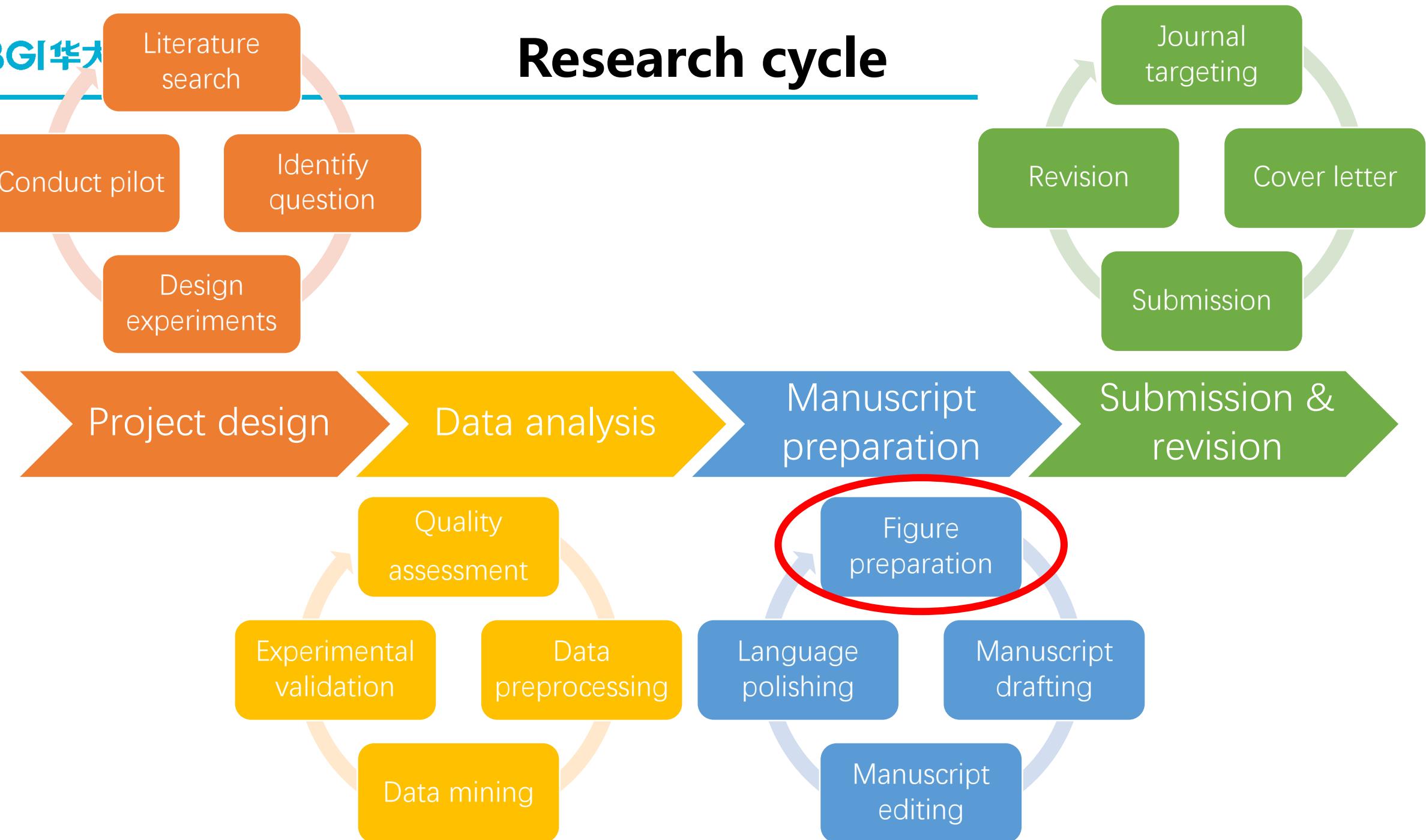


# Scientific figure

魏桐

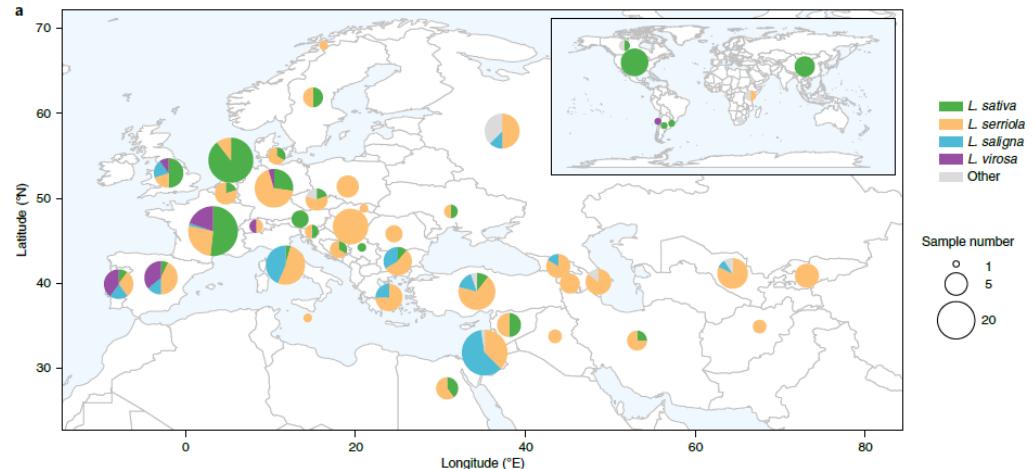
3/22/2023

# Research cycle



# What is scientific figure

Provide **an efficient way to communicate** your research



VS

**Plant materials and sequencing.** The collection of *Lactuca* SSD lines (<http://www.wur.eu/cgnsc002>) used in this study comprises a core subset of the regular collection of the Centre for Genetic Resources, the Netherlands (CGN) and includes all crop types of cultivated lettuce and main wild relatives used in plant breeding<sup>32</sup>. The total study set of 445 SSD lines included 131 cultivated lettuce (*L. sativa*) accessions collected worldwide, 201 *L. serriola* accessions, 57 *L. saligna* accessions, 37 *L. virosa* accessions and 19 lines from another eight *Lactuca* species (Supplementary Table 1).

# Part I

# Plotting by R

- R in brief
  - A programming language for **statistical computing and graphics**
  - **Open-source** implementation of the S language developed by Ross Ihaka and Robert Gentleman in 1990s
  - **A variety of packages** from an active community
- Usage
  - Data processing
  - Statistical analysis
  - Scientific figures

# Statistical analysis

```
> x <- rnorm(10)
> y <- rnorm(10, .5)
> t.test(x, y)

Welch Two Sample t-test
```

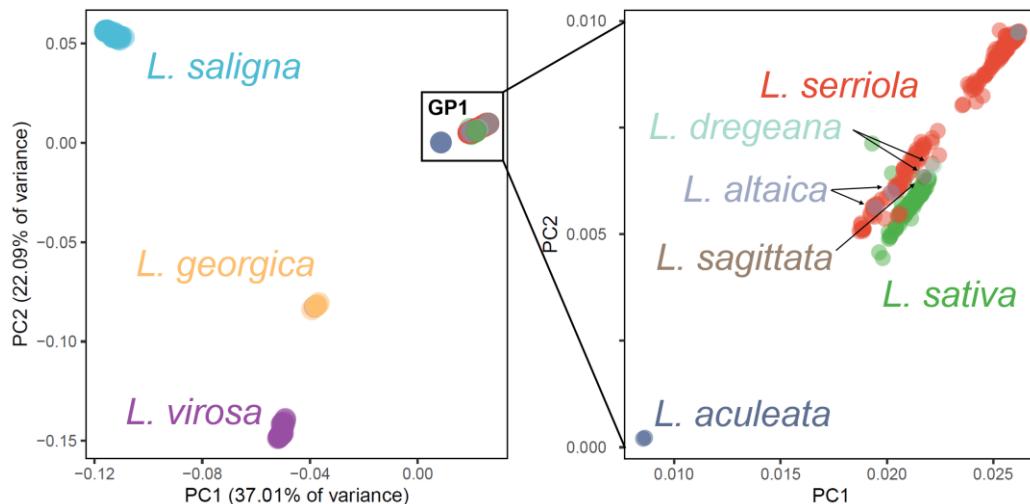
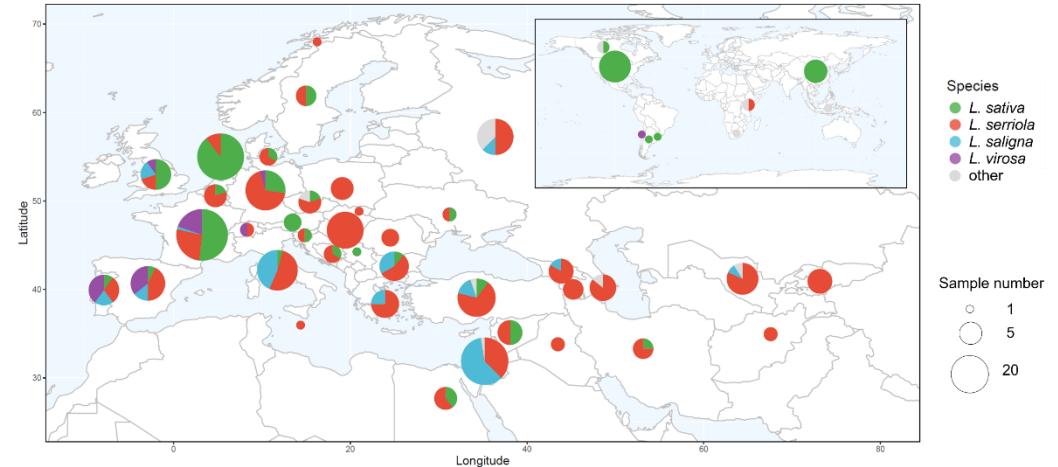
```
data: x and y
t = -0.72454, df = 15.418, p-value = 0.4796
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-1.1112264 0.5464145
sample estimates:
mean of x mean of y
0.02802619 0.31043212
```

```
> lm(y ~ x)

Call:
lm(formula = y ~ x)

Coefficients:
(Intercept)          x
0.3013        0.3253
```

# Scientific figure





[CRAN](#)  
[Mirrors](#)  
[What's new?](#)  
[Task Views](#)  
[Search](#)

[About R](#)  
[R Homepage](#)  
[The R Journal](#)

[Software](#)  
[R Sources](#)  
[R Binaries](#)  
[Packages](#)  
[Other](#)

[Documentation](#)  
[Manuals](#)  
[FAQs](#)  
[Contributed](#)

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 (2018-07-02, Feather Spray) [R-3.5.1.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](#)

Questions About R

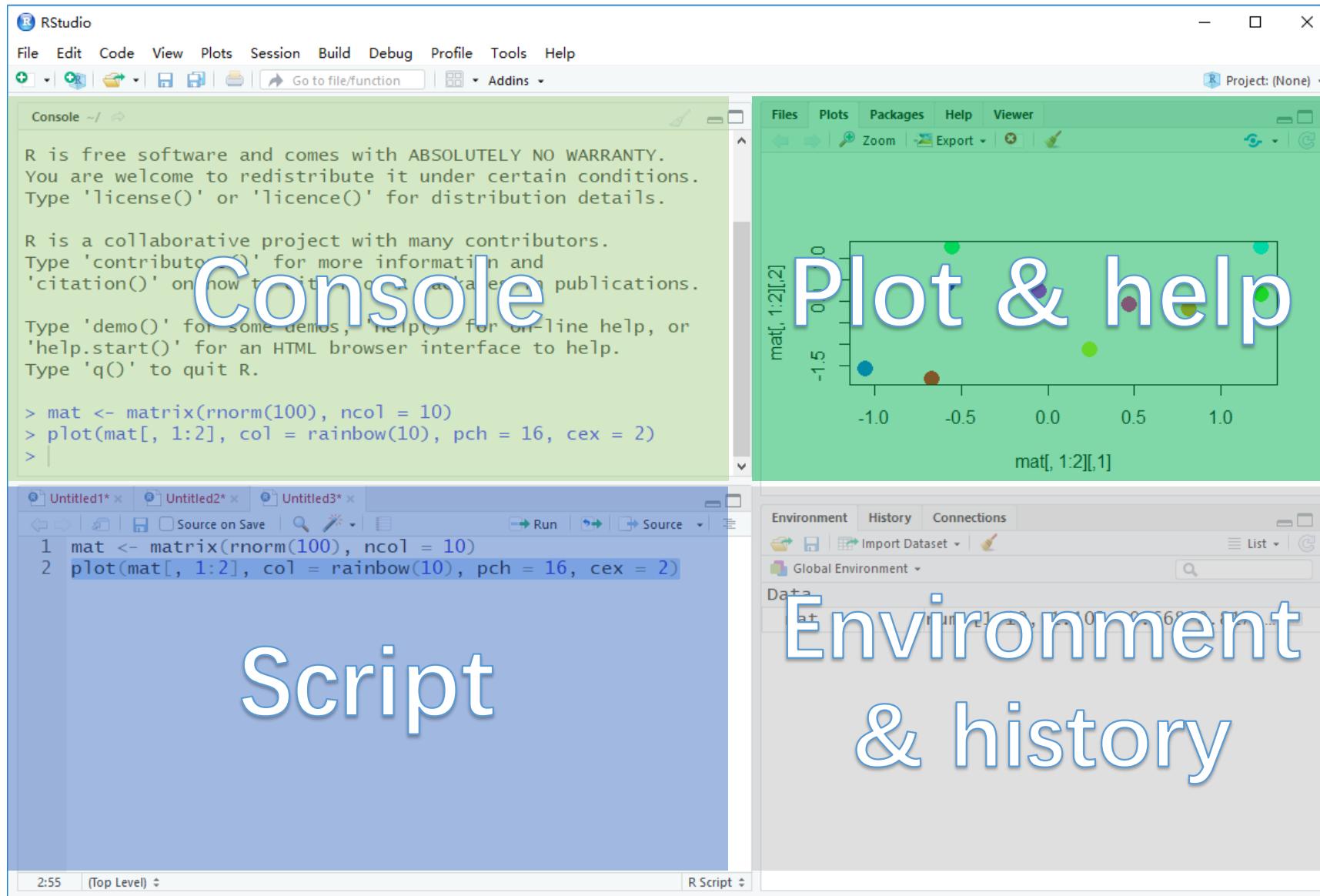
- If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#) before you send an email.

# R Studio download

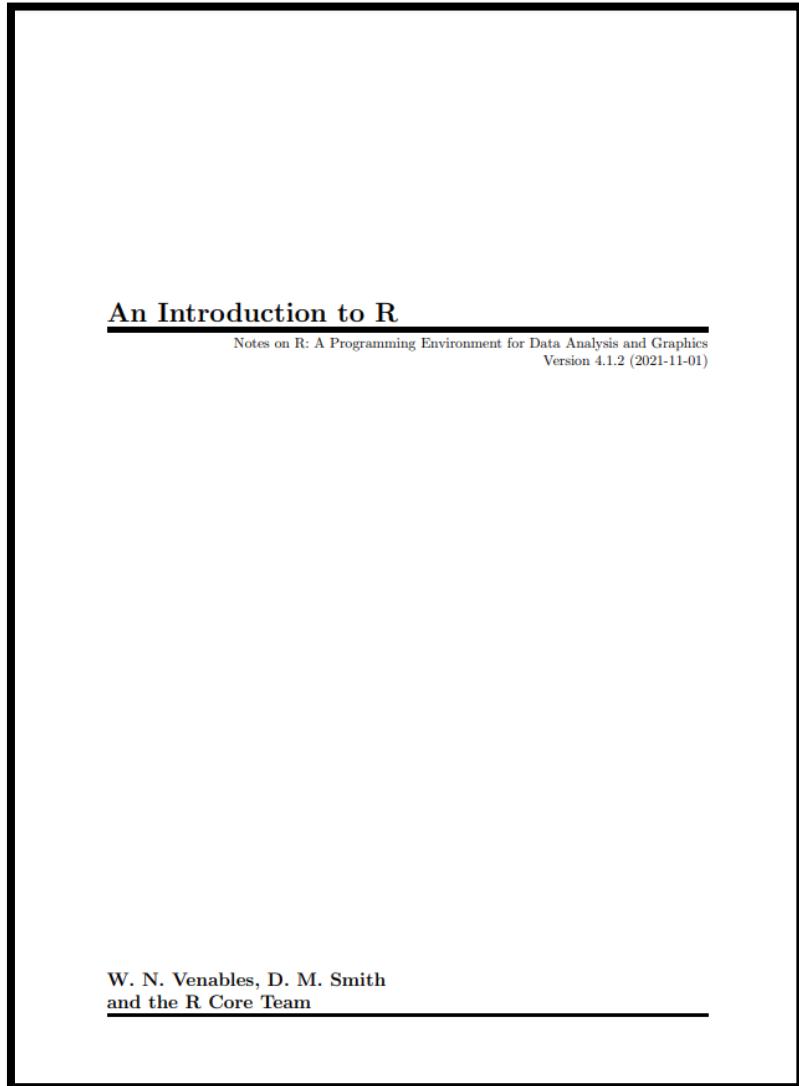
The screenshot shows the RStudio website's header. On the left is the RStudio logo. To the right are navigation links: Products, Solutions, Customers, Resources, About, and Pricing. The "DOWNLOAD" link is highlighted with a red oval. A magnifying glass icon is next to the "Pricing" link. Below the header, there's a large text block about RStudio Connect, followed by a blue "LEARN MORE" button.

**RStudio Connect** is the best way to host  
and share data products made with  
R and Python.

**LEARN MORE**



# Learning material: R intro



1. Introduction and preliminaries
2. Simple manipulations; numbers and vectors
3. Objects, their modes and attributes
4. Ordered and unordered factors
5. Arrays and matrices
6. Lists and data frames
7. Reading data from files
8. Probability distributions
9. Grouping, loops and conditional execution
10. Writing your own functions
11. Statistical models in R
12. Graphical procedures
13. Packages
14. OS facilities

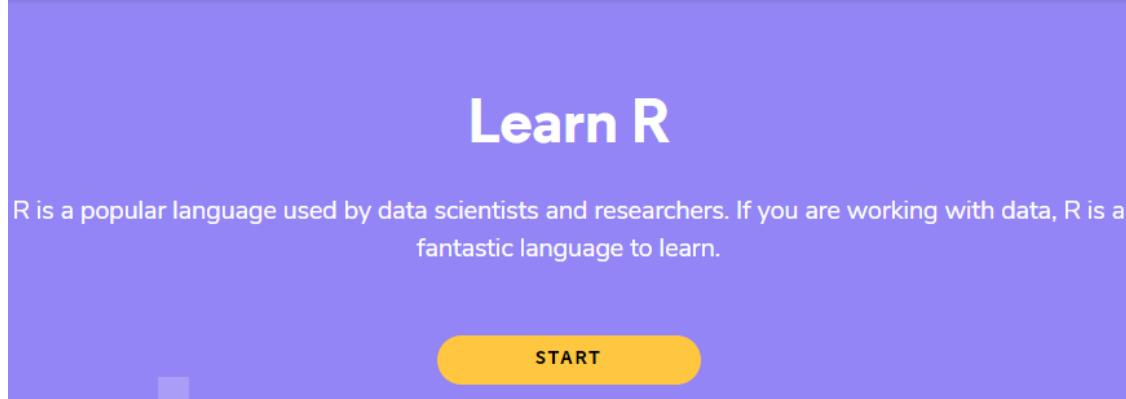
# Other online materials

## Coursera



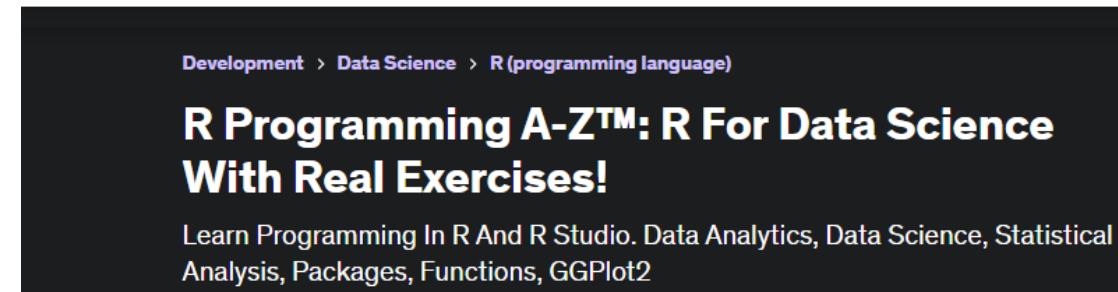
The image shows the Coursera Data Science Specialization landing page. It features a large banner with the title "Data Science Specialization" and a subtitle "Launch Your Career in Data Science. A nine-course introduction to data science, developed and taught by leading professors." Below the banner is a green button labeled "Try for Free". To the left of the banner is a sidebar with links: "About this Specialization", "Courses", "Pricing", "Creators", and "FAQ".

## codecademy



The image shows the codecademy Learn R landing page. The main heading is "Learn R". Below it is a paragraph: "R is a popular language used by data scientists and researchers. If you are working with data, R is a fantastic language to learn." At the bottom is a yellow "START" button.

## Udemy



The image shows the Udemy course page for "R Programming A-Z™: R For Data Science With Real Exercises!". The course title is prominently displayed in red. Below it is a brief description: "Learn Programming In R And R Studio. Data Analytics, Data Science, Statistical Analysis, Packages, Functions, GGPlot2". The course is located under the categories "Development > Data Science > R (programming language)".

## edX



The image shows the edX R Programming courses landing page. It features a large heading "R Programming courses" and a subtext: "Real college courses from Harvard, MIT, and more of the world's leading universities". Above the text is the edX logo.

# R codes in the script panel

Comment lines



```
# Step 0, read data and metadata
# 0.1, read the count data into R environment
# 0.2, generate a design matrix

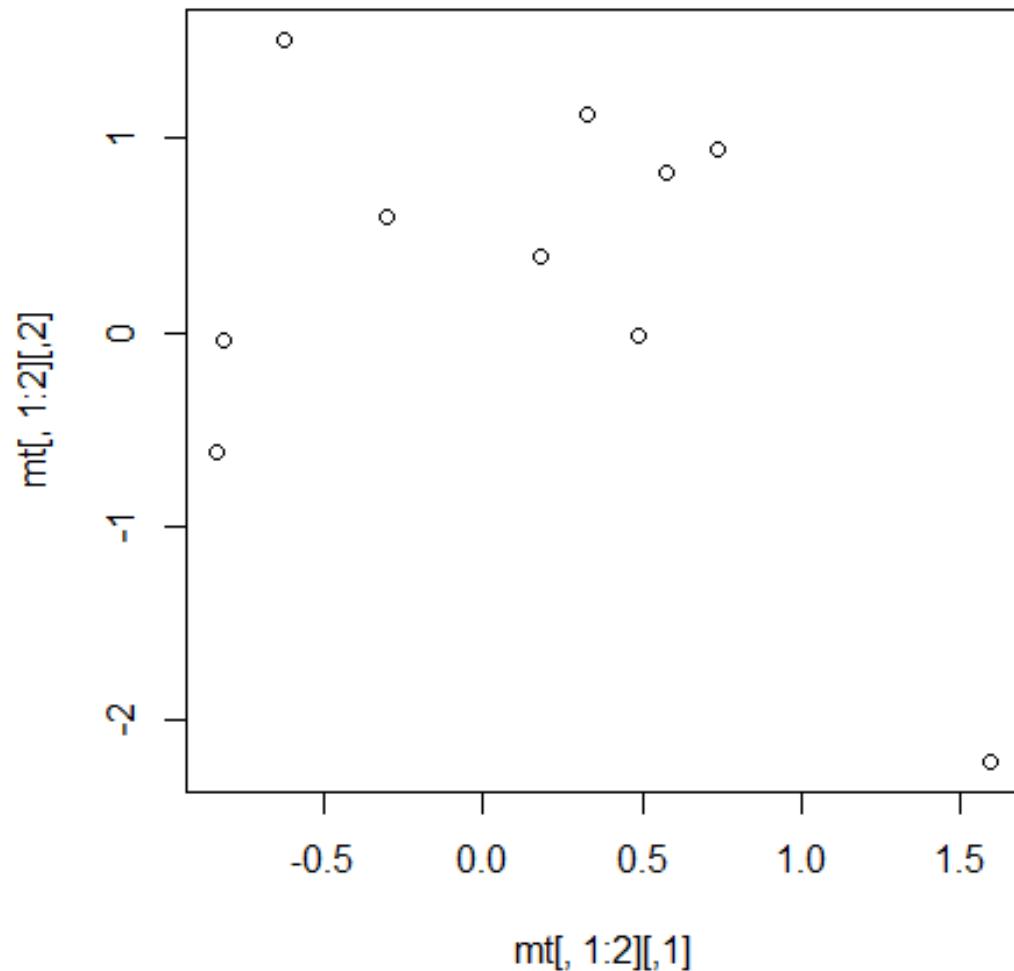
dat <- read.table("tophat_counts.txt", header = T, sep = "\t")
rownames(dat) <- dat$Gene_ID
dat <- dat[, -1]
dat.log2 <- log2(dat + 1)
# Labeling
sample.ID <- colnames(dat)
group.ID <- sub("_\\d", "", sample.ID)
sample.mat <- sapply(sample.ID, strsplit, split = "_|\\.\\.", simplify = T)
sample.mat <- t(as.data.frame(sample.mat))
colnames(sample.mat) <- c("genotype", "treatment", "time", "replicate")
```

object <- function(arguments)



# Function: plot()

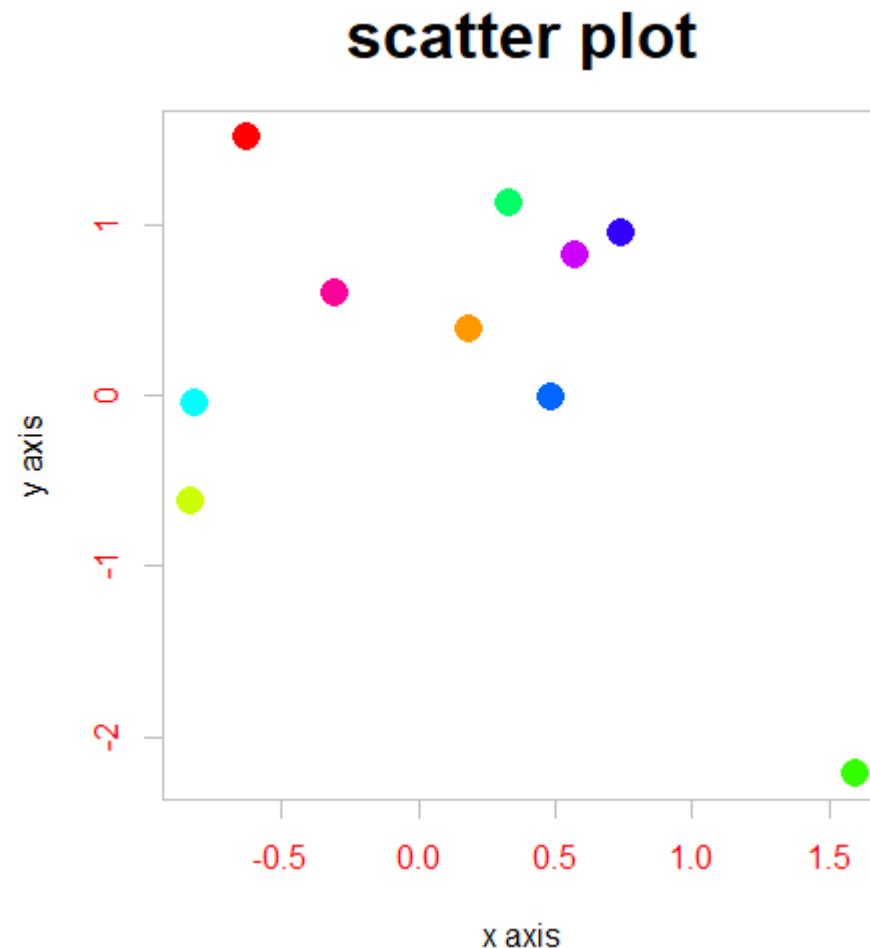
```
> set.seed(1)
> mt <- matrix(rnorm(20), ncol = 2)
> mt
      [,1]      [,2]
[1,] -0.6264538  1.51178117
[2,]  0.1836433  0.38984324
[3,] -0.8356286 -0.62124058
[4,]  1.5952808 -2.21469989
[5,]  0.3295078  1.12493092
[6,] -0.8204684 -0.04493361
[7,]  0.4874291 -0.01619026
[8,]  0.7383247  0.94383621
[9,]  0.5757814  0.82122120
[10,] -0.3053884  0.59390132
> plot(mt[, 1:2])
```



?graphics::plot for more information

# Function: plot()

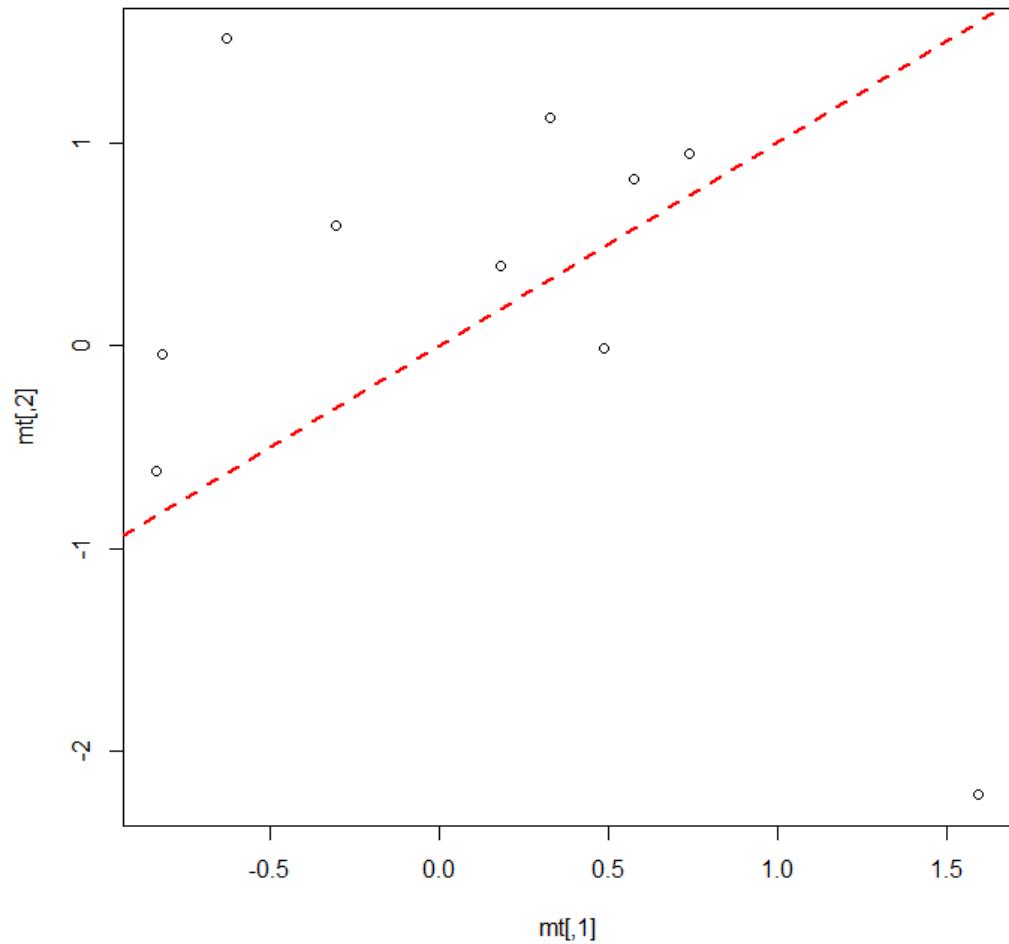
```
> plot(mt[, 1:2],  
+       pch = 16, cex = 2,  
+       col = rainbow(10),  
+       xlab = "x axis",  
+       ylab = "y axis",  
+       main = "scatter plot",  
+       cex.main = 2,  
+       fg = "grey",  
+       col.axis = "red")
```



?graphics::plot for more information

# Function: abline()

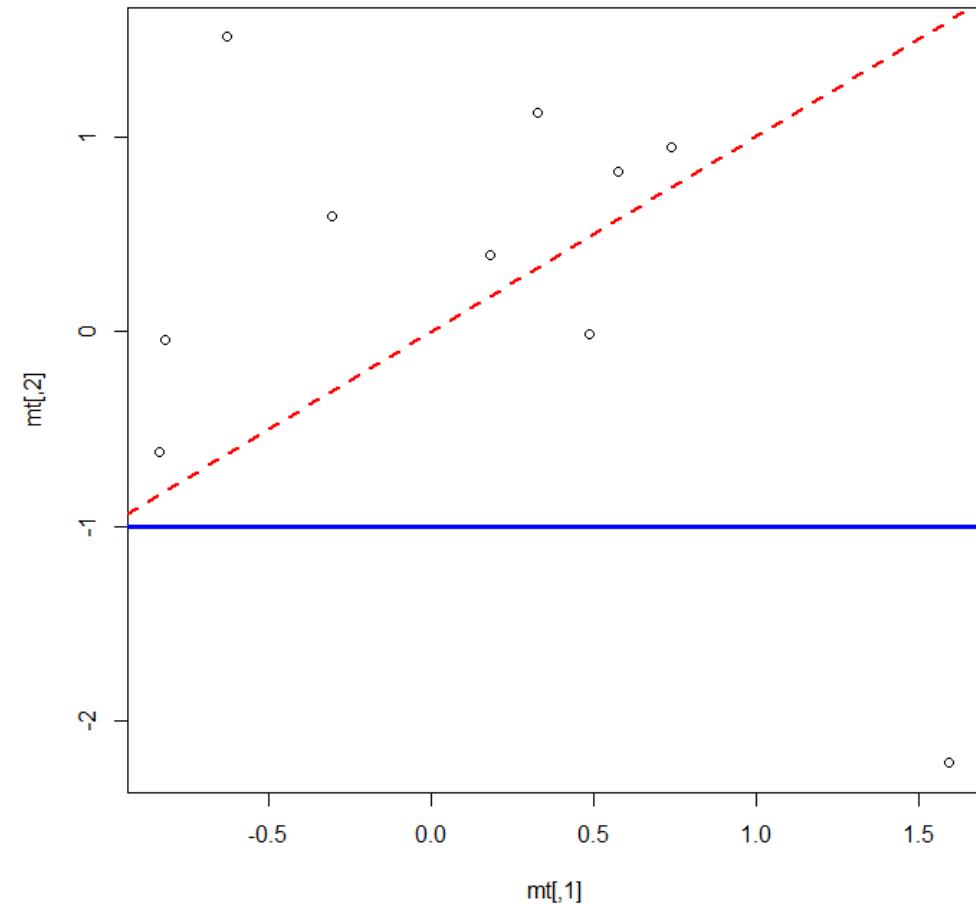
```
> plot(mt)
> abline(a = 0, b = 1,
+           col = "red", lwd = 2, lty = 2)
```



# Function: abline()

```
> plot(mt)
> abline(a = 0, b = 1,
+           col = "red", lwd = 2, lty = 2)
```

```
> abline(h = -1, col = "blue", lwd = 3)
```

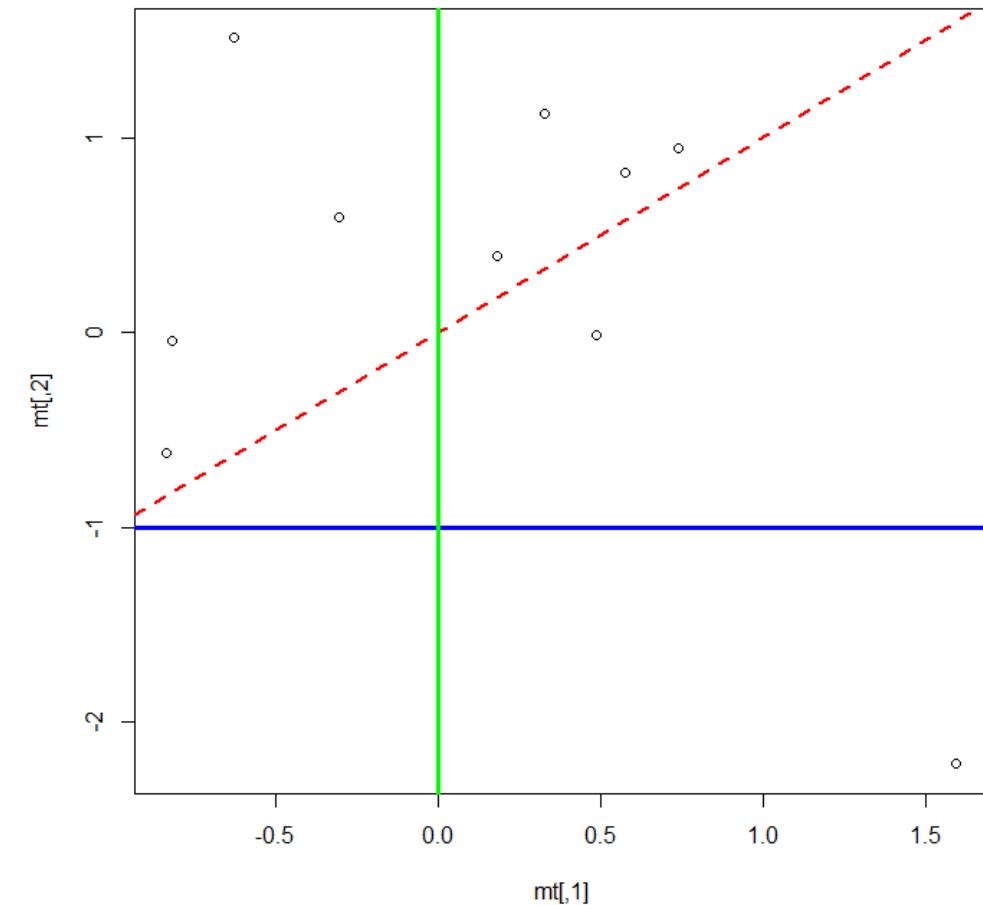


# Function: abline()

```
> plot(mt)
> abline(a = 0, b = 1,
+           col = "red", lwd = 2, lty = 2)
```

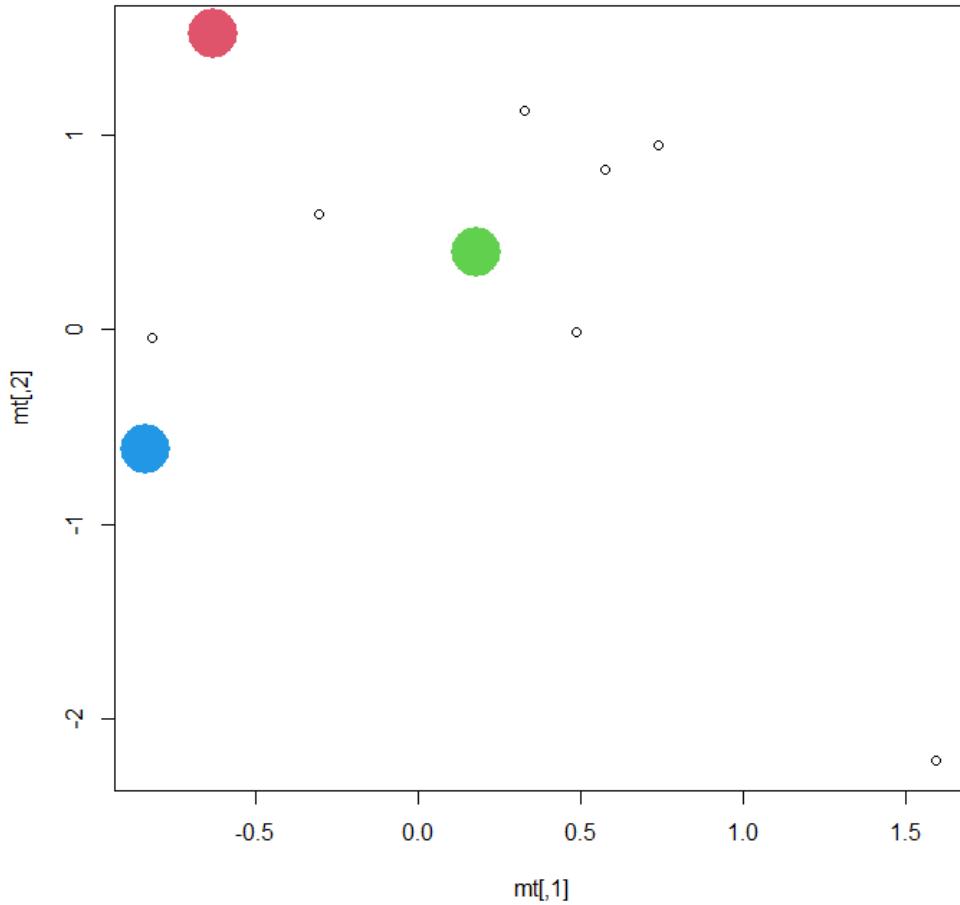
```
> abline(h = -1, col = "blue", lwd = 3)
```

```
> abline(v = 0, col = "green", lwd = 3)
```



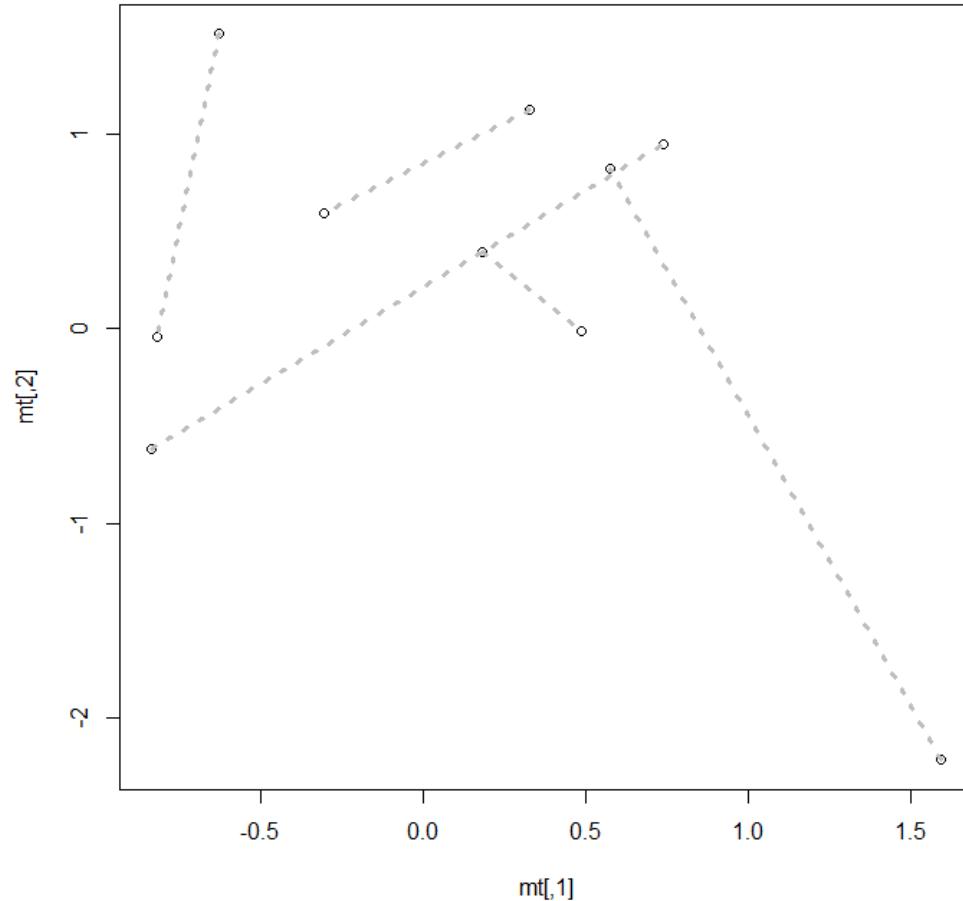
# Function: points()

```
> plot(mt)
> points(mt[1:3, ],
+          pch = 16, col = 2:4, cex = 5)
```



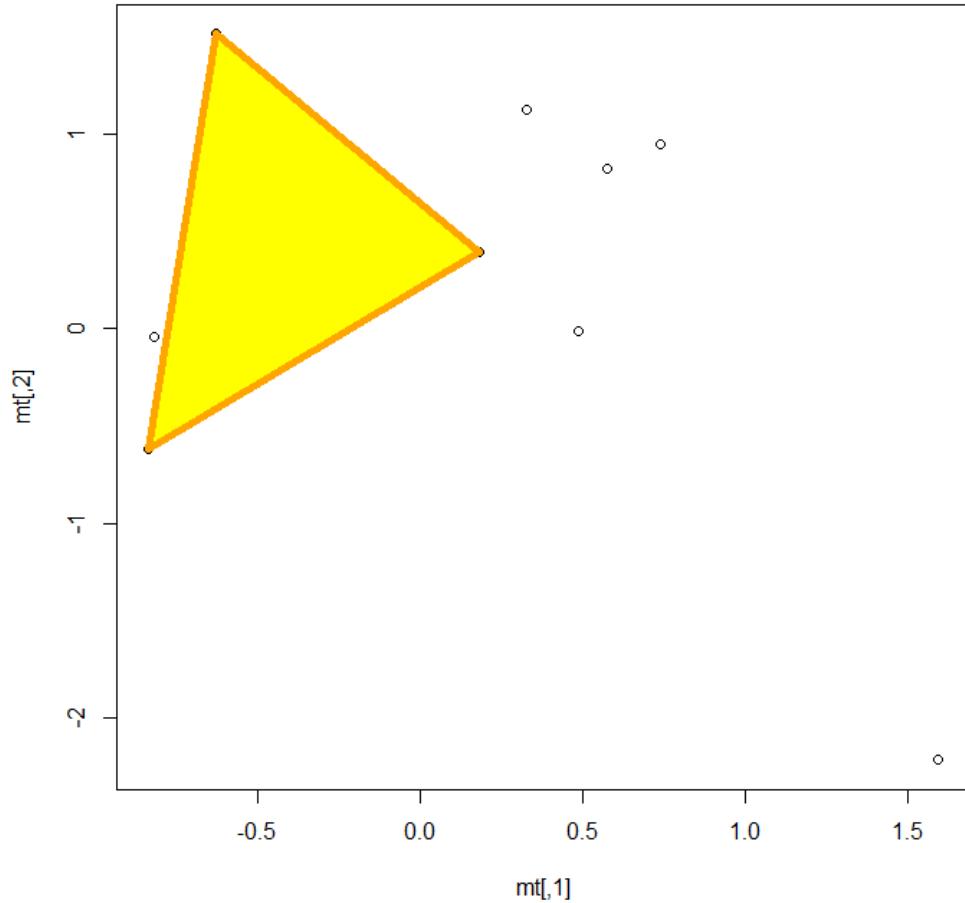
# Function: segments()

```
> plot(mt)
> segments(x0 = mt[1:5, 1],
+            y0 = mt[1:5, 2],
+            x1 = mt[6:10, 1],
+            y1 = mt[6:10, 2],
+            col = "grey",
+            lwd = 3, lty = 3)
```



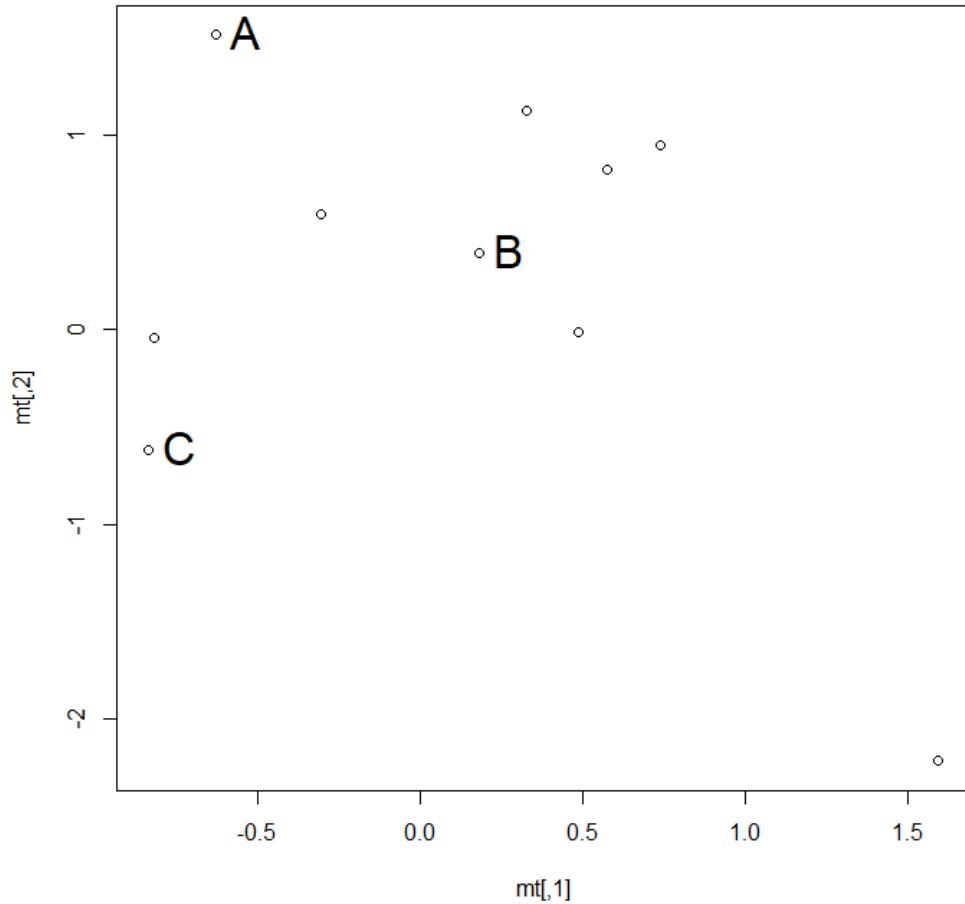
# Function: polygon()

```
> plot(mt)
> polygon(mt[1:3, ],
+           col = "yellow",
+           border = "orange",
+           lwd = 5)
```



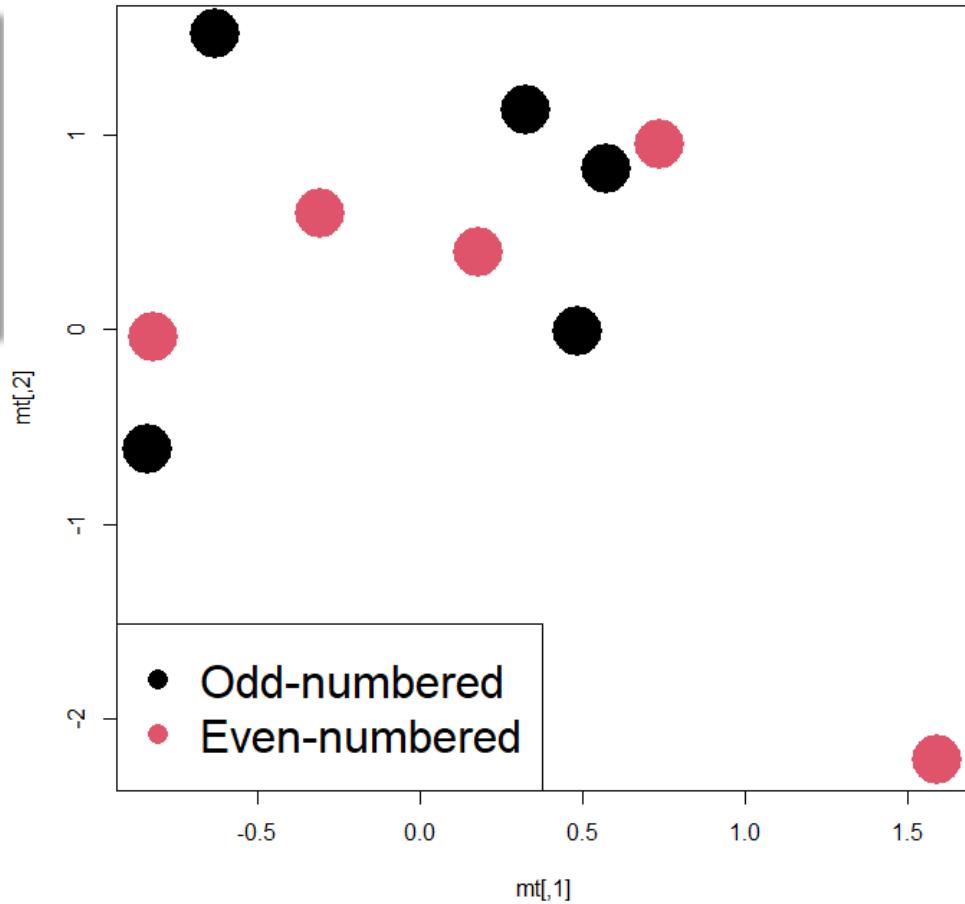
# Function: text()

```
> plot(mt)
> text(mt[1:3, ],
+       labels = LETTERS[1:3],
+       cex = 2,
+       pos = 4)
```



# Function: legend()

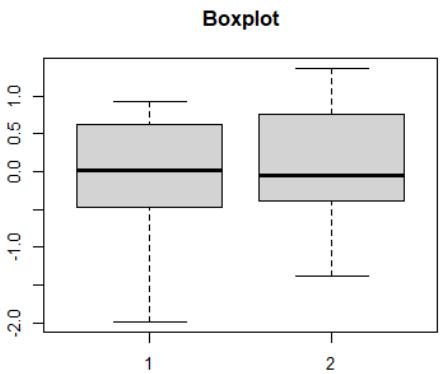
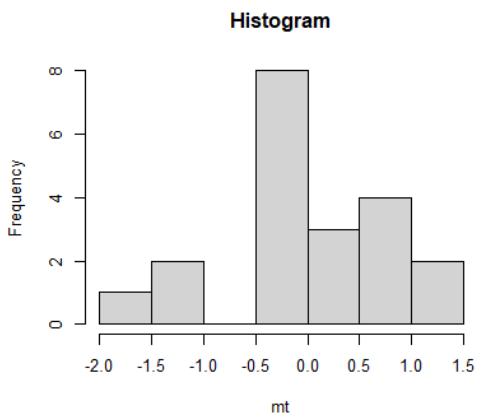
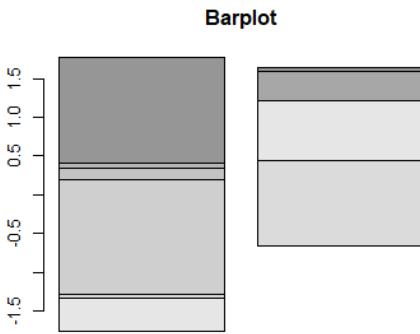
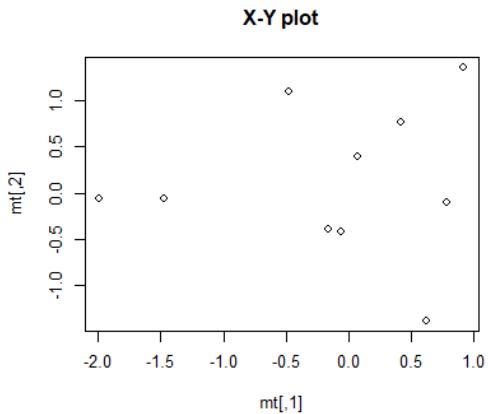
```
> plot(mt, col = 1:2, pch = 16, cex = 5)
> legend("bottomleft",
+         legend = c("Odd-numbered",
+                   "Even-numbered"),
+         col = 1:2, pch = 16, cex = 2)
```



# Function: par()

```
> par(mfrow = c(2, 2))
```

```
> plot(mt, main = "X-Y plot")
> barplot(mt, main = "Barplot")
> hist(mt, main = "Histogram")
> boxplot(mt, main = "Boxplot")
```



- High-level plotting functions
  - `plot(x, y)`
  - `hist()`; `barplot()`; `boxplot()`
- Low-level plotting functions
  - **Text:** `legend()`; `title()`; `text()`
  - **Points:** `points()`
  - **Lines:** `lines()`; `abline()`
  - **Shapes:** `rect()`; `polygon()`
- Parameter: `par()`

`help(package = "graphics")` for more information

# The ggplot2 package

- Layer by layer
- Geometric shapes, functions starting with `geom()`
- Aesthetics including color, arguments `aes()`
- Others: scales, facets, theme, statistical transformation, etc.



- Data types: numerical, categorical & ordinal
- 1-D\* numerical data: density lines (summary stats)
- 1-D categorical: barplot of counts
- 2-D numerical: scatterplot
- 2-D numerical vs categorical: boxplot/violin
- 2-D categorical: heatmap
- 3-D numerical/categorical: scatterplot w/ other aesthetics
- >3-D numerical/categorical: facet

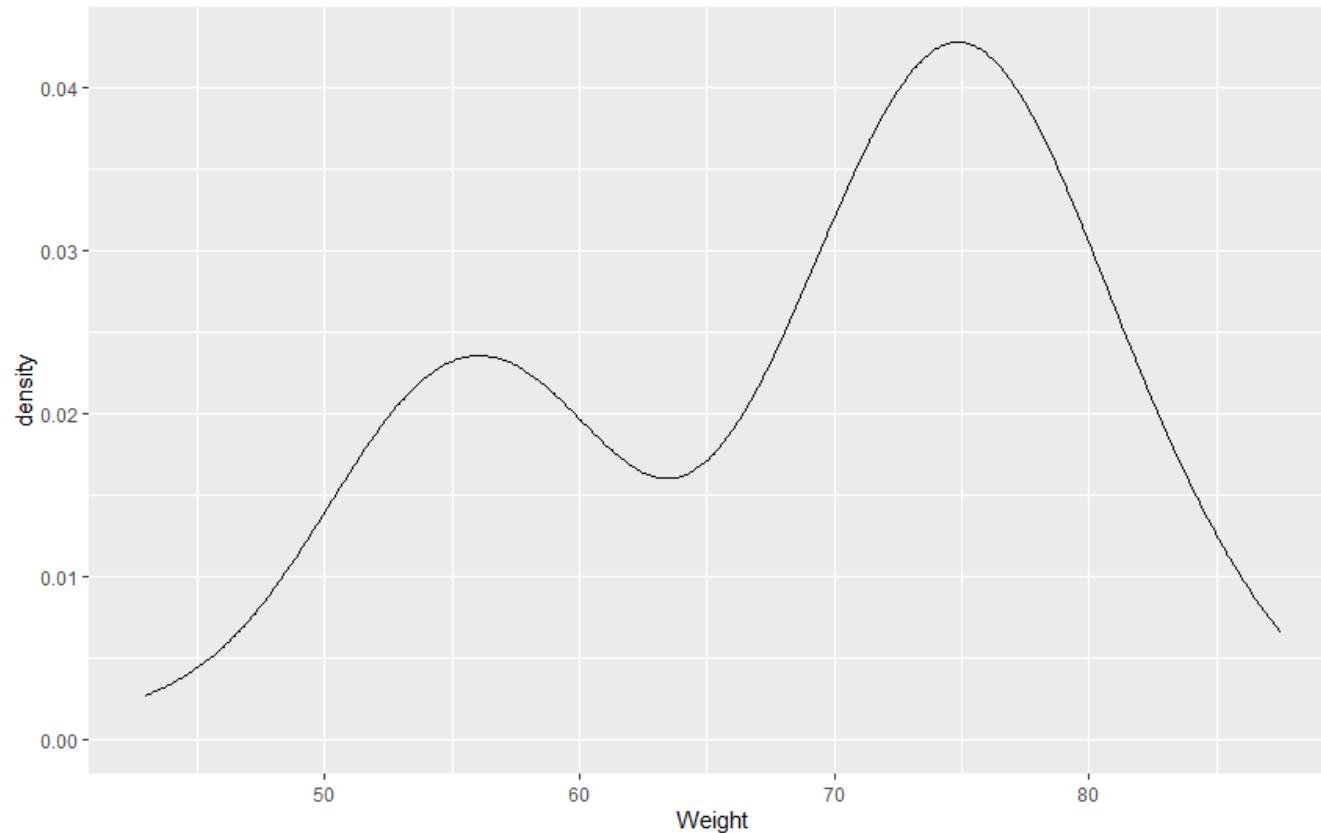
# A simulated data frame

```
set.seed(1)

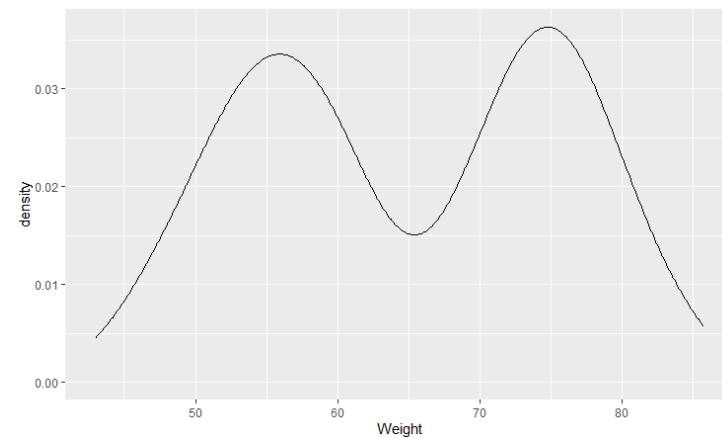
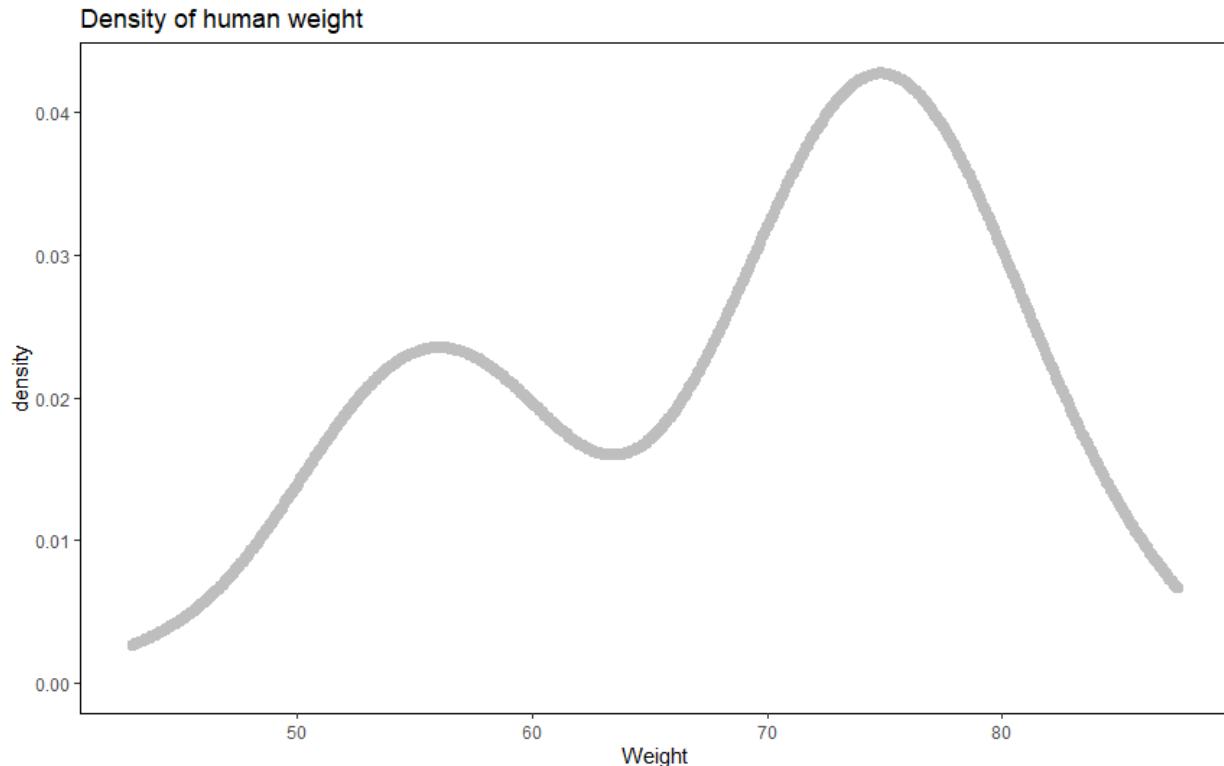
dat <- data.frame(Height = c(rnorm(100, 175, 10),
                             rnorm(50, 165, 10)),
                  Weight = c(rnorm(100, 75, 5),
                             rnorm(50, 55, 5)),
                  Sex = rep(c("Male", "Female"),
                            c(100, 50)),
                  Age = sample(rep(c("Young", "Old"),
                                  each = 75)))
```

```
> head(dat)
  Height  weight Sex Age
1 168.7355 77.25094 Male Old
2 176.8364 74.90720 Male Old
3 166.6437 73.40966 Male Old
4 190.9528 70.35319 Male Old
5 178.2951 67.56270 Male Old
6 166.7953 69.62404 Male Old
```

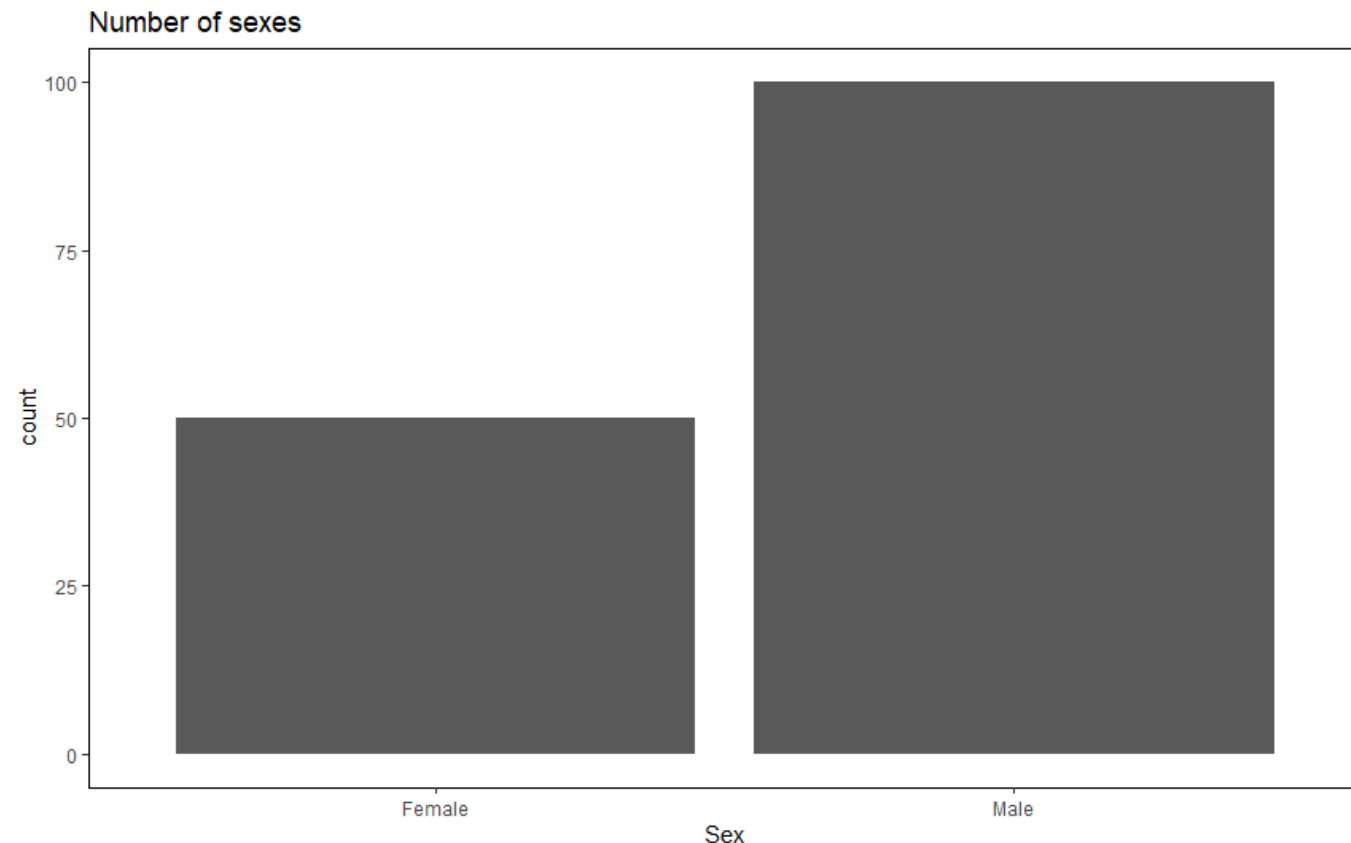
```
library(ggplot2)
g <- ggplot(dat)
g + geom_density(aes(x = Weight)) # density plot
```



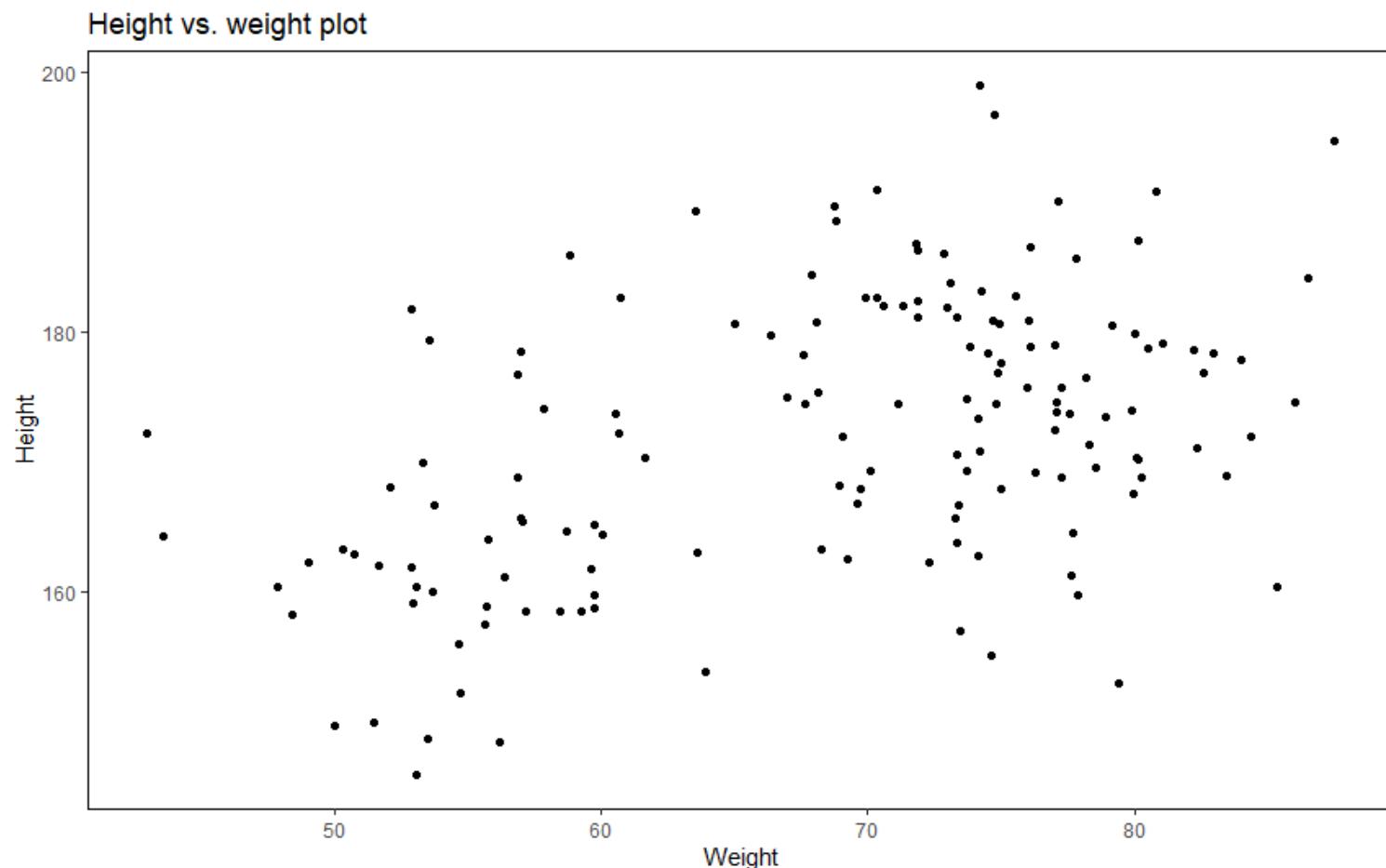
```
g + geom_density(aes(x = Weight),  
                  color = "grey", size = 3) +  
  labs(title = "Density of human weight") +  
  theme(panel.border = element_rect(color = "black", fill = NA),  
        panel.background = element_rect(fill = "white"))
```



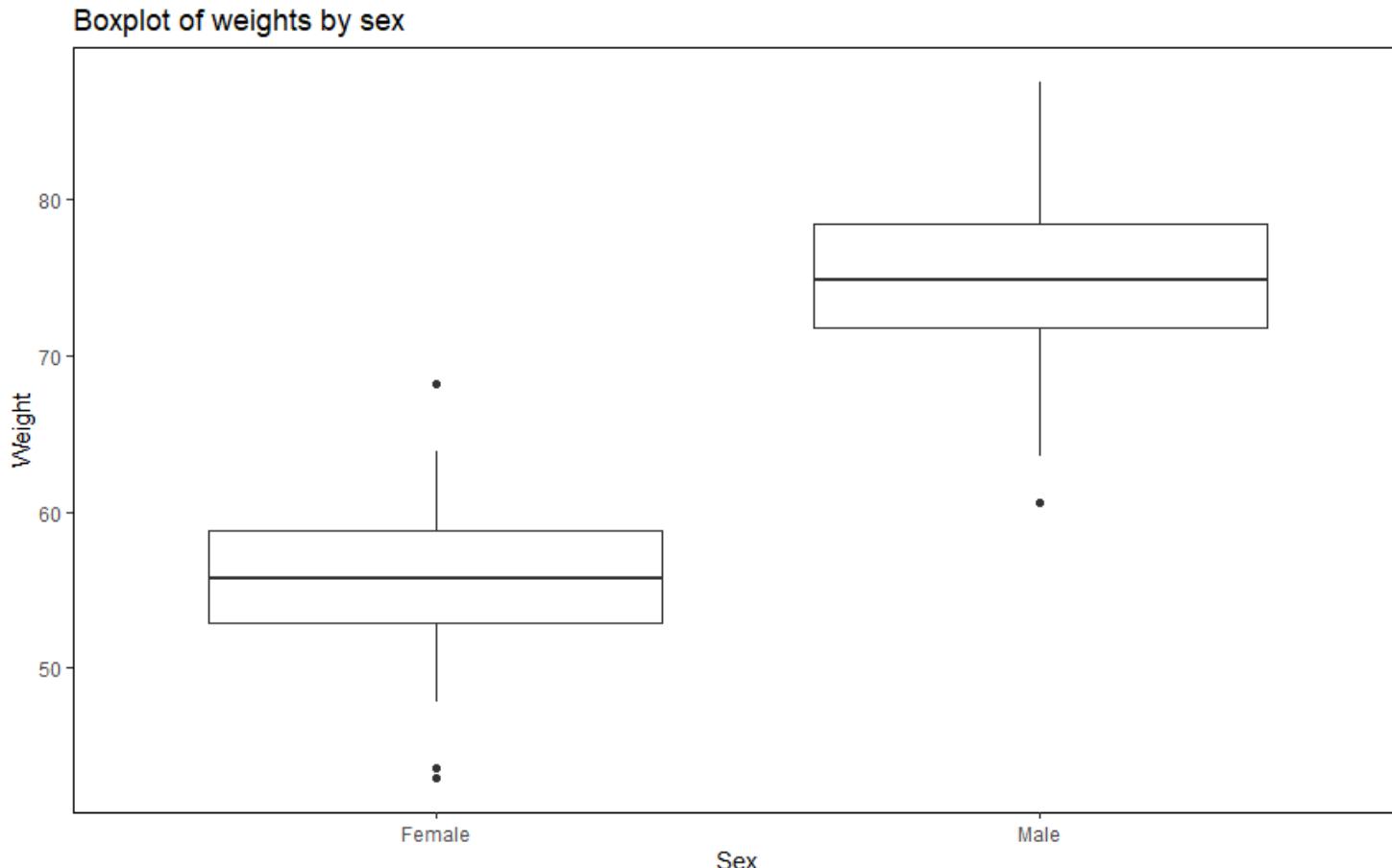
```
g <- ggplot(dat) +  
  theme(panel.border = element_rect(color = "black", fill = NA),  
        panel.background = element_rect(fill = "white"))  
g + geom_bar(aes(x = Sex)) + labs(title = "Number of sexes")
```



```
g + geom_point(aes(x = Weight, y = Height)) +  
  labs(title = "Height vs. weight plot")
```

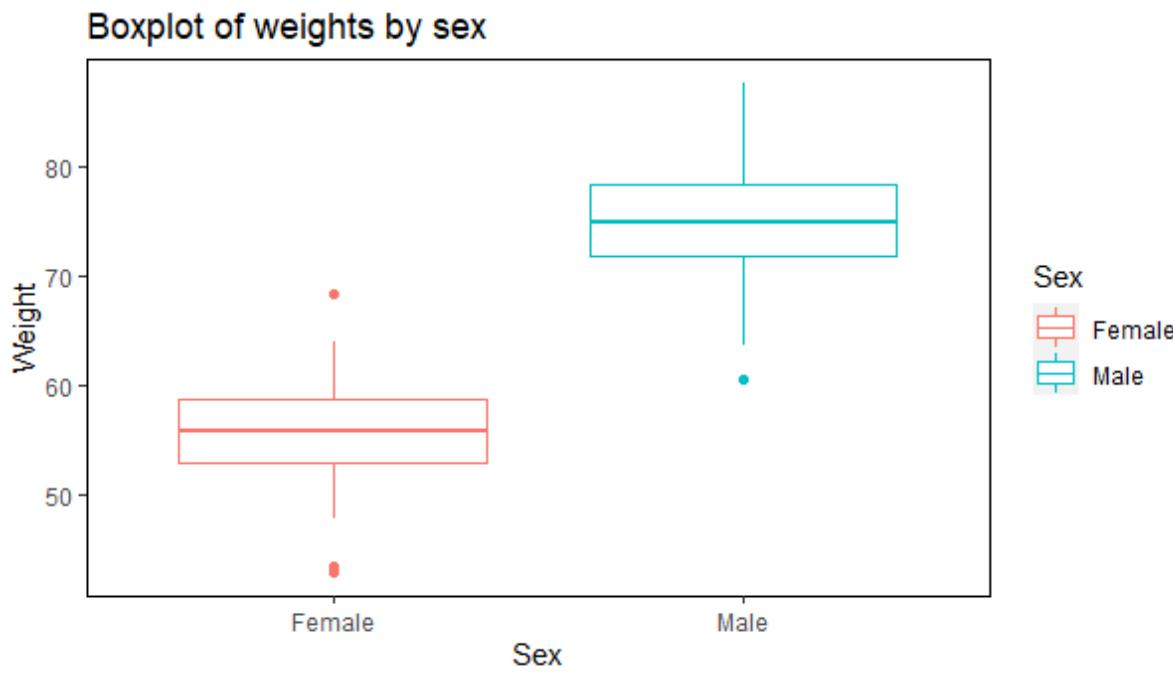


```
g + geom_boxplot(aes(x = Sex, y = Weight)) +  
  labs(title = "Boxplot of weights by sex")
```

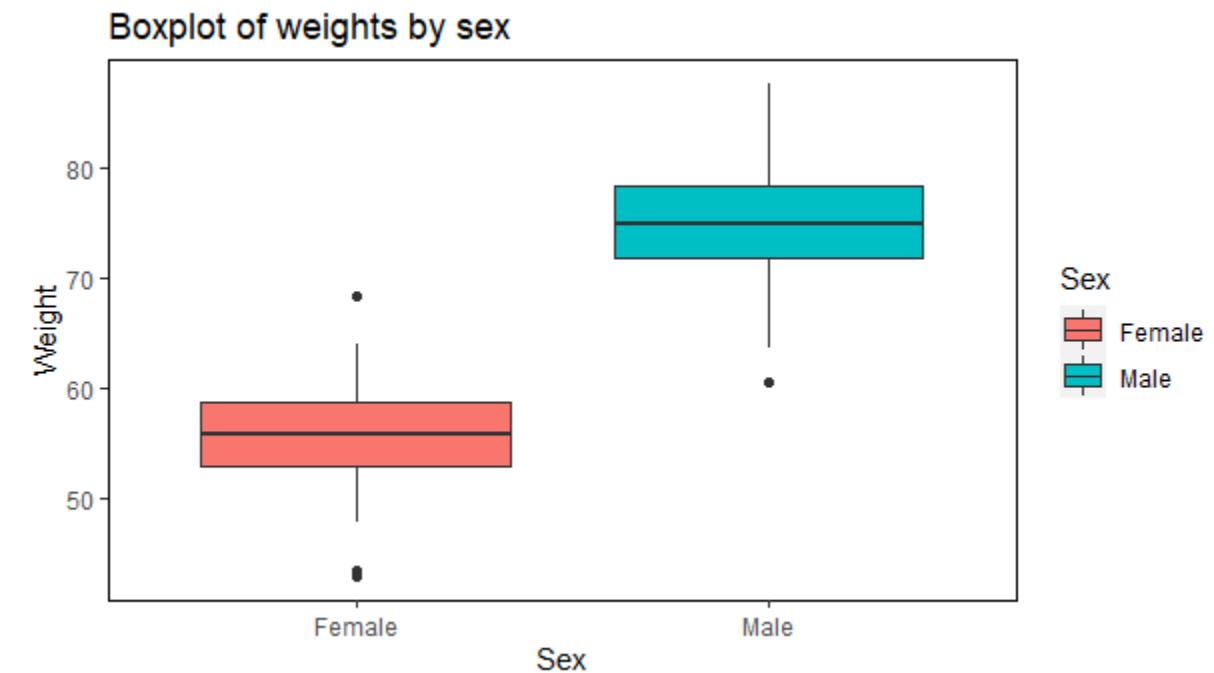


# Parameter: color vs. fill

```
g + geom_boxplot(aes(Sex, Weight, color = Sex)) +  
  labs(title = "Boxplot of weights by sex")
```



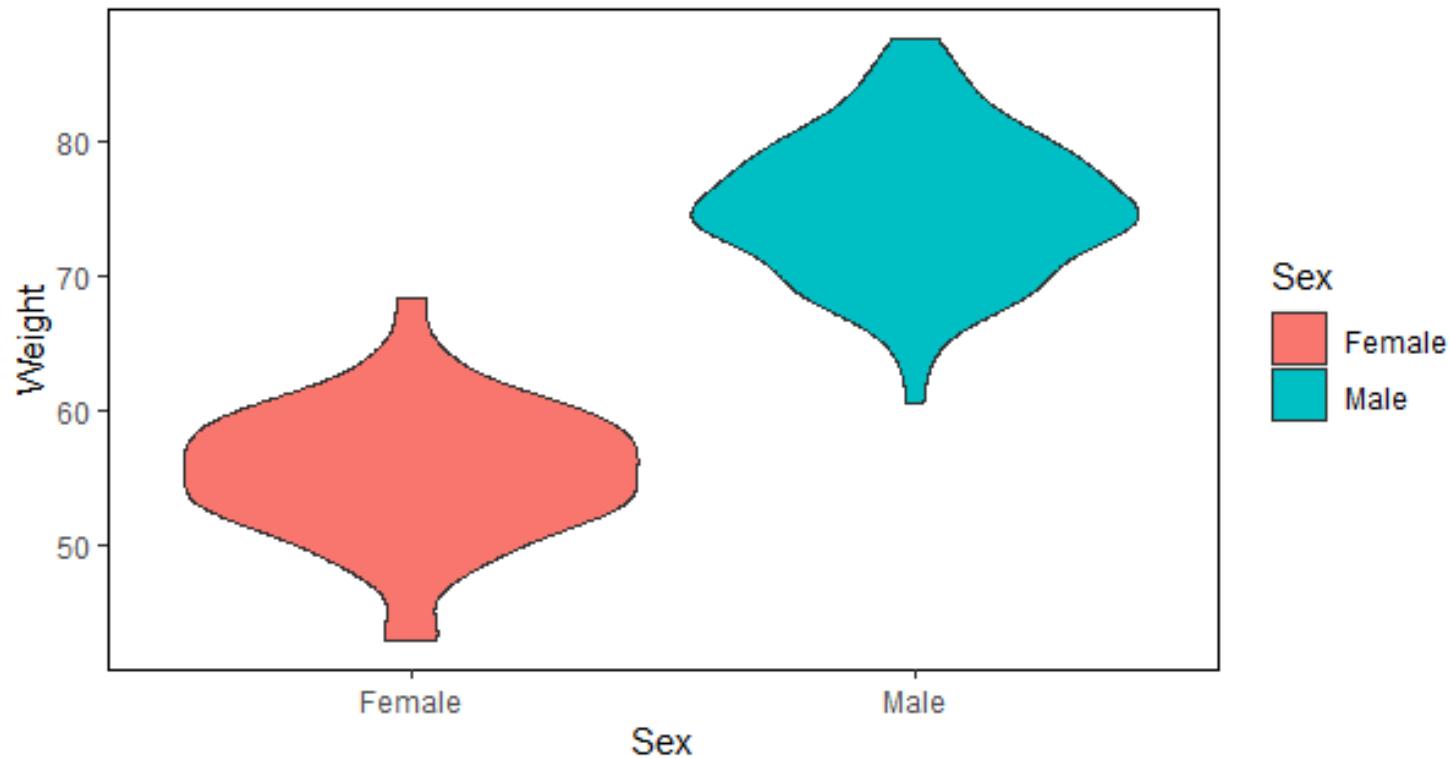
```
g + geom_boxplot(aes(Sex, Weight, fill = Sex)) +  
  labs(title = "Boxplot of weights by sex")
```



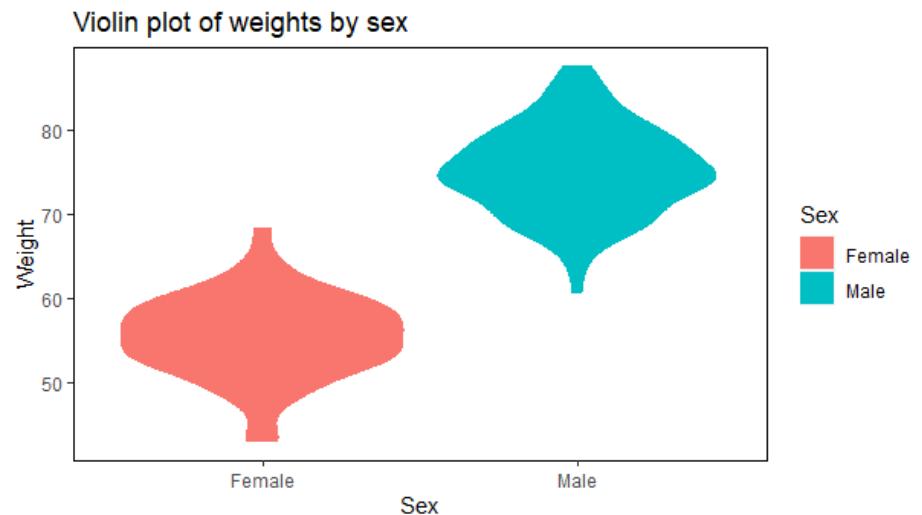
# Violin plot

```
g + geom_violin(aes(x = Sex, y = Weight, fill = Sex)) +  
  labs(title = "Violin plot of weights by sex")
```

Violin plot of weights by sex

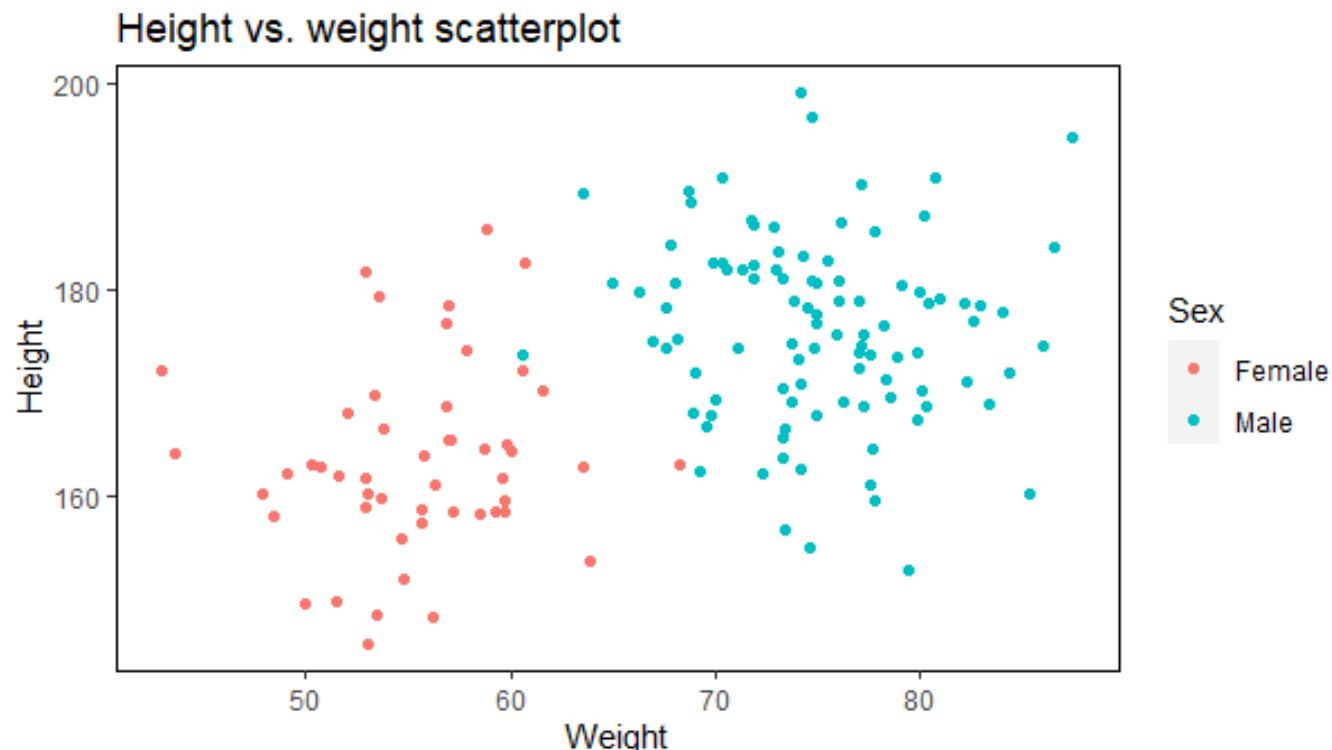


geom\_violin(..., color = NA)

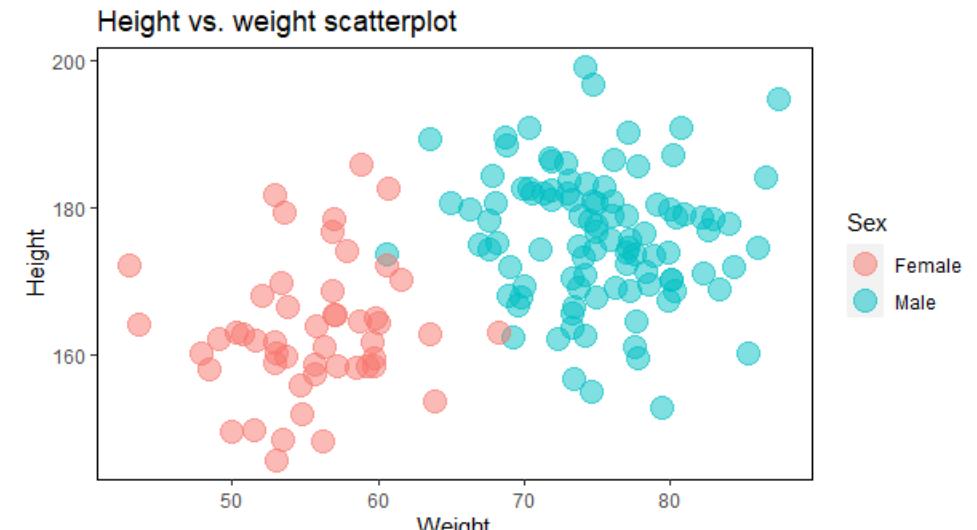


# Scatterplot w/ a 3<sup>rd</sup> variable

```
g + geom_point(aes(x = Weight, y = Height, color = Sex)) +  
  labs(title = "Height vs. weight scatterplot")
```

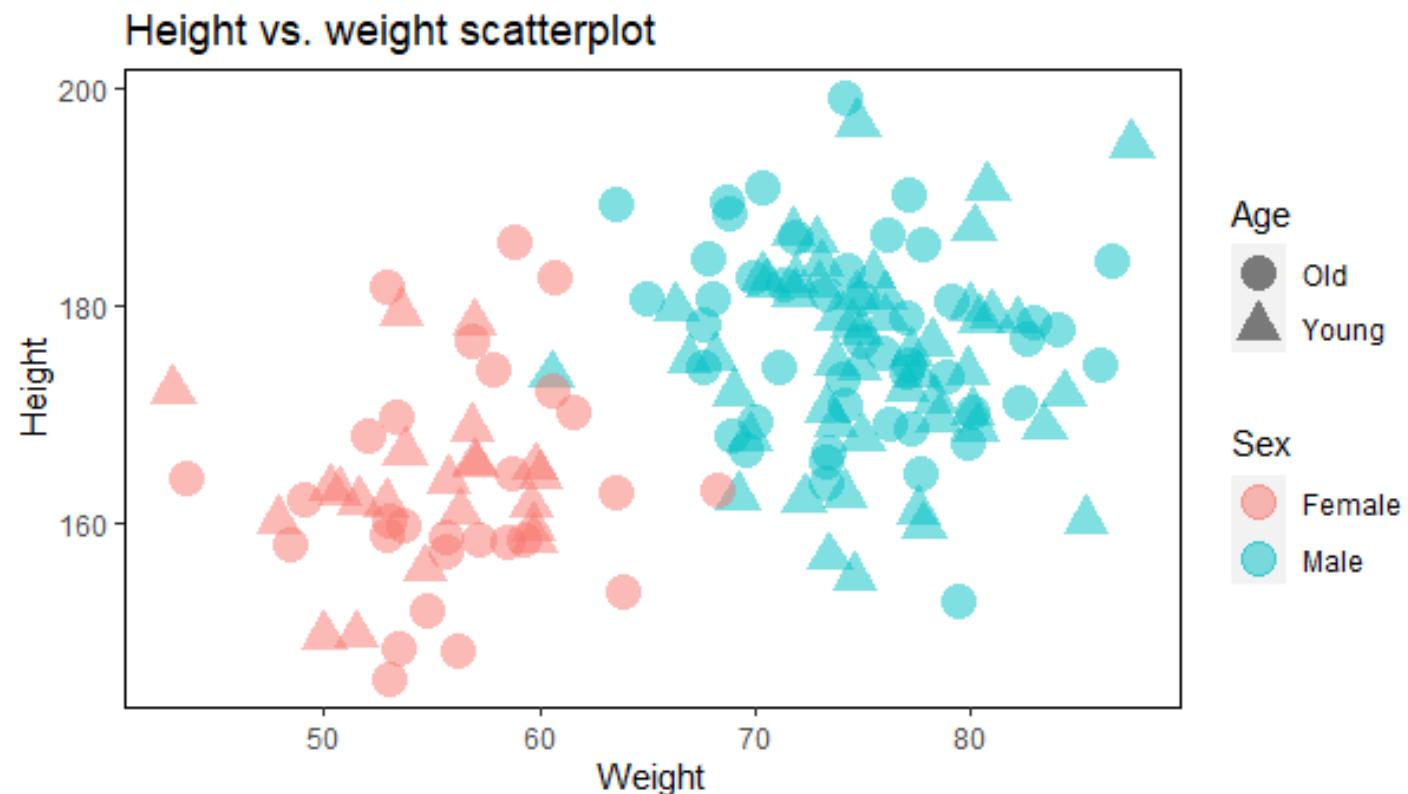


```
geom_point(..., size = 5, alpha = 0.5)
```

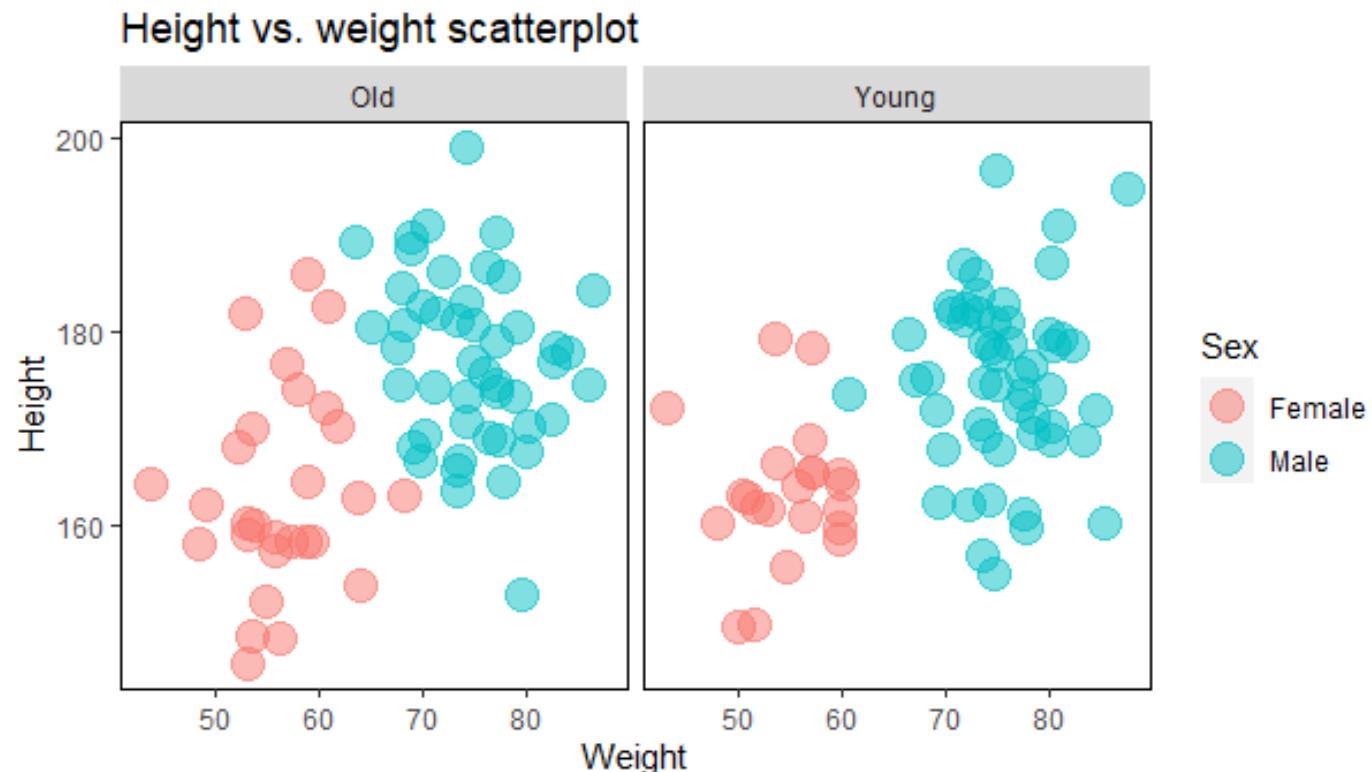


# 4-D data: aesthetics

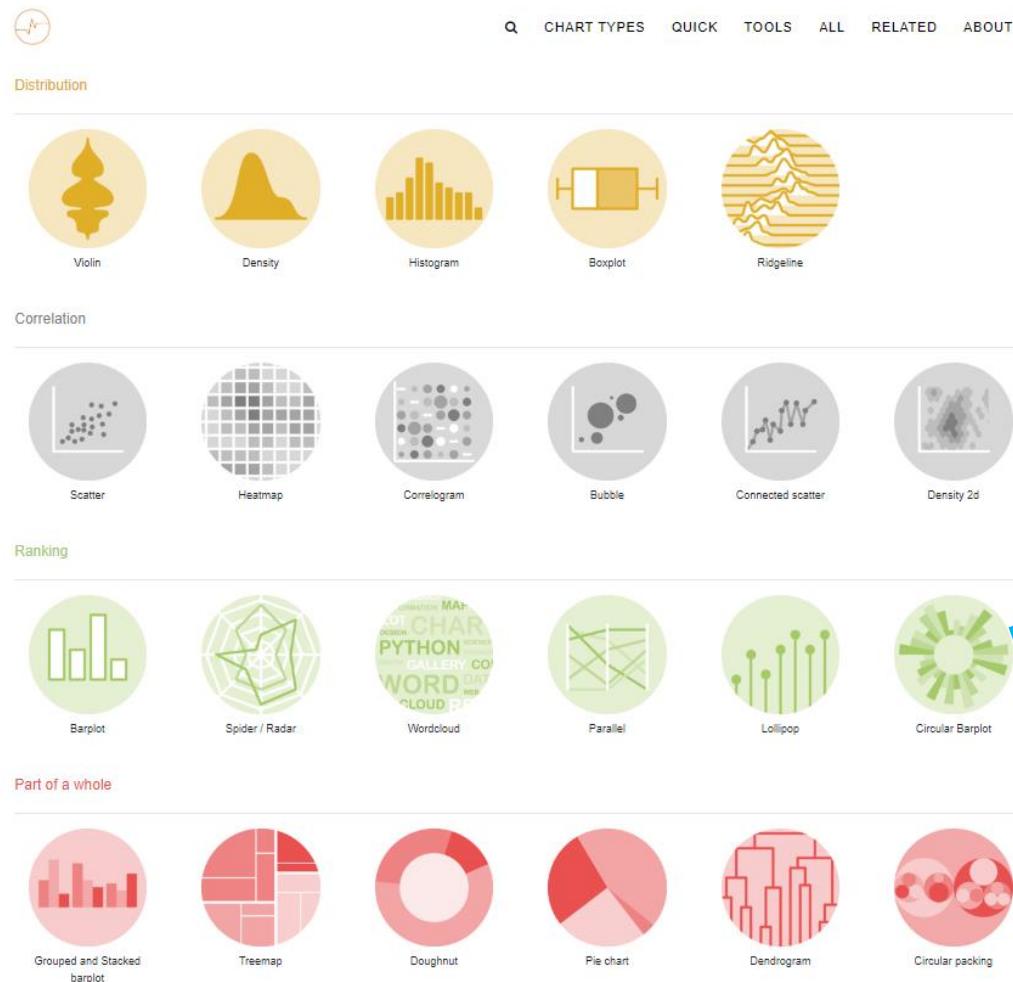
```
g + geom_point(aes(x = Weight, y = Height,  
                    color = Sex, shape = Age),  
                    size = 5, alpha = .5) +  
  labs(title = "Height vs. weight scatterplot")
```



```
g + geom_point(aes(x = Weight, y = Height, color = Sex),  
               size = 5, alpha = .5) +  
  labs(title = "Height vs. weight scatterplot") +  
facet_wrap(facets = vars(Age))
```



# More plots in R gallery



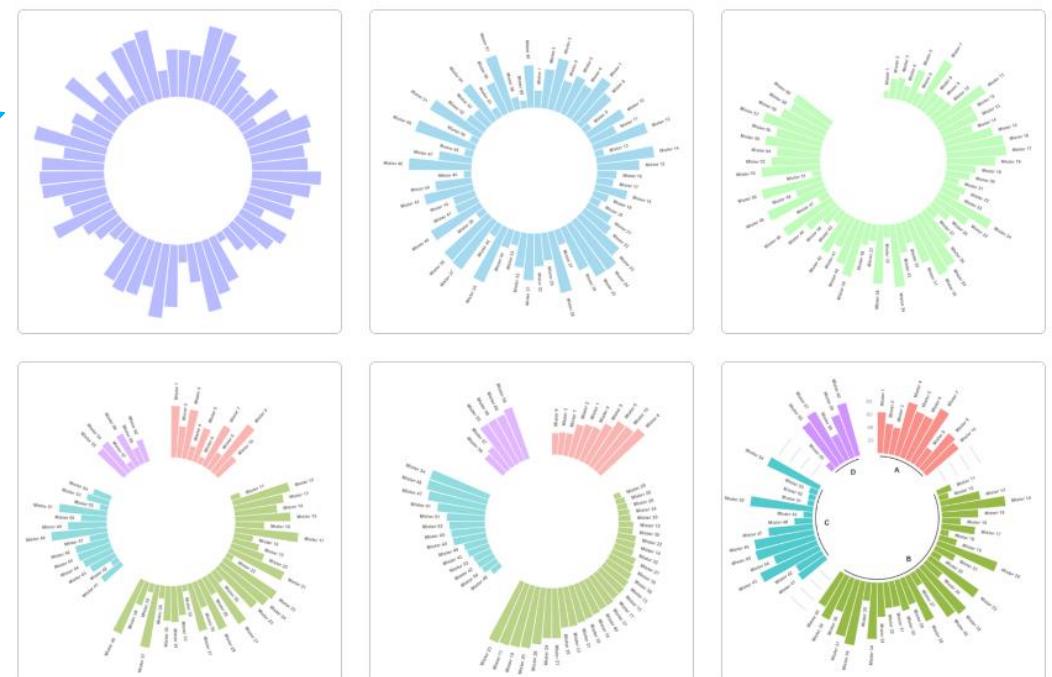
## Circular barplot



This is the [circular barplot](#) section of the gallery, a variation of the well known [barplot](#). Note that even if visually appealing, [circular barplot](#) must be used with care since groups [do not share the same Y axis](#). It is very adapted for cyclical data though. Visit [data-to-viz.com](#) for more info.

### STEP BY STEP

Here is a set of examples leading to a proper circular barplot, step by step. The first [most basic circular barchart](#) shows how to use `coord_polar()` to make the barchart circular. Next examples describe the next steps to get a proper figure: [gap](#) between groups, [labels](#) and customization.



# Part II

# A case study

# Case I. Lettuce domestication

nature  
genetics

ARTICLES

<https://doi.org/10.1038/s41588-021-00831-0> Check for updates

## Whole-genome resequencing of 445 *Lactuca* accessions reveals the domestication history of cultivated lettuce

Tong Wei<sup>1,11</sup>, Rob van Treuren<sup>1,2,11</sup>✉, Xinjiang Liu<sup>1,11</sup>, Zhaowu Zhang<sup>1,3</sup>, Jiongjiong Chen<sup>4</sup>, Yang Liu<sup>1</sup>, Shanshan Dong<sup>5</sup>, Peinan Sun<sup>4</sup>, Ting Yang<sup>1</sup>, Tianming Lan<sup>1,6</sup>, Xiaogang Wang<sup>7</sup>, Zhouquan Xiong<sup>7</sup>, Yaqiong Liu<sup>8</sup>, Jinpu Wei<sup>8</sup>, Haorong Lu<sup>1,8</sup>, Shengping Han<sup>8</sup>, Jason C. Chen<sup>8</sup>, Xuemei Ni<sup>1</sup>, Jian Wang<sup>1,9</sup>, Huanming Yang<sup>1,9</sup>, Xun Xu<sup>1,10</sup>, Hanhui Kuang<sup>4</sup>, Theo van Hintum<sup>2</sup>, Xin Liu<sup>1,11</sup>✉ and Huan Liu<sup>1</sup>✉

Lettuce (*Lactuca sativa*) is an important vegetable crop worldwide. Cultivated lettuce is believed to be domesticated from *L. serriola*; however, its origins and domestication history remain to be elucidated. Here, we sequenced a total of 445 *Lactuca* accessions, including major lettuce crop types and wild relative species, and generated a comprehensive map of lettuce genome variations. In-depth analyses of population structure and demography revealed that lettuce was first domesticated near the Caucasus, which was marked by loss of seed shattering. We also identified the genetic architecture of other domestication traits and wild introgressions in major resistance clusters in the lettuce genome. This study provides valuable genomic resources for crop breeding and sheds light on the domestication history of cultivated lettuce.

# Case I. Lettuce population

- As a vegetable crop
  - A rich source of vitamin K and vitamin A, and a moderate source of folate and iron
  - Consumed worldwide
- As a model Asteraceae plant
  - A short life cycle; easy for transformation
  - Well-maintained germplasms
  - Various agronomic traits
  - Potential bioreactor



var. *longifoliaf*



var. *angustanairish*

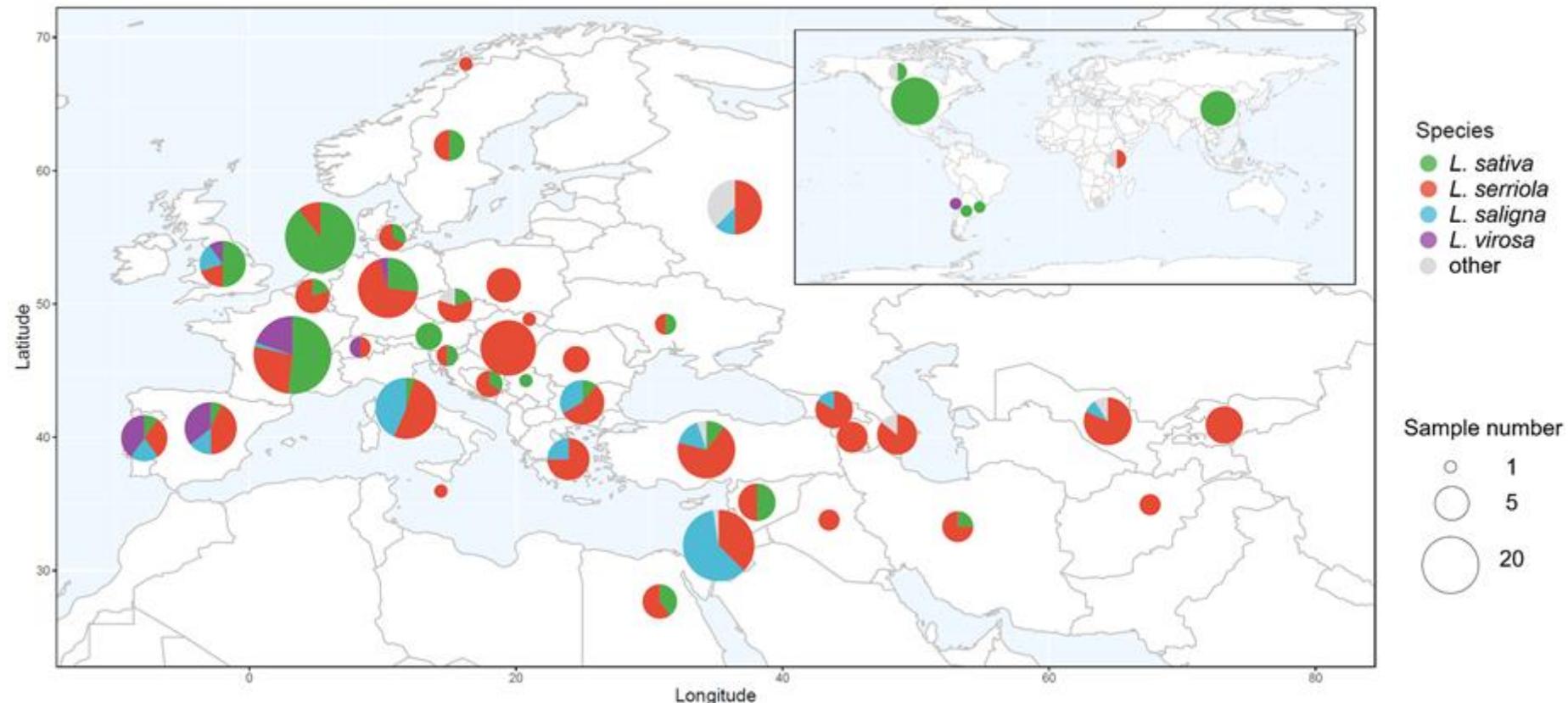
# Case I. Scientific questions



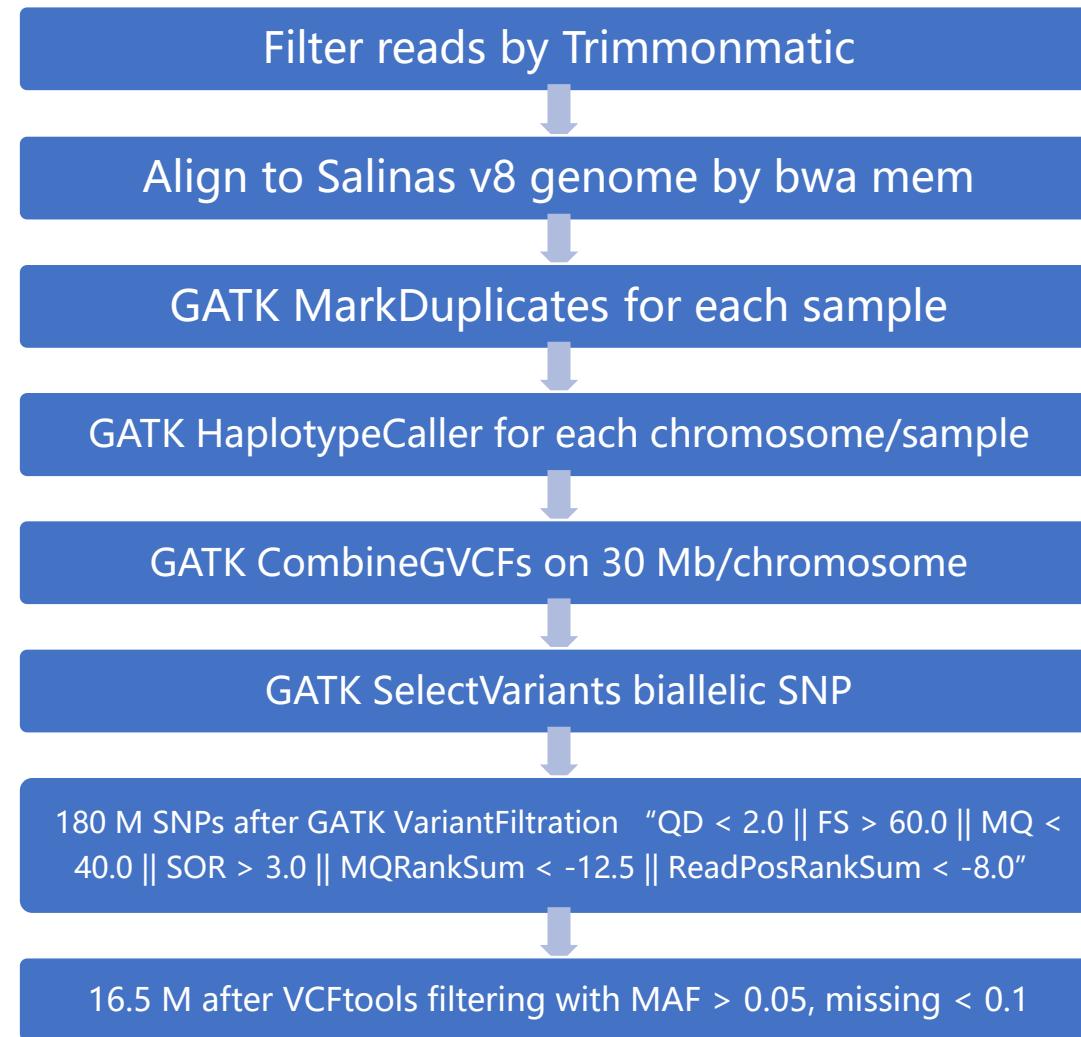
1. **Where** was lettuce domesticated?
2. **When** was lettuce domesticated?
3. What are the **key events** during lettuce domestication?
4. What are the **genetic determinants** for the domestication traits?

# Case I. Global representation in world map

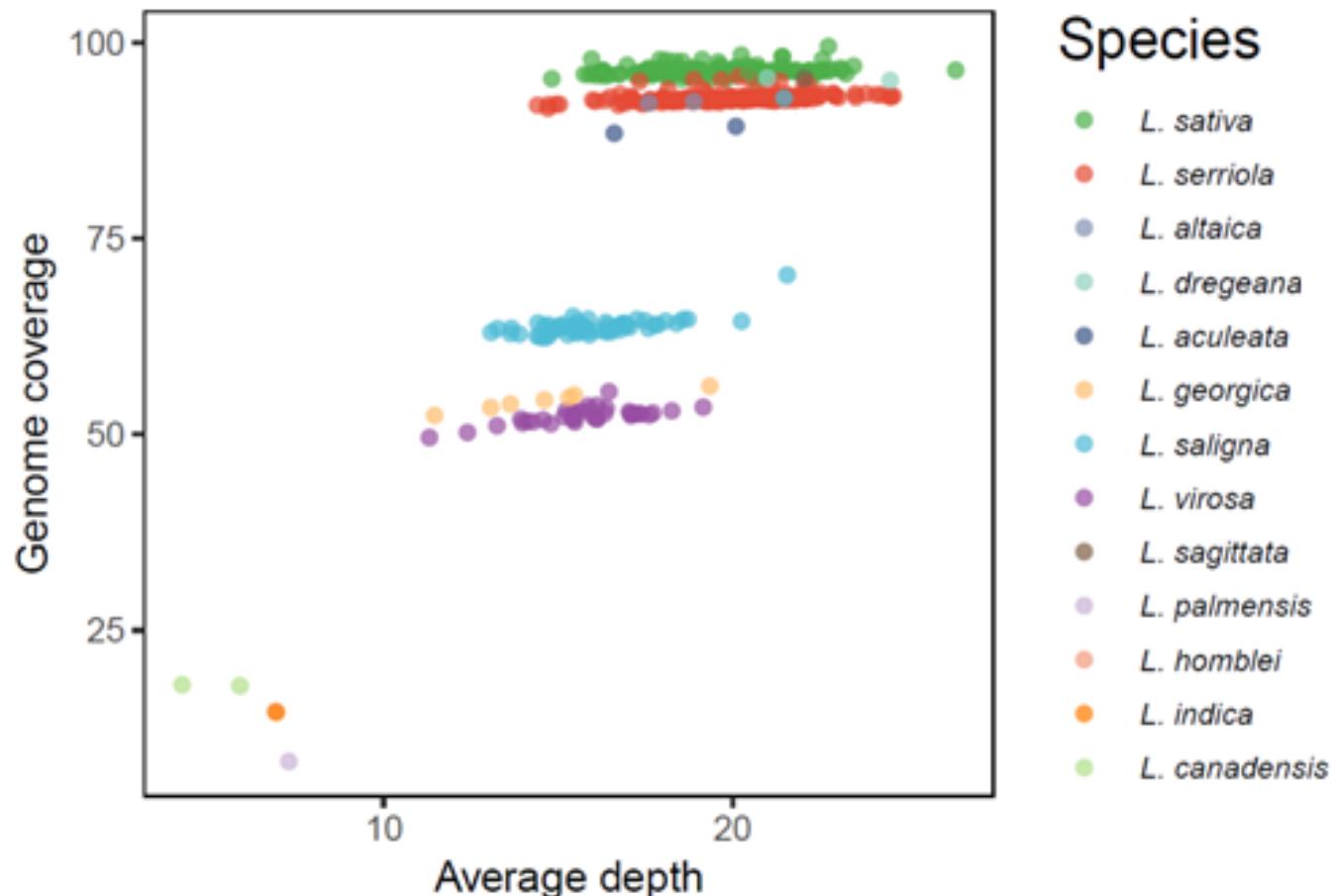
445 lines from 47 countries, mainly composed of GP breeding materials of *L. sativa*, *L. serriola*, *L. saligna* and *L. virosa*



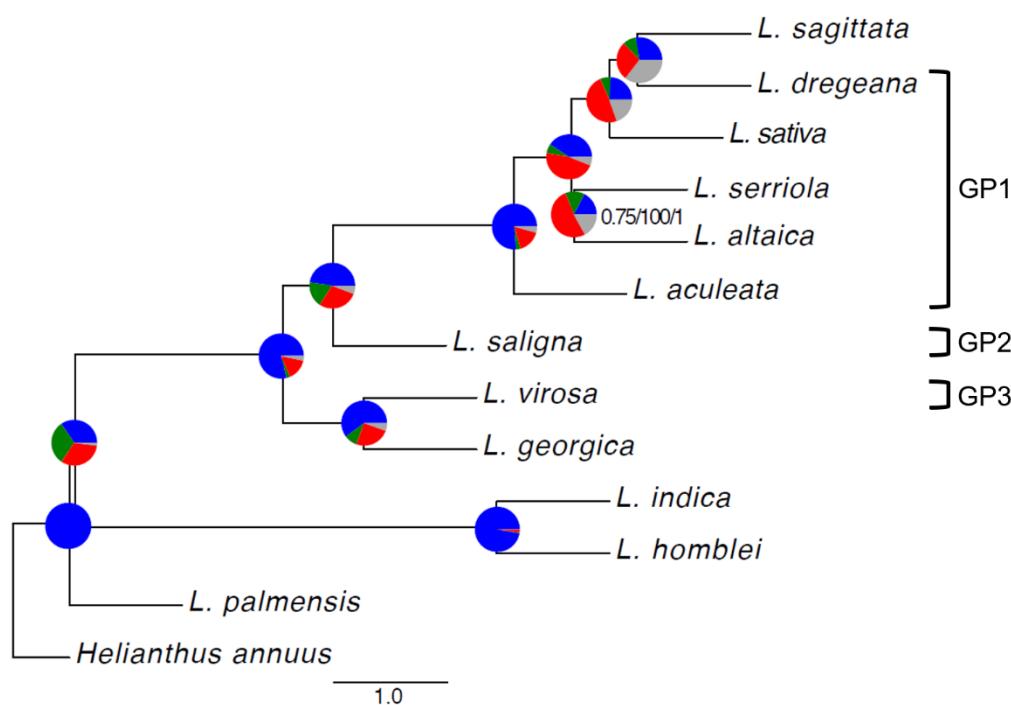
# Case I. Pipeline in flow chart



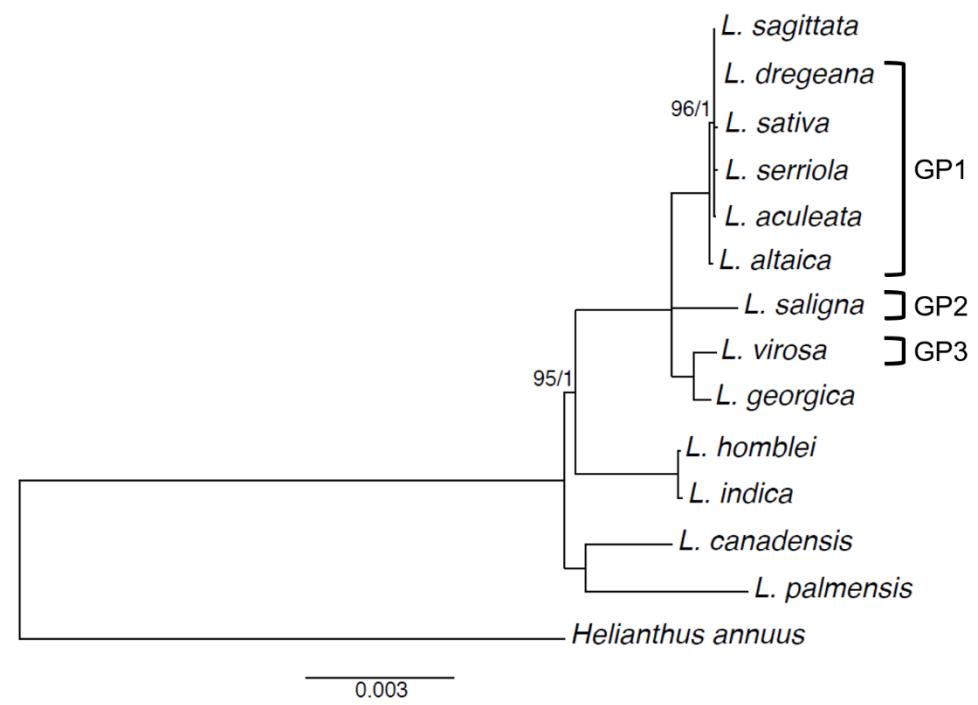
# Case I. Read alignment in scatter plot



# Case I. Phylogeny in tree plot

**a**

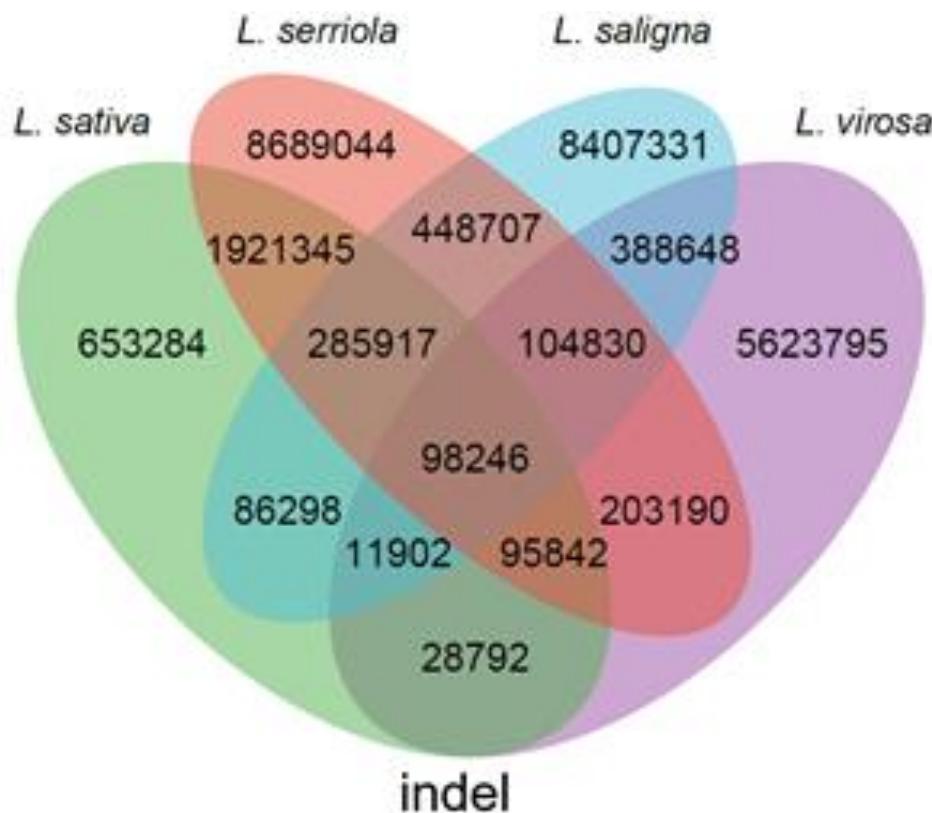
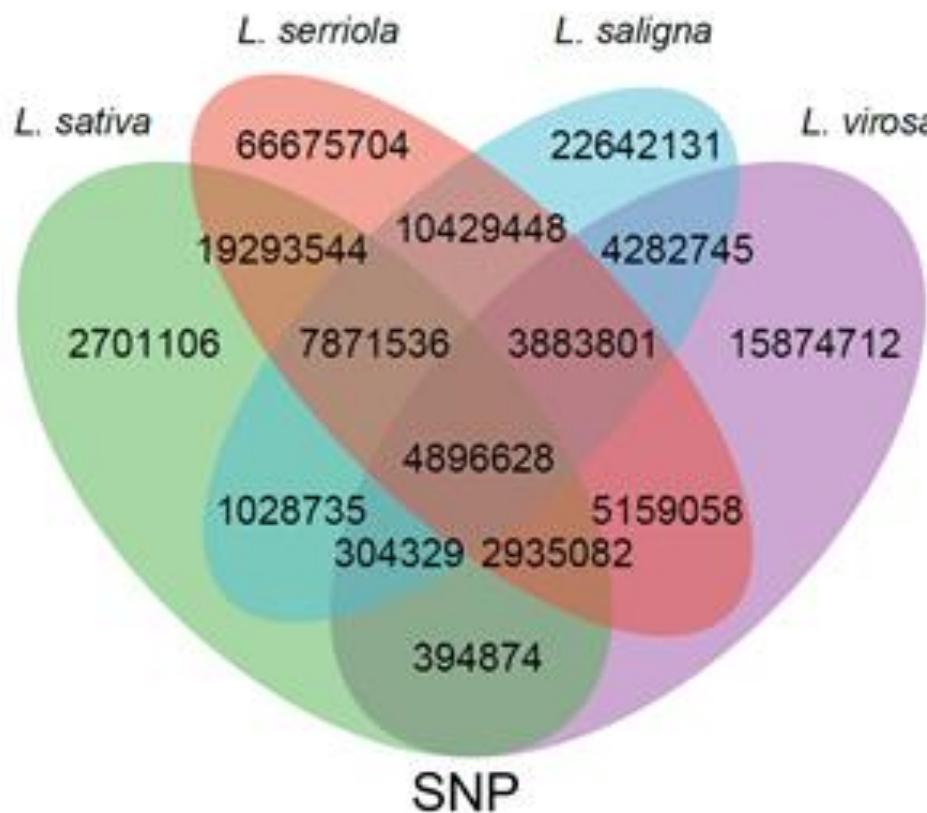
Coalescence-based tree from  
4,513 nuclear genes

**b**

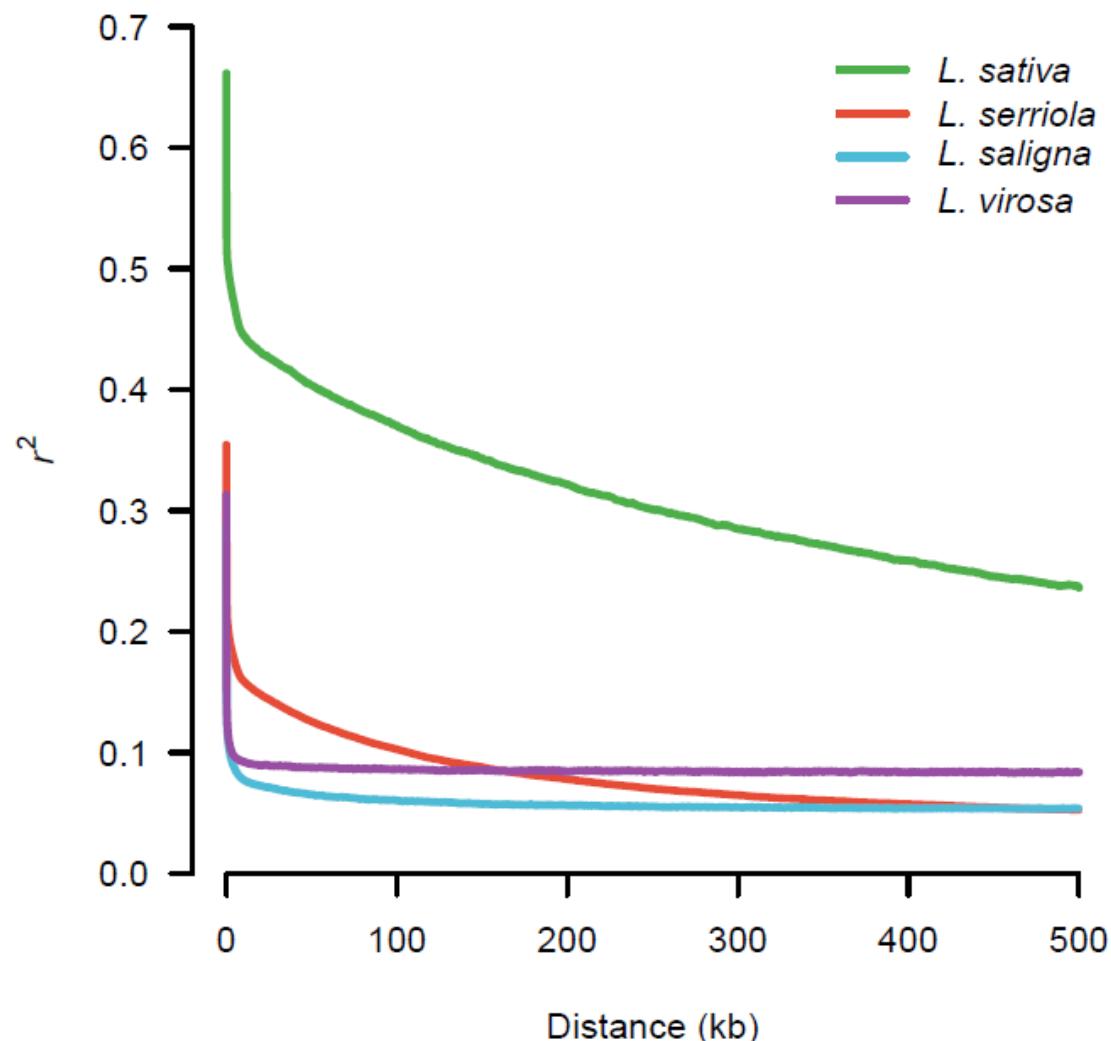
Concatenation-based tree from  
75 plastid genes

# Case I. SNP/indel stats in Venn diagram

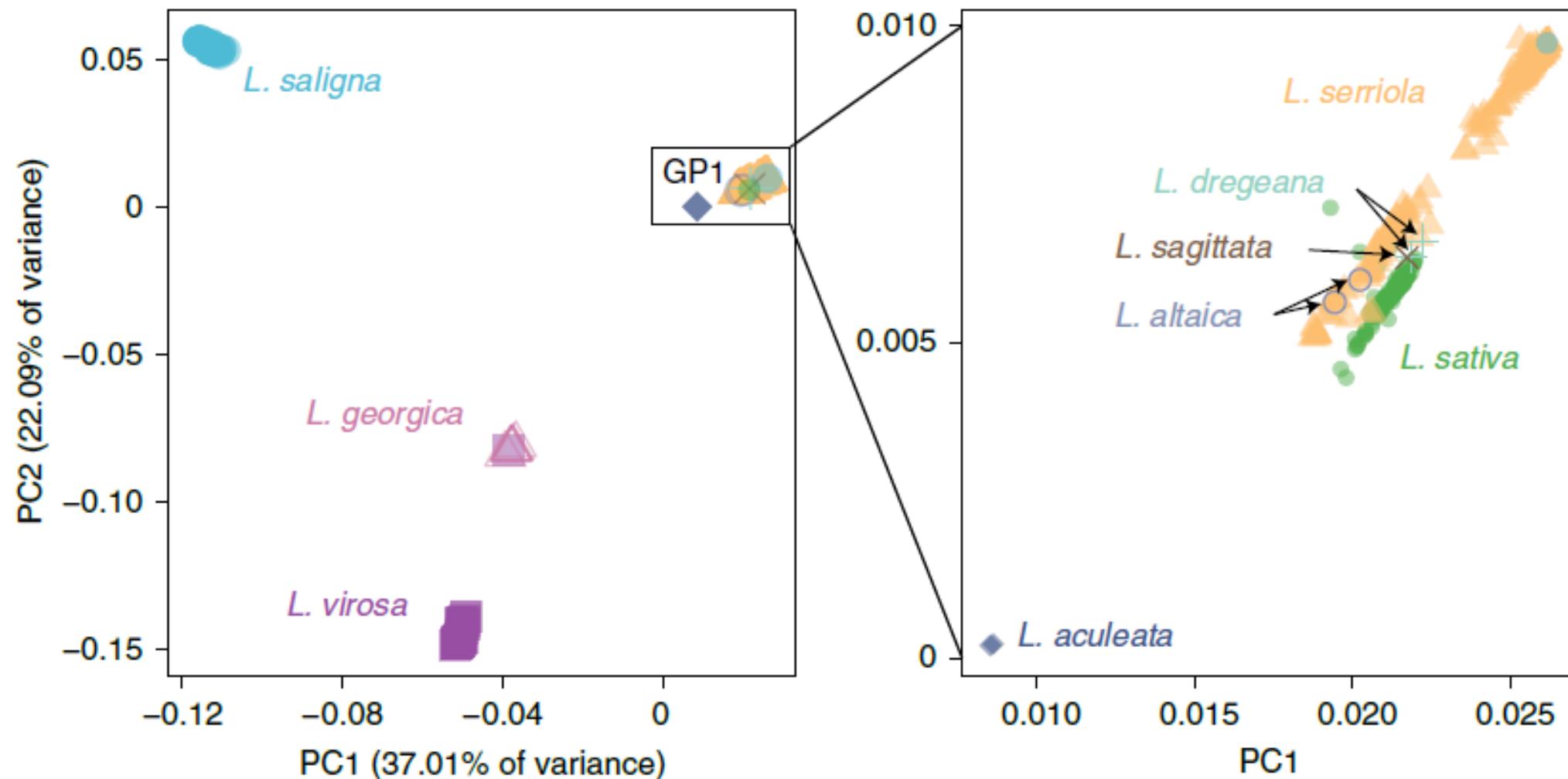
88.51% SNPs in *L. sativa* are shared with *L. serriola*.



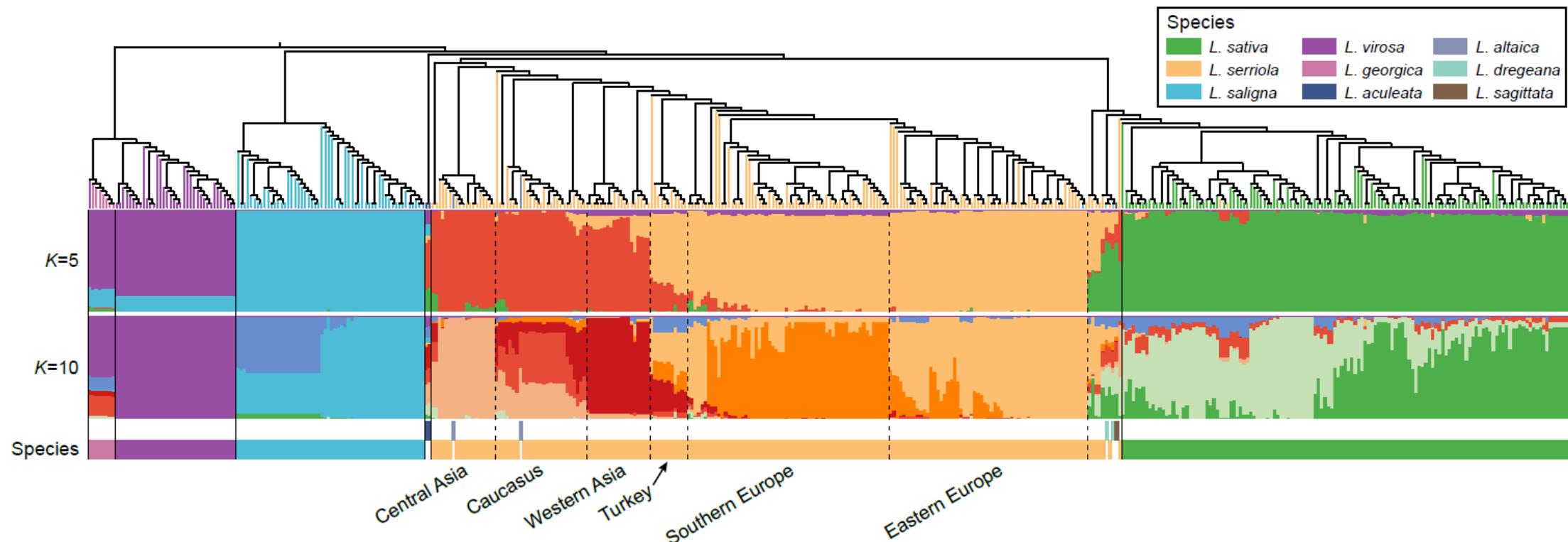
# Case I. Linkage disequilibrium in line chart



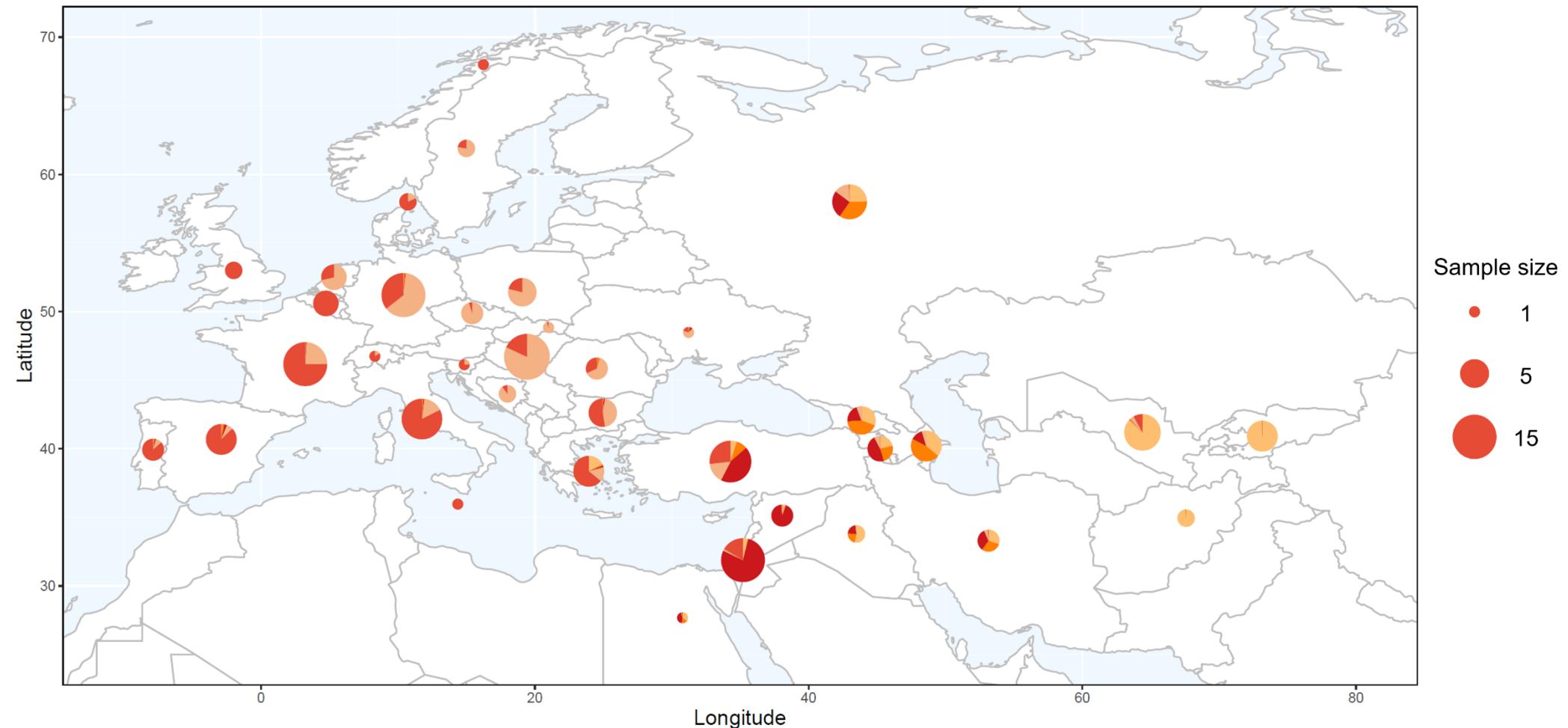
# Case I. PCA in scatter plot



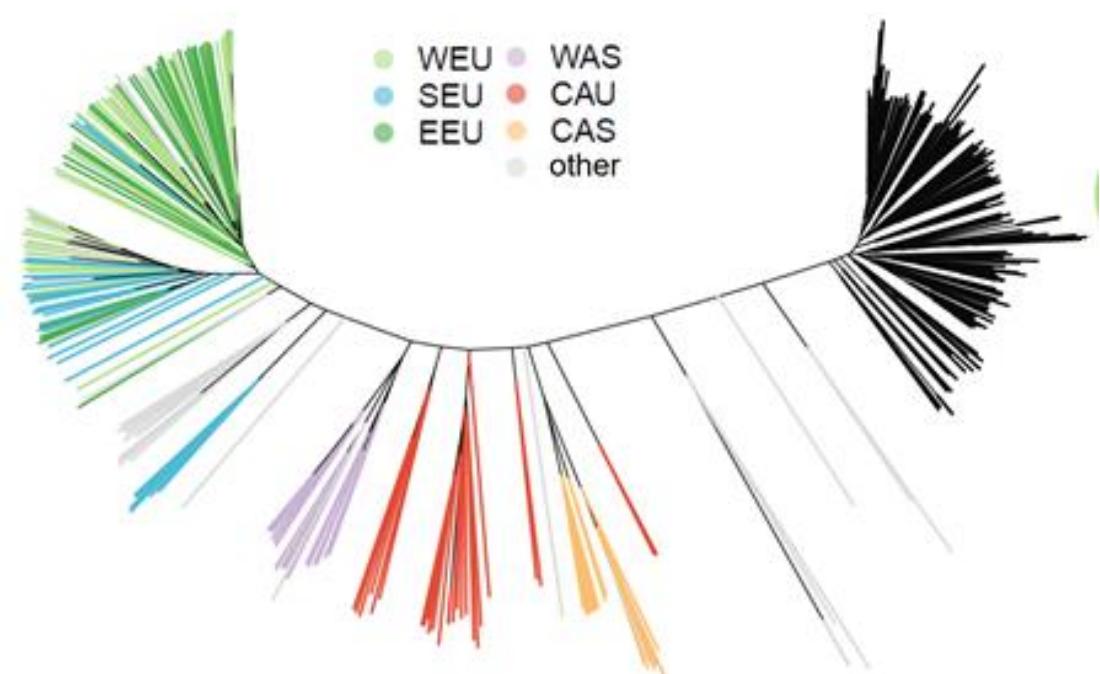
# Case I. Structure in barplot



# Case I. Genetic composition in piechart



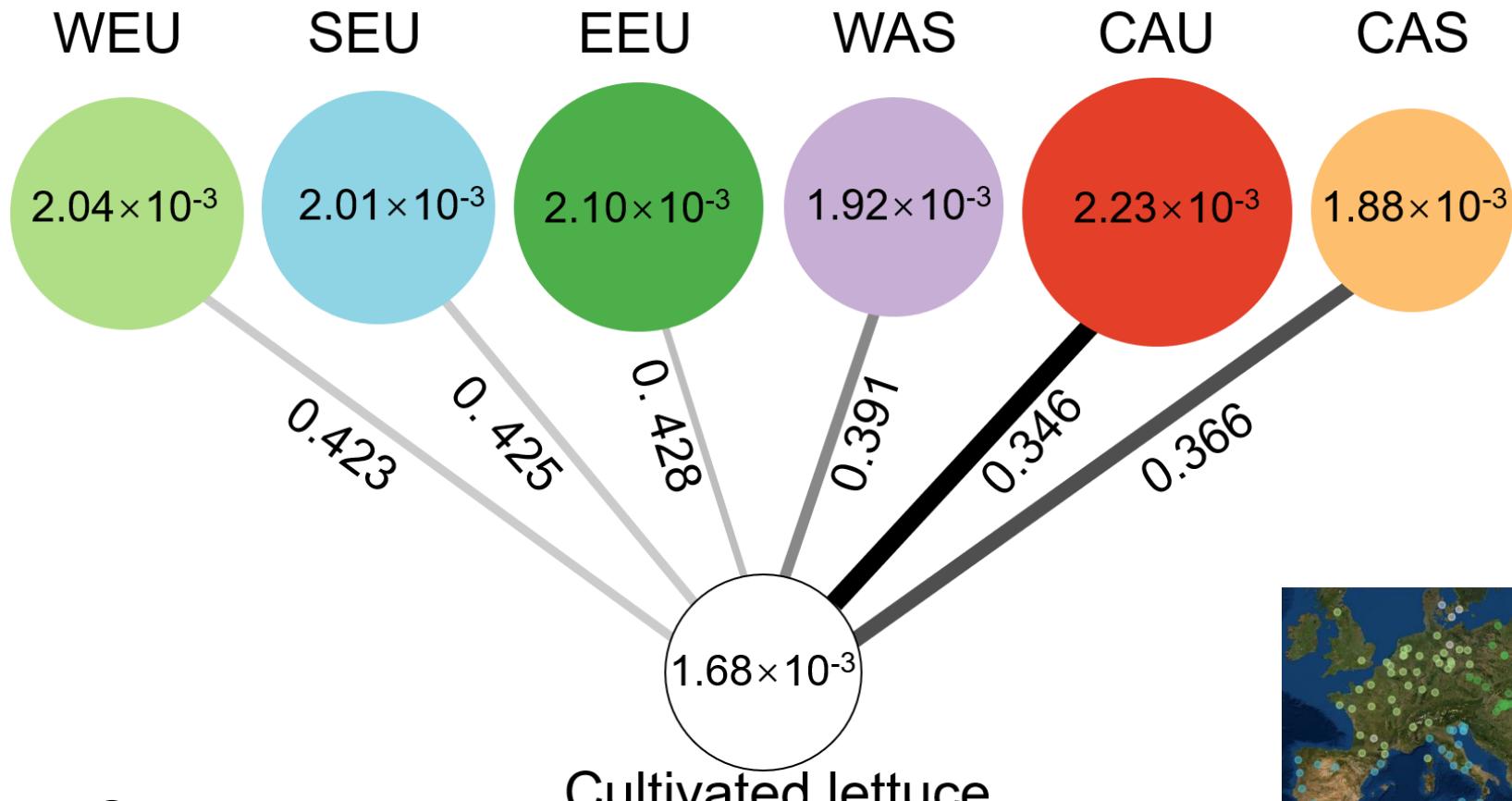
# Case I. Genetic & geographic grouping



3 European groups  
WEU: western Europe  
SEU: southern Europe  
EEU: eastern Europe

3 Asian groups  
WAS: western Asian  
CAU: Caucasian  
CAS: central Asian

# Case I. Population statistic in infographics

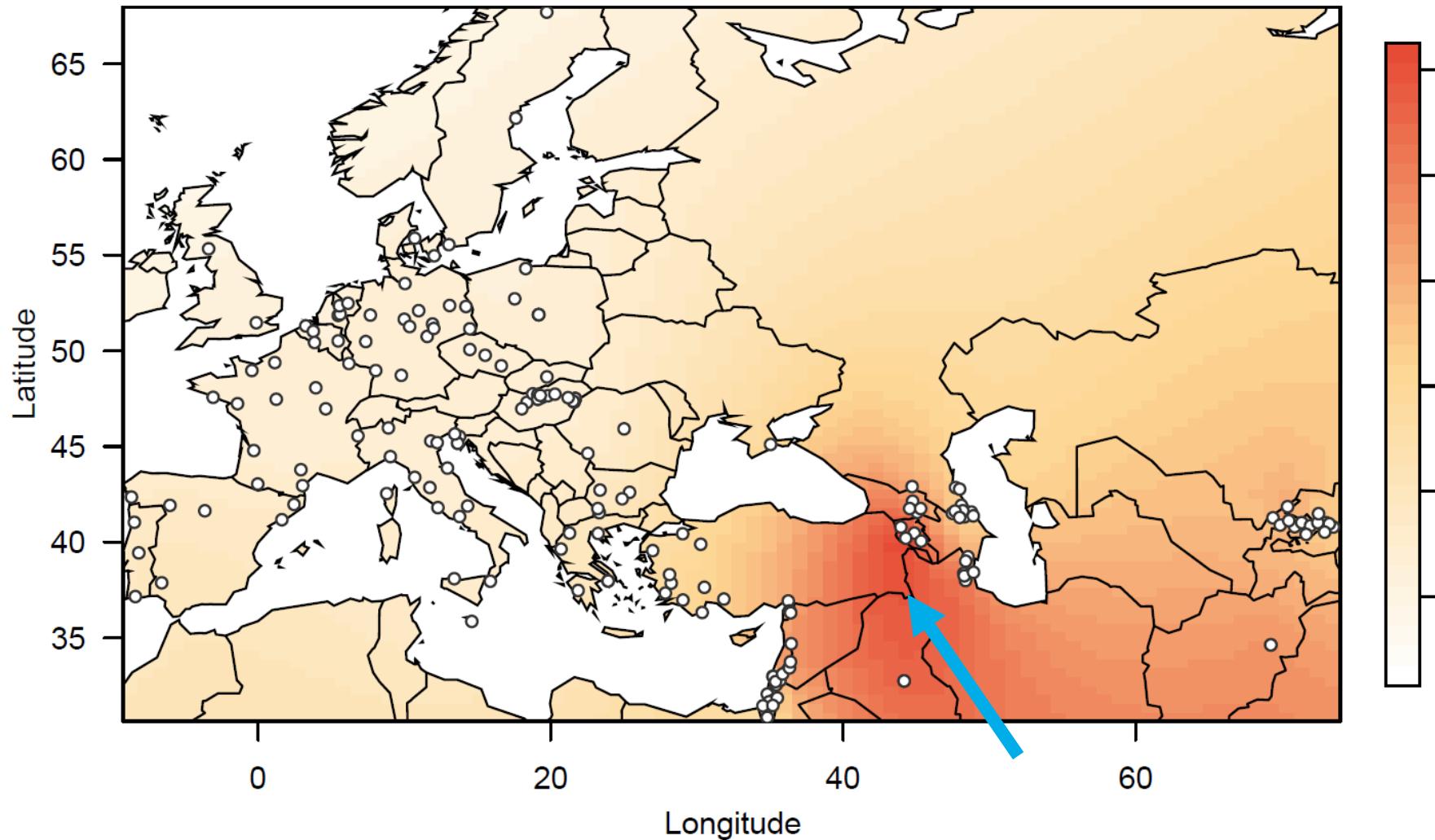


The highest  $\pi$  in CAU, and  
the smallest  $F_{ST}$  between CAU & lettuce



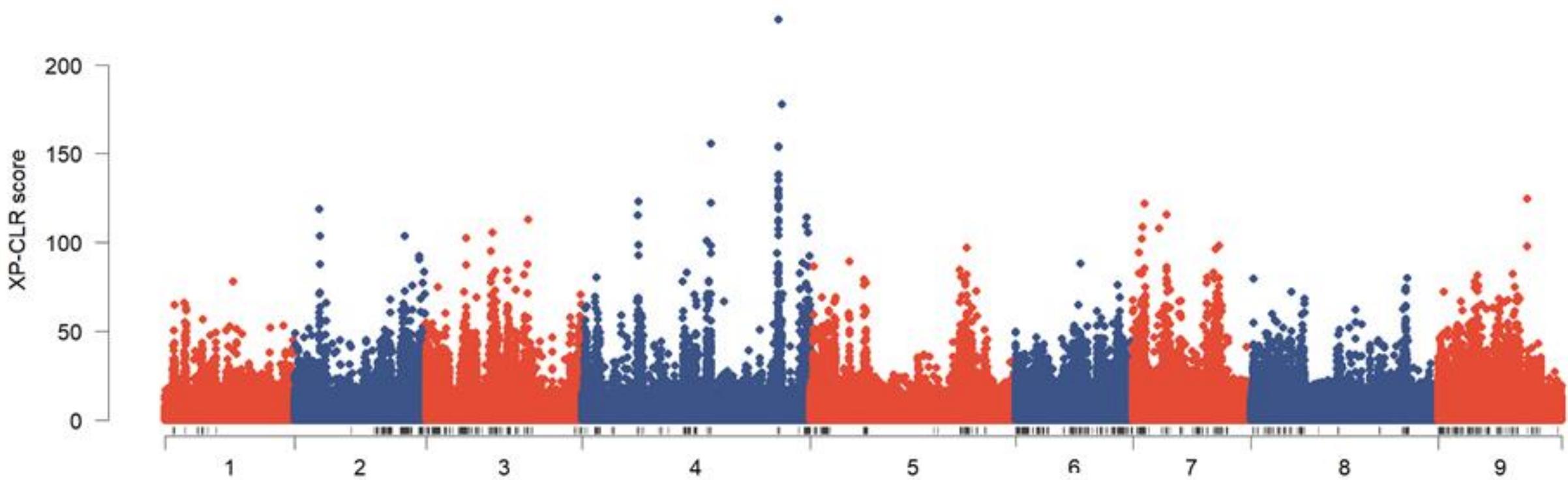
# Case I. Singleton heatmap on map

The highest singletons in CAU



# Case I. Selective sweeps in Manhattan plot

107.7 Mb and 2,304 genes within 4,089 selective sweeps



# Case I. GWAS in Manhattan plot

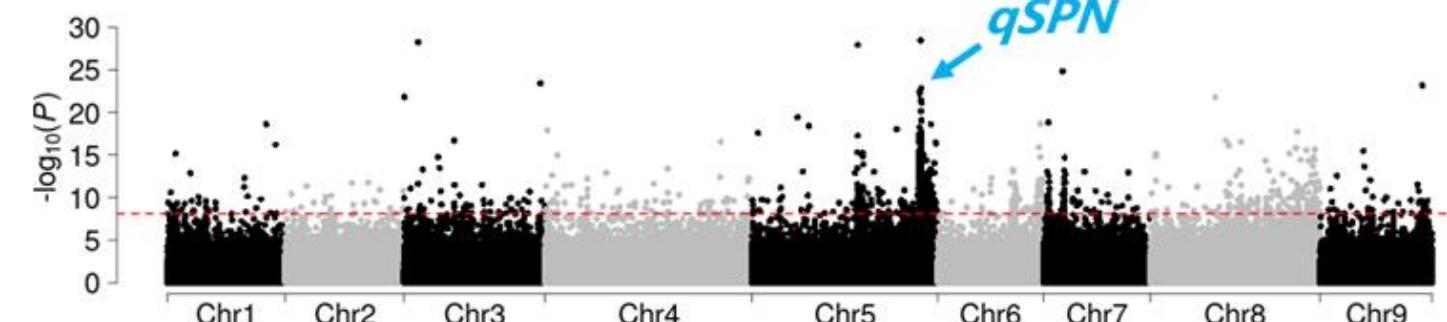
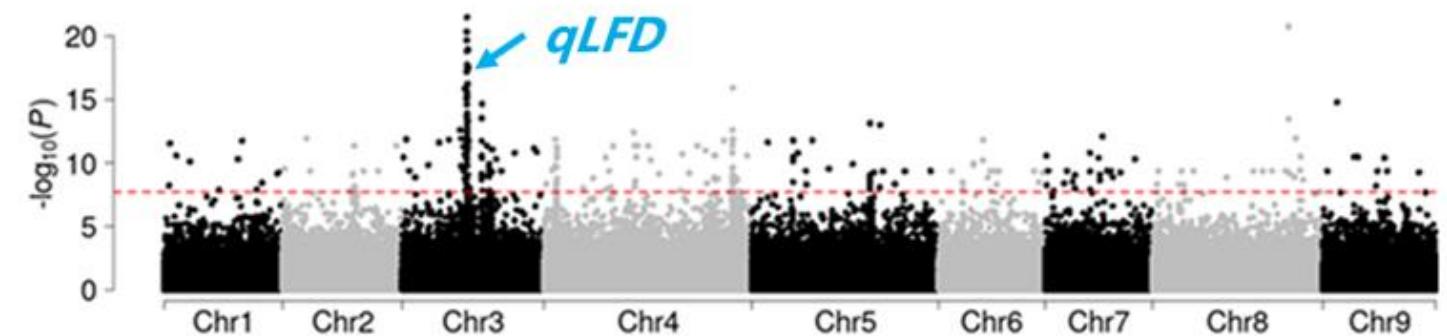
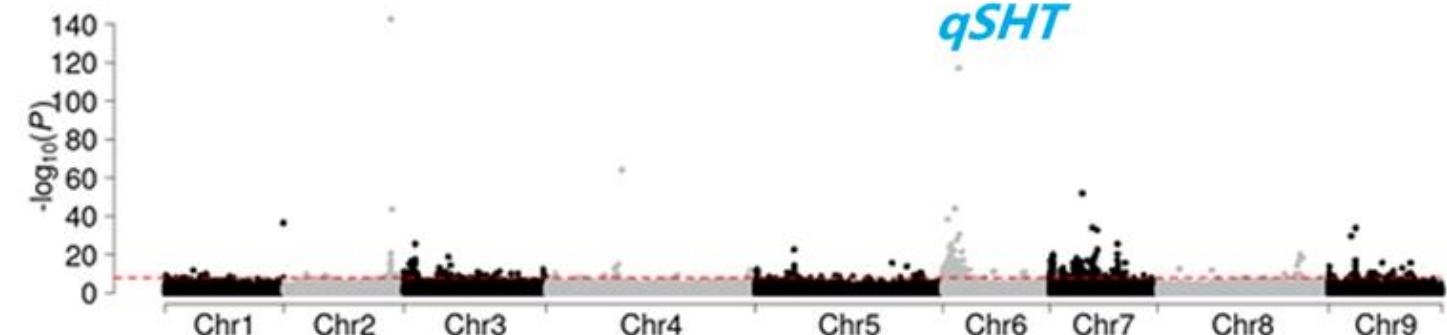
Seed shattering



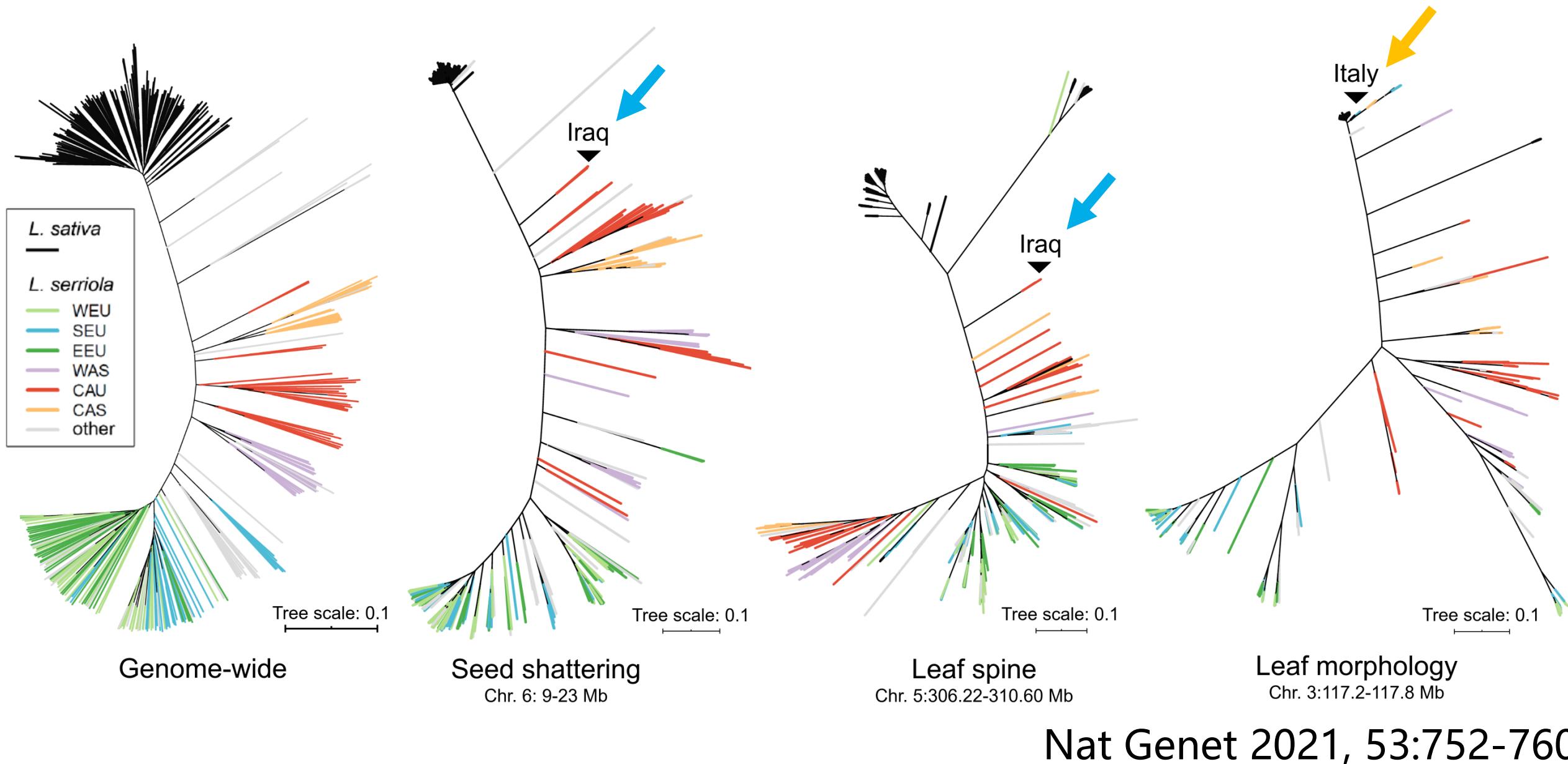
Leaf morphology



Leaf spine

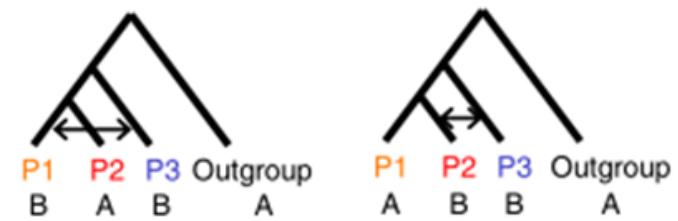
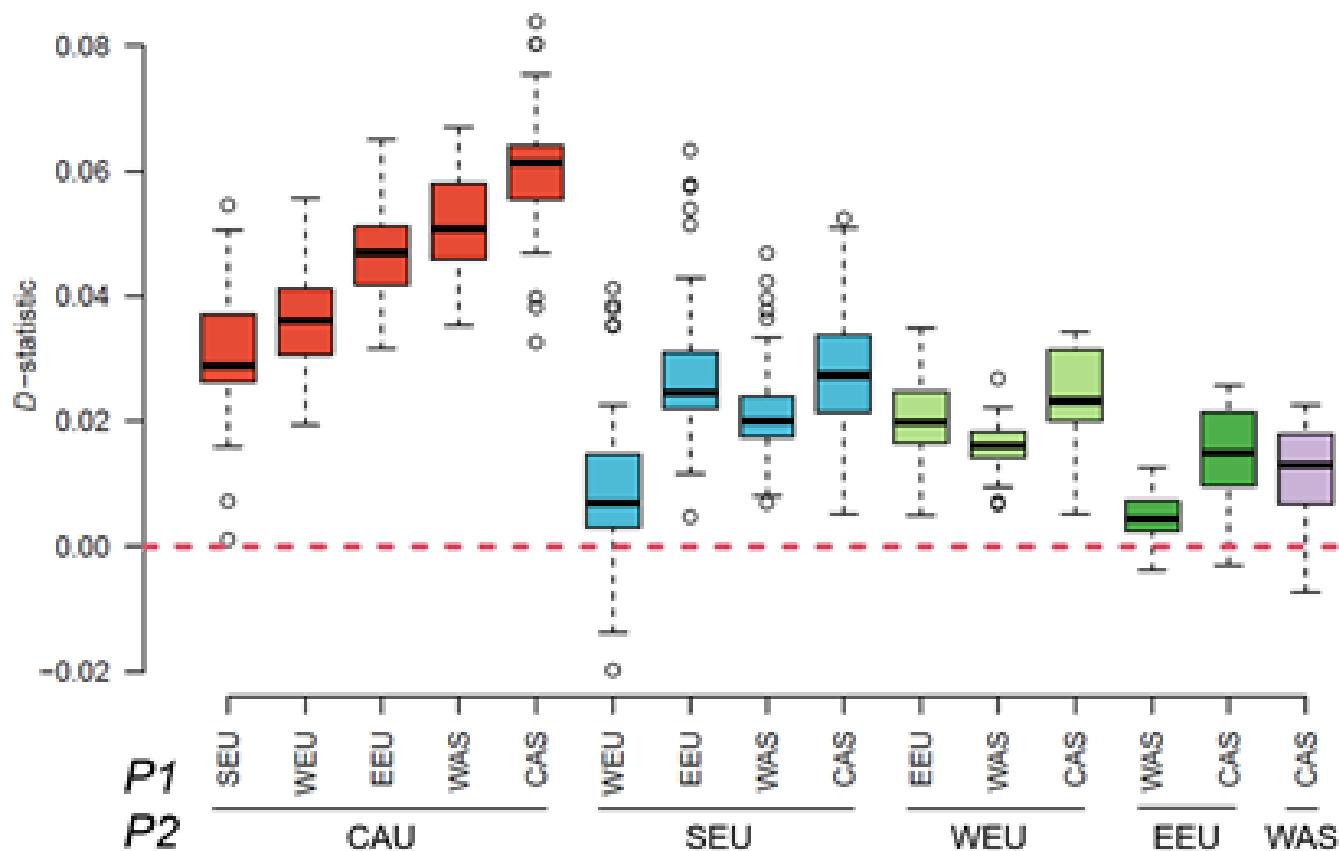
*L. sativa**L. serriola*

# Case I. Trait origins deduced in tree plot



# Case I. ABBA replicates in boxplot

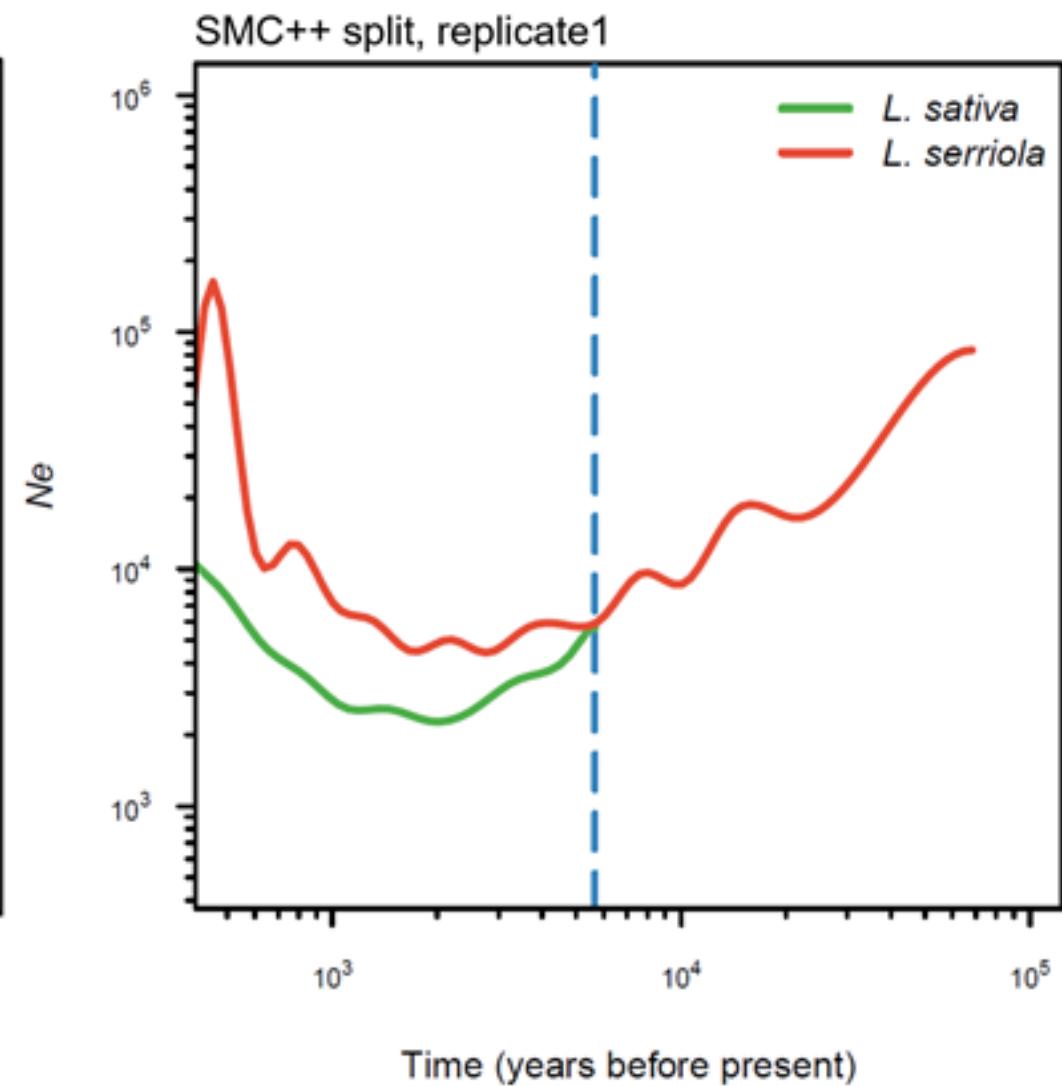
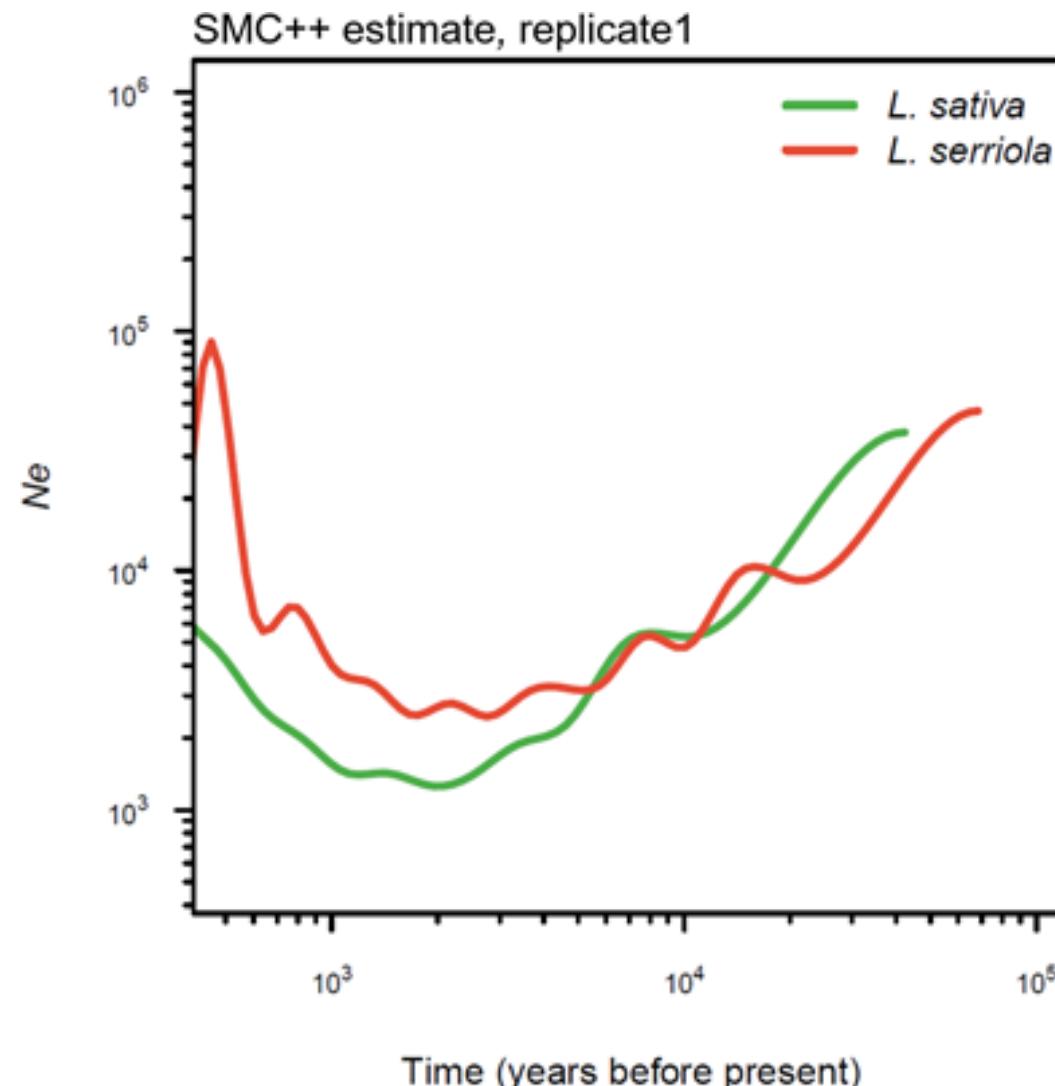
A close relationship in CAU and SEU



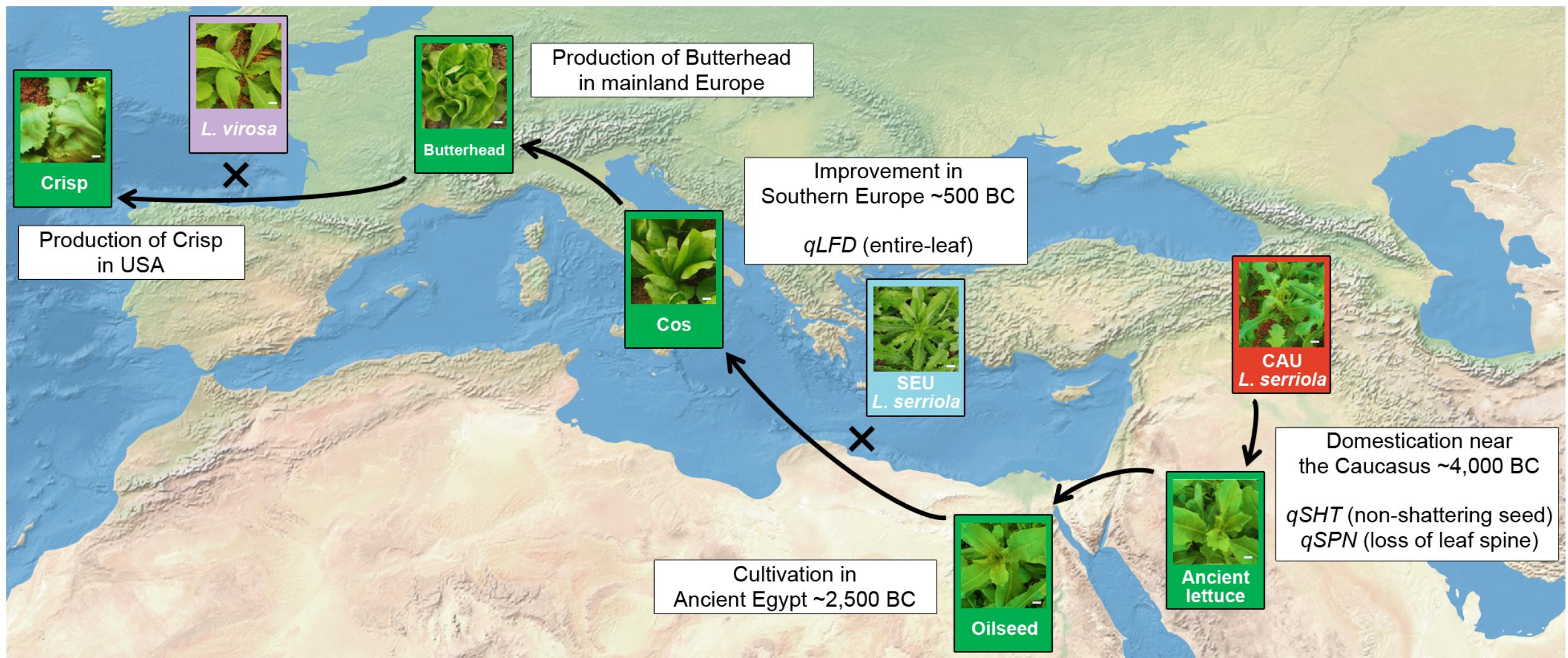
Outgroup: *L. saligna*  
P3: *L. sativa*  
P1 & P2: *L. serriola* groups



# Case I. Effective population size in line chart



# Case I. Domestication history in infographics

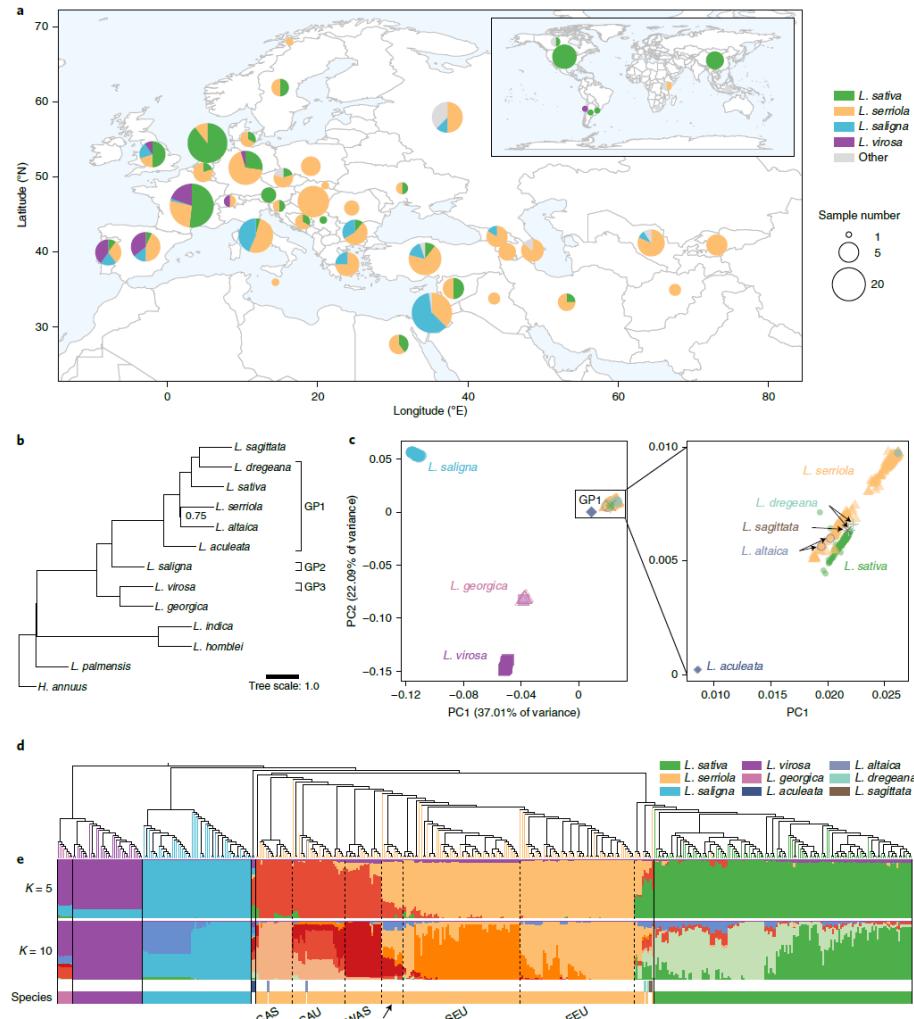


# Case I. Answers to scientific questions

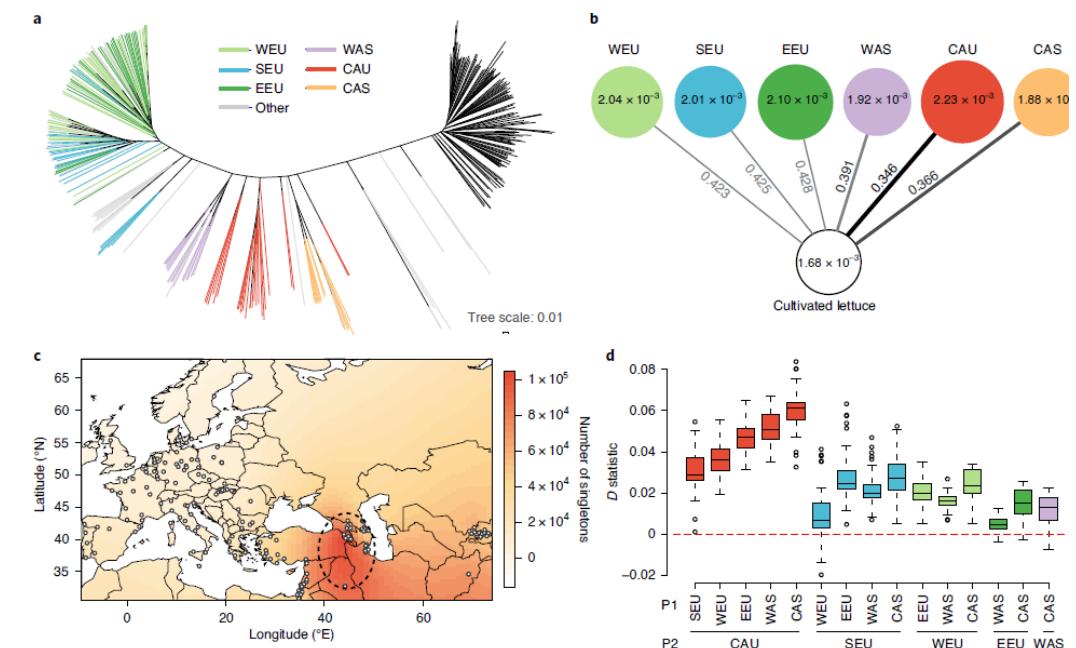
- Q: **Where and when** was lettuce was domesticated?
- A: Lettuce was domesticated **near Caucasus approximately 6,000 years ago.**
- Q: What are the **major events** and traits during lettuce domestication and improvement?
- A: **Non-shattering seeds** marked lettuce domestication; the **entire-leaf trait was introduced later** from a Southern European wild population.

# Case I. Main figures about population

**Fig. 1 |** Phylogeny and population structure of cultivated lettuce and its wild relative species.

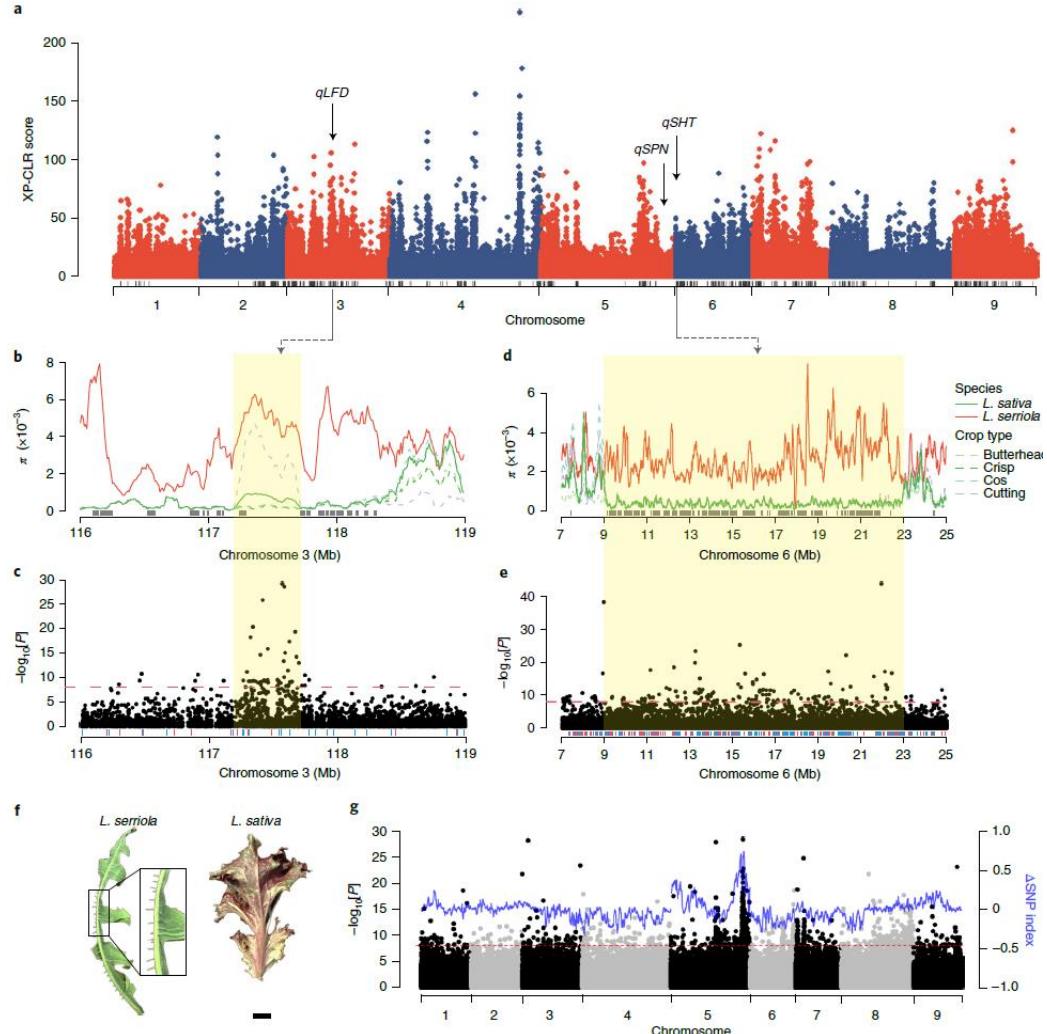


**Fig. 2 |** Proposed domestication center of cultivated lettuce near the Caucasus.

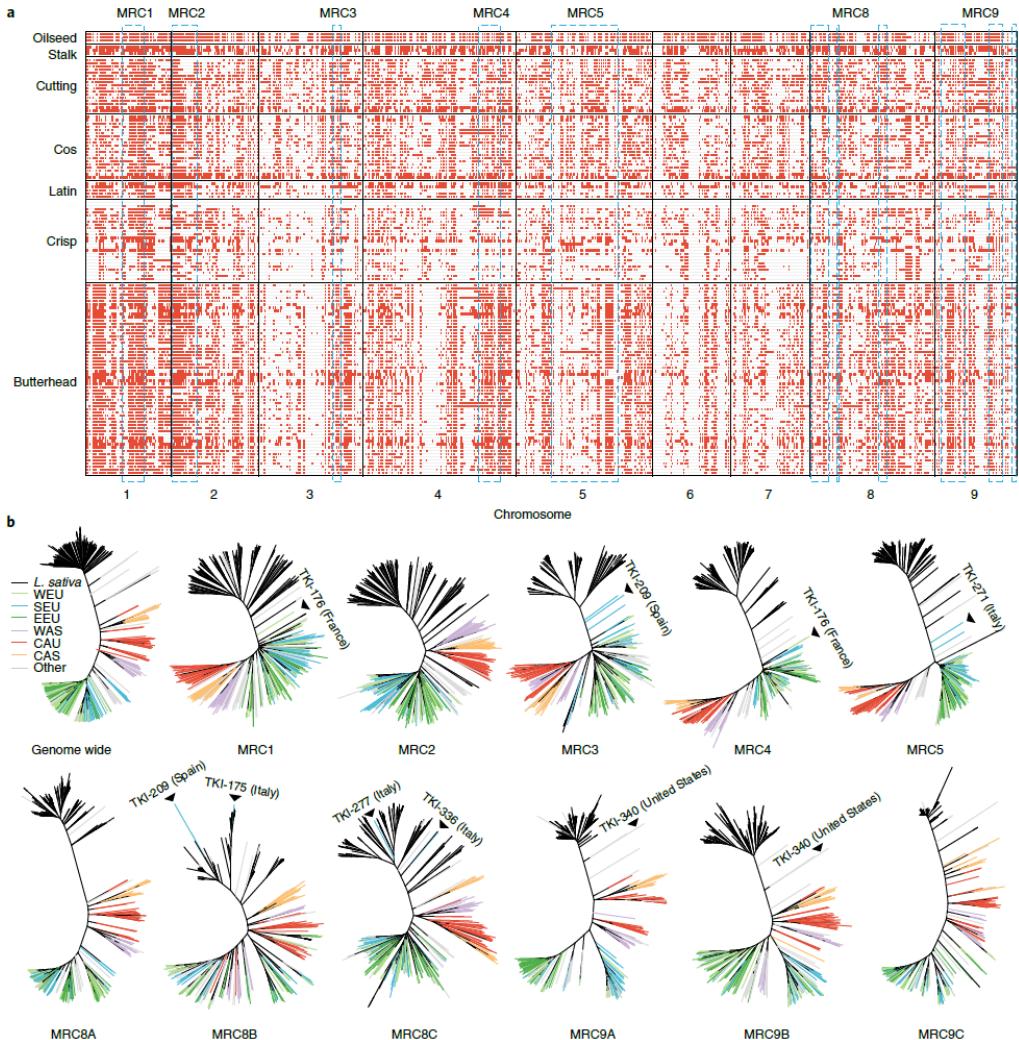


# Case I. Main figures about function

**Fig. 3 | Identification of selective sweeps associated with domestication traits in cultivated lettuce.**

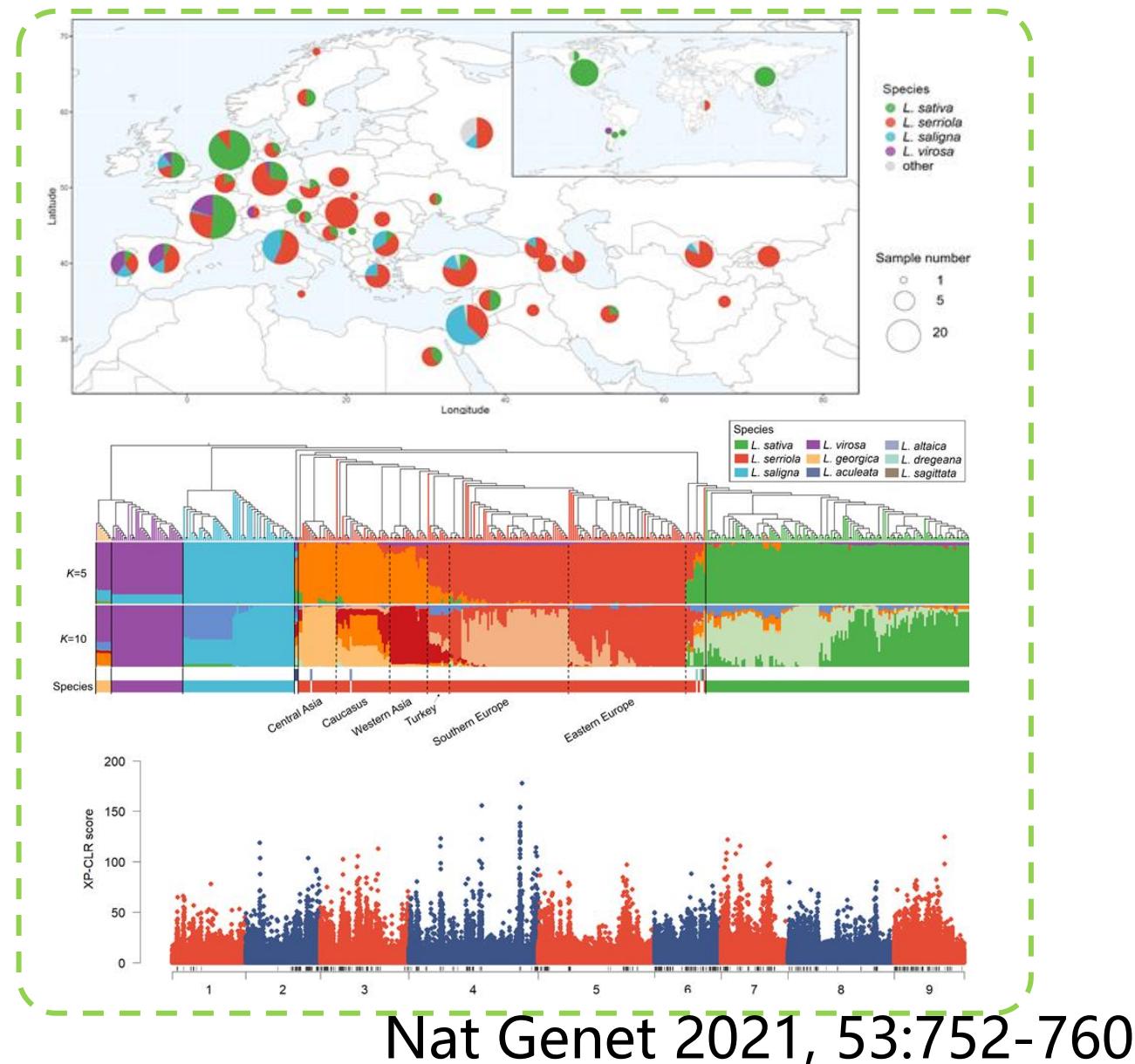


## **Fig. 4 | Introgressive contribution of *L. serriola* to lettuce resistance breeding.**

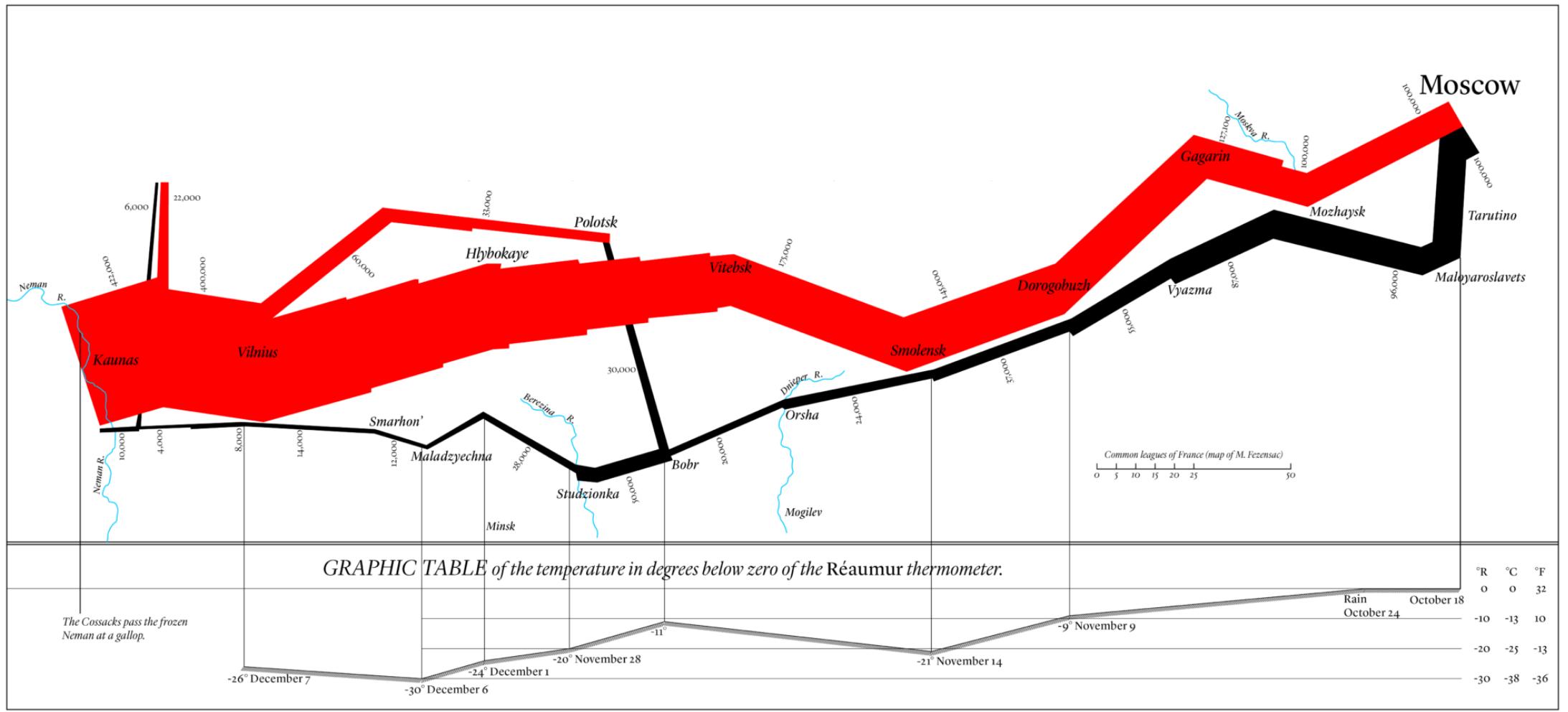


# Case I summary

- Population genomics of 445 cultivated and wild lettuce lines.
- Generated a comprehensive variation map.
- Curated GP status.
- Identified selective sweeps and major genetic determinants of domestication traits.
- Revealed lettuce domestication history.



# May the infographics work for you



Charles Minard' s infographic of Napoleon' s 1812 march on Moscow

# Questions?