

# The State of HPC in the Open Source R Ecosystem

Drew Schmidt

November 12, 2016



## Support and Disclaimer

This material is based upon work supported by the National Science Foundation Division of Mathematical Sciences under Grant No. 1418195.

The findings and conclusions in this presentation have not been formally disseminated by the U.S. Department of Health & Human Services nor by the U.S. Department of Energy, and should not be construed to represent any determination or policy of University, Agency, Administration and National Laboratory.



## Speaker Bio

- M.S. in mathematics.
- Former statistics consultant.
- Former full-time university researcher.
- Now a miserable grad student.
- Prolific complainer on twitter.



## Goals of This Talk

- Convince you that R has a legitimate place in HPC.
- Give a broad overview of the R package landscape.
- Make some very safe predictions.



# Contents

- 1 Background and Motivation
- 2 A Little History
- 3 Packages
- 4 A Closer Look at HPC and R
- 5 Concluding Remarks



## 1 Background and Motivation

- R Is Weird
- R Is Popular



## 1 Background and Motivation

- R Is Weird
- R Is Popular



## Types

- logical (“boolean”)
- integer (32-bit int)
- numeric (double)
- complex (double complex)
- character (string)

Also raw and external pointer

## Data Structures

- Vectors (matrices, n-dim arrays)
- Lists (arrays of pointers)
- Dataframes (lists with constraints)
- Environments (hash tables?!)
- That’s it.





## Happy Opposite Day!

```
1 T
2 ## [1] TRUE
3 F
4 ## [1] FALSE
5
6 T <- FALSE
7 F <- TRUE
8
9 T
10 ## [1] FALSE
11 F
12 ## [1] TRUE
```

## Odd Conventions

- `.` has no semantic meaning (except when it does
  - `t.test()`
  - `t.data.frame()`
- A *package* is installed in a *library*.

## Package or Library?

- I wrote a library.
- I put that library into a package.
- I installed the package ... into a library.
- I load the package with `library()` ???



\*BOOM\*



## 1 Background and Motivation

- R Is Weird
- R Is Popular























## Part Programming Language, Part Data Analysis Package

*“R is a shockingly dreadful language for an exceptionally useful data analysis environment.”*









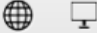
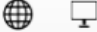
— Tim Smith, from **aRrgh: a newcomer’s (angry) guide to R**.



## IEEE Spectrum's 2014 Ranking of Programming Languages

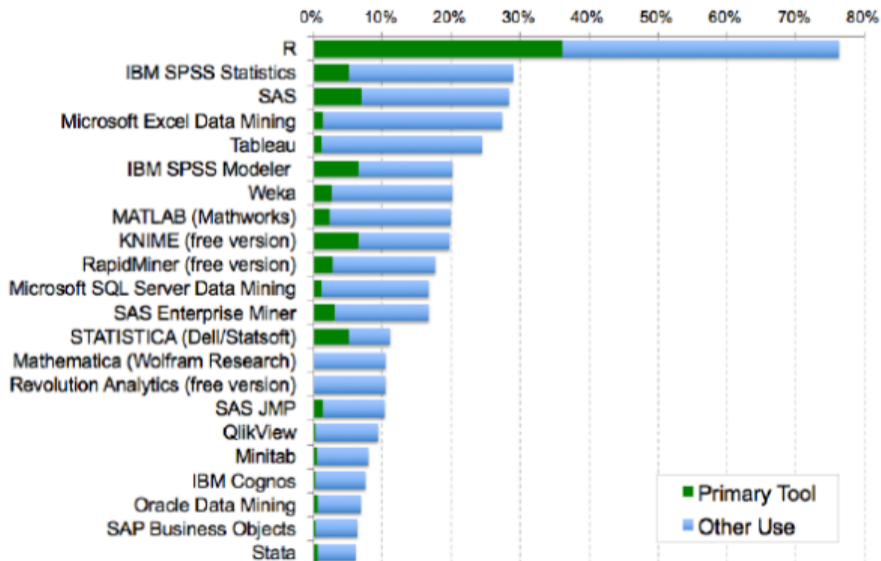
Language Rank	Types	Spectrum Ranking
1. Java	  	100.0
2. C	  	99.3
3. C++	  	95.5
4. Python	 	93.4
5. C#	  	92.4
6. PHP		84.7
7. Javascript	 	84.4
8. Ruby		78.8
9. R		74.2
10. MATLAB		72.9

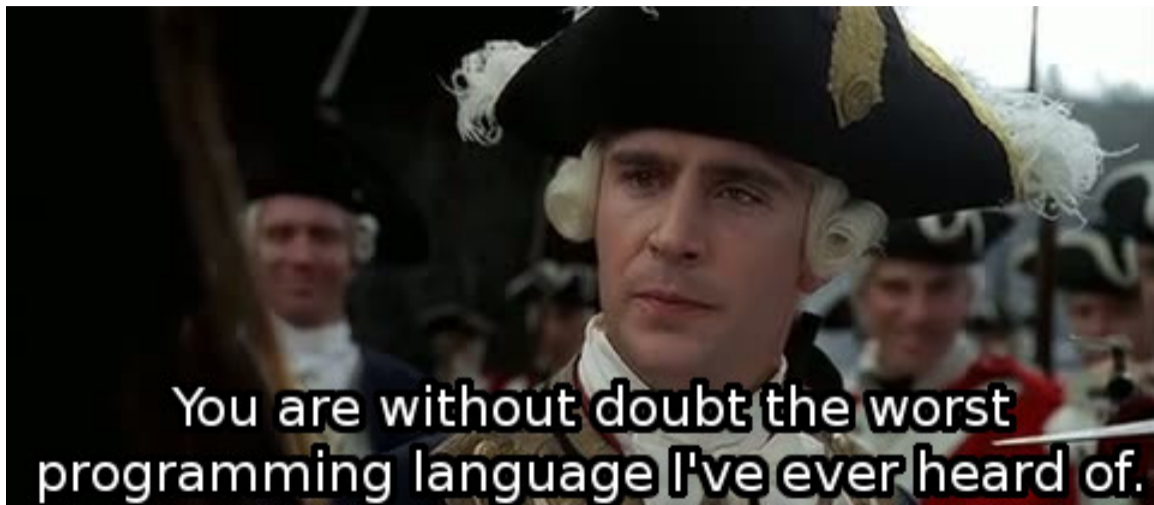
## IEEE Spectrum's 2016 Ranking of Programming Languages

Language Rank	Types	Spectrum Ranking
1. C		100.0
2. Java		98.1
3. Python		98.0
4. C++		95.9
5. R		87.9
6. C#		86.7
7. PHP		82.8
8. JavaScript		82.2
9. Ruby		74.5
10. Go		71.9



## Rexer 2015 data scientist survey







# Computer Scientists Hate It!



Looks **Weird**  
Is **Useful**

Esoteric language topples industry standards using one weird trick. Click to learn its stunning secrets.

**LEARN THE TRUTH NOW**



## Why use R at all?

- Most diverse set of statistical methods available.
- Rapid prototyping.
- CRAN (and increasingly GitHub) packages.
- *Awesome* community.
- Syntax is designed for analysis of data.



## 2 A Little History

- Statistics, Data Science, Big Data, and So On
- Enter R



## 2 A Little History

- Statistics, Data Science, Big Data, and So On
- Enter R



## HPC: Not Just for PDE'S Anymore!

- R's use in HPC.
- No traditional HPC...
- Lots of interesting work





# About Traditional HPC...

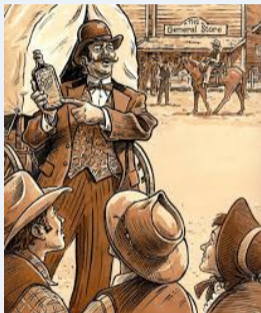
The image shows a screenshot of a web browser window. The address bar displays the URL [www.dursi.ca/hpc-is-dying-and-mpi-is-killing-it/](http://www.dursi.ca/hpc-is-dying-and-mpi-is-killing-it/). The browser's navigation bar includes icons for back, forward, and search, along with a star icon for bookmarks and a red 'ABP' icon. The website's header is a dark blue bar with the navigation links 'HOME', 'ABOUT', and 'CONTACT' in white text. The main content area features a dark, stylized illustration of a scientist in a white lab coat and glasses, sitting at a computer workstation. The scientist is holding a purple flask with a glowing purple liquid inside. In the background, there are computer monitors and a rack of server units with red labels 'A7', 'A8', and 'A9'. A large, glowing green and purple DNA double helix is superimposed over the scene. The title 'HPC is dying, and MPI is killing it' is written in large, white, sans-serif font across the center of the illustration. Below the title, on the left, is a small circular profile picture of Jonathan Dursi and his name 'Jonathan Dursi'. On the right, there is a small white icon of a hand holding a mouse cursor and the text 'hpc'.

## Changing Landscape of HPC

- “non-traditional” HPC: everybody but physics.
- What kind of software do they need?
- Can we leverage any existing HPC stuff?



## Problems with "Big Data" Software



- *Many* frameworks; what do they all do?
- Don't always play nice with HPC systems.
- Often not as "high level" as advertised.
- **Almost exclusively batch!**

Data Analysis Is An Interactive Activity

# Data analysis is an interactive activity<sup>a</sup>

---

<sup>a</sup>Data analysis is an interactive activity



# Data science in action



## 2 A Little History

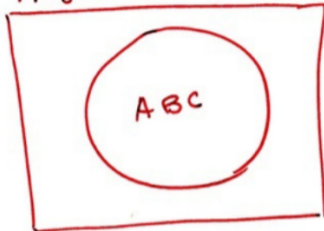
- Statistics, Data Science, Big Data, and So On
- Enter R



JHC  
①

## Algorithm Interface

5/5/76



XABC

ABC: general  
(FORTRAN)  
algorithm

XABC: FORTRAN  
subroutine to  
provide interface  
between ABC &  
language and/or  
utility programs

XABC (INSTR, OUTSTR)

### 3 Packages

- Advanced Compute Packages
- HPC Packages
- Hadoop and Applications
- Ok, So What?





## Where to Begin?

- Many packages of varying scope and quality.
- 1 core package (parallel)
- Over 100 contributed packages  
<https://cran.r-project.org/web/views/HighPerformanceComputing.html>
- Even more on GitHub.





### 3 Packages

- Advanced Compute Packages
  - HPC Packages
  - Hadoop and Applications
  - Ok, So What?



## Out of Core Packages

- ff, bigmemory and friends
- R is very “copy happy”
- Many statisticians don't know about things like XSEDE.
- Others hear “Linux” and run away screaming.
- Bizarrely, cloud computing is changing this.



## Rcpp

- Rcpp
- RcppArmadillo, RcppEigen
- RcppParallel
- ...

### 3 Packages

- Advanced Compute Packages
- HPC Packages
- Hadoop and Applications
- Ok, So What?



## Accelerator Packages

- gputools, Magma, HiPLARM, a few others.
- Accessibility mostly from things like nvblas and Intel MKL.



## Distributed Packages

- Rmpi
- snow
- pbdMPI and friends



## Remote Evaluation Packages

- rzmq, pbdZMQ
- remoter, future

### 3 Packages

- Advanced Compute Packages
- HPC Packages
- **Hadoop and Applications**
- Ok, So What?



## Hadoop et al Packages

- RHadoop, RHIPE
- SparkR
- sparklyr
- h2o

## “Applications”

- dplyr and data.table
- caret
- randomForest
- xgboost



### 3 Packages

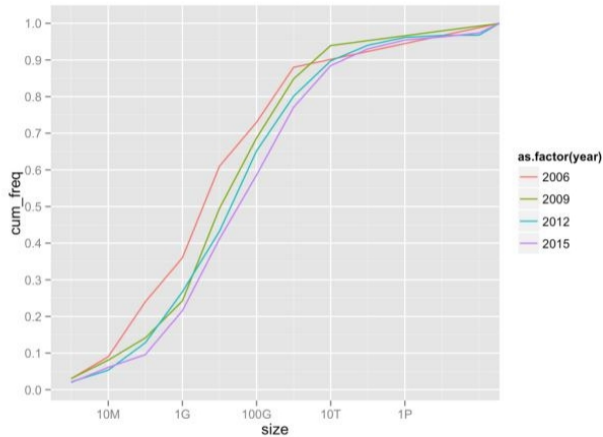
- Advanced Compute Packages
- HPC Packages
- Hadoop and Applications
- Ok, So What?



## Is the R community using this stuff?

- Short answer: yes.
- Long answer: mostly single-node parallelism.
- Hard truth: in addition to hype and buzzwords — fear and distrust

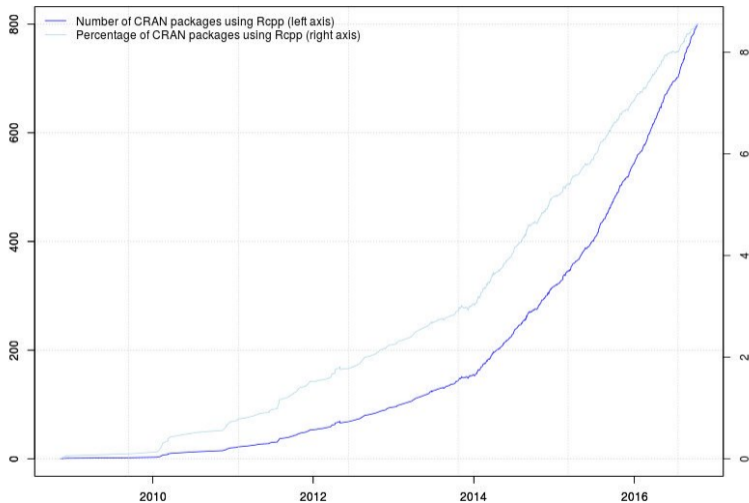




## BIG RAM IS EATING BIG DATA – SIZE OF DATASETS USED FOR ANALYTICS

NOVEMBER 18, 2015 EDUCATION SZILARD PAKKA 1 COMMENT 3

## Growth of Rcpp usage on CRAN



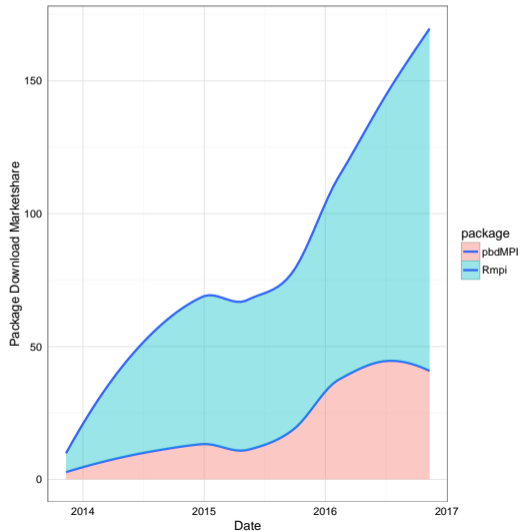
Source <https://twitter.com/eddelbuettel/status/787740983433854977>



## 4 A Closer Look at HPC and R



# HPC may be dying, but we're behind the times




## “OLCF Researchers Scale R to Tackle Big Science Data Sets”



- A problem that takes several hours on Apache Spark [was analyzed] in less than a minute using R on OLCF high-performance hardware.
- “... *for situations where one needs interactive near-real-time analysis, the **pbdR** approach is much better.*”

[https://www.hpcwire.com/2016/07/06/  
olcf-researchers-scale-r-tackle-big-science-data-sets/](https://www.hpcwire.com/2016/07/06/olcf-researchers-scale-r-tackle-big-science-data-sets/)

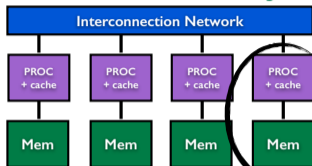


"R? in *my*  
HPC?"

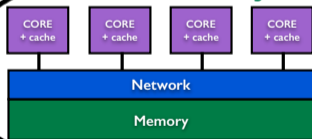
It's more likely than you think.

**FREE PC CHECK!**

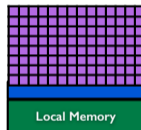
## Distributed Memory



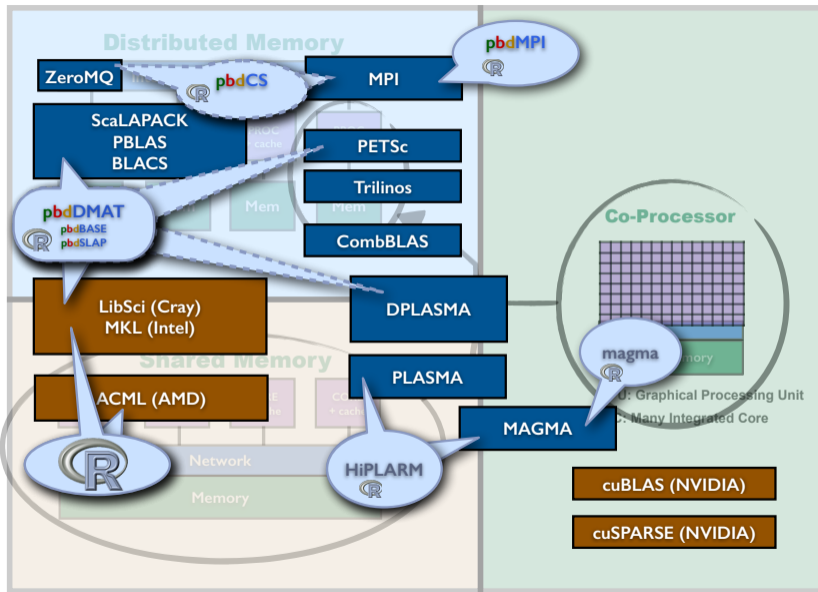
## Shared Memory



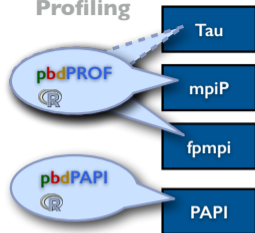
## Co-Processor



GPU: Graphical Processing Unit  
MIC: Many Integrated Core



Profiling



I/O



Learning



— Released      - - - - - Under Development

## 5 Concluding Remarks



## The Future?

- Better dplyr backends.
- More threading + accelerator usage in packages (Rcpp + RcppParallel).
- Astronomical amounts of buzz in the Hadoop/Spark-and-friends space — will ultimately hurt us in the MPI space.





~Thanks!~

# Questions?



Email: [wmathematics@gmail.com](mailto:wmathematics@gmail.com)



GitHub: <https://github.com/wmathematics>



Web: <http://wmathematics.info>



Twitter: [@wmathematics](https://twitter.com/wmathematics)