

# Analyzing Analytics: Advanced Performance Analysis Tools for R

Drew Schmidt\*, Wei-Chen Chen, Christian Heckendorf, George Ostrouchov

November 14, 2016



# The pbdR Core Team

Wei-Chen Chen<sup>1</sup>

George Ostrouchov<sup>2,3</sup>

Drew Schmidt<sup>3</sup>



## Support

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.

<sup>1</sup>FDA  
Washington, DC, USA

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

<sup>3</sup>University of Tennessee  
Knoxville TN, USA



# Contents

- 1 Why R?
- 2 Advanced Profilers
- 3 A Few Examples
- 4 Concluding Remarks



- 1 Why R?
- 2 Advanced Profilers
- 3 A Few Examples
- 4 Concluding Remarks

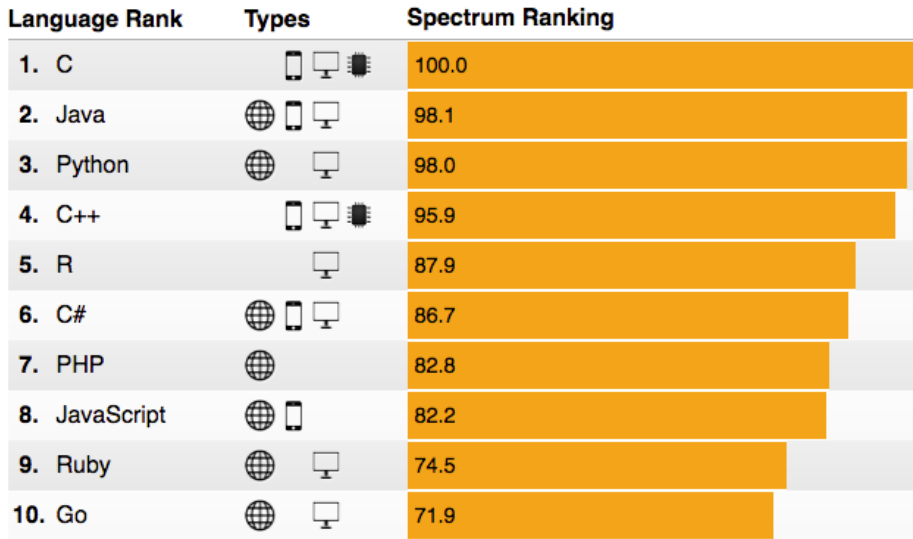


# "R? in *my* HPC?"

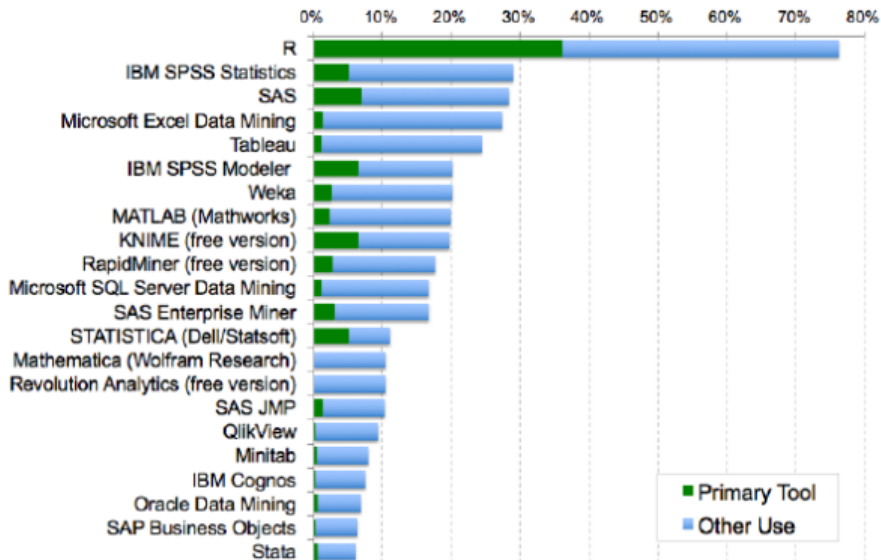
It's more likely than you think.

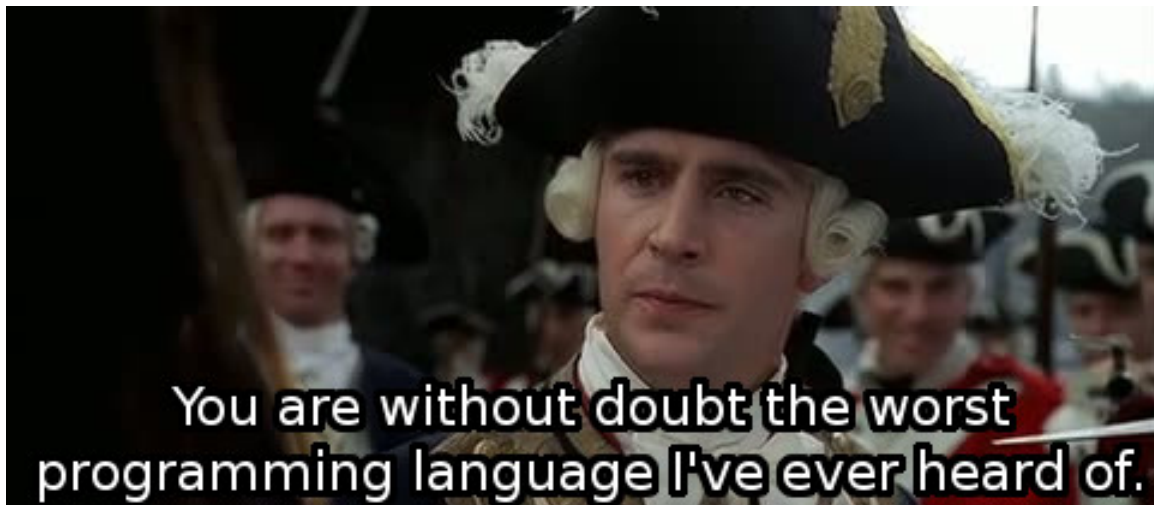
**FREE PC CHECK!**

## IEEE Spectrum's 2016 Ranking of Programming Languages



## Rexer 2015 data scientist survey

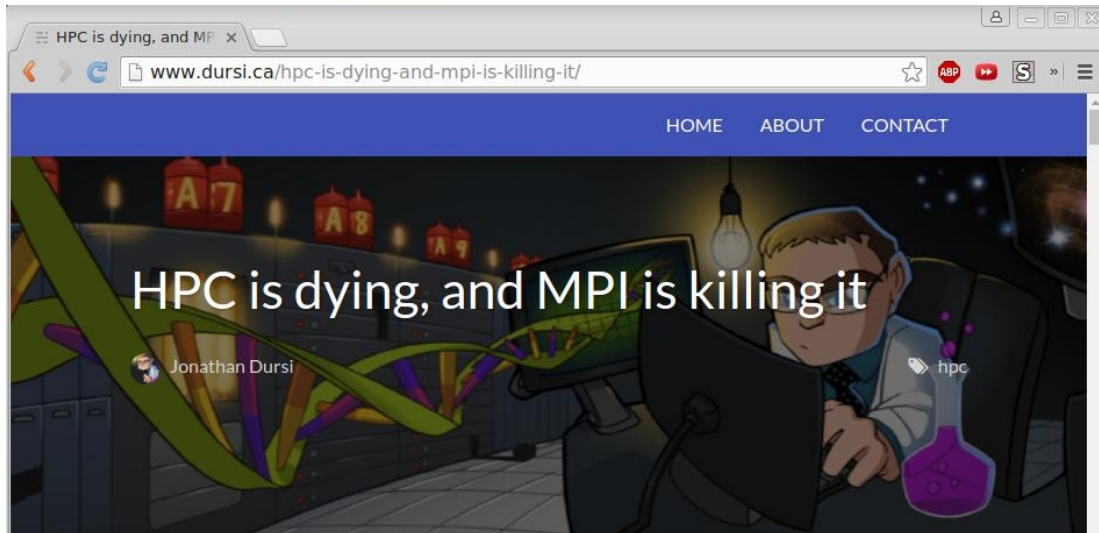




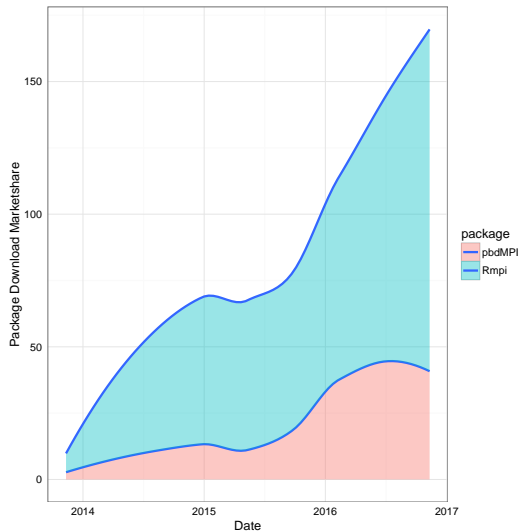




# About Traditional HPC...



# 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/>

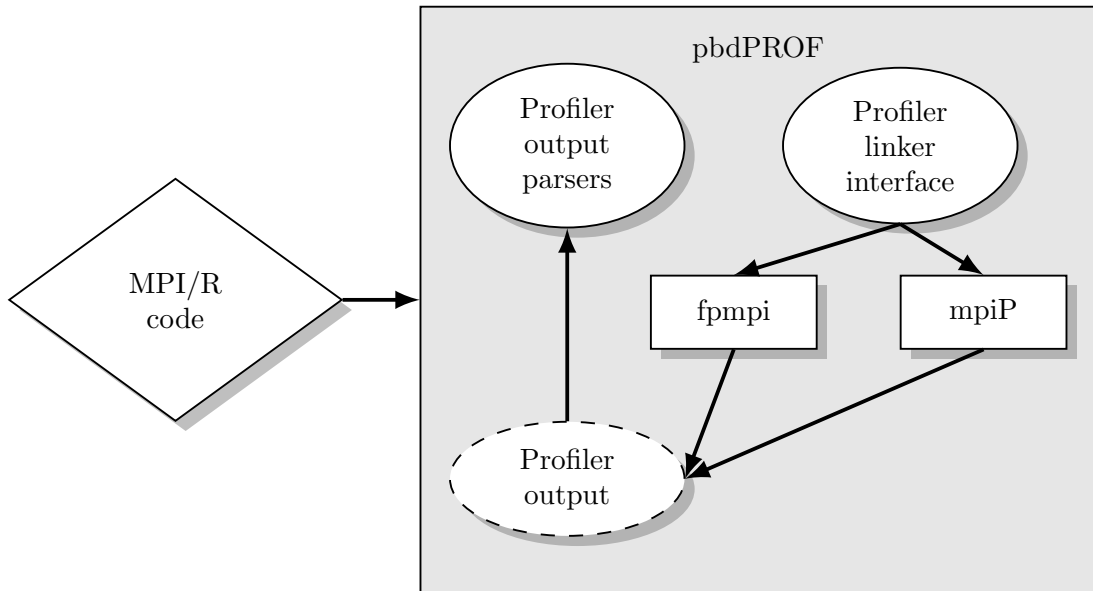
- 1 Why R?
- 2 Advanced Profilers
- 3 A Few Examples
- 4 Concluding Remarks

## Profiling in R Versus Traditional HPC

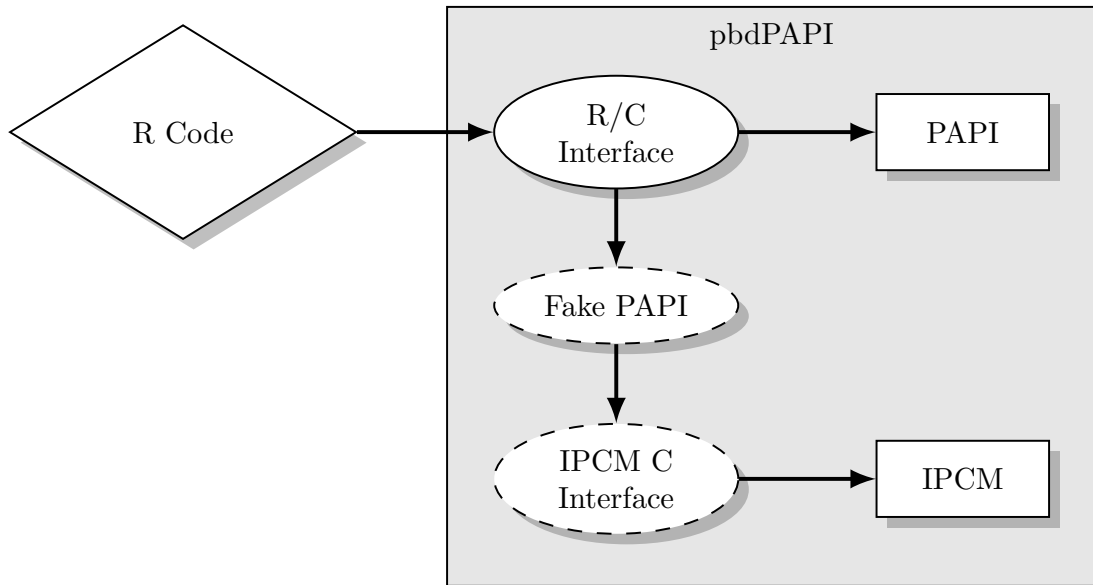
- Profiling in R (core): Stack sampler, some memory profilers.
- Every (other) profiling package in R: wall clock times.
- What's missing: MPI profilers, performance counters.

## Three New R Packages

- pbdPROF - MPI profiler support (fpmpi and mpiP).
- pbdPAPI - PAPI bindings and utilities.
- hpcvis - Plot methods for the above 2.

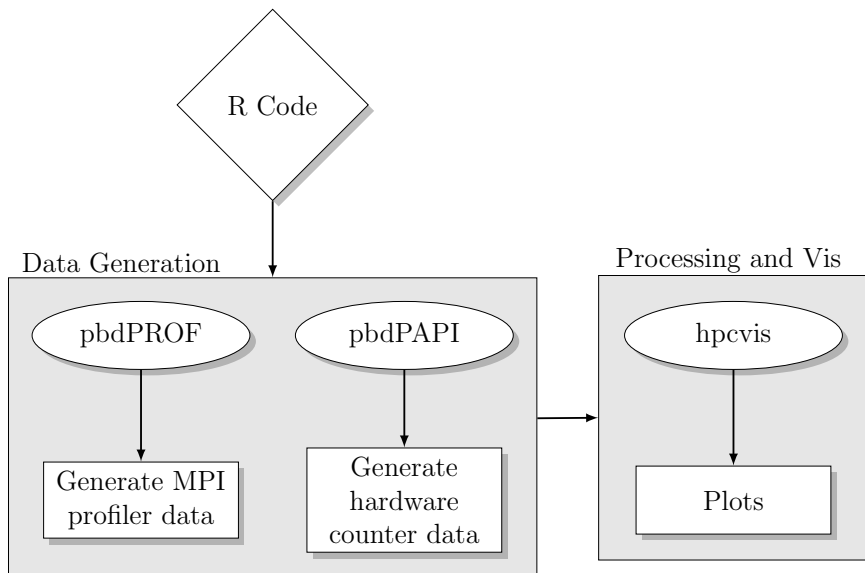






## pbdPAPI High-Level Interface

Function	Measurement
<code>system.cache()</code>	Cache misses, hits, accesses
<code>system.flops()</code>	MFLOPS
<code>system.idle()</code>	Idle cycles
<code>system.cpuormem()</code>	Classify CPU/RAM bound
<code>system.utilization()</code>	CPU utilization



## pbdPROF + hpcvis Syntax Example

```
1 library(pbdPROF)
2 library(hpcvis)
3
4 prof_file <- "R.4.26802.1.mpiP"
5 x <- read.prof(prof_file)
6
7 profplot(x, color=FALSE, stacked=FALSE, title=NULL, which=c(2,4))
8 profplot(x, plot.type="stats", color=FALSE, stacked=TRUE, title=NULL)
```

## pbdPAPI + hpcvis Syntax Example

```
1 library(pbdPAPI)
2 library(hpcvis)
3
4 cb <- cachebench(foo(x), bar(x))
5 papiplot(cb, label.angle=15, levels=1:2)
```

## Note

All plots for the remainder of this presentation were made with hpcvis.

1 Why R?

2 Advanced Profilers

3 A Few Examples

- Clustering
- K-Means Clustering Algorithms
- Distributed Clustering with pmclust

4 Concluding Remarks

## Clustering Approaches

- K-means and GMM
- Two unsupervised learning approaches.
- Trying to discover latent groupings in the data.
- K-means
  - Cover data with  $k$  *identical spheres*.
- GMM
  - Can use different ellipsoids, pointing in different directions.

## Now With Math

- Partition  $N$  data points in  $p$  dimensions into  $K$  clusters

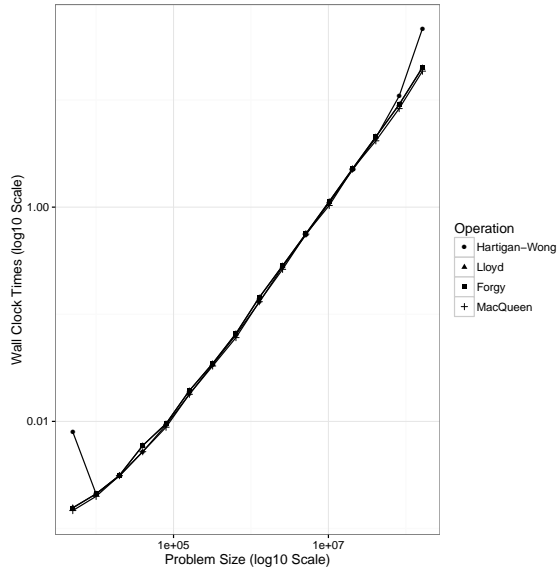
- k-means model (stated as Gaussian mixture):

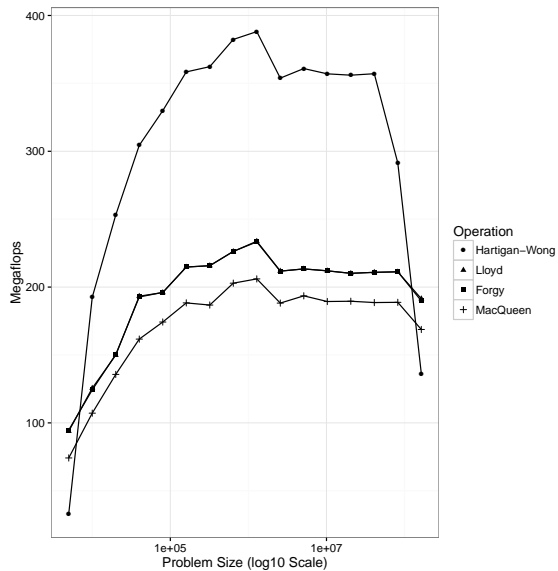
- $X_1, \dots, X_N \stackrel{iid}{\sim} \phi(\cdot | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ 
  - $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_k\}$
  - $\boldsymbol{\Sigma} = \{\sigma_1, \dots, \sigma_k\}$
  - $\phi(\cdot | \boldsymbol{\mu}, \sigma_k) = \sum_{k=1}^K \eta_k \phi_p(\cdot | \mu_k, \sigma_k)$
  - $2Kp + 1$  parameters

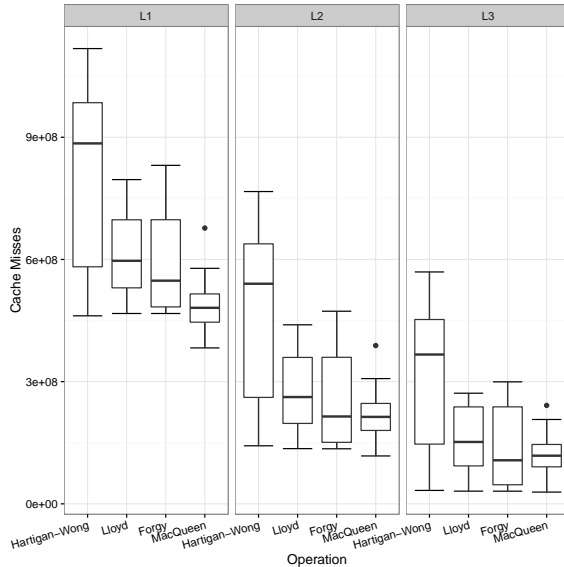
- **pmclust**: Gaussian mixture model (GMM):

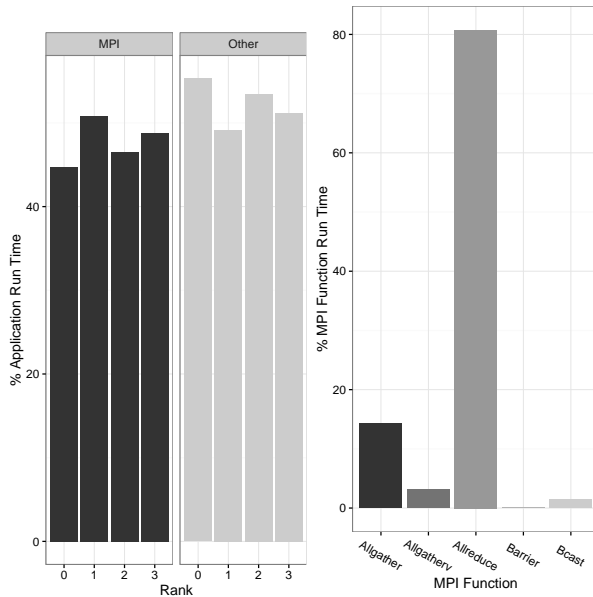
- $X_1, \dots, X_N \stackrel{iid}{\sim} \phi(\cdot | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ 
  - $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_k\}$
  - $\boldsymbol{\Sigma} = \{\Sigma_1, \dots, \Sigma_K\}$
  - $\phi(\cdot | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{k=1}^K \eta_k \phi_p(\cdot | \mu_k, \Sigma_k)$
  - $2Kp + K \frac{p \times (p+1)}{2}$  parameters

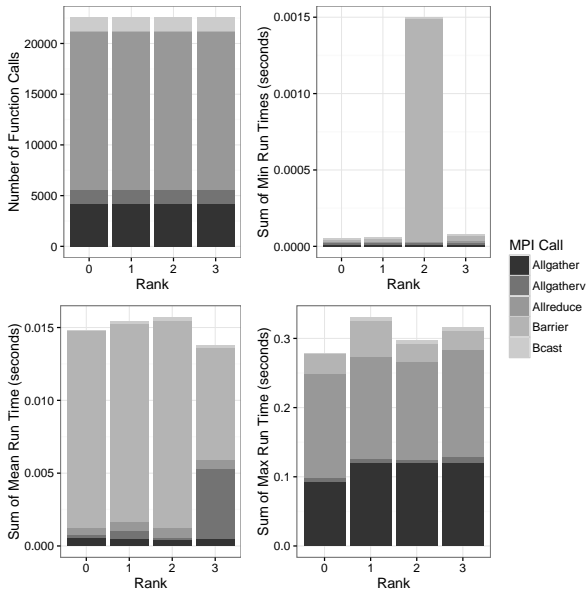












- 1 Why R?
- 2 Advanced Profilers
- 3 A Few Examples
- 4 Concluding Remarks

## Challenges

- pbdPAPI still difficult for R users to install.
- Mostly Linux-only.

## Future Work

- Build a web interface?
- Integrate with tau?

~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)