod(X) / (N-1)
6
7 print(Cov.X)
```

### Parallel Code

```
1 N <- allreduce(nrow(X.gbd), op="sum")
2 mu <- allreduce(colSums(X.gbd) / N, op="sum")
3
4 X.gbd <- sweep(X.gbd, STATS=mu, MARGIN=2)
5 Cov.X <- allreduce(crossprod(X.gbd), op="sum") / (N-1)
6
7 comm.print(Cov.X)
```

## 6 Basic Statistics Examples

- pbdMPI Example: Monte Carlo Simulation
- pbdMPI Example: Sample Covariance
- pbdMPI Example: Linear Regression
- Summary

### Example 3: Linear Regression

Find  $\beta$  such that

$$y = X\beta + \epsilon$$

When  $X$  is full rank,

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

### Example 3: Linear Regression GBD Algorithm

- 1 Locally, compute  $tx = x^T$
- 2 Locally, compute  $A = tx * x$ . Query every other processor for this result and sum up all the results.
- 3 Locally, compute  $B = tx * y$ . Query every other processor for this result and sum up all the results.
- 4 Locally, compute  $A^{-1} * B$

## Example 3: Linear Regression Code

### Serial Code

```
1 tX <- t(X)
2 A <- tX %*% X
3 B <- tX %*% y
4
5 ols <- solve(A) %*% B
```

### Parallel Code

```
1 tX.gbd <- t(X.gbd)
2 A <- allreduce(tX.gbd %*% X.gbd, op = "sum")
3 B <- allreduce(tX.gbd %*% y.gbd, op = "sum")
4
5 ols <- solve(A) %*% B
```

## 6 Basic Statistics Examples

- pbdMPI Example: Monte Carlo Simulation
- pbdMPI Example: Sample Covariance
- pbdMPI Example: Linear Regression
- Summary

## Summary

- SPMD programming is (often) a natural extension of serial programming.
- More **pbdMPI** examples in **pbdDEMO**.



# Contents

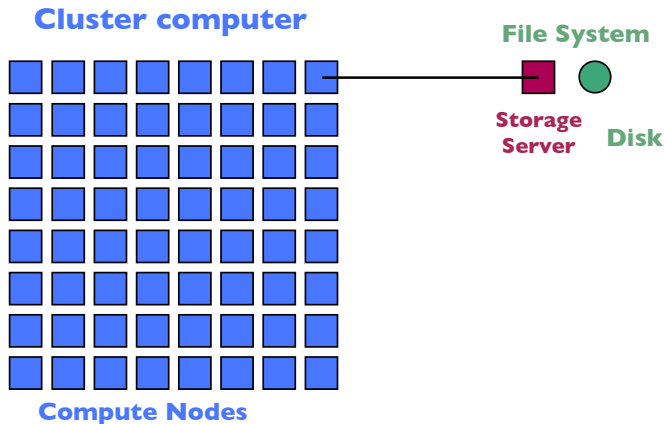
## 7 Data Input

- Cluster Computer and File System
- Serial Data Input
- Parallel Data Input
- Summary

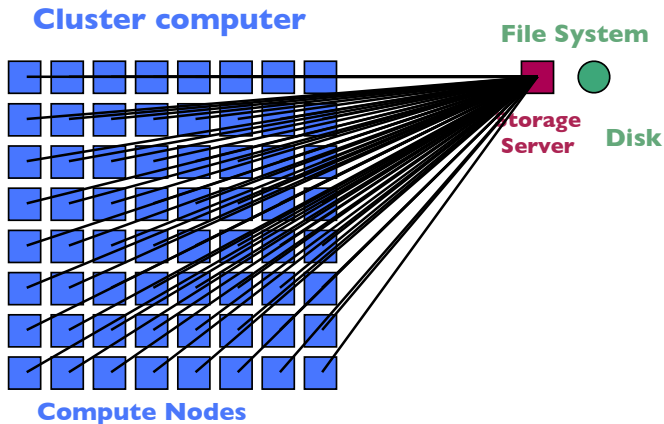
## 7 Data Input

- Cluster Computer and File System
- Serial Data Input
- Parallel Data Input
- Summary

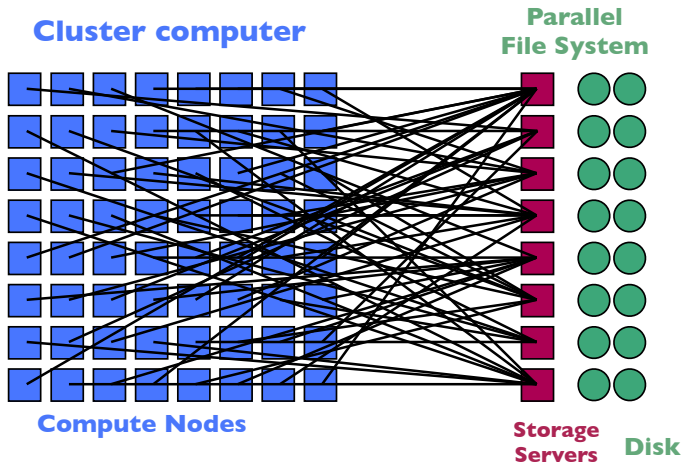
# One Node to One Storage Server



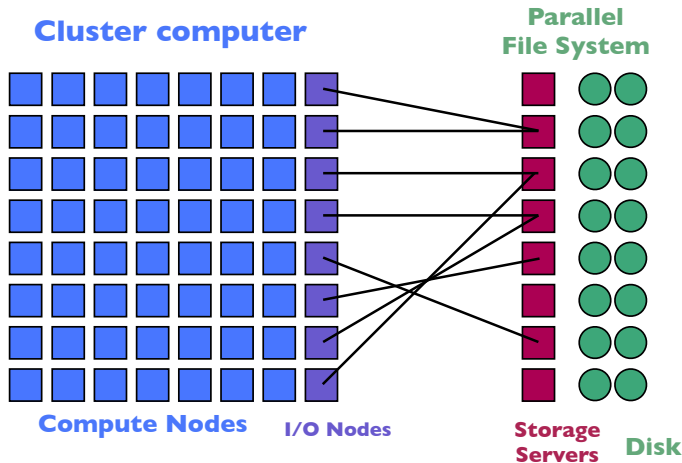
# Many Nodes to One Storage Server



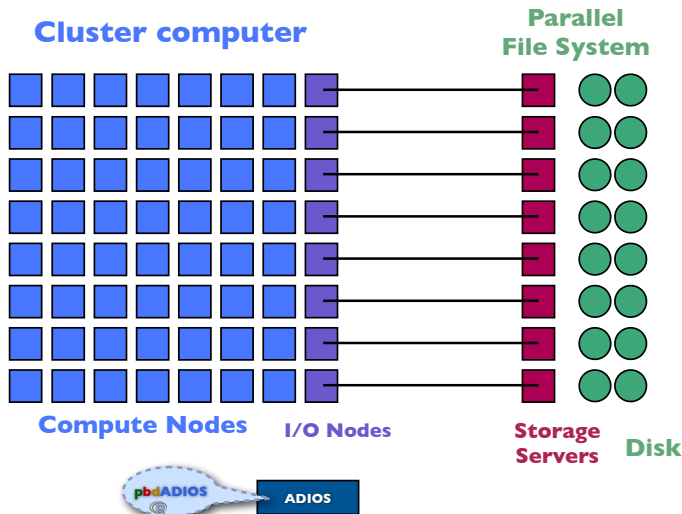
# Many Nodes to Few Storage Servers



# Few Nodes to Few Storage Servers - Default



# Few Nodes to Few Storage Servers - Coordinated



ADIOS

## 7 Data Input

- Cluster Computer and File System
- **Serial Data Input**
- Parallel Data Input
- Summary



Separate manual: <http://r-project.org/>

- `scan()`
- `read.table()`
- `read.csv()`
- `socket`

## CSV Data: Read Serial then Distribute

### Listing:

```
1 library(pbdDMAT)
2 if(comm.rank() == 0) { # only read on process 0
3   x <- read.csv("myfile.csv")
4 } else {
5   x <- NULL
6 }
7
8 dx <- as.ddmatrix(x)
```

## 7 Data Input

- Cluster Computer and File System
- Serial Data Input
- **Parallel Data Input**
- Summary

## New Issues

- How to read in parallel?
- CSV, SQL, NetCDF4, HDF, ADIOS, custom binary
- How to partition data across nodes?
- How to structure for scalable libraries?
- Read directly into form needed or restructure?
- ...
- A lot of work needed here!

## CSV Data

### Serial Code

```
1 x <- read.csv("x.csv")
2
3 x
```

### Parallel Code

```
1 library(pbdDEMO, quiet = TRUE)
2 init.grid()
3
4 dx <- read.csv.ddmatrix("x.csv", header=TRUE, sep=",",
5                          nrows=10, ncols=10, num.rdrs=2, ICTXT=0)
6
7 dx
8
9 finalize()
```

## Binary Data: Vector

```
1 ## set up start and length for reading a vector of n doubles
2 size <- 8 # bytes
3
4 my_ids <- get.jid(n, method="block")
5
6 my_start <- (my_ids[1] - 1)*size
7 my_length <- length(my_ids)
8
9 con <- file("binary.vector.file", "rb")
10 seekval <- seek(con, where=my_start, rw="read")
11 x <- readBin(con, what="double", n=my_length, size=size)
```

## Binary Data: Matrix

```
1 ## read an nrow by ncol matrix of doubles split by columns
2 size <- 8 # bytes
3
4 my_ids <- get.jid(ncol, method="block")
5 my_ncol <- length(my_ids)
6 my_start <- (my_ids[1] - 1)*nrow*size
7 my_length <- my_ncol*nrow
8
9 con <- file("binary.matrix.file", "rb")
10 seekval <- seek(con, where=my_start, rw="read")
11 x <- readBin(con, what="double", n=my_length, size=size)
12
13 ## glue together as a column-block ddmatrix
14 gdim <- c(nrow, ncol)
15 ldim <- c(nrow, my_ncol)
16 bldim <- c(nrow, allreduce(my_ncol, op="max"))
17 X <- new("ddmatrix", Data=matrix(x, nrow, my_ncol),
18       dim=gdim, ldim=ldim, bldim=bldim, ICTXT=1)
19
20 ## redistribute for ScaLAPACK's block-cyclic
21 X <- redistribute(X, bldim=c(2, 2), ICTXT=0)
22 Xprc <- prcomp(X)
```

## NetCDF4 Data

```
1 ### parallel read after determining start and length
2 nc <- nc_open_par(file_name)
3
4 nc_var_par_access(nc, "variable_name")
5 new.X <- ncvar_get(nc, "variable_name", start, length)
6 nc_close(nc)
7
8 finalize()
```



## 7 Data Input

- Cluster Computer and File System
- Serial Data Input
- Parallel Data Input
- Summary

## Summary

- Mostly “do it yourself”
- Parallel file system for big data
  - Binary files for true parallel reads
  - Know number of readers vs number of storage servers
- Redistribution help from `ddmatrix` functions
- More help under development

# Contents

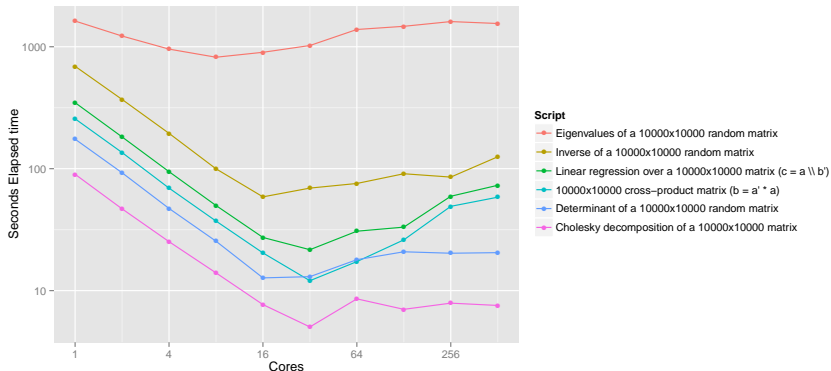
- 8 Introduction to pbdDMAT and the ddmatrix Structure
  - Introduction to Distributed Matrices
  - pbdDMAT
  - Summary

## 8 Introduction to pbdDMAT and the ddmatrix Structure

- Introduction to Distributed Matrices
  - pbdDMAT
  - Summary

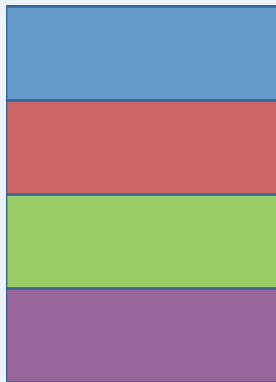
## Distributed Matrices

You can only get so far with one node. . .

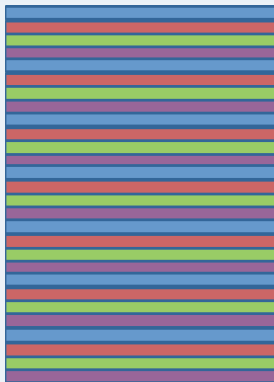


The solution is to distribute the data.

## Distributed Matrices



(a) Block



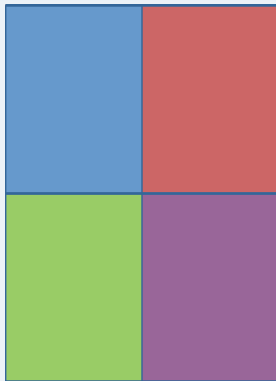
(b) Cyclic



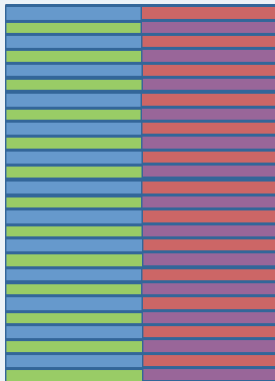
(c) Block-Cyclic

Figure: Matrix Distribution Schemes

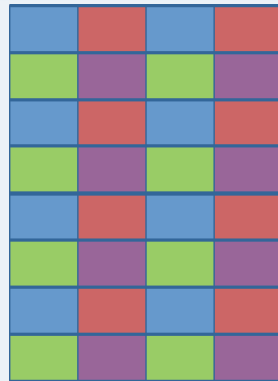
## Distributed Matrices



(a) 2d Block



(b) 2d Cyclic



(c) 2d Block-Cyclic

Figure: Matrix Distribution Schemes Onto a 2-Dimensional Grid

## Processor Grid Shapes

$$\begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}^T$$

(a)  $1 \times 6$

$$\begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{bmatrix}$$

(b)  $2 \times 3$

$$\begin{bmatrix} 0 & 1 \\ 2 & 3 \\ 4 & 5 \end{bmatrix}$$

(c)  $3 \times 2$

$$\begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}$$

(d)  $6 \times 1$

Table: Processor Grid Shapes with 6 Processors



## The ddmatrix Class

For **d**istributed **d**ense **m**atrix objects, we use the special S4 class `ddmatrix`.

<code>ddmatrix</code> =	{	<b>Data</b>	The local submatrix (an R matrix)
		<b>dim</b>	Dimension of the global matrix
		<b>ldim</b>	Dimension of the local submatrix
		<b>bldim</b>	Dimension of the blocks
		<b>ICTXT</b>	MPI Grid Context

## Understanding ddmatrix: Global Matrix

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

## ddmatrix: 1-dimensional Row Block

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

$$\text{Processor grid} = \begin{vmatrix} 0 \\ 1 \\ 2 \\ 3 \end{vmatrix} = \begin{vmatrix} (0,0) \\ (1,0) \\ (2,0) \\ (3,0) \end{vmatrix}$$

## ddmatrix: 2-dimensional Row Block

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

$$\text{Processor grid} = \left| \begin{array}{cc} 0 & 1 \\ 2 & 3 \end{array} \right| = \left| \begin{array}{cc} (0,0) & (0,1) \\ (1,0) & (1,1) \end{array} \right|$$

## ddmatrix: 1-dimensional Row Cyclic

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

$$\text{Processor grid} = \begin{vmatrix} 0 \\ 1 \\ 2 \\ 3 \end{vmatrix} = \begin{vmatrix} (0,0) \\ (1,0) \\ (2,0) \\ (3,0) \end{vmatrix}$$

## ddmatrix: 2-dimensional Row Cyclic

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

$$\text{Processor grid} = \begin{vmatrix} 0 & 1 \\ 2 & 3 \end{vmatrix} = \begin{vmatrix} (0,0) & (0,1) \\ (1,0) & (1,1) \end{vmatrix}$$

## ddmatrix: 2-dimensional Block-Cyclic

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

$$\text{Processor grid} = \begin{vmatrix} 0 & 1 \\ 2 & 3 \end{vmatrix} = \begin{vmatrix} (0,0) & (0,1) \\ (1,0) & (1,1) \end{vmatrix}$$

## 8 Introduction to pbdDMAT and the ddmatrix Structure

- Introduction to Distributed Matrices
- pbdDMAT
- Summary



## The ddmatrix Data Structure

The more complicated the processor grid, the more complicated the distribution.

## ddmatrix: 2-dimensional Block-Cyclic with 6 Processors

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

$$\text{Processor grid} = \begin{vmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{vmatrix} = \begin{vmatrix} (0,0) & (0,1) & (0,2) \\ (1,0) & (1,1) & (1,2) \end{vmatrix}$$

## Understanding ddmatrix: Local View

$\begin{bmatrix} X_{11} & X_{12} & X_{17} & X_{18} \\ X_{21} & X_{22} & X_{27} & X_{28} \\ X_{51} & X_{52} & X_{57} & X_{58} \\ X_{61} & X_{62} & X_{67} & X_{68} \\ X_{91} & X_{92} & X_{97} & X_{98} \end{bmatrix}_{5 \times 4}$	$\begin{bmatrix} X_{13} & X_{14} & X_{19} \\ X_{23} & X_{24} & X_{29} \\ X_{53} & X_{54} & X_{59} \\ X_{63} & X_{64} & X_{69} \\ X_{93} & X_{94} & X_{99} \end{bmatrix}_{5 \times 3}$	$\begin{bmatrix} X_{15} & X_{16} \\ X_{25} & X_{26} \\ X_{55} & X_{56} \\ X_{65} & X_{66} \\ X_{95} & X_{96} \end{bmatrix}_{5 \times 2}$
$\begin{bmatrix} X_{31} & X_{32} & X_{37} & X_{38} \\ X_{41} & X_{42} & X_{47} & X_{48} \\ X_{71} & X_{72} & X_{77} & X_{78} \\ X_{81} & X_{82} & X_{87} & X_{88} \end{bmatrix}_{4 \times 4}$	$\begin{bmatrix} X_{33} & X_{34} & X_{39} \\ X_{43} & X_{44} & X_{49} \\ X_{73} & X_{74} & X_{79} \\ X_{83} & X_{84} & X_{89} \end{bmatrix}_{4 \times 3}$	$\begin{bmatrix} X_{35} & X_{36} \\ X_{45} & X_{46} \\ X_{75} & X_{76} \\ X_{85} & X_{86} \end{bmatrix}_{4 \times 2}$

$$\text{Processor grid} = \begin{vmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{vmatrix} = \begin{vmatrix} (0,0) & (0,1) & (0,2) \\ (1,0) & (1,1) & (1,2) \end{vmatrix}$$

## The ddmatrix Data Structure

- ① ddmatrix is *distributed*. No one processor owns all of the matrix.
- ② ddmatrix is *non-overlapping*. Any piece owned by one processor is owned by no other processors.
- ③ ddmatrix can be row-contiguous or not, depending on the processor grid and blocking factor used.
- ④ ddmatrix is locally column-major and globally, it depends...
- ⑥ GBD is a generalization of the one-dimensional block ddmatrix distribution. Otherwise there is no relation.
- ⑦ ddmatrix is confusing, but very robust.

X <sub>11</sub>	X <sub>12</sub>	X <sub>13</sub>	X <sub>14</sub>	X <sub>15</sub>
X <sub>21</sub>	X <sub>22</sub>	X <sub>23</sub>	X <sub>24</sub>	X <sub>25</sub>
X <sub>31</sub>	X <sub>32</sub>	X <sub>33</sub>	X <sub>34</sub>	X <sub>35</sub>
X <sub>41</sub>	X <sub>42</sub>	X <sub>43</sub>	X <sub>44</sub>	X <sub>45</sub>
X <sub>51</sub>	X <sub>52</sub>	X <sub>53</sub>	X <sub>54</sub>	X <sub>55</sub>
X <sub>61</sub>	X <sub>62</sub>	X <sub>63</sub>	X <sub>64</sub>	X <sub>65</sub>
X <sub>71</sub>	X <sub>72</sub>	X <sub>73</sub>	X <sub>74</sub>	X <sub>75</sub>
X <sub>81</sub>	X <sub>82</sub>	X <sub>83</sub>	X <sub>84</sub>	X <sub>85</sub>
X <sub>91</sub>	X <sub>92</sub>	X <sub>93</sub>	X <sub>94</sub>	X <sub>95</sub>

## Pros and Cons of This Data Structure

### Pros

- Robust for matrix computations.

### Cons

- Confusing layout.

*This is why we hide most of the distributed details.*

The details are there if you want them (you don't want them).

## Methods for class ddmatrix

**pbdDMAT** has over 100 methods with *identical* syntax to R:

- ``[, rbind(), cbind(), ...`
- `lm.fit(), prcomp(), cov(), ...`
- ``%*%`, solve(), svd(), norm(), ...`
- `median(), mean(), rowSums(), ...`

### Serial Code

```
1 cov(x)
```

### Parallel Code

```
1 cov(x)
```

## Comparing pbdMPI and pbdDMAT

### pbdMPI:

- MPI + sugar.
- GBD not the only structure **pbdMPI** can handle (just a useful convention).

### pbdDMAT:

- Distributed matrices + statistics.
- The `ddmatrix` structure *must* be used for **pbdDMAT**.
- If the data is not 2d block-cyclic compatible, `ddmatrix` will *definitely* give the wrong answer.

- 8 Introduction to pbdDMAT and the ddmatrix Structure
  - Introduction to Distributed Matrices
  - pbdDMAT
  - Summary



## Summary

- 1 Start by loading the package:

```
1 library(pbdDMAT, quiet = TRUE)
```

- 2 Always initialize before starting and finalize when finished:

```
1 init.grid()  
2  
3 # ...  
4  
5 finalize()
```

# Contents

## 9 Examples Using pbdDMAT

- RandSVD
- Summary

## 9 Examples Using pbdDMAT

- RandSVD
- Summary

Randomized SVD<sup>1</sup>

## PROTOTYPE FOR RANDOMIZED SVD

Given an  $m \times n$  matrix  $A$ , a target number  $k$  of singular vectors, and an exponent  $q$  (say,  $q = 1$  or  $q = 2$ ), this procedure computes an approximate rank- $2k$  factorization  $U\Sigma V^*$ , where  $U$  and  $V$  are orthonormal, and  $\Sigma$  is nonnegative and diagonal.

**Stage A:**

- 1 Generate an  $n \times 2k$  Gaussian test matrix  $\Omega$ .
- 2 Form  $Y = (AA^*)^q A\Omega$  by multiplying alternately with  $A$  and  $A^*$ .
- 3 Construct a matrix  $Q$  whose columns form an orthonormal basis for the range of  $Y$ .

**Stage B:**

- 4 Form  $B = Q^* A$ .
- 5 Compute an SVD of the small matrix:  $B = \tilde{U}\Sigma V^*$ .
- 6 Set  $U = Q\tilde{U}$ .

**Note:** The computation of  $Y$  in step 2 is vulnerable to round-off errors. When high accuracy is required, we must incorporate an orthonormalization step between each application of  $A$  and  $A^*$ ; see Algorithm 4.4.

## ALGORITHM 4.4: RANDOMIZED SUBSPACE ITERATION

Given an  $m \times n$  matrix  $A$  and integers  $\ell$  and  $q$ , this algorithm computes an  $m \times \ell$  orthonormal matrix  $Q$  whose range approximates the range of  $A$ .

- 1 Draw an  $n \times \ell$  standard Gaussian matrix  $\Omega$ .
- 2 Form  $Y_0 = A\Omega$  and compute its QR factorization  $Y_0 = Q_0 R_0$ .
- 3 **for**  $j = 1, 2, \dots, q$
- 4     Form  $\tilde{Y}_j = A^* Q_{j-1}$  and compute its QR factorization  $\tilde{Y}_j = \tilde{Q}_j \tilde{R}_j$ .
- 5     Form  $Y_j = A \tilde{Q}_j$  and compute its QR factorization  $Y_j = Q_j R_j$ .
- 6 **end**
- 7  $Q = Q_q$ .

## Serial R

```

1 randSVD <- function(A, k, q=3)
2 {
3   ## Stage A
4   Omega <- matrix(rnorm(n*2*k),
5                 nrow=n, ncol=2*k)
6   Y <- A %*% Omega
7   Q <- qr.Q(qr(Y))
8   At <- t(A)
9   for(i in 1:q)
10    {
11      Y <- At %*% Q
12      Q <- qr.Q(qr(Y))
13      Y <- A %*% Q
14      Q <- qr.Q(qr(Y))
15    }
16
17   ## Stage B
18   B <- t(Q) %*% A
19   U <- La.svd(B)$u
20   U <- Q %*% U
21   U[, 1:k]
22 }
```

<sup>1</sup>Halko N, Martinsson P-G and Tropp J A 2011 Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions *SIAM Rev.* 53 217–88

## Randomized SVD

## Serial R

```

1 randSVD <- function(A, k, q=3)
2 {
3   ## Stage A
4   Omega <- matrix(rnorm(n*2*k),
5     nrow=n, ncol=2*k)
6   Y <- A %*% Omega
7   Q <- qr.Q(qr(Y))
8   At <- t(A)
9   for(i in 1:q)
10    {
11      Y <- At %*% Q
12      Q <- qr.Q(qr(Y))
13      Y <- A %*% Q
14      Q <- qr.Q(qr(Y))
15    }
16
17   ## Stage B
18   B <- t(Q) %*% A
19   U <- La.svd(B)$u
20   U <- Q %*% U
21   U[, 1:k]
22 }

```

## Parallel pbdR

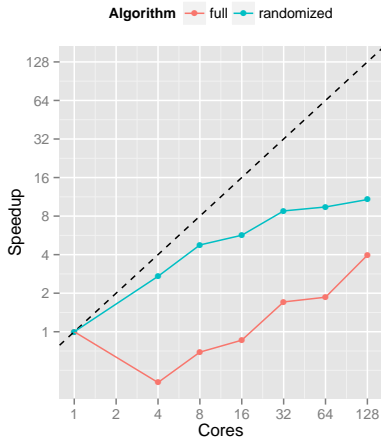
```

1 randSVD <- function(A, k, q=3)
2 {
3   ## Stage A
4   Omega <- ddmatrix("rnorm",
5     nrow=n, ncol=2*k)
6   Y <- A %*% Omega
7   Q <- qr.Q(qr(Y))
8   At <- t(A)
9   for(i in 1:q)
10    {
11      Y <- At %*% Q
12      Q <- qr.Q(qr(Y))
13      Y <- A %*% Q
14      Q <- qr.Q(qr(Y))
15    }
16
17   ## Stage B
18   B <- t(Q) %*% A
19   U <- La.svd(B)$u
20   U <- Q %*% U
21   U[, 1:k]
22 }

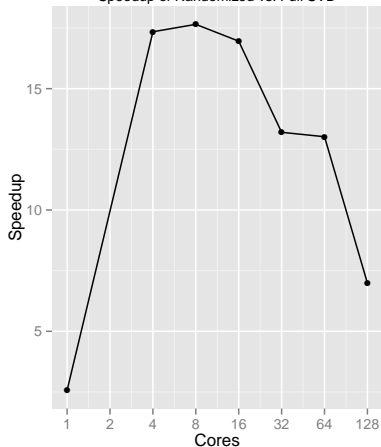
```

# Randomized SVD

30 Singular Vectors from a 100,000 by 1,000 Matrix



30 Singular Vectors from a 100,000 by 1,000 Matrix  
Speedup of Randomized vs. Full SVD



## 9 Examples Using pbdDMAT

- RandSVD
- Summary

## Summary

- **pbdDMAT** makes distributed (dense) linear algebra easier.
- Can enable rapid prototyping at large scale.



# Contents

## 10 MPI Profiling

- Profiling with the pbdPROF Package
- Installing pbdPROF
- Example
- Summary

## 10 MPI Profiling

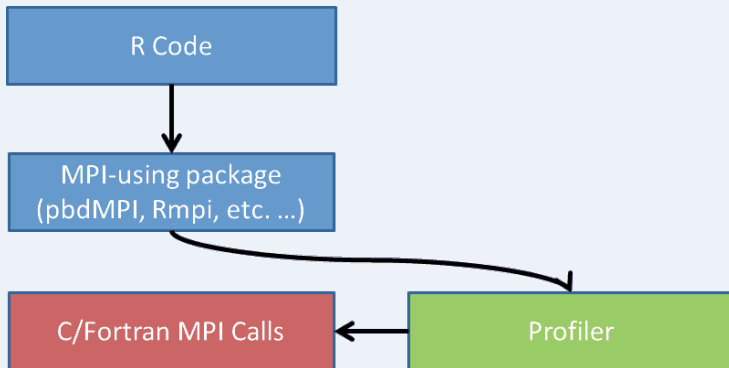
- Profiling with the pbdPROF Package
- Installing pbdPROF
- Example
- Summary

## Introduction to **pbdPROF**

- Successful Google Summer of Code 2013 project.
- Available on the CRAN.
- Enables profiling of MPI-using R scripts.
- **pbdR** packages officially supported; can work with others. . .
- Also reads, parses, and plots profiler outputs.

## How it works

MPI calls get hijacked by profiler and logged:



## Introduction to **pbdPROF**

- Currently supports the profilers **fpmpi** and **mpiP**.
- **fpmpi** is distributed with **pbdPROF** and installs easily, but offers minimal profiling capabilities.
- **mpiP** is fully supported also, but you have to install and link it yourself.

## 10 MPI Profiling

- Profiling with the pbdPROF Package
- Installing pbdPROF
- Example
- Summary

## Installing **pbdPROF**

- 1 Build **pbdPROF**.
- 2 Rebuild **pbdMPI** (linking with **pbdPROF**).
- 3 Run your analysis as usual.
- 4 Interactively analyze profiler outputs with **pbdPROF**.

This is explained at length in the **pbdPROF** vignette.

## Rebuild pbdMPI

```
R CMD INSTALL pbdMPI_0.2-2.tar.gz  
  --configure-args="--enable-pbdPROF"
```

- Any package which explicitly links with an MPI library must be rebuilt in this way (**pbdMPI**, **Rmpi**, ...).
- Other **pbdR** packages link with **pbdMPI**, and so do not need to be rebuilt.
- See **pbdPROF** vignette if something goes wrong.



## 10 MPI Profiling

- Profiling with the pbdPROF Package
- Installing pbdPROF
- Example
- Summary

## An Example from **pbDMMAT**

- Compute SVD in **pbDMMAT** package.
- Profile MPI calls with **mpiP**.

## Example Script

my\_svd.r

```
1 library(pbdMPI, quietly=TRUE)
2 library(pbdDMAT, quietly=TRUE)
3 init.grid()
4
5
6 n <- 1000
7 x <- ddmatrix("rnorm", n, n)
8
9 my.svd <- La.svd(x)
10
11
12 finalize()
```

## Example Script

Run example with 4 ranks:

```
$ mpirun -np 4 Rscript my_svd.r
mpiP:
mpiP: mpiP: mpiP V3.3.0 (Build Sep 23 2013/14:00:47)
mpiP: Direct questions and errors to
      mpiP-help@lists.sourceforge.net
mpiP:
Using 2x2 for the default grid size

mpiP:
mpiP: Storing mpiP output in [./R.4.5944.1.mpiP].
mpiP:
```

## Read Profiler Data into R

Interactively (or in batch) Read in Profiler Data

```
1 library(pbdPROF)
2 prof.data <- read.prof("R.4.28812.1.mpiP")
```

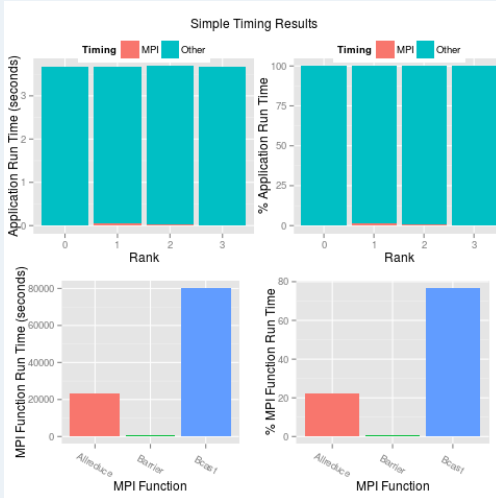
Partial Output of Example Data

```
> prof.data
An mpip profiler object:
[[1]]
  Task AppTime MPITime MPI.
1     0    5.71  0.0387 0.68
2     1    5.70  0.0297 0.52
3     2    5.71  0.0540 0.95
4     3    5.71  0.0355 0.62
5     *   22.80  0.1580 0.69

[[2]]
  ID Lev File.Address Line_Parent_Funct MPI_Call
1   1  0 1.397301e+14      [unknown] Allreduce
2   2  0 1.397301e+14      [unknown]   Bcast
```

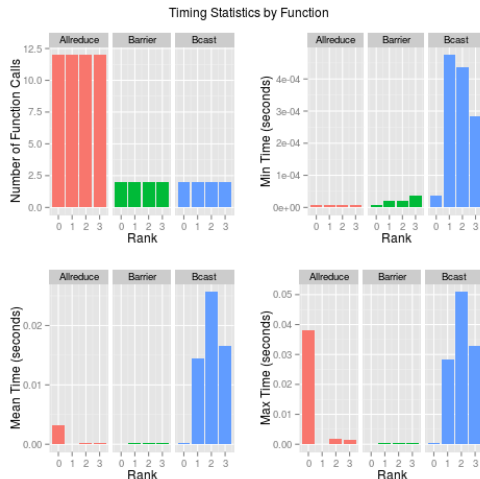
## Generate plots

```
1 plot(prof.data)
```



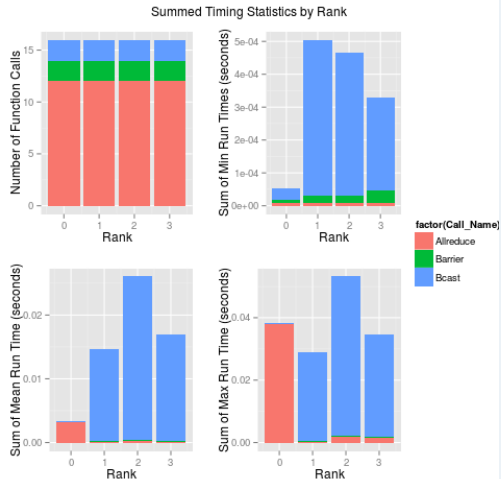
## Generate plots

```
1 plot(prof.data, plot.type="stats1")
```



## Generate plots

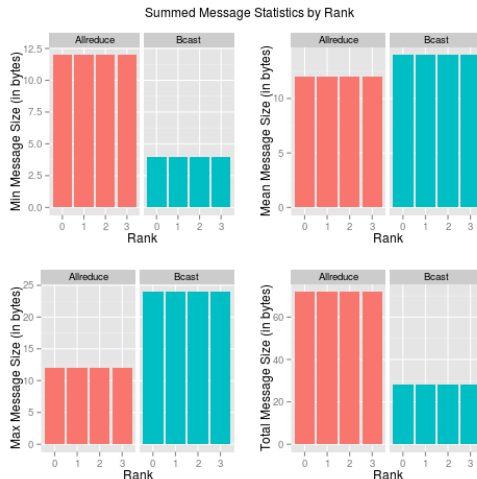
```
1 plot(prof.data, plot.type="stats2")
```





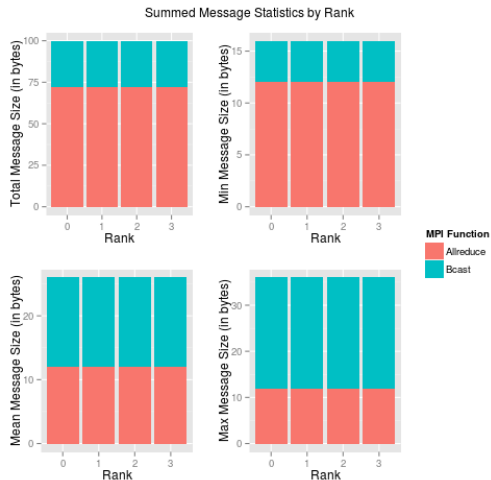
## Generate plots

```
1 plot(prof.data, plot.type="messages1")
```



## Generate plots

```
1 plot(prof.data, plot.type="messages2")
```



## 10 MPI Profiling

- Profiling with the pbdPROF Package
- Installing pbdPROF
- Example
- Summary

## Summary

- **pbdPROF** offers tools for profiling R-using MPI codes.
- Easily builds **fpmi**; also supports **mpiP**.

# Contents

## 11 Wrapup

## Summary

- Profile your code to understand your bottlenecks.
- **pb**dR makes distributed parallelism with R easier.
- Distributing data to multiple nodes
- For truly large data, I/O must be parallel as well.

## The pbdR Project

- Our website: <http://r-pbd.org/>
- Email us at: [RBigData@gmail.com](mailto:RBigData@gmail.com)
- Our google group: <http://group.r-pbd.org/>

## Where to begin?

- The **pbdDEMO** package  
<http://cran.r-project.org/web/packages/pbdDEMO/>
- The **pbdDEMO** Vignette: <http://goo.gl/HZkRt>

Thanks for coming!

# Questions?



<http://r-pbd.org/>

Come see our poster on Wednesday at 5:30!