# DSC aids in reproducible research

### Motivation

- **★** Statistical comparisions typically performed is suboptimal in many ways
- **→ DSC** is an attempt to make statistical comparisions extensible & reproducible

### Goal

- \* A platform to make it simple, even fun to perform DSCs in research [1]
- \* A repository of DSCs for projects in & outside Stephens Lab

Tell DSC how benchmark data is generated

simulate:

exec: datamaker.R

Nsamp: 2, 10, 50

**Ngene: 10000** 

params:

tissue: Adipose-Subcutaneous,

breaksample: FALSE, TRUE

return: data, meta = R(data\$meta)

Command interface

> dsc settings.dsc -j8 --host midway.rcc.uchicago.edu

INFO: DSC script exported to settings.html ...

INFO: Constructing DSC from settings.dsc ...

INFO: Running DSC jobs ...

simulate\_1+correction\_1+transform\_1+test\_3+score\_1: 100% [========] 5 20.8/s

simulate\_1+correction\_1+transform\_1+test\_4+score\_1: 100% [========] 5 20.8/s

simulate\_1+correction\_1+transform\_2+test\_3+score\_1: 100% [========] 5 20.8/s

simulate\_1+correction\_1+transform\_2+test\_4+score\_1: 100% [========] 5 20.8/s

simulate\_1+correction\_1+transform\_3+test\_3+score\_1: 100% [========] 5 20.8/s

simulate\_1+correction\_1+transform\_3+test\_4+score\_1: 100% [=========] 5 20.8/s

simulate\_1+correction\_1+transform\_4+test\_3+score\_1: 100% [========] 5 20.8/s

simulate 1+correction 1+transform 4+test 4+score 1:100% [========] 5 20.8/s

INFO: Building output database rna\_seq.rds ...

INFO: DSC complete!

INFO: Elapsed time 254.829 seconds.

.alias: args = Pack()

simulate\_normalized(simulate):

voom.normalize: TRUE

(Adipose-Subcutaneous, Lung)

**seed:** R(1:50)

Tell DSC we can generate data differently based on what we've done

## Illustration

Comparison of methods for differential gene expression analysis of RNA sequencing data

# Dynamic Statistical Comparisons

Gao Wang

Matthew Stephens

gaow@uchicago.edu mstephens@uchicago.edu University of Chicago http://github.com/stephenslab

## Simulation data set A.I

DSC

scripts

MD oracle K + MST

MD arbitrary K + MST

ositive PMD oracle K + MST

Positive PMD arbitrary K + MST

prowser

**Simulation** data set A.II **Simulation** data set B.I

Benchmark

results

browser

**SVA** myrna **RUV** 

voom quasibinom edgeRGIm DESeq2Glm<sup>1</sup>

limma ash jointash edgeR

(correction \* transform \* test[3:4], correction \* transform[1] \* test[5], test[1:2]) \* score by connecting various computational routines

# DSC vs. scientific workflow systems

- **★** DSC harnesses workflow systems [3] yet is tailored for use with methods development in data science
- **★** Compact syntax to configure families of methods
- **★** With steps linked by *variables* not *files* DSC makes it easy to build benchmarks from existing scripts

### DSCR (a.k.a. DSC1) DSC2: a multilingual system

**★** The pilot phase of DSC

design implemented in

complex and capacity

for large scale problems

the R language

**★** Readily usable [2]

Lacks flexibility for

- \* Seamless integration of routines written in R, Matlab and Python
  - **★** Flexible assembly of multiple statistical procedures at will
  - Designed with modern workflow management system standards
  - **Focused on experience of** simplicity
  - Works in three major computational environments

Tell DSC to consider these counfounder



correction: exec: SVA.R, RUV.R, These are statistical myrna.R methods tackling params: data: \$data the same problem .alias: args = Pack() return: data

... and these normalization methods for RNA-Seg data

Give DSC a

DE methods

evaluate

Finally, construct DSC pipelines

metric to

score

control

methods

exec: voom.R, quasibinom.R, edgeRglm.R, DESeq2glm.R params: data: \$data .alias: args = Pack() return: data

edgeR.R, DESeq2.R, ash.R, jointash.R, limma.R ... and these differential data: \$data expression (DE) .alias: args = Pack() analysis

exec(1,2): glm: TRUE, FALSE exec(5): robust: TRUE, FALSE return: output

Amethod may have different flavors (parameters)

exec: score.R params: data: \$meta work\_dir: \$output return: result



## Reference & Links

[1] DSC2 URL / QR code (scan above) github.com/stephenslab/dsc2

[2] DSCR URL

qithub.com/stephens999/dscr

[3] DSC2 uses pysos library for workflow management and redis cluster for job distribution github.com/bopeng/sos redis.io/topics/cluster-tutorial

[4] The RNA-Seq illustration is adapted from a DSCR implemention by Mengyin Lu github.com/mengyin/dscr-gtex-total