

# Evomics 2026

## R & ggplot2



# Outline

- Short introduction
  - Why is R useful
  - RStudio
  - R Markdown
  - Data structures
- Dataset for practical
- Practical

# What is R?

A free software environment  
(and language) for statistical  
computing and graphics

<http://www.r-project.org>

```
R Console

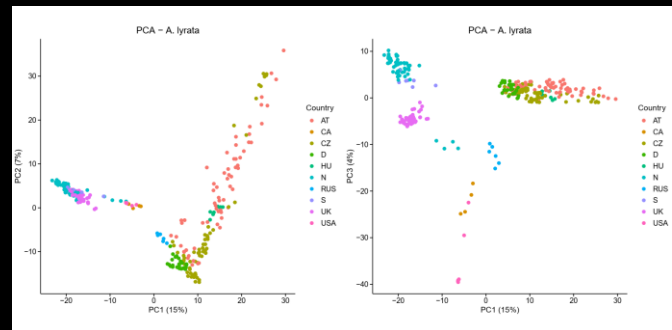
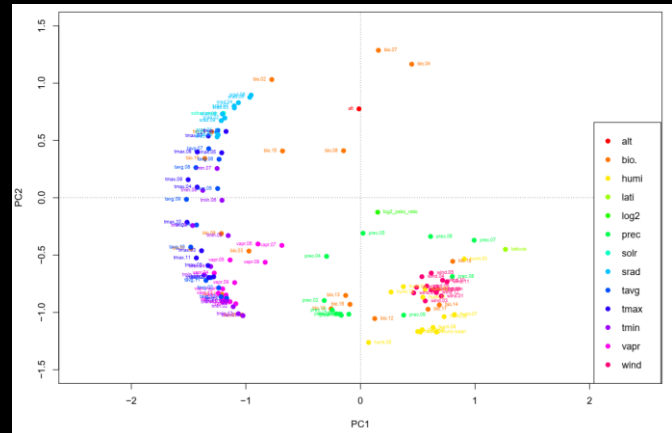
R version 4.3.3 (2024-02-29 ucrt) -- "Angel Food Cake"
Copyright (C) 2024 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

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.

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.

> |
```



# Why is R useful?

- **Open source**
- **Data management and manipulation**
  - Importing data in various formats (like text files, excel files, etc.)
  - Manipulating data (subsetting and filtering tables, merging, transposing, etc.)
- Cutting-edge **graphical data visualization**
- Support for rich **statistical simulation and modeling**
- Well established system of **packages and documentation**
- **Active development** and dedicated **community**

# When do we use R?

Big data (*fastq etc.*)



UNIX tools

Small data (*count tables etc.*)



R / Python

Final results, visualisations

# Why R and not Excel?



# Why R and not Excel?

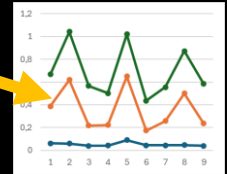


```
"genotype" "cell.width" "X"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 38
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084 NA
"control" 35.888557069 NA
"XA53" 39.0640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 28
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.053
"XA53" 13.9409304767847 74
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			



# Why R and not Excel?

## Scenario 1: Data changed

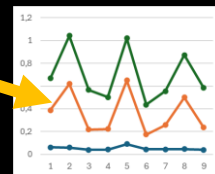


```
"genotype" "cell.width" "X"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 38
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084 NA
"control" 35.888557069 NA
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.053
"XA53" 13.9409304767847 74
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

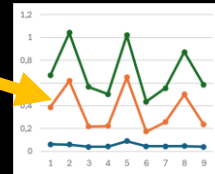


```
"genotype" "cell.width" "X"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 38
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084 NA
"control" 35.888557069 NA
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.053
"XA53" 13.9409304767847 74
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			





# Why R and not Excel?

Scenario 1: Data changed



# Why R and not Excel?

Scenario 1: Data changed

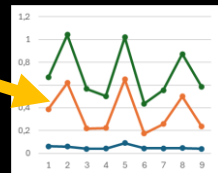


```
"genotype" "cell.width" "4"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

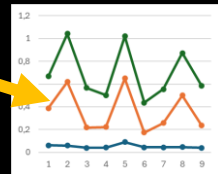


```
"genotype" "cell.width" "4"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

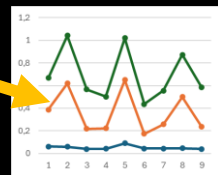


```
"genotype" "cell.width" "4"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

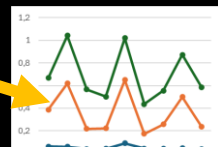


```
"genotype" "cell.width" "4"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```

	A	B	C
1			
2			
3			
4			
5			
6			

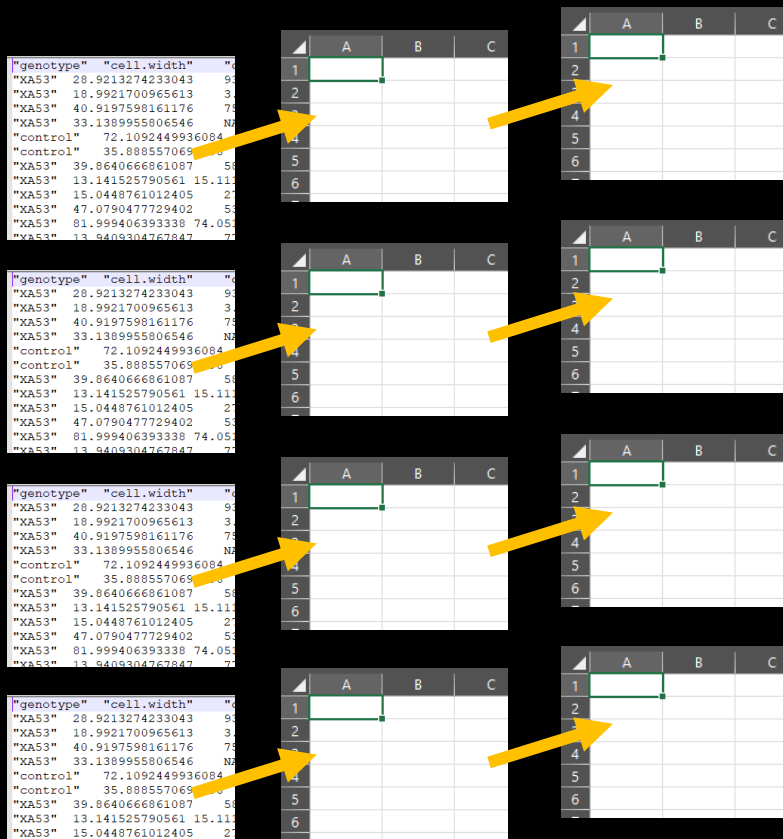
	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			



# Why R and not Excel?

## Scenario 2: Analysis changed



# Why R and not Excel?

## Scenario 2: Analysis changed

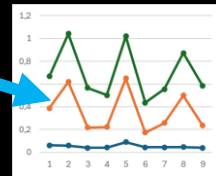


```
"genotype" "cell.width" "4"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 38
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

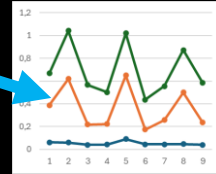


```
"genotype" "cell.width" "4"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 38
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

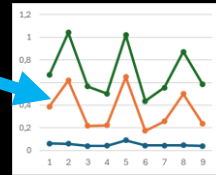


```
"genotype" "cell.width" "4"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 38
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			



```
"genotype" "cell.width" "4"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 38
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 27
"XA53" 47.0790477729402 53
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			



# Why R and not Excel?

Scenario 3: Many plots needed

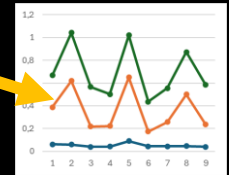


```
"genotype" "cell.width" "A"
"XA53" 28.9213274233043 98
"XA53" 18.9921700965613 38
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084 75
"control" 35.888557069 75
"XA53" 39.0640666861087 58
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 28
"XA53" 47.0790477729402 58
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 75
```

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

	A	B	C
1			
2			
3			
4			
5			
6			

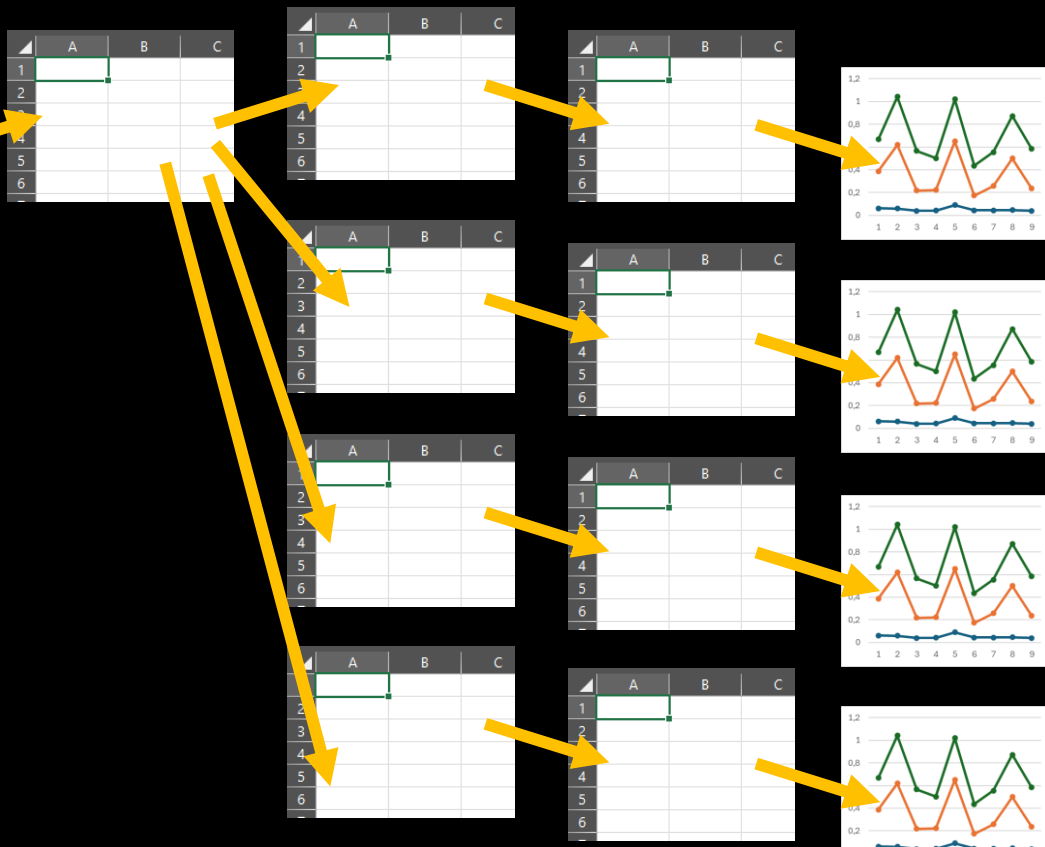


# Why R and not Excel?

Scenario 3: Many plots needed



```
"genotype" "cell.width" "X"
"XA53" 28.9213274233043 9
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 7
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888557069
"XA53" 39.0640666861087 5
"XA53" 13.141525790561 15.11
"XA53" 15.0448761012405 2
"XA53" 47.0790477729402 5
"XA53" 81.999406393338 74.05
"XA53" 13.9409304767847 7
```



# Why R and not Excel?

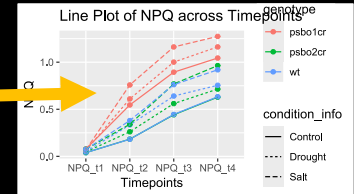


genotype	cell.width	
"XA53"	28.9213274233043	93
"XA53"	18.9921700965613	31
"XA53"	40.9197598161176	79
"XA53"	33.1389955806546	NA
"control"	72.1092449936094	
"control"	35.8885570000000	
"XA53"	39.8640666861087	58
"XA53"	13.141525790561	15.111
"XA53"	15.0448761012405	27
"XA53"	47.0790477729402	53
"XA53"	81.999406393338	74.052
"XA53"	13.8409304767847	7

```
##{r}
setwd("D:/!ecolgen/resources/orthofinder/
brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/Results_brassicaceae
_2/Comparative_Genomics_Statistics/Orthog
roups_speciesoverlaps.tsv")

## heatmap with values
pdf("R_analysis/Orthogroups_Speciesoverl
aps_heatmap.pdf", width=14, height=7,
onfile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
# input data
gdata <- spec.overlap
```



# Why R and not Excel?

## Scenario 1: Data changed



```
"genotype" "cell.width" "6"
"XA53" 28.9213274233043 93
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888550691898
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 27
"XA53" 47.079047729402 53
"XA53" 81.999406393338 74.052
"XA53" 13.9409304767847 7
```

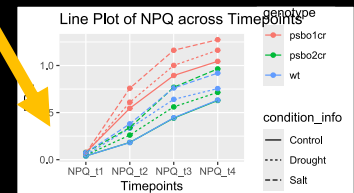
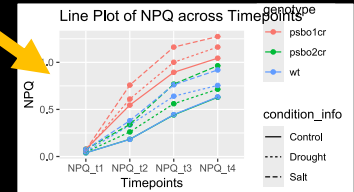
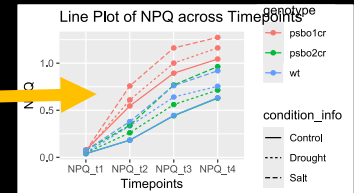
```
"genotype" "cell.width" "6"
"XA53" 28.9213274233043 93
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888550691898
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 27
"XA53" 47.079047729402 53
"XA53" 81.999406393338 74.052
"XA53" 13.9409304767847 7
```

```
"genotype" "cell.width" "6"
"XA53" 28.9213274233043 93
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888550691898
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 27
"XA53" 47.079047729402 53
"XA53" 81.999406393338 74.052
"XA53" 13.9409304767847 7
```

```
{r}
setwd("D:/!ecolgen/resources/orthofinder/brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/Results_brassicaceae_2/Comparative_Genomics_Statistics/Orthogroups_speciesoverlaps.tsv")

# heatmap with values
pheatmap("R_analysis/Orthogroups_Species_overlap_heatmap.pdf", width=14, height=7,
onefile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
input data
data <- spec.overlap
```





# Why R and not Excel?

## Scenario 2: Analysis changed



```
"genotype" "cell.width" "6"
"XA53" 28.9213274233043 93
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888550691898
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 27
"XA53" 47.079047729402 58
"XA53" 81.999406393338 74.052
"XA53" 13.9409304767847 7
```

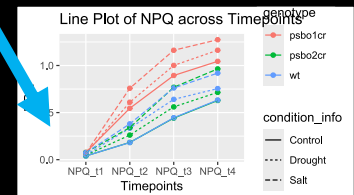
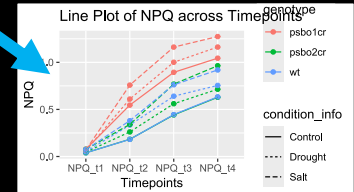
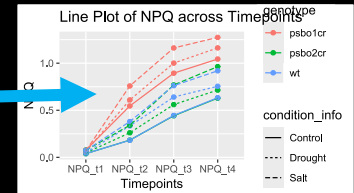
```
"genotype" "cell.width" "6"
"XA53" 28.9213274233043 93
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888550691898
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 27
"XA53" 47.079047729402 58
"XA53" 81.999406393338 74.052
"XA53" 13.9409304767847 7
```

```
"genotype" "cell.width" "6"
"XA53" 28.9213274233043 93
"XA53" 18.9921700965613 3
"XA53" 40.9197598161176 75
"XA53" 33.1389955806546 NA
"control" 72.1092449936084
"control" 35.888550691898
"XA53" 39.8640666861087 58
"XA53" 13.141525790561 15.111
"XA53" 15.0448761012405 27
"XA53" 47.079047729402 58
"XA53" 81.999406393338 74.052
"XA53" 13.9409304767847 7
```

```
##{r}
setwd("D:/!ecolgen/resources/orthofinder/brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/Results_brassicaceae_2/Comparative_Genomics_Statistics/Orthogroups_speciesoverlaps.tsv")

## Heatmap with values
pheatmap::pheatmap("R_analysis/Orthogroups_Species_overlap_heatmap.pdf", width=14, height=7,
onefile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
input data
data <- spec.overlap
```



# Why R and not Excel?

## Scenario 3: Many plots needed

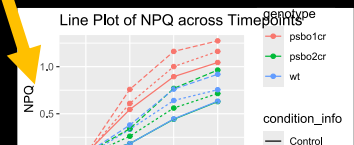
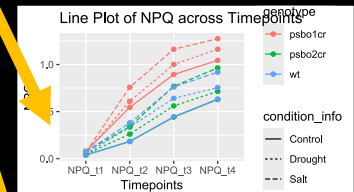
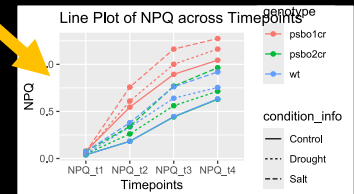
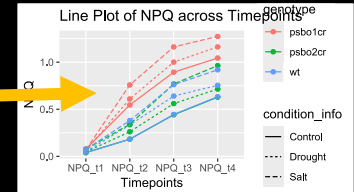


genotype	cell.width	
"XA53"	28.9213274233043	93
"XA53"	18.9921700965613	32
"XA53"	40.9197598161176	75
"XA53"	33.1389955806546	NA
"control"	72.1092449936094	
"control"	35.8885570000000	
"XA53"	39.8640666861087	58
"XA53"	13.141525790561	15.111
"XA53"	15.0448761012405	27
"XA53"	47.0790477729402	53
"XA53"	81.999406393338	74.052
"XA53"	13.8409304767847	7

```
{r}
setwd("D:/!ecolgen/resources/orthofinder/brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/brassicaceae_2/comparative_genomics_statistics/orthogroups/speciesoverlaps.tsv")

## heatmap with values
pdf("R_analysis/Orthogroups_Species_overlap_heatmap.pdf", width=14, height=7,
onefile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
# input data
gdata <- spec.overlap
```



# Why R and not Excel?



Scenario 4: Someone wants to understand or repeat the analysis

```
## {r}
setwd("b:/!ecolgen/resources/orthofinder/
brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/Results_brassicaceae
_2/Comparative_Genomics_Statistics/Orthog
roups_Speciesoverlaps.tsv")

## heatmap with values
pdf ("R_analysis/Orthogroups_Speciesoverl
aps_heatmap.pdf", width=14, height=7,
onefile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
# input data
gdata <- spec.overlap
```

# Why R and not Excel?



Scenario 4: Someone wants to understand or repeat the analysis

```
## {r}
setwd("b:/!ecolgen/resources/orthofinder/brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/Results_brassicaceae_2/Comparative_Genomics_Statistics/Orthogroups_Speciesoverlaps.tsv")

## heatmap with values
pdf ("R_analysis/Orthogroups_Speciesoverlaps_heatmap.pdf", width=14, height=7,
onefile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
# input data
gdata <- spec.overlap
```



Me in 2 months

# Why R and not Excel?



Scenario 4: Someone wants to understand or repeat the analysis

```
## {r}
setwd("b:/!ecolgen/resources/orthofinder/brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/Results_brassicaceae_2/Comparative_Genomics_Statistics/Orthogroups_SpeciesOverlaps.tsv")

## heatmap with values
pdf ("R_analysis/Orthogroups_SpeciesOverlaps_heatmap.pdf", width=14, height=7,
onefile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
# input data
gdata <- spec.overlap
```



Me in 2 months



Me in 2 years

# Why R and not Excel?



Scenario 4: Someone wants to understand or repeat the analysis

```
## {r}
setwd("b:/!ecolgen/resources/orthofinder/brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/Results_brassicaceae_2/Comparative_Genomics_Statistics/Orthogroups_Speciesoverlaps.tsv")

## heatmap with values
pdf ("R_analysis/Orthogroups_Speciesoverlaps_heatmap.pdf", width=14, height=7,
onefile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
# input data
gdata <- spec.overlap
```



Me in 2 months



Me in 2 years



Collaborator

# Why R and not Excel?



Scenario 4: Someone wants to understand or repeat the analysis

```
## {r}
setwd("b:/!ecolgen/resources/orthofinder/brassicaceae_2/")
old.par<-par(no.readonly = T)

spec.overlap <- read.table(file =
"orthofinder_results/Results_brassicaceae_2/Comparative_Genomics_Statistics/Orthogroups_Speciesoverlaps.tsv")

## heatmap with values
pdf ("R_analysis/Orthogroups_Speciesoverlaps_heatmap.pdf", width=14, height=7,
onefile = T)
par(mar = c(2, 12, 12, 2) + 0.1)
# input data
gdata <- spec.overlap
```



Me in 2 months



Me in 2 years



Collaborator



Paper reader

# Why R and not Excel?



Scenario 4: Someone wants to understand or repeat the analysis

```
## {r}  
setwd("D:/!ecolgen/resources/orthofinder/  
brassicaceae_2/")  
old.par<-par(no.readonly = T)  
spec.overlap <- read.table(file =
```

## Reproducibility

```
pdf ("R_analysis/orthogroups_speciesoverlap  
aps_heatmap.pdf", width=14, height=7,  
onefile = T)  
par(mar = c(2, 12, 12, 2) + 0.1)  
# input data  
gdata <- spec.overlap
```



Me in 2 months



Me in 2 years



Collaborator



Paper reader





BY JAMES MONTGOMERY FLAGG

**I WANT YOU  
FOR REPRODUCIBLE  
SCIENCE**

# Why R and not Excel?

- Automation
  - Many plots in one loop
  - Easily repeated if the data changes
- Reproducibility and transparency
  - You will know later what you did with the data
  - Other people will know what you did with the data
  - You can publish your code with your paper
- Excel tends to change some numbers to dates etc.

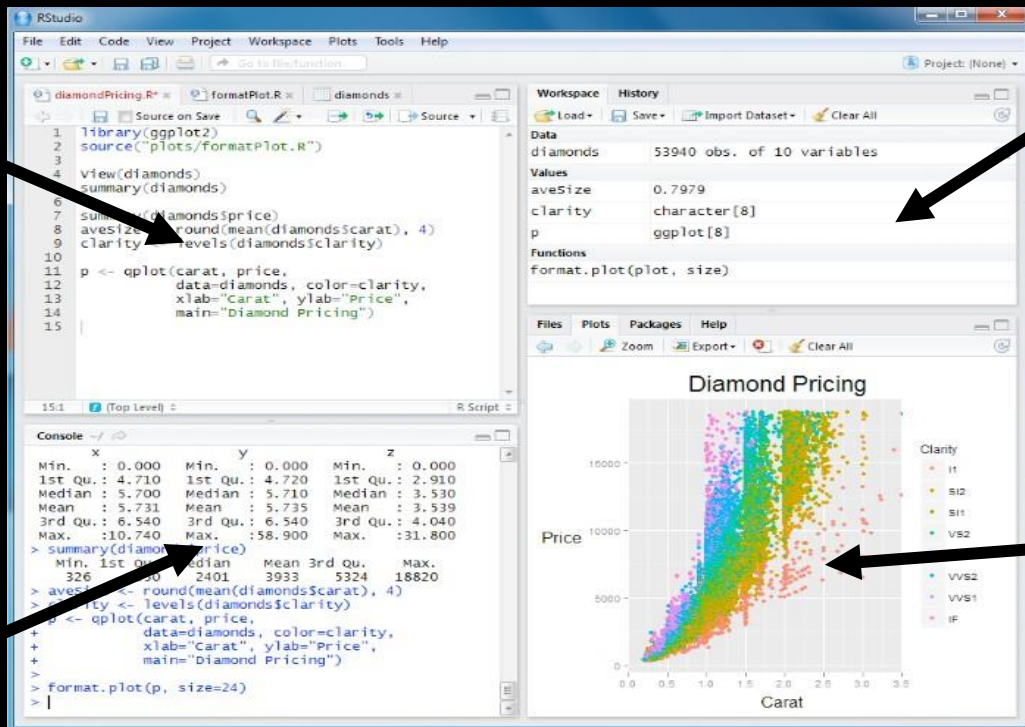
wt	10.233333	1007.22
psbo1cr	12.566666	71.56
psbo2cr	18.111111	516.33
wt	20.733333	1666.67
psbo1cr	23.166666	72.34

# R Studio

Integrated development environment (IDE) for R

Script (enter  
commands here)

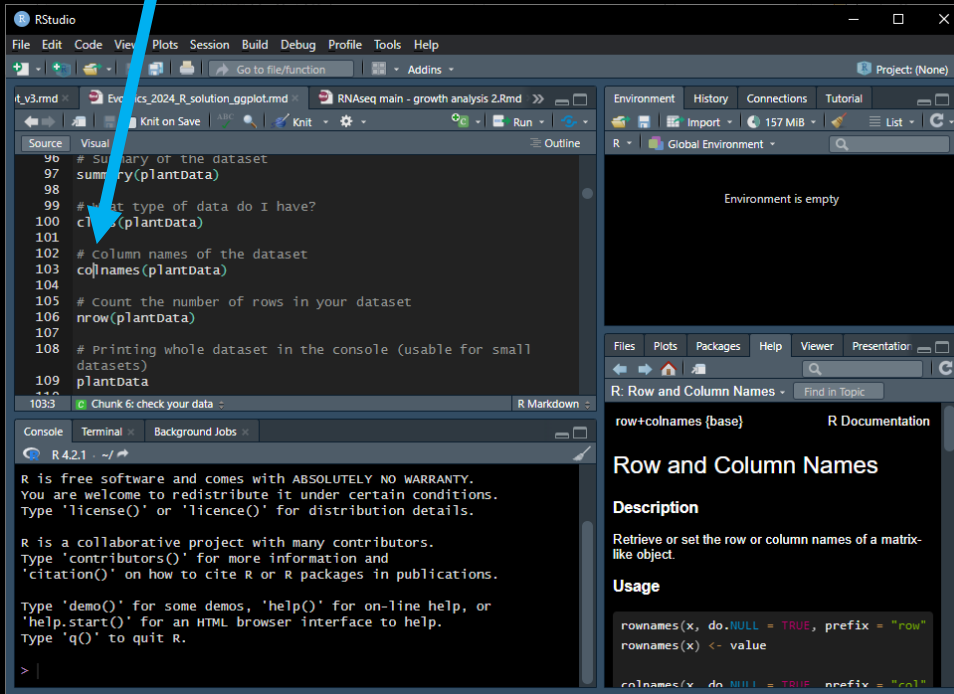
Workspace  
data



Help and Plots  
viewer

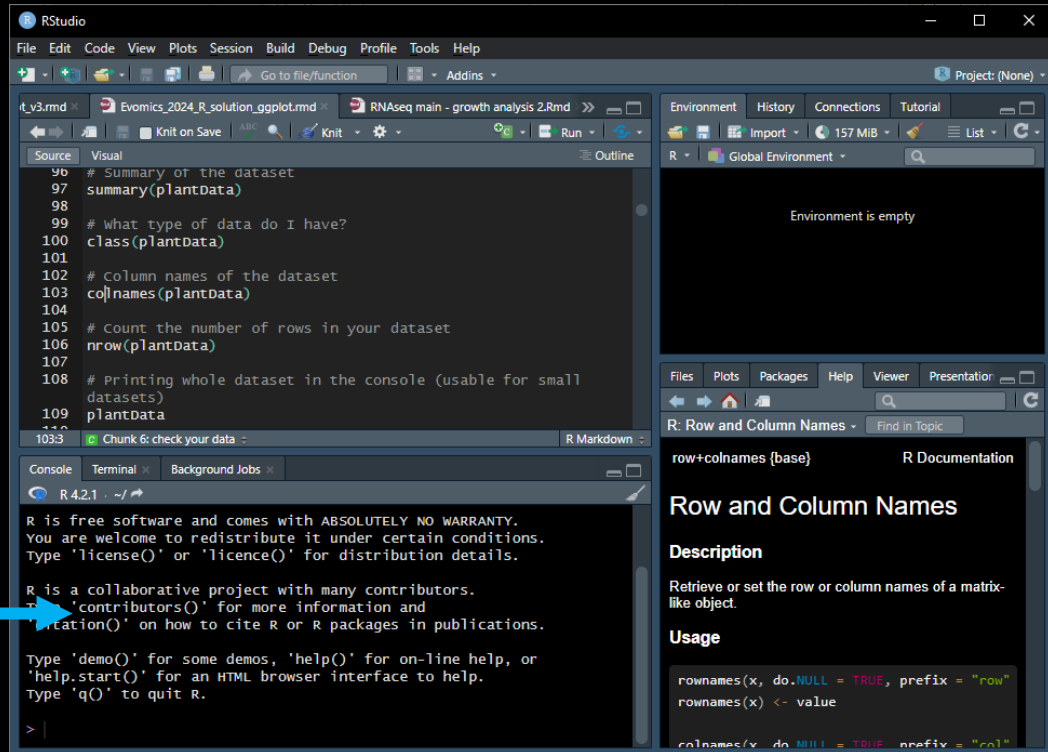
# Help in R Studio

Press **F1** when the cursor is in the name of the function



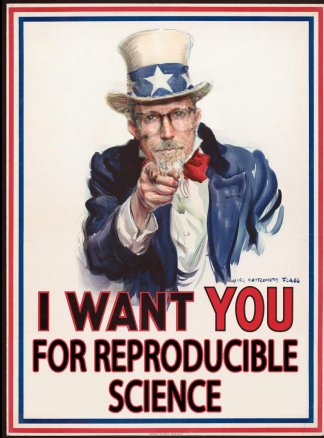
# Where to write the code?

Console? Not good for reproducibility.



# Where to write the code?

Console? Not good for reproducibility.



Console



```
# Summary of the dataset
summary(plantData)

# what type of data do I have?
class(plantData)

# Column names of the dataset
colnames(plantData)

# Count the number of rows in your dataset
nrow(plantData)

# Printing whole dataset in the console (usable for small
# datasets)
print(plantData)
```

1033 Chunk 6: check your data : R Markdown :

Console Terminal Background Jobs

R 4.2.1 ~ /

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.

R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' for 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.

> |

Environment History Connections Tutorial

R Global Environment

Environment is empty

Files Plots Packages Help Viewer Presenter

R: Row and Column Names - Find in Topic

row+colnames [base] R Documentation

### Row and Column Names

#### Description

Retrieve or set the row or column names of a matrix-like object.

#### Usage

```
rownames(x, do.NULL = TRUE, prefix = "row")
rownames(x) <- value

colnames(x, do.NULL = TRUE, prefix = "col")
```

# Where to write the code?

R Script / R Markdown

Source



The screenshot displays the RStudio IDE interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. Below the menu is a toolbar with icons for file operations and execution. The main editor area is divided into two panes: the left pane is the 'Source' editor, and the right pane is the 'Environment' pane. The 'Source' editor shows a script with R code comments and functions like `summary()`, `class()`, `colnames()`, `nrow()`, and `print()`. The 'Environment' pane shows the current environment, which is empty. Below the 'Source' editor is the 'Console' pane, which displays the R startup message and the prompt `> |`. To the right of the 'Environment' pane is the 'Documentation' pane, which shows the 'Row and Column Names' section of the R documentation, including a description and usage examples.

```
# Summary of the dataset
summary(plantData)

# what type of data do I have?
class(plantData)

# Column names of the dataset
colnames(plantData)

# Count the number of rows in your dataset
nrow(plantData)

# Printing whole dataset in the console (usable for small
# datasets)
print(plantData)
```

Environment

History

Connections

Tutorial

R

Global Environment

Environment is empty

Files

Plots

Packages

Help

Viewer

Presentation

R: Row and Column Names

Find in Topic

row+colnames [base]

R Documentation

## Row and Column Names

### Description

Retrieve or set the row or column names of a matrix-like object.

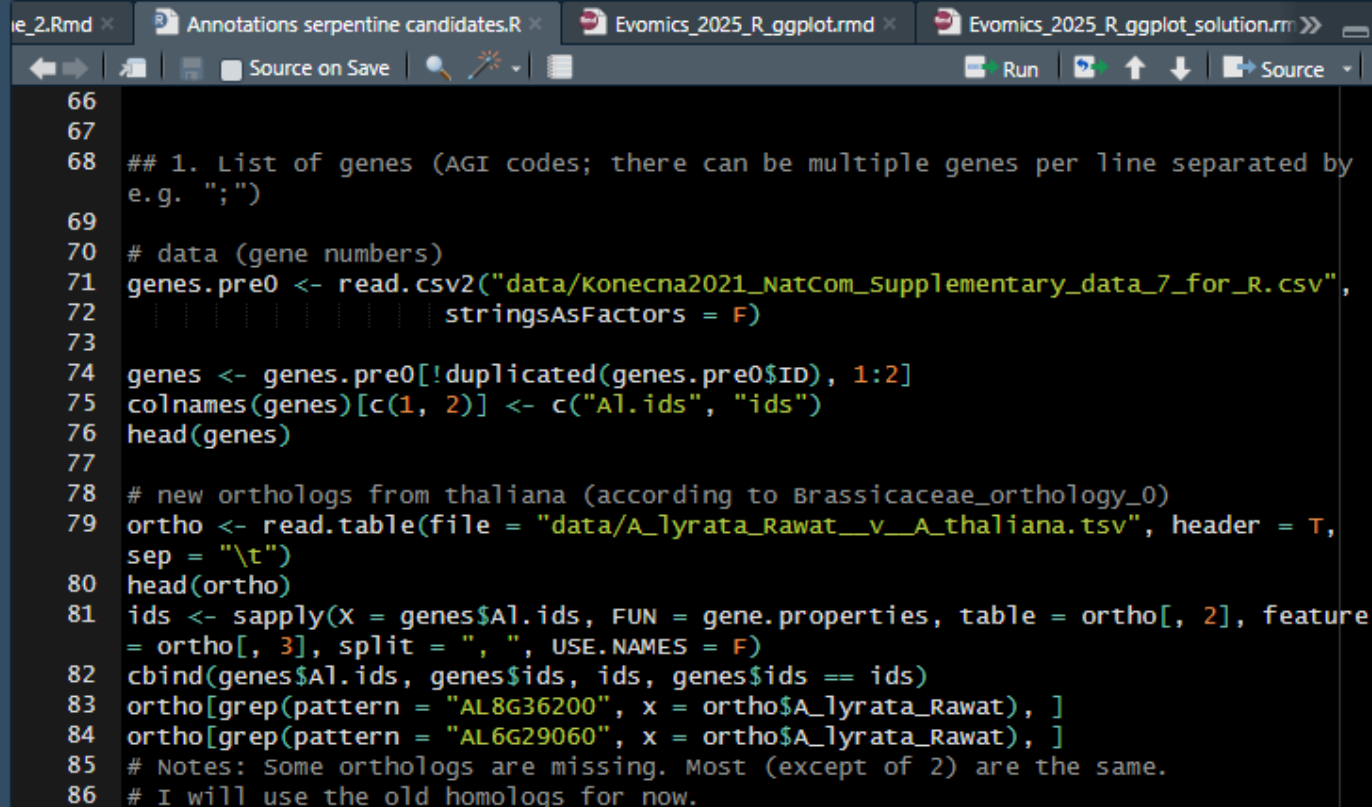
### Usage

```
rownames(x, do.NULL = TRUE, prefix = "row")
rownames(x) <- value

colnames(x, do.NULL = TRUE, prefix = "col")
```

# Where to write the code?

## R Script: Code + # Comments



The screenshot shows an RStudio interface with four tabs: 'ie\_2.Rmd', 'Annotations serpentine candidates.R', 'Evomics\_2025\_R\_ggplot.rmd', and 'Evomics\_2025\_R\_ggplot\_solution.rmd'. The active tab is 'Evomics\_2025\_R\_ggplot\_solution.rmd'. The script contains the following R code:

```
66
67
68 ## 1. List of genes (AGI codes; there can be multiple genes per line separated by
   e.g. ";")
69
70 # data (gene numbers)
71 genes.pre0 <- read.csv2("data/Konecna2021_NatCom_Supplementary_data_7_for_R.csv",
72                        stringsAsFactors = F)
73
74 genes <- genes.pre0[!duplicated(genes.pre0$ID), 1:2]
75 colnames(genes)[c(1, 2)] <- c("Al.ids", "ids")
76 head(genes)
77
78 # new orthologs from thaliana (according to Brassicaceae_orthology_0)
79 ortho <- read.table(file = "data/A_lyrata_Rawat_v_A_thaliana.tsv", header = T,
80                    sep = "\t")
81 head(ortho)
82 ids <- sapply(x = genes$Al.ids, FUN = gene.properties, table = ortho[, 2], feature
83             = ortho[, 3], split = "", ", ", USE.NAMES = F)
84 cbind(genes$Al.ids, genes$ids, ids, genes$ids == ids)
85 ortho[grep(pattern = "AL8G36200", x = ortho$A_lyrata_Rawat), ]
86 ortho[grep(pattern = "AL6G29060", x = ortho$A_lyrata_Rawat), ]
87 # Notes: Some orthologs are missing. Most (except of 2) are the same.
88 # I will use the old homologs for now.
```



# Where to write the code?

R Markdown: Formatted text + ````Code chunks````

```
ie_2.Rmd × Evomics_2025_R_ggplot_solution.rmd × Annotations serpentine candidates.R × Evomics_2025_R_ggplot.rmd ×
← → | | | Knit on Save | ABC | | Knit | | Run |
Source Visual Outline

259
260 After that, we will remove plants that have died during the experiment. These
plants have *NA* values in columns `size_mm2` and `QY_max`. The NPQ values were
not measured for all plants, and we are not going to plot those columns yet, so it
is fine for now that there are some *NA* values there.

261
262 ```{r clean your data}
263 # Remove all NA values from size_mm2 & QY_max column, as this indicates that
plants died during the experiment and we do not have data for them
264 pd_clean <- plantData %>%
265   filter(!is.na(size_mm2), !is.na(QY_max))
266 ```

267
268
269 **Exercise 3:** check how many lines we have removed.

270
271 Hint: use the `nrow()` function or `dim()` function on the original `plantData`
data.frame and on the new `pd_clean` data.frame.

272
273 ```{r solution dimensions, class.source= 'fold-hide', eval=FALSE}
274 # Dimensions of the original data.frame (number of rows and number of columns)
275 dim(plantData)
276 # Dimensions of the cleaned data.frame
```

# R Markdown

Can be  
“knitted” to  
produce report  
in html, pdf,  
docx, GitHub  
md etc.

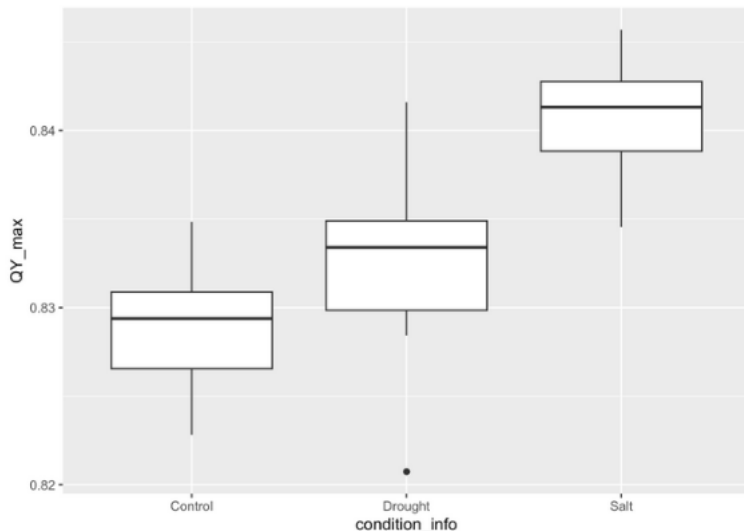
## 4.3 Modify your graph aesthetics

We will now make our box plot a bit fancier. Although the defaults often work well, you can modify almost everything within the `ggplot2` package.

Here you can see how to modify various things in the plot.

Hide

```
# Original box plot of QY_max by condition_info
p1 <- ggplot(pD_clean_wt, aes(x=condition_info, y=QY_max)) +
  geom_boxplot()
p1
```



Hide

*# Now let's get fancy with this plot. We'll start with our p1 plot and sequentially add layers to it.*

```
p1_fancy <- ggplot(pD_clean_wt, aes(x=condition_info, y=QY_max)) +
  geom_boxplot() + # add a boxplot layer (same as before)
  geom_point() + # add points to the boxplot
```

CC BY SA Posit Software, PBC • [info@posit.co](mailto:info@posit.co) • [posit.co](https://posit.co) • Learn more at [rmarkdown.rstudio.com](https://rmarkdown.rstudio.com) • HTML cheatsheets at [posit.co/cheatsheets](https://posit.co/cheatsheets) • rmarkdown 2.23 • Updated: 2023-07-10

# General data structures

- **Vector** - ordered collection of data

```
vector_1 <- c(2, 3, 4, 10)
```

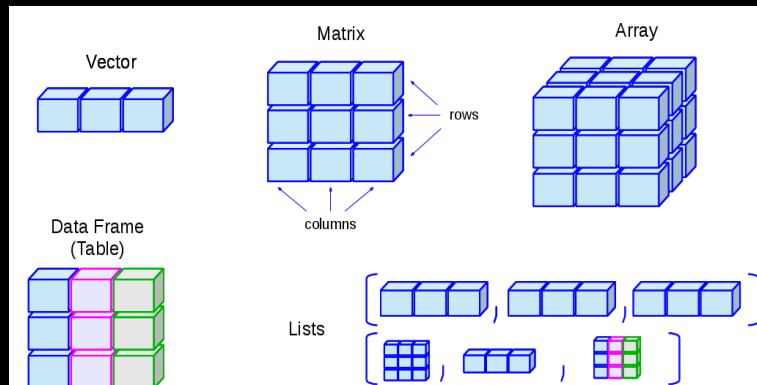
```
vector_2 <- c("potato", "lemonade", "avocado")
```

- **Matrix** - 2D collection of vectors with same data type

- **Array** - multiple dimension collection of vectors

- **Dataframe** - matrix-like with multiple data types (like an excel table with text and numbers)

- **Lists** - ordered collection of any objects (can contain also other lists inside it)



**But...**

**which dataset should we use to try R?**

# *Arabidopsis thaliana* mutants *psbo1* and *psbo2*

WT



*psbo1*

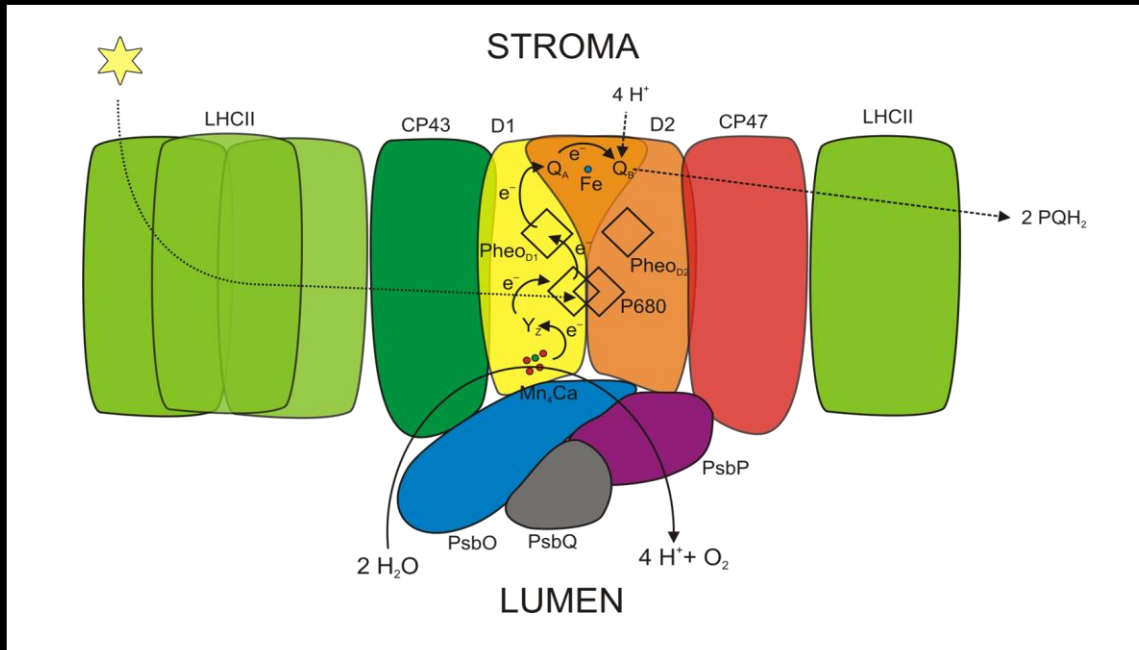


*psbo2*



# PsbO protein

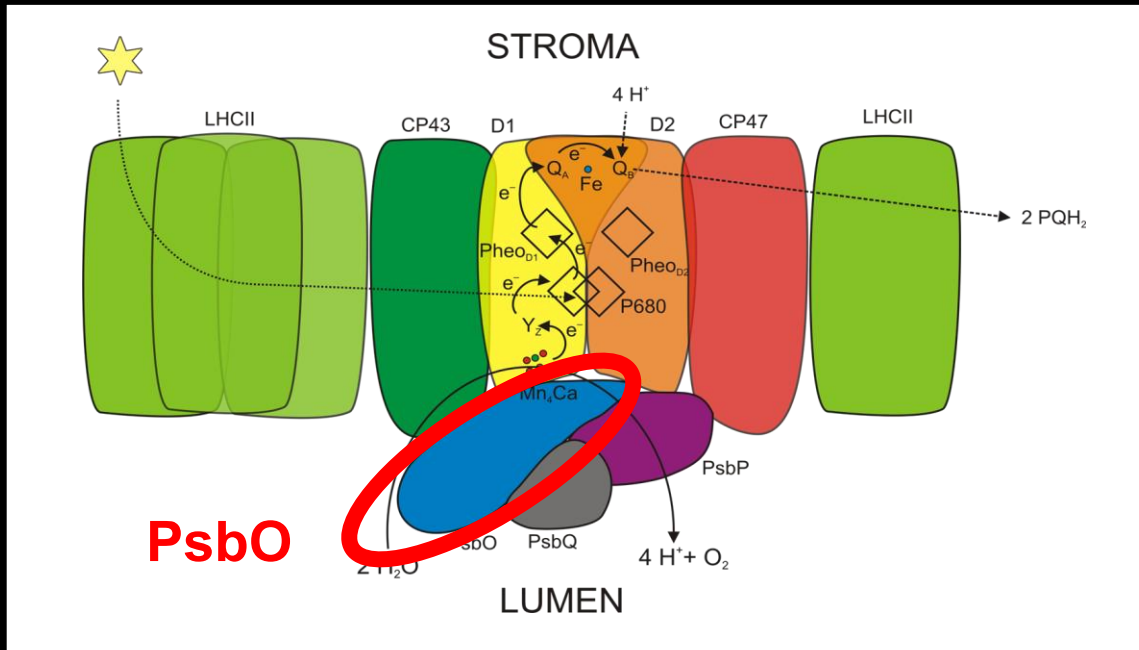
- Subunit of photosystem II
- Important for water splitting
- *Arabidopsis*: PsbO1 and PsbO2



Photosystem II

# PsbO protein

- Subunit of photosystem II
- Important for water splitting
- *Arabidopsis*: PsbO1 and PsbO2

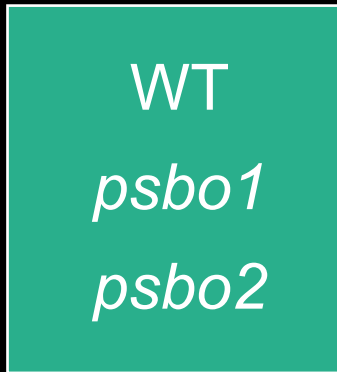


Photosystem II

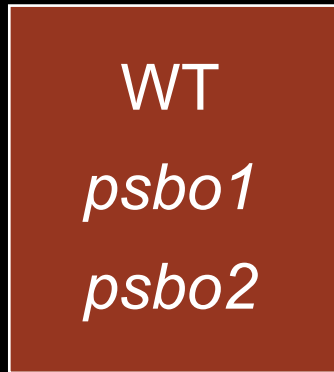


# Experimental design

Control

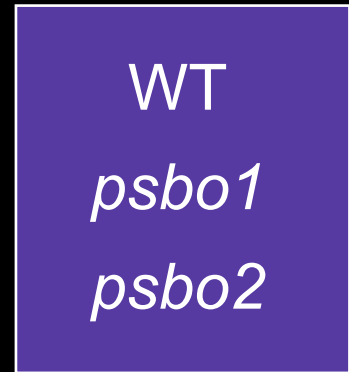


Drought



- water

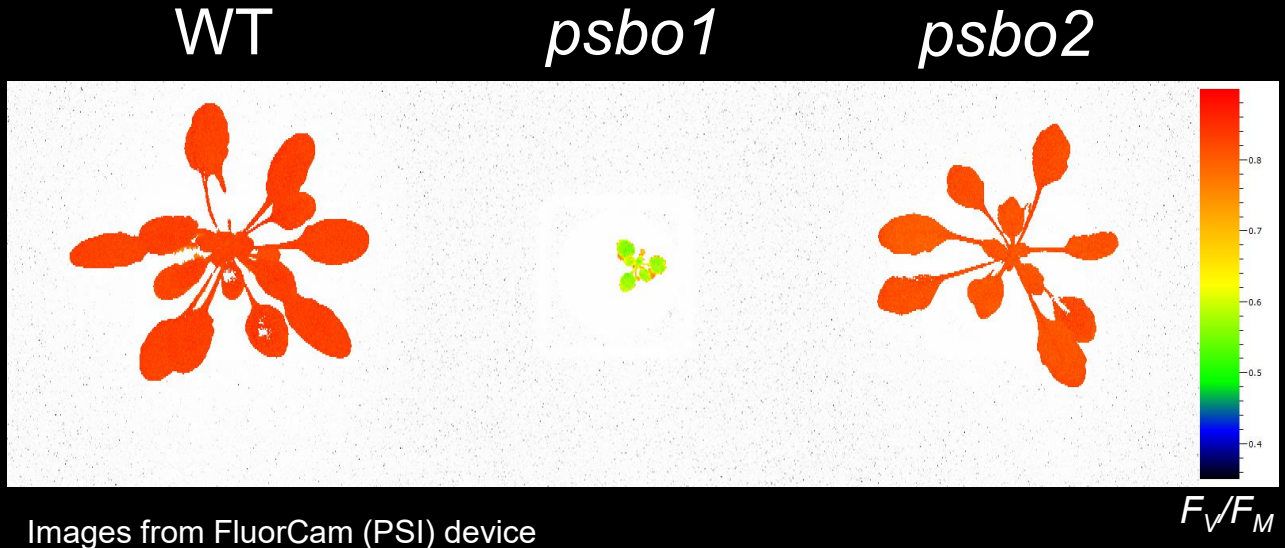
Salt stress



+ NaCl

# Measurement – chlorophyll fluorescence

- Leaf rosette area
- $F_v/F_M$  (QY\_max) – maximum quantum yield of photosystem II



# Let's start the practical!

Open the Rstudio server by typing in browser:  
<your IP>:8787



Remember:

- Practise makes the masters.
- Do sanity checks. Always.
- Use AI, but try to understand, check and improve the code.