# Large-scale data analysis in R

**Research Computing Center
and the Dept. of Human Genetics
University of Chicago**

Peter Carbonetto
pcarbo@uchicago.edu

*November 14, 2017*

# Getting started

**Initial setup:**

- Wifi

- Power outlets

- RCC cluster access

- Yubikeys (`rcc-guest-xxxx`)

- Helper(s)

- Handout / exercises

- Pace, questions (e.g., keyboard shortcuts)

# Getting started

1. Download workshop packet.

```
$ cd ~
$ wget https://tinyurl.com/yad63z9c \
    -O temp.tar.gz
$ tar zxvf temp.tar.gz
$ cd R-large-scale
$ cd code
```

2. Add these three lines to your `.bashrc` file.

```
# SLURM.
export SACCT_FORMAT="jobid,partition,user,account%12,alloccpus,node%12,elapsed,totalcpu,maxRSS,ReqM"
export SQUEUE_FORMAT="%13i %12j %10P %10u %12a %8T %9r %10l %.11L %5D %4C %8m %N"
```

# Getting started

## 3. Start up three midway2 sessions.

# Getting started

4. Request a compute node.

```
$ screen -S rcc_workshop
$ sinteractive --partition=broadwl \
    --time=2:00:00 --account=rcc-guest \
    --reservation=rworkshop2
$ echo $HOSTNAME
```

5. Start up R programming environment.

```
$ module load R/3.4.1
$ R
R> getwd()
```

# Aims of workshop

- ~~Do an analysis of a large-scale data set in R.~~

- Develop some useful skills for large-scale data analysis in R *within a high-performance computing environment*, and apply these skills to a medium-scale data set.

# Outline

- Initial setup.

- Brief introduction.

- Part 1: Warm-up with PCA of RegMap data.

- Part 2: Implementing multithreaded computation in R for analysis of genetic adaptation to climate.

- Brief recap.

# Part 1: PCA of the RegMap data



**Adaptation to Climate Across the _Arabidopsis thaliana_ Genome**

Angela M. Hancock,[1] Benjamin Brachi,[2] Nathalie Faure,[2] Matthew W. Horton,[1]
Lucien B. Jarymowycz,[1] F. Gianluca Sperone,[1] Chris Toomajian,[3] Fabrice Roux,[2] Joy Bergelson[1]*

www.sciencemag.org   **SCIENCE**   VOL 334   7 OCTOBER 2011

1. Inspect RegMap data.

2. Run PCA interactively and assess time & memory needs.

3. Examine PCA results.

4. Run PCA using Rscript.

5. Run PCA using SLURM job engine, and inspect output.

# *Exercise:* Run PCA using SLURM job engine, and inspect output.

```
sbatch pca.sbatch
```

- Useful commands while job is running:

```
squeue --user=<cneitd> | less -S
ssh midway2-xxxx
```
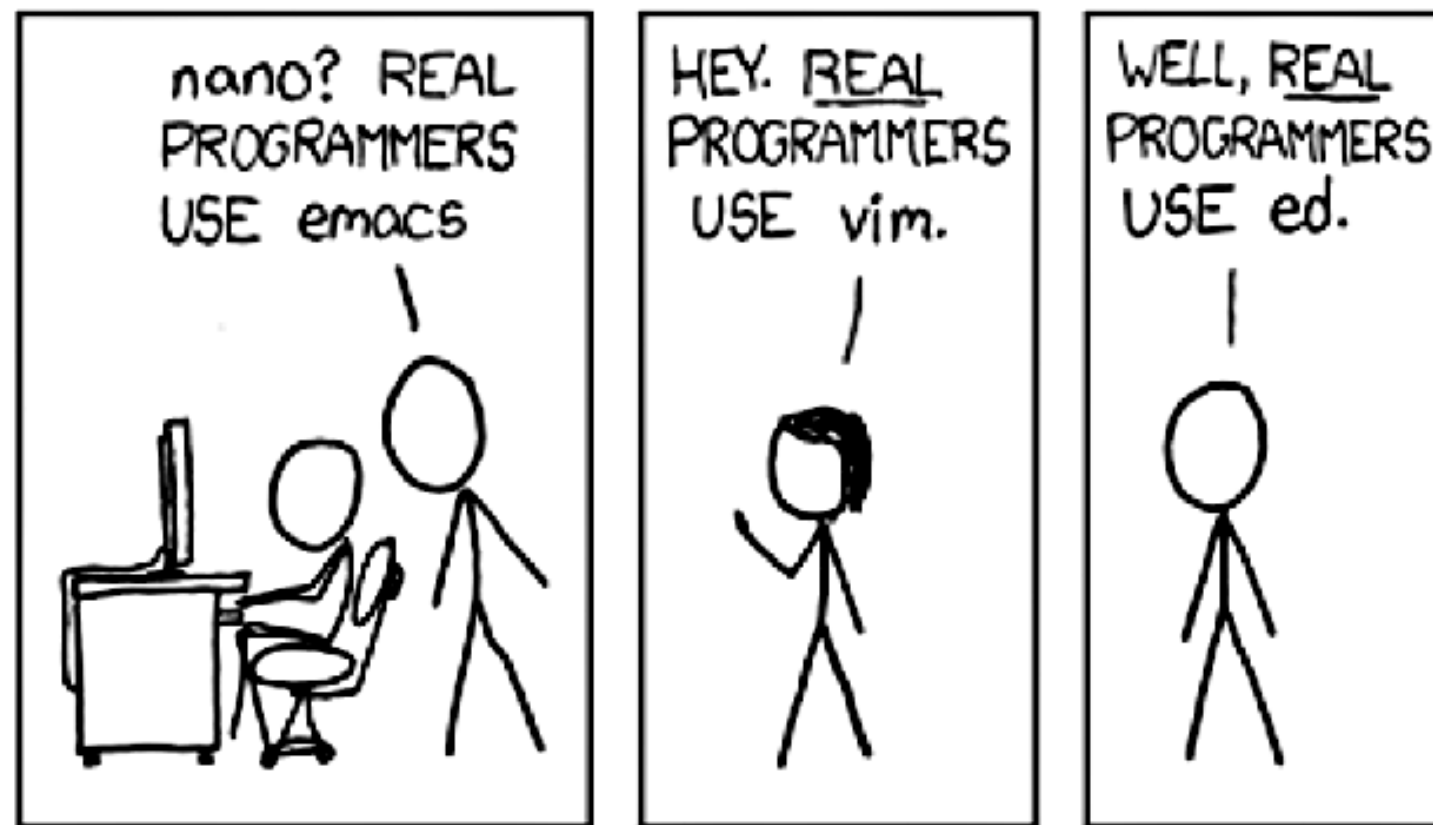
- Useful commands after job has completed:

```
htop --user=<cnetid>
sacct --user=<cnetid> --units=G | less -S
less ../output/pca_err.txt
less ../output/pca_out.txt
```

# Part 2: Genetic analysis of climate variables in RegMap data

1. Inspect climate ("phenotype") data.

2. Run analysis interactively and assess time & memory needs.

3. Experiment with multithreaded matrix operations to improve computation time.

4. Experiment with parallel computation of the weights using `mclapply`.

5. Automate `climate.R` analysis using command-line arguments.

6. Run analysis using SLURM job engine, and inspect output.

7. Use SLURM job engine to automate analysis of all 48 climate variables.

# There is no best tool—use *whatever works for you*.

# Some general advice

1. `help(package = cool_package)`.

2. Use **midway2**, not midway1.

3. Email [help@rcc.uchicago.edu](mailto:help@rcc.uchicago.edu) R help on the RCC cluster.

4. Learn to avoid loops as much as possible; e.g., use `apply()`, `lapply()`, `tapply()`, `do.call()`.

5. Document your setup—start with `sessionInfo()`.

# After the workshop

- Feedback.

- Suggestions for future R workshop topics (e.g., Rcpp).