

What is



?

Principally a collaborative software development project

But it is also:

- a software repository
- a bioinformatics support site
- data repository
- publisher for supplementary materials
- source for tutorials and instructional documentation

Managed and maintained by a core team of 6 people,
with contributions coming from all over the world

What is



?

Principally a collaborative software development project

But it is

- a software
- a bioinformatics
- data science
- public
- source

Managed
with



What is Bioconductor?

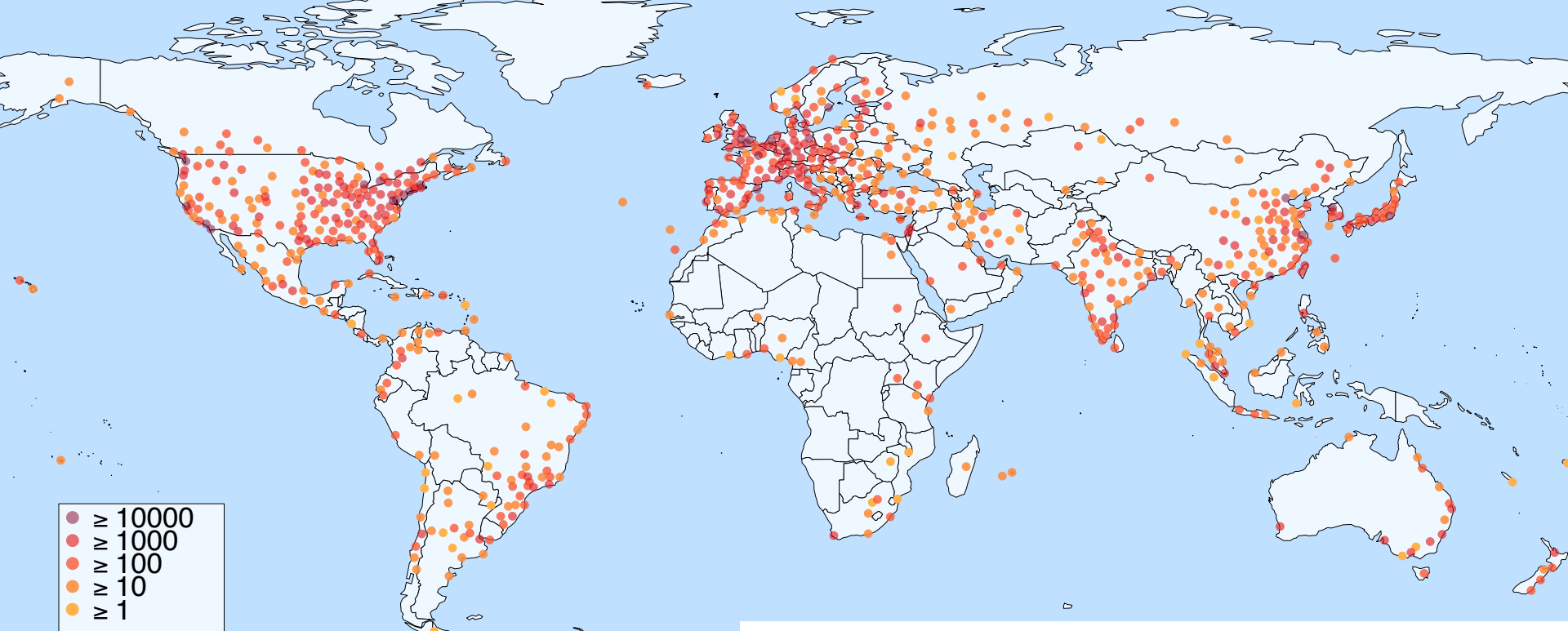
Started 2002 as a platform for analysis & understanding of microarray data

More than 1,300 packages. Domains of expertise:

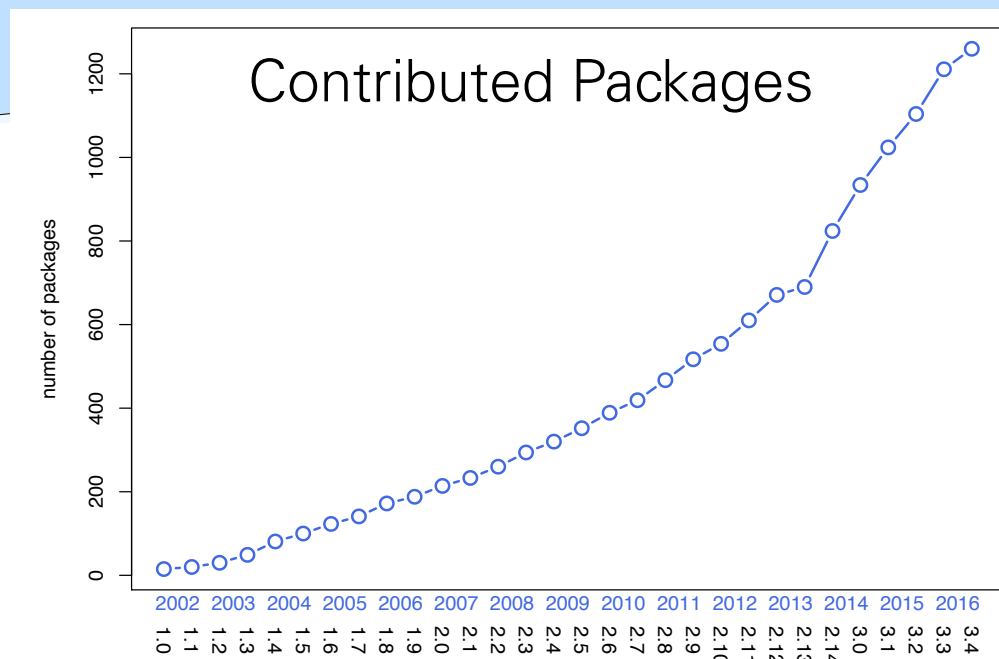
- Sequencing (RNASeq, ChIPSeq, single-cell, called variants, ...)
- Microarrays (methylation, expression, copy number, ...)
- Flow cytometry
- Proteomics
- Multi-Omics data integration

Important themes

- Reproducible research
- Interoperability between packages & workflows
... even from different authors
- Usability



World largest bioinformatics project
 10,000s users
 >18,000 papers in PubmedCentral



Collaborative
and distributed development
Open source

Lower barrier of entry
Training
Turn users into
developers

Data import,
preprocessing
Integration of data types
Based on R

Interoperable
components
Rapid development
Code re-use

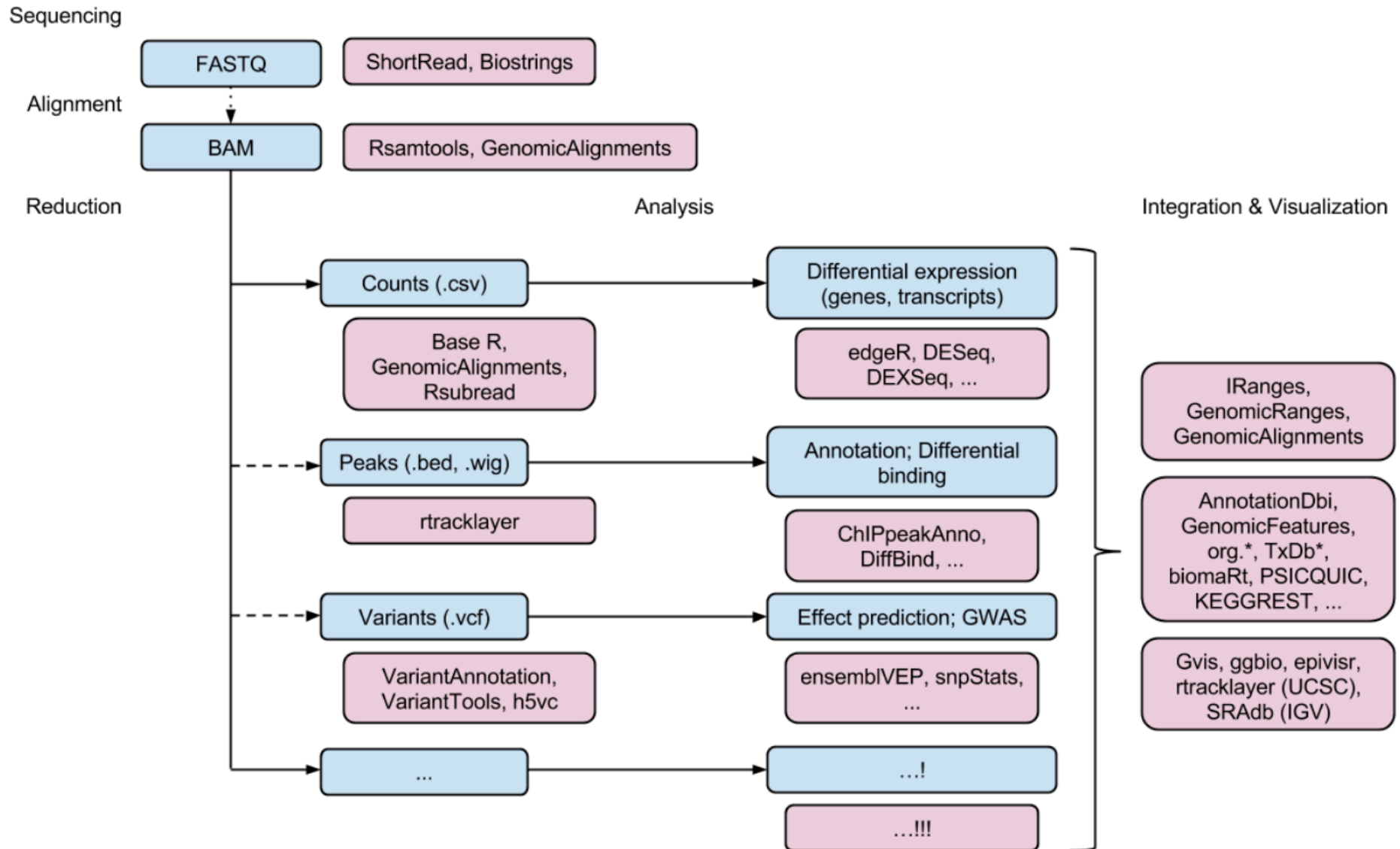
Publication of software
Computational reproducibility

Motivating principles

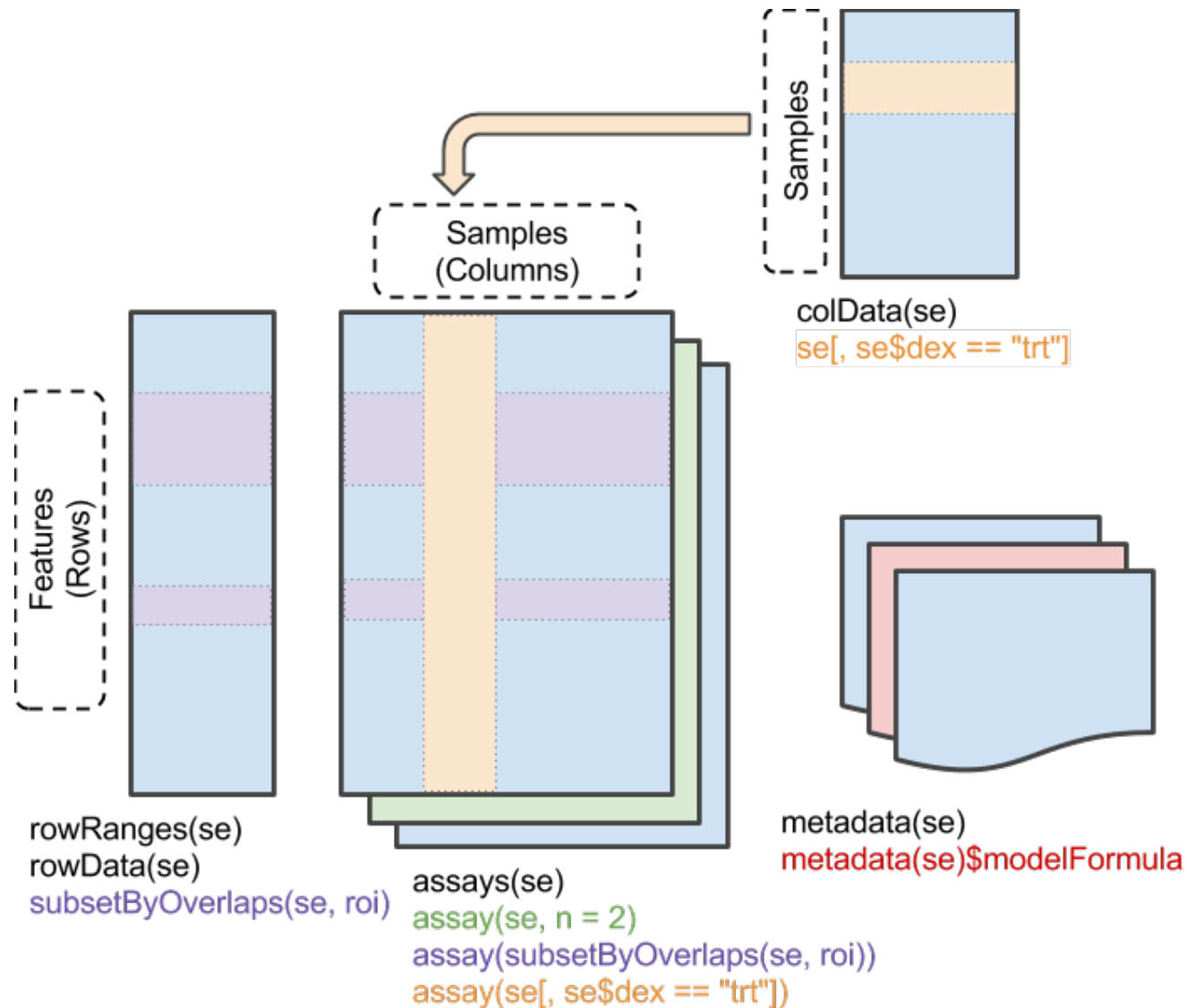
- Provide a compelling **user** experience:
documentation, demos, tutorials
 - workflows
 - package vignettes
 - function manual pages
- Support active & open **developer** community
 - training on software development & programming techniques
 - distributed development by domain experts (→ interoperability)
 - common data structures that enable workflows integrating multiple data types and disciplines



Workflows for HT Sequencing



Bioconductor Classes: e.g., Summarized Experiment



Bioconductor Classes: e.g., Summarized Experiment



- Validity checking: enforces contracts with the user
- Synchronized handling of multiple tables / matrices (subsetting, dangling pointers)
- Encapsulation: separation of interface from implementation
- Specialized highly efficient methods for manipulation (e.g. GRanges class)

```
subsetByOverlaps(se, roi) assays(se)  
assay(se, n = 2) assay(subsetByOverlaps(se, roi))  
assay(se[, se$dex == "trt"])
```

Annotation & Datasets

Annotation of genes, transcripts, proteins, pathways, metabolites; GO, Reactome, Pubmed, ...

Sequences

You don't have to download text files from NCBI / EBI and parse them into R - use ready made packages with nice interfaces.

[ExperimentHub](#): published datasets already curated into efficient R objects, with documentation

Modes of documentation

- Manual Pages (for each function)
- Vignettes: Narrative overviews on what you can do with a package
- Workflows: end-to-end descriptions of a scientific question
- F1000Research papers, Bioinformatics application notes: peer-reviewed, citable

Support Forum -

<http://support.bioconductor.org>

The screenshot shows the Bioconductor Support Forum interface. At the top, there's a navigation bar with links like 'My: messages', 'votes', 'posts', 'tags', 'following', and 'bookmarks'. The main header features the Bioconductor logo and navigation links: 'ASK QUESTION', 'LATEST', 'NEWS', 'JOBS', 'TUTORIALS', 'TAGS', and 'USERS'. Below the header, there's a search bar and a list of questions. Each question entry includes the number of votes, answers, and views, the question title, tags, and the author's name and reputation. The questions listed are:

- Testing for DE genes** (0 votes, 1 answer, 51 views) by Lennart.Vermader. Tags: limma, differential gene expression, affymetrix microarrays.
- read.maimages(targets, source="genepix") Error in `[.data.frame`(obj, , columns[[a]]) : undefined columns selected** (0 votes, 0 answers, 21 views) by sankhwar.madhu. Tags: limma, bioconductor, agilent microarrays.
- Feature request: syntax for dropping unused seqlevels in GRanges objects** (2 votes, 1 answer, 45 views) by Charles Plessy. Tags: seqlevels, granges, feature request.
- Job: Research Fellow or Junior Research Fellow, Statistical Genetics and Genomics, Australia** (0 votes, 0 answers, 28 views) by hong.lee. Tags: genetics, job, biostatistics.
- Bioconductor for predicting transcription factors from a gene list?** (0 votes, 0 answers, 33 views) by hkarakurt. Tags: transcription, transcriptdb.
- DESeq2 DESeqDataSetFromTximport design formula** (0 votes, 0 answers, 30 views) by stephensmith. Tag: deseq2.
- check splice isoforms** (0 votes, 0 answers, 26 views) by ry82722. Tags: deseq2, tximport.
- Experiment design contrasts any time edgeR** (5 votes, 2 answers, 63 views) by delfino.pietro. Tags: maseq, edgeR, ma, dge.
- binnedaverage - is it possible to add an na.rm option?** (0 votes, 0 answers, 37 views) by Janet Young. Tag: genomics.

On the right side, there's a 'Recent...' section with a list of replies and a 'Votes' section with a list of votes. At the bottom, there's a 'Locations' section showing the location of the user and the time taken to load the page.

Traffic: 123 users visited in the last hour

Support Forum -

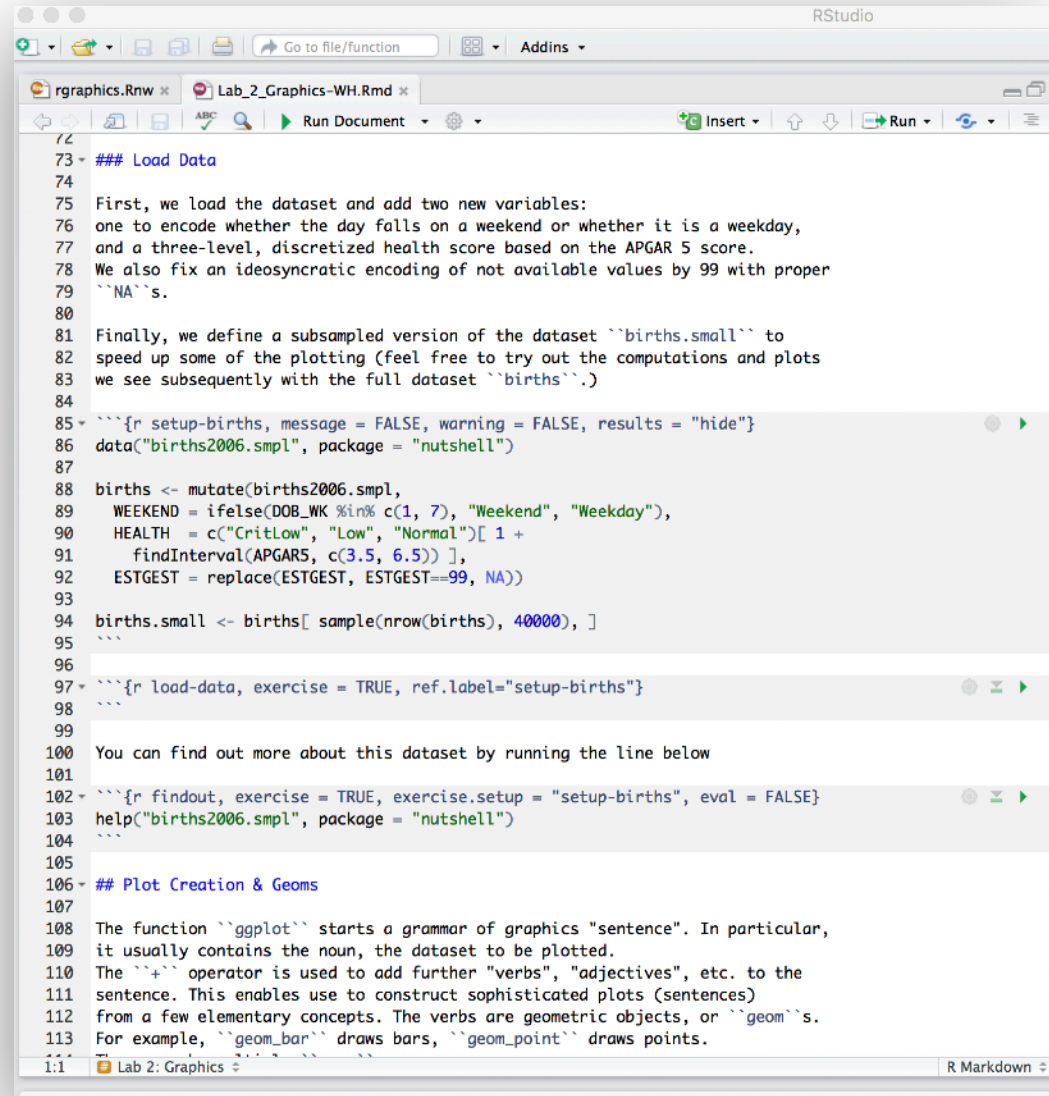
Etiquette (Posting Guide)

- Make sure you use most recent versions
- Read the documentation
- Use Google to see if a similar question has already been asked
- Prepare a minimal working example and post its code
- Remember your manners when reporting “bugs” or “missing features”
- Use descriptive subject line and precise language
- Post `'devtools::session_info()'`

<https://www.bioconductor.org/help/support/posting-guide/>

Scientific software should be assessed by similar criteria as a scientific publication

- Reproducible
- Peer-reviewed
- Easy to access by other researchers & society
- Builds on the work of others
- Others will build their work on top of it



The screenshot shows the RStudio interface with a script editor open. The script is titled 'Lab_2_Graphics-WH.Rmd' and contains R code for loading data, creating a subsampled dataset, and setting up for plotting. The code is as follows:

```
## Load Data

First, we load the dataset and add two new variables:
one to encode whether the day falls on a weekend or whether it is a weekday,
and a three-level, discretized health score based on the APGAR 5 score.
We also fix an idiosyncratic encoding of not available values by 99 with proper
`NA`s.

Finally, we define a subsampled version of the dataset `births.small` to
speed up some of the plotting (feel free to try out the computations and plots
we see subsequently with the full dataset `births`.)

```{r setup-births, message = FALSE, warning = FALSE, results = "hide"}
data("births2006.smpl", package = "nutshell")

births <- mutate(births2006.smpl,
 WEEKEND = ifelse(DOB_WK %in% c(1, 7), "Weekend", "Weekday"),
 HEALTH = c("CritLow", "Low", "Normal")[1 +
 findInterval(APGAR5, c(3.5, 6.5))],
 ESTGEST = replace(ESTGEST, ESTGEST==99, NA))

births.small <- births[sample(nrow(births), 40000),]
```

```{r load-data, exercise = TRUE, ref.label="setup-births"}
```

You can find out more about this dataset by running the line below

```{r findout, exercise = TRUE, exercise.setup = "setup-births", eval = FALSE}
help("births2006.smpl", package = "nutshell")
```

## Plot Creation & Geoms

The function `ggplot` starts a grammar of graphics "sentence". In particular,
it usually contains the noun, the dataset to be plotted.
The `+` operator is used to add further "verbs", "adjectives", etc. to the
sentence. This enables use to construct sophisticated plots (sentences)
from a few elementary concepts. The verbs are geometric objects, or `geom`s.
For example, `geom_bar` draws bars, `geom_point` draws points.
```

Code re-use

Writing good software is hard

Existing, well-used and maintained software contains fewer bugs

Common problems are already solved

Avoid re-implementation — produce interfaces

Focus on new things

→ Lots of package interdependencies (>1000 packages, 100s developers)

Don't reinvent the wheel

Shared code base, maintained by core team

Bioconductor already has code to:

- Read common file formats
- Represent common data types e.g. Genomic Ranges, Summarized Experiments
- Load genomes and annotation
- etc.

Let users become
developers

What are the benefits from using the Bioconductor development environment?

Standardised and powerful data structures for representing datasets incl. metadata

Many tools for data I/O and preprocessing. Access to databases of primary data and annotation (ExperimentHub)

Support for writing good documentation

Support for supporting your users



What are the benefits from using the Bioconductor development environment?

Free code review

Package system

Daily checks - continuous integration

Version control system

Release & devel branches

Six monthly release cycle

Stable version for most users,
but easy to make new features
public



Why R?

- high-level, interpreted programming language
- rapid prototyping, creativity, flexibility and reproducibility
- scientific and statistical computing capabilities
- publication quality graphics system ('grammar of graphics')
- convenient data I/O & wrangling
- mature package management system
- inter-language interfaces (C, C++, Java, JavaScript)
- lots of momentum with recent language innovations (RStudio, tidyverse, Jupyter, commercial adaptations, ...)

LISP/Scheme inside