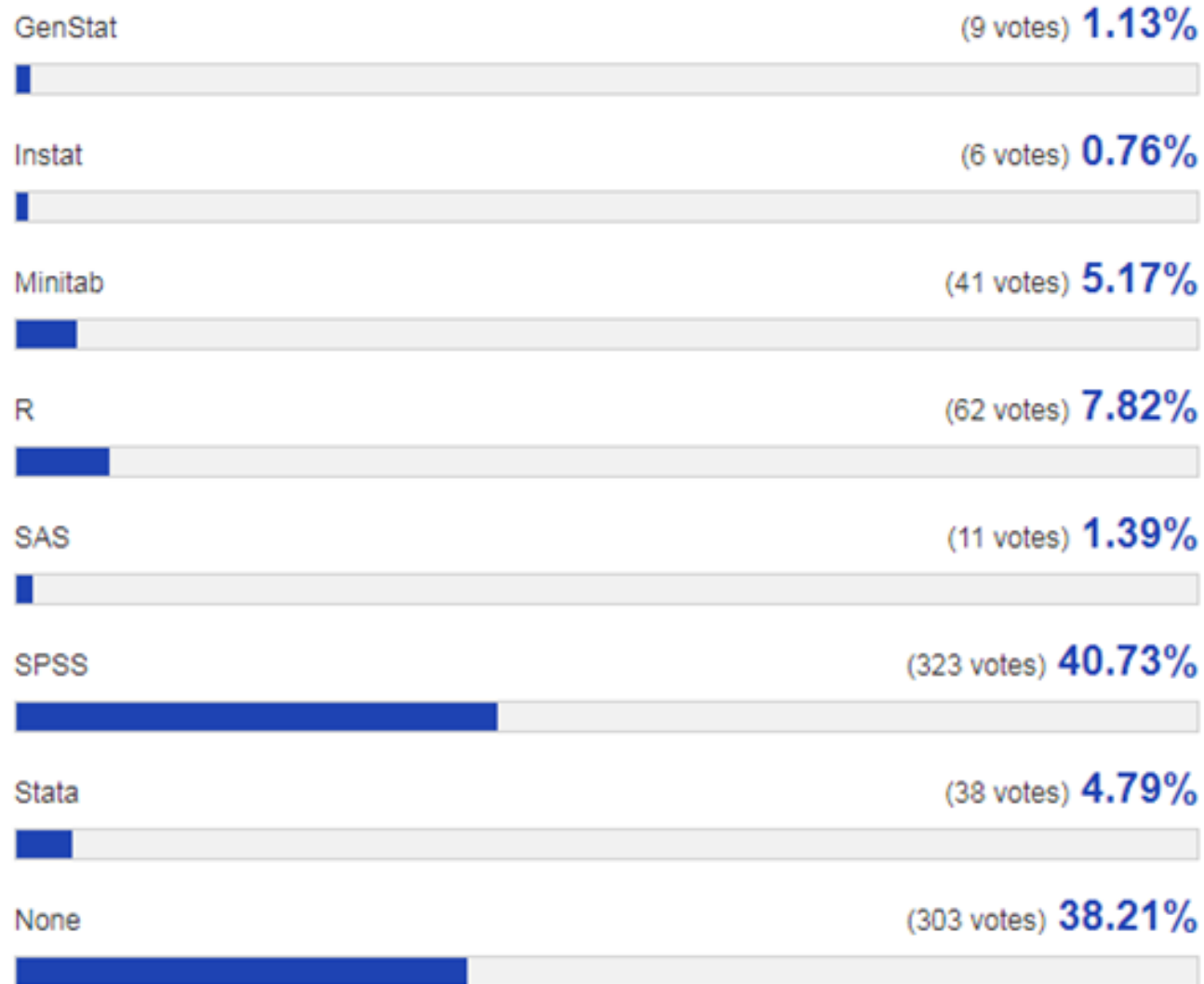# Workshop 1:
# The tidyverse and beyond

```
- take a parachute and jump
(into the tidyverse)
```
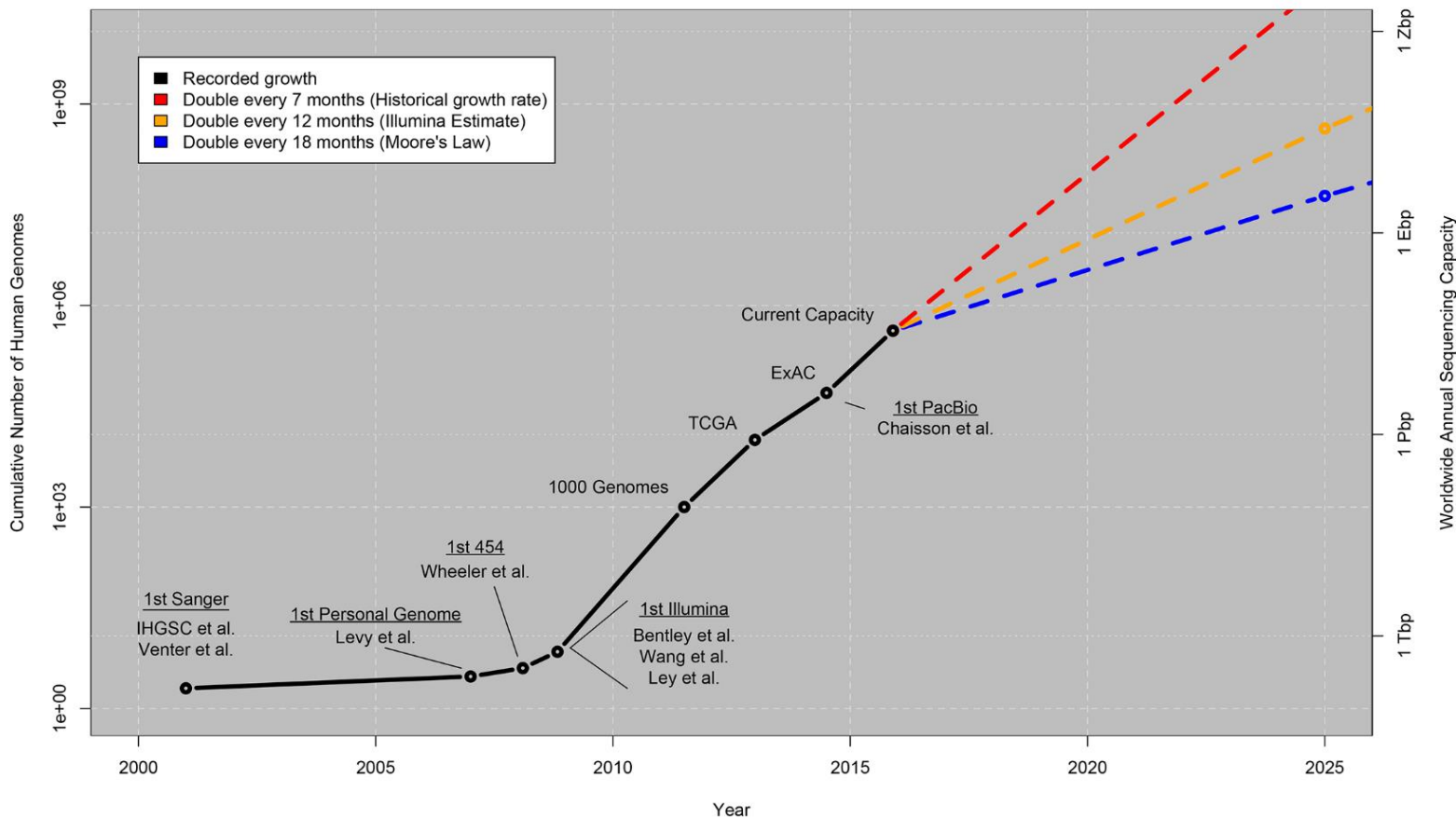
**Brendan Palmer,**

**Statistics & Data Analysis Unit,**

**Clinical Research Facility - Cork**

# Snapshot of data analysis tools used in UCC

GenStat            (9 votes) **1.13%**

Instat            (6 votes) **0.76%**

Minitab            (41 votes) **5.17%**

R            (62 votes) **7.82%**

SAS            (11 votes) **1.39%**

SPSS            (323 votes) **40.73%**

Stata            (38 votes) **4.79%**

None            (303 votes) **38.21%**

# The explosion of data in the life sciences

# R is worth learning

It's free +

    accessible +

        reproducible +

            widely used +
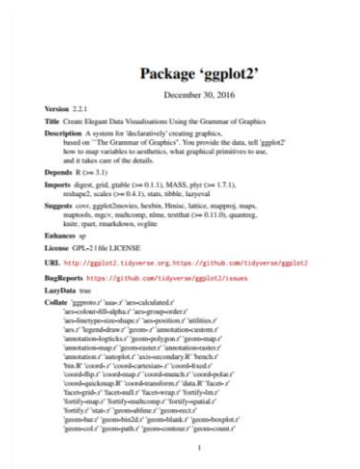
                broadly applicable +

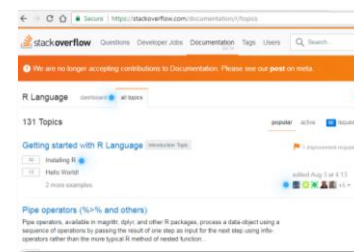                    works across platforms +

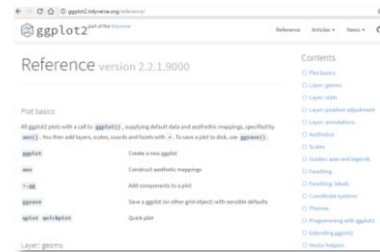                  and …………

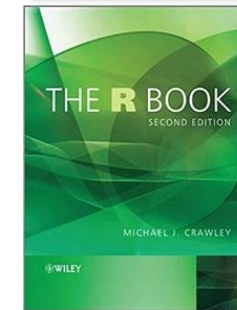# …help comes in many forms

**Vignettes**

**Webpages**

**eBooks**

**Twitter**

# Cheatsheets

Base R

The Caret Package

Data Transformation with dplyr : : CHEAT SHEET

R Markdown : : CHEAT SHEET

R Markdown Reference Guide

R Syntax Comparison : : CHEAT SHEET

Data Science in Spark with Sparklyr : : CHEAT SHEET

Data Import : : CHEAT SHEET

RStudio IDE

Data Visualization with ggplot2 : : CHEAT SHEET

Shiny : : CHEAT SHEET

Advanced R

# R functions

# R packages

functions

documentation

sample data

R comes pre-loaded with ~30 other packages (e.g. base, stats, graphics etc.)

Other packages:
Install once
Update regularly
Load each session

tidyverse

# Data analysis in a nutshell

"The good news about computers is that they do what you tell them to do. The bad news is that they do what you tell them to do." - Ted Nelson

The tidyverse makes R code more human readable
    - it is easier to write, run and read

# Data analysis in a tidyverse nutshell

# You could write a book on that!!

# RStudio



Screenshot of the RStudio interface with four regions labeled: **Code editor** (top left), **Workspace** (top right), **R console** (bottom left), and **Files/Plots/Help** (bottom right).

# Worksheet 1
## Open ws1_script1_navigating_R_packages.R

# Basics of R code

| Symbol | What it does | Example 1 | Example 2 |
|---|---|---|---|
| <- | Assign operator Creates new objects | > x <- 5<br>> x<br>[1] 5 | > y <- "This"<br>> y<br>[1] "This" |
| c() | Helps create objects with more than one element | > v <- c(5,6,7,8)<br>> v<br>[1] 5 6 7 8 | > w <- c("This", "is", "easy! ")<br>> w<br>[1] "This" "is" "easy!" |
| # | Computer ignores what is written. Used for adding notes to code | > #print("hello")<br>> | > print("hello")<br>[1] "hello" |
| %>% | Literally translates as "then do this" | > data %>%<br>   do.something.to(data) | |
| %in% | returns a logical vector indicating if there is a match | > "x" %in% c("x", "y", "z")<br>[1] TRUE | > c("x", "y", "z") %in% "x"<br>[1]  TRUE FALSE FALSE |
| ? | Access information | > ?mean() | > ?geom_point() |

**FYI: R is case sensitive!!  Name.of.data ≠ name.of.data**

# The tidyverse package 1.2.1

```
> library(tidyverse)
-- Attaching packages ----------------------------------------- tidyverse 1.2.1 --
v ggplot2 2.2.1      v purrr   0.2.4
v tibble  1.4.2      v dplyr   0.7.4
v tidyr   0.8.0      v stringr 1.2.0
v readr   1.1.1      v forcats 0.2.0
-- Conflicts ---------------------------------------------- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
```

```
ggplot2 <- data visualisation      purrr   <- functional programming
tibble  <- data frames             dplyr   <- data manipulation
readr   <- data import             stringr <- string manipulation
tidyr   <- data tidying            forcats <- categorical variables
```

# Tidy data should satisfy the following:

Each variable forms a column

Each observation forms a row

**In Brauer et al., 2008:**
- column headers are values not variable names
- multiple variable are stored in one column

e.g. 1: the column "NAME" contains information such as;
SFB2 || ER to Golgi transport || molecular function unknown || YNL049C || 1082129
- these need to be split into new columns

e.g. 2: columns G0.05 to U0.03 identify the limiting nutrient (letter) and the growth rate (number) combinations

# Try to limit "uninformative" data

"GWEIGHT" contains the same information in every cell
       - this isn't going to add to our analysis

"GID" and "YORF" appear to be study specific IDs

"NAME" column contains a lot of information

Going back to the previous example;
SFB2 || ER to Golgi transport || molecular function unknown || YNL049C || 1082129

SFB2: Gene names, but not present in all cases
ER to Golgi transport: Biological process
molecular function unknown: Molecular function
YNL049C: Gene ID listed on public repositories
1082129: Another identifier that does not appear to be useful

# Code structure example

```
separated_gene_df <- separate(raw_gene_df, NAME, c("name", "BP", "MF", "systematic_name", "number"), sep = "\\|\\|")
```

`separated_gene_df` —the new object you will create

`<-` —the assign operator

`separate` —the function you are calling

`(raw_gene_df,` —the input data

`NAME,` —the column to be separated

`c("name", "BP", "MF", "systematic_name", "number"),` —new columns IDs

`sep = "\\|\\|")` —how the data will be split

# Worksheet 2
## Open ws1_script2_stepwise_Brauer_analysis.R

# How to plot in ggplot

```
Template:

ggplot(data = <DATA>) +

    <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>))    +

      linear model  +

      axes formatting  +

      legend formatting  +

      title   + etc. etc.
```

# Worksheet 3

**Open ws1_script3_piped_Brauer_analysis.R**

# Introductory R Workshops

~~**Week 1 (13ᵗʰ February):**~~
~~**Take a parachute and jump (into the tidyverse)**~~
~~- tidying and visualisation of NGS data~~
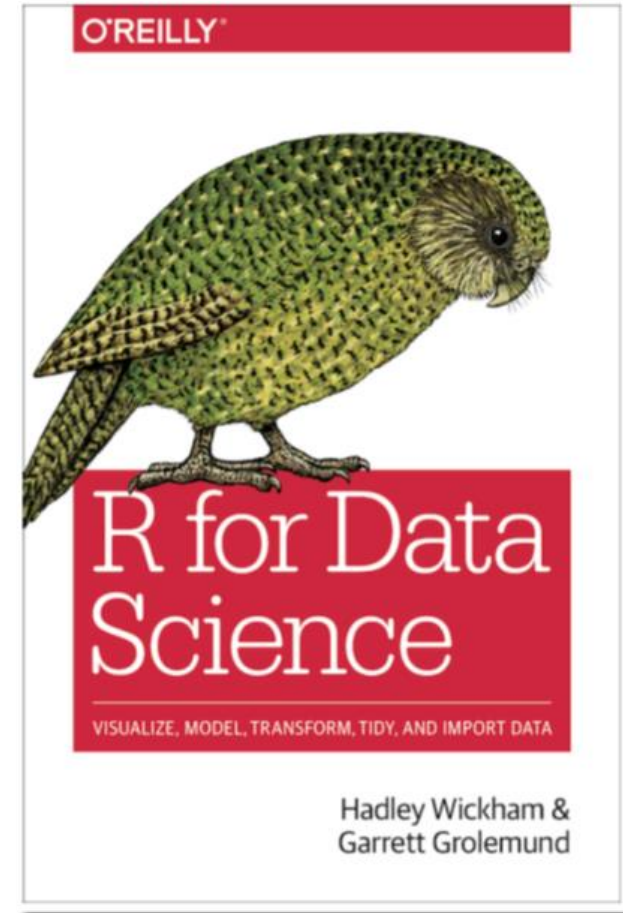~~using sample R scripts~~

**Week 2 (20ᵗʰ February):**
**We built this software on base R code**
    - overview and structure of R syntax

**Week 3 (27ᵗʰ February):**
**Sending an SOS to the world**
    - how to identify with errors in your code and get help

# Introductory R Workshops

**Week 4 (6ᵗʰ March):**
**It's the end of base R as you know it**
    - introduction to the tidyverse packages tidyr and dplyr

**Week 5 (13ᵗʰ March):**
**Welcome to the ggungle**
    - analysis and visualisation of data

**Week 6 (20ᵗʰ March):**
**Yesterday**
    - writing clear code and making your work reproducible