

# Bootcamp outline

**Day #1:** Basic command line interface,  
reproducible research, and sequence alignment

**Day #2:** Intro to R and RStudio, data input, and  
cleaning

**Day #3:** R and data manipulation, processing

**Day #4:** R and data plotting, presentations

# Illumina sequencing generates FASTQ files



**Each paired-end lane has  
four FASTQ files**

**FASTQ files contain the sequence and quality data, along with sequencer, index, and lane information.**

# Illumina sequencing analysis pipeline

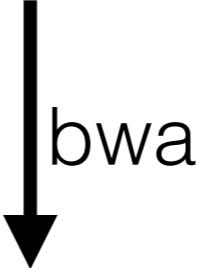
```
data — bash — 93x22
~/Desktop/IBiS-Bootcamp/data 😊 head -20 DL238.chrIII.fq
@FCC1GWRACXX:6:1103:2609:96765#AGTCAGAA/2
TCCCTAACGCTTAAGCCTAACGCTAACGCTAACGCCCTAAGCTGCTCTTACACATCTGTGGCGGAGG
+
^_a\caacgefegefihfhi[efefghifbgffffheggfhiiifggf^cgdbaaebe_bbehghhhgdgddcdbbbedbb^aa
@FCC1GWRACXX:6:1302:2023:49027#AGTCAGAA/1
GCCCTAACGCTAACGCTAACGCTAACGCCAACGCTAACGCCCTACACCTAACGCTAACGCCAACGCA
+
JPJJacJQQ\\bcdd_dd]X[^Re^J[J`QIYYdYPI^I^IIIYIH000acXBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@FCC1GWRACXX:6:1310:17686:75253#AGTCAGAA/2
TCCCTAACGCTAACGCTAACGCTAACGCCCTAACGCTAACGCCCTAACCTGTCTTACACATCTGTGGCGGAGG
+
^PZcce^c]YbRY0`bY`Z^`dfebg__`Xde[b[XX^c``cfffec^Z^`aaRW\acZMH]edgfhhdggyb]`_cZVzbBBBBBBB
@FCC1GWRACXX:6:2101:8675:28447#AGTCAGAA/1
TCCCTAACGCTAACGCTAACGCCAACGCTAACGCCCTAACGCTAACGCCAACGCTAACGCCAACGCTAAC
+
JP\ceedeggfoggfhhiiiiifdhfiibfhifhiihhiihhiihhiihhii_gbddgdd`ggee_bcdd]]]^b`zz_
~/Desktop/IBiS-Bootcamp/data 😊
```

FASTQ



```
>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus]
LCLYTHIGRNIYYGSYLYSETWNTGIMLLITMATAFMGYVLPWGQMSFWGATVITNLFSAIPYIGTNLV
EWIWGGFSVDKATLNRRFAFHFILEPLFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPFPYYTIKDFLG
LLILLLLLLALLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRSVPNKLGGVLALFSLIVL
GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFSIILAFLPIAGX
IENY
```

FASTA



Binary alignment data  
BAM file



samtools/bcftools

Variant Call File  
VCF

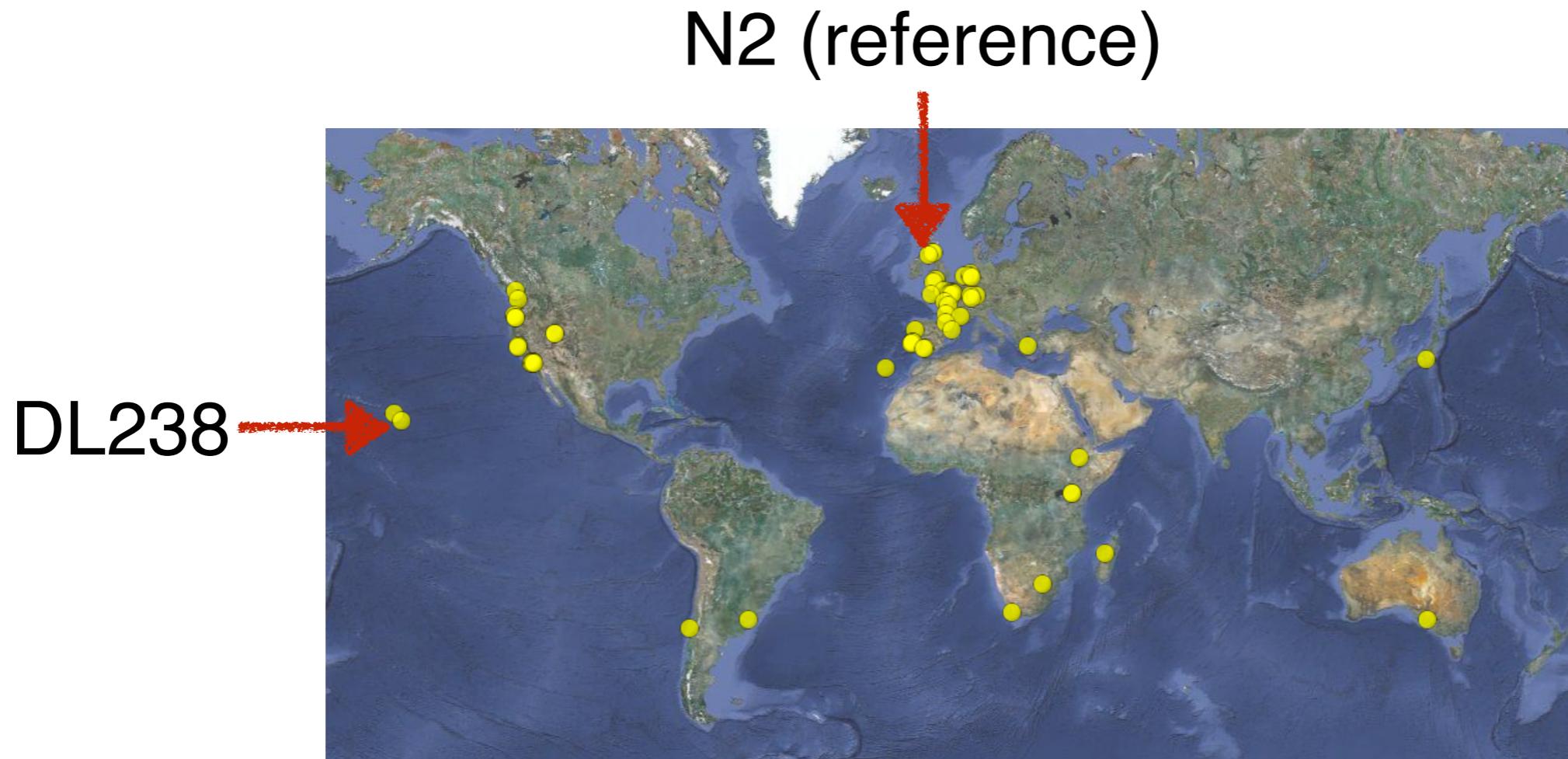
# **Let's analyze the sequence data for variation.**

Go to scripts folder in IBiS\_bootcamp and open  
DL238.Variant.Calling.sh with Atom

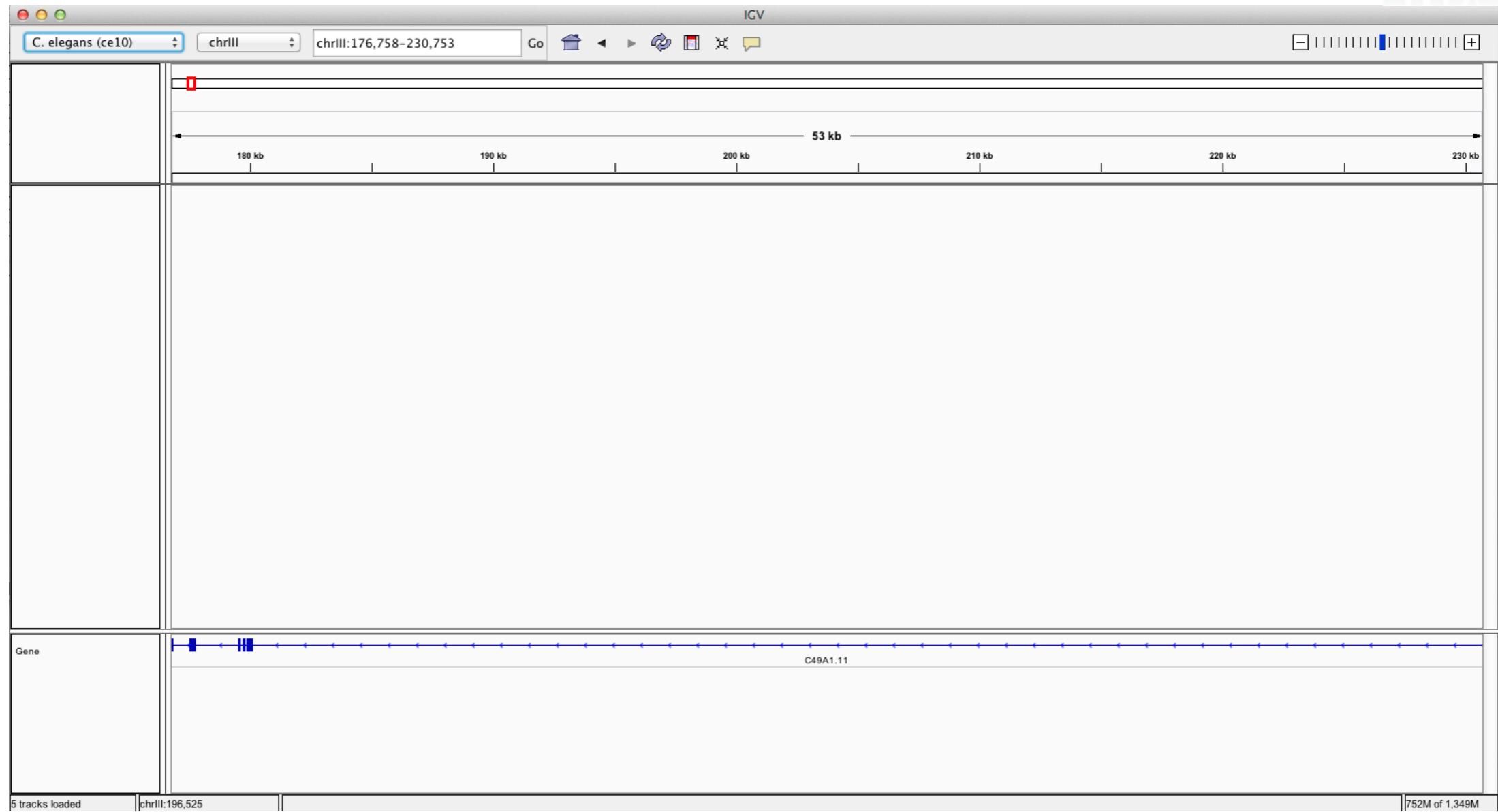
# Let's use the script to analyze the sequence data for variation.

Go to the command line  
and then go to the data folder in IBiS\_Bootcamp

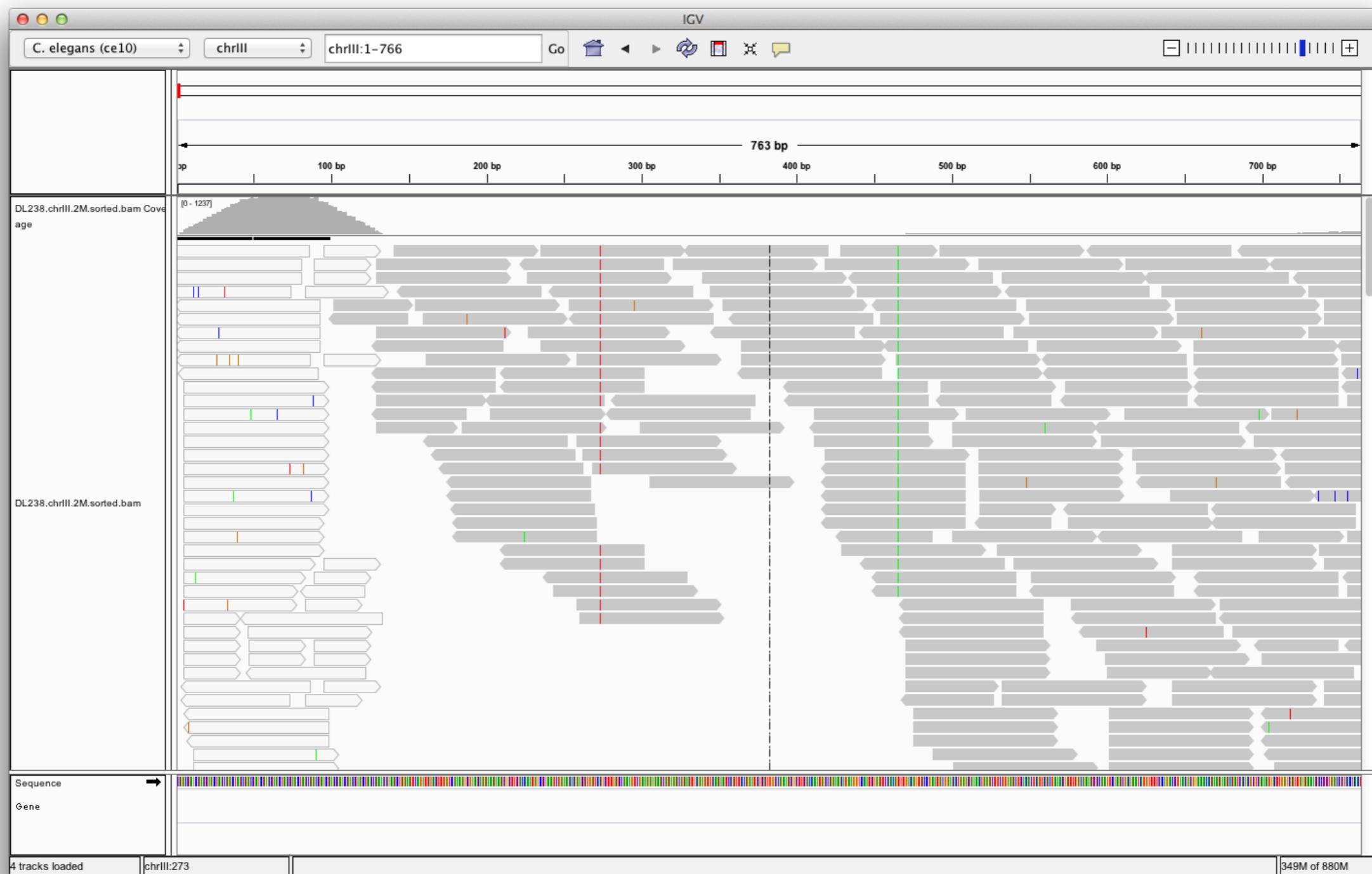
Type: bash ../scripts/DL238.variant.Calling.sh



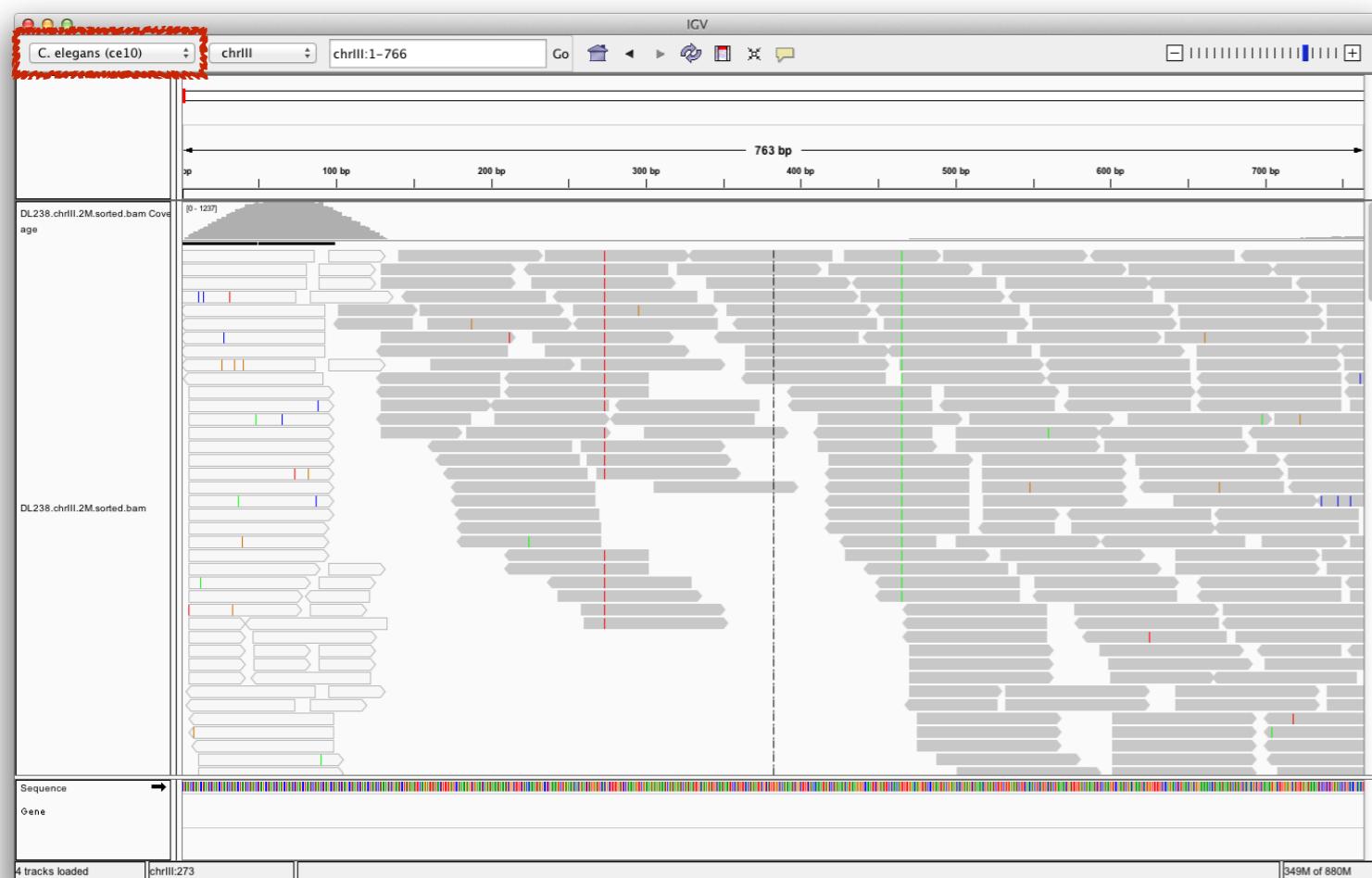
# Let's visualize the sequence data



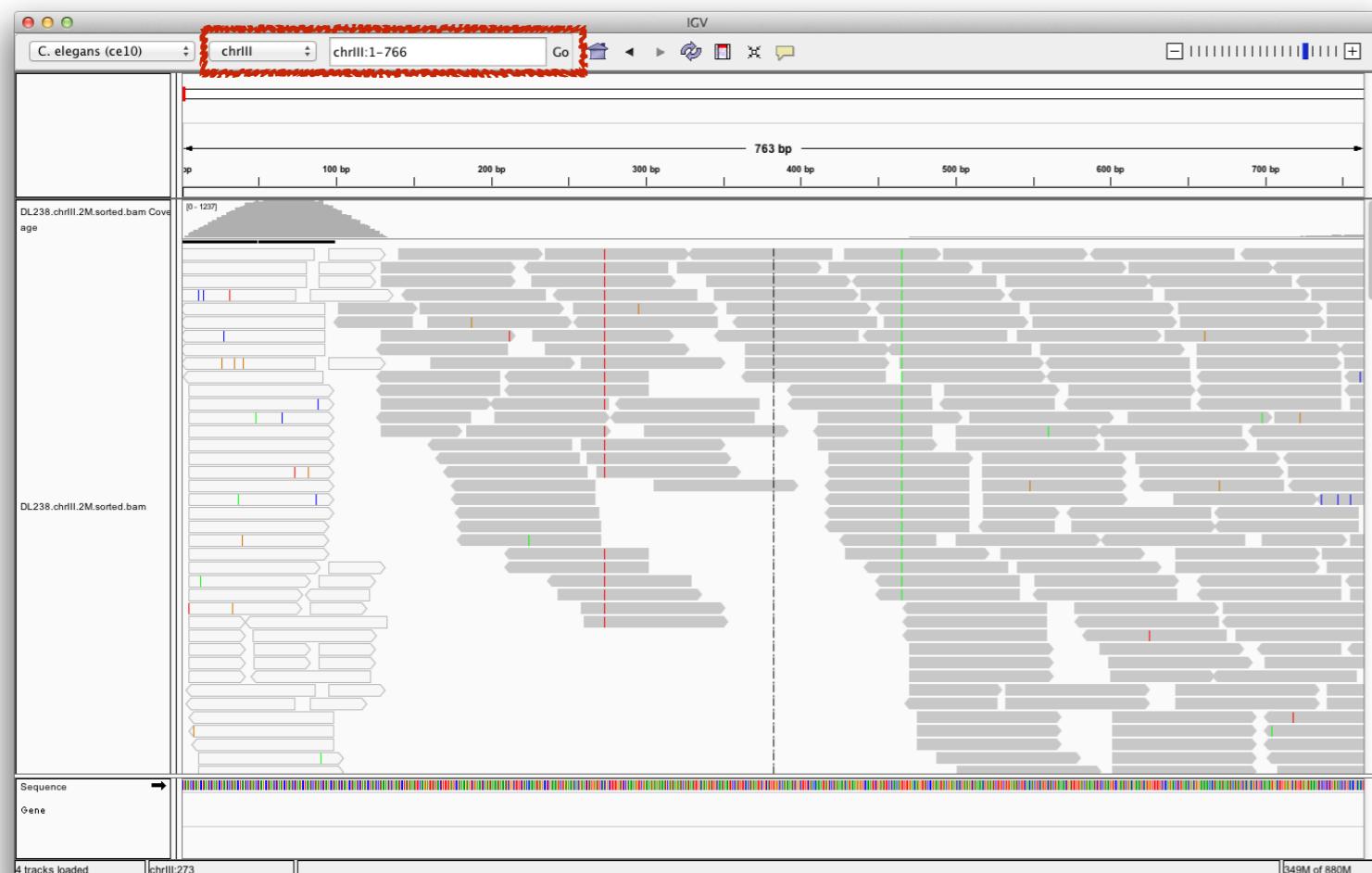
Type `igv` from your command line. IGV will start



# A Genome browser for sequence data



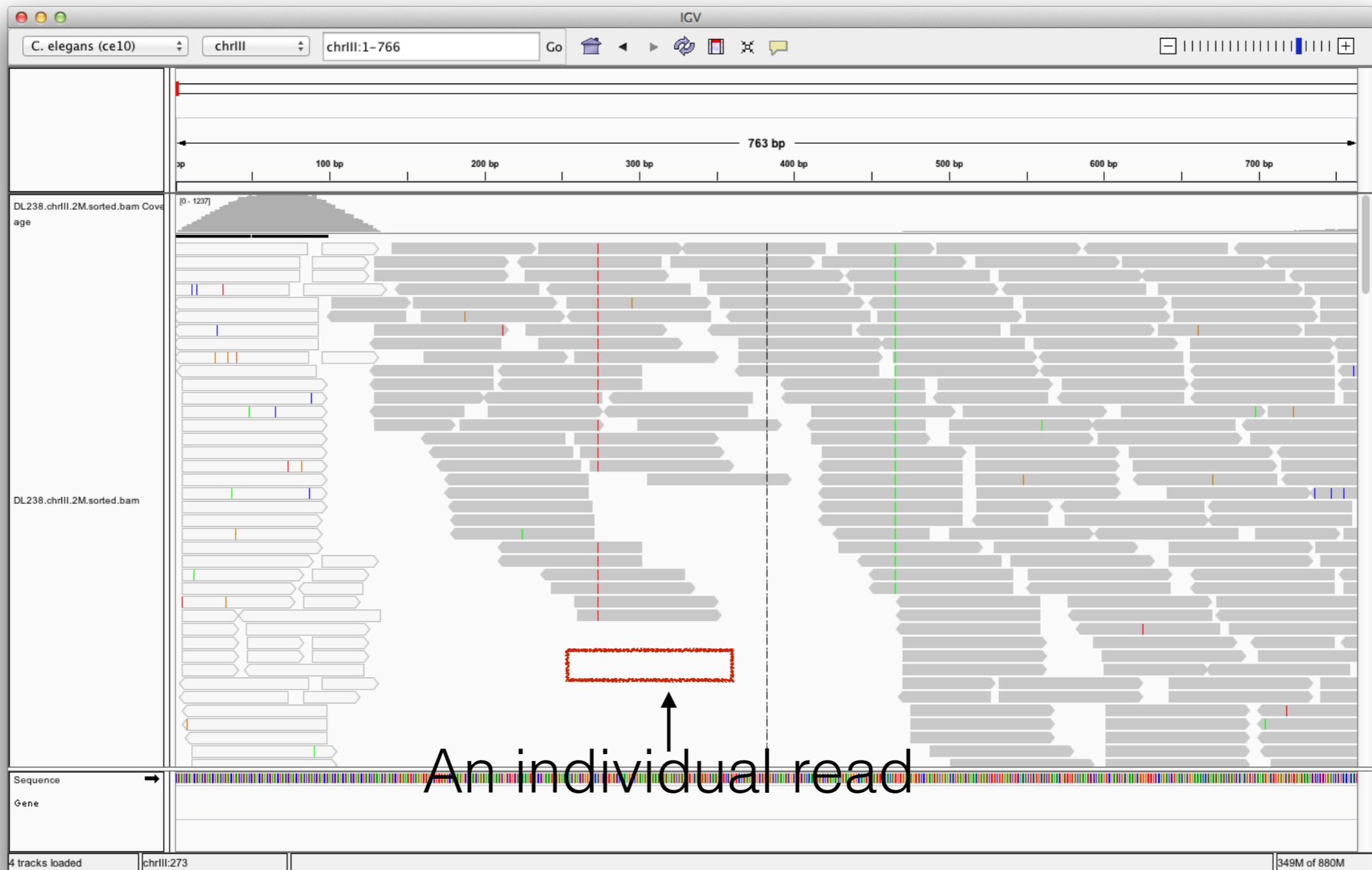
Download *C. elegans* (ce10)



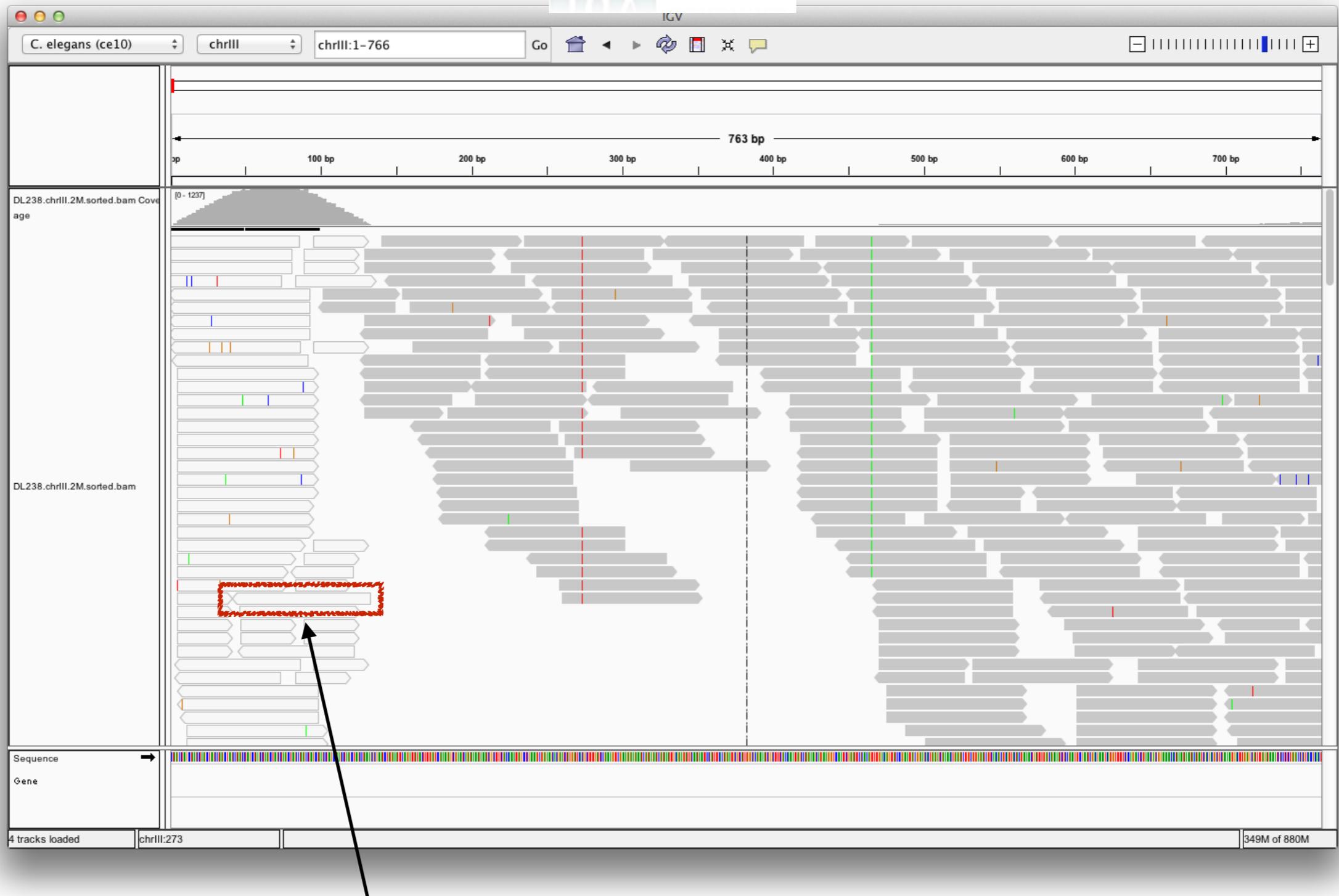
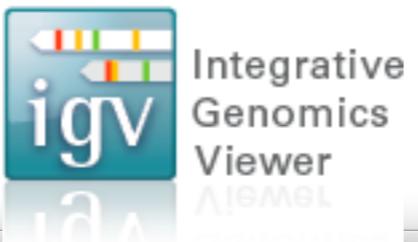
1. Drag [DL238.chrIII.2M.sorted.bam](#) into IGV
2. Navigate to chrIII:1-766



Integrative  
Genomics  
Viewer



A Genome browser for sequence data

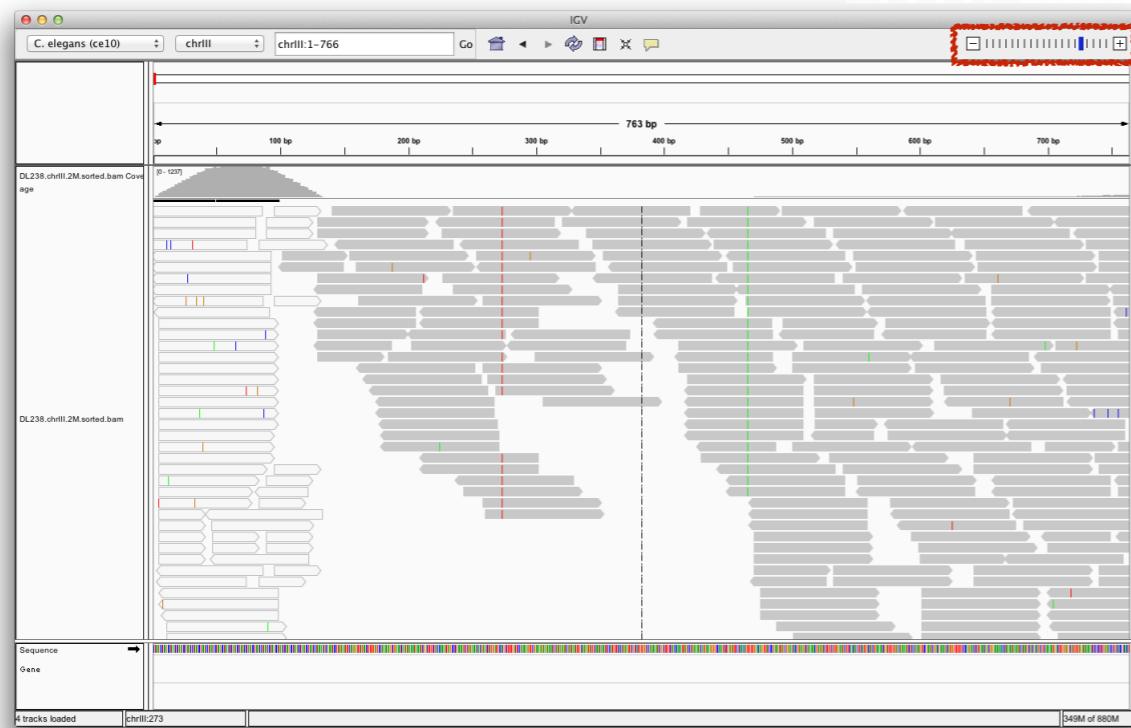


Light Grey = Mapping quality of 0.

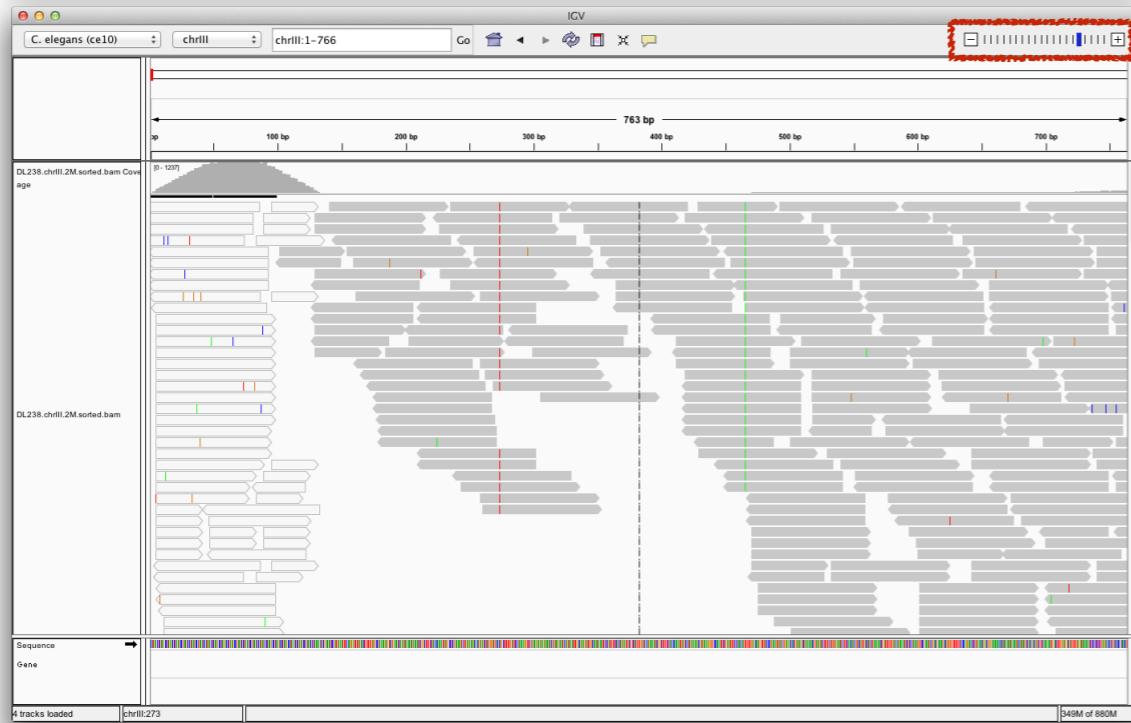
Why might the reads on the left side all be gray?



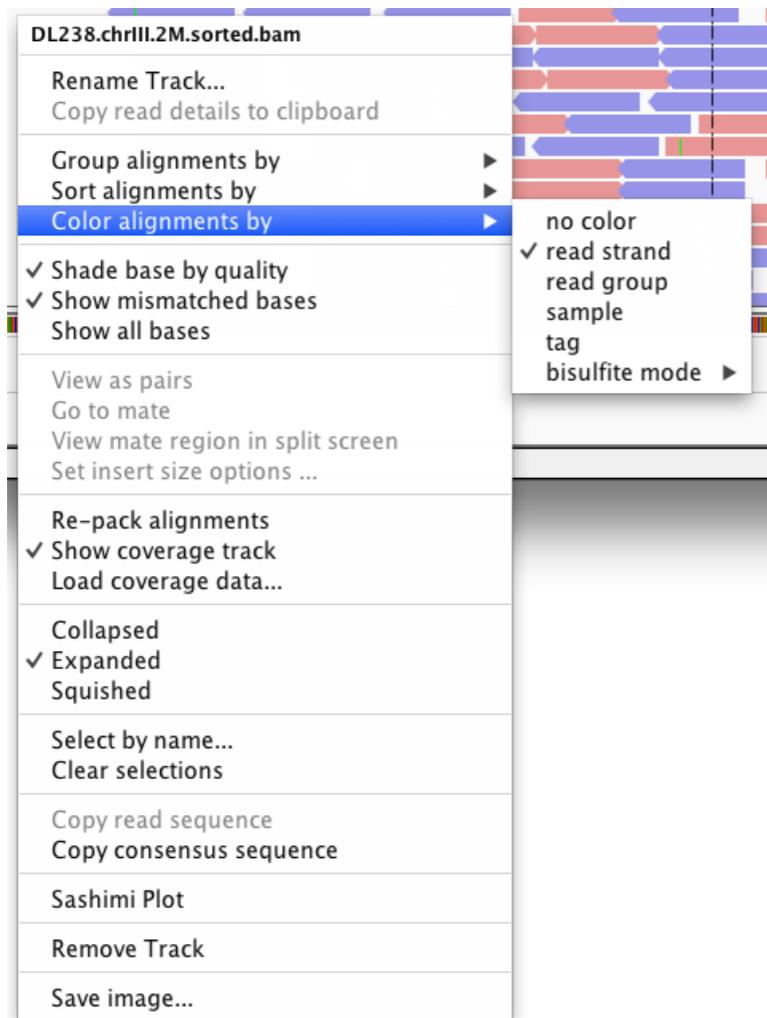
Integrative  
Genomics  
Viewer



Zoom out



Pan left and right



Right click - lots of ways  
to examine data.

Try setting size to 'squished'



We also called variants using samtools+bcftools.

1. Drag DL238.fixed.vcf.gz into IGV



2. Drag the sorted.bam file down so the vcf is on top.



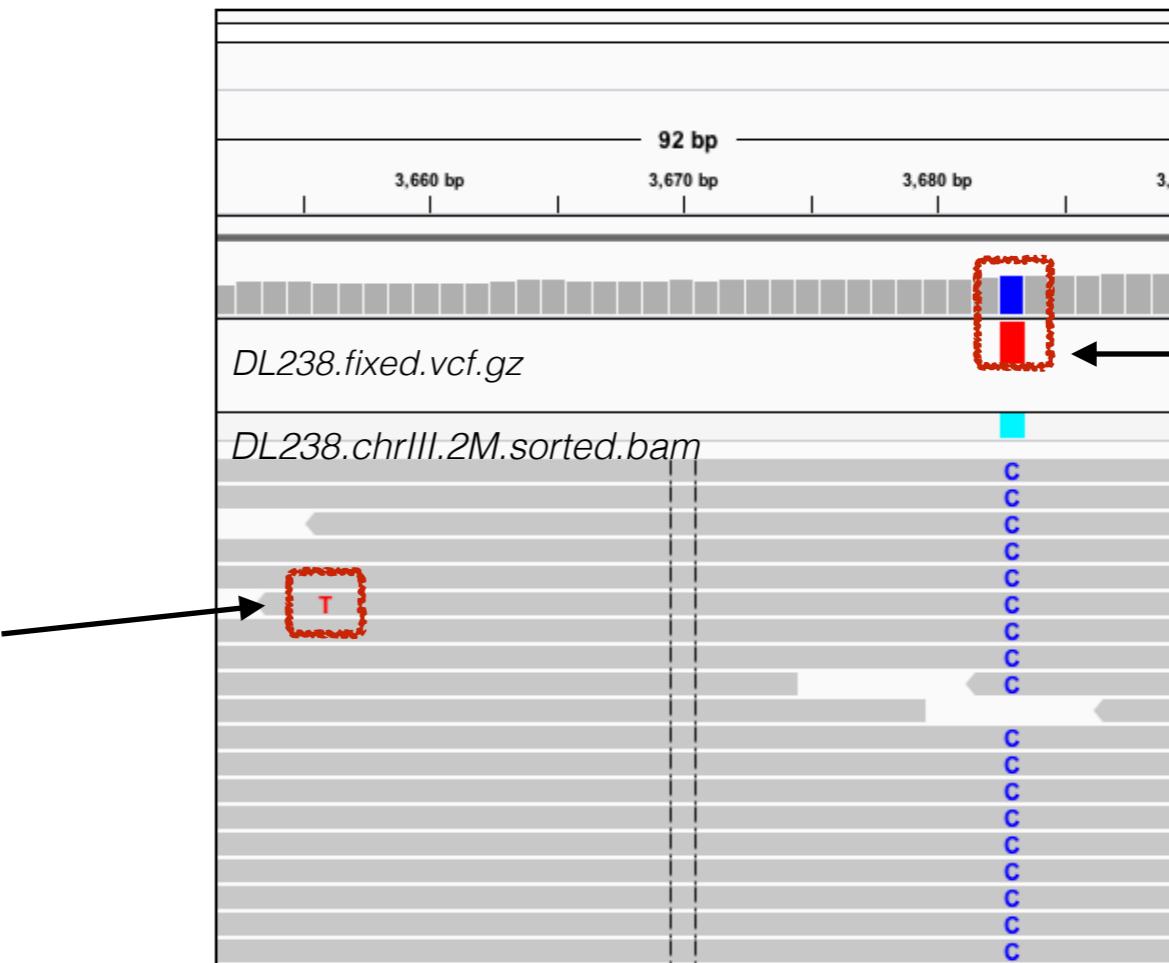
# IGV Scavenger Hunt

Within the first 5,000 base pairs, find:

Three SNPs

Two Indels

Sequencing Error



Legitimate SNP,  
Called correctly  
and supported by BAM

# Bootcamp outline

**Day #1:** Basic command line interface,  
reproducible research, and sequence alignment

**Day #2:** Intro to R and RStudio, data input, and  
cleaning

**Day #3:** R and data manipulation, processing

**Day #4:** R and data plotting, presentations

# The computational biology triumvirate



**Free and open source are core principles of science**

**You ARE learning to code!**



**is a software environment  
for statistical computing**



**is free!**

**Software**



**Cost**

\$8,700 - \$140,000 / year

\$2150 + \$1,000s for modules

**\$0**



**is managed by users!**

## **Comprehensive R Archive Network**

<http://cran.us.r-project.org>

Over 5000 free add-on *packages*

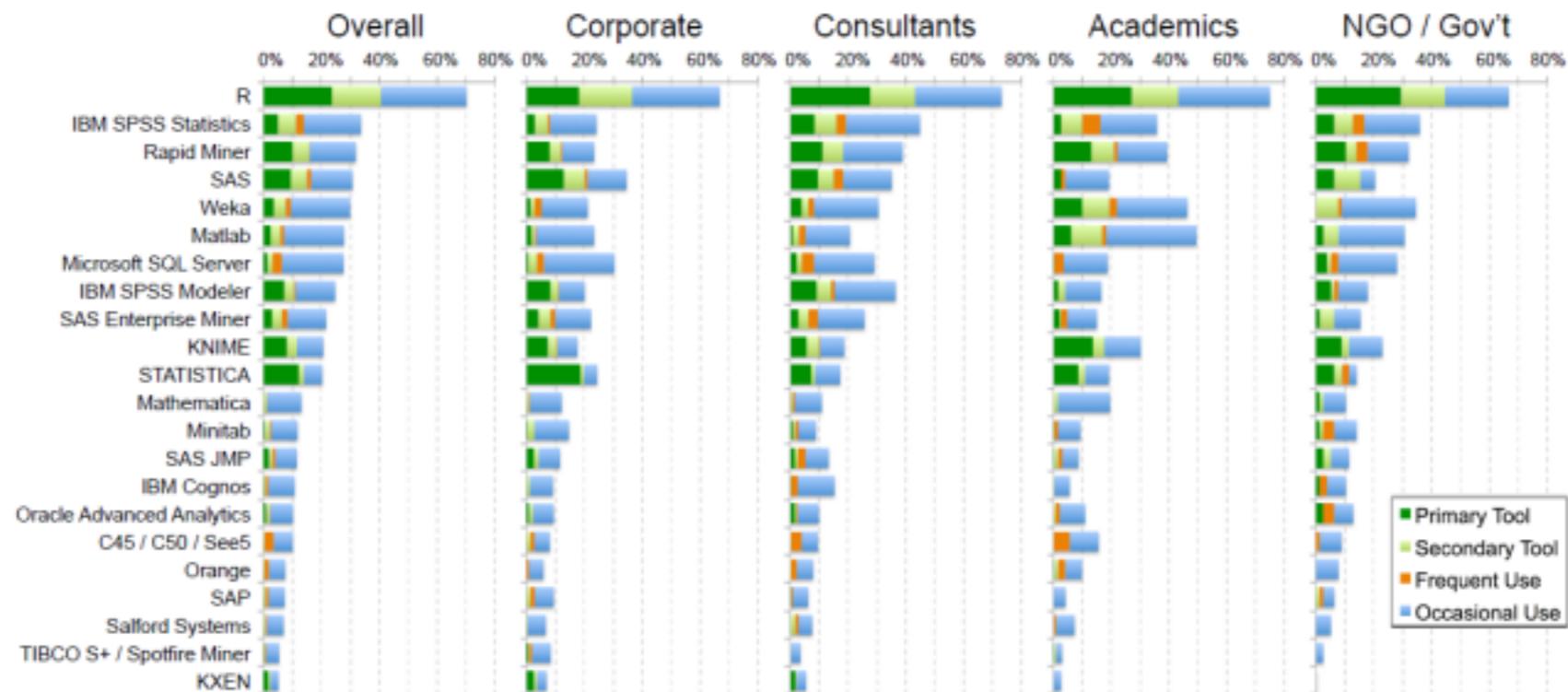
## **Stack Overflow**

<http://stackoverflow.com>

Free help



# is popular for data analysis!





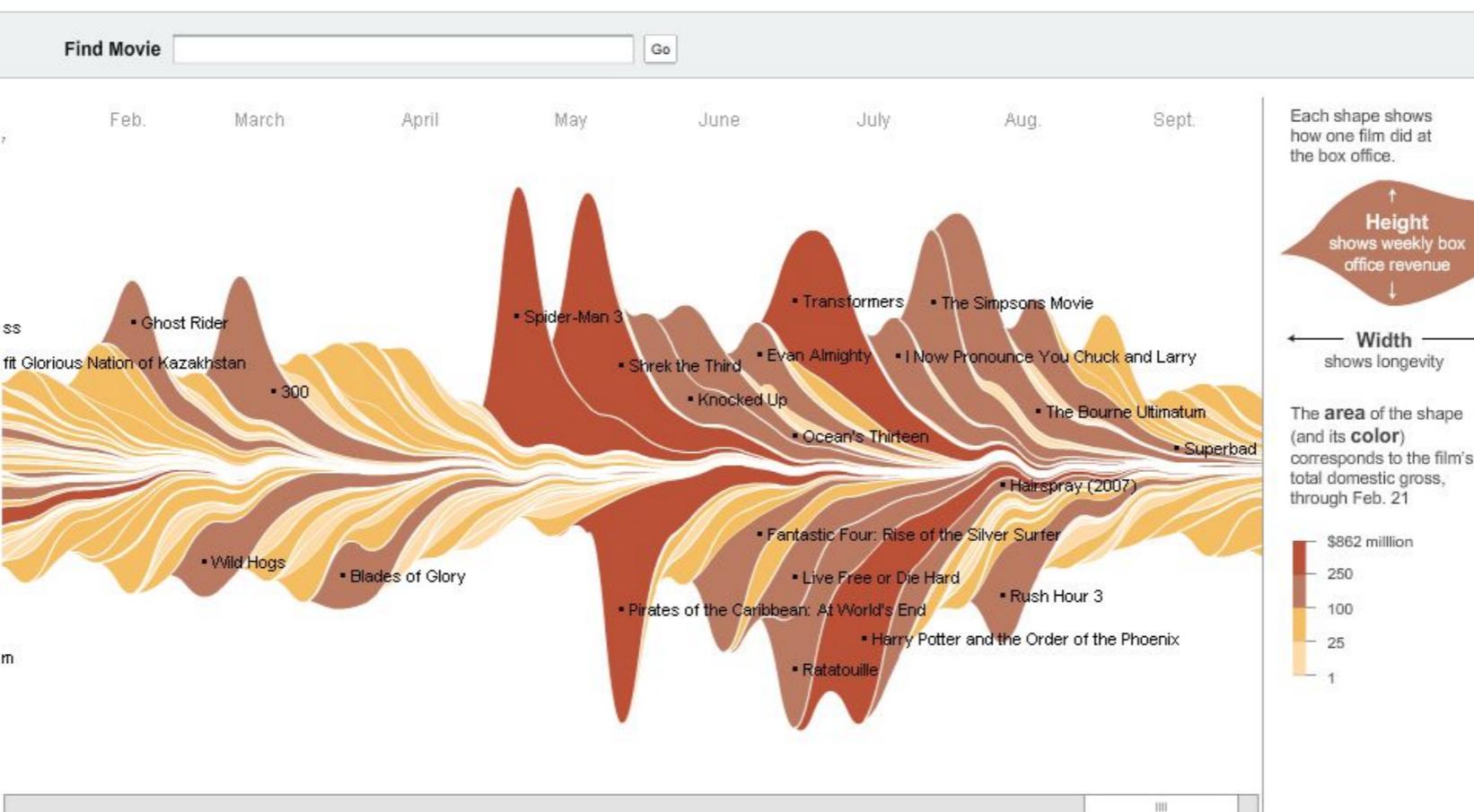
# has amazing graphics!

February 23, 2008

[SIGN IN TO E-MAIL OR SAVE THIS](#) | [FEEDBACK](#)

## The Ebb and Flow of Movies: Box Office Receipts 1986 — 2008

Summer blockbusters and holiday hits make up the bulk of box office revenue each year, while contenders for the Oscars tend to attract smaller audiences that build over time. Here's a look at how movies have fared at the box office, after adjusting for inflation.



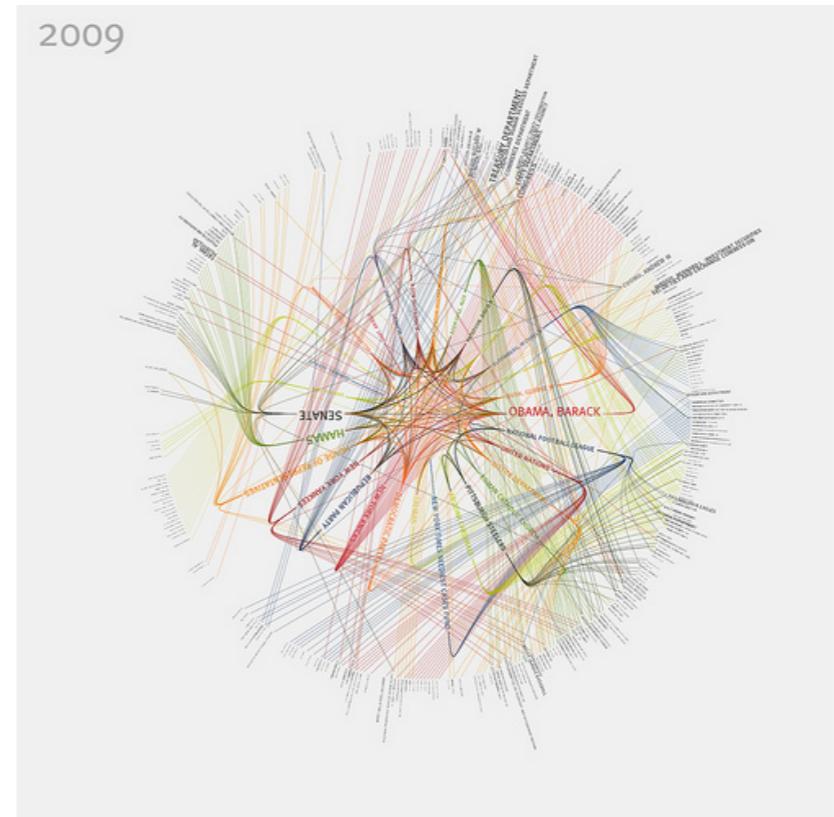
Sources: Baseline StudioSystems; Box Office Mojo

Mathew Bloch, Lee Byron, Shan Carter and Amanda Cox

## New York Times favorite tool



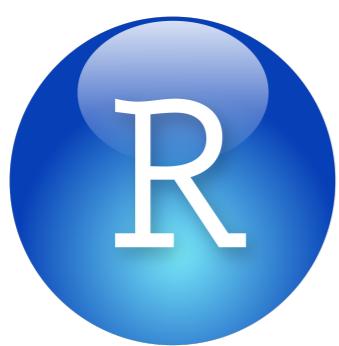
has amazing graphics!





# **is a programming language!**

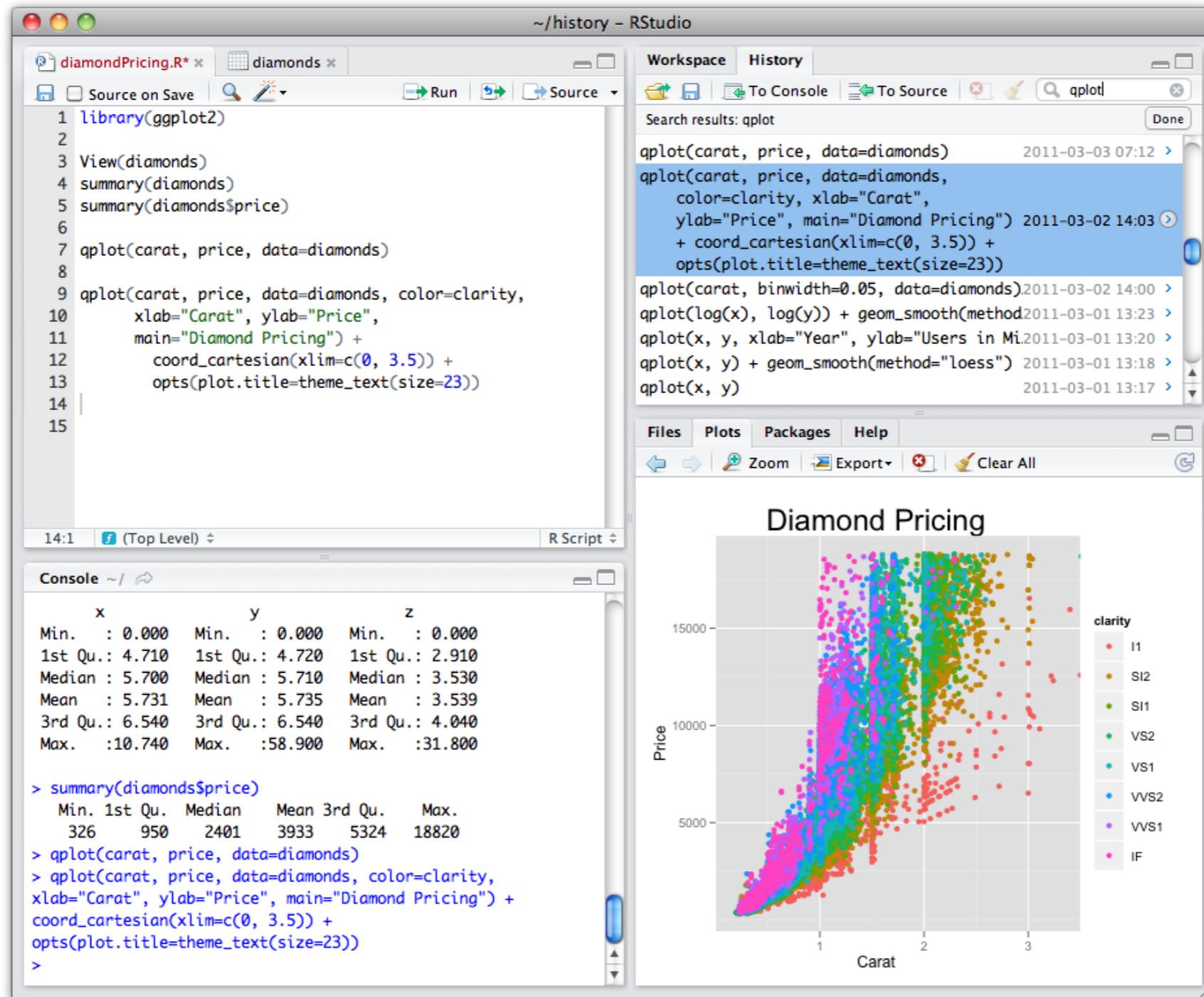
- String handling, functions, object oriented, file i/o
- Parallel programming
- Statistically orientated data structures
- Runs command line utilities, python, C, etc.
- Integrates with a variety of database systems



# RStudio makes R even easier!

Editor

Console



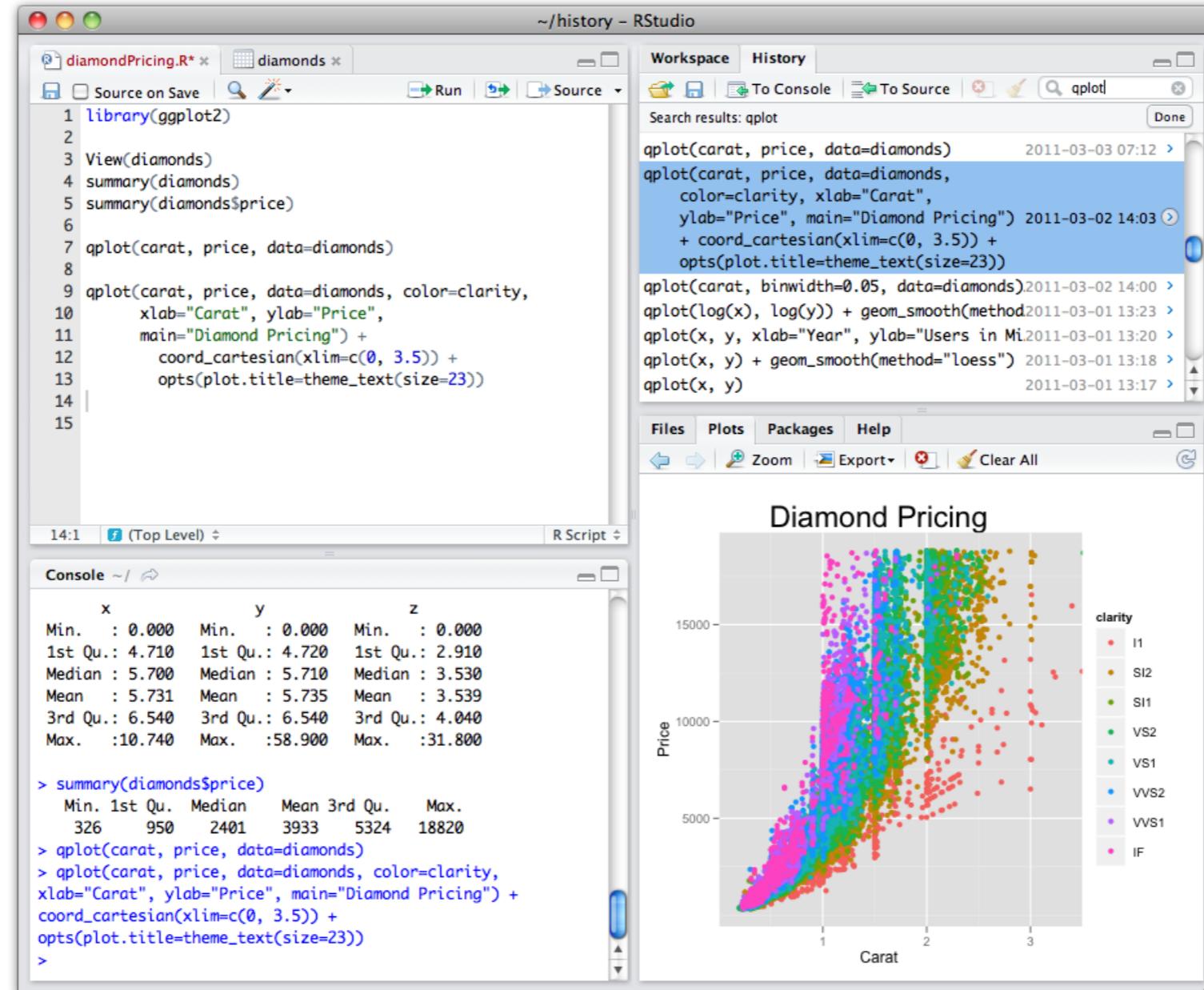
Workspace

Graphics

# Two main ways to interact with R

Editor

Console

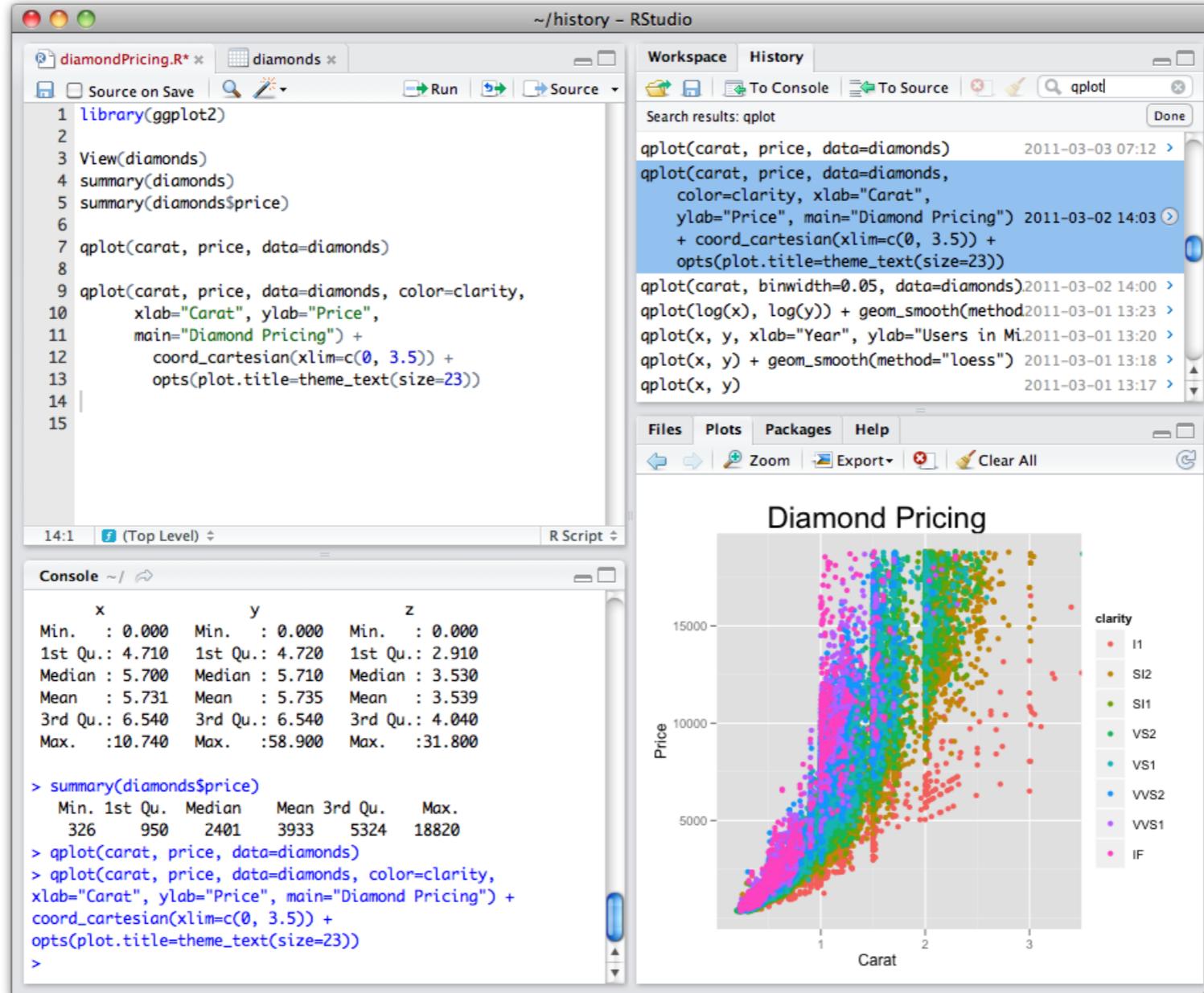


Think of the console as asking R questions about data.

# Two main ways to interact with R

Editor

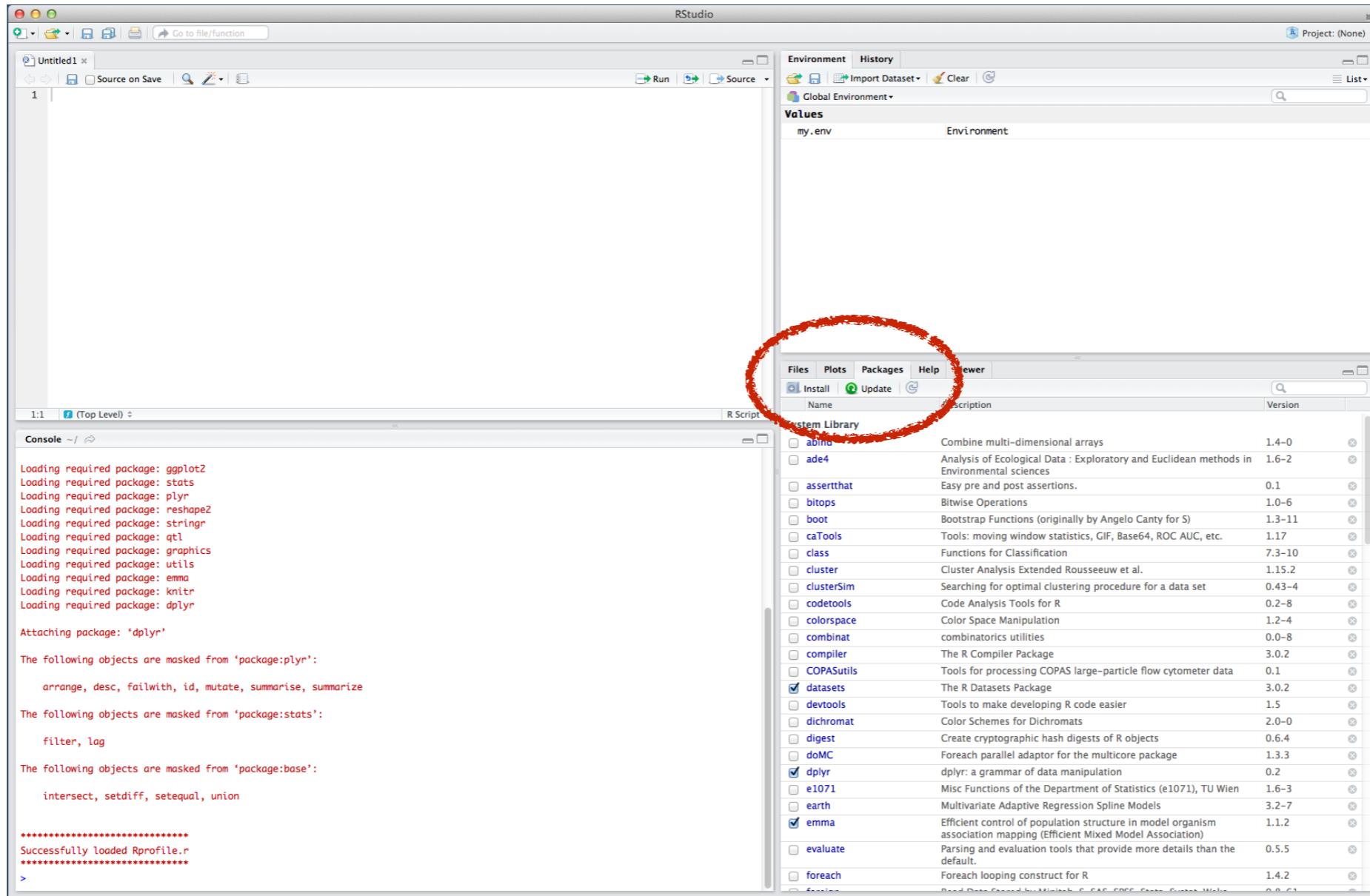
Console



The editor is a plain text file editor (like Atom)

For scripting and saving work (markdown too)

# Let's install some packages



1. Click on Packages in bottom right pane
2. Click on Install
3. Make sure the Install from drop down is CRAN
4. Enter lubridate

# Basic R syntax

> The prompt

+ at the prompt means that something isn't complete

Functions are words followed by (), i.e. `mean(x)` for calculating the mean of x

Equations are entered as assignments, i.e. `y <- mean(x)` for y is equal to the mean of x.

If you can't remember a function, use `apropos()`. For example, `apropos("var")` will display every item in the workspace with "var"

Same goes for `??var`

You can also use ?function. For example, `?mean`

Try `????mean`

# is used to enter comments not read by R

# Basic Data Types in R

## Atomic vector types

character

numeric

integer

logical

complex

raw

*atomic* = only holds data of one type

```
> # Assign 1 into y
```

```
> y <- 1
```

```
> y
```

```
> # Assign 1 to 10 into y
```

```
> y <- 1:10
```

```
> y
```

```
> str(y)
```

```
> View(y)
```

# Basic Data Types in R

## Atomic vector types

character  
numeric  
integer  
logical  
complex  
raw

*atomic* = only holds data of one type

```
> # Assign 1 into y                      > length(y)
> y <- 1                                > y <- as.numeric(y)
> y                                         > y
                                              > str(y)
                                              > as.character(y)

> # Assign 1 to 10 into y
> y <- 1:10
> y

> str(y)
> View(y)
```

# Basic Data Types in R

## Atomic vector types

character  
numeric  
integer  
logical  
complex  
raw

*atomic* = only holds data of one type

```
> length(y)          > y <- logical(5)
> y <- as.numeric(y) > y
> y                  > as.numeric(y)
> str(y)
> as.character(y)    > y <- c(1,2,3)
> str(y)
> y <- c("TRUE", "TRUE", "FALSE")
> str(y)
> y <- c("Bob", "Rick", "Ishwar")
> str(y)
```

# Basic data manipulation

Adding elements:

```
> y <- c(y, "Carole")
> y
> y <- c("Greg", y)
> y
```

Missing data (NA):

```
> z <- c(0.5, NA, 0.7)
> z
> is.na(z)
```

Sequences:

```
> series <- 1:10
> series
> seq(1:10)
> seq(from =1, to =10, by = 0.1)
> LETTERS
> letters
> as.roman(2015)
```

Special values:

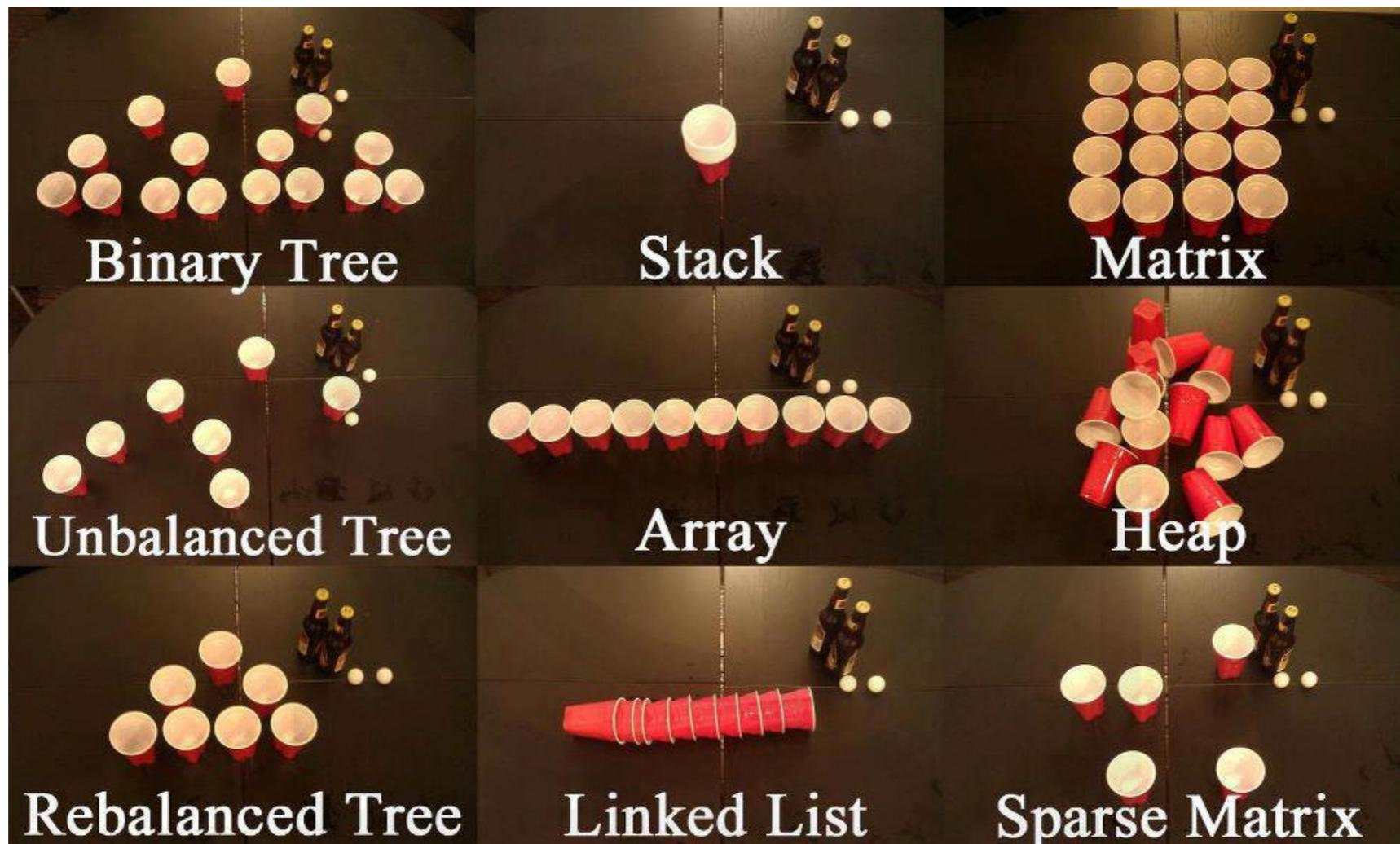
```
> 1/0
> 0/0
```

# Data structures in R

Vectors are one-dimensional data sets of any one data type

Data frames are two-dimensional data sets of any data type

Lists are groups of vectors, data frames, or other lists.



# Data structures in R

Vectors are one-dimensional data sets of any one data type

Data frames are two-dimensional data sets of any data type

Lists are groups of vectors, data frames, or other lists.

# How do we explore data objects?

`names(df)` = gives the names of columns of df

`rownames(df)` = gives the names of rows of df

`colnames(df)` = gives the names of columns of df

`dim(df)` = outputs the number of rows then columns of df

`length(df)` = outputs the length of df

`summary(df)` = outputs summary statistics of df

# Let's make our own functions

```
> add20 <- function(x) {x + 20}

> add20(10)

> # This function will convert Fahrenheit to Celsius

> FtoC <- function(fahrenheit) {
  celsius = (fahrenheit - 32) * 5 / 9
  return(celsius)}
```

# What do we do with “big data”?



# Flow of data analysis

1. Read                          Today
  2. Tidy
  3. Process
  4. Plot
  5. Present
- 
- ```
graph TD; A[1. Read] --> B[2. Tidy]; B --> C[3. Process]; C --> D[4. Plot]; D --> E[5. Present]; B --- F["{ Tomorrow }"]; C --- F; D --- G["{ Friday }"]; E --- G;
```

# We all owe Hadley Wickham a beer.



- Chief Scientist at RStudio
- Made tidyverse package to clean up data
- Made dplyr package for grammar of data manipulation
- ggplot2 package for beautiful R graphics
- Made reshape2, stringr, lubridate, devtools, plyr, rvest, etc.

# Choose your own data set for analysis

Baby names, life tables, births

```
devtools::install_github("hadley/babynames")
```

All 2015 police shootings

<https://github.com/washingtonpost/data-police-shootings>

Electric power consumption

Go To: <https://goo.gl/rGS5GN>

All UFO sightings

Go To: <https://goo.gl/R6KmP2>

Every Daily Show guest with Jon Stewart

Google: "analysis of guests Jon Stewart Stephen Turner", go to RPubs link

Query IMDB and pull data

```
devtools::install_github("hrbrmstr/omdbapi")
```

All LEGO sets from 1970-2014

```
devtools::install_github("seankross/lego")
```

**Others...you find them...the web is filled with data  
Check out:**

**[http://koaning.io/posts/fun\\_datasets.html](http://koaning.io/posts/fun_datasets.html)**

**<https://data.cityofchicago.org/>**

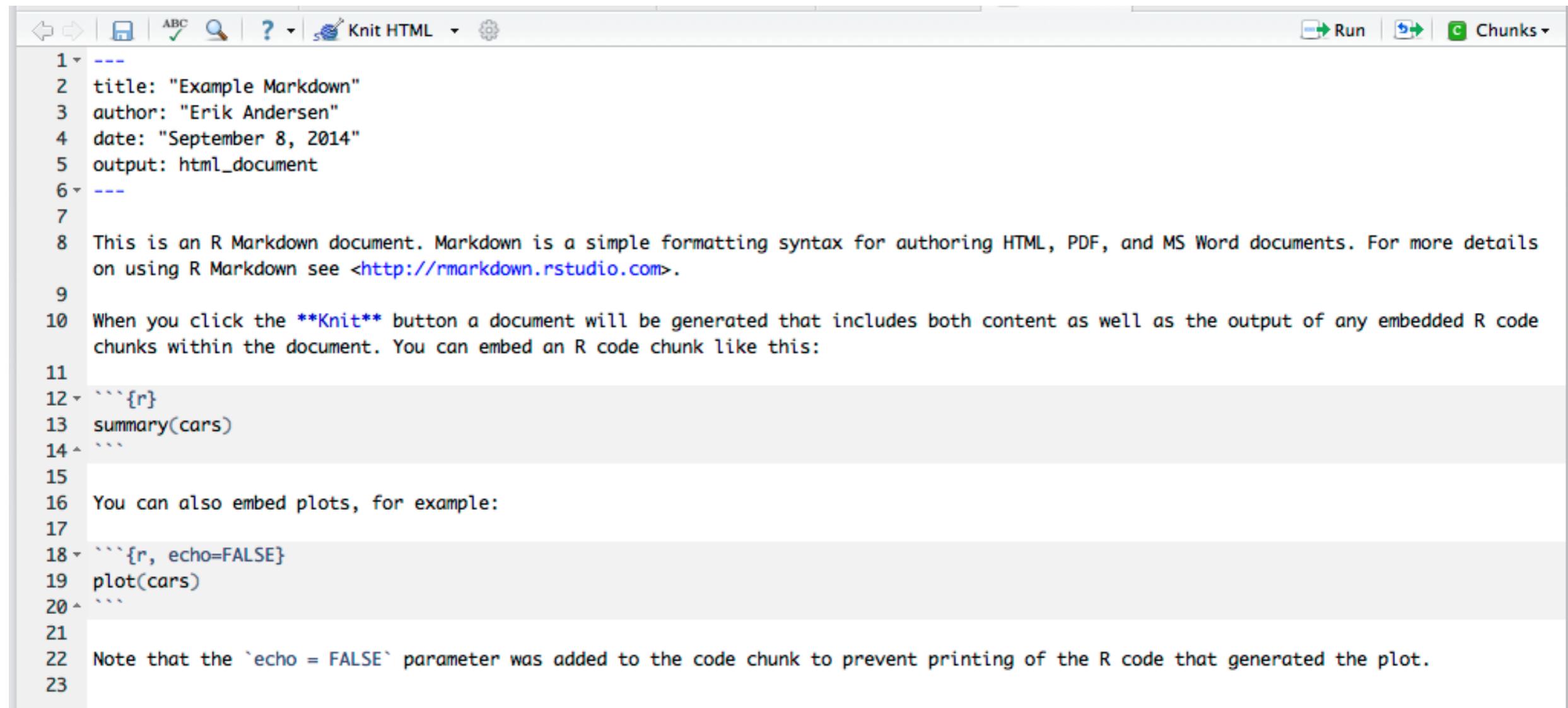
**<https://github.com/datasets>**

**<https://github.com/caesar0301/awesome-public-datasets>**

# Time to find some data



# RStudio can generate markdown reports

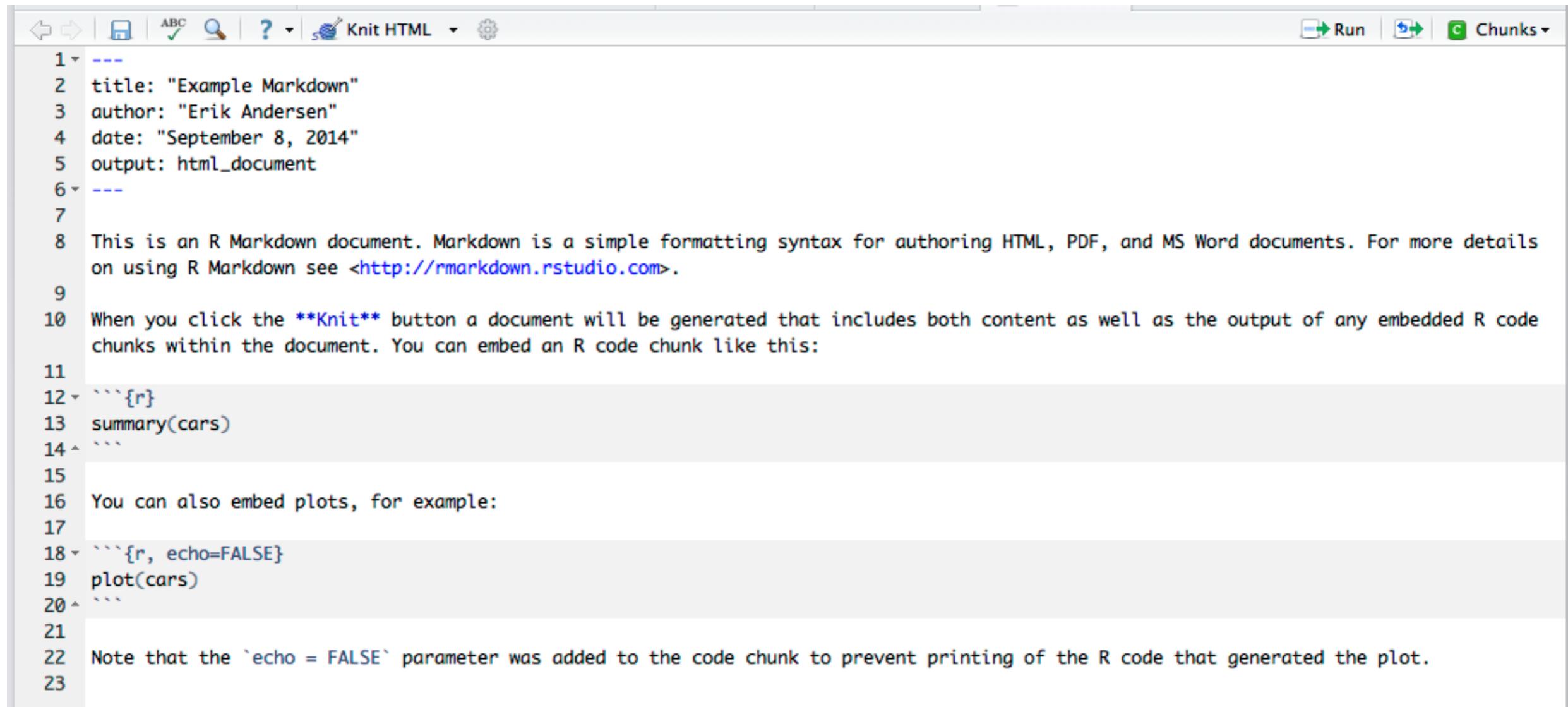


The screenshot shows the RStudio IDE interface with an R Markdown file open. The toolbar at the top includes icons for back, forward, file, search, help, and knit HTML. The main editor area displays the following R Markdown code:

```
1 ---  
2 title: "Example Markdown"  
3 author: "Erik Andersen"  
4 date: "September 8, 2014"  
5 output: html_document  
6 ---  
7  
8 This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.  
9  
10 When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:  
11  
12 ```{r}  
13 summary(cars)  
14```  
15  
16 You can also embed plots, for example:  
17  
18 ```{r, echo=FALSE}  
19 plot(cars)  
20```  
21  
22 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.  
23
```

- 1. Go to File, New File, R Markdown**
- 2. Install the necessary packages**
- 3. Write your R code and Markdown text**
- 4. Knit an HTML report**

# Make a markdown for your data analysis project



The screenshot shows the RStudio interface with an R Markdown file open. The toolbar at the top includes icons for back, forward, ABC, search, help, Knit HTML (which is currently selected), and run. The code editor displays the following R Markdown content:

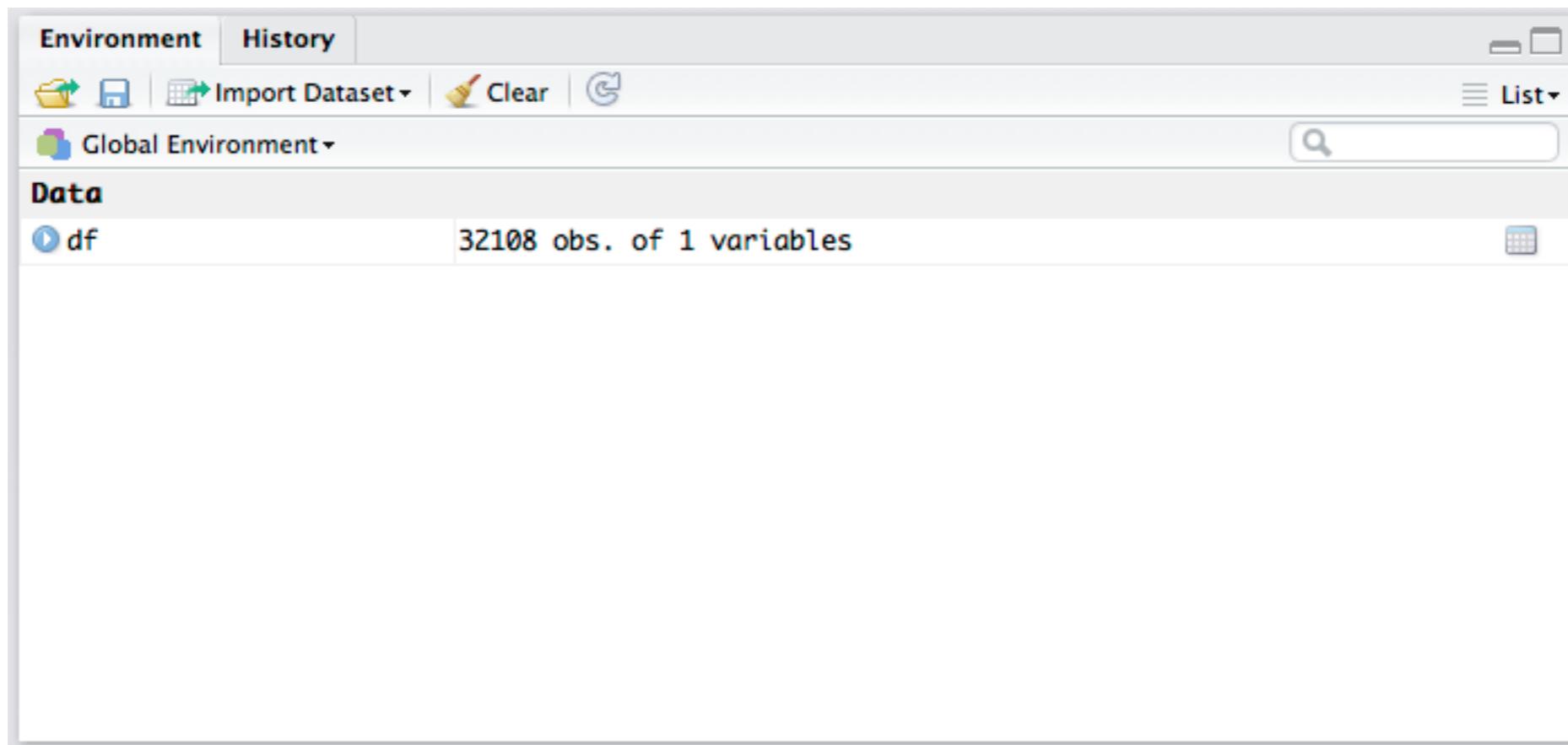
```
1 ---  
2 title: "Example Markdown"  
3 author: "Erik Andersen"  
4 date: "September 8, 2014"  
5 output: html_document  
6 ---  
7  
8 This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details  
on using R Markdown see <http://rmarkdown.rstudio.com>.  
9  
10 When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code  
chunks within the document. You can embed an R code chunk like this:  
11  
12 ```{r}  
13 summary(cars)  
14 ````  
15  
16 You can also embed plots, for example:  
17  
18 ```{r, echo=FALSE}  
19 plot(cars)  
20 ````  
21  
22 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.  
23
```

- 1. Write a short description of the data and ideas you might have**
- 2. Take a chunk to read in the data**
- 3. Think of analyses you might want to do tomorrow  
(time series, averages, etc.)**
- 4. Prepare to clean the data!**

# Reading in data

**read.delim()** is a generic function to read in simple text files, like csv or tsv files.

```
df <- read.delim(file("~/Desktop/IBiS-Bootcamp/data/sorter_data.txt", header=TRUE, sep="\t")
```



A new data object is in your workspace. Click on it. Type `df` to display all of it.

Try `head(df)` and `tail(df)`.

Try `summary(df)`.

# Reading in data

Other packages have much more powerful data reading capabilities.

Try `readr`, `read_csv` or `read_tsv()`

Try `rio`, `import()`

A new data object is in your workspace. Click on it. Type `df` to display all of it.

Try `head(df)` and `tail(df)`.

Try `summary(df)`.

# Time to read in data



# Observing data

It looks like the tail has some data that do not match the pattern of the rest of the file.

```
df <- read.delim(file "~/Desktop/IBiS-Bootcamp/data/sorter_data.txt", header=TRUE, sep="\t", nrows=32093)
```

Type `str(df)` to look at the structure of the file

```
Console ~/ 
' ends field 1 on line 1 when detecting types: Green PMT voltage 600
> str(df)
'data.frame': 32093 obs. of 27 variables:
 $ Id      : int 0 1 2 3 4 5 6 7 8 9 ...
 $ Plate   : int 1 1 1 1 1 1 1 1 1 1 ...
 $ Row     : chr "A" "A" "A" "A" ...
 $ Column  : int 1 1 1 1 1 1 1 1 1 1 ...
 $ Clog    : chr "N" "N" "N" "N" ...
 $ Scan.rate: int 2500 2500 2500 2500 2500 2500 2500 2500 2500 2500 ...
 $ Status.sort: int 0 0 0 0 0 0 0 0 0 0 ...
 $ Status.sel : int 40 40 40 40 40 40 40 40 40 40 ...
 $ TOF      : int 52 556 88 639 531 34 247 152 96 79 ...
 $ EXT      : int 120 2375 130 2697 2520 36 534 474 210 178 ...
 $ Green    : int 2 148 3 205 151 1 16 66 20 4 ...
 $ Yellow   : int 3 193 5 241 170 3 21 26 7 5 ...
 $ Red      : int 0 90 2 102 71 2 14 11 4 2 ...
 $ PH.Ext   : int 0 0 0 0 0 0 0 0 0 ...
 $ PW.Ext   : int 0 0 0 0 0 0 0 0 0 ...
 $ PC.Ext   : int 0 0 0 0 0 0 0 0 0 ...
 $ PH.Green : int 0 0 0 0 0 0 0 0 0 ...
 $ PW.Green : int 0 0 0 0 0 0 0 0 0 ...
 $ PCGreen  : int 0 0 0 0 0 0 0 0 0 ...
 $ PH.Yellow: int 0 0 0 0 0 0 0 0 0 ...
 $ PW.Yellow: int 0 0 0 0 0 0 0 0 0 ...
 $ PCYellow : int 0 0 0 0 0 0 0 0 0 ...
 $ PH.Red   : int 0 0 0 0 0 0 0 0 0 ...
 $ PW.Red   : int 0 0 0 0 0 0 0 0 0 ...
 $ PCRed   : int 0 0 0 0 0 0 0 0 0 ...
 $ Time.Stamp: int 2933417 2933417 2933432 2933432 2933463 2933463 2933463 2933495 2933495 2933495 ...
 $ X       : logi NA NA NA NA NA NA ...> |
```

# Data structures in R

Data.frames are data organized into rows and columns.

Data frames are organized by [row,column]. You must use the comma!

You can look at a vector of the data frame using \$, i.e. `df$TOF`

Or you can specify the row or column to output, i.e. `df[, 9]` outputs the ninth column and `df[, "TOF"]` outputs the column named “TOF”

```
Console ~ / 
> 
> df[,"TOF"]
 [1] 52 556 88 639 531 34 247 152 96 79 54 133 64 471 417 21 491 400 504 163 25 436 57 36
[25] 190 29 574 476 394 51 584 74 172 373 26 49 92 31 25 521 177 108 266 59 57 435 26 135
[49] 62 197 52 38 236 97 545 438 82 28 129 76 67 48 105 140 147 22 88 74 24 265 162 21
[73] 31 122 50 20 70 53 329 27 71 80 158 517 523 448 40 30 86 49 519 134 26 59 49 78
[97] 65 46 91 219 259 21 85 128 244 216 59 188 182 56 43 75 51 60 136 114 80 77 52 58
[121] 219 26 203 81 25 127 190 53 63 54 37 40 468 197 172 59 85 196 45 23 24 245 58 209
[145] 314 400 229 488 99 35 260 46 103 199 421 116 481 133 112 40 55 164 51 161 212 161 128 212
[169] 64 23 491 47 20 22 69 471 23 120 28 51 524 124 265 20 69 22 62 481 103 32 53 143
[193] 48 37 203 116 89 31 45 255 73 281 296 105 472 49 20 43 20 46 113 92 143 194 46 87
[217] 194 23 70 525 54 52 51 46 36 59 71 27 506 22 331 73 847 438 58 160 50 55 402 31
[241] 208 524 58 50 78 479 54 102 55 73 295 516 109 38 121 142 127 132 486 44 101 335 66 393
[265] 94 205 77 172 205 565 29 133 54 443 511 34 97 139 182 25 39 27 182 50 60 101 56 59
[289] 146 44 206 48 118 324 36 310 127 537 194 59 46 111 411 61 439 74 32 51 91 115 46 250
[313] 467 229 47 58 181 25 103 35 184 42 295 185 406 119 32 61 549 635 84 186 111 179 49 142
```

# **Take some time to move around through your data**

1. Are there missing values?
2. What are the column names?
3. What are the dimensions?
4. What are the row values?
5. Are the data organized long or wide?

# Some thoughts about organizing raw data

- Be consistent.
- Write dates as YYYY-MM-DD or YYYYMMDD.
- Fill in all of the cells.
- Put just one thing in a cell.
- Make it a rectangle.
- Create a data dictionary.
- No calculations in the raw data files.
- Don't use font color or highlighting as data.
- Choose good names for things.
- Make backups.
- Use data validation to avoid data entry mistakes. Unit tests!
- Save the data in plain text files.

from Karl Broman (UW Madison), check out his fantastic blog  
Other great rules here: <https://github.com/jtleek/datasharing>

**The most important part of data analysis  
is to think about your data**

What do you want to test?

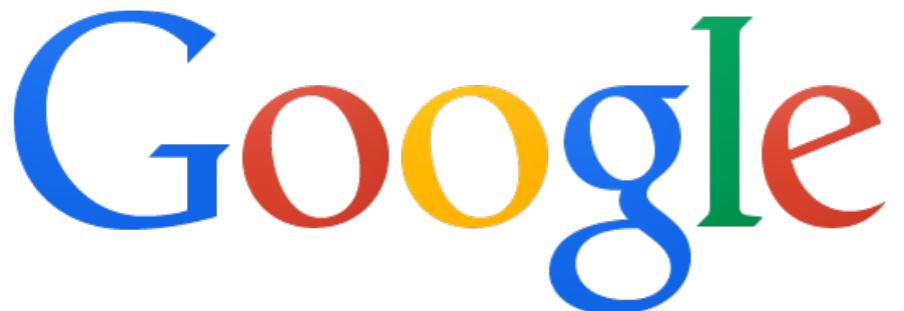
How will you show that conclusion?

Put the goal at the top of your markdown report

# **Data exploration time; define goal(s)**



# What to do when you are lost?



<http://www.r-bloggers.com>

<http://gettinggeneticsdone.blogspot.com>

<http://rseek.org>

# Day #2 Homework

1. Install the swirl package
2. Type `swirl()`
3. Take the introduction to R programming class
4. Play with your new data set

{swirl}

Learn R, in R.