# Automated RNA-seq Differential Expression Validation

Center for Health Bioinformatics, Harvard School of Public Health

# Pipeline proliferation

# Complexity

- **Installation**
  - **Third party tools**
  - **Bizarre environments**
- **Choices**
  - **Tools, parameters**
- **Data**
- **Glue**



3

# Development goals of bcbio-nextgen

▶ Community developed and driven

▶ Scalable

▶ Easy to install. Easy to use and extend.

▶ Well-documented

▶ Quantifiable

# Installation

Tools
 compatible
 versioned
 no sudo, no problem
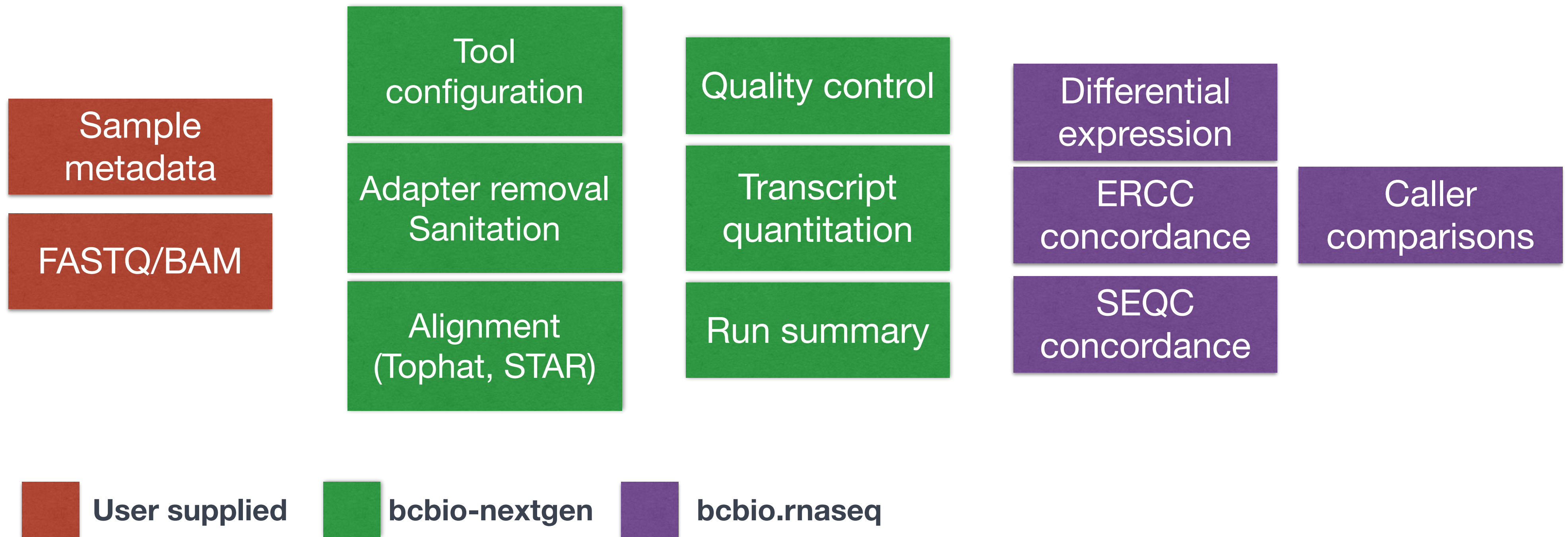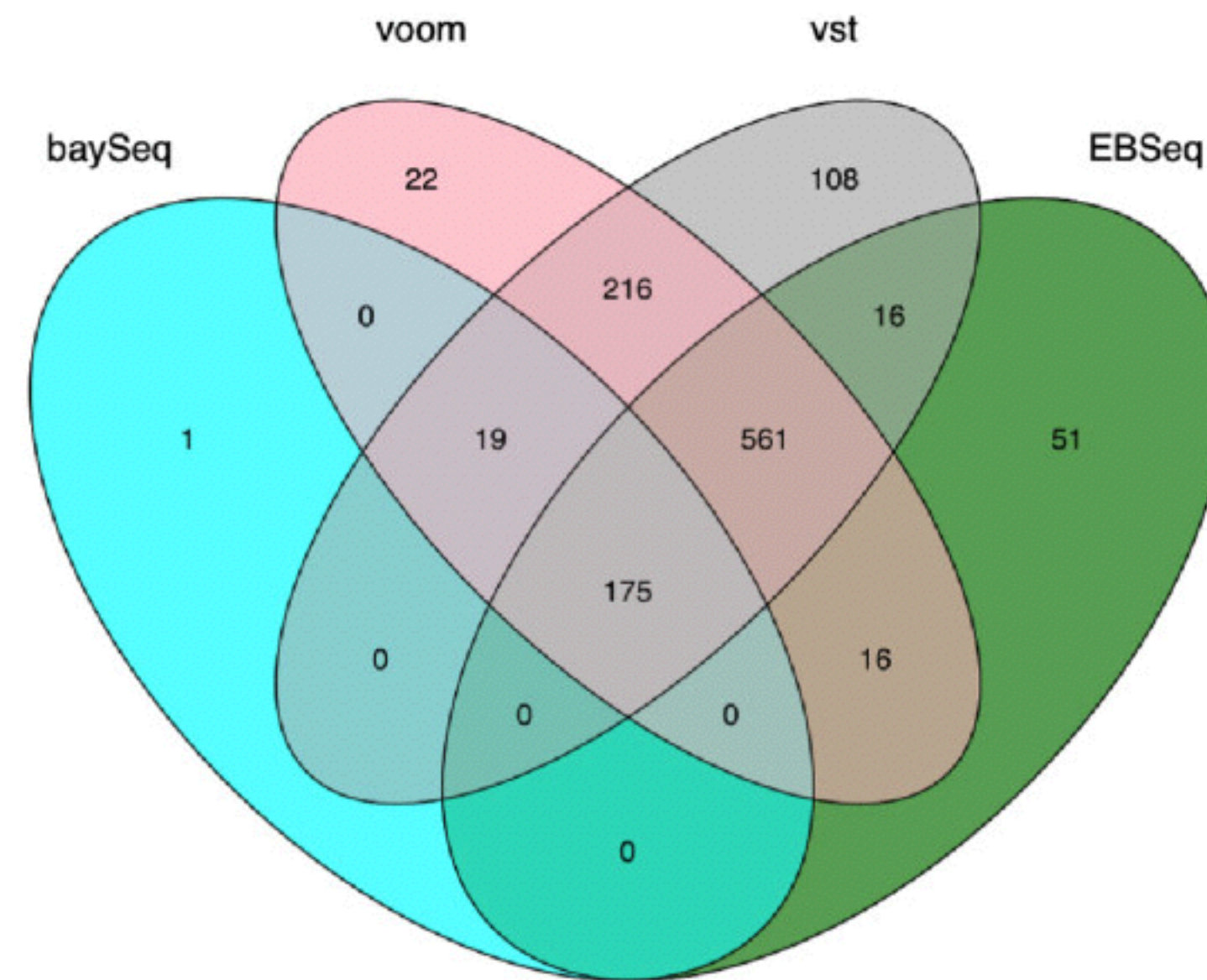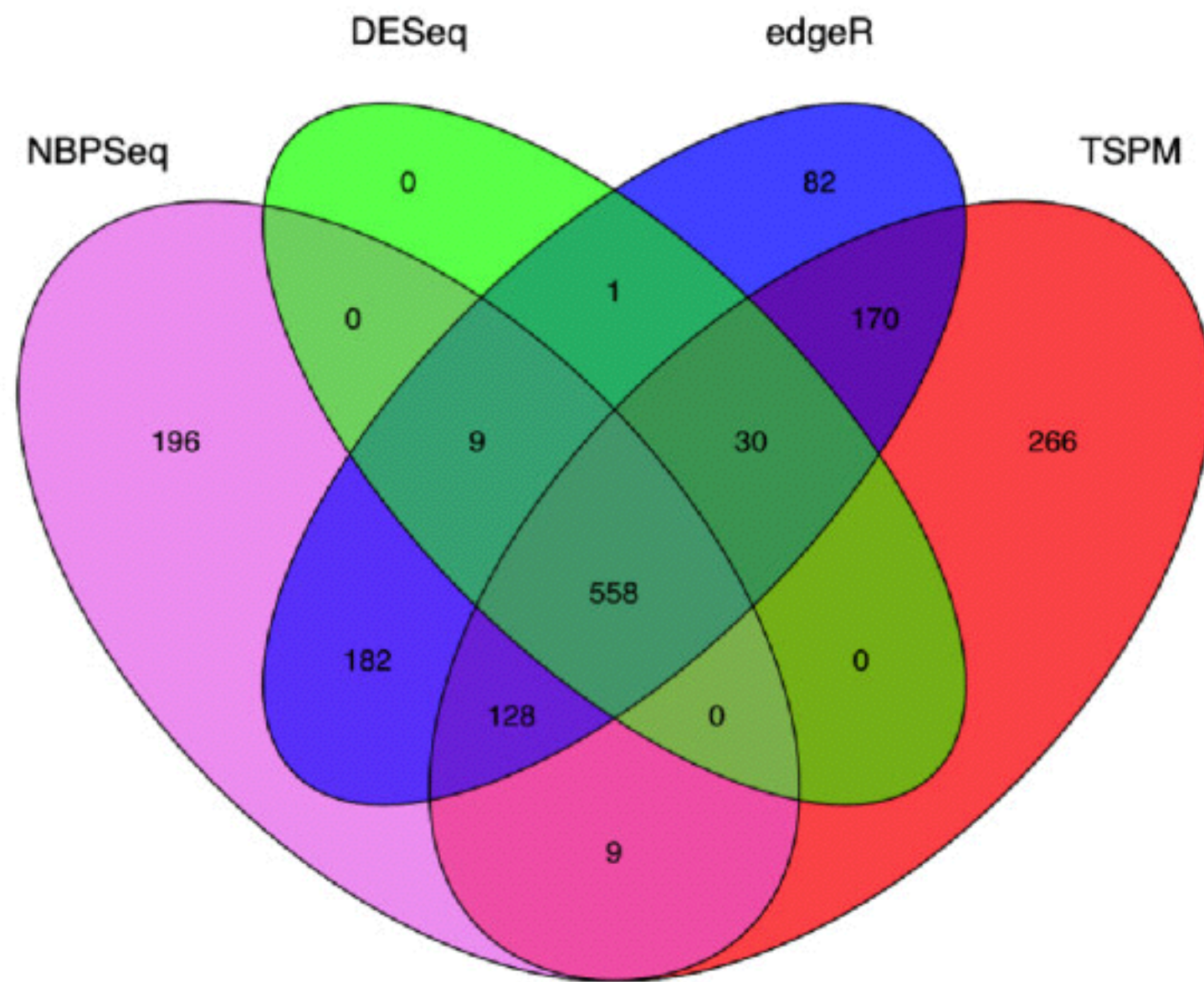 sandboxed

Data
 coherent
 versioned

5

# Ease of use

▶ Tools come pre-configured

▶ Analysis involves

  ▶ Putting FASTQ/BAM files in a directory

  ▶ Creating a CSV metadata file describing the samples

  ▶ Editing a small configuration file

```
samplename,description,panel
SRR950078,UHRR_rep1,UHRR
SRR950079,HBRR_rep1,HBRR
SRR950080,UHRR_rep2,UHRR
SRR950081,HBRR_rep2,HBRR
SRR950082,UHRR_rep3,UHRR
SRR950083,HBRR_rep3,HBRR
SRR950084,UHRR_rep4,UHRR
SRR950085,HBRR_rep4,HBRR
SRR950086,UHRR_rep5,UHRR
SRR950087,HBRR_rep5,HBRR
```

```
details:
  - analysis: RNA-seq
    genome_build: GRCh37
    algorithm:
      aligner: star
      quality_format: Standard
      trim_reads: read_through
      adapters: [truseq, polya]
      strandedness: unstranded
```

RNA-seq pipeline overview

**A comparison of methods for differential expression analysis of RNA-seq data**
Charlotte Soneson[1][*] and Mauro Delorenzi[1][2]

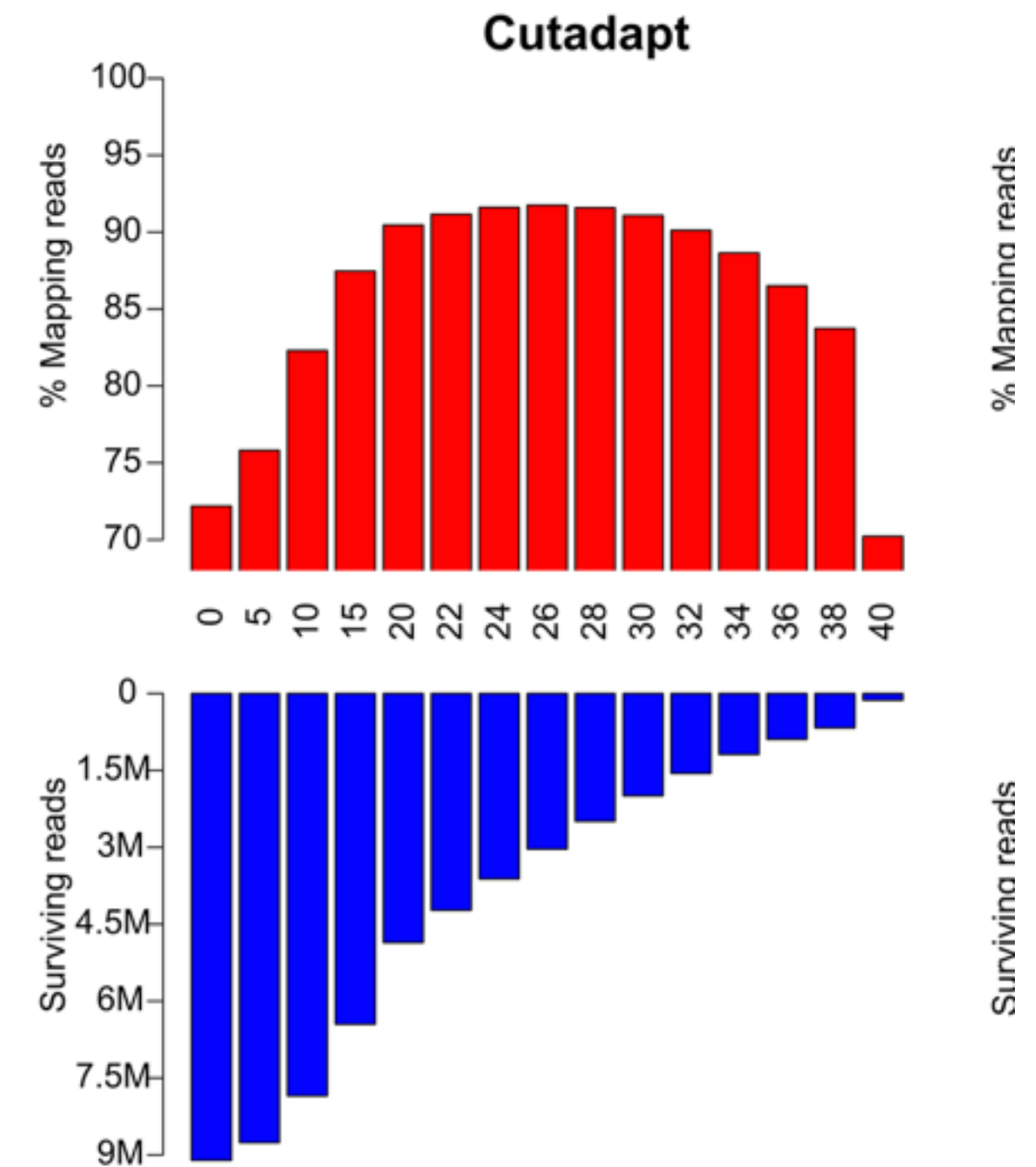# Varying DE calls between methods

# Simulation

▶ SEQC data set not a great set

▶ Count based simulation

  ▶ More complicated models

  ▶ Model biological variability

▶ Which algorithm is best?

▶ Plug in and go

# Is trimming beneficial in RNA-seq?



**Cutadapt**

An Extensive Evaluation of Read Trimming Effects on Illumina NGS Data Analysis

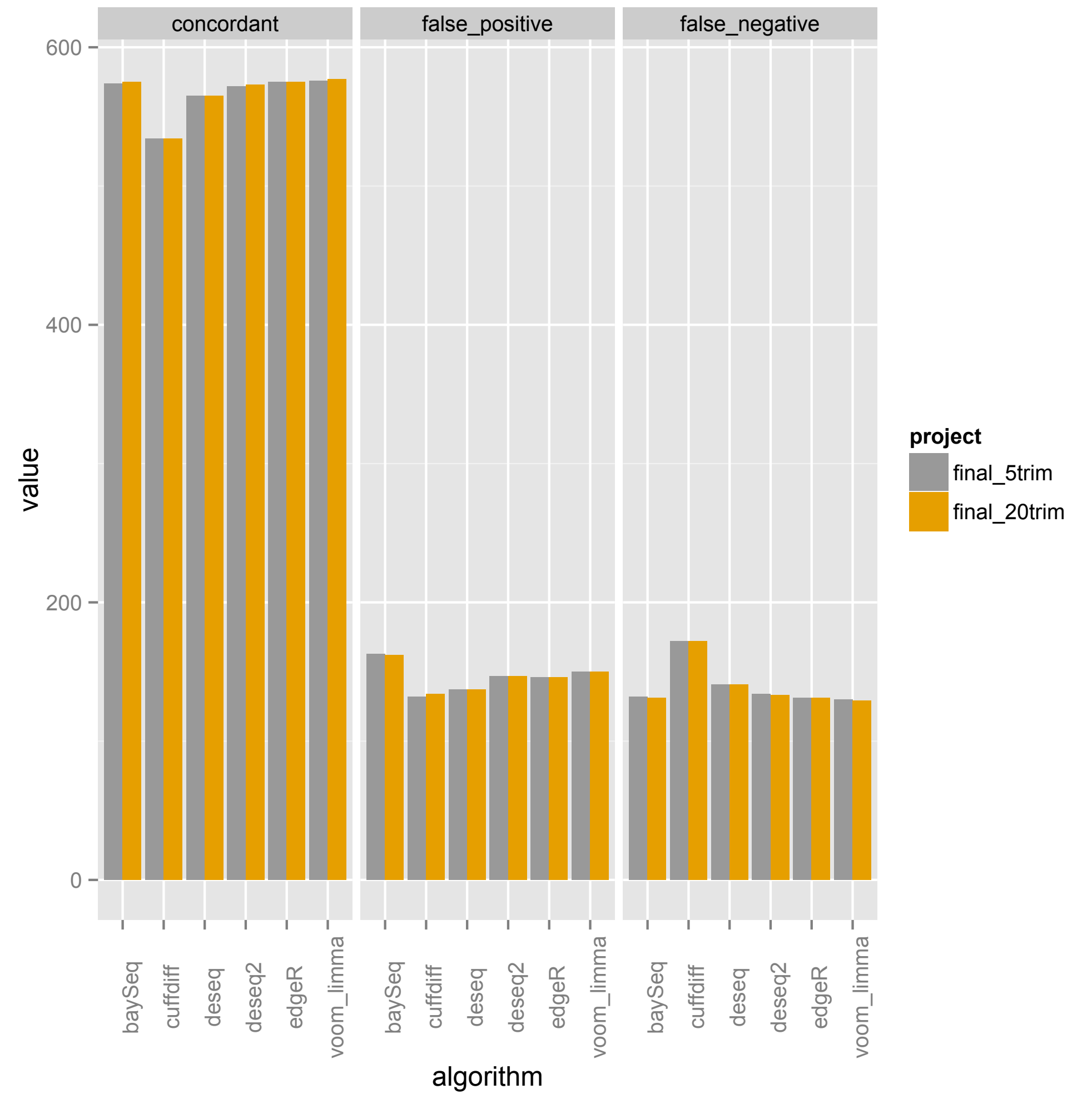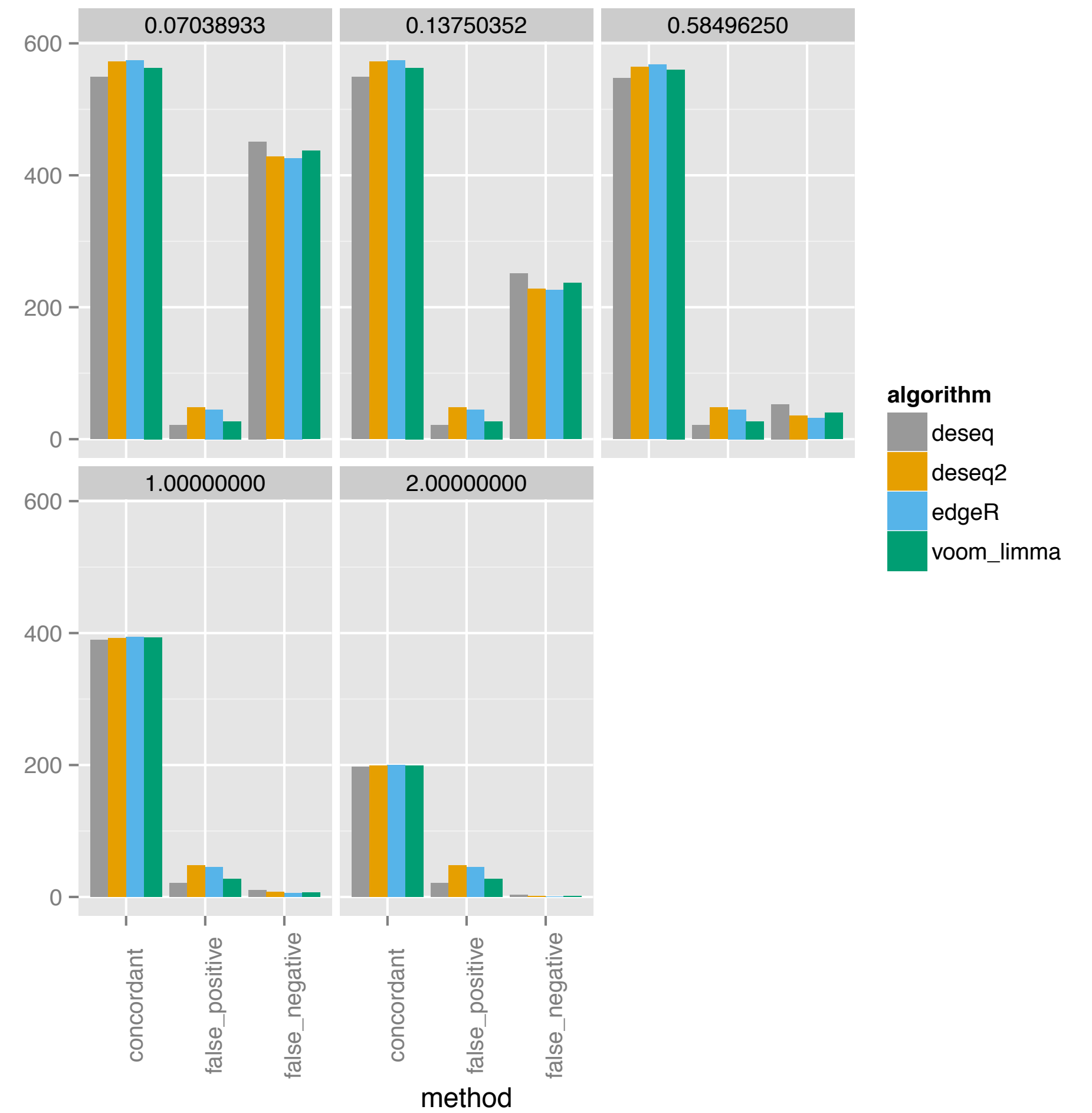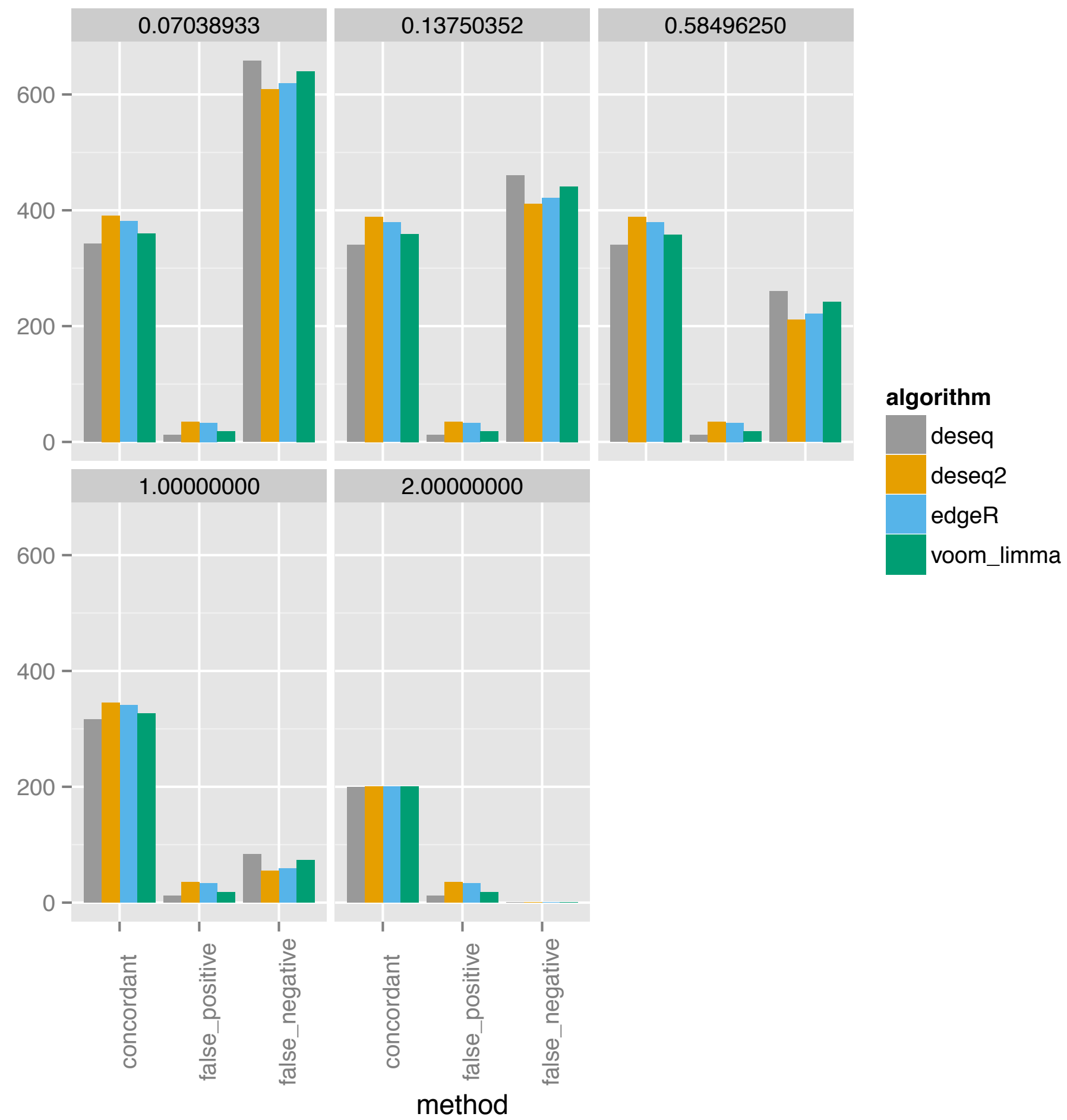Cristian Del Fabbro, Simone Scalabrin, Michele Morgante, Federico M. Giorgi

# Concordance

concordant/false positive/false negative

Jaccard index

Fold change

3 replicates, 100M

15 replicates, 20M

# Get, install, develop

**Get**

wget https://raw.github.com/chapmanb/bcbio-nextgen/master/scripts/bcbio_nextgen_install.py

**Install**

python bcbio_nextgen_install.py /usr/local/share/bcbio-nextgen —tooldir=/usr/local

**Develop**

**https://github.com/chapmanb/bcbio-nextgen**     (Python)

**https://github.com/roryk/bcbio.rnaseq**          (Clojure, R)
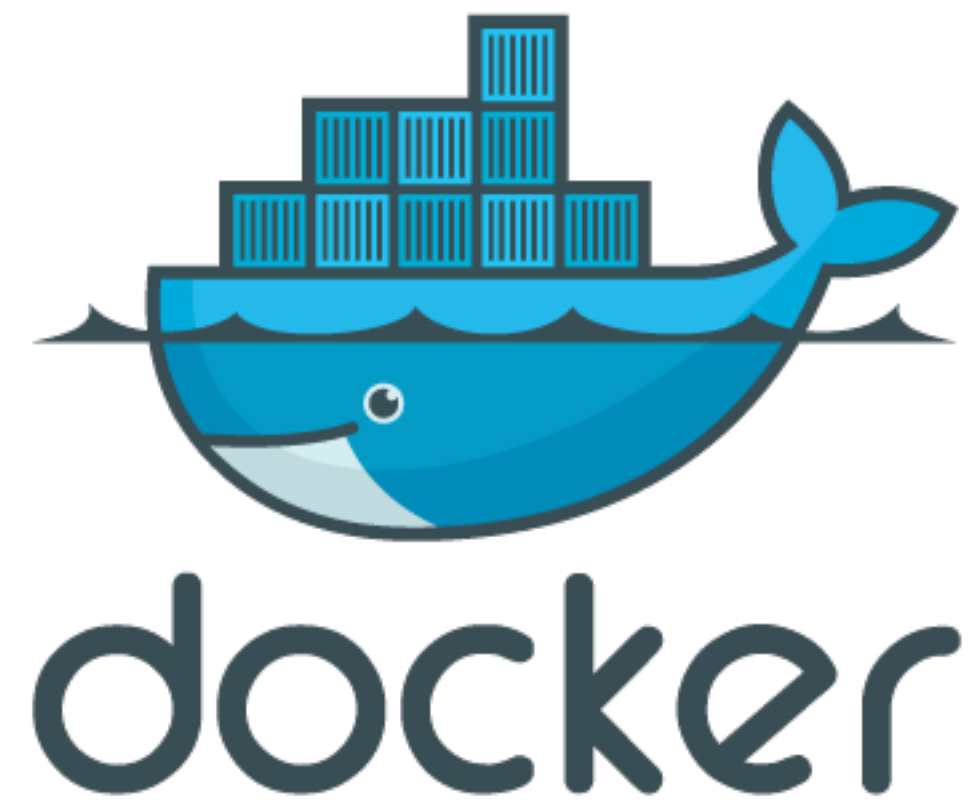
# Current target environment

- Cluster scheduler
  - Torque
  - SLURM
  - SGE
  - LSF
- Shared filesystem
  - NSF
  - Lustre
- Local temporary disk
  - SSD

14

# Virtualization and reproducibility

# Differential expression callers

edgeR

DESeq

DESeq2

BaySeq

voom + limma

Cuffdiff

NOISeq*

DERFinder*

```
# deseq analysis
# Soneson, C. & Delorenzi, M. A comparison of methods for differential expression
# analysis of RNA-seq data. BMC Bioinformatics 14, 91 (2013).

library(DESeq)
library(limma)
library(HTSFilter)
library(tools)
count_file = {{{count-file}}}
out_file = {{{out-file}}}
class = {{{class}}}
project = {{{project}}}
normalized_file = paste(strsplit(out_file, file_ext(out_file)[[1]][[1]]),
    "counts", sep="")
counts = read.table(count_file, header=TRUE, row.names="id")
DESeq.cds = newCountDataSet(countData = counts, conditions = class)
DESeq.cds = estimateSizeFactors(DESeq.cds)
DESeq.cds = estimateDispersions(DESeq.cds, method = "per-condition",
                                fitType = "local")
#DESeq.cds <- HTSFilter(DESeq.cds, s.len=25)$filteredData
res = nbinomTest(DESeq.cds, levels(class)[1], levels(class)[2])

comparison = paste(levels(class)[1], "_vs_", levels(class)[2], sep="")
out_table = data.frame(id=res$id, expr=res$baseMean, logFC=res$log2FoldChange,
          pval=res$pval, padj=res$padj, algorithm="deseq", project=project)
out_table$pval[is.na(out_table$pval)] = 1
out_table$padj[is.na(out_table$padj)] = 1
write.table(out_table, file=out_file, quote=FALSE, row.names=FALSE,
            sep="\t")
write.table(counts(DESeq.cds, normalized=TRUE), file=normalized_file,
            quote=FALSE, sep="\t")
```

# Community

Nick Loman
@pathogenomenick

Follow

Loman's law of bioinformatics: If you haven't found at least one bug in someone's pipeline then you don't understand it properly yet.

# Community