

Programming with Big Data in R

Drew Schmidt and George Ostrouchov

July 8, 2013



Affiliations and Support

The pbdR Core Team

<http://r-pbd.org>

Wei-Chen Chen¹, George Ostrouchov^{1,2}, Pragneshkumar Patel², Drew Schmidt¹

Ostrouchov, Patel, and Schmidt were supported in part by the project “NICS Remote Data Analysis and Visualization Center” funded by the Office of Cyberinfrastructure of the U.S. National Science Foundation under Award No. ARRA-NSF-OCI-0906324 for NICS-RDAV center.

Chen and Ostrouchov were supported in part by the project “Visual Data Exploration and Analysis of Ultra-large Climate Data” funded by U.S. DOE Office of Science under Contract No. DE-AC05-00OR22725.

¹Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN

²Remote Data Analysis and Visualization Center, University of Tennessee, Knoxville, TN

About This Presentation

Downloads

This presentation and supplemental materials are available at:

<http://r-pbd.org/user2013>

About This Presentation

Speaking Serial R with a Parallel Accent

The content of this presentation is based in part on the **pbdDEMO** vignette *Speaking Serial R with a Parallel Accent*

<https://github.com/wrathematics/pbdDEMO/blob/master/inst/doc/pbdDEMO-guide.pdf?raw=true>

It contains more examples, and sometimes added detail.

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.

About This Presentation

Conventions

We use:

- “.” as a decimal mark
- “,” as order of magnitude separator

Example	Yes	No
One million	1,000,000	1.000.000
One half	0.5	0,5
One thousand and one half	1,000.5	1.000,5

Introduction	pbdR	pbdMPI	GBD	Break	Stats eg's	pbdDMAT	pbdDMAT eg's	Wrapup
oooooo oooooo oooooo	oooo oooo oooo	oooo oooooooo oooooo	ooo ooo ooo		oooo ooo ooo	ooooo oooooo oooooooo	ooo oooo oo	

Contents

- 1 Introduction
- 2 pbdR
- 3 Introduction to pbdMPI
- 4 The Generalized Block Distribution
- 5 Brief Intermission
- 6 Basic Statistics Examples
- 7 Introduction to pbdDMAT
- 8 Examples Using pbdDMAT
- 9 Wrapup

Contents

- 1 Introduction
 - A Concise Introduction to Parallelism
 - Common Terminology
 - R and Parallelism

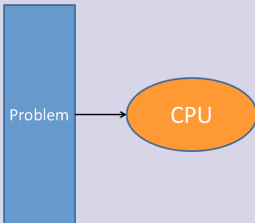
What is Parallelism?

Broadly, *doing more than one thing at a time.*

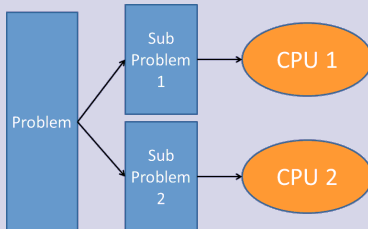
The simultaneous use of multiple compute resources to solve a computational problem:

Parallelism

Serial Programming

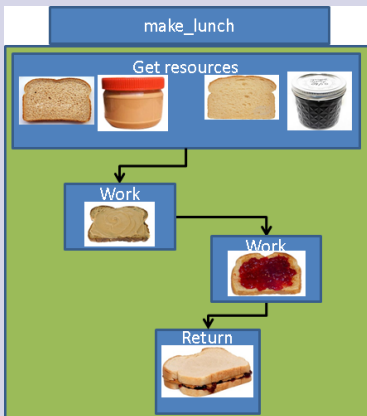


Parallel Programming

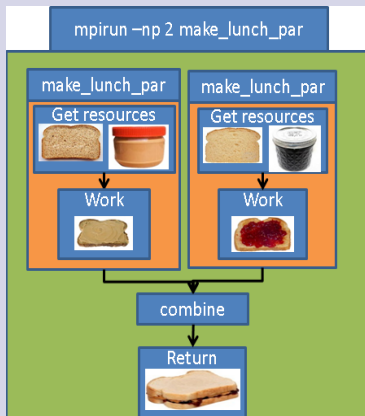


Parallelism

Serial Programming



Parallel Programming



Kinds of Parallelism

- *Data Parallelism*: Data is distributed
- *Task Parallelism*: Tasks are distributed

pbdR Paradigms: Data Parallelism

With data parallelism:

- No one processor/node owns all the data.
- Processors own local pieces of a (conceptually) global object

```

○○○○○●
○○○○○
○○○○○

```

```

○○○○
○○○○

```

```

○○○○
○○○○○○○○
○○○○○○

```

```

○○○
○○○
○○○

```

```

○○○○
○○○
○○○

```

```

○○○○○
○○○○○○
○○○○○○○○

```

```

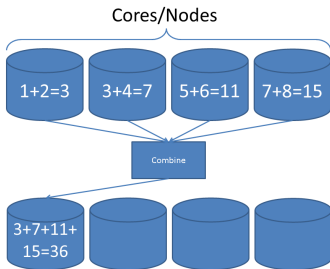
○○○
○○○○
○○

```

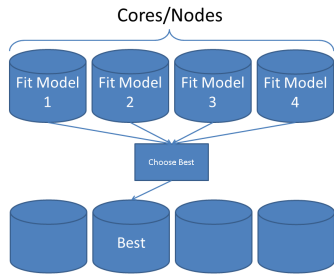
A Concise Introduction to Parallelism

Data vs Task Parallelism

Data Parallelism



Task Parallelism



Difficulty

- 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.

Introduction	pbdR	pbdMPI	GBD	Break	Stats eg's	pbdDMAT	pbdDMAT eg's	Wrapup
○○○○○○○ ○●○○○ ○○○○○	○○○○ ○○○○	○○○○ ○○○○○○○○○ ○○○○○○○	○○○ ○○○ ○○○		○○○○ ○○○ ○○○	○○○○○ ○○○○○○○ ○○○○○○○○○	○○○ ○○○○ ○○	

Common Terminology

Scalability

Scalability: unitless measure of performance;

$$\frac{\tau_i}{\tau_0}$$


```

○○○○○○
○●○○○
○○○○

```

```

○○○○
○○○

```

```

○○○○
○○○○○○○
○○○○○

```

```

○○○
○○○
○○○

```

```

○○○○
○○○
○○○

```

```

○○○○
○○○○○
○○○○○○○

```

```

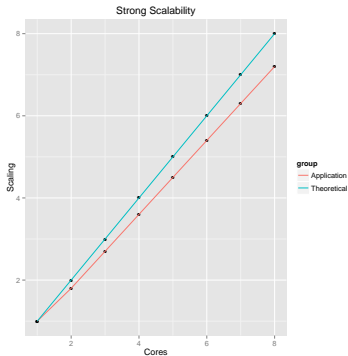
○○○
○○○
○○

```

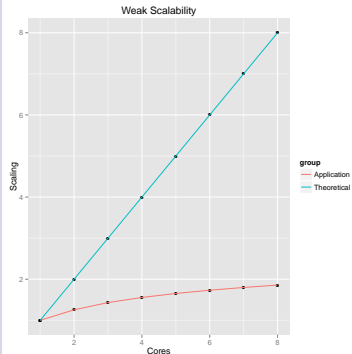
Common Terminology

Types of Scalability: Strong and Weak

Strong

Fix *total* data size

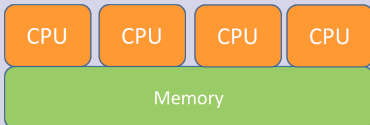
Weak

Fix *local* data size

Shared and Distributed Memory Machines

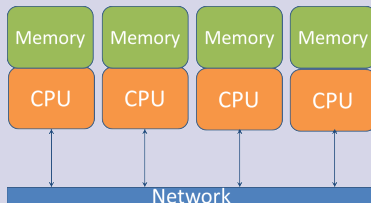
Shared Memory

Different processors can directly access and modify each others' memory. There is only one node.



Distributed

Different processors/nodes can not directly access/modify different processors'/nodes' memory.



○○○○○○
○○○○●
○○○○○

○○○○
○○○○

○○○○
○○○○○○○○
○○○○○○○

○○○
○○○
○○○

○○○○
○○○
○○○

○○○○○
○○○○○○
○○○○○○○○

○○○
○○○○
○○

Common Terminology

Shared and Distributed Memory Machines

Shared Memory Machines

Thousands of cores



Nautilus, University of Tennessee

1024 cores

Distributed Memory Machines

Hundreds of thousands of cores



Kraken, University of Tennessee

112,896 cores

○○○○○○○
 ○○○○○○
 ●○○○○

○○○○
 ○○○○

○○○○
 ○○○○○○○○
 ○○○○○○

○○○
 ○○○
 ○○○

○○○○
 ○○○
 ○○○

○○○○○
 ○○○○○○
 ○○○○○○○○

○○○
 ○○○○
 ○○

R and Parallelism

What about R?

```

○○○○○○
○○○○○
○●○○○

```

```

○○○○
○○○○

```

```

○○○○
○○○○○○○○
○○○○○○

```

```

○○○
○○○
○○○

```

```

○○○○
○○○
○○○

```

```

○○○○○
○○○○○○
○○○○○○○○

```

```

○○○
○○○○
○○

```

Problems with Serial R

- ❶ Slow.
- ❷ If you don't know what you're doing, it's *really* slow.
- ❸ Performance improvements usually for small machines.
- ❹ Very ram intensive.
- ❺ Chokes on big data.

Parallel R Packages

Shared Memory

- 1 **foreach**
- 2 **parallel**
- 3 **snow**
- 4 **multicore**

Distributed

- 1 **Rmpi**
- 2 **R+Hadoop**
- 3 **pbdR**

```

○○○○○○
○○○○○
○○●○○

```

```

○○○○
○○○○

```

```

○○○○
○○○○○○○○
○○○○○○

```

```

○○○
○○○
○○○

```

```

○○○○
○○○
○○○

```

```

○○○○○
○○○○○○
○○○○○○○○

```

```

○○○
○○○○
○○

```

R and Parallelism

The solution to many of R's problems is parallelism. However ...

What we have

- ① Mostly serial.
- ② Mostly not distributed
- ③ Data parallelism mostly explicit

What we want

- ① Mostly parallel.
- ② Mostly distributed.
- ③ Mostly implicit.

Why We Need Parallelism

- 1 Saves time (long term).
- 2 Data size is skyrocketing.
- 3 Necessary for many problems.
- 4 Like it or not, it's coming.
- 5 *It's really cool.*