



# Technology is fickle



**Illumina, Inc.**  
5200 Illumina Way  
San Diego, CA 92122 USA  
tel 858.202.4500  
fax 858.202.4545  
[www.illumina.com](http://www.illumina.com)

24 March 2015

## Product Obsolescence Notification

Dear Valued Customer,

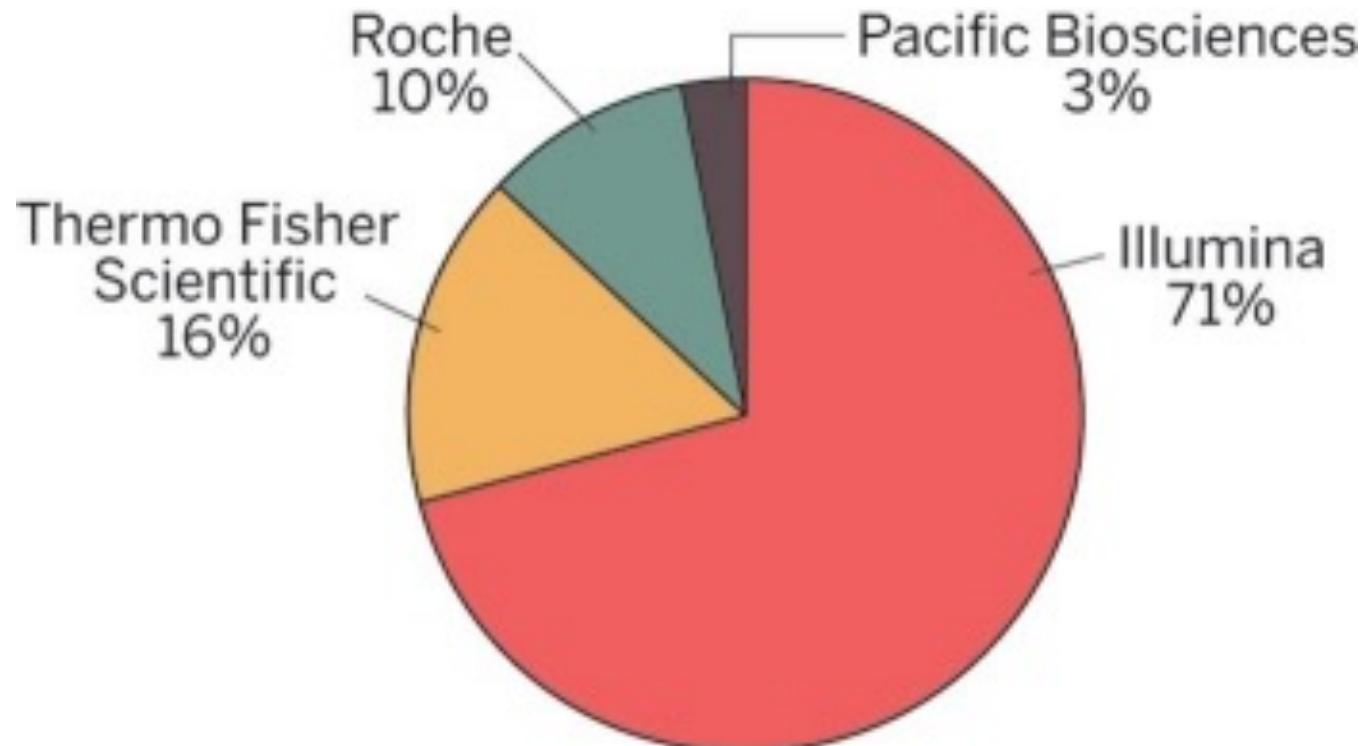
Due to unexpected demand for Illumina's Mouse Gene Expression BeadArray™, the inventory of beads for this product has been completely consumed. We are unable to synthesize exact replacement beadpools due to the nature and complexity of the product. As a result we are forced to discontinue and no longer accept orders for the products listed below.

You are receiving this notice because our records indicate that you have ordered one or more of these products or have orders pending.

### Products Affected

Product Description	Catalog Number
MouseRef-8 v2.0 BeadChip (16 Samples)	BD-202-0202
MouseRef-8 v2.0 BeadChip (48 Samples)	BD-202-0602
MouseWG-6 v2.0 BeadChip (12 Samples)	BD-201-0202
MouseWG-6 v2.0 BeadChip (36 Samples)	BD-201-0602

# The sequencing landscape is dominated by Illumina



World market in 2013 = \$1.3 billion

# The year of sequencing

In 2007, the next-generation sequencing technologies have come into their own with an impressive array of successful applications. Kelly Rae Chi reports.

## SPECIAL FEATURE | METHOD OF THE YEAR

### CREATING THE GENOME ANALYZER

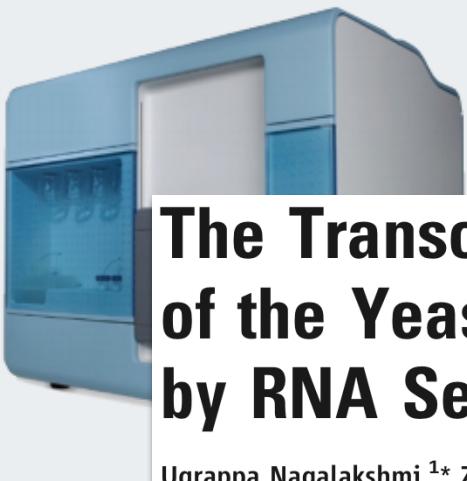
When John West started as CEO of Solexa Ltd. in August of 2004, the longest stretch of DNA that the company could sequence was only six bases long. "That was a little bit intimidating," recalls West, now the vice president and general manager of Illumina's DNA sequencing business unit following the acquisition of Solexa Ltd. by Illumina. "The problem was we never had a commercial platform to be able to sequence it on."

The next-generation sequencing platform they developed and commercialized by June 2006, West says, has been a huge success.

Since the commercial release of the platform they have sold 100 instruments and increased the scale of what they have tackled using the technology—from a 5,300-base-pair viral genome to a 150-million-base-pair human X chromosome. But the machine was a challenge to develop. The developers had to bring together key elements of chemistry, people and technology to make it work.

By the time Solexa Ltd. announced its plans to merge with Lynx Therapeutics in August 2004, a lot of the core sequencing-by-synthesis chemistry and molecular biology had already been done, West says. The early chem  
work, done by Solexa Ltd. in the United Kingdom, led to t  
creation of a reversible terminator nucleotide and a polym  
that would incorporate it. Solexa Ltd. and Lynx Therapeut  
had bought cluster technology for solid phase amplificatio  
together from the Swiss company Manteia SA, and did som  
instrumentation design.

Still, the company was in a state of flux in late 2004 and needed to figure out how to combine the chemistry and technology into a complete system. The researchers needed to meet and brainstorm, operating on an eight-hour time difference between the two branches, one in the UK and one in the US. Ligation chemistry developed by Lynx Therapeutics was an option. There were differing opinions about which chemistry, ligation or polymerase would work best, but ultimately they took a gamble on the Solexa polymerase because it could



## The Transcriptional Landscape of the Yeast Genome Defined by RNA Sequencing

Ugrappa Nagalakshmi,<sup>1\*</sup> Zhong Wang,<sup>1\*</sup> Karl Waern,<sup>1</sup> Chong Shou,<sup>2</sup> Debasish Raha,<sup>1</sup> Mark Gerstein,<sup>2,3</sup> Michael Snyder<sup>1,2,3†</sup>

The identification of untranslated regions, introns, and coding regions within an organism remains challenging. We developed a quantitative sequencing-based method called RNA-Seq for mapping transcribed regions, in which complementary DNA fragments are subjected to high-throughput sequencing and mapped to the genome. We applied RNA-Seq to generate a high-resolution transcriptome map of the yeast genome and demonstrated that most (74.5%) of the nonrepetitive sequence of the yeast genome is transcribed. We confirmed many known and predicted introns and demonstrated that others are not actively used. Alternative initiation codons and upstream open reading frames also were identified for many yeast genes. We also found unexpected 3'-end heterogeneity and the presence of many overlapping genes. These results indicate that the yeast transcriptome is more complex than previously appreciated.

## RNA Transcription

ChIRP-Seq



CHART



RAP



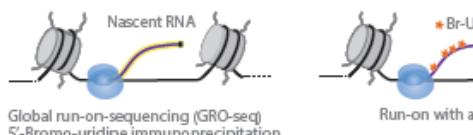
GRO-seq

BRIC-Seq

Bru-Seq

BruChase

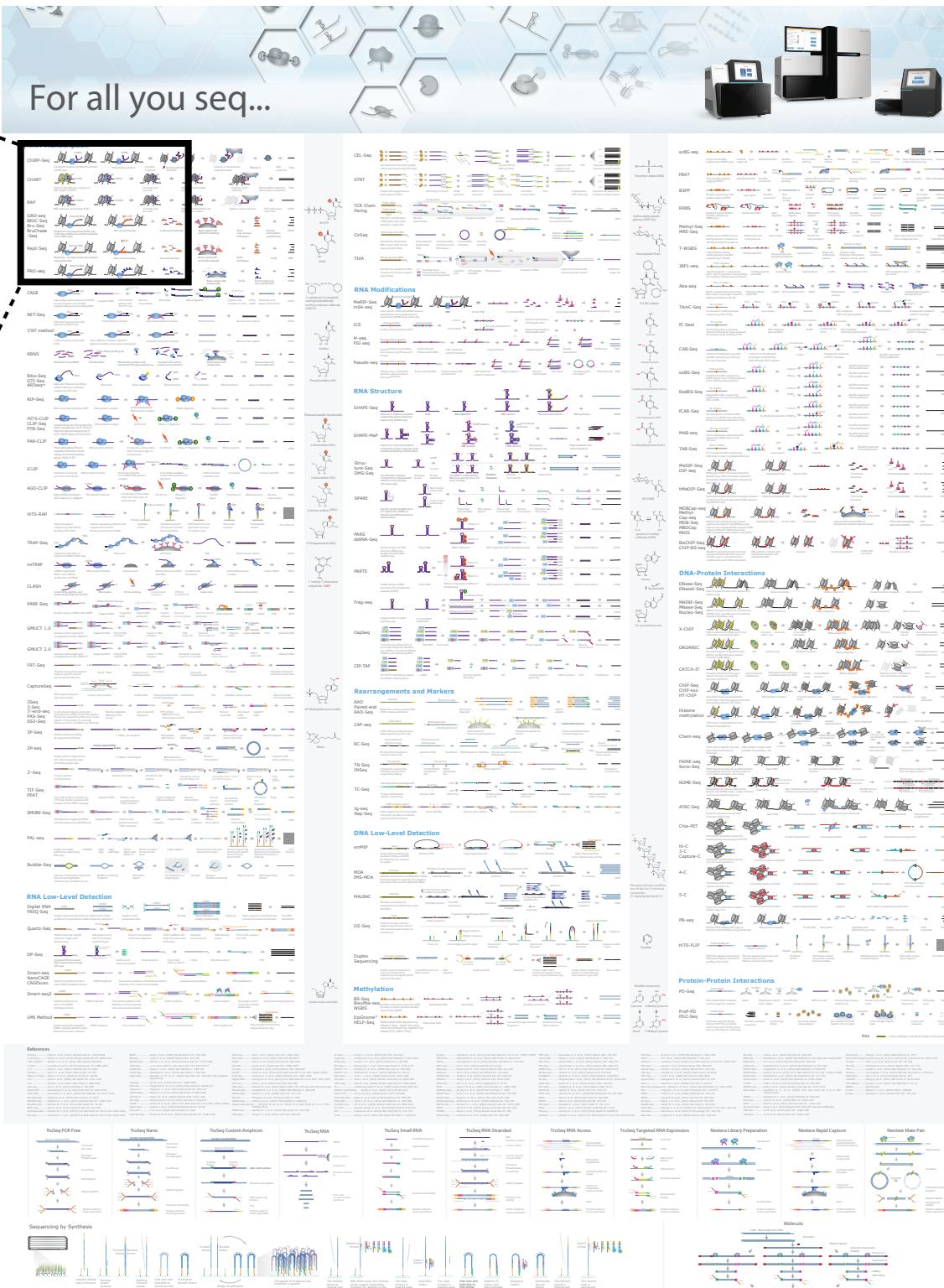
-Seq



Repli Seq



PRO-seq



<http://bit.ly/1NBxOJH>

illumina®

# Illumina RNAseq

MiSeq



14-20 million reads  
250bp paired-end  
24 hr run time

- well-suited for a lab, but not sufficient for transcriptomics
- Best suited for:
  - ✓ amplicon sequencing (16S)
  - ✓ bacterial genome seq
  - ✓ QC of RNAseq libraries

NextSeq



400 million reads  
150bp paired-end  
24 hr run time

HiSeq

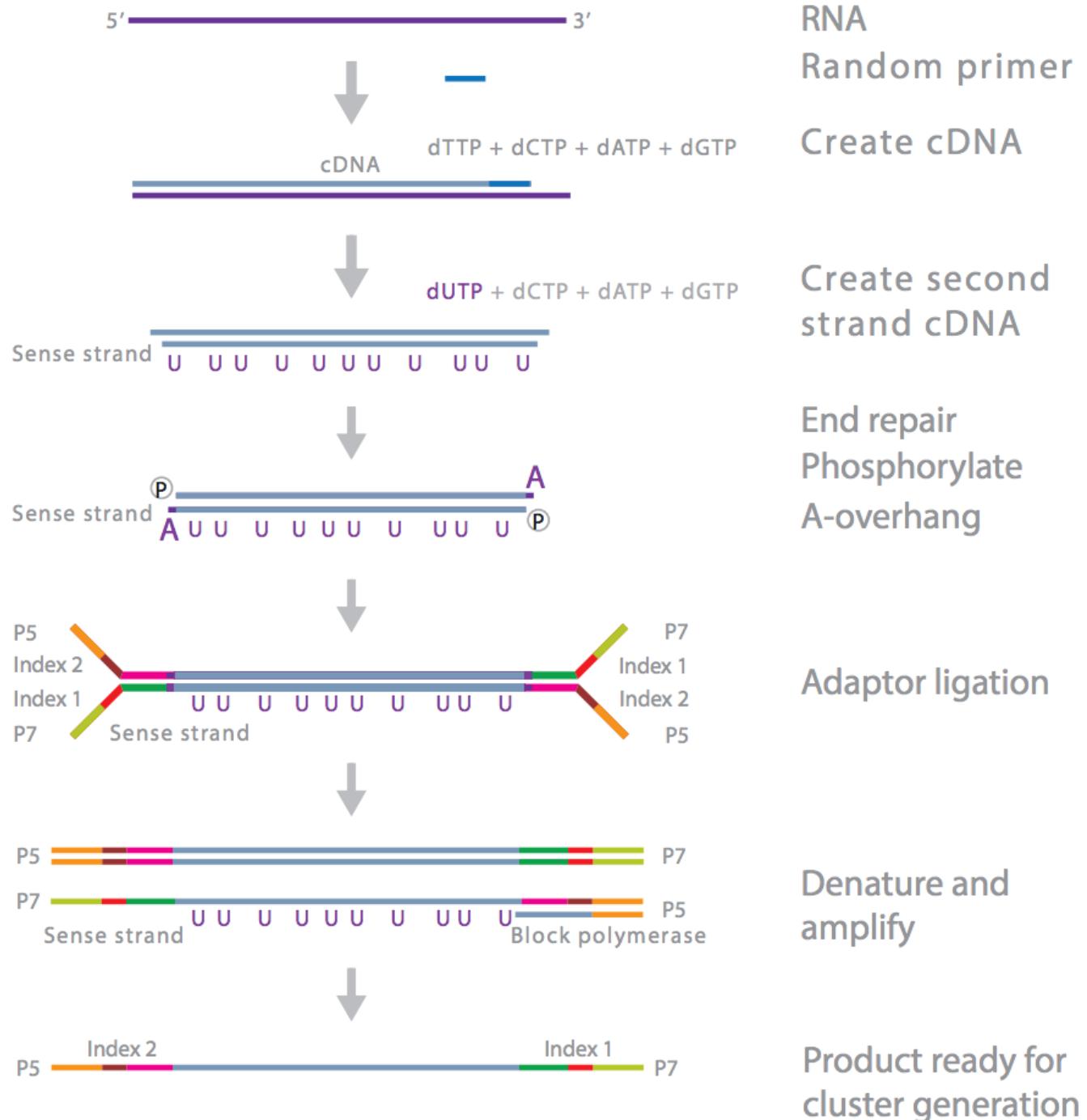


0.3 - 2 billion reads  
150-250bp paired-end  
5 day run time

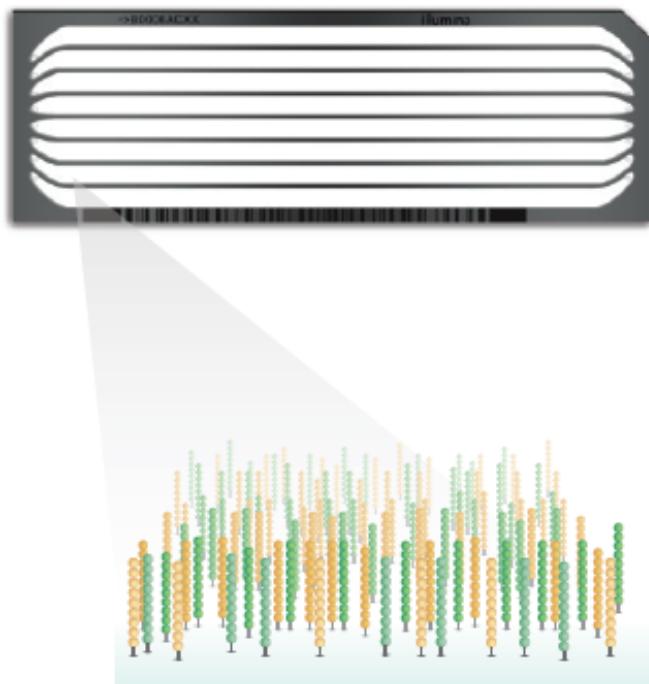
- best fit for a large core facility
- requires a dedicated technician
- complicated set-up
- Best suited for:
  - ✓ RNAseq and DNAseq
  - ✓ metagenomics

All sequencing starts with a  
‘library prep’

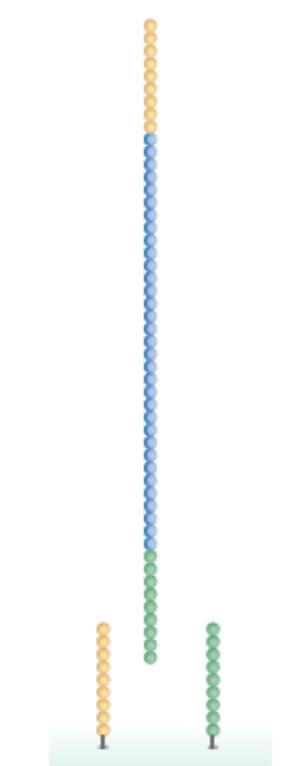
# TruSeq RNA Stranded



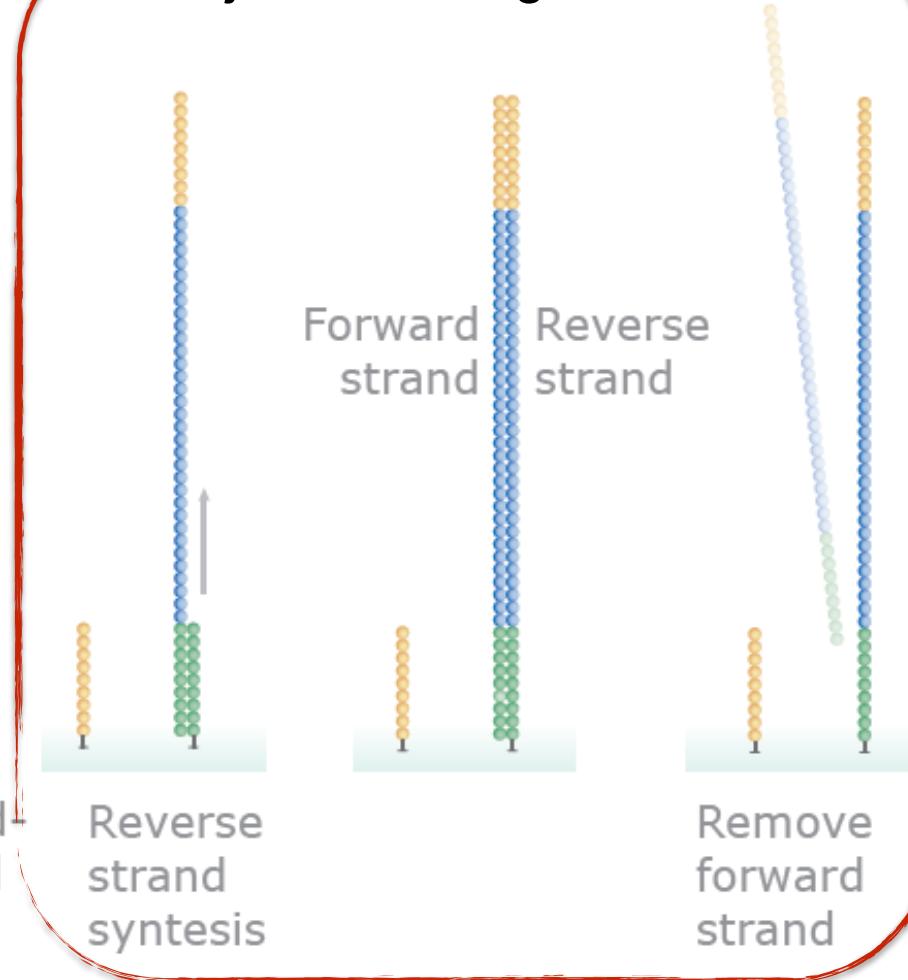
# Sequencing by Synthesis



Adapter hybridizes to flowcell



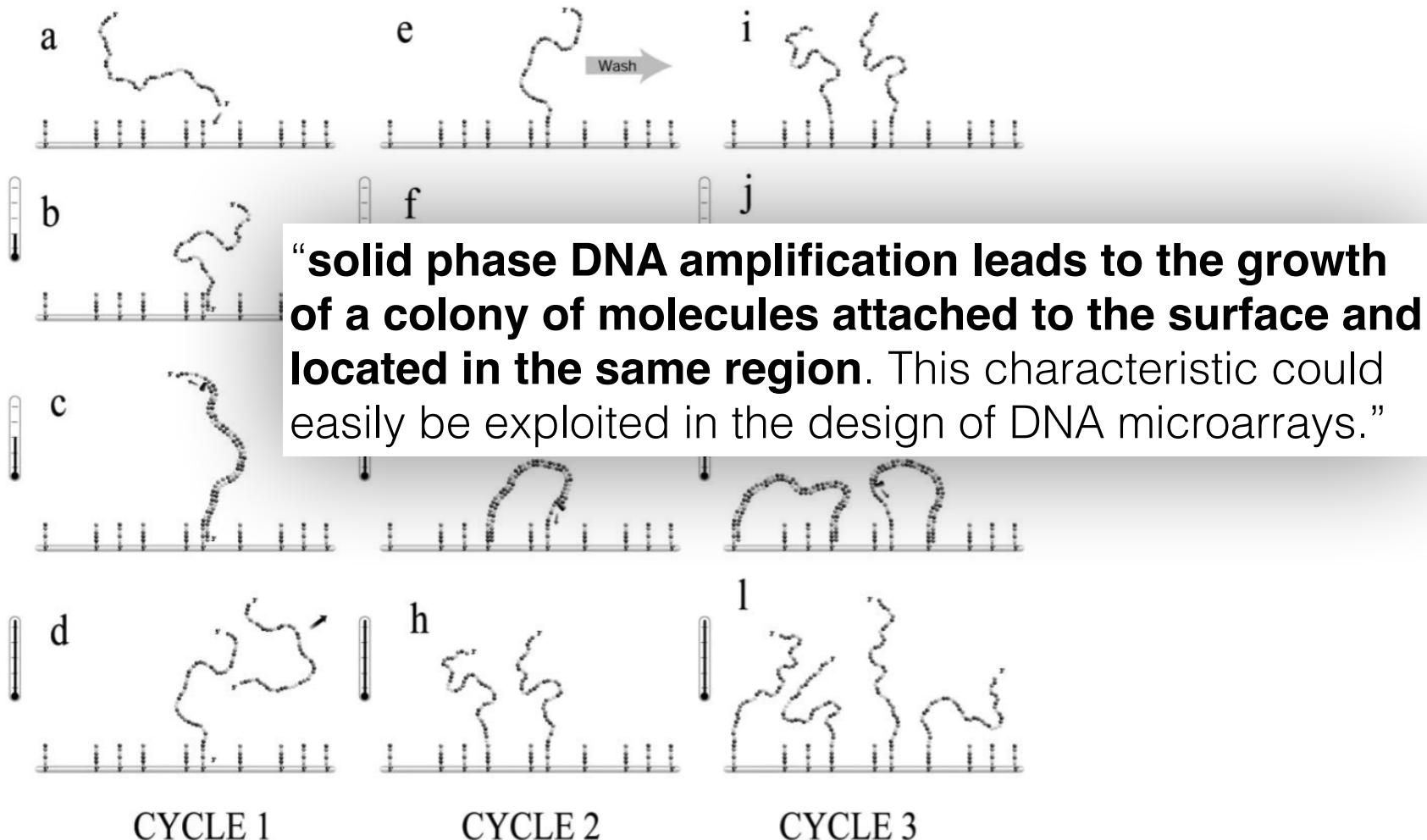
major technological hurdle



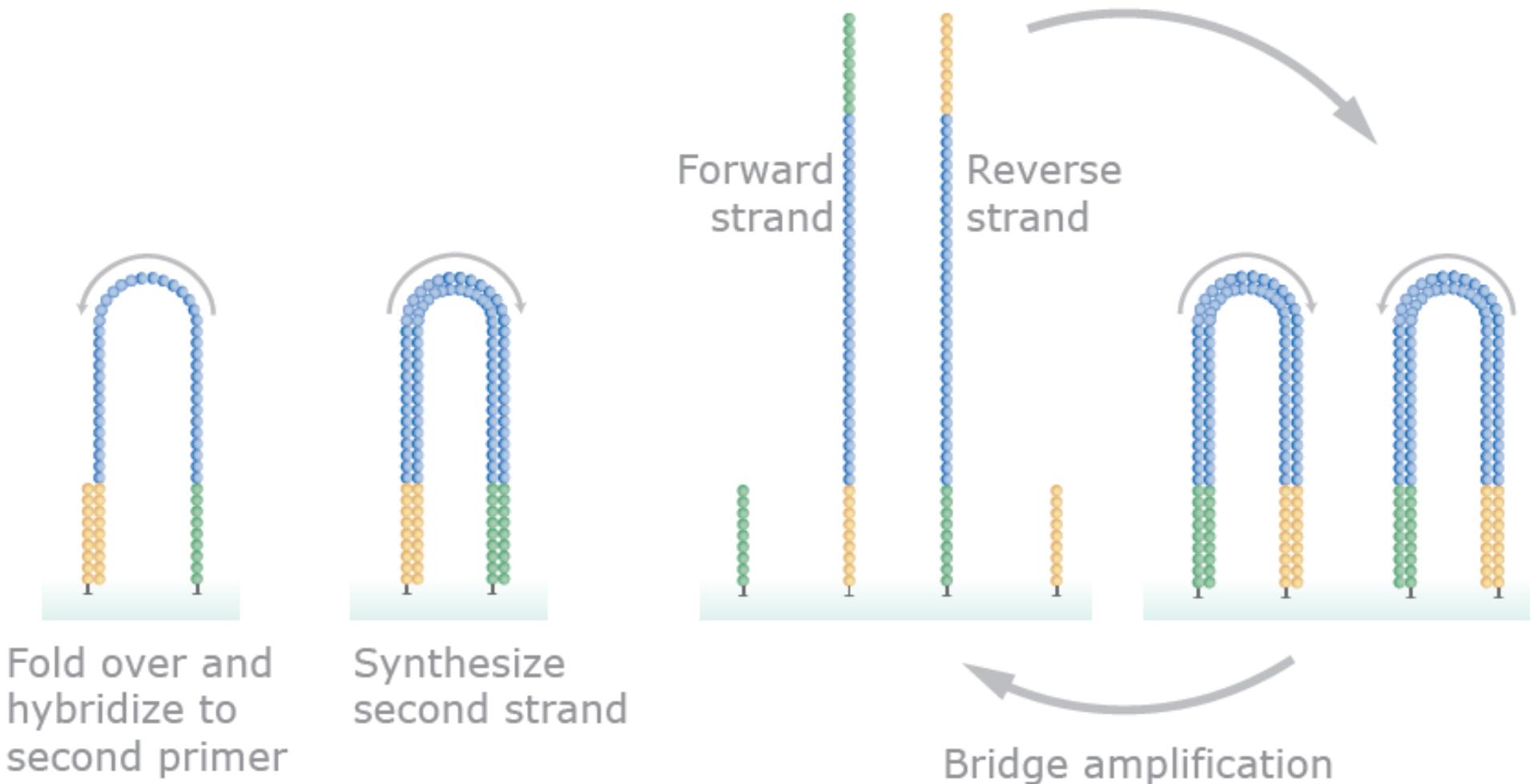
# Solid Phase DNA Amplification: A Simple Monte Carlo Lattice Model

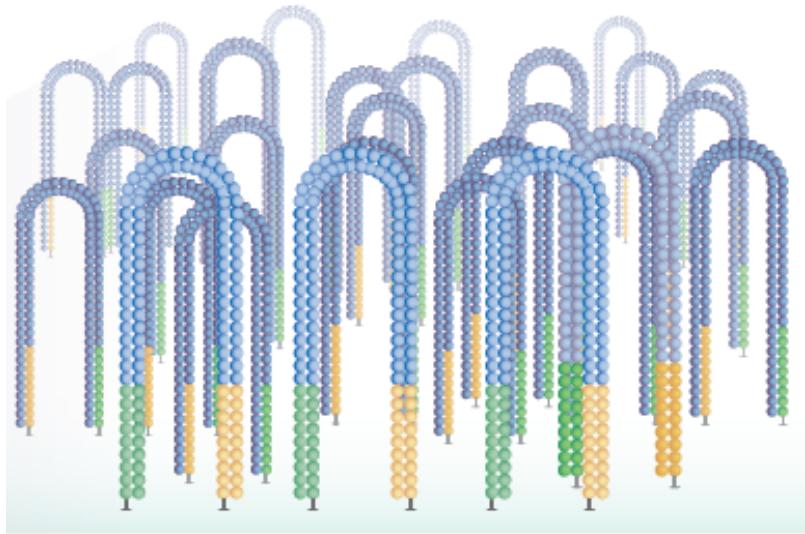
Jean-Francois Mercier,\* Gary W. Slater,\* and Pascal Mayer<sup>†</sup>

\*Department of Physics, University of Ottawa, Ottawa, Ontario, Canada; and <sup>†</sup>Manteia Predictive Medicine S.A., Coinsins, Switzerland

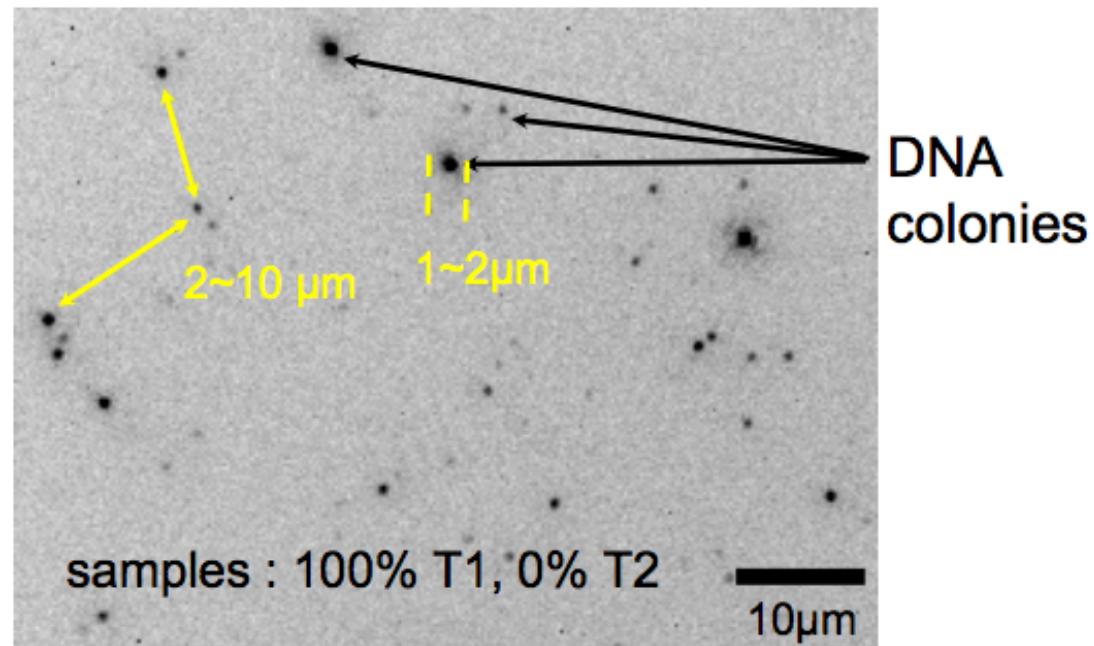


# Sequencing by Synthesis (SBS)

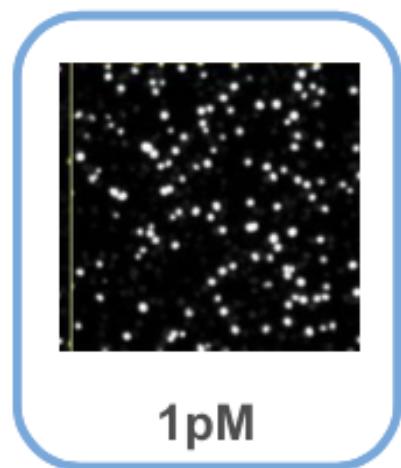
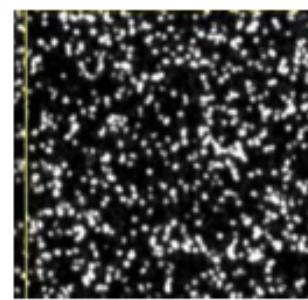
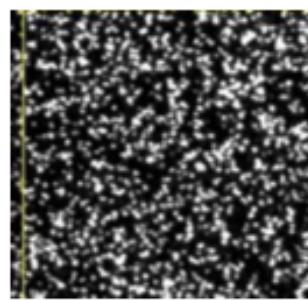
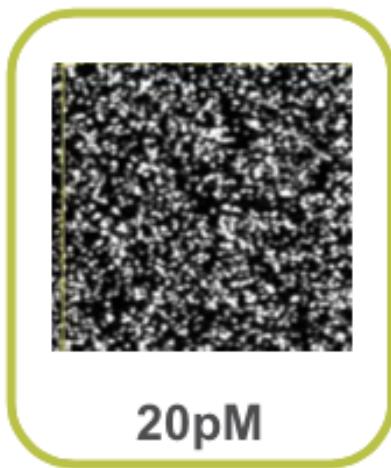




Thousands of molecules are amplified in parallel



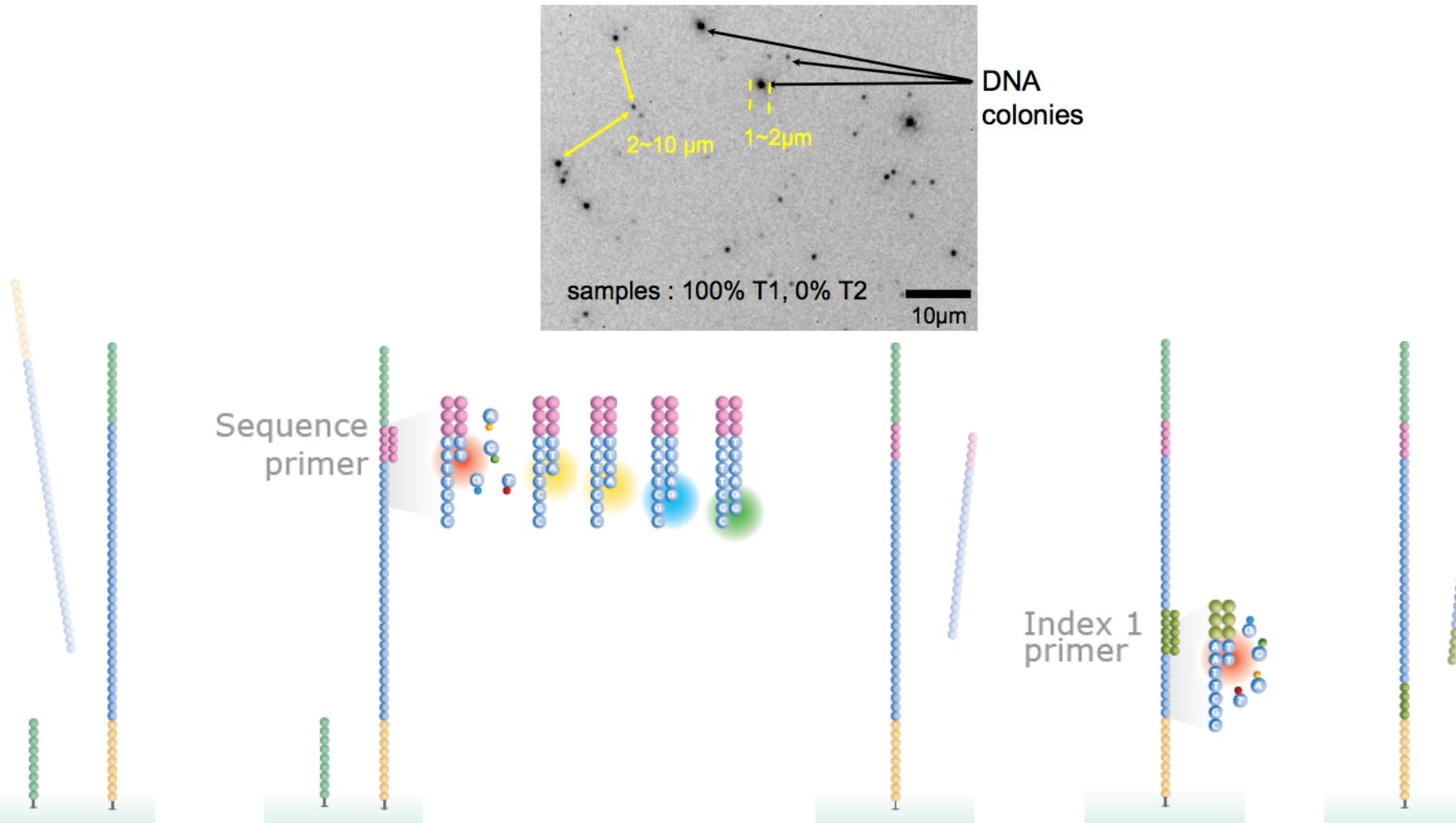
Mayer et al., Presentation 1998



\*HiSeq: 600-800K clusters/mm<sup>2</sup>

\*NextSeq: 130-160K clusters/mm<sup>2</sup>

# Sequencing by Synthesis (SBS)



The reverse strand is cleaved and washed away

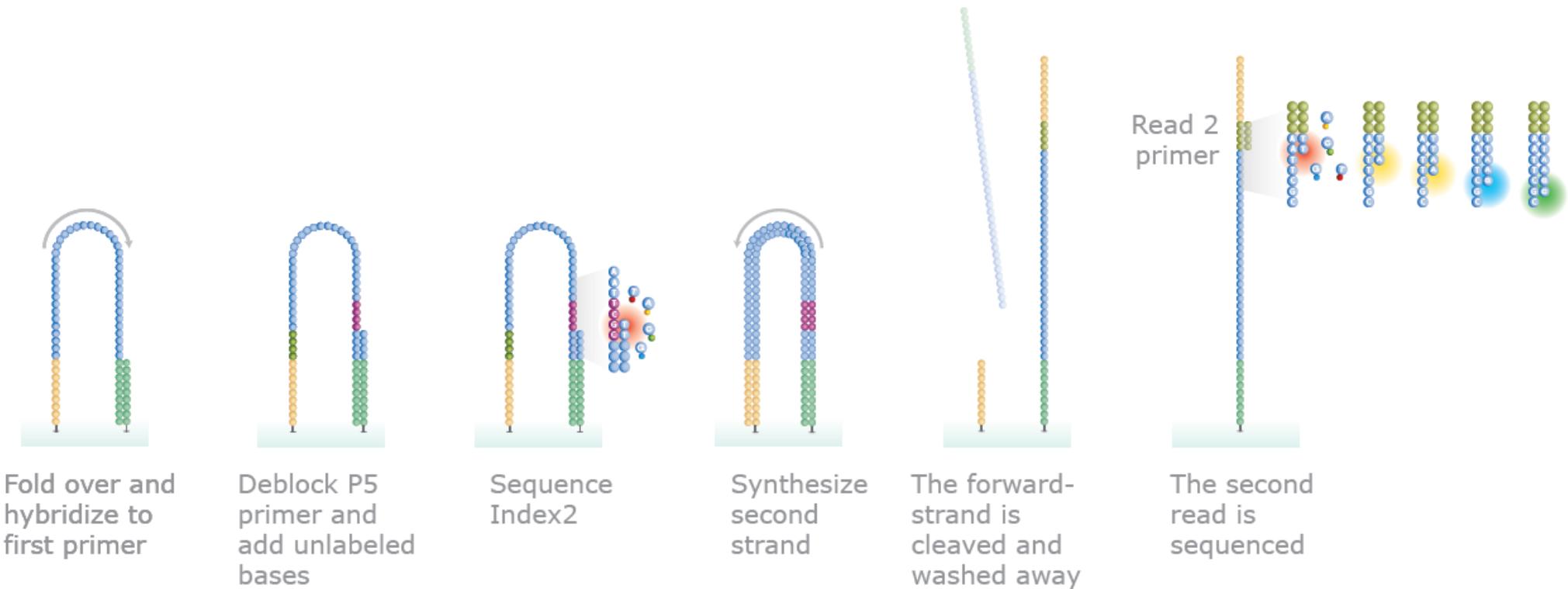
With each cycle, four fluorescently tagged nucleotides compete for addition to the growing chain. Only one is incorporated based on the sequence of the template.

The read product is washed away

Sequence Index1

The read product is washed away

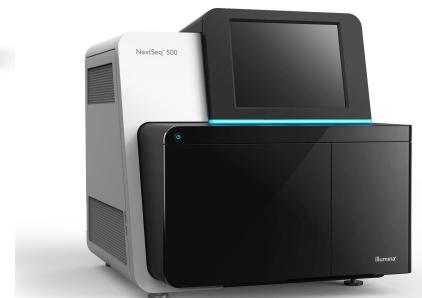
# Sequencing by Synthesis (SBS)



<https://www.youtube.com/watch?v=HMyCqWhwB8E>

# experimental design

1. What is the question that you want to address with sequencing?
2. Do you have independent evaluation that your experiment worked?
3. What are the appropriate controls?
4. What are the questions that will drive your data analysis?
5. How many replicates are needed?
6. How deep should I sequence (reads)?
7. How much will it cost?



# Deeper sequencing or more replicates?

BIOINFORMATICS

DISCOVERY NOTE

Vol. 30 no. 3 2014, pages 301–304  
doi:10.1093/bioinformatics/btt688

Gene expression

Advance Access publication December 6, 2013

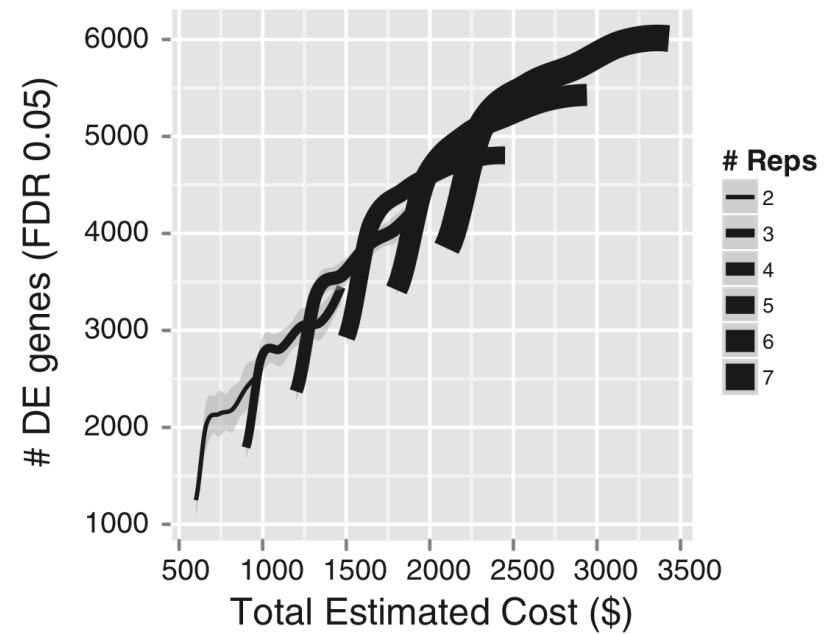
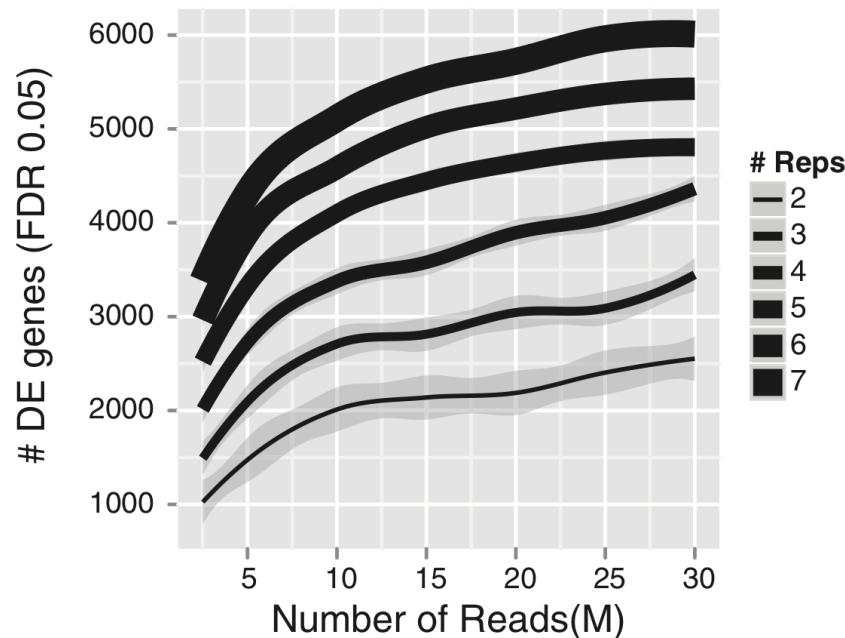
## RNA-seq differential expression studies: more sequence or more replication?

Yuwén Liu<sup>1,2</sup>, Jie Zhou<sup>1,3</sup> and Kevin P. White<sup>1,2,3,\*</sup>

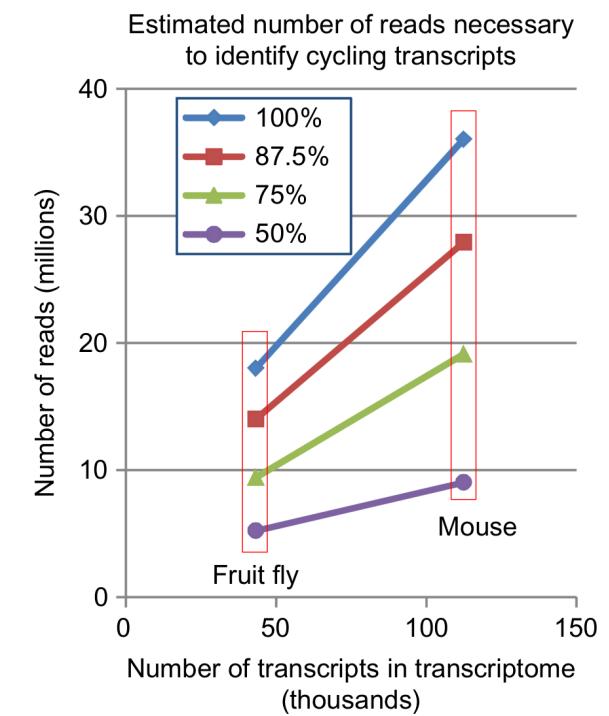
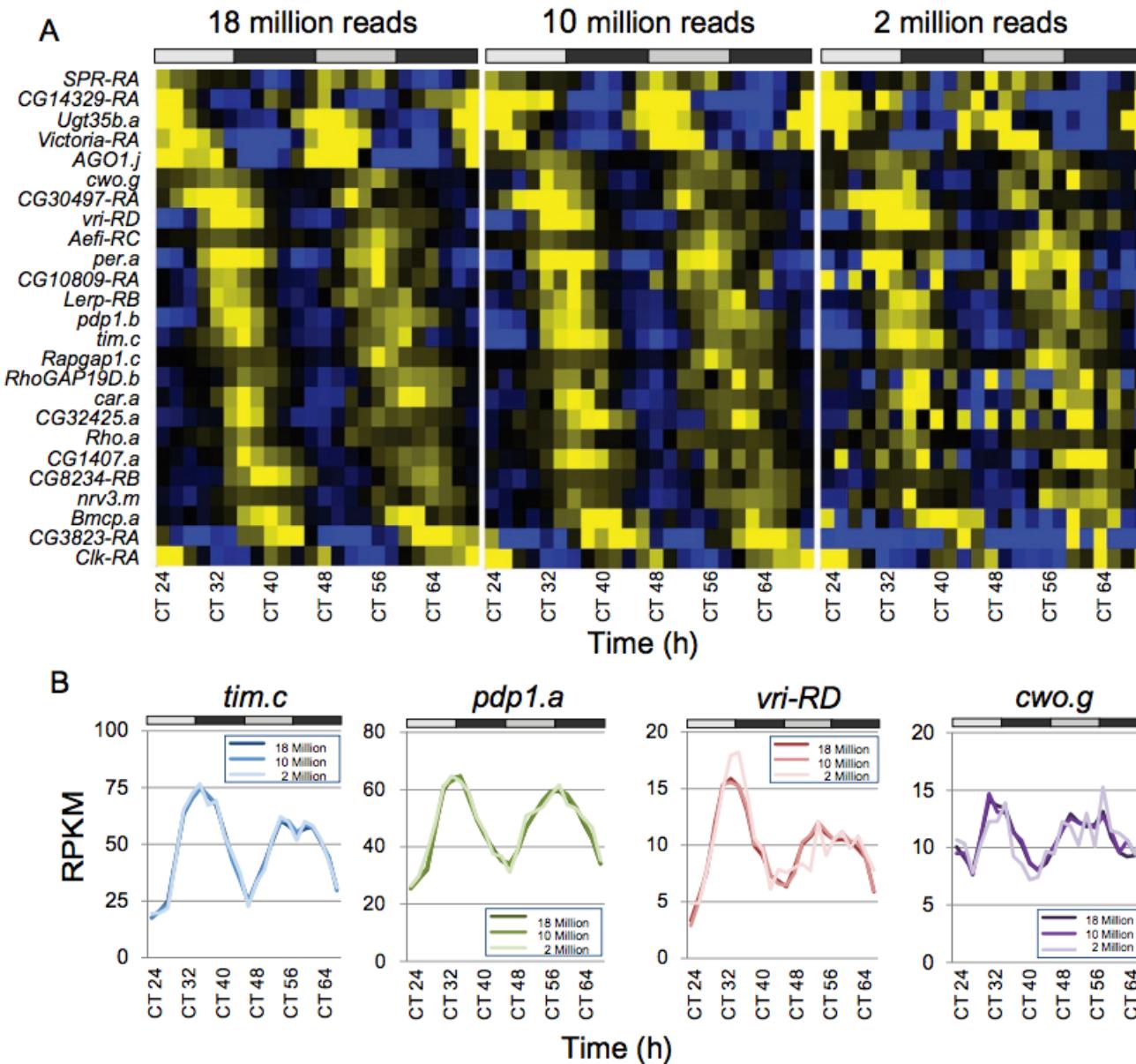
<sup>1</sup>Institute of Genomics and Systems Biology, <sup>2</sup>Committee on Development, Regeneration, and Stem Cell Biology and

<sup>3</sup>Department of Human Genetics, University of Chicago, Chicago, IL 60637, USA

Associate Editor: Janet Kelso



# low read depth can still reveal biology



# experimental design

1. What is the question that you want to address with sequencing?
2. Do you have independent evaluation that your experiment worked?
3. What are the appropriate controls?
4. What are the questions that will drive your data analysis?
5. How many replicates are needed?
6. How deep should I sequence (reads)?
7. How much will it cost?



# Library Preparation



# Sequencing



\*mid output  
130M clusters

\*high output  
400M clusters

TruSeq kit	cat#	cost
<b>stranded total RNA LT (w/ RiboZero)</b>	RS-122-2201	\$5,520
<b>stranded mRNA LT</b>	RS-122-2101	\$2,640

\*all kits process 48 samples

cycles	cat#	cost
300	FC-404-2004	\$4,120
150	FC-404-2002	\$2,575
75	FC-404-2005	\$1,340

# total transcriptome, 150 bp PE

## Library Preparation



## Sequencing



\*mid output  
130M clusters

\*high output  
400M clusters

TruSeq kit	cat#	cost	cycles	cat#	cost
stranded total RNA LT (w/ RiboZero)	RS-122-2201	\$5,520	300	FC-404-2004	\$4,120
stranded mRNA LT	RS-122-2101	\$2,640	150	FC-404-2002	\$2,575
			75	FC-404-2005	\$1,340

\*all kits process 48 samples

**12 samples**

library prep = \$115/sample

sequencing = \$343/sample

data output = 30M reads/sample

# mRNAseq, 150 bp SE

## Library Preparation



## Sequencing



\*mid output  
130M clusters

\*high output  
400M clusters

TruSeq kit	cat#	cost
stranded total RNA LT (w/ RiboZero)	RS-122-2201	\$5,520
stranded mRNA LT	RS-122-2101	\$2,640

cycles	cat#	cost
300	FC-404-2004	\$4,120
150	FC-404-2002	\$2,575
75	FC-404-2005	\$1,340

\*all kits process 48 samples

### 12 samples

library prep = \$55/sample  
sequencing = \$215/sample  
data output = 30M reads/sample

# mRNaseq, 75 bp SE

## Library Preparation



## Sequencing



\*mid output  
130M clusters

\*high output  
400M clusters

TruSeq kit	cat#	cost
stranded total RNA LT (w/ RiboZero)	RS-122-2201	\$5,520
stranded mRNA LT	RS-122-2101	\$2,640

cycles	cat#	cost
300	FC-404-2004	\$4,120
150	FC-404-2002	\$2,575
75	FC-404-2005	\$1,340

\*all kits process 48 samples

**12 samples**

library prep = \$55/sample

sequencing = \$112/sample

data output = 30M reads/sample

# mRNAseq, 75 bp SE (24 samples)

## Library Preparation



## Sequencing



\*mid output  
130M clusters

\*high output  
400M clusters

TruSeq kit	cat#	cost
stranded total RNA LT (w/ RiboZero)	RS-122-2201	\$5,520
stranded mRNA LT	RS-122-2101	\$2,640

cycles	cat#	cost
300	FC-404-2004	\$4,120
150	FC-404-2002	\$2,575
75	FC-404-2005	\$1,340

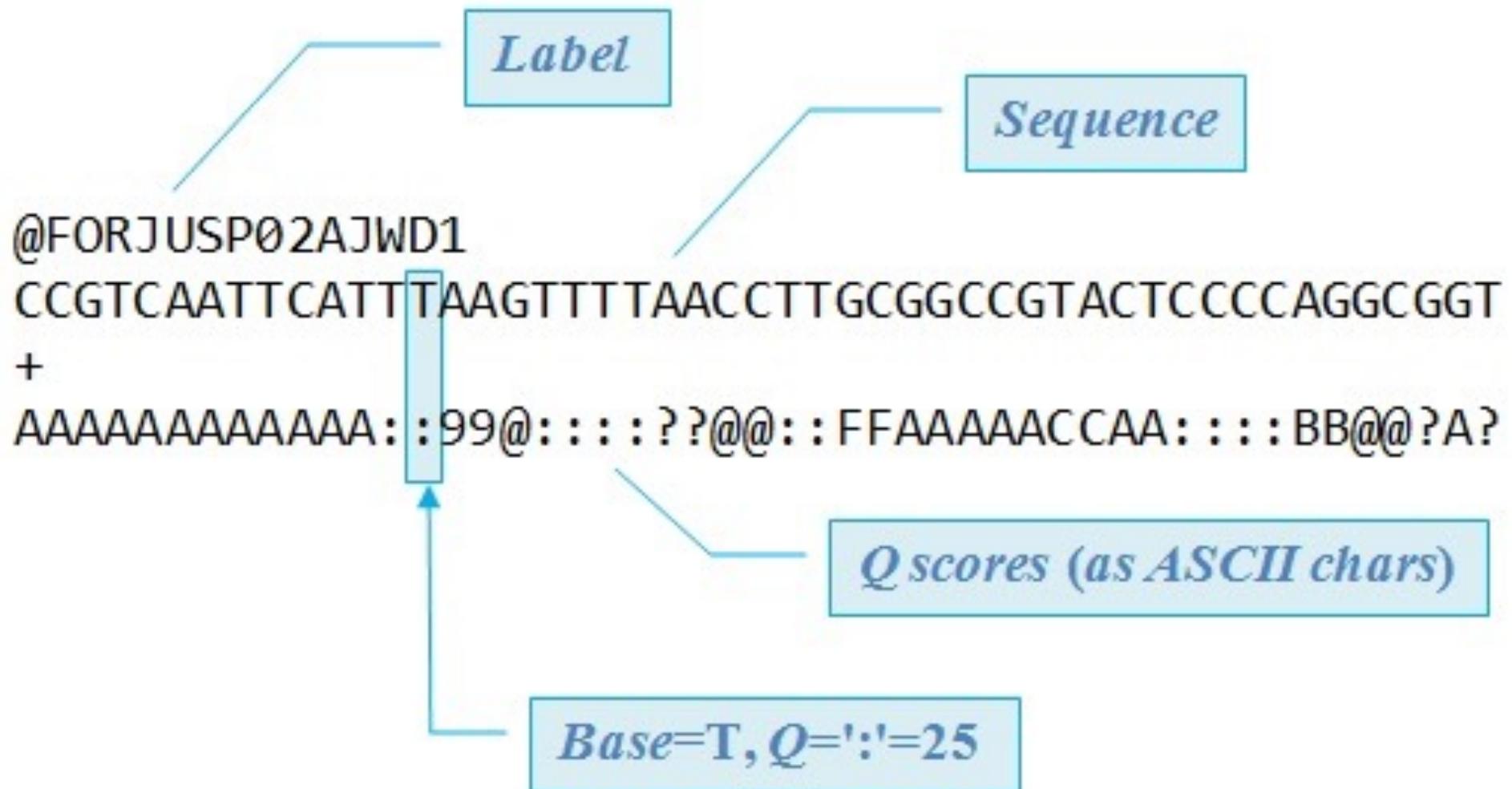
\*all kits process 48 samples

### 24 samples

library prep = \$55/sample  
sequencing = \$56/sample  
data output = 15M reads/sample

# A word about raw data

RNAseq (.fastq.gz)



5-10Gb/sample



# Workshop website

## Toxoplasma Pre-conference Workshop

### Introduction

This is the GitHub [repository](#) for the Toxoplasma Pre-Conference Workshop on data analysis and career perspectives. Goals of this workshop include:

- to interact with scientists who have pursued careers outside of the traditional academic tenure-track
- to gain an understanding of how RNAseq data is generated and how the raw data is processed (for practical reasons, we will discuss but not demonstrate this)
- to learn to use the R programming language to manage and explore the data (hands-on with example dataset)
- to learn to use ToxoDB to interpret RNAseq data in the context of a genome browser (also hands-on)

The workshop will have two main components: 1) career perspectives from scientists outside of the traditional tenure-track. 2) hands-on workshop on analysing RNAseq data. The latter will be broken a morning and afternoon workshop. The morning workshop will cover RNAseq analysis in R/bioconductor (Dan Beiting); and the afternoon will focus on using the genome browser to view RNAseq data (Omar Harb). We encourage all attendees to bring their laptops to allow participation in the hands-on exercises.

# What you will need to participate

**There are four main things you will need in order to participate in the hands-on aspects of the workshop**

## 1. Your laptop computer

- Everyone is encouraged to bring their own internet-enabled laptop equipped with either the recent Mac or Windows OS.

## 2. A “toolbox” of free software

- Download and install the appropriate version of the [R Programming Language](#) for your operating system
- Download and install the graphical user interface for R, called [RStudio](#)
- Download a text editor. I use [TextWrangler](#) (Mac only) and [Sublime](#) (Mac and PC)

## 3. Sample dataset

- We will be analyzing RNAseq data from this [2014 PLOS One](#) published by Kami Kim's lab.
- The original raw data is available on the Short Read Archive [here](#), but you **do not** need to download this (files are huge). Instead, please download the 'digital gene expression list' (DGEList) from [here](#)
- In addition to the data, please download this simple [text file](#) that describes the design of the study

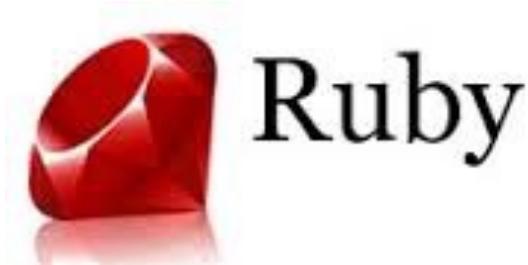
## 4. An analysis script

- the R script we'll use in the workshop can be downloaded [here](#).



## Bioinformatics toolbox

<b>Text editor</b>	TextWrangler; Sublime	Emacs; sublime
<b>R/bioconductor</b>	R	R
<b>GUI for R</b>	RStudio	RStudio
<b>Terminal</b>	terminal	terminal
<b>Network analysis</b>	Cytoscape	Cytoscape
<b>Graphing</b>	R	R



Unix (shell)

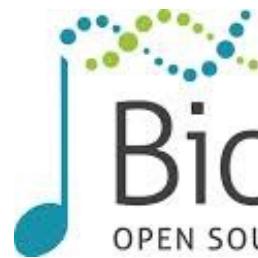
A screenshot of a terminal window titled "Terminal". It displays a file listing in a standard Unix-style format. The output includes file names, permissions, sizes, dates, and paths. The terminal window has a dark background with light-colored text.



# Perl



# Ruby



# Bioconductor

OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

# C++



# Unix (shell)

```
-rwxr-x--x 1 sys 52850 Jun 8 1979 httpunix  
drwxr-x--x 2 bin 320 Sep 22 05:33 lib  
drwxr-x--x 2 root 96 Sep 22 05:46 adec  
-rwxr-x--x 1 root 50950 Jun 8 1979 rkmunix  
-rwxr-x--x 1 root 51982 Jun 8 1979 rl2unlx  
-rwxr-x--x 1 sys 51750 Jun 8 1979 rphunix  
-rwxr-x--x 1 sys 51274 Jun 8 1979 roptunix  
drwxr-x--x 2 root 48 Sep 22 05:50 tmp  
drwxr-x--x 2 root 192 Sep 22 05:48 var  
# ls -l /usr  
Total 11  
drwxr-x--x 3 bin 128 Sep 22 05:45 dict  
drwxr-x--x 2 dmr 32 Sep 22 05:48 dmr  
drwxr-x--x 5 bin 416 Sep 22 05:46 games  
drwxr-x--x 3 sys 496 Sep 22 05:42 include  
drwxr-x--x 10 bin 528 Sep 22 05:43 lib  
drwxr-x--x 11 bin 176 Sep 22 05:45 man  
drwxr-x--x 3 bin 208 Sep 22 05:46 adec  
drwxr-x--x 2 bin 80 Sep 22 05:46 pub  
drwxr-x--x 6 root 96 Sep 22 05:45 spool  
drwxr-x--x 13 root 208 Sep 22 05:42 src  
# ls -l /usr/dmr  
Total 0
```

# Commerical Solutions

The screenshot displays the GeneSpring 11.9 software interface, which includes:

- Project Navigator:** Shows a "Demo Project" with an experiment titled "HeLa cells treated with compound X".
- Analysis Tools:** A tree view of analysis flags, including:
  - All Entities
  - Filtered on Flags [P, M]
  - T-test, p<0.05
  - Fold change >= 2.0
  - GO Analysis, p < 0.1
  - molecular\_function
  - cellular\_component
  - K-Means on Treatment
  - Significant Pathways, p < 0.1
  - MT-HeavyMetal-Path
  - Fold change >= 5.0
- Data Visualizations:**
  - A line plot titled "HeLa cells treated with compound X" showing "Normalized Intensity Values" over time.
  - A heatmap showing gene expression levels across samples.
  - A network diagram titled "MT-HeavyMetal-Pathway X" illustrating interactions between genes like FLMC, MTIF, MTIE, TNF, MTIX, TAFI, MTIA, and CO2.
- Workflow:** A sidebar listing experimental setup, quality control, analysis (Statistical Analysis, Filter on Volcano Plot, Fold Change, Clustering, Find Similar Entities, Filter on Parameters, Principal Component Analysis), and Class Prediction.

**Annotations:**

- Easy to use Wizard-Driven Workflows**
- Powerful data visualization**
- Integrated Pathway Analysis**
- Intuitive Data Management**

**price: \$4846/yr  
not open source  
not reproducible  
power vs simplicity**

Open R

# What you need to participate...

**There are four main things you will need in order to participate in the hands-on aspects of the workshop**

## 1. Your laptop computer

- Everyone is encouraged to bring their own internet-enabled laptop equipped with either the recent Mac or Windows OS.

## 2. A “toolbox” of free software

- Download and install the appropriate version of the [R Programming Language](#) for your operating system
- Download and install the graphical user interface for R, called [RStudio](#)
- Download a text editor. I use [TextWrangler](#) (Mac only) and [Sublime](#) (Mac and PC)

## 3. Sample dataset

- We will be analyzing RNAseq data from this 2014 PLOS One published by Kami Kim's lab.
- The original raw data is available on the Short Read Archive [here](#), but you **do not** need to download this (files are huge). Instead, please download the 'digital gene expression list' (DGEList) from [here](#)
- In addition to the data, please download this simple [text file](#) that describes the design of the study

## 4. An analysis script

- the R script we'll use in the workshop can be downloaded [here](#).

...you need R

## The Comprehensive R Archive Network

### Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

### Source Code for all Platforms

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (2015-04-16, Full of Ingredients) [R-3.2.0.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in time periods before a planned release).
- Daily snapshots of current patched and development versions are [available here](#). Please read about [new features and bug fixes](#) before filing corresponding feature requests or bug reports.
- Source code of older versions of R is [available here](#).
- Contributed extension [packages](#)



R Console



[ ~ ] Q Help Search

R version 3.1.2 (2014-10-31) -- "Pumpkin Helmet"  
Copyright (C) 2014 The R Foundation for Statistical Computing  
Platform: x86\_64-apple-darwin13.4.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

[R.app GUI 1.65 (6833) x86\_64-apple-darwin13.4.0]

[Workspace restored from /Users/danielbeiting/.RData]  
[History restored from /Users/danielbeiting/.Rhistory]

&gt;



# Why R/bioconductor?

*A powerful and flexible tool for bioinformatics*

## Pros

- R supports missing values - handles real-world data
- R's package system - simple method for handling user-contributed functions so that others understand them (standardized, modular, and extensible)
- used for over a decade to analyze genomic data
- often don't need to write code *de novo*
- huge support community
- publication quality graphics
- not just arrays - ChIPseq, RNAseq, genetic screens, interaction networks

# Why R/bioconductor?

*A powerful and flexible tool for bioinformatics*

## Cons

- steep learning curve
- no single ‘right way’ to do anything
- not as fast as C++
- historically, R was lacking in terms of user interface

Welcome to RStudio - Open source  
and enterprise-ready professional  
software for R

[Download RStudio](#)[Discover Shiny](#)

# What you need to participate...

**There are four main things you will need in order to participate in the hands-on aspects of the workshop**

## 1. Your laptop computer

- Everyone is encouraged to bring their own internet-enabled laptop equipped with either the recent Mac or Windows OS.

## 2. A “toolbox” of free software

- Download and install the appropriate version of the [R Programming Language](#) for your operating system
- Download and install the graphical user interface for R, called [RStudio](#)
- Download a text editor. I use [TextWrangler](#) (Mac only) and [Sublime](#) (Mac and PC)

## 3. Sample dataset

- We will be analyzing RNAseq data from this [2014 PLOS One](#) published by Kami Kim's lab.
- The original raw data is available on the Short Read Archive [here](#), but you **do not** need to download this (files are huge). Instead, please download the 'digital gene expression list' (DGEList) from [here](#)
- In addition to the data, please download this simple [text file](#) that describes the design of the study

## 4. An analysis script

- the R script we'll use in the workshop can be downloaded [here](#).

# Open RStudio

The screenshot shows the RStudio interface with a red box highlighting the left pane, which contains the R console output. The right pane shows the 'Environment' and 'Global Environment' tabs, with the 'Global Environment' tab selected. The environment is currently empty.

R version 3.1.1 (2014-07-10) -- "Sock it to Me"  
Copyright (C) 2014 The R Foundation for Statistical Computing  
Platform: x86\_64-apple-darwin13.1.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |

# Console

Environment History

Import Dataset Clear

Global Environment

Environment is empty

Files Plots Packages Help Viewer

Install Update

Name	Description	Version
acepack	ace() and avas() for selecting regression transformations	1.3-3.3
affxparser	Affymetrix File Parsing SDK	1.36.0
affy	Methods for Affymetrix Oligonucleotide Arrays	1.42.3
affyio	Tools for parsing Affymetrix data files	1.32.0
annotate	Annotation for microarrays	1.42.1
AnnotationDbi	Annotation Database Interface	1.26.1
assertthat	Easy pre and post assertions.	0.1
base64	Base 64 encoder/decoder	1.1
base64enc	Tools for base64 encoding	0.1-2
BatchJobs	Batch computing with R.	1.5
BBmisc	Miscellaneous Helper Functions for B. Bischl	1.9
beanplot	Visualization via Beanplots (like Boxplot/Stripchart/Violin Plot)	1.2
BH	Boost C++ Header Files	1.55.0-3
Biobase	Biobase: Base functions for Bioconductor	2.24.0
BiocGenerics	S4 generic functions for Bioconductor	0.10.0
BiocInstaller	Install/Update Bioconductor and CRAN Packages	1.14.3
BiocParallel	Bioconductor facilities for parallel evaluation	0.6.1
biomaRt	Interface to BioMart databases (e.g. Ensembl, COSMIC, Wormbase and Gramene)	2.20.0
Biostrings	String objects representing biological sequences, and matching algorithms	2.32.1
biovizBase	Basic graphic utilities for visualization of genomic data.	1.12.3
bit	A class for vectors of 1-bit booleans	1.1-12

Console ~/ Go to file/function

R version 3.1.1 (2014-07-10) -- "Sock it to Me"  
Copyright (C) 2014 The R Foundation for Statistical Computing  
Platform: x86\_64-apple-darwin13.1.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |

# Console

Environment History

Import Dataset Clear

Global Environment

# Workspace

Files Plots Packages Help Viewer

Install Update

Name	Description	Version
acepack	ace() and avas() for selecting regression transformations	1.3-3.3
affxparser	Affymetrix File Parsing SDK	1.36.0
affy	Methods for Affymetrix Oligonucleotide Arrays	1.42.3
affyio	Tools for parsing Affymetrix data files	1.32.0
annotate	Annotation for microarrays	1.42.1
AnnotationDbi	Annotation Database Interface	1.26.1
assertthat	Easy pre and post assertions.	0.1
base64	Base 64 encoder/decoder	1.1
base64enc	Tools for base64 encoding	0.1-2
BatchJobs	Batch computing with R.	1.5
BBmisc	Miscellaneous Helper Functions for B. Bischl	1.9
beanplot	Visualization via Beanplots (like Boxplot/Stripchart/Violin Plot)	1.2
BH	Boost C++ Header Files	1.55.0-3
Biobase	Biobase: Base functions for Bioconductor	2.24.0
BiocGenerics	S4 generic functions for Bioconductor	0.10.0
BiocInstaller	Install/Update Bioconductor and CRAN Packages	1.14.3
BiocParallel	Bioconductor facilities for parallel evaluation	0.6.1
biomaRt	Interface to BioMart databases (e.g. Ensembl, COSMIC, Wormbase and Gramene)	2.20.0
Biostrings	String objects representing biological sequences, and matching algorithms	2.32.1
biovizBase	Basic graphic utilities for visualization of genomic data.	1.12.3
bit	A class for vectors of 1-bit booleans	1.1-12

Project: (None)

Console ~/ Go to file/function

R version 3.1.1 (2014-07-10) -- "Sock it to Me"  
Copyright (C) 2014 The R Foundation for Statistical Computing  
Platform: x86\_64-apple-darwin13.1.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |

# Console

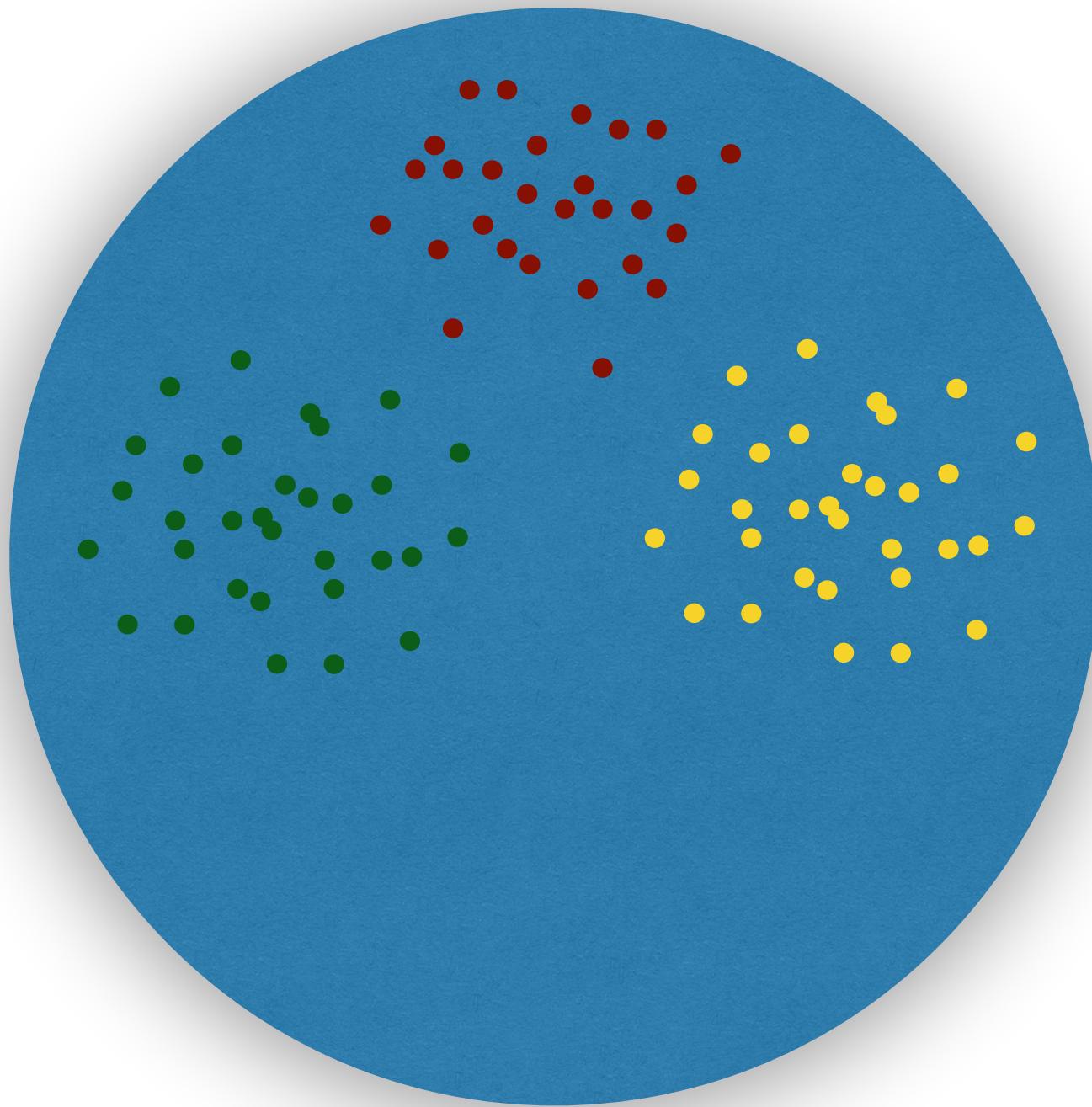
# Workspace

# Library and plots

Environment History Import Dataset Clear List Global Environment

Name	Description	Version
acepack	ace() and avas() for selecting regression transformations	1.3-3.3
affxparser	Affymetrix File Parsing SDK	1.36.0
affy	Methods for Affymetrix Oligonucleotide Arrays	1.42.3
affyio	Tools for parsing Affymetrix data files	1.32.0
annotate	Annotation for microarrays	1.42.1
AnnotationDbi	Annotation Database Interface	1.26.1
assertthat	Easy pre and post assertions.	0.1
base64	Base 64 encoder/decoder	1.1
base64enc	Tools for base64 encoding and decoding	0.1-2
BatchJobs	Batch Computing with R	0.5
BBmisc	Miscellaneous Helper Functions for B. Bischl	1.9
beanplot	Visualization via Beanplots (like Boxplot/Stripchart/Violin Plot)	1.2
BH	Boost C++ Header Files	1.55.0-3
Biobase	Biobase: Base functions for Bioconductor	2.24.0
BiocGenerics	S4 generic functions for Bioconductor	0.10.0
BiocInstaller	Install/Update Bioconductor and CRAN Packages	1.14.3
BiocParallel	Bioconductor facilities for parallel evaluation	0.6.1
biomaRt	Interface to BioMart databases (e.g. Ensembl, COSMIC, Wormbase and Gramene)	2.20.0
Biostrings	String objects representing biological sequences, and matching algorithms	2.32.1
biovizBase	Basic graphic utilities for visualization of genomic data.	1.12.3
bit	A class for vectors of 1-bit booleans	1.1-12

# The wonderful world of R



Package type

- statistics
- graphing
- modeling

~6500 packages

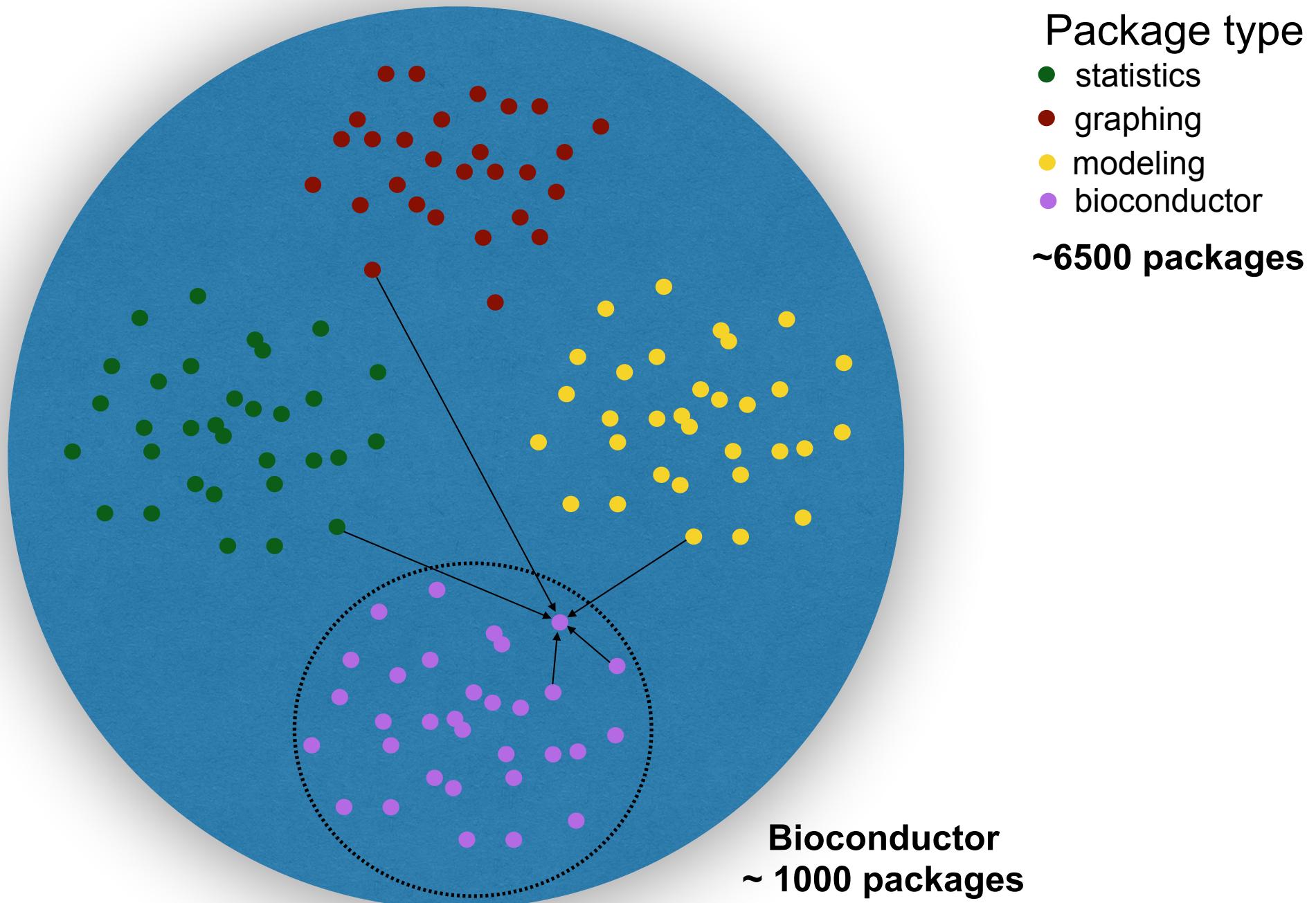
# Orchestrating high-throughput genomic analysis with Bioconductor

Wolfgang Huber<sup>1</sup>, Vincent J Carey<sup>2,3</sup>, Robert Gentleman<sup>4</sup>, Simon Anders<sup>1</sup>, Marc Carlson<sup>5</sup>, Benilton S Carvalho<sup>6</sup>, Hector Corrada Bravo<sup>7</sup>, Sean Davis<sup>8</sup>, Laurent Gatto<sup>9</sup>, Thomas Girke<sup>10</sup>, Raphael Gottardo<sup>11</sup>, Florian Hahne<sup>12</sup>, Kasper D Hansen<sup>13,14</sup>, Rafael A Irizarry<sup>3,15</sup>, Michael Lawrence<sup>4</sup>, Michael I Love<sup>3,15</sup>, James MacDonald<sup>16</sup>, Valerie Obenchain<sup>5</sup>, Andrzej K Oleś<sup>1</sup>, Hervé Pagès<sup>5</sup>, Alejandro Reyes<sup>1</sup>, Paul Shannon<sup>5</sup>, Gordon K Smyth<sup>17,18</sup>, Dan Tenenbaum<sup>5</sup>, Levi Waldron<sup>19</sup> & Martin Morgan<sup>5</sup>

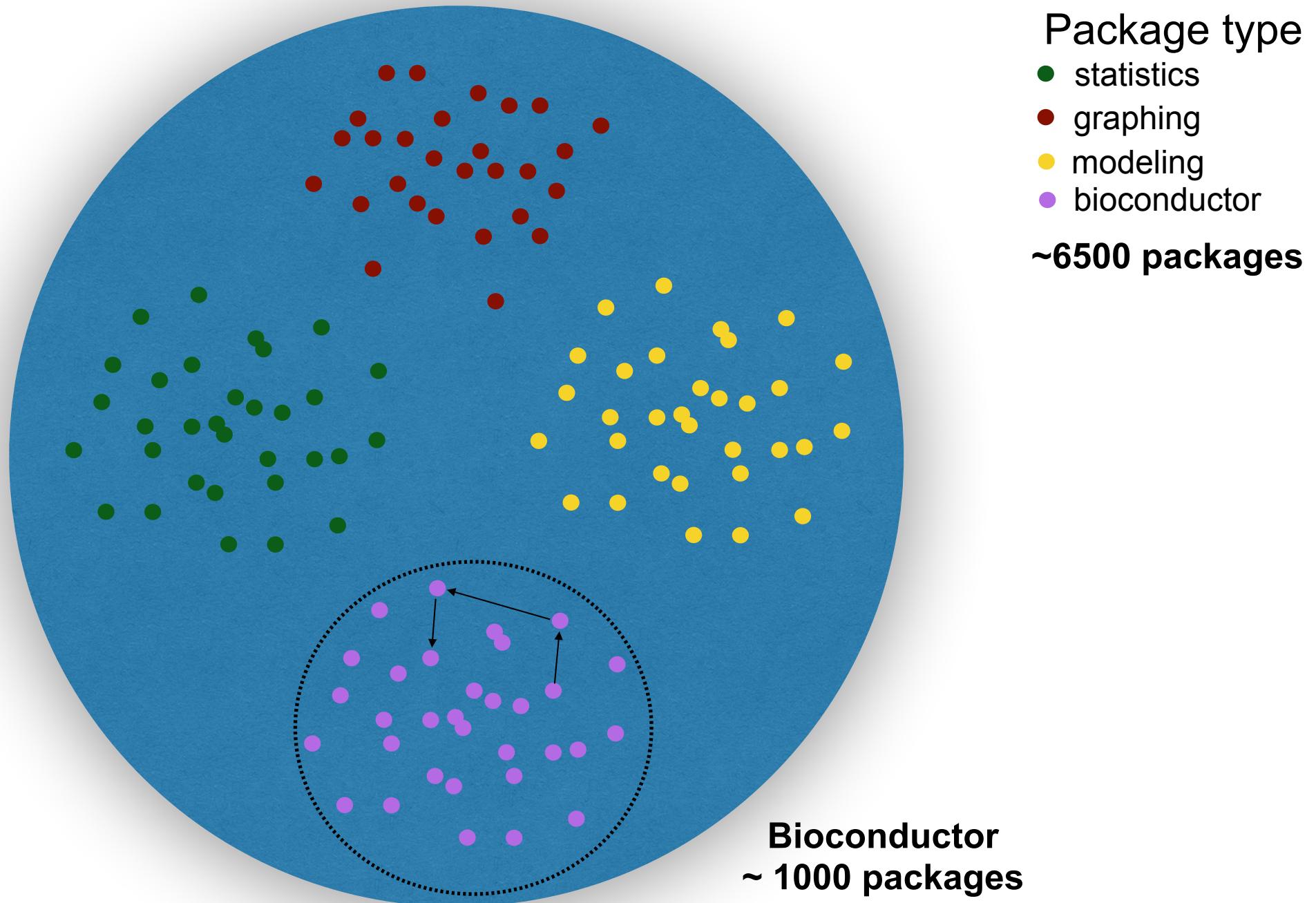
NATURE METHODS | VOL. 12 NO. 2 | FEBRUARY 2015 | 115

“We have embraced R for its scientific and statistical computing capabilities, for its graphics facilities and for the convenience of an interpreted language. R also interfaces with low-level languages including C and C++ for computationally intensive operations, Java for integration with enterprise software and JavaScript for interactive web-based applications and reports.”

# Packages are modular and can work together



# Workflows utilize multiple packages to step through a complex process



*how do I actually get  
bioconductor?*

R version 3.1.1 (2014-07-10) -- "Sock it to Me"  
Copyright (C) 2014 The R Foundation for Statistical Computing  
Platform: x86\_64-apple-darwin13.1.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> source("http://bioconductor.org/biocLite.R")  
**biocLite()**

Environment History

Import Dataset Clear

Global Environment

Environment is empty

Files Plots Packages Help Viewer

Install Update

Name	Description	Version
acepack	ace() and avas() for selecting regression transformations	1.3-3.3
affxparser	Affymetrix File Parsing SDK	1.36.0
affy	Methods for Affymetrix Oligonucleotide Arrays	1.42.3
affyio	Tools for parsing Affymetrix data files	1.32.0
annotate	Annotation for microarrays	1.42.1
AnnotationDbi	Annotation Database Interface	1.26.1
assertthat	Easy pre and post assertions.	0.1
base64	Base 64 encoder/decoder	1.1
base64enc	Tools for base64 encoding	0.1-2
BatchJobs	Batch computing with R.	1.5
BBmisc	Miscellaneous Helper Functions for B. Bischl	1.9
beanplot	Visualization via Beanplots (like Boxplot/Stripchart/Violin Plot)	1.2
BH	Boost C++ Header Files	1.55.0-3
Biobase	Biobase: Base functions for Bioconductor	2.24.0
BiocGenerics	S4 generic functions for Bioconductor	0.10.0
BiocInstaller	Install/Update Bioconductor and CRAN Packages	1.14.3
BiocParallel	Bioconductor facilities for parallel evaluation	0.6.1
biomaRt	Interface to BioMart databases (e.g. Ensembl, COSMIC, Wormbase and Gramene)	2.20.0
Biostrings	String objects representing biological sequences, and matching algorithms	2.32.1
biovizBase	Basic graphic utilities for visualization of genomic data.	1.12.3
bit	A class for vectors of 1-bit booleans	1.1-12

# Getting help

[Home](#) » [Help](#)

## Bioconductor Release »

Packages in the stable, semi-annual release:

- [Software](#)
- [Annotation Data](#) (Genome, Array, etc.)
- [Experiment Data](#)
- [Latest Release Announcement](#)

Bioconductor is also available as an [Amazon Machine Image](#) and a series of [Docker images](#).

## Support »

Please read the [posting guide](#). Post questions about Bioconductor to one of the following locations:

- [Support site](#) - for questions about Bioconductor packages
- [Bioc-devel mailing list](#) - for package developers

## Courses & Conferences »

Bioconductor courses and conferences explore R programming and the analysis of genomic data. View catalog of a [recent course material](#).

- [2015](#)
- [2014](#)
- [2013](#)
- [2012](#)
- [2011](#)
- [2010](#)
- [2009](#)
- [2008](#)
- [2007](#)
- [2006](#)
- [2005](#)
- [2004](#)
- [2003](#)
- [2002](#)

## Workflows »

Common Bioconductor workflows include:

- [Oligonucleotide Arrays](#)
- [High-throughput Sequencing](#)
- [Counting Reads for Differential Expression \(parathyroideSE vignette\)](#)
- [Annotation](#)
- [Annotating Variants](#)
- [Annotating Ranges](#)
- [Flow Cytometry](#) and other assays
- [Candidate Binding Sites for Known Transcription Factors](#)
- [Cloud-enabled cis-eQTL search and annotation](#)
- [RNA-Seq workflow: gene-level exploratory analysis and differential expression](#)
- [Changing genomic coordinate systems with rtracklayer:::liftOver](#)
- [Mass spectrometry and proteomics data analysis](#)

## Community Help Resources »

Community-contributed resources for learning Bioconductor:

- [Book](#) and [lab](#) from MOOC: PH525x Data Analysis for Genomics.
- [Video Lectures by Rafael Irizarry and guests](#)
- [R/Bioconductor Notes by Thomas Girke](#).
- [Using Bioconductor to Analyze your 23andme Data by Vince Buffalo](#).
- [Bioconductor Code Search by Itoshi NIKAIDO](#).
- [R/Bioconductor Blog By Sean Davis](#).

## Newsletters »

- [April 2014](#)
- [July 2014](#)
- [October 2014](#)
- [January 2015](#)

[Home](#) » [Help](#)

## Bioconductor Release »

Packages in the stable, semi-annual release:

- [Software](#)
- [Annotation Data](#) (Genome, Array, etc.)
- [Experiment Data](#)
- [Latest Release Announcement](#)

Bioconductor is also available as an [Amazon Machine Image](#) and a series of [Docker images](#).

## Support »

Please read the [posting guide](#). Post questions about Bioconductor to one of the following locations:

- [Support site](#) - for questions about Bioconductor packages
- [Bioc-devel mailing list](#) - for package developers

## Courses & Conferences »

Bioconductor courses and conferences explore R programming and the analysis of genomic data. View catalog of a [recent course material](#).

- [2015](#)
- [2014](#)
- [2013](#)
- [2012](#)
- [2011](#)
- [2010](#)
- [2009](#)
- [2008](#)
- [2007](#)
- [2006](#)
- [2005](#)
- [2004](#)
- [2003](#)
- [2002](#)

## Workflows »

Common Bioconductor workflows include:

- [Oligonucleotide Arrays](#)
- [High-throughput Sequencing](#)
- [Counting Reads for Differential Expression \(parathyroideSE vignette\)](#)
- [Annotation](#)
- [Annotating Variants](#)
- [Annotating Ranges](#)
- [Flow Cytometry](#) and other assays
- [Candidate Binding Sites for Known Transcription Factors](#)
- [Cloud-enabled cis-eQTL search and annotation](#)
- [RNA-Seq workflow: gene-level exploratory analysis and differential expression](#)
- [Changing genomic coordinate systems with rtracklayer:::liftOver](#)
- [Mass spectrometry and proteomics data analysis](#)

## Community Help Resources »

Community-contributed resources for learning Bioconductor:

- [Book](#) and [lab](#) from MOOC: PH525x Data Analysis for Genomics.
- [Video Lectures by Rafael Irizarry and guests](#)
- [R/Bioconductor Notes by Thomas Girke](#).
- [Using Bioconductor to Analyze your 23andme Data by Vince Buffalo](#).
- [Bioconductor Code Search by Itoshi NIKAIDO](#).
- [R/Bioconductor Blog By Sean Davis](#).

## Newsletters »

- [April 2014](#)
- [July 2014](#)
- [October 2014](#)
- [January 2015](#)

[Home](#) » [Help](#)

## Bioconductor Release »

Packages in the stable, semi-annual release:

- [Software](#)
- [Annotation Data](#) (Genome, Array, etc.)
- [Experiment Data](#)
- [Latest Release Announcement](#)

Bioconductor is also available as an [Amazon Machine Image](#) and a series of [Docker images](#).

## Support »

Please read the [posting guide](#). Post questions about Bioconductor to one of the following locations:

- [Support site](#) - for questions about Bioconductor packages
- [Bioc-devel mailing list](#) - for package developers

## Courses & Conferences »

Bioconductor courses and conferences explore R programming and the analysis of genomic data. View catalog of a [recent course material](#).

- [2015](#)
- [2014](#)
- [2013](#)
- [2012](#)
- [2011](#)
- [2010](#)
- [2009](#)
- [2008](#)
- [2007](#)
- [2006](#)
- [2005](#)
- [2004](#)
- [2003](#)
- [2002](#)

## Workflows »

Common Bioconductor workflows include:

- [Oligonucleotide Arrays](#)
- [High-throughput Sequencing](#)
- [Counting Reads for Differential Expression \(parathyroideSE vignette\)](#)
- [Annotation](#)
- [Annotating Variants](#)
- [Annotating Ranges](#)
- [Flow Cytometry](#) and other assays
- [Candidate Binding Sites for Known Transcription Factors](#)
- [Cloud-enabled cis-eQTL search and annotation](#)
- [RNA-Seq workflow: gene-level exploratory analysis and differential expression](#)
- [Changing genomic coordinate systems with rtracklayer:::liftOver](#)
- [Mass spectrometry and proteomics data analysis](#)

## Community Help Resources »

Community-contributed resources for learning Bioconductor:

- [Book](#) and [lab](#) from MOOC: PH525x Data Analysis for Genomics.
- [Video Lectures by Rafael Irizarry and guests](#)
- [R/Bioconductor Notes by Thomas Girke](#).
- [Using Bioconductor to Analyze your 23andme Data by Vince Buffalo](#).
- [Bioconductor Code Search by Itoshi NIKAIDO](#).
- [R/Bioconductor Blog By Sean Davis](#).

## Newsletters »

- [April 2014](#)
- [July 2014](#)
- [October 2014](#)
- [January 2015](#)



ASK QUESTION

LATEST 70

NEWS 2

JOBS

TAGS

USERS

Limit ▾

Sort ▾

Search

0 votes  
0 answers

4 views

[UPC and Sequence file to remove GC bias](#)

upc\_generic\_expressionset brainarray

written 1 hour ago by sanj • 10

0 votes  
0 answers

2 views

[UPC and Sequence file to remove GC bias](#)

upc\_generic\_expressionset brainarray

written 1 hour ago by sanj • 10

0 votes  
4 answers

44 views

[SICER and Diffbind](#)

difffbind sicer

written 13 days ago by sergio.espeso-gil • 0

0 votes  
0 answers

8 views

[Support for MouseDivGeno arrays in Bioconductor](#)

snpchip affy microarray

written 4 hours ago by Sean Davis • 20k

3 votes  
7 answers

244 views

[reactome.db is not updated?](#)

pathways geneset enrichment reactome.db

written 4 months ago by Guangchuang Yu • 120 • updated 7 hours ago by willem.ligtenberg • 0

0 votes  
4 answers

277 views

[reactome.db: reactome IDs not mapped to pathway names](#)

pathways annotation

written 5 months ago by johannes.rainer • 90 • updated 7 hours ago by willem.ligtenberg • 0

1 vote  
2 answers

35 views

[SCAN UPC normalisation](#)

normalization scan.upc

written 13 days ago by sanj • 10

0 votes  
0 answers

14 views

[browseGenome rtracklayer, save / export UCSC genome browser image e.g. as pdf](#)

rtracklayer browsegenome image save file ucsc

written 7 hours ago by James Perkins • 60

0 votes  
1 answer

18 views

[Error while reading in peaksets into DiffBind.](#)

reading peaks difffbind

written 19 hours ago by surjray • 0 • updated 8 hours ago by Gord Brown • 180

0 votes  
1 answer

27 views

[How to analyze a known fusion gene expression in RNAseq](#)

fusion genes rnaseq

written 9 hours ago by stephen66 • 30 • updated 9 hours ago by b.nota • 10

0 votes  
0 answers

27 views

[Why I don't use the MethyAnalisis?](#)

methyanalysis

written 13 hours ago by 452519402 • 0

0 votes  
1 answer

25 views

[riboSeqR: plotTranscript error](#)

riboseqr

written 23 hours ago by antonvila.s • 0 • updated 18 hours ago by Thomas J Hardcastle • 160

0 votes  
2 answers

34 views

[Annotation information for ecoli2.db](#)

ecoli2.db

## Recent...

## Replies

- [C: SCAN UPC normalisation](#) by Stephen Piccolo • 390
- [C: Can't find any different...](#) by James W. MacDonald • 34k
- [C: reactome.db: reactome ID...](#) by Dan Du • 210
- [C: reactome.db: reactome ID...](#) by willem.ligtenberg • 0
- [A: SICER and Diffbind](#) by sergio.espeso-gil • 0

## Votes

- [C: Can't find any different...](#)
- [Bioconductor Newsletter - A...](#)
- [C: How to use brainarray cu...](#)
- [A: Problem installing XPS](#)
- [A: How to Find Common Dysre...](#)

## Awards • All »

- [Scholar](#)  to Martin Morgan • 14k
- [Teacher](#)  to Aaron Lun • 1.3k
- [Scholar](#)  to Aaron Lun • 1.3k
- [Teacher](#)  to Gordon Smyth • 22k
- [Autobiographer](#)  to Gordon Smyth • 22k
- [Scholar](#)  to James W. MacDonald • 34k

## Locations • All »

- United States, just now
- United States, 26 minutes ago
- European Molecular Biology Laboratory (Heidelberg, Germany), 28 minutes ago
- United States, 29 minutes ago
- United States, 45 minutes ago

[Home](#) » [Bioconductor 3.0](#) » [Software Packages](#) » limma

## limma

### Linear Models for Microarray Data

Bioconductor version: Release (3.0)

Data analysis, linear models and differential expression for microarray data.

Author: Gordon Smyth [cre,aut], Matthew Ritchie [ctb], Jeremy Silver [ctb], James Wettenhall [ctb], Natalie Thorne [ctb], Davis McCarthy [ctb], Di Wu [ctb], Yifang Hu [ctb], Wei Shi [ctb], Belinda Phipson [ctb], Alicia Oshlack [ctb], Carolyn de Graaf [ctb], Mette Langaas [ctb], Egil Ferkingstad [ctb], Marcus Davy [ctb], Francois Pepin [ctb], Dongseok Choi [ctb]

Maintainer: Gordon Smyth <smyth at wehi.edu.au>

Citation (from within R, enter `citation("limma")`):

Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W and Smyth GK (2015). "limma powers differential expression analyses for RNA-sequencing and microarray studies." *Nucleic Acids Research*, **43**, pp. doi: 10.1093/nar/gkv007.

### Installation

To install this package, start R and enter:

```
source("http://bioconductor.org/biocLite.R")
biocLite("limma")
```

### Documentation

To view documentation for the version of this package installed in your system, start R and enter:

```
browseVignettes("limma")
```

<a href="#">PDF</a>	<a href="#">R Script</a>	Limma One Page Introduction
<a href="#">PDF</a>		usersguide.pdf
<a href="#">PDF</a>		Reference Manual
<a href="#">Text</a>		NEWS

### Workflows »

Common Bioconductor workflows include:

- [Oligonucleotide Arrays](#)
- [High-throughput Sequencing](#)
- [Counting Reads for Differential Expression](#) (*parathyroideSE* vignette)
- [Annotation](#)
- [Annotating Variants](#)
- [Annotating Ranges](#)
- [Flow Cytometry](#) and other assays
- [Candidate Binding Sites for Known Transcription Factors](#)
- [Cloud-enabled cis-eQTL search and annotation](#)
- [RNA-Seq workflow: gene-level exploratory analysis and differential expression](#)
- [Changing genomic coordinate systems with rtracklayer::liftOver](#)
- [Mass spectrometry and proteomics data analysis](#)

### Mailing Lists »

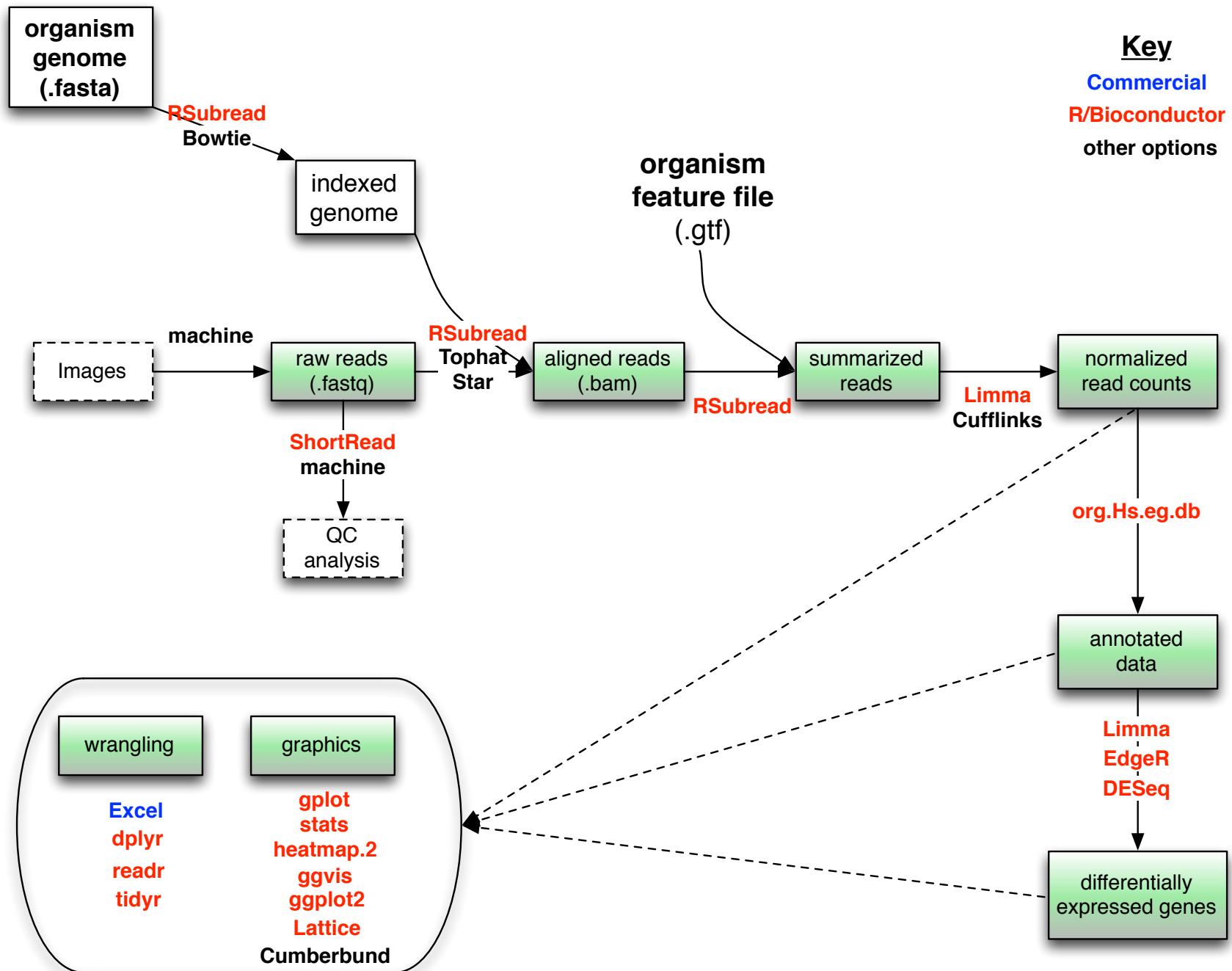
Post questions about Bioconductor packages to our mailing lists. Read the [posting guide](#) before posting!

- [bioconductor](#)
- [bioc-devel](#)

help for individual packages

How is R/Bioconductor used to analyze RNAseq data?

# RNAseq analysis workflow



# RNAseq: Subread aligner + FeatureCounts

## The Subread aligner: fast, accurate and scalable read mapping by seed-and-vote

Yang Liao<sup>1,2</sup>, Gordon K. Smyth<sup>1,3</sup> and Wei Shi<sup>1,2,\*</sup>

*Nucleic Acids Research*, 2013, Vol. 41, No. 10 e108  
doi:10.1093/nar/gkt214

<sup>1</sup>Division of Bioinformatics, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, Victoria 3052, Australia, <sup>2</sup>Department of Computing and Information Systems, The University of Melbourne, Parkville, Victoria 3010, Australia and <sup>3</sup>Department of Mathematics and Statistics, The University of Melbourne, Parkville, Victoria 3010, Australia

BIOINFORMATICS

ORIGINAL PAPER

Vol. 30 no. 7 2014, pages 923–930  
doi:10.1093/bioinformatics/btt656

Sequence analysis

Advance Access publication November 13, 2013

## featureCounts: an efficient general purpose program for assigning sequence reads to genomic features

Yang Liao<sup>1,2</sup>, Gordon K. Smyth<sup>1,3</sup> and Wei Shi<sup>1,2,\*</sup>

<sup>1</sup>Bioinformatics Division, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, VIC 3052, <sup>2</sup>Department of Computing and Information Systems and <sup>3</sup>Department of Mathematics and Statistics, The University of Melbourne, Parkville, VIC 3010, Australia

Associate Editor: Martin Bishop

# What you need to participate...

**There are four main things you will need in order to participate in the hands-on aspects of the workshop**

## 1. Your laptop computer

- Everyone is encouraged to bring their own internet-enabled laptop equipped with either the recent Mac or Windows OS.

## 2. A “toolbox” of free software

- Download and install the appropriate version of the [R Programming Language](#) for your operating system
- Download and install the graphical user interface for R, called [RStudio](#)
- Download a text editor. I use [TextWrangler](#) (Mac only) and [Sublime](#) (Mac and PC)

## 3. Sample dataset

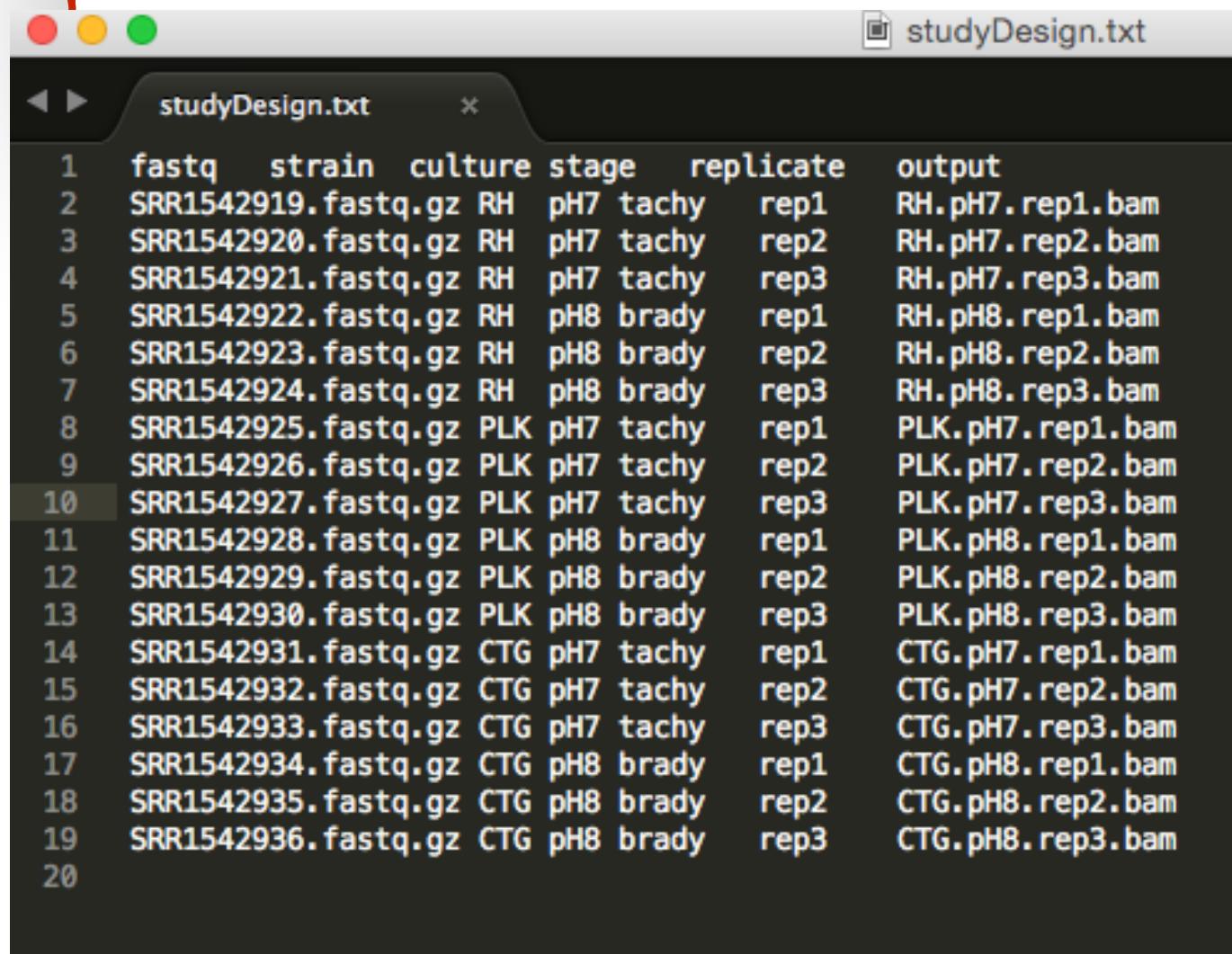
- We will be analyzing RNAseq data from this [2014 PLOS One published by Kami Kim's lab](#).
- The original raw data is available on the Short Read Archive [here](#), but you **do not** need to download this (files are huge). Instead, please download the 'digital gene expression list' (DGEList) from [here](#)
- In addition to the data, please download this simple [text file](#) that describes the design of the study

## 4. An analysis script

- the R script we'll use in the workshop can be downloaded [here](#).

# Distinct Strains of *Toxoplasma gondii* Feature Divergent Transcriptomes Regardless of Developmental Stage

Matthew McKnight Croken<sup>1</sup>, Yanfen Ma<sup>2</sup>, Lye Meng Markillie<sup>3</sup>, Ronald C. Taylor<sup>4</sup>, Galya Orr<sup>3</sup>, Louis M. Weiss<sup>2,5\*</sup>, Kami Kim<sup>1,2,5\*</sup>

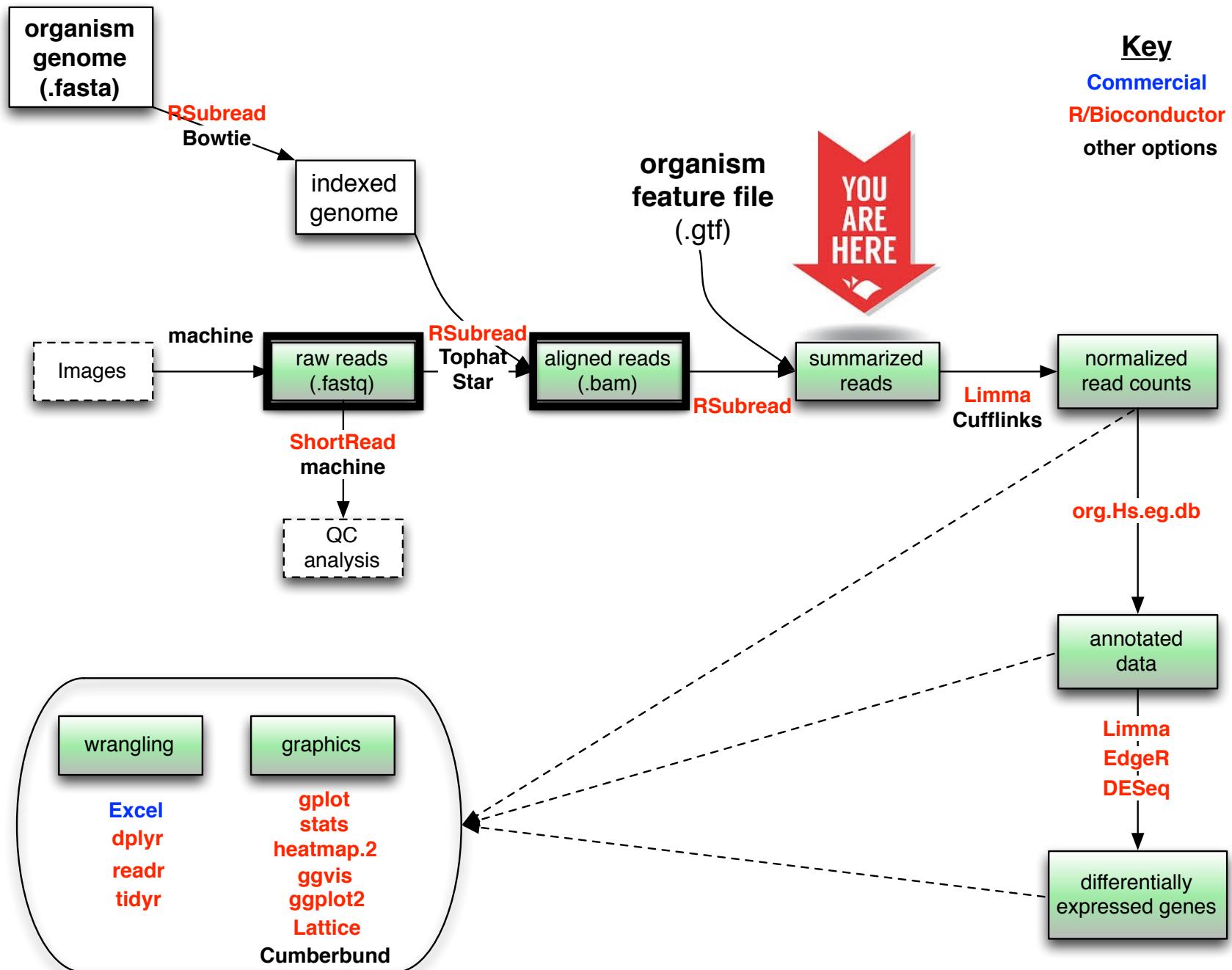


studyDesign.txt

	fastq	strain	culture	stage	replicate	output
1	SRR1542919.fastq.gz	RH	pH7	tachy	rep1	RH.pH7.rep1.bam
2	SRR1542920.fastq.gz	RH	pH7	tachy	rep2	RH.pH7.rep2.bam
3	SRR1542921.fastq.gz	RH	pH7	tachy	rep3	RH.pH7.rep3.bam
4	SRR1542922.fastq.gz	RH	pH8	brady	rep1	RH.pH8.rep1.bam
5	SRR1542923.fastq.gz	RH	pH8	brady	rep2	RH.pH8.rep2.bam
6	SRR1542924.fastq.gz	RH	pH8	brady	rep3	RH.pH8.rep3.bam
7	SRR1542925.fastq.gz	PLK	pH7	tachy	rep1	PLK.pH7.rep1.bam
8	SRR1542926.fastq.gz	PLK	pH7	tachy	rep2	PLK.pH7.rep2.bam
9	SRR1542927.fastq.gz	PLK	pH7	tachy	rep3	PLK.pH7.rep3.bam
10	SRR1542928.fastq.gz	PLK	pH8	brady	rep1	PLK.pH8.rep1.bam
11	SRR1542929.fastq.gz	PLK	pH8	brady	rep2	PLK.pH8.rep2.bam
12	SRR1542930.fastq.gz	PLK	pH8	brady	rep3	PLK.pH8.rep3.bam
13	SRR1542931.fastq.gz	CTG	pH7	tachy	rep1	CTG.pH7.rep1.bam
14	SRR1542932.fastq.gz	CTG	pH7	tachy	rep2	CTG.pH7.rep2.bam
15	SRR1542933.fastq.gz	CTG	pH7	tachy	rep3	CTG.pH7.rep3.bam
16	SRR1542934.fastq.gz	CTG	pH8	brady	rep1	CTG.pH8.rep1.bam
17	SRR1542935.fastq.gz	CTG	pH8	brady	rep2	CTG.pH8.rep2.bam
18	SRR1542936.fastq.gz	CTG	pH8	brady	rep3	CTG.pH8.rep3.bam
19						
20						

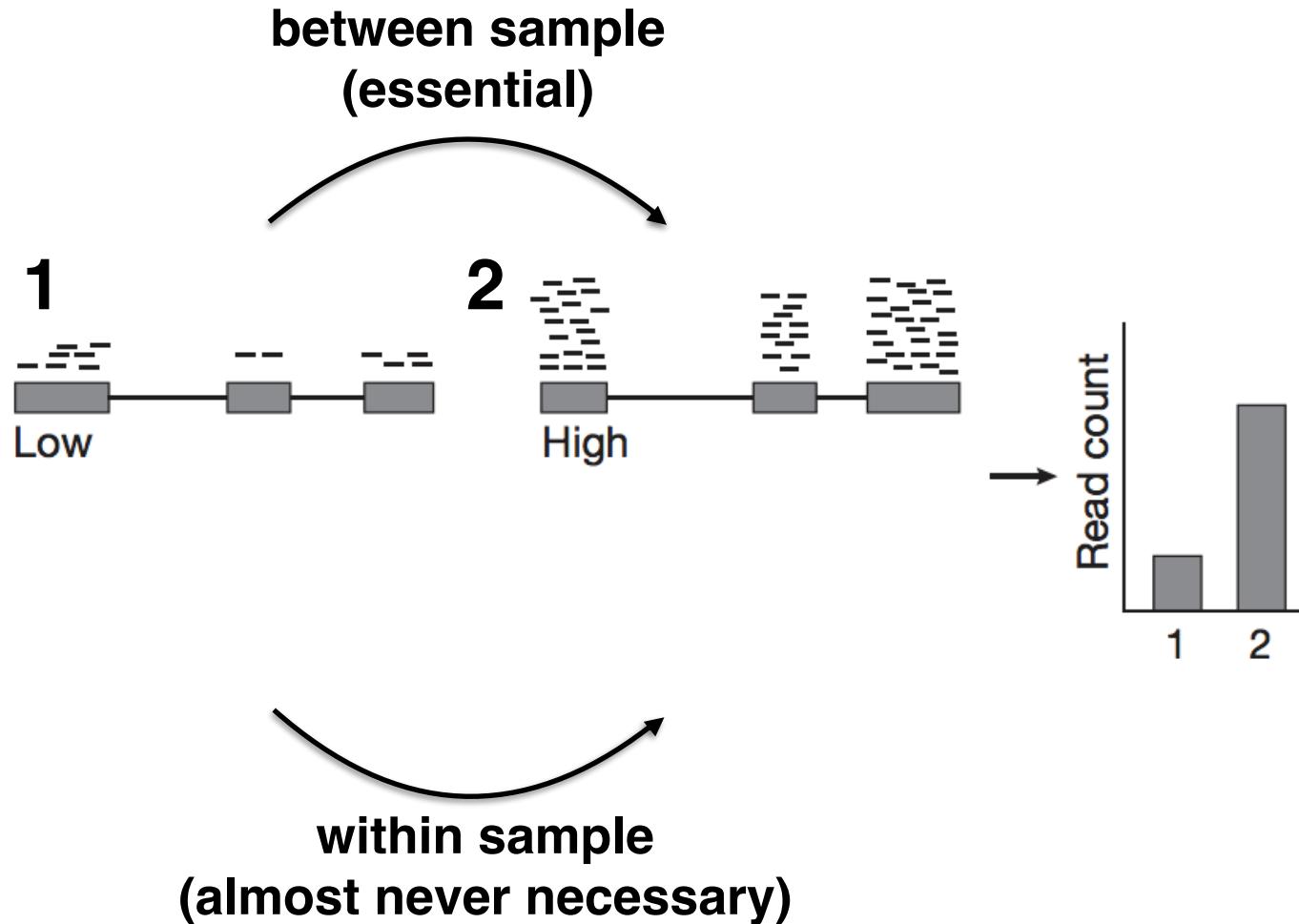
s of America, **2** Department of Pathology, Albert  
y, Pacific Northwest National Laboratory, Richland,  
, Pacific Northwest National Laboratory, Richland,  
rk, United States of America

# RNAseq analysis workflow

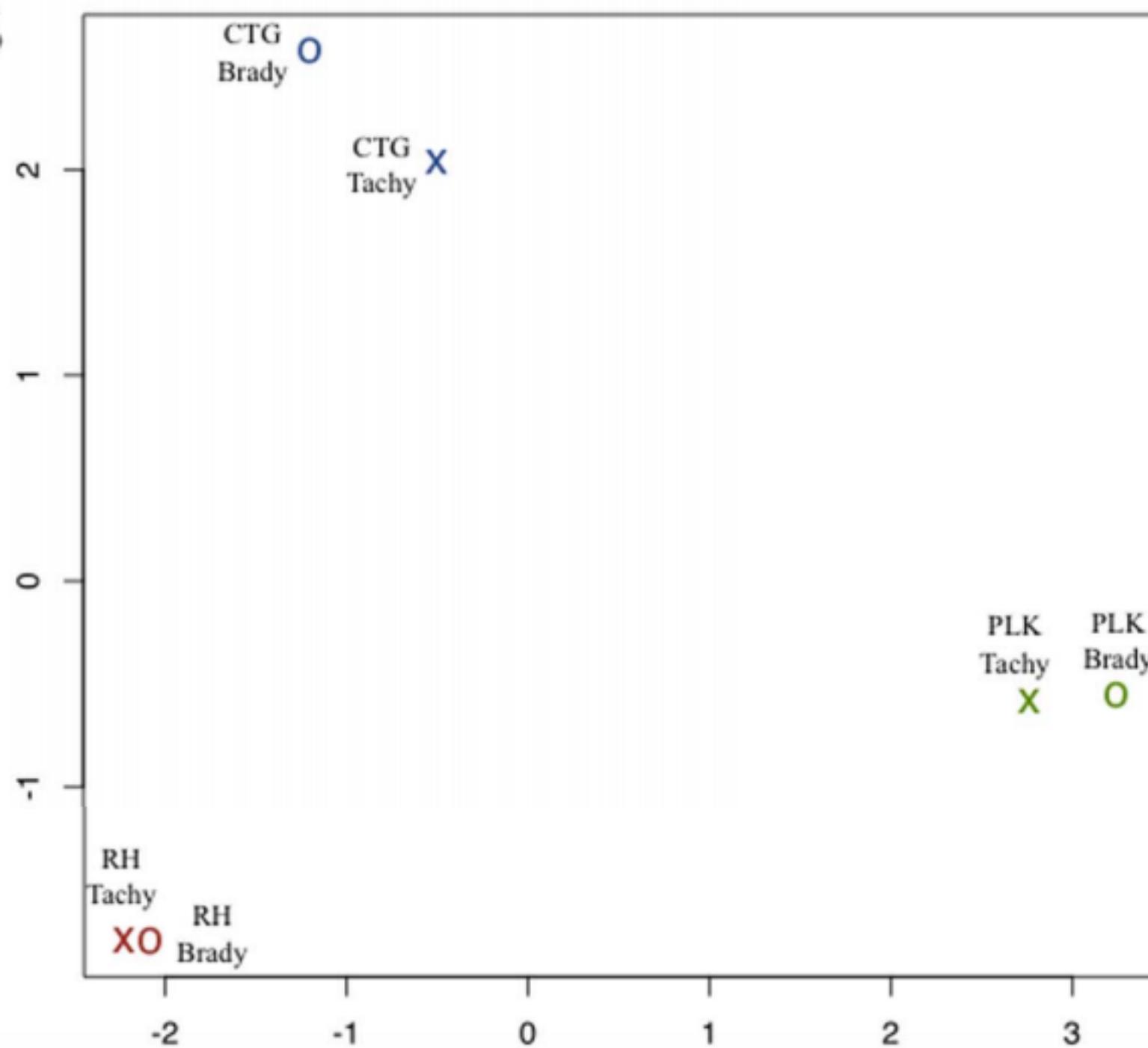


# Demo analysis

# Normalizing RNAseq data



*modified from Garber et al., Nature Methods, 2011*



# Comparison of the transcriptional landscapes between human and mouse tissues

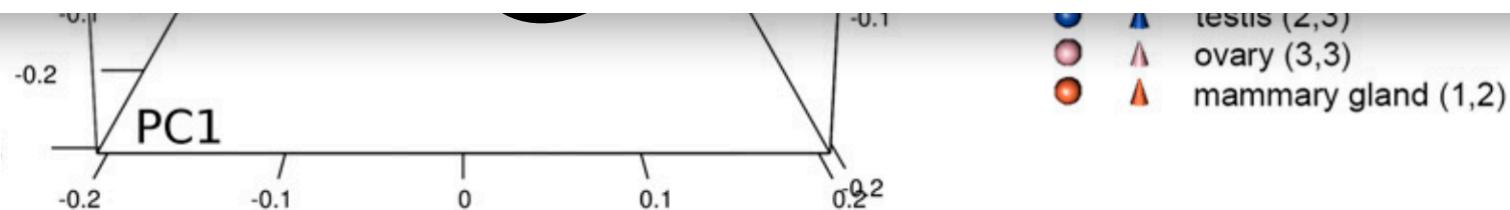
Shin Lin<sup>a,b,1</sup>, Yiing Lin<sup>c,1</sup>, Joseph R. Nery<sup>d</sup>, Mark A. Urich<sup>d</sup>, Alessandra Breschi<sup>e,f</sup>, Carrie A. Davis<sup>g</sup>, Alexander Dobin<sup>g</sup>, Christopher Zaleski<sup>g</sup>, Michael A. Beer<sup>h</sup>, William C. Chapman<sup>c</sup>, Thomas R. Gingeras<sup>g,i</sup>, Joseph R. Ecker<sup>d,j,2</sup>, and Michael P. Snyder<sup>a,2</sup>

<sup>a</sup>Department of Genetics, Stanford University, Stanford, CA 94305; <sup>b</sup>Division of Cardiovascular Medicine, Stanford University, Stanford, CA 94305;

<sup>c</sup>Department of Surgery, Washington University School of Medicine, St. Louis, MO 63110; <sup>d</sup>Genomic Analysis Laboratory, The Salk Institute for Biological Studies, La Jolla, CA 92037; <sup>e</sup>Centre for Genomic Regulation and UPF, Catalonia, 08003 Barcelona, Spain; <sup>f</sup>Departament de Ciències Experimentals i de la Salut, Universitat Pompeu Fabra, 08003 Barcelona, Spain; <sup>g</sup>Functional Genomics, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11742; <sup>h</sup>McKusick-Nathans Institute of Genetic Medicine and the Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21205; <sup>i</sup>Affymetrix, Inc., Santa Clara, CA 95051; and <sup>j</sup>Howard Hughes Medical Institute, The Salk Institute for Biological Studies, La Jolla, CA 92037

## Significance

To date, various studies have found similarities between humans and mice on a molecular level, and indeed, the murine model serves as an important experimental system for biomedical science. In this study of a broad number of tissues between humans and mice, high-throughput sequencing assays on the transcriptome and epigenome reveal that, in general, differences dominate similarities between the two species. These findings provide the basis for understanding the differences in phenotypes and responses to conditions in humans and mice.





## RESEARCH ARTICLE

# A reanalysis of mouse ENCODE comparative gene expression data [v1; ref status: awaiting peer review, <http://f1000r.es/5ez>]

Yoav Gilad, Orna Mizrahi-Man

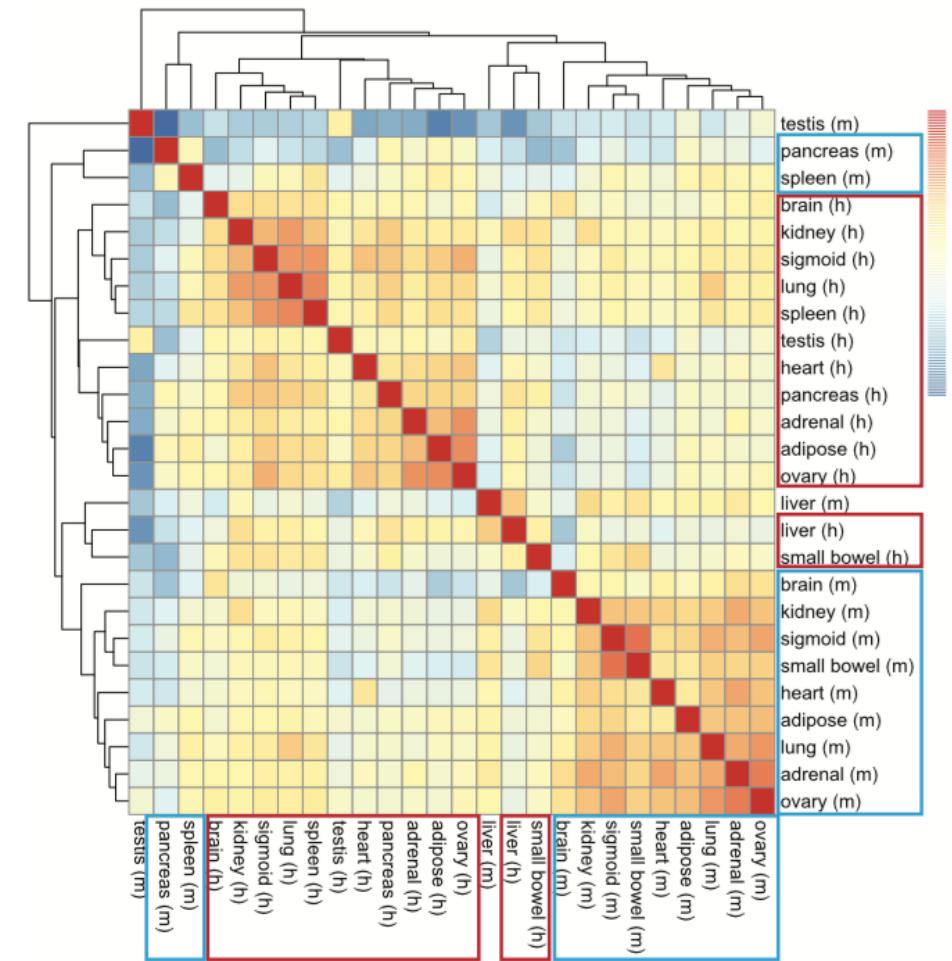
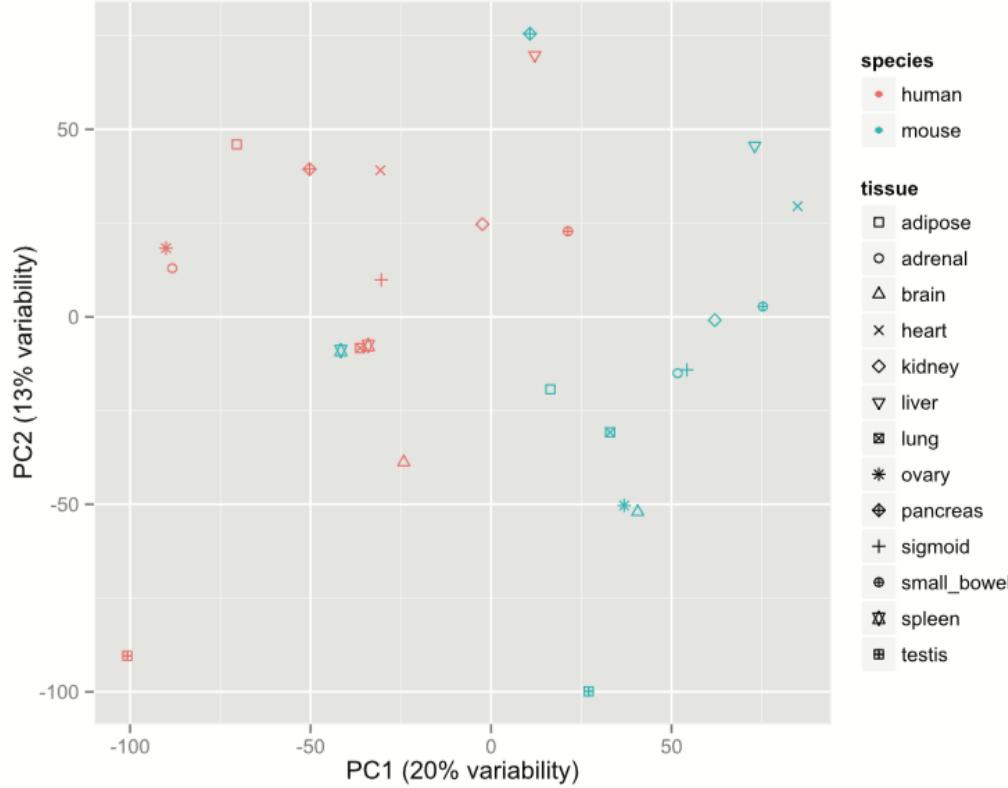
Department of Human Genetics, University of Chicago, Chicago, IL, 60637, USA

“The implications of the observation that human and mouse gene expression data may be clustering by species more than by tissues can be profound. To a large degree, modern biology is built upon the empirical observation that homologous gene regulatory networks establish the identities of homologous cell-types, tissues, and organs across species – the results of Lin *et al.*, if true, challenge these observations and the biological basis of homology. From a more practical perspective, the mouse is arguably the most important animal model for biomedical research. If gene regulation in any mouse tissue is markedly more representative of a general mouse regulatory network than the regulatory network of a corresponding human tissue, this would call into question the utility of the mouse, and perhaps any other non-human animal, as a useful model system for biomedical research.”

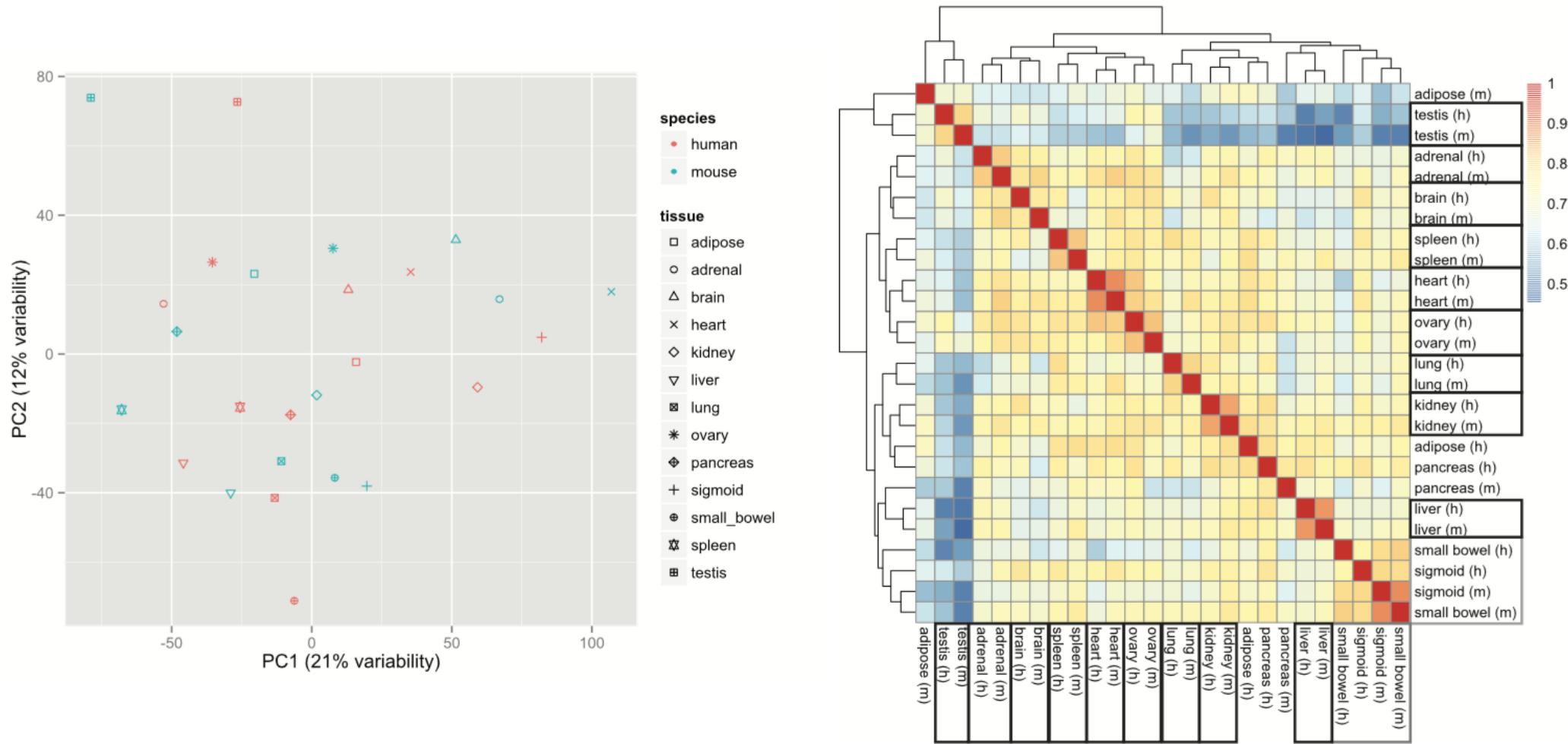
@<machine\_id>:<runnumber>:<flowcellID>:<lane>:<tile>:<x-pos>: <y-pos>  
 <read>:<is filtered>:<control number>:<index sequence>

@D4LHBFN1:276:C2HKJACXX:4:1101:3448:12374 1:N:0:AGTTCC

D87PMJN1 (run 253, flow cell D2GUAACXX, lane 7)	D87PMJN1 (run 253, flow cell D2GUAACXX , lane 8)	D4LHBFN1 (run 276, flow cell C2HKJACXX , lane 4)	MONK (run 312, flow cell C2GR3ACXX , lane 6)	HWI-ST373 (run 375, flow cell C3172ACXX , lane 7)
heart	adipose	adipose	heart	brain
kidney	adrenal	adrenal	kidney	pancreas
liver	sigmoid colon	sigmoid colon	liver	brain
small bowel	lung	lung	small bowel	spleen
spleen	ovary	ovary	testis	<span style="color:red;">●</span> Human
testis		pancreas		<span style="color:blue;">●</span> Mouse



**Figure 2. Recapitulating the patterns reported by the mouse ENCODE papers.** **a.** Two-dimensional plots of principal components calculated by performing PCA of the transposed log-transformed FPKM values (from 14,744 orthologous gene pairs) for the 26 samples, after removal of invariant columns (genes). **b.** Heatmap based on pairwise Pearson correlation of expression data used in panel **a**. We used Euclidean distance and complete linkage as distance measure and clustering method, respectively.



**Figure 3. Clustering of data once batch effects are accounted for.** **a.** Two-dimensional plots of principal components calculated by applying PCA to the transposed matrix of batch-corrected log-transformed normalized fragment counts (from 10,309 orthologous gene pairs that remained after the exclusion steps described in the results) for the 26 samples, after removal of invariant columns (genes). **b.** Heatmap based on pairwise Pearson correlation of the expression data used in panel **a**. We used Euclidean distance and complete linkage as distance measure and clustering method, respectively.

## Discuss this Article

Reader Comment 20 May 2015

**Shin Lin**, Department of Genetics, Stanford University, USA

“There remains the issue of our study design with respect to confounding of lane effect and species. It should be noted that our study design minimized library preparation and primer index effect. Although our experience is consistent with [reports] showing that lane effect is not a large contributor to variance, we recognize it is better to have data. Thus, we are sequencing under a new pattern of pooled libraries, and soon, we will post the results for the community.” - *authors of PNAS study*

Author Response 20 May 2015

**Yoav Gilad**, Human Genetics, University of Chicago, USA

Why are you calling these 'lane effects'? The samples were sequenced on different instruments (different sequencers) and by default - different flow cells. Can you explicitly acknowledge that fact please?

# Data wrangling with the dplyr package

*4 essential ‘verbs’ encompass  
most of the functionality of dplyr*

1. **Mutate** - add new columns as a function of existing data
2. **Filter** - picks rows based on rules you apply
3. **Select** - picks columns based on rules you apply
4. **Arrange** - reorders rows by sorting based on column(s)