

# Scientific Reading and Presentation

## I. Introduction

魏桐

9/11/2021

- About this course, *Scientific Reading and Presentation*
  - Course goals, learning objectives, assessments & scoring
- Scientific reading
  - Why reading is important?
  - How to read a scientific paper?
- Scientific presentation, i.e. slide show
  - What is scientific presentation?
  - How to prepare and give a presentation?
  - A presentation of my recent work

# Section I

# About this course

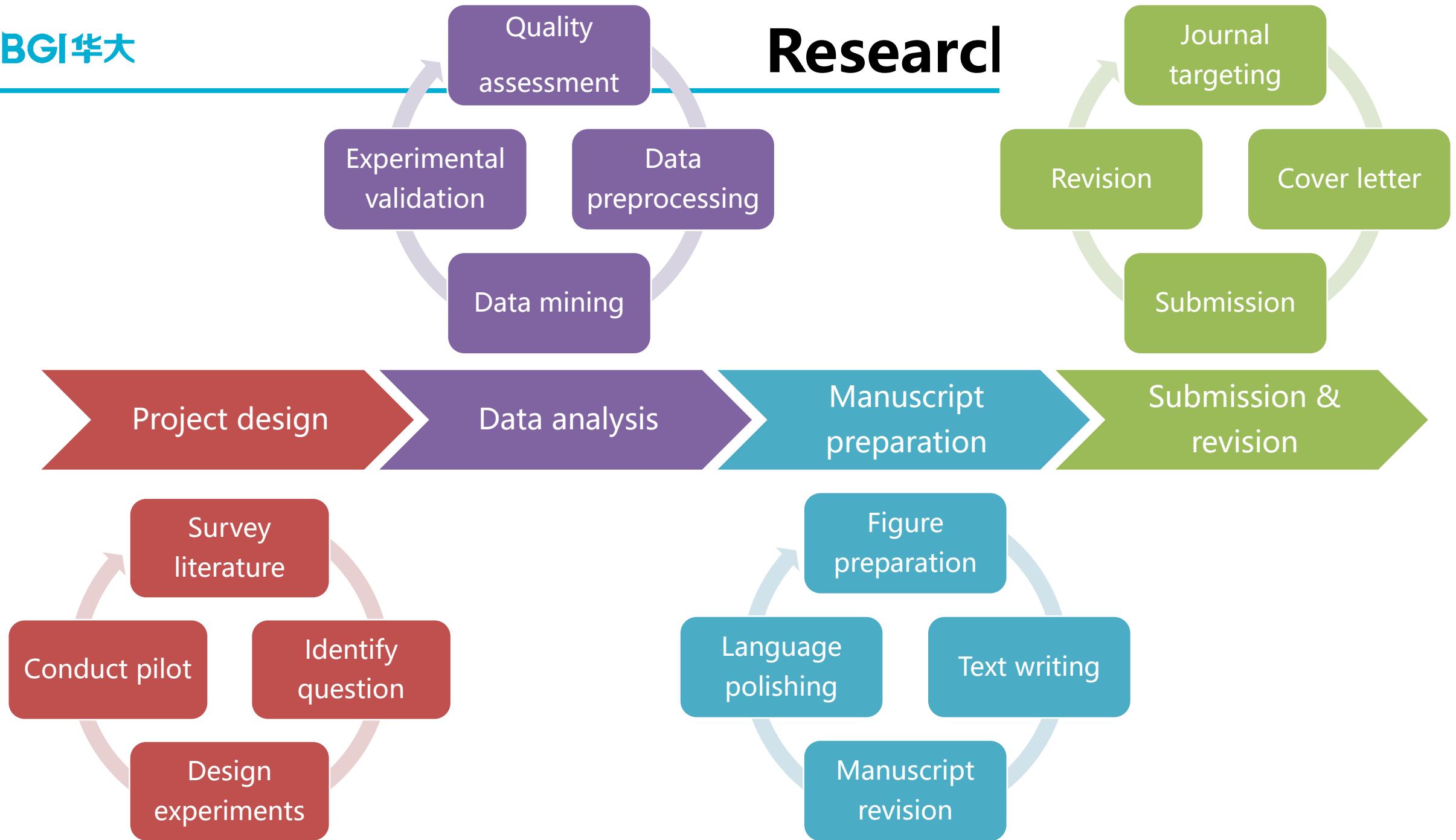
# Scientific method



An empirical method of **acquiring knowledge** from observations.

It **applies to all** scientific enterprises.

# Research



# Abilities required for life science research

---

- The ability to **read and write scientific papers**
- The ability to **think independently**, and solve problems
- The ability to **employ a variety of techniques**
- The ability to **engage in dialogue** using appropriate scientific language
  
- Others skills like communication, interpersonal, practical, self-management and so on

- The goal is to learn,
  - How to **read scientific literature**
  - How to **present scientific reports/papers**
  - How to **engage into scientific discussions**
- The objectives include,
  - The **methodology** of scientific reading
  - A **habit** of periodic literature reading and reviewing

# About this course (continued)

---

- Organization
  - A **45-minute talk** by the tutor or a guest
  - Three **35-45-minute presentation** by student groups + 10-15 Q&A
  - Personal presentation in the final week
- Learn by practice, including
  - Reading **50+ papers**
  - Giving **3 talks**
  - Engaging into the talks

- Introduction of scientific reading and presentation
- The history of sequencing technology
- Literature about human genomics
- Literature about evolutionary genomics
- Literature about population genomics
- Literature about functional genomics
- Literature about recent omics approaches
- Final presentation

Focus on  
human  
research

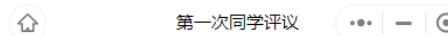
# Contents: example

---

- DNA sequencing with **chain-terminating inhibitors**. Proceedings of the National Academy of Sciences, USA (1977) 74: 5463-7.
- Initial sequencing and analysis of **the human genome**. Nature (2001) 409, 860-921.
- **Ancient human genome** sequence of an extinct Palaeo-Eskimo. Nature (2010) 463, 757-762.
- A map of human genome variation from **population-scale sequencing**. Nature (2010) 467, 1061–1073
- Identification and analysis of **function elements** in 1% of the human genome by the ENCODE pilot project. Nature (2007) 447, 799-816.
- Construction of a human cell landscape at **single-cell level**. Nature (2020) 581, 303-309.

- Weekly seminar (25% each)
  - ~40-minute talk + 10-minute question
  - Tutor + peer evaluation (on-line survey)
  - Style/format (5%) + presenting (5%) + background (5%)  
+ results (10%)
- Discussion (25%)
  - Q & A, reporting in 2 days (1-2% each)
- Final presentation (25%)
  - 8-minute talk + 1-2 questions

## Peer evaluation – 问卷星



### 第一次同学评议

\*1. 姓名

\*2. ppt制作: 中英文是否统一、图表是否变形、字体是否清晰、设计是否美观 (满分5分)

1	2	3	4	5
<input type="radio"/> ①	<input type="radio"/> ②	<input type="radio"/> ③	<input type="radio"/> ④	<input type="radio"/> ⑤

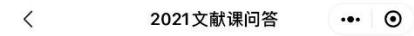
\*3. 展示技巧: 表达是否流畅、与观众是否互动、回答问题是否切题 (满分5分)

1	2	3	4	5
<input type="radio"/> ①	<input type="radio"/> ②	<input type="radio"/> ③	<input type="radio"/> ④	<input type="radio"/> ⑤

\*4. 背景介绍: 背景是否介绍清楚、科学问题是是否清晰、科学意义是否点明 (满分5分)

1	2	3	4	5
<input type="radio"/> ①	<input type="radio"/> ②	<input type="radio"/> ③	<input type="radio"/> ④	<input type="radio"/> ⑤

## Q & A – 腾讯文档



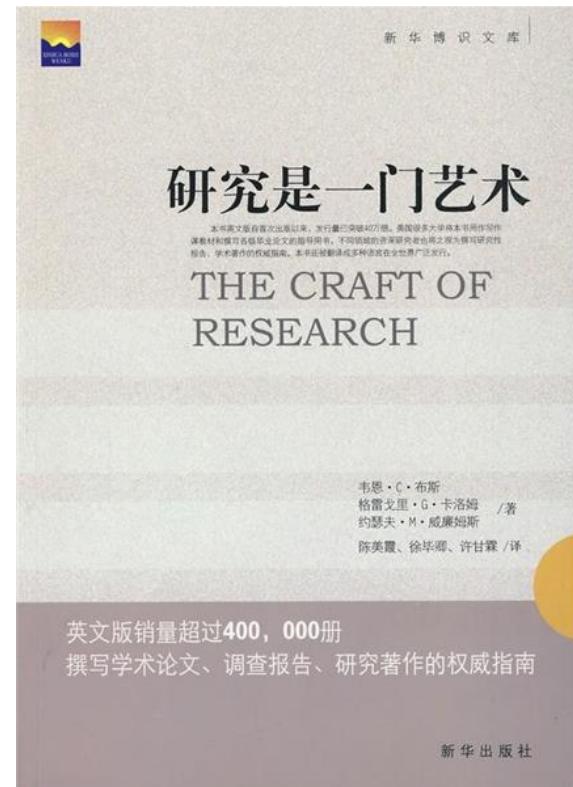
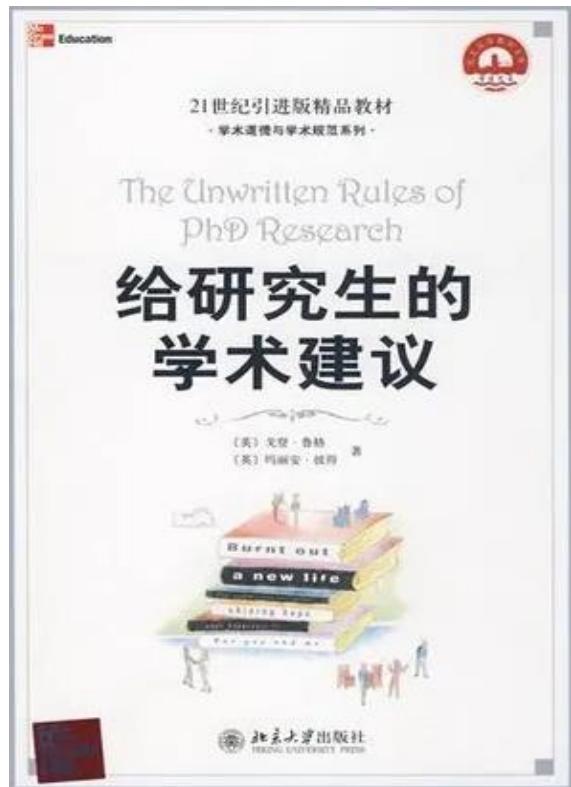
	A	B	C	D	E
1	姓名	提问对象	问题	回答	你的思考
2					
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					
14					
15					
16					
17					
18					
19					
20					
21					
22					
23					
24					
25					



# **Grading (continued)**

# Learning materials: books

- The Craft of Research
- The Unwritten Rules of PhD Research



- About the course
  - <https://github.com/popgenome/Reading2021/blob/main/README.md>
- Q & A, comment under the following link
  - <https://github.com/popgenome/Reading2021/blob/main/Q&A.md>

# Learning materials: papers

---

<https://pan.genomics.cn/ucdisk/s/rYBVVf>



# Learning materials: websites

---

- Nature milestones on genomics sequencing:  
<https://www.nature.com/immersive/d42859-020-00099-0/index.html>
- Human genome project, <https://www.genome.gov/human-genome-project>
- 1000 Genomes, <https://www.internationalgenome.org/>
- ENCODE, <https://www.encodeproject.org/>
- Epigenome Roadmap,  
<http://www.roadmapepigenomics.org/>
- The Cancer Genome Atlas, <https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga>
- ...

# Skills trained in this course

---

- Intellectual skills
  - **Read and summarize literature** critically
  - **Identify and frame key problems**
  - Recognize moral and ethical issues
- Practical skills
  - Undertake sufficient work
  - Design, plan, conduct, and **report investigations**
  - Obtain, record, and analyze data
  - Conduct work in a safe and responsible manner

# Skills trained (continued)

---

- Communication/presentation skills
  - Present research in a scientific manner
  - Cite others work appropriately
- Interpersonal/teamwork skills
  - Identify individual and collective goals
  - Debate in a scientific and respectful manner
  - Evaluate individual and team performance
- Self-management skills
  - Think and work independently and analytically
  - Develop resilience
  - Develop employable skills
  - Develop lifelong learning skills

# How to thrive in this course

---

- During reading
  - **Read** as many as possible
  - Try the best to **understand**
  - **Extend** reading if possible
- During presentation
  - Follow **the logic**
  - **Cite** appropriately
  - **Engage** in talks
- Attitude MATTERS

# Topics for the first 3 weeks

- Topic 1: Sanger sequencing
- Topic 2: Sequencing-by-synthesis
- Topic 3: Single-molecule sequencing

CourseReading2021>CourseReading2021>reading1technology

来自: 魏桐(Weitong) | 分享时间: 2021/08/31 11:17:19 | 到期时间: 永不过期

[转存到...](#) [批量下载](#)

文件名称	更新时间	文件大小
1 1977 (pnas) A new method for sequencing DNA.pdf	2021/08/31 10:20:24	1.76MB
1 1977 (pnas) DNA sequencing with chain-terminating inhibitors.pdf	2021/08/31 10:20:25	2.05MB
1953 (nature) Molecular structure of nucleic acids a structure for deoxyribose nucleic acid.pdf	2021/08/31 10:20:26	143.42 KB
2 2005 (nature) Genome sequencing in microfabricated high-density picolitre reactors.pdf	2021/08/31 10:20:26	1.11MB
2 2005 (science) Accurate Multiplex Polony Sequencing of an Evolved Bacterial Genome.pdf	2021/08/31 10:20:27	564.24 KB
2012 US20120160687A1 Characterization of individual polymer molecules based on monomer interface interactions.pdf	2021/08/31 10:20:28	2.87MB
2016 (genome biol) The Oxford Nanopore MinION.pdf	2021/08/31 10:20:28	920.72 KB
2021 (nat methods) Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm.pdf	2021/08/31 10:20:29	3.55MB
3 2003 (science) Zero-Mode Waveguides for Single-Molecule Analysis at High Concentrations.pdf	2021/08/31 10:20:31	672.44 KB
3 2016 (nat biotechnol) Three decades of nanopore sequencing.pdf	2021/08/31 10:20:34	7.48MB

已全部加载,共10个

- Topic 4: Human Genome Project
- Topic 5: HapMap
- Topic 6: Telomere-to-telomere genome

CourseReading2021>CourseReading2021>reading2genome

来自: 魏桐(Weitong) | 分享时间: 2021/08/31 11:17:19 | 到期时间: 永不过期

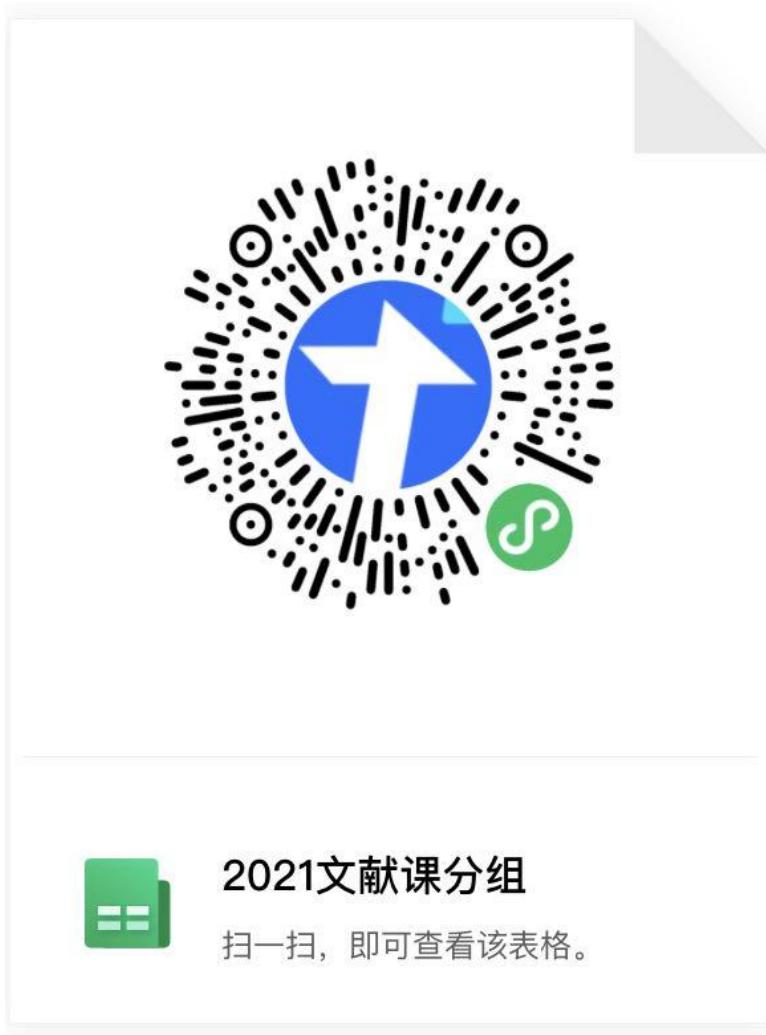
[转存到...](#) [批量下载](#)

文件名称	更新时间	文件大小
1 2001 (nature) Initial sequencing and analysis of the human genome.pdf	2021/08/31 10:20:35	2.32MB
1 2001 (science) The Sequence of the Human Genome.pdf	2021/08/31 10:20:36	3.35MB
1 2004 (nature) Finishing the euchromatic sequence of the human genome.pdf	2021/08/31 10:20:36	1011.70 KB
2 2003 (nature) The international hapmap project.pdf	2021/08/31 10:20:37	554.36 KB
2 2005 (nature) A haplotype map of the human genome.pdf	2021/08/31 10:20:38	1.22MB
2008 (nature) The diploid genome sequence of an Asian individual.pdf	2021/08/31 10:20:38	445.81 KB
2020 (genome biol) Genotyping structural variants in pangenome graphs using the vg toolkit.pdf	2021/08/31 10:20:38	1.69MB
2020 (genome biol) The design and construction of reference pangenome graphs with minigraph.pdf	2021/08/31 10:20:39	1.54MB
2020 (nat biotechnol) Fully phased human genome assembly without parental data using single-cell strand sequencing ...	2021/08/31 10:20:40	4.93MB
2020 Haplotype-resolved de novo assembly with phased assembly graphs.pdf	2021/08/31 10:20:41	2.33MB
3 2020 (nature) Telomere-to-telomere assembly of a complete human X chromosome.pdf	2021/08/31 10:20:43	8.61MB

- Topic 7: Ancient human genome
- Topic 8: Ancient human population
- Topic 9: Adaptation

The screenshot shows a digital library interface with a search bar at the top. Below the search bar is a list of documents, each with a thumbnail, title, update time, and file size. The titles are categorized into three main sections: 1. Ancient human genome sequence, 2. Population genomics and ancestry, and 3. Adaptation. The first two sections are grouped together and highlighted with a blue border.

文件名称	更新时间	文件大小
1 2010 (nature) Ancient human genome sequence of an extinct Palaeo-Eskimo.pdf	2021/08/31 10:48:07	1.48MB
1 2015 (nature) Population genomics of Bronze Age.pdf	2021/08/31 10:48:09	6.79MB
2 2014 (nature) The genomic landscape of Neanderthal ancestry in present-day humans.pdf	2021/08/31 10:48:10	1.86MB
2 2014 (science) Resurrecting surviving Neandertal lineages from modern human genomes.pdf	2021/08/31 10:48:11	1.87MB
2 2020 (nature) Ancient West African foragers in the context of African population history.pdf	2021/08/31 10:48:11	3.42MB
2015 (nature) Genomic evidence for the Pleistocene and recent population history of Native Americans.pdf	2021/08/31 10:48:13	2.96MB
2018 (nature) Terminal Pleistocene Alaskan genome reveals first founding population of Native Americans.pdf	2021/08/31 10:48:14	2.34MB
2019 (nature) The population history of northeastern Siberia since the Pleistocene.pdf	2021/08/31 10:48:18	14.98MB
2019 (science) Ancient Rome A genetic crossroads of Europe and the Mediterranean.pdf	2021/08/31 10:48:19	1.48MB
2020 (nature) Population genomics of the Viking world.pdf	2021/08/31 10:48:25	25.64MB
2020 (science) Ancient DNA indicates human population shifts and admixture in northern and southern China.pdf	2021/08/31 10:48:26	1.07MB
2021 (nature) Genomic insights into the formation of human populations in East Asia.pdf	2021/08/31 10:48:30	19.02MB
2021 (nature) Pleistocene sediment DNA reveals hominin and faunal turnovers at Denisova Cave.pdf	2021/08/31 10:48:32	12.45MB
3 2014 (nature) Altitude adaptation in Tibetans caused by introgression of Denisovan-like DNA.pdf	2021/08/31 10:48:34	4.93MB
3 2020 (nature) Origin of complexity in haemoglobin evolution.pdf	2021/08/31 10:48:35	5.66MB



The image shows two separate screenshots of Microsoft Excel spreadsheets, both titled "个人信息" (Personal Information) and "分组" (Grouping).  
  
The top screenshot displays a table with columns A through F. The first row contains headers: 姓名 (Name), 组别 (Group), 年级 (Grade), and 邮箱 (Email). The second row contains data: 姓名 (Name) and 邮箱 (Email). The third row contains data: 姓名 (Name) and 邮箱 (Email). The fourth row contains data: 姓名 (Name) and 邮箱 (Email). The fifth row contains data: 姓名 (Name) and 邮箱 (Email). The sixth row contains data: 姓名 (Name) and 邮箱 (Email). The seventh row contains data: 姓名 (Name) and 邮箱 (Email). The eighth row contains data: 姓名 (Name) and 邮箱 (Email). The ninth row contains data: 姓名 (Name) and 邮箱 (Email). The tenth row contains data: 姓名 (Name) and 邮箱 (Email).  
  
The bottom screenshot displays a table with columns A through G. The first row contains headers: 组别 (Group), 选题1 (Topic 1), 姓名1 (Name 1), 姓名2 (Name 2), 姓名3 (Name 3), 姓名4 (Name 4), and 姓名5 (Name 5). The second row contains data: 选题1 (Topic 1) and 姓名1 (Name 1). The third row contains data: 选题1 (Topic 1) and 姓名2 (Name 2). The fourth row contains data: 选题1 (Topic 1) and 姓名3 (Name 3). The fifth row contains data: 选题1 (Topic 1) and 姓名4 (Name 4). The sixth row contains data: 选题1 (Topic 1) and 姓名5 (Name 5). The seventh row contains data: 选题1 (Topic 1) and 姓名1 (Name 1). The eighth row contains data: 选题1 (Topic 1) and 姓名2 (Name 2). The ninth row contains data: 选题1 (Topic 1) and 姓名3 (Name 3). The tenth row contains data: 选题1 (Topic 1) and 姓名4 (Name 4).

# Grouping and topic selection



华大云盘

# Section II

# Scientific reading

# What is scientific reading?

---

- Scientific reading is
  - to **read and use** literature with **a full and critical understanding**,
  - while **addressing such questions** as content, context, objectives, and its interpretation and application.

- Know **the structure**
  - Title page, Abstract
  - **Main text in IMRaD** : Introduction, Methods & Materials, Results, and Discussion
  - Other materials: supplementary/supporting materials, data and code, peer review information, etc.
  
- Understand **the logic** in,
  - Introduction
  - Results + Figures/Tables
  - Discussion

# Methodology in reading a paper

---

- **Why** did they do it?
- **How** did they do it?
- **What** did they get?
- **So what** did it mean?

# It applies to Abstract/Conclusion

---

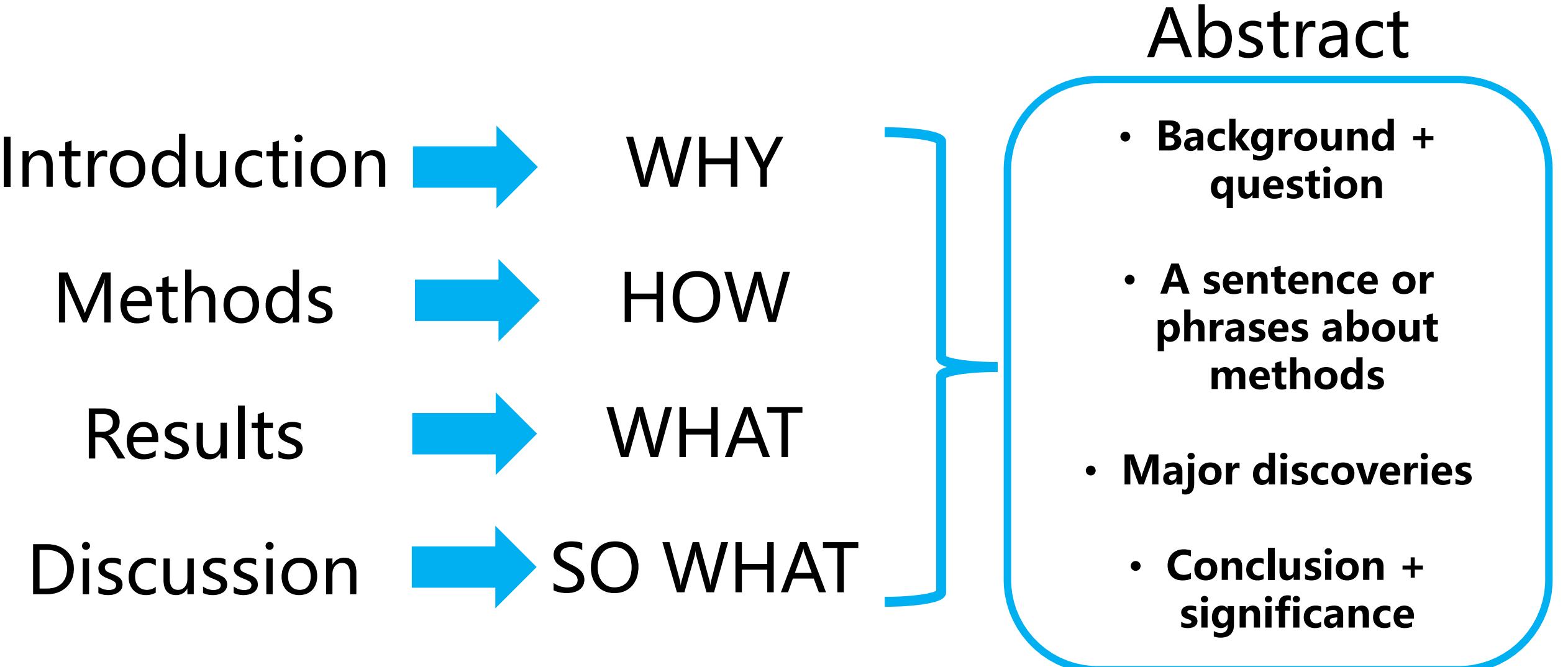
- Why did they do it?
  - ... is important in organ development/crop science; however, the mechanism remains unclear.
- How did they do it?
  - We carried out a multi-omics approach ...
- What did they get?
  - The results showed that genes were associated with ...
- So what did it mean?
  - Our work discovered key players and shed light on ...

# It also applies to Results/Discussion

---

- Why did they do it?
  - To investigate/reveal/study the mechanism of [your research], ..
  - To further explore the population structure, ...
- How did they do it?
  - We conducted a transcriptomic approach ...
  - We analyzed the SNPs from two populations ...
- What did they get?
  - The results showed that genes were differentially expressed ...
  - The population analysis revealed additional structure ..
- So what did it mean?
  - Our data indicate the transcriptomic reprogramming during ...
  - The results suggest the domestication history of ...

## It connects sections



- A **concise and meaningful** title with the emphasis on the main discoveries
- Author list
- Affiliations
- Corresponding author(s)



DONE →

FOUND →

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

Tong Wei<sup>1,11</sup>, Rob van Treuren<sup>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✉</sup> and Huan Liu<sup>1✉</sup>

<sup>1</sup>State Key Laboratory of Agricultural Genomics, BGI-Shenzhen, Shenzhen, China. <sup>2</sup>Centre for Genetic Resources, the Netherlands, Wageningen, the Netherlands. <sup>3</sup>BGI Education Center, University of Chinese Academy of Sciences, Shenzhen, China. <sup>4</sup>Huazhong Agricultural University, Wuhan, China.

<sup>5</sup>Fairy Lake Botanical Garden, Shenzhen, China. <sup>6</sup>University of Copenhagen, Copenhagen, Denmark. <sup>7</sup>BGI-Laos, Vientiane, Laos. <sup>8</sup>China National GeneBank, Shenzhen, China. <sup>9</sup>James D. Watson Institute of Genome Sciences, Hangzhou, China. <sup>10</sup>Guangdong Provincial Key Laboratory of Genome Read and Write, Shenzhen, China. <sup>11</sup>These authors contributed equally: Tong Wei, Rob van Treuren, Xinjiang Liu. ✉e-mail: robbert.vantreuren@wur.nl; liuxin@genomics.cn; liuhuan@genomics.cn

- Structure
  - One sentence for **background**
  - One sentence for **scientific question**
  - Two or three sentences for **major discoveries**
  - One sentence for **conclusion and significance**
- Read Abstract as **a miniature article**
- Use as a **guideline**

# Abstract example

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.

# Abstract example

WHY →  
HOW →

WHAT →

SO WHAT

1. Lettuce (*Lactuca sativa*) is an important vegetable crop worldwide.
2. Cultivated lettuce is believed to be domesticated from *L. serriola*; however, its origins and domestication history remain to be elucidated.
3. 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.
4. 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.
5. We also identified the genetic architecture of other domestication traits and wild introgressions in major resistance clusters in the lettuce genome.
6. This study provides valuable genomic resources for crop breeding and sheds light on the domestication history of cultivated lettuce.

- Structure
  - The 1<sup>st</sup>, 2<sup>nd</sup>, ... paragraphs are background, which describes **the importance** of your area of study, and review **the major findings related** to this work
  - In the end of the 2<sup>nd</sup> last paragraph, raise **the scientific question(s)** in a logical way
  - The last paragraph state the **major discoveries**
- Read and understand the general background
- Pay attention on **how the questions are raised**
- Take a glance at the major findings

# Introduction example: background

Lettuce (*Lactuca sativa* L.) is an important vegetable crop in the Compositae (also known as Asteraceae) family and is widely consumed as salad greens in many countries. Lettuce was first depicted on wall paintings of Egyptian tombs around 2,500 BC<sup>1,2</sup>, making it one of the oldest known vegetable crops. It is believed that cultivated lettuce was domesticated from its progenitor *L. serriola*, and several hypotheses were proposed regarding the domestication center of lettuce, including Egypt, the Mediterranean area, the Middle East and Southwest Asia<sup>1-3</sup>. Modern lettuce varieties are classified based on morphological characteristics into leaf lettuce types (namely cos, butterhead, crisp, Latix and cutting) and non-leaf types (stalk and oilseed)<sup>4</sup>. Oilseed, mostly grown in Egypt for seed oil, is considered the most primitive type, while cos lettuce represents the predecessor of leaf types<sup>2,5</sup>. Despite the morphological variations, different lettuce types share common agronomic traits, such as entire leaf morphology, loss of seed shattering and an absence of spines along the leaf midvein, which are recognized as the domestication syndrome in cultivated lettuce<sup>2</sup>.

# Introduction example: question(s)

Advances in DNA sequencing technology make it feasible to study the genetic architecture in such germplasm collections. A previous RNA sequencing (RNA-seq) study of 240 lettuce accessions demonstrated that different crop types of cultivated lettuce were derived from a single domestication event<sup>9</sup>. However, the domestication history of cultivated lettuce and the genetic basis of human selection remain largely unknown.

# Introduction example: findings

In this study, we sequenced 445 *Lactuca* accessions from 47 countries, comprising the major lettuce crop types and wild relative species. A comprehensive variation map, including 179 million single-nucleotide polymorphisms (SNPs), 30 million insertions/deletions (indels) and 244,866 structural variants, was constructed, from which we analyzed the phylogenetic relationship within the gene pool species and the domestication history of cultivated lettuce. The genetic architecture of domestication traits and introgression regions in resistance clusters were also identified. These sequencing results provide a valuable resource for lettuce research and breeding in the future.

- Organized by pipelines/experiments
- Read useful sections
- Pay attention **in details**, software, parameters, filtering criteria, etc

## ARTICLES

## NATURE GENETICS

## Methods

**Plant materials and sequencing.** The collection of *Lactuca* SSD lines (<http://www.wur.eu/cgnss002>) 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>27</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). Ten seeds were sown for each accession in December 2017 and transplanted in January 2018 in a greenhouse at BGI-Laos. Leaf samples were harvested from representative plants in March for genomic DNA extraction using the cetyl trimethylammonium bromide method<sup>28</sup>. Libraries with an insert size of 250 bp were constructed and paired-end reads (2×100 bp) were produced on a BGISEQ-500 platform at BGI-Shenzhen following the manufacturer's procedures<sup>29</sup>. All of the samples were sequenced with 20× depth except 12 wild species sequenced with higher depth for de novo assembly (Supplementary Table 2).

**Genome assembly of wild *Lactuca* species.** To construct genome assemblies for 12 wild *Lactuca* species, raw reads were filtered using Trimmomatic (version 0.27)<sup>30</sup> with the parameters ILLUMINACLIP:adapter.fa:2:3:5:4:2:trt LEADING:3 TRAILING:3 SLIDINGWINDOW:5:15 MINLEN:50. The genome size for each species was estimated using KmerFreq (version 5.0)<sup>31</sup> with *K*=17, and genome heterozygosity and repeat ratios were assessed using GCE (version 1.0)<sup>32</sup>. The DNA C-values were retrieved from the Plant DNA C-values Database of Kew Gardens (<https://cvalues.science.kew.org/>)<sup>33</sup>. For each species, genome assembly was run using SOAPdenovo (version 2.04) with ten different *K*-mer values, and completeness was assessed by BUSCO (Supplementary Table 3). The assemblies with the highest BUSCO scores were selected for genome annotation with multiple pipelines, as previously reported<sup>34</sup>. Briefly, transposable elements were identified using a combination of homology-based and de novo approaches. RepeatMasker (version open-4.0.6)<sup>35</sup> and RepeatProteinMask (version 1.5.0)<sup>36</sup> were used to identify transposable elements with the known repeats from Repbase (release 25.03)<sup>37</sup> and custom repeat libraries annotated by RepeatModeler (version open-1.0.8)<sup>38</sup>. Tandem Repeats Finder (version 4.07b)<sup>39</sup> was used to find tandem repeats. All of the repeats identified by different approaches were masked before gene prediction. Ab initio prediction was carried out using AUGUSTUS (version 3.2.3)<sup>40</sup> and GeneMark (version 1.0)<sup>41</sup>. RNA-seq data obtained from the National Center for Biotechnology Information (NCBI) Sequence Read Archive database (details in Supplementary Table 3) were assembled into transcripts using Bridger (version r2014-12-01)<sup>42</sup>. The resulting transcripts and expressed sequence tag assemblies downloaded from PlantGDB<sup>43</sup> were aligned against the genome assemblies using BLASTN<sup>44</sup>. Homology-based prediction was performed by BLASTN (*e*-value < 1 × 10<sup>-3</sup>) using five published plant genomes downloaded from the NCBI (<https://www.ncbi.nlm.nih.gov/>), including *L. sativa* (version 8.0)<sup>45</sup>, *H. annuus* (HanXRQ1.0)<sup>46</sup>, *Cynara cardunculus* (CerdV1)<sup>47</sup>, *Arabidopsis thaliana* (TAIR10)<sup>48</sup> and *Solanum lycopersicum* (SL3.0)<sup>49</sup>. All evidences were combined into a final consensus gene set using MAKER (version 2.31.8)<sup>50</sup>.

**Phylogenetic inference of *Lactuca* species.** To construct the phylogenetic relationships of the investigated *Lactuca* species, 4,513 single-copy genes were defined by OrthoMCL (version 5)<sup>51</sup> from the genome of an outgroup *H. annuus*, the lettuce reference genome and the genome assemblies of 11 wild *Lactuca* species, excluding the tetraploid species *L. canadensis*. TranslatorX (<http://translatorx.co.uk/>)<sup>52</sup> was used to translate DNA sequences into amino acids using the standard genetic code and to create amino acid alignments using MAFFT (version 5.0)<sup>53</sup>, which was used for nucleotide alignments after trimming ambiguous alignment portions using Gblocks (version 0.91b)<sup>54</sup>. Individual gene trees were constructed using RAxML (version 7.2.3)<sup>55</sup> with the GTR+GAMMA model, and the species tree was summarized by ASTRAL-III<sup>56</sup>. The Phytools software (version 0.0.1; <https://bitbucket.org/blackrim/phytools>) was used to demonstrate topological concordances and conflicts between individual gene trees and the species tree. Two additional phylogenetic trees were constructed

reads were filtered by Trimomatic using the same parameter as for the genome assembly, and aligned to the *L. sativa* cv. Salinas reference genome (version 8.0)<sup>52</sup> using BWA-MEM with default parameters (version 0.7.12)<sup>50</sup>. Twelve wild accessions with higher sequencing depth were downsampled to about 20×. Five samples from four distantly related species, *L. canadensis*, *L. hombleti*, *L. indica* and *L. palmeensis*, were removed from variant calling because of the low mapping rates on the lettuce reference genome. The alignment bam files were then sorted and PCR duplicates were marked by MarkDuplicates, and HaplotypeCaller was run on each bam file in a genomic variant call format (GVCF) mode. The GVCF files from 440 accessions were consolidated into a single GVCF file, from which SNPs and small indels were identified using a joint calling approach. The SNPs and indels were further filtered using the following criteria: (1) SNPs were filtered with QD < 2.0 | FS > 60.0 | MQ < 40.0 | SOR > 3.0 | MQRankSum < -12.5 | ReadPosRankSum < -8.0, and indels with "QD < 2.0 | FS > 200.0 | SOR > 10.0 | MQRankSum < -12.5 | ReadPosRankSum < -8.0"; (2) genotype calls with a depth < 2 or > 50, (3) variants with more than two alleles and (4) variants with a missing rate of >10% or a minor allele frequency (MAF) of <0.05 were removed, resulting in a set of 13 million filtered SNPs used for population genetic analyses; and (5) linkage disequilibrium pruning was performed with PLINK (version 1.9) using a window size of 10 kb with a step size of one SNP and an *r*<sup>2</sup> threshold of 0.5, resulting a set of 2.77 million pruned SNPs for clustering analysis. The variants from 332 *L. sativa* and *L. serriola* accessions were extracted from the quality-filtered VCF file and filtered with a missing rate of <10% and a MAF of >0.05, resulting in a set of 25.6 million SNPs for demography and GWAS.

Structural variant calling was performed on PCR-duplicate-marked bam files using three programs: Delly (version 0.8.1)<sup>57</sup>, Manta (version 1.5.0)<sup>58</sup> and BreakDancer (1.3.6)<sup>59</sup>. Delly and Manta were run with the default parameters and structural variant calls with imprecise breakpoints were removed (flag IMPRECISE). BreakDancer was performed with the parameters -m 10000000 -q 30 -y 30 -r 2. Structural variants identified in each accession were integrated from different programs and then merged among all 440 accessions using SURVIVOR (version 1.0)<sup>70</sup>. The structural variants identified by at least two programs were kept for further analysis.

All of the variants were annotated using Snpeff (version 4.3g)<sup>71</sup>. SNPs, indels and structural variants were categorized based on their positions on the chromosome (including intergenic regions, exons, introns, splicing sites, untranslated regions and 1-kb upstream and downstream regions) and on their effects (including missense, start codon gain or loss, stop codon gain or loss and splicing mutations).

**Population genetic analysis.** PCA was performed on the filtered SNP set using GCTA (version 1.91.4beta3)<sup>72</sup>. A neighbor-joining tree was constructed with 100 bootstraps using PHYLIP (version 3.69)<sup>73</sup> and the tree layout was generated using the online tool iTOL (<http://itol.embl.de>). The population structure was analyzed with the cluster number *K* ranging from 1–20 by ADMIXTURE (version 1.3.0)<sup>74</sup> using a default fivefold cross-validation (-cv=5). Each *K* was run with 20 replicates and the outputs were aligned by CLUMPP (version 1.1.2)<sup>75</sup>. Considering a cultivated ancestry of 20% (*K*=10), we labeled one *L. sativa*, two *L. dracena* and seven *L. serriola* accessions as admixed samples (Supplementary Table 5).

Genetic differentiation (*F<sub>ST</sub>*) and nucleotide diversity (*π*) were calculated within a nonoverlapping 100-kb window using VCFtools (version 0.1.13)<sup>76</sup>. For a given species, only biallelic SNPs with a missing rate of <0.1 and a MAF of >0.05 were used. Linkage disequilibrium was calculated on SNP pairs within a 500-kb window using PopLDdecay (version 3.31; <https://github.com/BGI-shenzhen/PopLDdecay>). Linkage disequilibrium decay measured the distance at which the Pearson's correlation coefficient (*r*<sup>2</sup>) dropped to half of the maximum. Singleton SNPs in each accession were calculated from the hard-filtering SNP set using VCFtools (version 0.1.13)<sup>76</sup>, and the geographic distribution of singletons was assessed using the kriging regression function Krig implemented in the R package fields<sup>77</sup>.

**Demographic analysis of lettuce evolutionary history.** *D* statistics<sup>78</sup> were calculated within a nonoverlapping 100-kb window to detect asymmetric gene

# Method example

**Read mapping and variant calling.** Variant calling was carried out following the Genome Analysis Toolkit (GATK version 4.0.3.0) Best Practices<sup>64,65</sup>. Raw reads were filtered by Trimmomatic using the same parameter as for the genome assembly, and aligned to the *L. sativa* cv. Salinas reference genome (version 8.0)<sup>12</sup> using BWA-MEM with default parameters (version 0.7.12)<sup>66</sup>. Twelve wild accessions with higher sequencing depth were downsampled to about 20×. Five samples from four distantly related species, *L. canadensis*, *L. homblei*, *L. indica* and *L. palmensis*, were removed from variant calling because of the low mapping rates on the lettuce reference genome. The alignment bam files were then sorted and PCR duplicates were marked by MarkDuplicates, and HaplotypeCaller was run on each bam file in a genomic variant call format (GVCF) mode. The GVCF files from 440 accessions were consolidated into a single GVCF file, from which SNPs and small indels were identified using a joint calling approach.

# Method example: details in pipeline

- Variant calling was carried out following the Genome Analysis ToolKit (GATK version 4.0.3.0) Best Practices <sup>41</sup>. Raw reads were filtered by Trimmomatic using the same parameter as in the genome assembly procedure, and aligned to the *L. sativa* cv. Salinas reference genome (v8.0, downloaded from <http://genomevolution.org/coge>)<sup>14</sup> using BWA mem with default parameters (version 0.7.12)<sup>42</sup>.
- Twelve wild accessions with higher sequencing depth were downsampled to obtain a similar sequencing depth to the rest of the samples. The alignment bam files were then sorted and PCR duplicates were marked, and a GATK tool HaplotypeCaller was run on each bam file in a GVCF (genomic variant call format) mode. The gVCF files from 440 accessions were consolidated into a single gVCF file, from which single nucleotide polymorphisms (SNPs) and small indels were identified using a joint calling approach.

- Organized under subtitles
- Read section by section
- Find and follow **the story line**, follow-up analyses, functional validation, etc.

## NATURE GENETICS

## ARTICLES

into crisp lettuce detected here agrees with the pedigree record of using an *L. vitrosa* line to breed crisp varieties in the United States<sup>16</sup>.

As little archeological evidence has been discovered for lettuce, we estimated the domestication time by assessing the change of effective population sizes ( $N_e$ ) of *L. sativa* and *L. serriola*. Our result showed that both *L. sativa* and *L. serriola* experienced a gradual decline of  $N_e$  starting from 10,000 years ago, but the  $N_e$  contraction in *L. sativa* was stronger than that in *L. serriola* and continued until 2,000 years ago (Extended Data Fig. 5a,b). The  $N_e$  of six *L. serriola* groups and four *L. sativa* crop types displayed a similar pattern of population contraction and expansion to wild and cultivated lettuce, respectively (Supplementary Fig. 5). The divergence time between *L. sativa* and *L. serriola* was estimated to be around 6,000 years ago (Extended Data Fig. 5c,d), denoting the domestication of cultivated lettuce. All of these results indicate that lettuce was first domesticated near the Caucasus, and that the European *L. serriola* population played an important role in crop improvement after lettuce was spread to mainland Europe.

**Human selection of domestication traits in lettuce.** In addition to a reduced genetic diversity caused by population bottlenecks, domestication often results in purifying selection of genomic regions that control agronomic traits favored by humans. In search for the signature of selection, we scanned the lettuce genome using a cross-population composite likelihood ratio test (XP-CLR) and identified 4,089 selective sweeps covering 4.66% (107.7 Mb) of the assembled genome and harboring 2,304 genes in total (Fig. 3a and Supplementary Table 8). Genes involved in fatty acid and carbohydrate metabolism are enriched among those under selection, probably caused by human selection on nutritional values (Supplementary Table 9).

To better understand the impact of human selection on the lettuce genome, we carried out genome-wide association studies (GWAS) of the domestication traits in cultivated lettuce (Supplementary Table 10). As entire and lobed leaf morphology were both recorded in cultivated and wild lettuce, we analyzed the leaf morphology in *L. sativa* and *L. serriola* accessions separately. The major GWAS signals were detected within a 600-kb region on chromosome 3 in *L. sativa* (leading SNP  $P=2.45 \times 10^{-22}$ ), which overlapped with the major signals detected in the same region in *L. serriola* (leading SNP  $P=6.62 \times 10^{-30}$ ; Fig. 3c, Supplementary Fig. 6 and Supplementary Tables 11 and 12). Our results align with two previous studies using quantitative trait locus (QTL) mapping and bulked segregant analysis<sup>7,18</sup>, which linked the same region to lobed leaf morphology. The nucleotide diversity of this interval was markedly reduced in cultivated lettuce and in butterhead, crisp and cos crop types (Fig. 3b). However, the nucleotide diversity in the cutting type was comparable to that in *L. serriola*, consistent with four out of 12 cutting accessions developing lobed leaves. We then carried out a phylogenetic analysis using the SNPs within this region to explore the underlying genetic architecture. The entire leaf lettuce accessions, including three primitive oilseed samples, formed a single clade and clustered with entire leaf *L. serriola* accessions from SEU (Supplementary

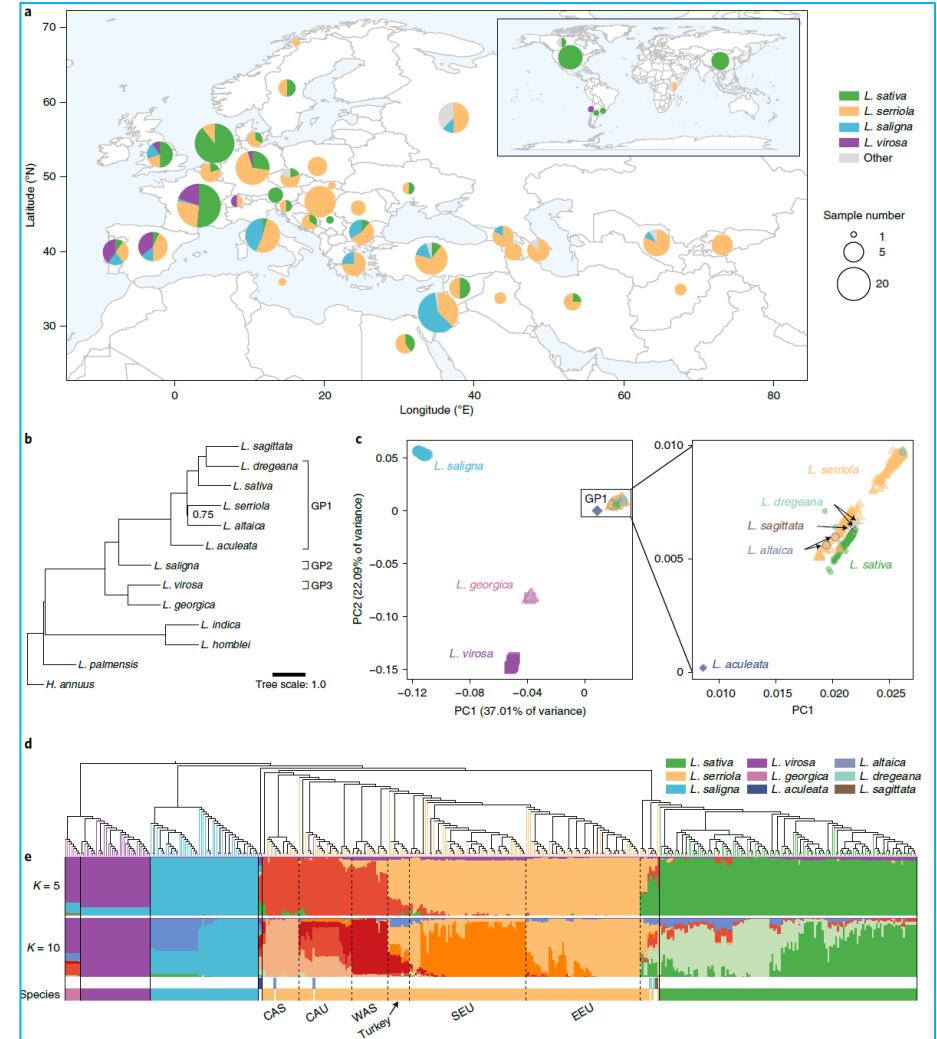
detected several major signals across the lettuce chromosomes (Fig. 3e and Supplementary Fig. 8a) and the signals on chromosome 6 agree with a previously identified QTL, *qSHT*, which explained 85% of the phenotypic variation of seed shattering<sup>17</sup>. The associated region was overlapped with an extended selective region spanning 9–23 Mb, within which the nucleotide diversity was markedly reduced in *L. sativa* compared with *L. serriola* (Fig. 3d). A transcription factor gene, encoding a homolog to NAC SECONDARY WALL THICKENING PROMOTING FACTOR 1 (NST1), which controls pod dehiscence in *Arabidopsis*<sup>22</sup>, was found within this region as an intriguing candidate (Supplementary Fig. 8b,c). The phylogeny within this region showed that all of the cultivated lettuce accessions and nonshattering mixed samples formed a tight clade that is sister to the CAU group of *L. serriola* (Extended Data Fig. 6c). These results suggest that the loss of seed shattering in cultivated lettuce was probably caused by spontaneous mutation(s) derived from a wild ancestor near the Caucasus.

Wild plants develop spinose leaf structures such as spines and thorns as essential physical defense traits against herbivores, but these primitive traits in vegetables, especially leafy ones, are undesirable for human consumption. To identify the associated genomic regions, we carried out an association analysis in the investigated *L. serriola* samples that were recorded with or without spines in leaf midveins. The major GWAS signals were detected within 306.22–310.60 Mb on chromosome 5 (leading SNP with  $P=3.45 \times 10^{-29}$ ; Fig. 3g). Additionally, we sequenced two bulked F2 populations with or without leaf spines, which derived from a TKI-143×GLHZ cross (an *L. serriola* accession with spines crossed with a cutting lettuce accession without spines; Fig. 3f). We calculated the differences in SNP indices ( $\Delta$ SNP index) between the bulked samples and identified a single region overlapped with the GWAS signals on chromosome 5. Our result was also supported by a previously reported QTL, *qSPN*, which explained 82% of the phenotypic variation of spine presence<sup>17</sup>. Among the 70 candidate genes in this region, 21 were differentially regulated between *L. sativa* and *L. serriola*, which can be investigated further by genetic approaches (Supplementary Table 12). The phylogeny in this region suggests an early divergence from the CAU group of *L. serriola* (Extended Data Fig. 6d), although no selection signals were detected within it.

**Identification of loci related to agronomic traits.** Modern lettuce cultivars display diverse characteristics in many agronomic traits, such as flowering time, anthocyanin biosynthesis and leaf development. Among them, prolonged vegetative growth and delayed flowering have been recognized as important agronomic traits in cultivated lettuce. Our GWAS analysis detected a strong signal around 164.5 Mb on chromosome 7 ( $P=3.45 \times 10^{-14}$ ), where a PHYTOCHROME C (*PHYC*) gene resides (Extended Data Fig. 7). Two major haplotypes of *PHYC* were discovered in cultivated lettuce, and the accessions carrying the reference G allele displayed a significant delay in flowering date accompanied by reduced *PHYC* expression (Extended Data Fig. 7d and Supplementary Tables 12 and 14), resembling a wheat *PHYC* mutant that showed delayed

# Result example

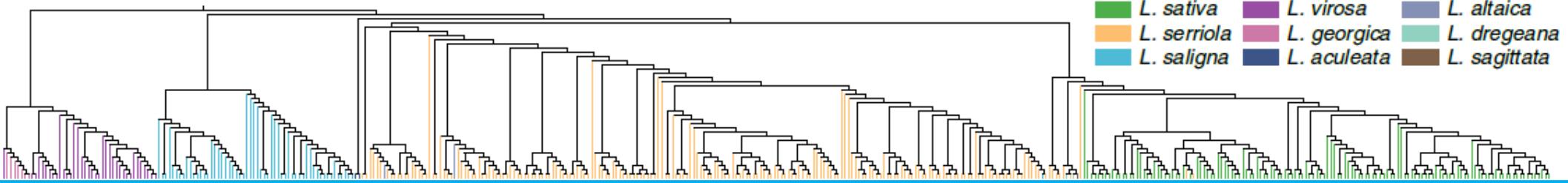
We explored the phylogenetic relationships among the 440 *Lactuca* accessions using 13 million high-quality SNPs. Our result showed that all of the *L. sativa* accessions formed a monophyletic clade, suggesting a single domestication event for cultivated lettuce (Fig. 1d). The phylogenetic positions of other wild species are consistent with the species tree, with most GP1 species having a close relationship with cultivated lettuce except that *L. georgica* was found close to the GP3 species *L. virosa*. Model-based clustering analysis revealed additional inter- and intraspecies relationships. Asia- and Europe-originated accessions formed two groups in both *L. serriola* and *L. saligna*, reflecting the spatial genetic variations within these species (Fig. 1e and Extended Data Fig. 2). Subgroups were further revealed in *L. serriola* when a higher number of  $K$  was assumed. These consisted of accessions from Central Asia, the Caucasus, Western Asia and Southern and Eastern Europe. Admixture between cultivated lettuce and wild species was detected in seven *L. serriola* accessions as well as two *L. dregeana* and one *L. sagittata* samples (Supplementary Table 5). We also observed admixture between Western Asian and Southern European compositions in Turkish accessions (Extended Data Fig. 2 and Supplementary Fig. 3). The phylogenetic relationships were revealed by the principal component analysis (PCA) as well, in which the GP1 species except *L. georgica* grouped together and *L. serriola* accessions were split into two distinct groups with Asian and European origins (Fig. 1c and Extended Data Fig. 3a,b).



# Result example: figure + legend

Panel

d



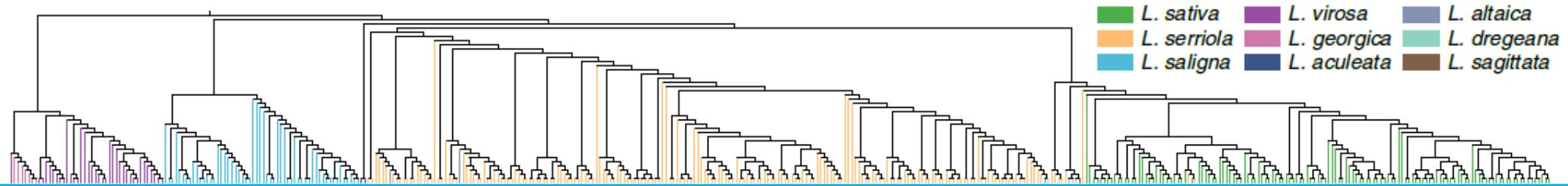
Legend

d, A neighbor-joining tree of 440 *Lactuca* accessions. Branch colors denote species.

# Result example: figure + text

## Panel

d



## Result

1. We explored the phylogenetic relationships among the 440 *Lactuca* accessions using 13 million high-quality SNPs.
2. Our result showed that all of the *L. sativa* accessions formed a monophyletic clade,
3. suggesting a single domestication event for cultivated lettuce

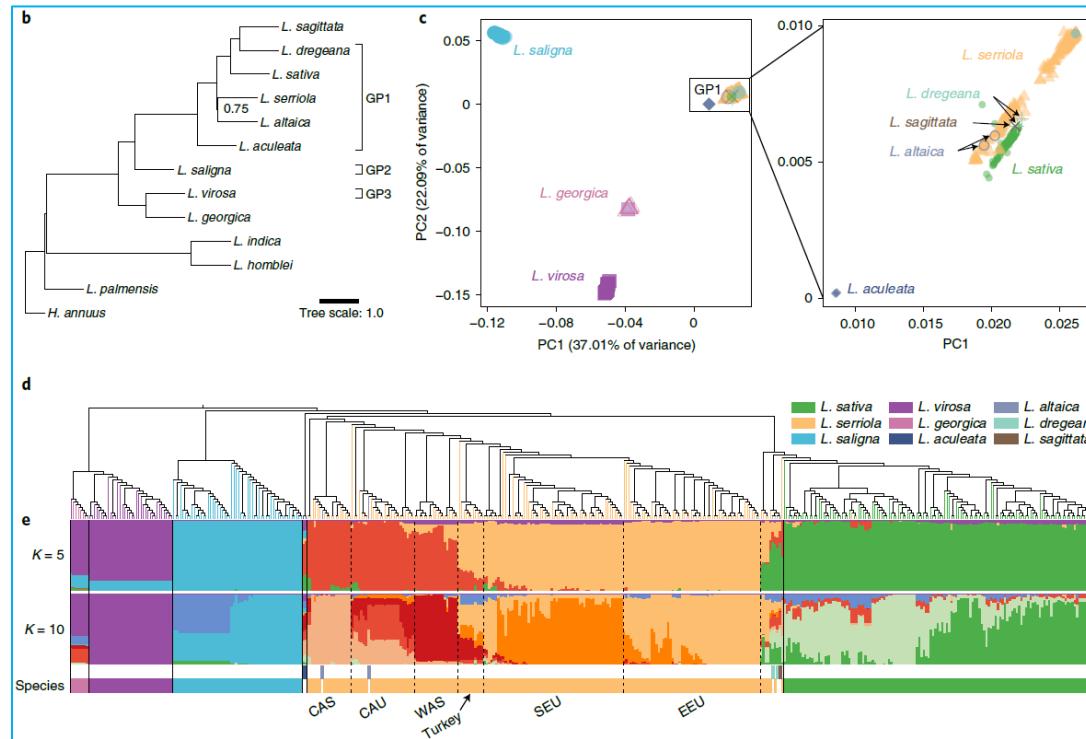
# Result example: logic within paragraph

We explored the phylogenetic relationships among the 440 *Lactuca* accessions using 13 million high-quality SNPs. Our result showed that all of the *L. sativa* accessions formed a monophyletic clade, suggesting a single domestication event for cultivated lettuce (Fig. 1d). The phylogenetic positions of other wild species are consistent with the species tree, with most GP1 species having a close relationship with cultivated lettuce except that *L. georgica* was found close to the GP3 species *L. virosa*. Model-based clustering analysis revealed additional inter- and intraspecies relationships. Asia- and Europe-originated accessions formed two groups in both *L. serriola* and *L. saligna*, reflecting the spatial genetic variations within these species (Fig. 1e and Extended Data Fig. 2). Subgroups were further revealed in *L. serriola* when a higher number of  $K$  was assumed. These consisted of accessions from Central Asia, the Caucasus, Western Asia and Southern and Eastern Europe. Admixture between cultivated lettuce and wild species was detected in seven *L. serriola* accessions as well as two *L. dregeana* and one *L. sagittata* samples (Supplementary Table 5). We also observed admixture between Western Asian and Southern European compositions in Turkish accessions (Extended Data Fig. 2 and Supplementary Fig. 3). The phylogenetic relationships were revealed by the principal component analysis (PCA) as well, in which the GP1 species except *L. georgica* grouped together and *L. serriola* accessions were split into two distinct groups with Asian and European origins (Fig. 1c and Extended Data Fig. 3a,b).

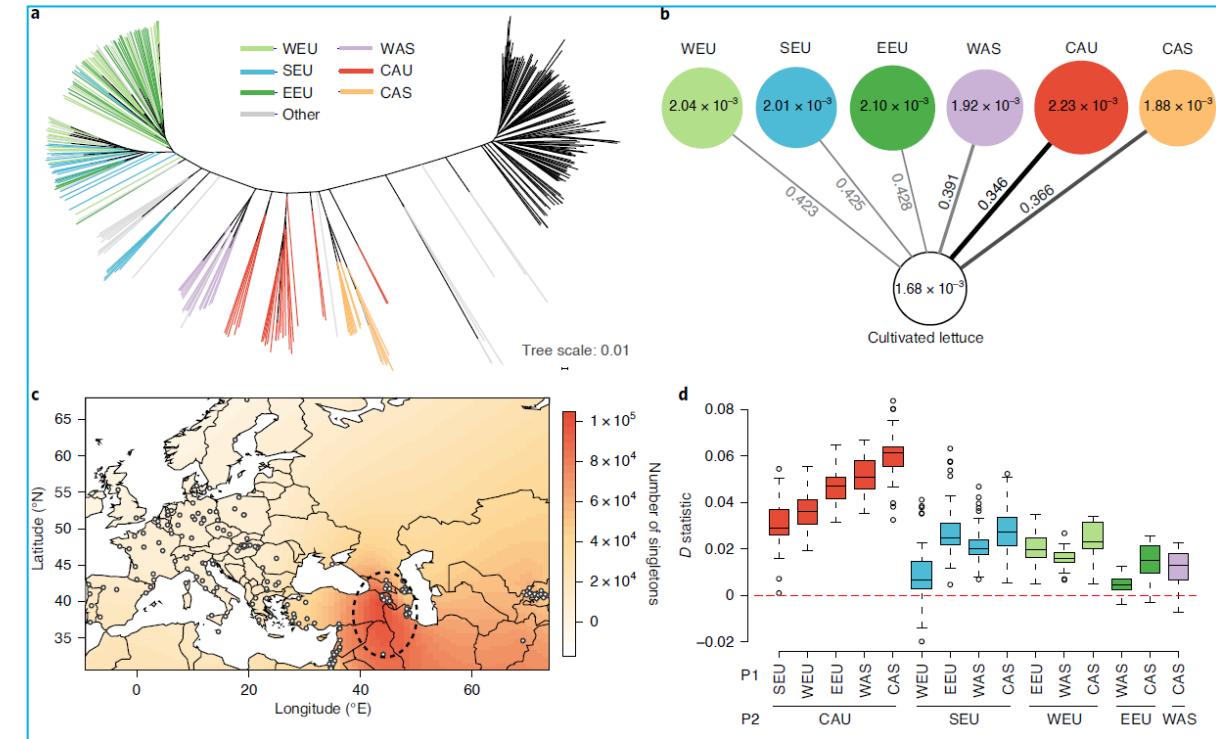
1. Fig. 1d: Our result showed that all of the *L. sativa* accessions formed a monophyletic clade, ... (Fig. 1d).
2. Fig. 1e: Model-based clustering analysis revealed additional inter- and intraspecies relationships. Asia- and Europe-originated accessions formed two groups ... (Fig. 1e and Extended Data Fig. 2).
3. Fig. 1c: The phylogenetic relationships were revealed by the principal component analysis (PCA) as well, ... (Fig. 1c and Extended Data Fig. 3a,b).

# Result example: logic between sections

**Fig. 1. Population structure of *Lactuca* accessions**



**Fig. 2. Lettuce domestication center near the Caucasus**



- Structure
  - The 1<sup>st</sup> paragraph, restates **the major findings**.
  - The 2<sup>nd</sup>, 3<sup>rd</sup>, and ..., discuss in details **topic by topic**, and bring out **the significance in the context**.
  - The last one ends with a **conclusion**.
- Read topic by topic
- Pay attention to the interpretation of results **in the context**

# Discussion example: major findings

## Discussion

In this work, we analyzed the genome sequences of 445 *Lactuca* accessions representing lettuce crop types and its wild gene pool species. More than 208 million sequence variants were identified, from which we revealed the population structure of the genebank collection and the domestication history of cultivated lettuce.

# Discussion example: one topic

WHY

As germplasms of major crops are maintained as genebank collections, understanding the population structure and phylogenetic relationships is of great importance for genebank management and utilization. In lettuce breeding, the GP1 species are used widely as

WHAT

there is no reproductive barrier within the group<sup>27</sup>. Our phylogenetic analyses clarified several issues regarding the taxonomic status of these GP1 species (see the Supplementary Note for a detailed discussion). First, the presumed GP1 species *L. georgica* should be reassigned as it clustered with the GP3 species *L. virosa*. The *L. dregeana* and *L. sagittata* samples are not to be considered as true wild species. Another GP1 species, *L. altaica*, has been considered as conspecific with *L. serriola*<sup>29,30</sup>, but the plastid phylogeny implied an introgression and fixation of a distantly related plastid haplotype in *L. altaica*. Phylogenetic analyses with additional samples will clarify these taxonomic issues in wild species. Our study also pointed out future directions in germplasm collection and utilization. Among the investigated samples, *L. serriola* from the Caucasus represents the most promising resource because the population from this area showed the highest nucleotide diversity. *L. aculeata* represents another potentially important gene pool, as its phylogenetic position distinct from other GP1 species suggests a different genetic repertoire. Thus, our study provided new insights regarding accession identity and genetic resources for crop improvement, demonstrating the value of whole-genome sequencing in the management of crop collections and the utilization thereof.

HOW

SO WHAT

# Discussion example: conclusion

Overall, our study constructed phylogenetic relationships within lettuce gene pool species and revealed the genetic basis of human selection during lettuce domestication. The genome sequences and the variation map generated in this study will serve as a valuable resource for lettuce research and breeding in the future.

- Extended Data
- Supplementary information
- Source Data
- Data and code availability
- Peer review information
- Accepted time

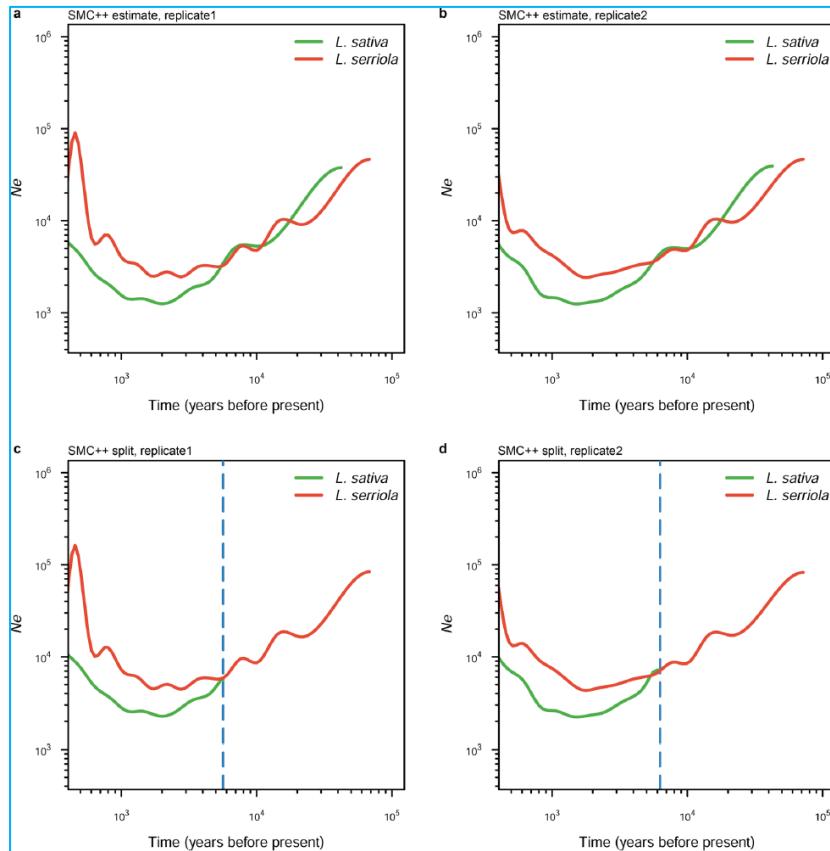
## Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41588-021-00831-0>.

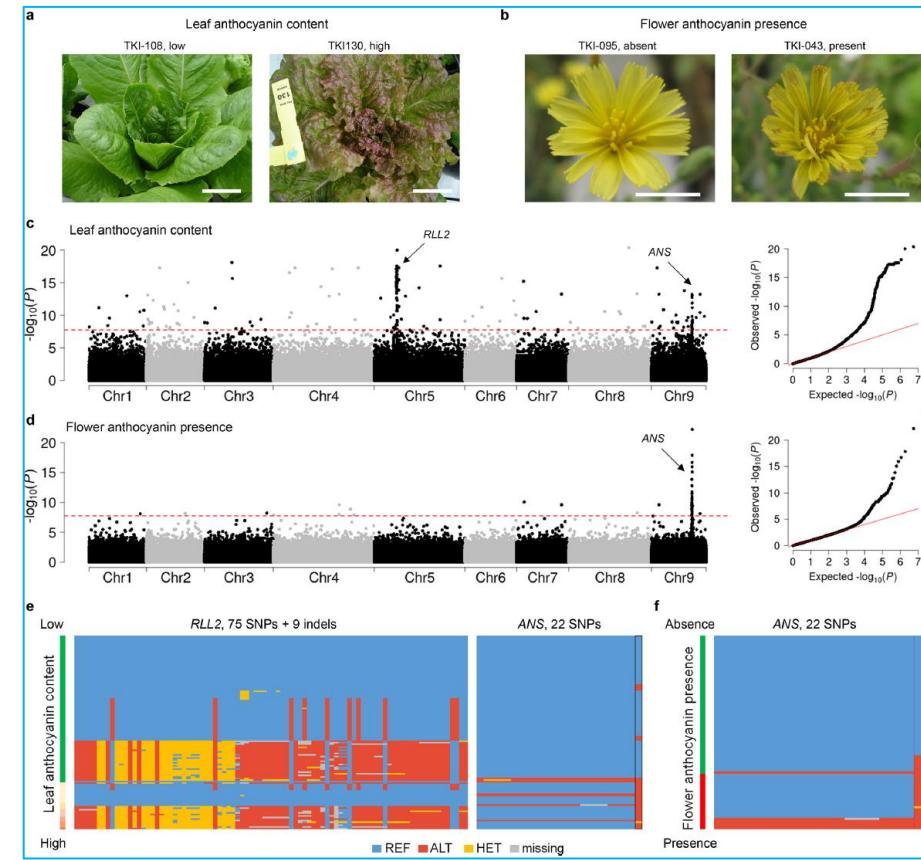
Received: 26 December 2019; Accepted: 1 March 2021;  
Published online: 12 April 2021

# Extended Data example

## Extended Data Fig. 5



## Extended Data Fig. 8



# Supplementary Information example

## Supplementary Note

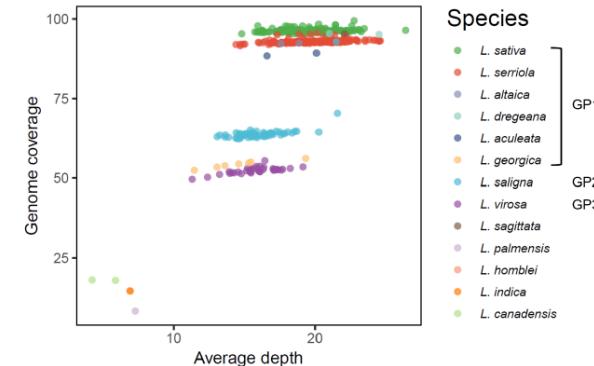
### Supplementary Note

#### 1. Lettuce Gene Pool species

Wild relative species are commonly used as a source of novel traits in lettuce breeding, such as in variety improvement for biotic and abiotic stress and for quality characters. The International Plant Name Index (IPNI) includes several hundreds of species of the genus *Lactuca*, but the majority of these taxa refer to synonyms and basionyms or have been reassigned to other genera<sup>1</sup>. Generally, the *Lactuca* genus is considered to consist of about 100 species, of which approximately 20 are part of the lettuce gene pool<sup>2,3</sup>. The primary gene pool (GP1) consists of completely inter-fertile taxa, including the crop species *L. sativa* L. and its wild relatives *L. aculeata* Boiss. & Kotschy, *L. altaica* Fisch. & C.A. Mey., *L. azerbaijanica* Rech. f., *L. dregeana* DC., *L. georgica* Grossh., *L. scarioloides* Boiss. and *L. serriola* L. The secondary gene pool (GP2) is formed by *L. saligna* L. alone, while the tertiary gene pool (GP3) includes *L. acanthifolia* (Willd.) Boiss., *L. alpestris* (Gand.) Rech. f., *L. aurea* (Vis. & Pančić) Stebbins, *L. longidentata* DC., *L. orientalis* (Boiss.) Boiss., *L. quercina* L., *L. sibirica* (L.) Benth. ex Maxim., *L. tatarica* (L.) C.A. Mey., *L. viminea* (L.) J. Presl & C. Presl, *L. virosa* L. and *L. watsoniana* Trel. *L. serriola*, *L. saligna* and *L. virosa* are the main species that have been extensively used in breeding.

## Supplementary Figures

### Supplementary Figures



Supplementary Fig. 1. Average sequencing depth and genome coverage of 445 *Lactuca* accessions on the lettuce reference genome. The primary (GP1), secondary (GP2), and tertiary gene pool (GP3) species are indicated in the legend.

# Source Data example

## Source Data

### Source data

#### Source Data Fig. 1

Statistical source data.

#### Source Data Fig. 2

Statistical source data.

#### Source Data Fig. 3

Statistical source data.

#### Source Data Fig. 4

Statistical source data.

#### Source Data Extended Data Fig. 1

Statistical source data.

#### Source Data Extended Data Fig. 2

Statistical source data.

#### Source Data Extended Data Fig. 3

Statistical source data.

## Source Data Fig. 1

A Country	B L. sativa	C L. serriola	D L. saligna	E L. virosa	F other	G Total number
Afghanistan	0	2	0	0	0	2
Argentina	1	0	0	0	0	1
Armenia	0	4	0	0	0	4
Austria	3	0	0	0	0	3
Azerbaijan	0	6	0	0	1	7
Belgium	1	4	0	0	0	5
Bulgaria	1	5	3	0	0	9
Canada	1	0	0	0	1	2
Chile	0	0	0	1	0	1
China	7	0	0	0	0	7
Croatia	1	2	0	0	0	3
Czech Republic	1	3	0	0	1	5
Denmark	1	2	0	0	0	3
Egypt	2	3	0	0	0	5
France	29	15	1	11	0	56
Georgia	0	5	1	0	0	6
Germany	6	15	0	1	0	22
Greece	0	6	2	0	0	8
Hungary	0	17	0	0	0	17
Indonesia	0	0	0	0	1	1
Iran	1	3	0	0	0	4
Iraq	0	2	0	0	0	2
Israel	0	15	24	0	1	40
Italy	1	12	10	0	0	23

Fig1a Fig1b Fig1c Fig1d Fig1e +

## Data availability

All raw sequencing data were deposited into the Sequence Read Archive (under BioProject accession PRJNA693894) and CNGB Nucleotide Sequence Archive (CNSA; under the accession number CNP0000335). Variant files, genome assemblies and annotation files are stored in CNSA under the same accession number. Source data are provided with this paper.

## Code availability

All of the code used in this study is available at <https://github.com/popgenome/lettuce2020>.

- Find and follow **the logic/story line**
- Pay attention in **details** in figures, tables, and methods
- Ask self the four questions in main text
- Make notes

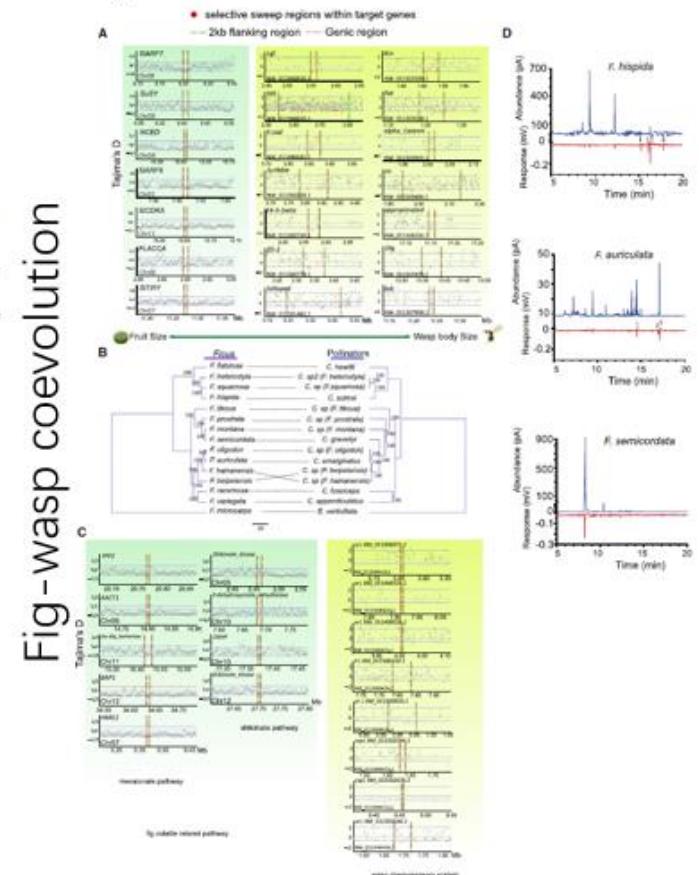
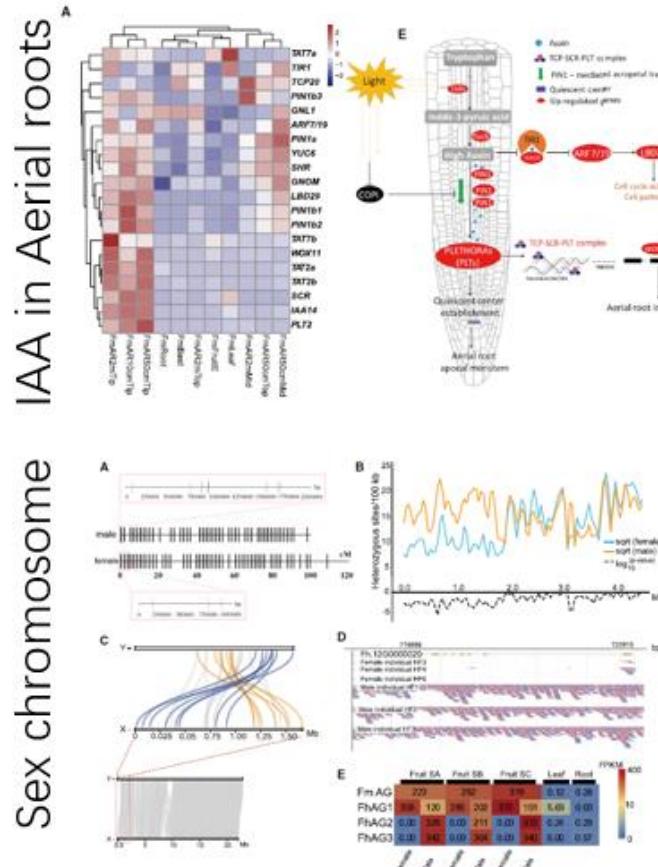
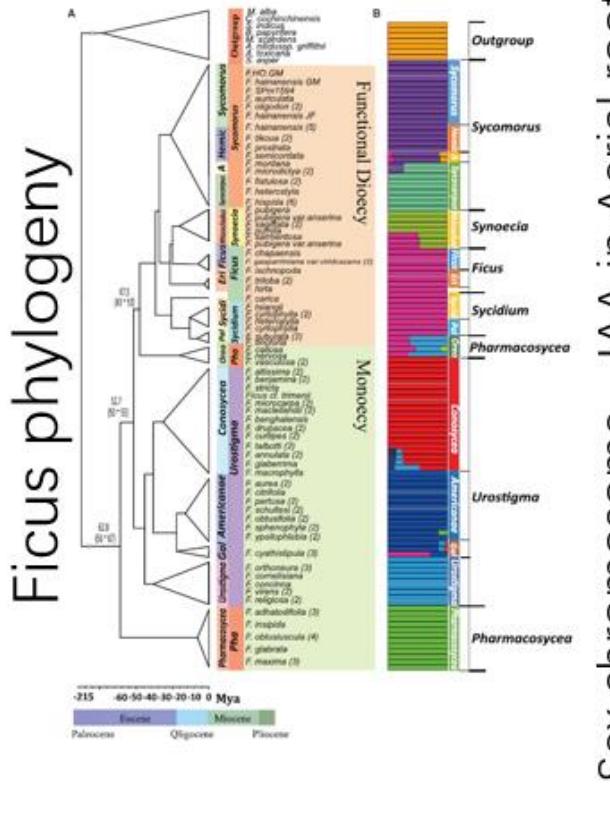
# Note example: methods

Genomes of the banyan tree and pollinator wasp provide insights into fig-wasp coevolution. Cell (2020) 183:1

- Library: 63 PacBio SMRT cells (30 for *Ficus macrocarpa* and 33 for *F. hispida*) sequenced on Pacbio RSII; paired-end libraries sequenced on Illumina X10 with 150-bp cycles and 300-500 bp insert size, Hi-C libraries sequenced on Illumina X10 (80% valid reads evaluated by HiC-Pro); *E. vorticellata* sequenced in 170x on PacBio and 84x on Illumina; RNA libraries sequenced on Illumina HiSeq 2500 or X10.
- Assembly: PacBio reads self-corrected by CANU; assembled by CANU, FALCON, SMARTdenovo and evaluated by N50/size/BUSCO; merge by Quickmerge; assemblies polished by Pilon; uniquely-mapped Hi-C reads corrected by 3D-DNA pipeline and used for scaffolding by ALLHiC; wasp genome assembled by CANU and redundancy removed by Redundans; evaluated by BUSCO, PASA transcripts, Illumina reads, chromatin interactions, and genetic maps based on resequencing of 30 male and female.
- Annotation: a repeat library constructed by RepeatModeler; TEs identified by RepeatMasker, tandem repeats by TRF; TE classified by TEclass; annotation using MAKER pipeline, including training on Trinity/PASA assembled transcripts by SNAP, GENEMARK, AUGUSTUS, a second round of training with AED<0.2 gene models, and combined with HISAT/StringTie assembled transcripts and homologous proteins from 6 species; miRNA predicted by mapping miRBase miRNAs to two genomes by bowtie.
- Comparative genomics: segmental duplication identified by the 1<sup>st</sup> round of blast of 400-kb unique segments with at least 88% identity and 500 bp and a 2<sup>nd</sup> round alignment of global alignments at least 90% identity and 1 kb; SVs identified by MUMmer alignment of MaSuRCA assembled contigs and a web-based tool Assemblytics; CNVs significantly different from genome average depth identified by count Illumina reads in 5-kb non-TE sequence.
- Phylogeny: concatenated tree constructed from MUSCLE alignment of single-copy genes 5 species; ML tree constructed on SNPs by IQ-Tree and RAxML/GTRCAT model; divergence time estimated by PAML/MCMCTREE; co-phylogeny analyzed by Jane4, codivergence analyzed by Tajima's D in 14 fig-wasp pair; selective sweep by SweeD.
- Population: ancestry identified by ADMIXTURE; coefficient of genetic relatedness calculated within a 200-kb sliding window using IBD; ABBA analyzed by ANGSD/doAbbababa2.
- Hormone: 50 mg fresh tissues ground in liquid nitrogen, extracted in 0.5 mL methanol/water/formic acid, evaporated in nitrogen gas, reconstituted in 80% methanol ultraphoniced and filtered for LC-ESI-MS/MS.
- Sex chromosome: phased SNPs called from minimap2 alignment of corrected PacBio alignment by GATK4 with error correction by WhatsHap; SNPs from 40x *F. hispida* WGS data used to determine sex-phased blocks with at least 70% of sex-specific SNPs; sex chromosome de novo assembly by CANU and scaffolded by ALLHiC; sex-determining regions identified based on heterozygous site density from 15 male and 11 female; expression called from RNAseq from flowers at 3 stages by bowtie-RSEM pipeline.
- Pollinator attraction: odor from fig syconia at receptive phase collected by solid-phase micro-extraction over the headspace for 1h and profiled by GC-MS; stimulus to wasps detected by gas chromatography-electroantennogram detection (GC-EAD) with 9 replicates.

# Note example: results

Genomes of the banyan tree and pollinator wasp provide insights into fig-wasp coevolution. *Cell* (2020) 183:1



# Note example: in-depth reading

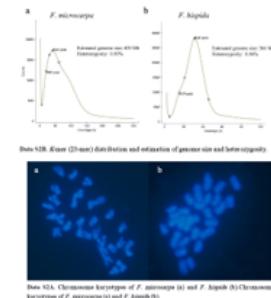
**Background: banyan tree from *Ficus* genus**

- Ecologically important: provide year-round fruit production in tropical forests
- The largest genus with ~800 species in the Moraceae family (桑科)
- Hemi-epiphytic habit; aerial roots
- A diverse sex determination system
- Fig-wasp obligate mutualism

4

## Genome survey

- Genome size, 366 Mb for *F. hispida*, 430 Mb for *F. microcarpa*, estimated from Illumina short reads using perl scripts
- Nuclear DNA estimated by flow cytometry analysis



5

## Genome assembly of *F. microcarpa*, *F. hispida* and *E. verticillata*

- 63 PacBio SMRT cells (30 for *Ficus macrocarpa* and 33 for *F. hispida*) sequenced on Pacbio RSII; Paired-end libraries sequenced on Illumina X10 with 150-bp cycles and 300-500 bp insert size; Hi-C libraries sequenced on Illumina X10 (80% valid reads evaluated by HiC-Pro)
- RNA libraries sequenced on Illumina HiSeq 2500 or X10.
- PacBio reads self-corrected by CANU; assembled by CANU, FALCON, SMARTdenovo
- Evaluated by N50/size/BUSCO; merge by Quickmerge; assemblies polished by Pilon;
- Uniquely-mapped Hi-C reads corrected by 3D-DNA pipeline and used for scaffolding by ALLHiC.

6

## Assembly validation by genetic map

- SNPs called from resequencing of 30 male and 30 female F1 individuals by GATK
- Maternal and paternal genetic maps determined from bin markers

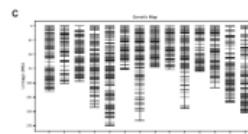
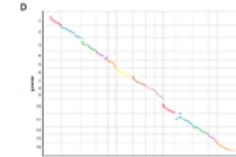


Figure S1. Assessment of Chromosome-Level Genome Assemblies of the Two *Ficus* Species, Related to Figure 1. Genome-wide analysis of chromatin interactions at 150-kb resolution in *F. microcarpa* (A) and *F. hispida* (B) genomes. (C) A high-density genetic map for *F. hispida* F1 population. (D) Comparison of F1 genetic map with Hi-C assembly in *F. hispida*.

7

## Assembly evaluation

- Completeness assessed by BUSCO on 1,440 conserved plant genes
- PASA assembled transcripts mapped to two assemblies
- Illumina reads mapped by bwa
- Genome-wide chromatin interaction



8

## Genome annotation

- A repeat library constructed by RepeatModeler; TEs identified by RepeatMasker; tandem repeats by TRF; TE classified by TEclass
- Annotation using MAKER pipeline, including training on Trinity/PASA assembled transcripts by SNAP, GENEMARK, AUGUSTUS, a second round of training with AED<0.2 gene models
- Combined with HISAT/StringTie assembled transcripts and homologous proteins from 6 species.

Table 1. Statistics for Assembly and Annotation of the <i>Ficus</i> genome		
Chromosome	No. of contigs	Total length of contigs (Mb)
Chromosome 1	408	34,791,363
Chromosome 2	408	36,202,493
Chromosome 3	205	32,700,637
Chromosome 4	205	35,161,573
Chromosome 5	205	34,415,657
Chromosome 6	360	35,213,584
Chromosome 7	443	35,486,046
Chromosome 8	268	35,171,889
Chromosome 9	193	35,152,172
Chromosome 10	463	35,733,657
Total no. of contigs	2,000	35,200,000
Total length of contigs (Mb)	4,790	3,674
Total length of assembled genome (Mb)	405	3,674
Assembly size (Mb)	39.62	3,674
No. of genes	26,716	23,271
Length of genes (Mb)	102	102

9

## Other article types

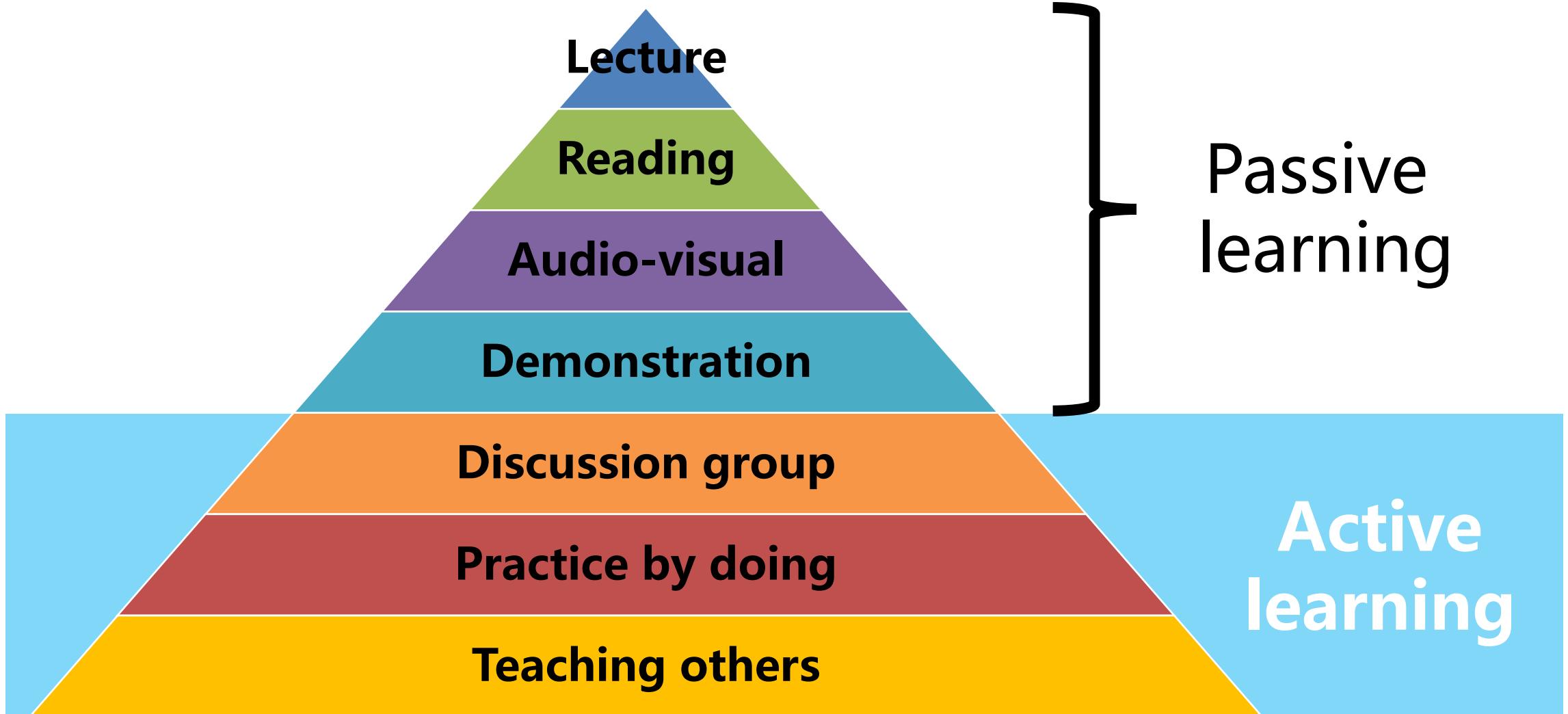
---

- Review
- Technical advances
- Data notes
- Editorial, Perspective, Comments, etc.
  
- Resources: PUBMED, google scholar, sci-hub

# Section III

# Scientific presentation

# Learning Pyramid



- Scientific presentation is
  - An **essential skill** of introducing personal work, communicating with peer, and seeking for supports, and
  - A **professional way** of describing observation, introducing hypotheses, demonstrating and interpreting results, and summarizing into conclusions.

# Prepare for presentations

---

- Prepare presentations in a logical way
- Elucidate details in methods and materials
- Summarize results in an effective manner
- Know your audience
  - **Use bare minimum** and provide only essential information
  - Stick to time limits (1 minute per slide)
- Know your story – **follow the logic**

# IMRaD structure for presentation

---

- IMRaD: **I**ntroduction, **M**ethods, **R**esults, and **D**iscussion
- Title slide includes the full title, the names, date, etc.
- Introduction
  - Provide a brief introduction of the **background**
  - State clearly the **aim or scientific questions**
- Methods
  - Use bullet points
  - Emphasize on **important details/parameters**

- Results
  - Present ONLY one thing in each slide
  - Use **figures** whenever possible
  - Enlarge text and numbers
- Discussion/Conclusion
  - Summarize the most important findings
- Other slides, acknowledgements, extra evidence,

- Follow a **logical structure**, with focus on major discoveries
  - Use **figures and tables** with an appropriate font size
  - Avoid using a big chunk of text
  - Use **own words** instead of copy from papers
  - Keep consistency in slides, font, size, color, etc.
- 
- **Practice, practice, practice...**

# Tips in presentation

---

- Plan for an impressive and informative start
  - Adjust **voice and pace** to ensure clarity
  - Make **eye contact** with your audience
  - Stay connected with your slides
  - Explain the details in each slide
- 
- Avoid reading each word in your slides

# Case: lettuce resequencing study



## 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.

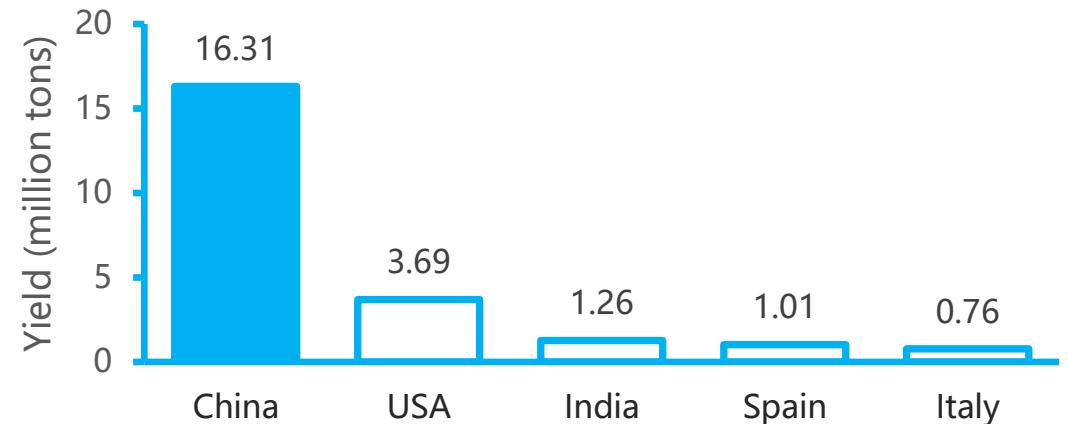
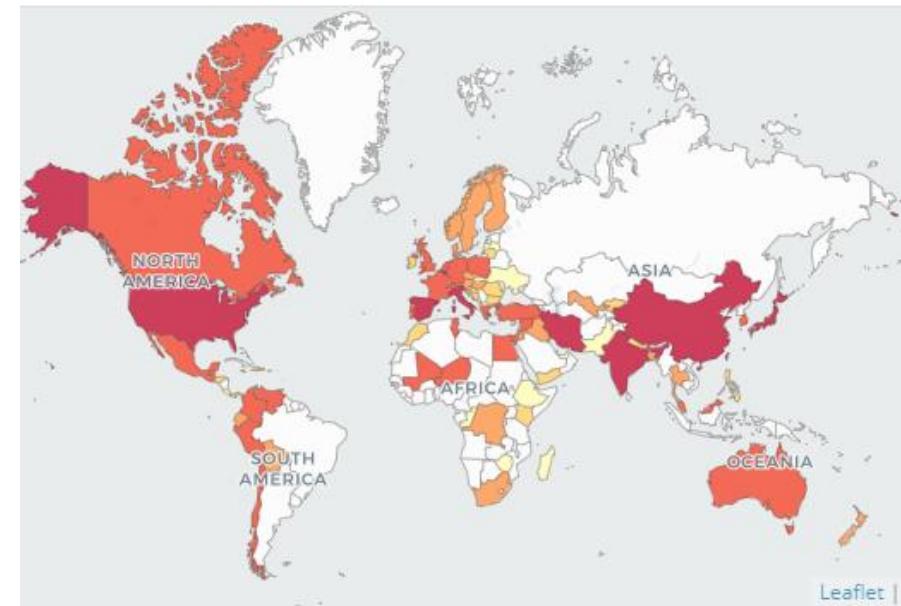
# Domestication history of cultivated lettuce revealed by resequencing 445 *Lactuca* accessions

魏桐

9月11日

# Lettuce is an important vegetable crop

- Consumed worldwide
- A rich source of vitamin K and vitamin A, and a moderate source of folate and iron
- Ranks **the 3<sup>rd</sup>** in leafy vegetable worldwide



# Lettuce is a model Asteraceae plant

- A short life cycle
- Easy for transformation
- Well-preserved  
germplasms
- Various agronomic traits
- Potential bioreactor



# Lettuce crop types



Crisp



Butterhead



Cutting



Cos



Stalk



Oilseed



*L. serriola*  
proposed wild progenitor

# Lettuce cultivation history



From a 14<sup>th</sup> century medical book in Europe



Painted stela of Tatiaset from the ancient Egypt

# The oldest archaeological evidence



This tomb dates to the V dynasty, 2494-2345 BC

# Scientific questions

---

- Population related
  - What is **the population structure** in lettuce germplasms?
  - What are **the phylogenetic relationships** between wild relatives and cultivated lettuce?
- Domestication related
  - **Where and when** was lettuce domesticated?
  - What are **the major events and selected traits** during lettuce domestication and improvement?
- Breeding related
  - What are **the genetic determinants** for phenotypic differences?
  - What are **the genetic determinants** for Bremia resistance?

## CGN lettuce collection

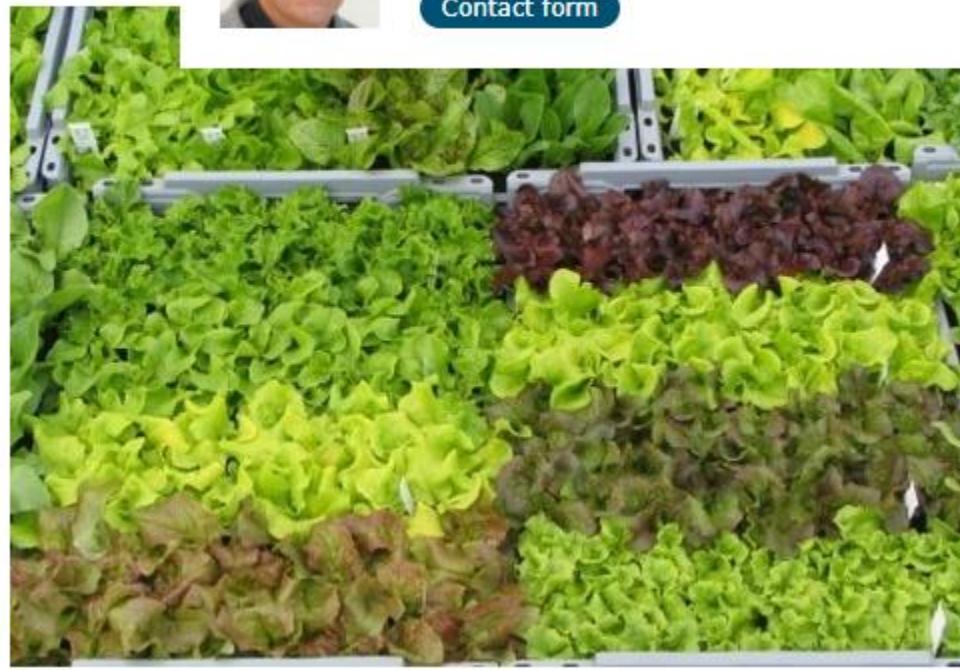
CGN maintains one of the largest genetic resources collections of lettuce in the world, currently (April 2020) including 2535 accessions. CGN's lettuce collection contains a fairly high representation of wild species, covering over 40% of the total collection. The collection is relatively well characterized by morphological and molecular descriptors, while also many trait data have been collected. More information about various aspects of the collection can be accessed below.



Curator

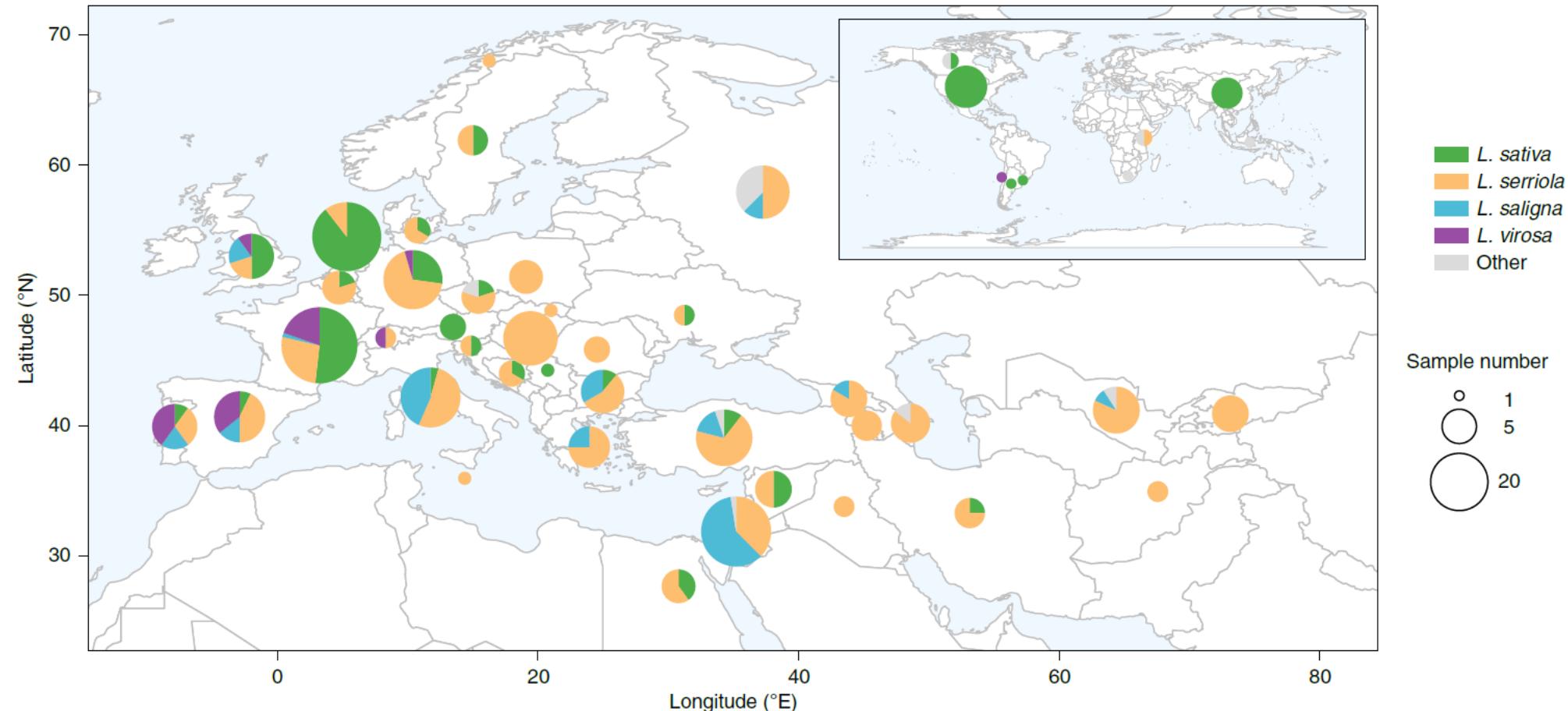
dr. R (Rob) van Treuren

[Contact form](#)

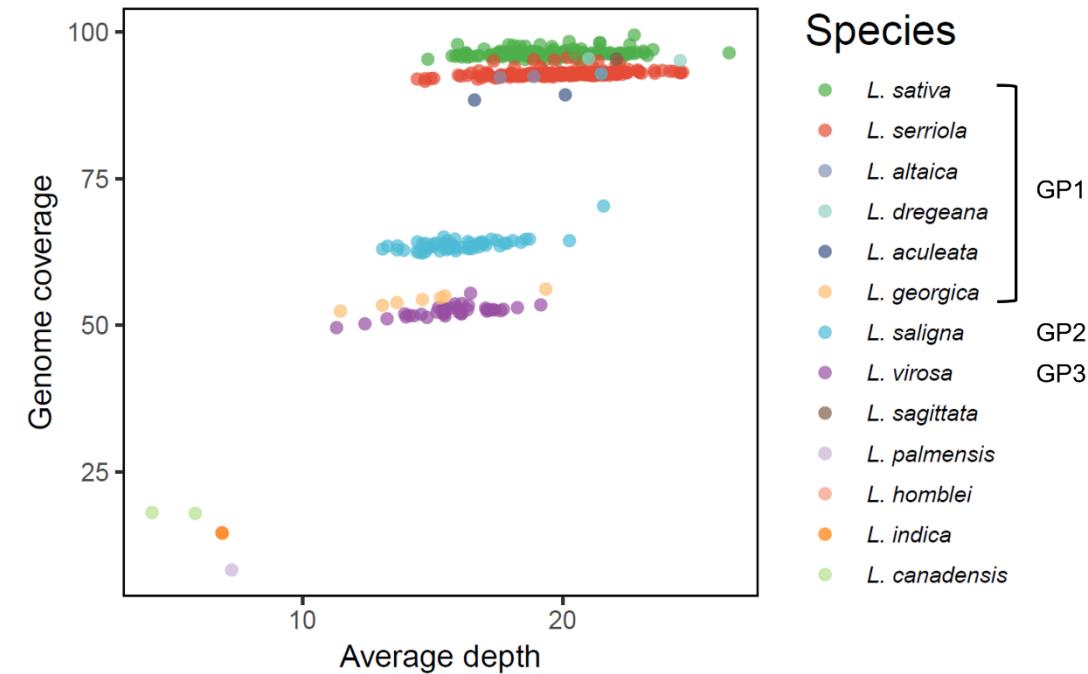
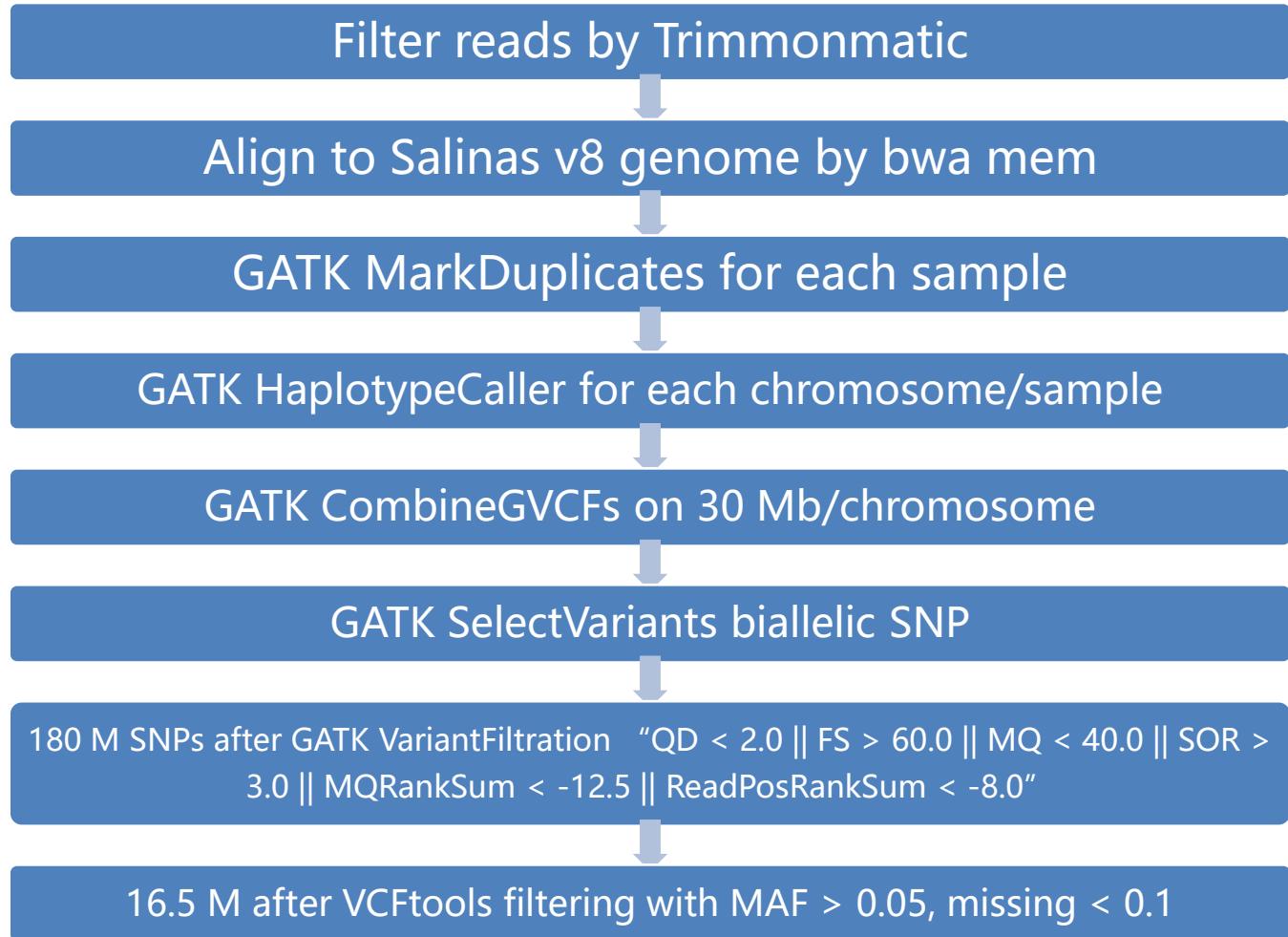


# Geographic distribution

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



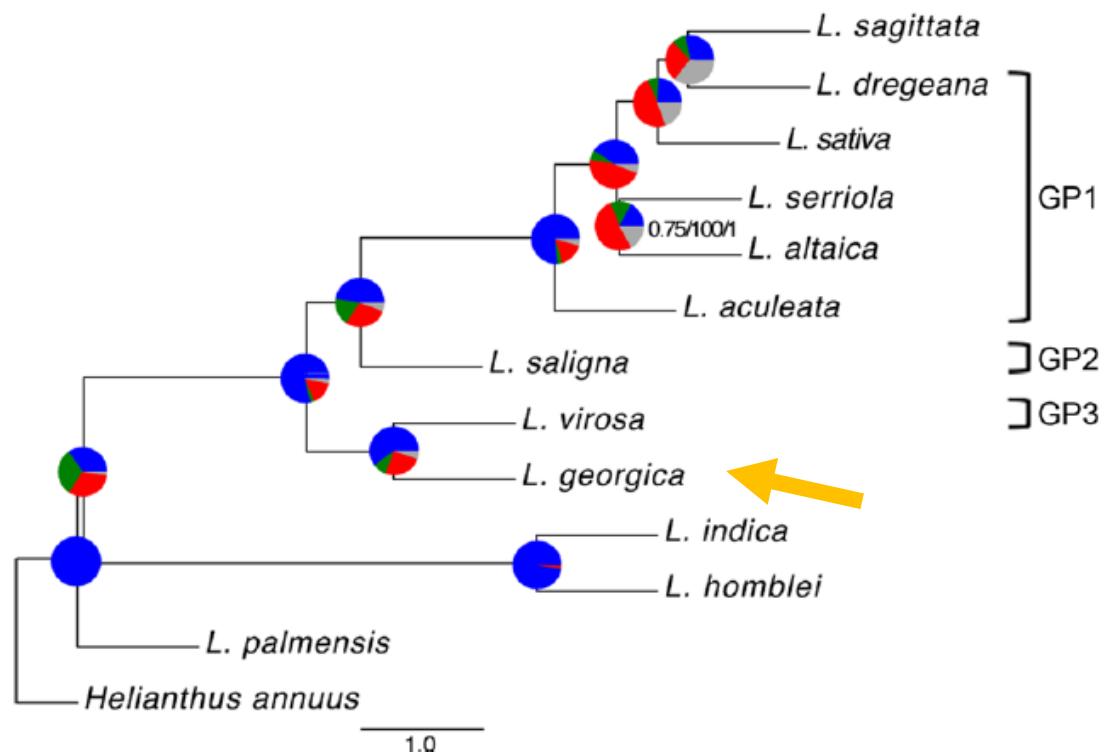
# Read alignment and variant calling



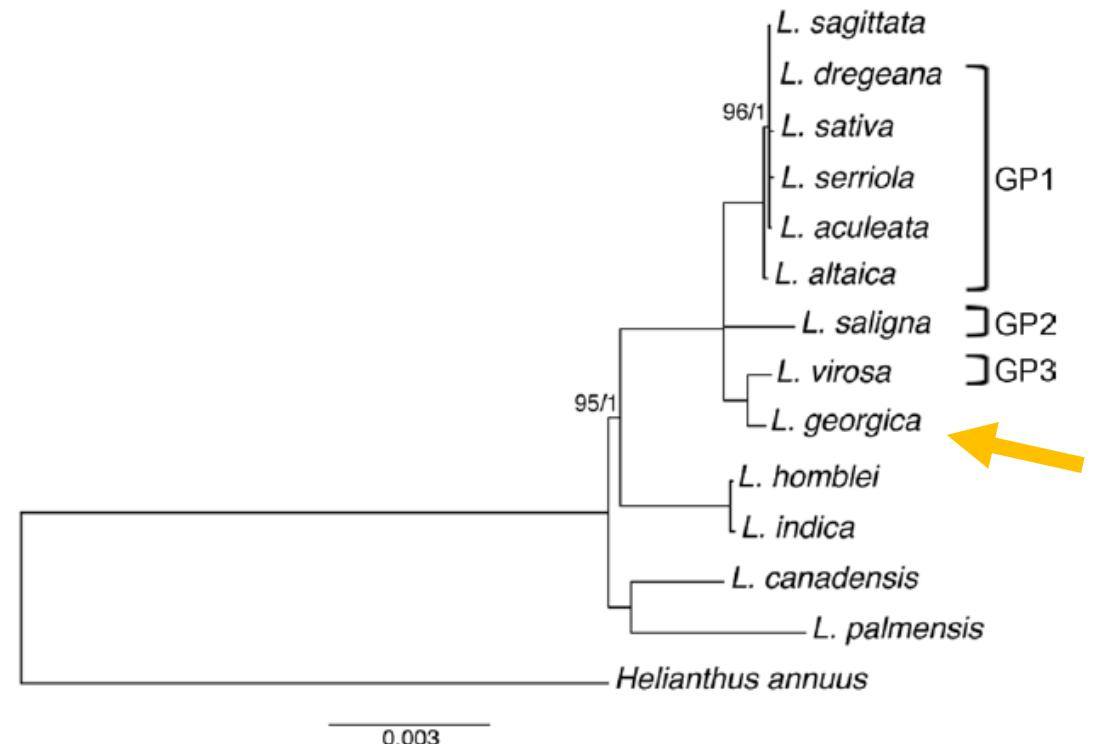
# *De novo assembly of 12 wild genome*

Species name	Gene pool	Est. size (Mb)	Heterozygosity (%)	Repeat (%)	Assembly (Mb)	Scaffold N50 (bp)	BUSCO	Chr. number	Ploidy level
<i>L. serriola</i>	GP1	2,848	0.17%	62.30%	2,069	6,406	88.50%	18	2
<i>L. georgica</i>	GP1	4,168	0.17%	66.41%	3,153	4,718	88.90%		
<i>L. aculeata</i>	GP1	2,739	0.05%	61.34%	2,194	9,169	89.40%	18	2
<i>L. altaica</i>	GP1	2,670	0.11%	61.18%	2,234	11,528	90.20%	18	2
<i>L. dregeana</i>	GP1	2,804	0.23%	62.00%	2,249	8,094	89.10%	18	2
<i>L. saligna</i>	GP2	2,413	0.15%	60.94%	1,761	8,104	88.30%	18	2
<i>L. virosa</i>	GP3	3,438	0.26%	65.46%	2,913	4,910	88.20%	18	2
<i>L. canadensis</i>		9,987	0.27%	61.11%	6,091	2,434	72.80%	34	4
<i>L. homblei</i>		6,068	0.22%	66.79%	3,705	5,643	85.50%	18	2
<i>L. indica</i>		6,200	0.19%	66.56%	3,299	3,485	83.10%	18	2
<i>L. sagittata</i>		2,762	0.13%	61.29%	1,924	6,815	89.20%		
<i>L. palmensis</i>		984	0.10%	67.20%	720	8,388	89.10%		2

# Phylogeny among *Lactuca* spp.



Coalescence-based tree from  
4,513 nuclear genes

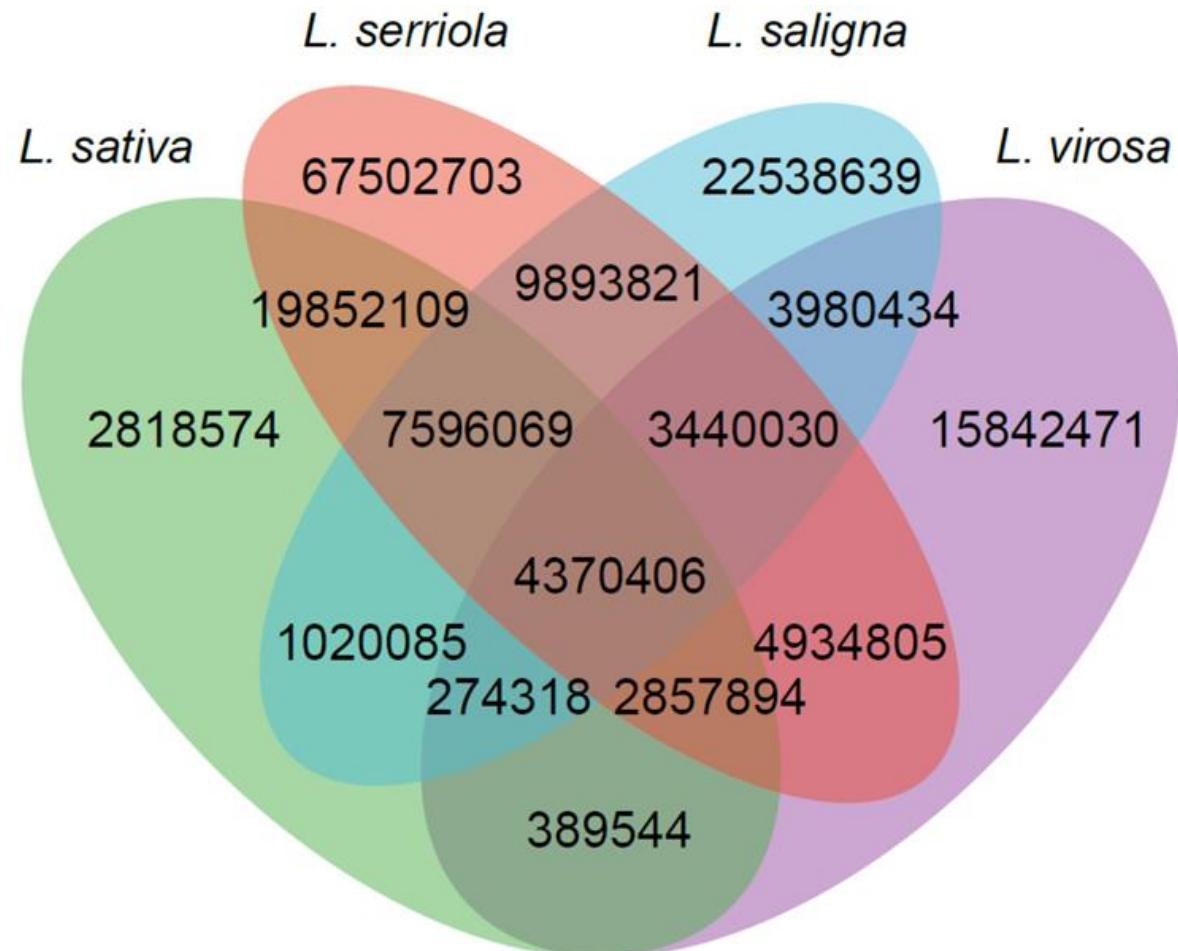


Concatenation-based tree from  
75 plastid genes

# A comprehensive variation map

	<i>L. sativa</i>	<i>L. serriola</i>	<i>L. saligna</i>	<i>L. virosa</i>	Total
Accessions	133	199	56	36	440
SNPs	39,178,999	120,447,837	53,113,802	36,089,902	178,807,215
indels	3,125,340	11,705,344	9,538,216	6,333,010	29,535,608
SVs	64,956	201,188	17,549	13,680	244,866
Nonsynonymous	219,770	700,293	567,159	447,516	1,580,355
Synonymous	196,874	541,084	625,850	501,114	1,399,417
5'-UTR	35,294	109,579	135,152	102,524	313,409
3'-UTR	31,409	97,492	112,224	85,743	262,579
Intron	499,831	1,472,748	1,129,693	845,087	3,010,148
Splicing	20,287	58,515	65,845	50,567	151,540
Intergenic	38,182,102	117,487,760	50,506,000	34,079,182	172,149,804

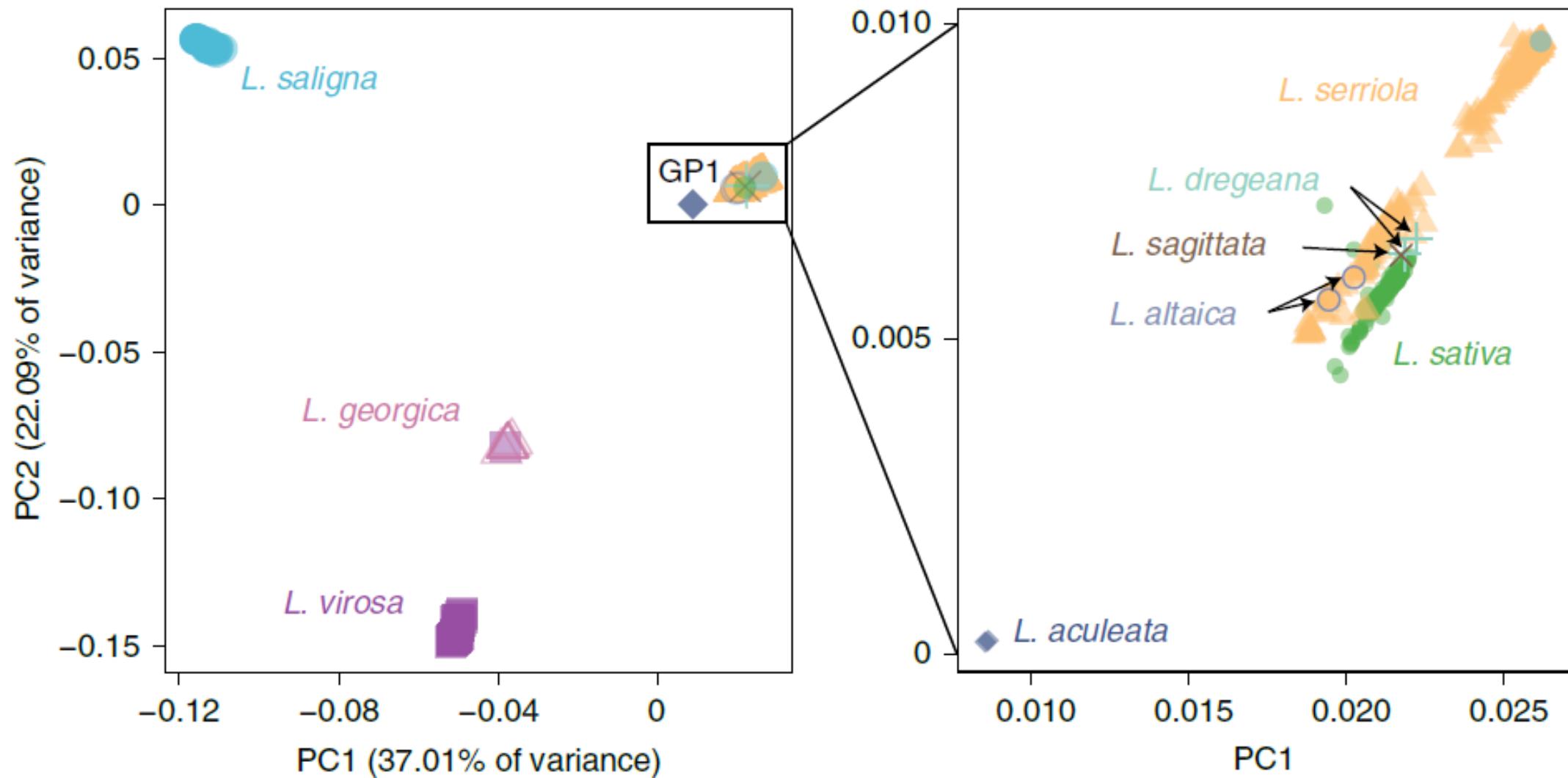
88.51% SNPs in *L. sativa* are shared with *L. serriola*.



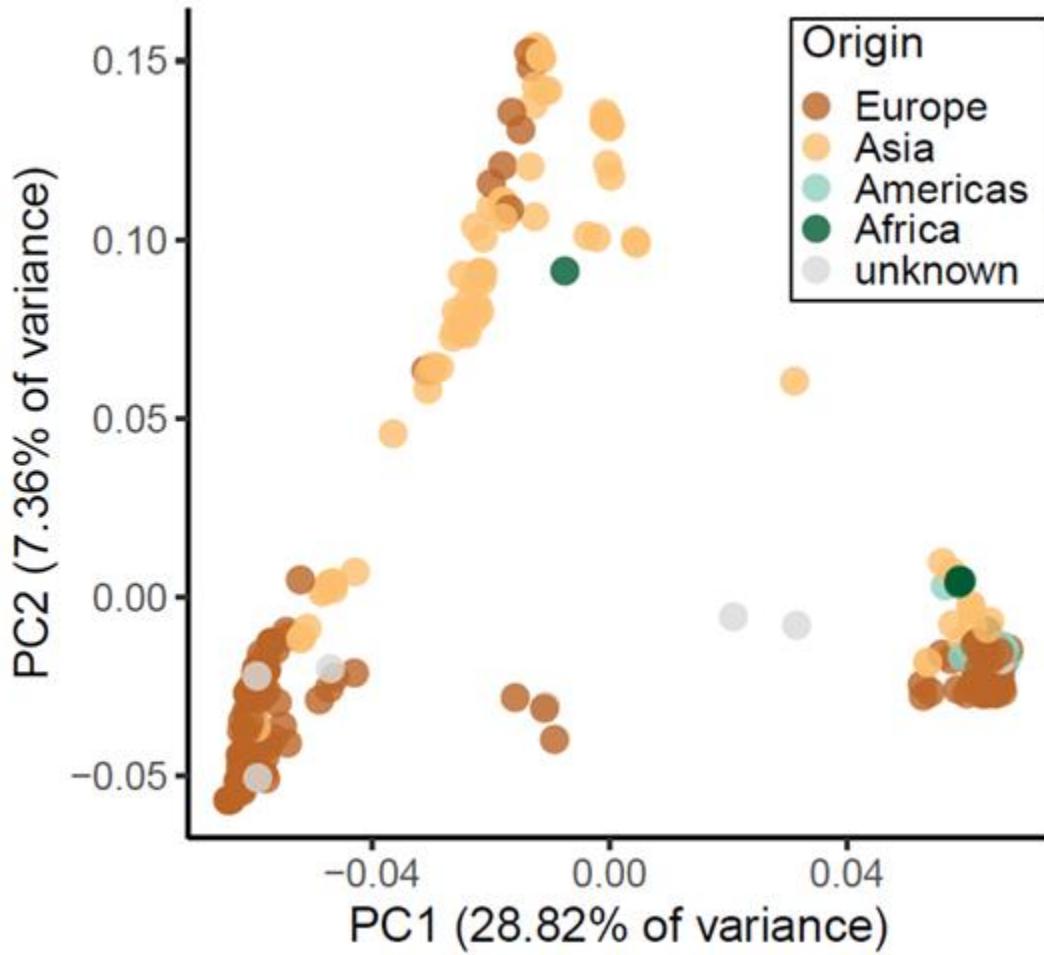
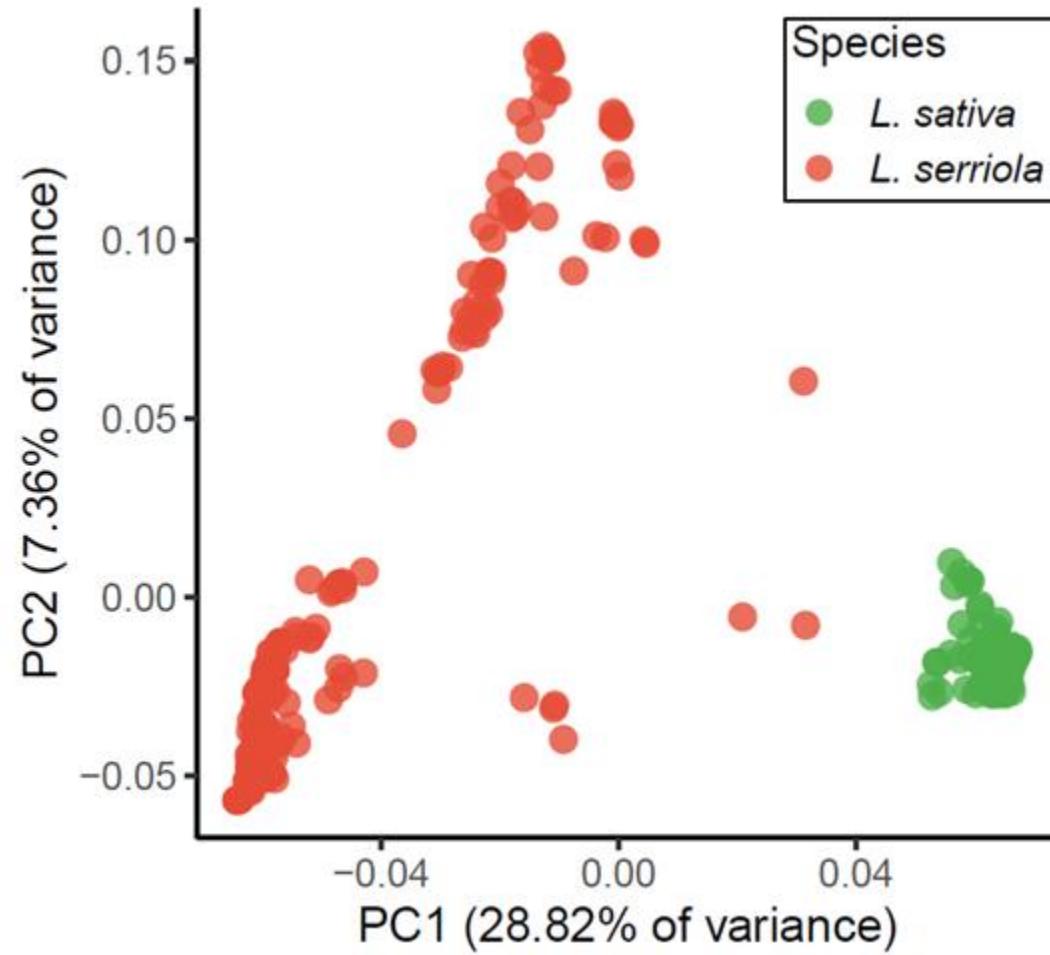
# Reduced diversity in cultivated lettuce

Species	Nucleotide diversity	$F_{ST}$ with lettuce
<i>L. sativa</i>	1.68	-
<i>L. serriola</i>	2.94	0.305
<i>L. saligna</i>	1.23	0.548
<i>L. virosa</i>	1.07	0.646

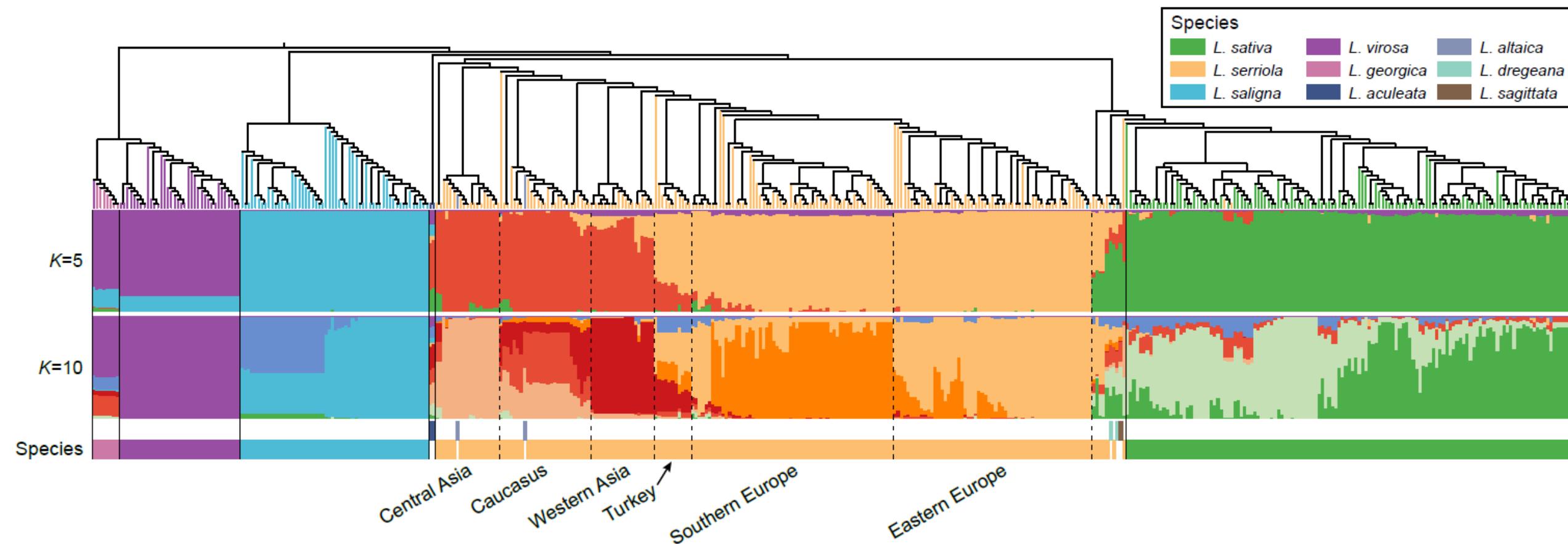
# GP1 samples cluster together



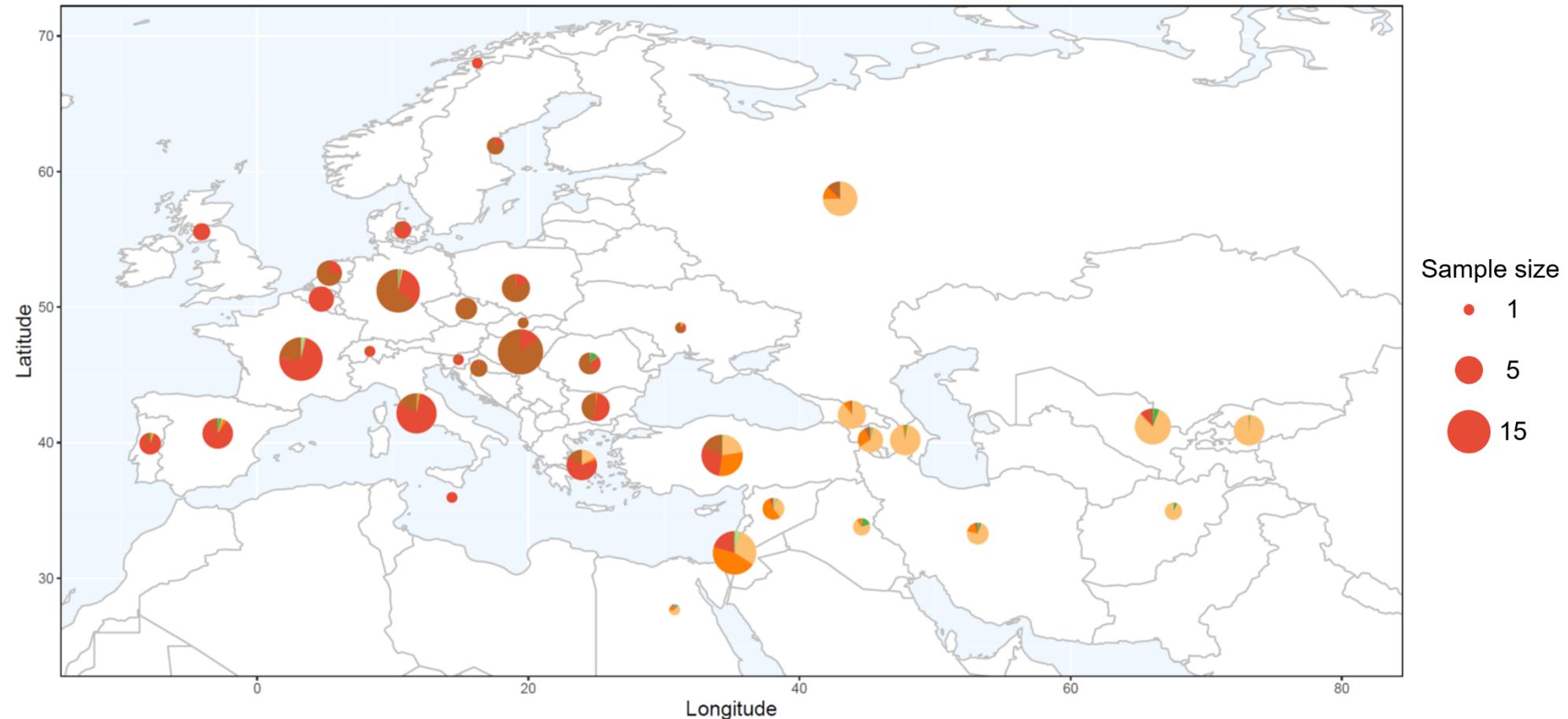
# Two groups in the *L. serriola* population



# Intra-specific structure in *L. serriola*



# Genetic composition in *L. serriola*

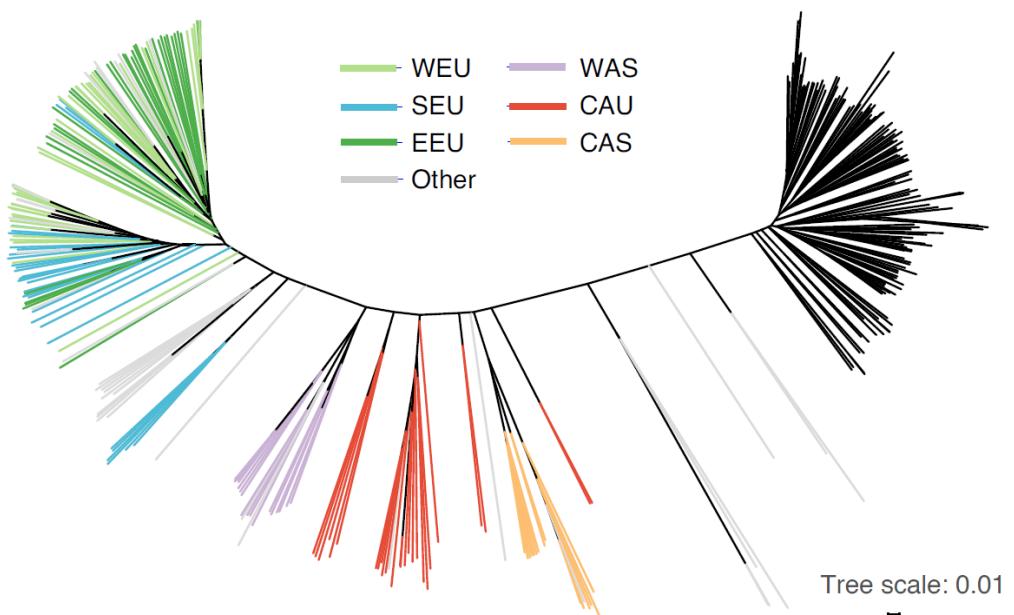


# Answers to the scientific questions

---

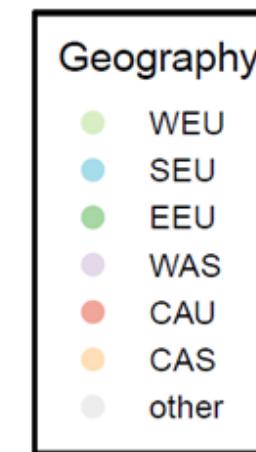
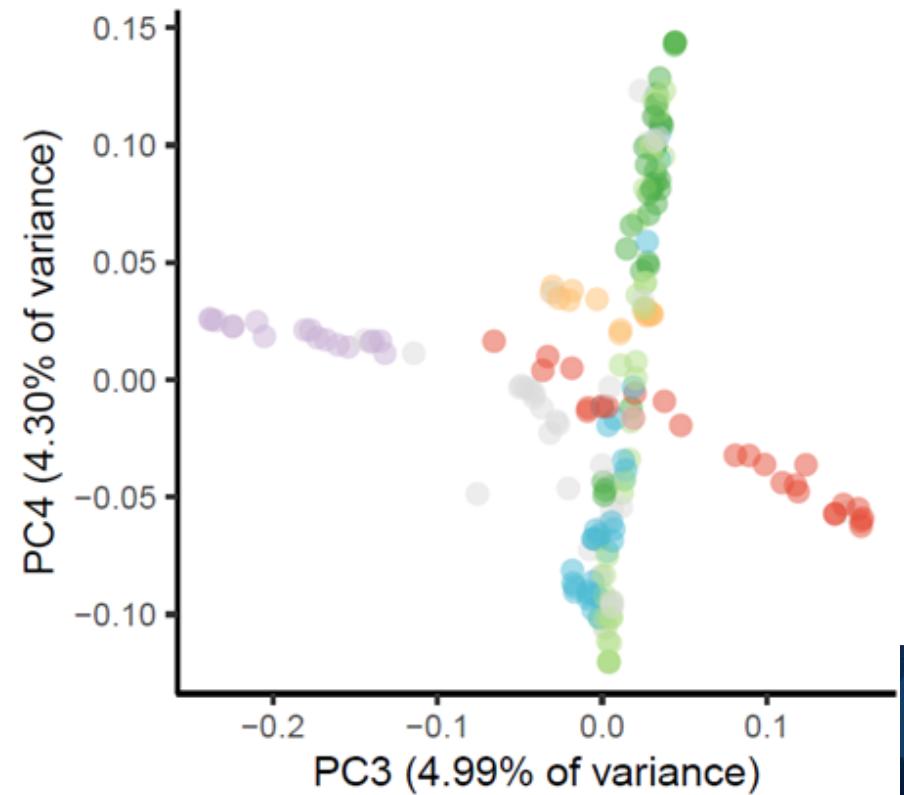
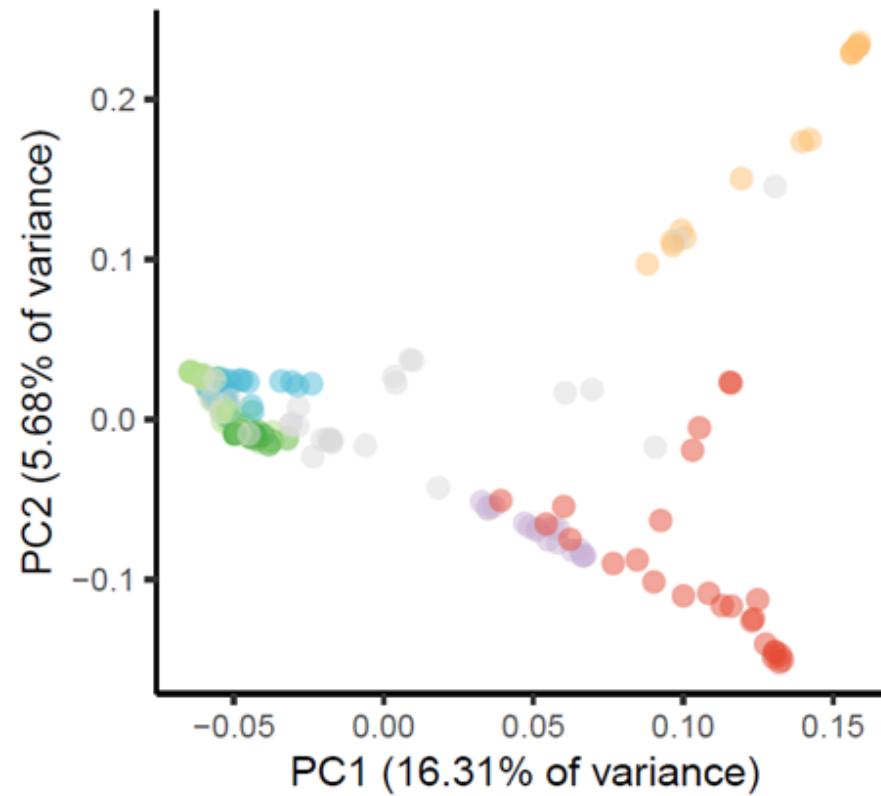
- Q1: What is the population structure in lettuce germplasms?
- A1: Cultivated lettuce clustered together, indicating a single domestication event; intra-specific structure exist in *L. sativa* and *L. serriola*.
- Q2: What are the phylogenetic relationships between wild relatives and cultivated lettuce?
- A2: *L. serriola*, especially the Asian group, is close to modern lettuce; the assumed GP1 *L. georgica* belongs to GP3.

# Phylo-geographic grouping of *L. serriola*

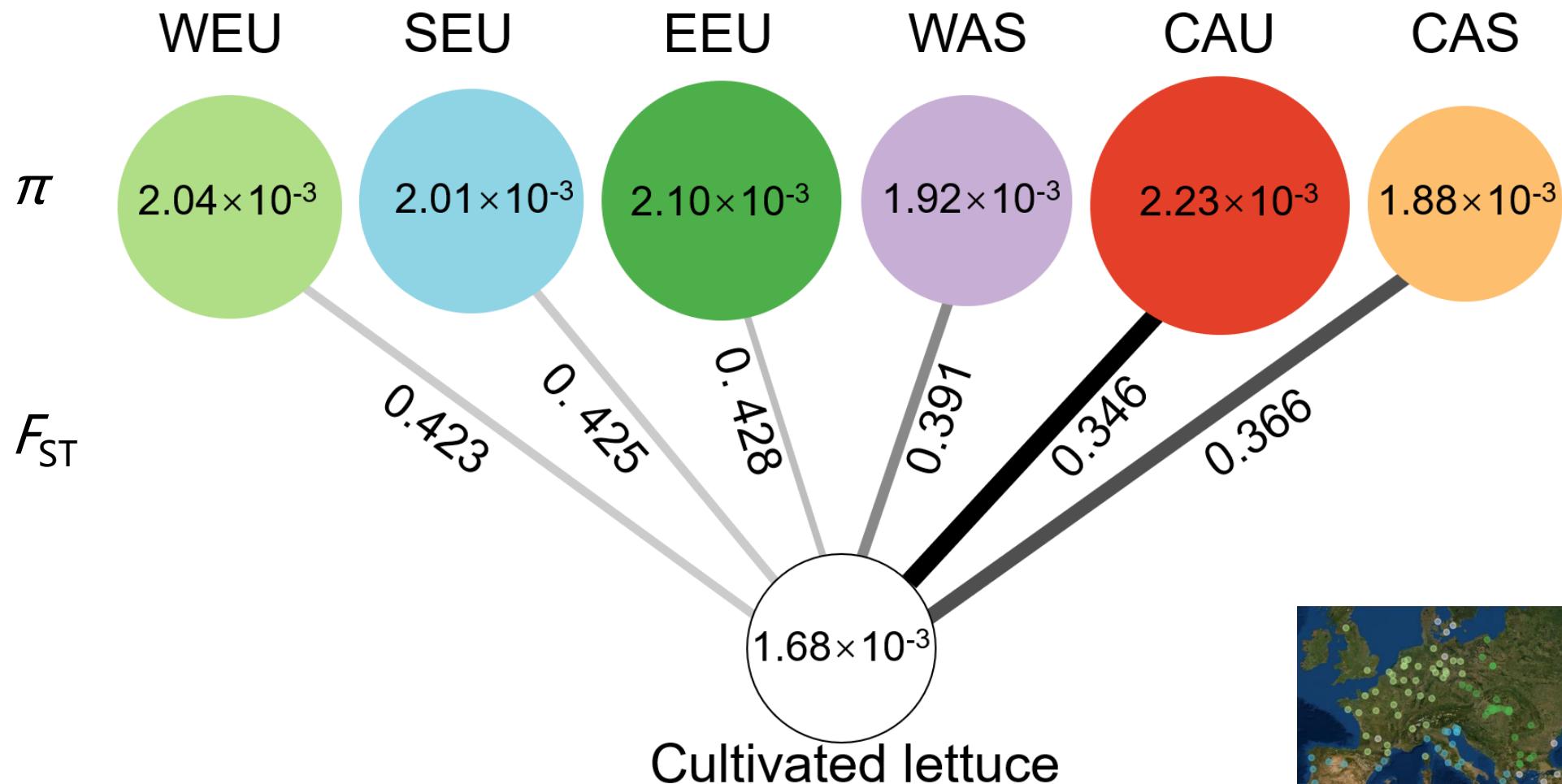


3 European groups  
WEU: western Europe  
SEU: southern Europe  
EEU: eastern Europe

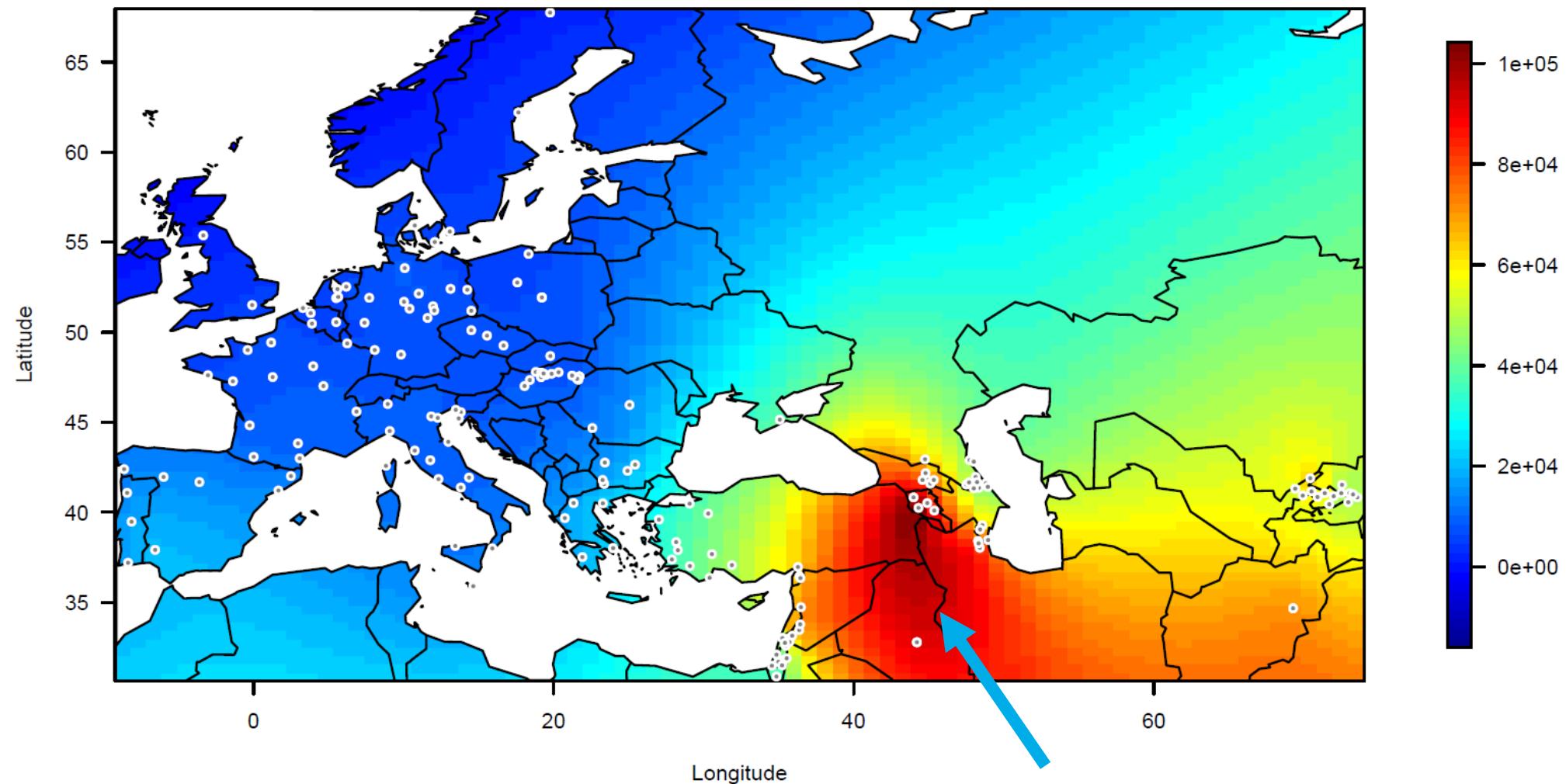
3 Asian groups  
WAS: western Asian  
CAU: Caucasian  
CAS: central Asian

PCA of *L. serriola*

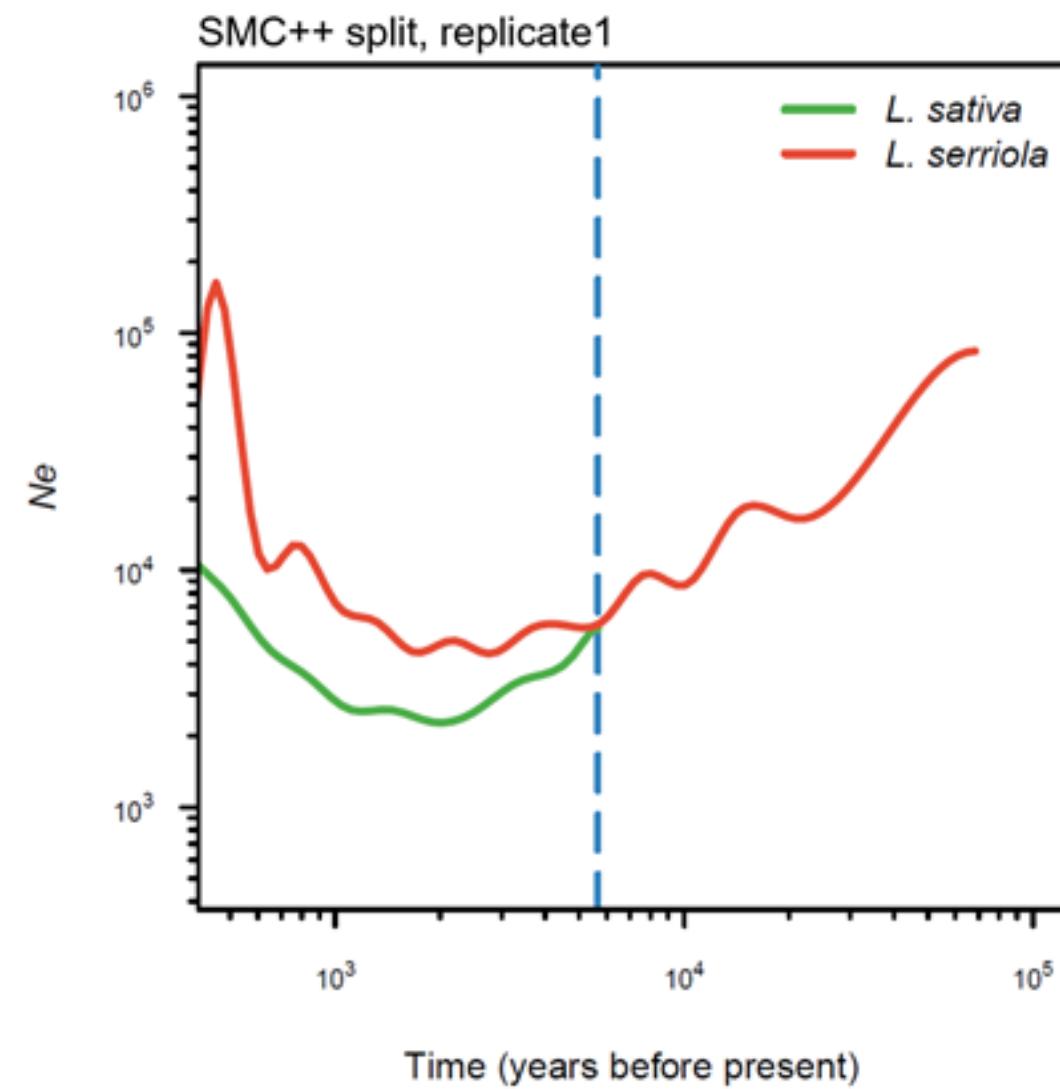
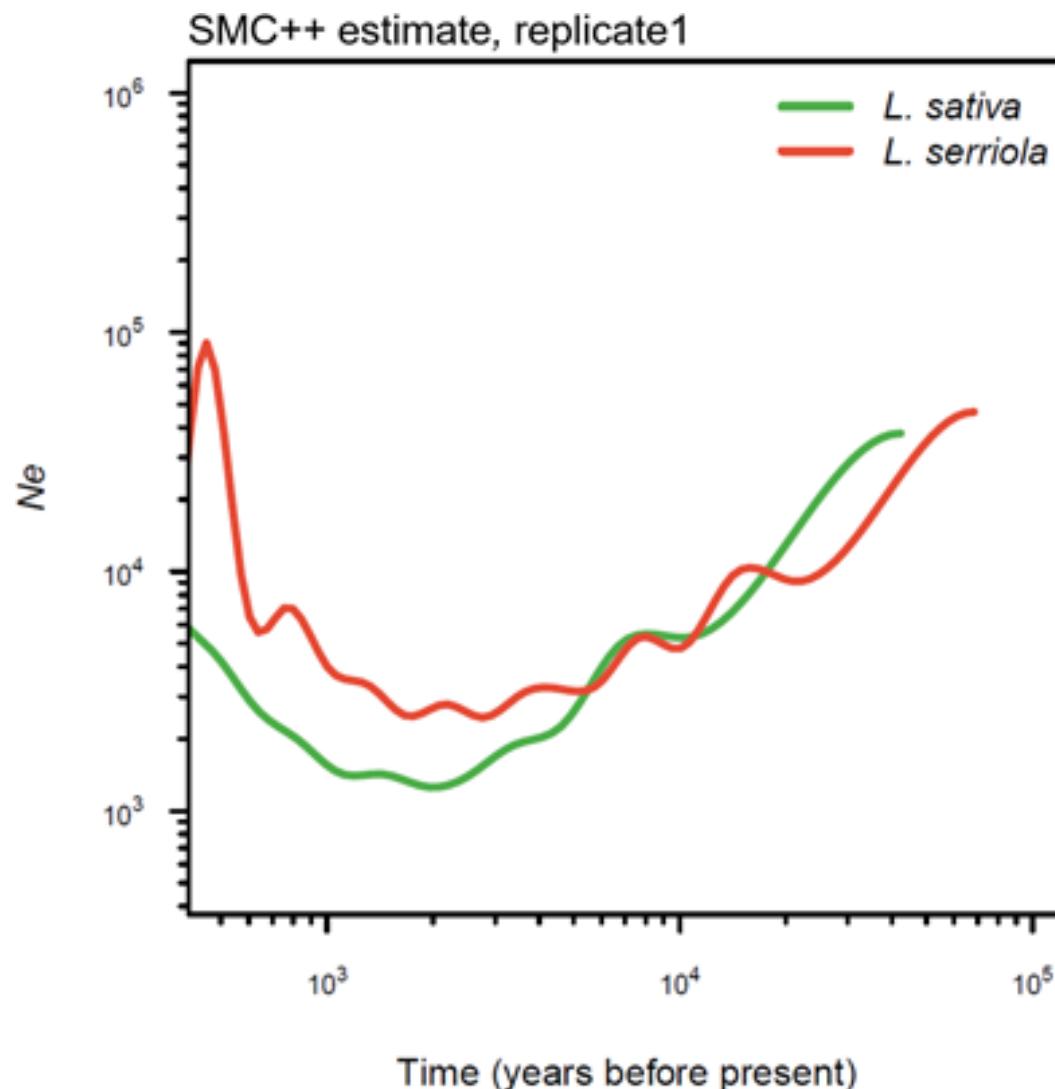
# The highest $\pi$ in the CAU group



# The most abundant singletons in CAU



# Divergence of *L. sativa* and *L. serriola*



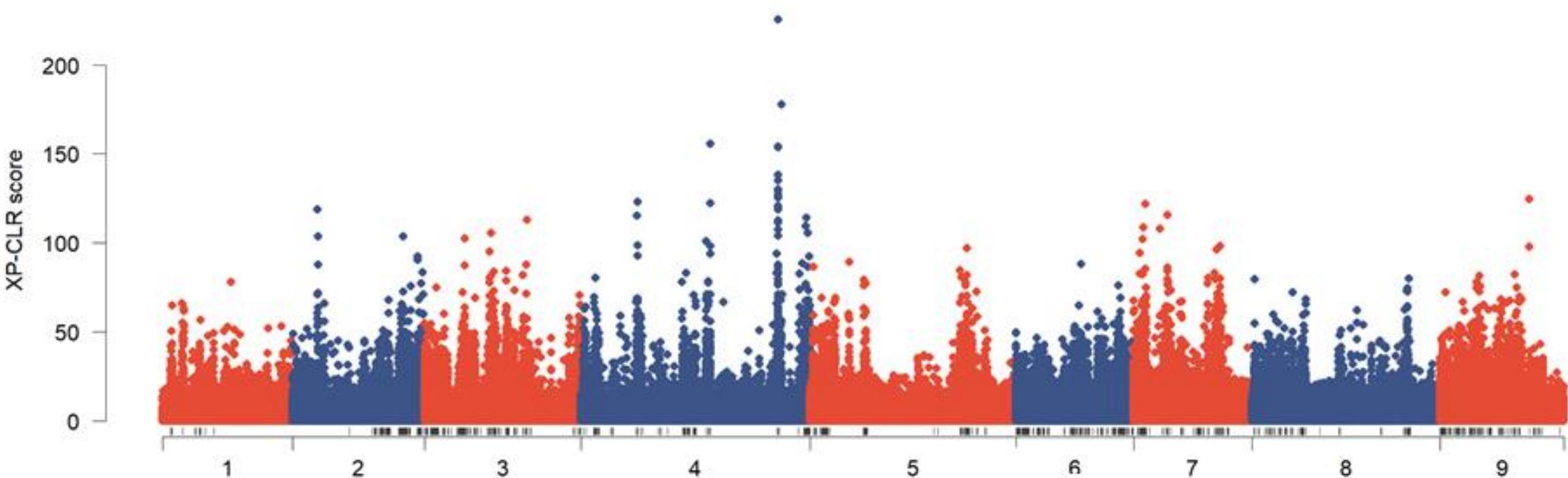
# Answers to the scientific questions

---

- Q3: Where and when was lettuce domesticated?
- A3: Lettuce was domesticated near the Caucasus approximately 6,000 years ago.

# Selective sweeps in the lettuce genome

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

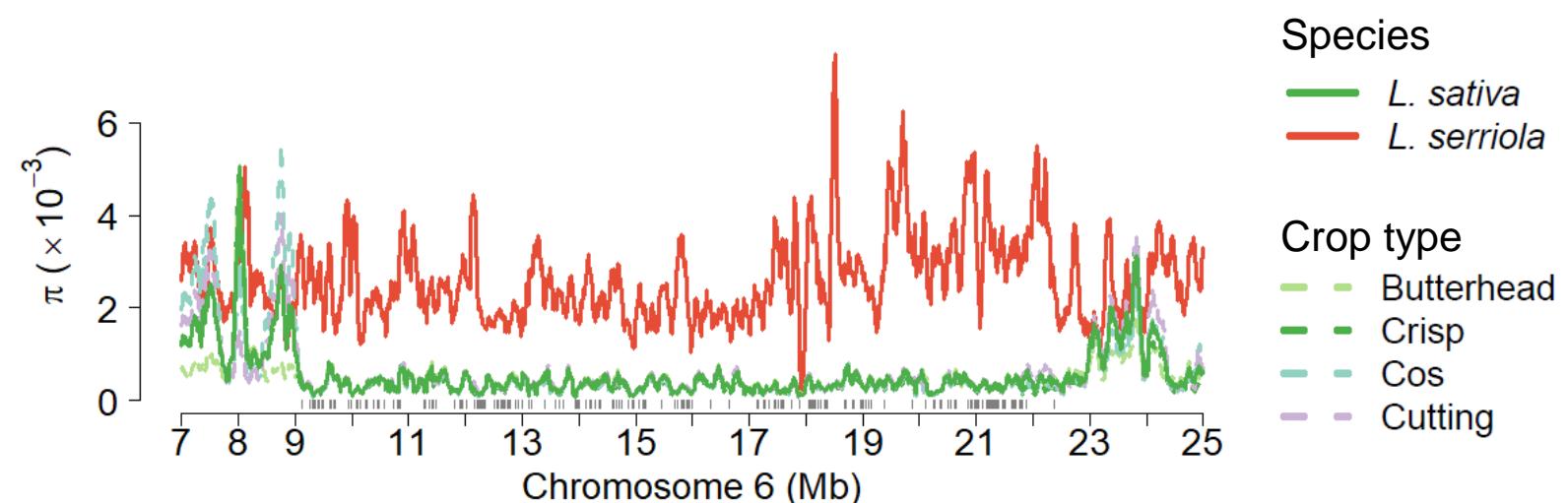
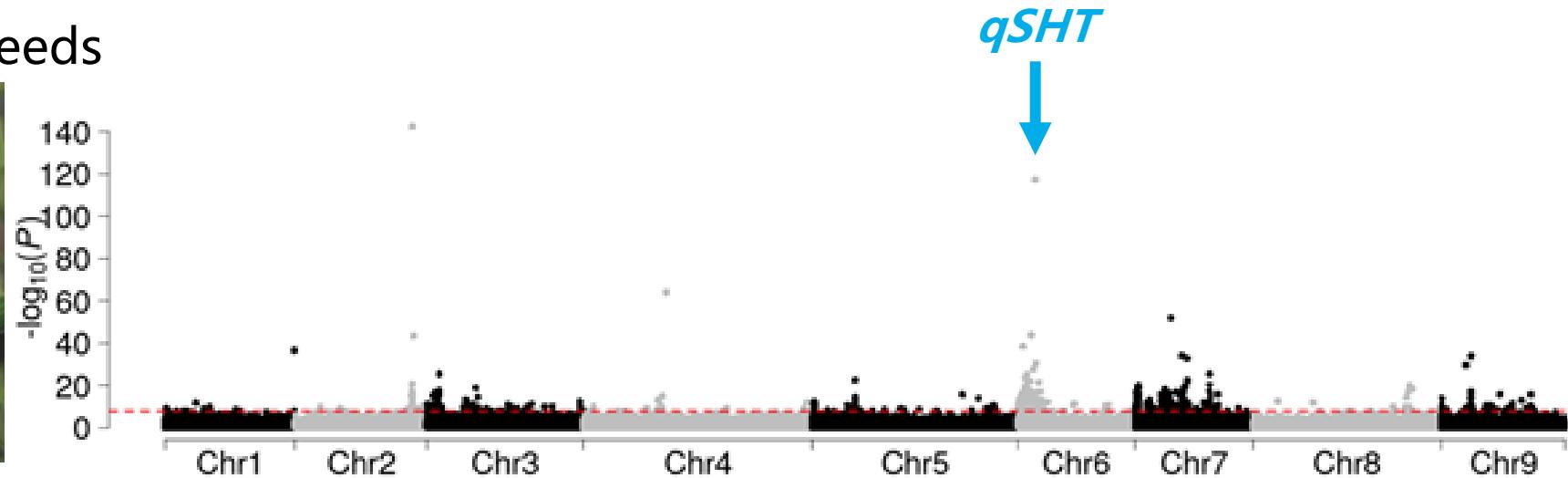


# Domestication trait – seed shattering

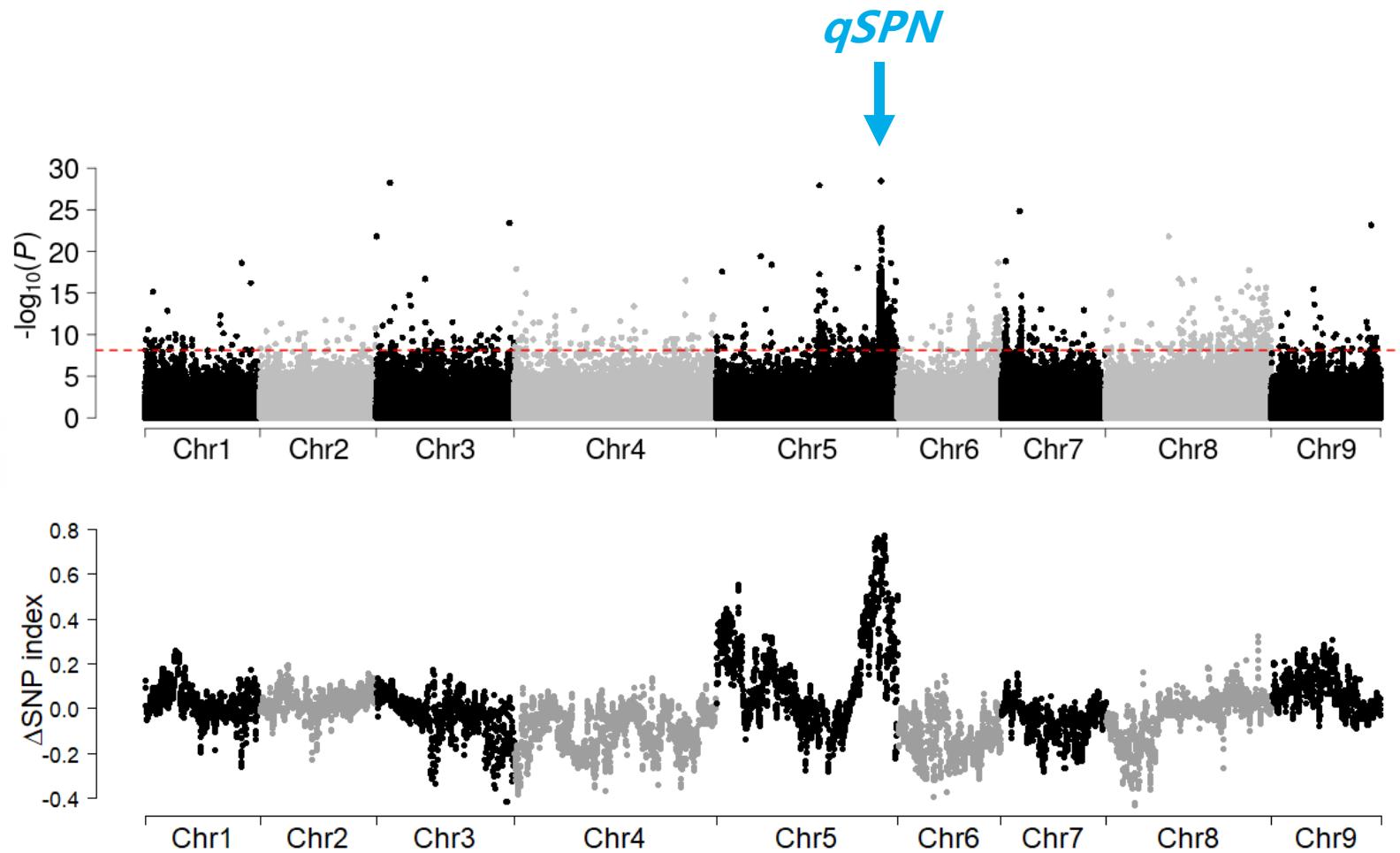
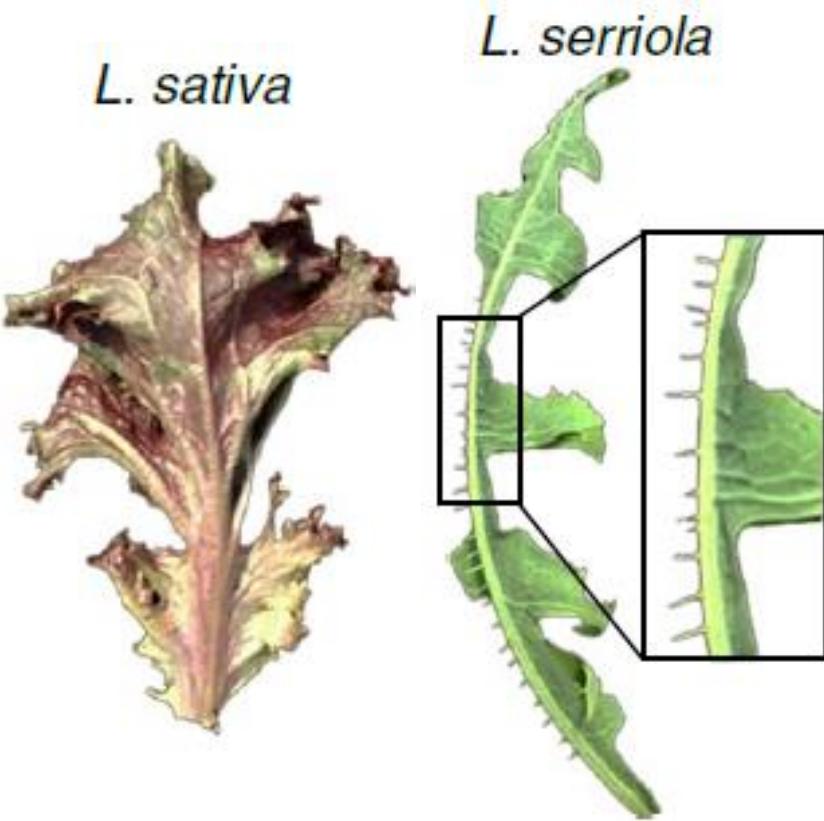
*L. sativa*, non-shattering seeds



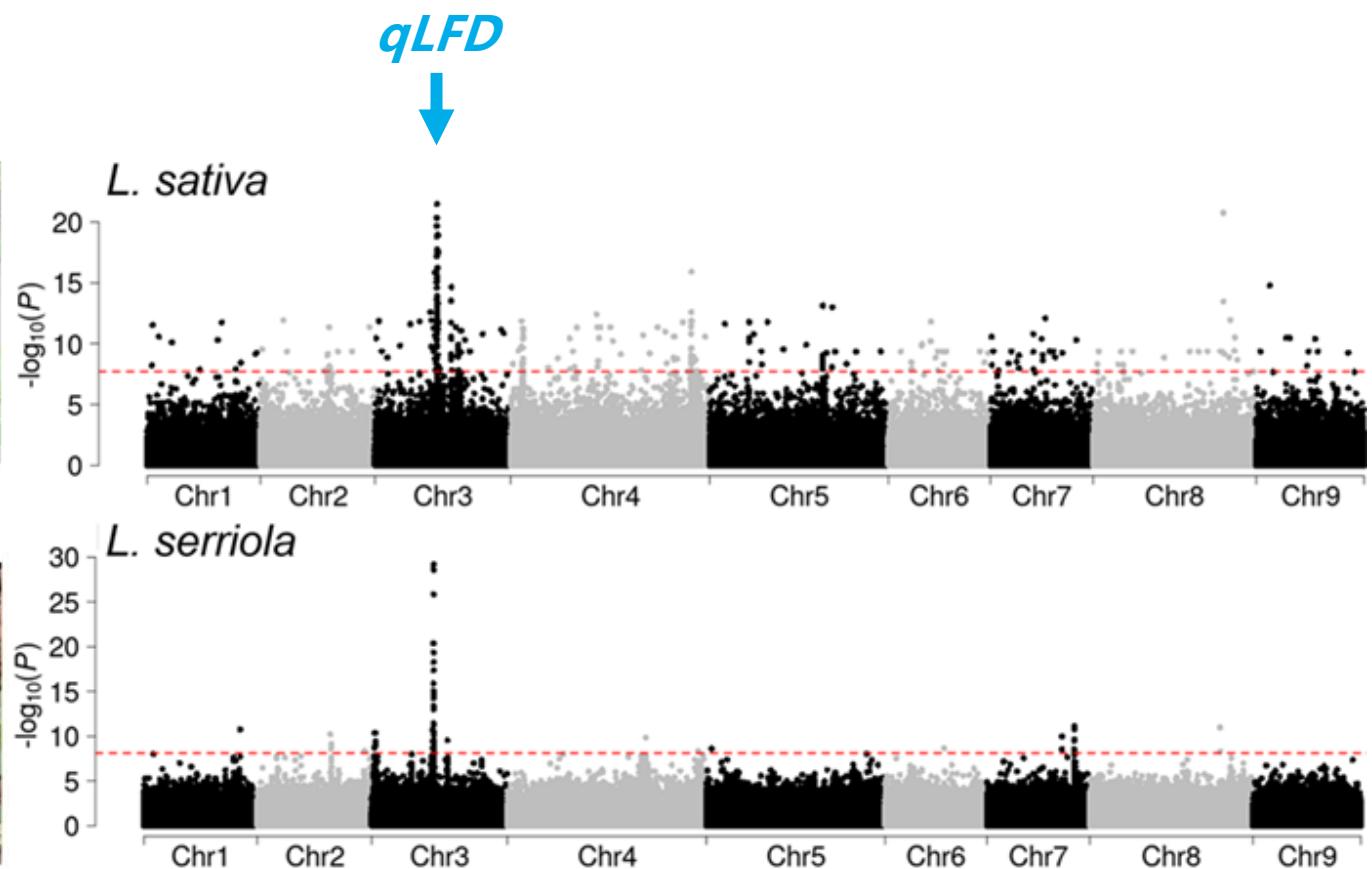
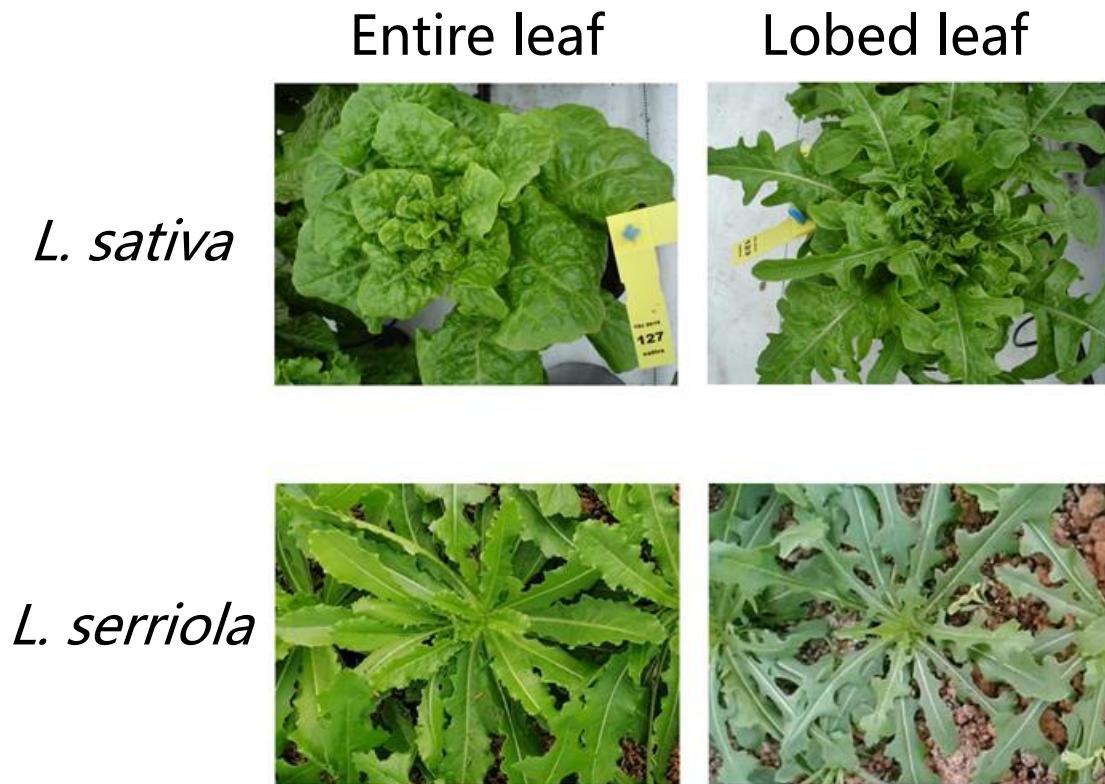
*L. serriola*, shattering seeds



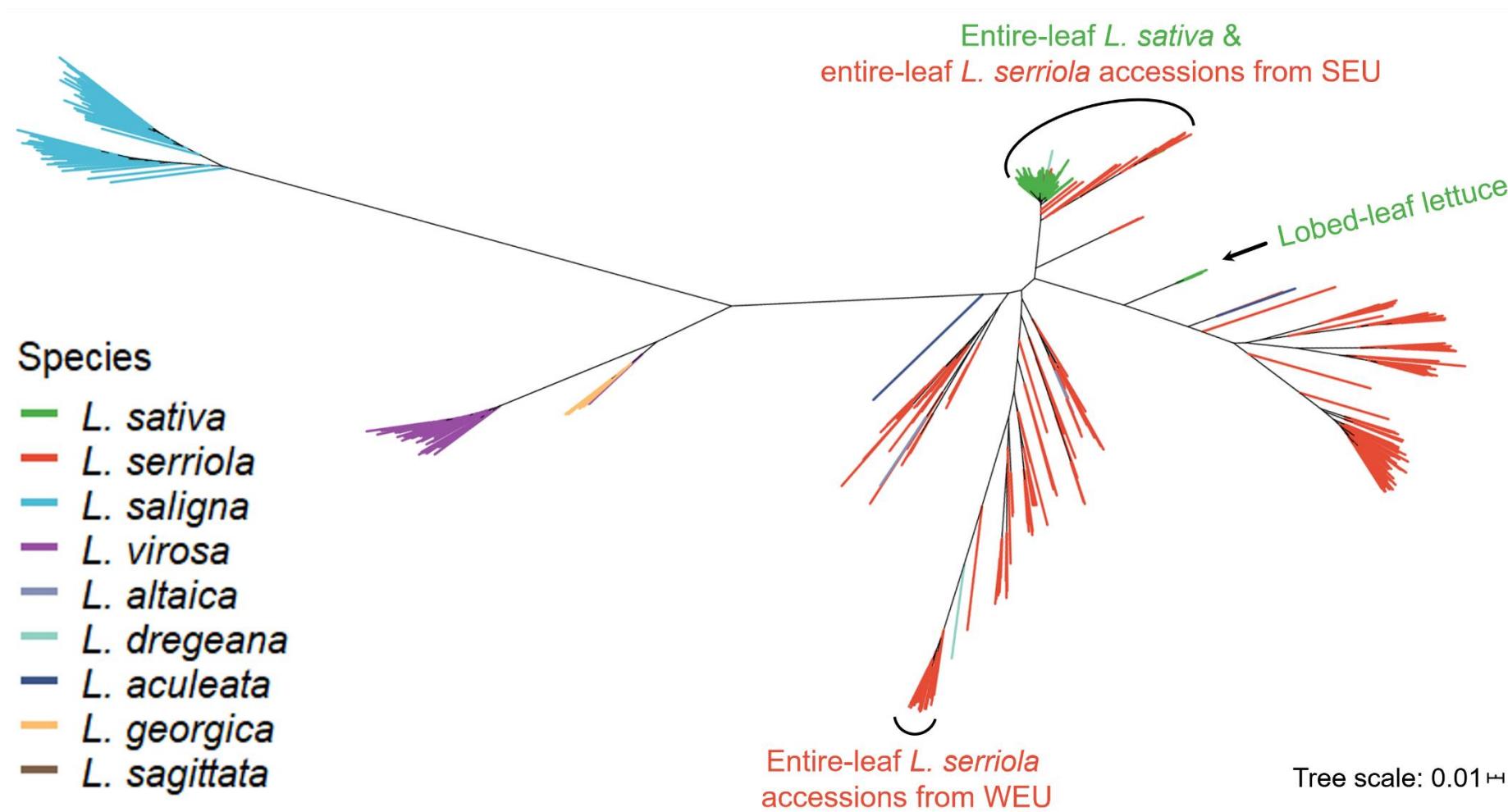
# Domestication trait – leaf spine



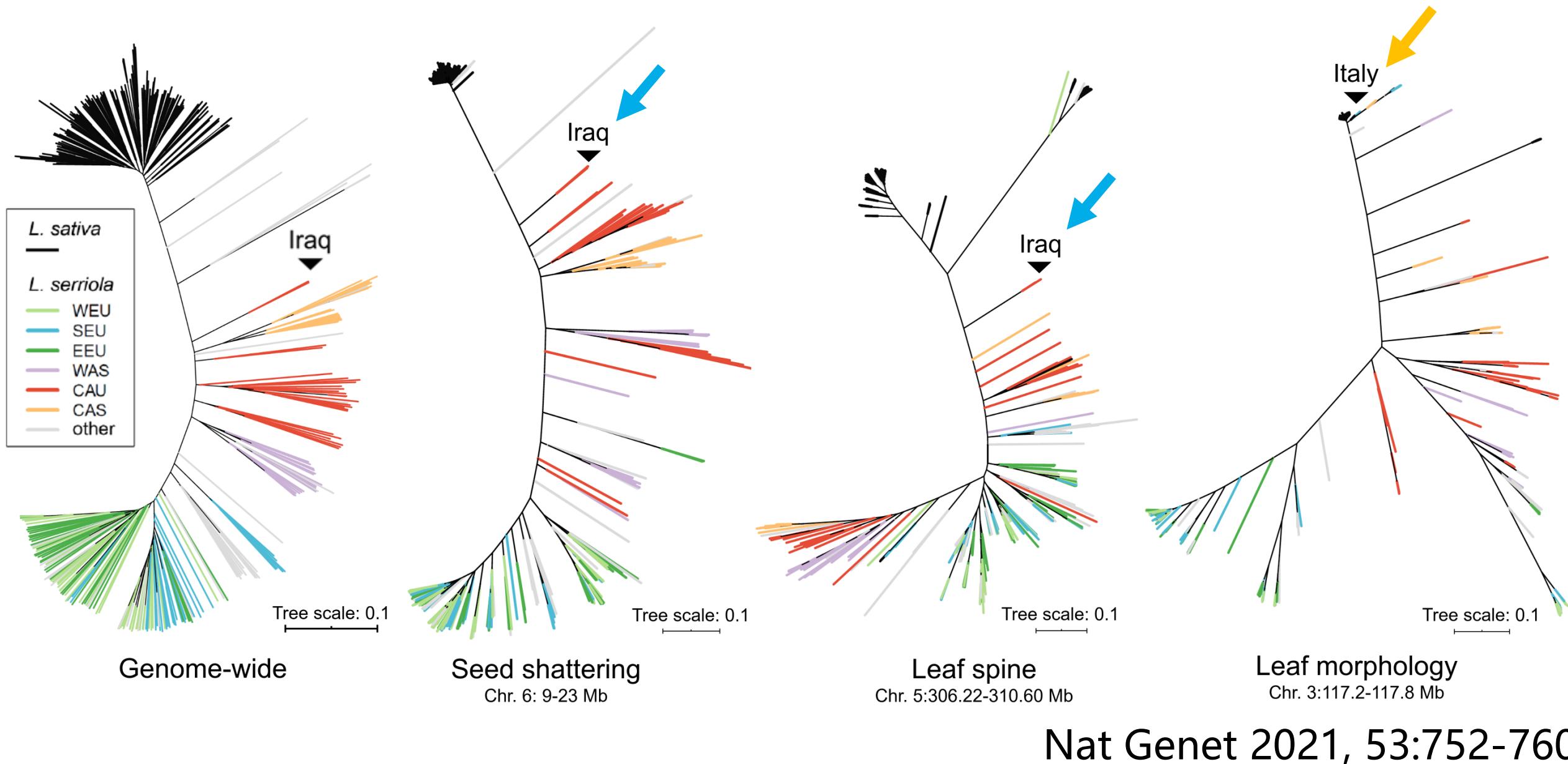
# Domestication trait – leaf morphology



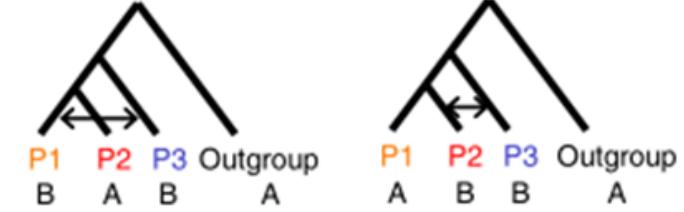
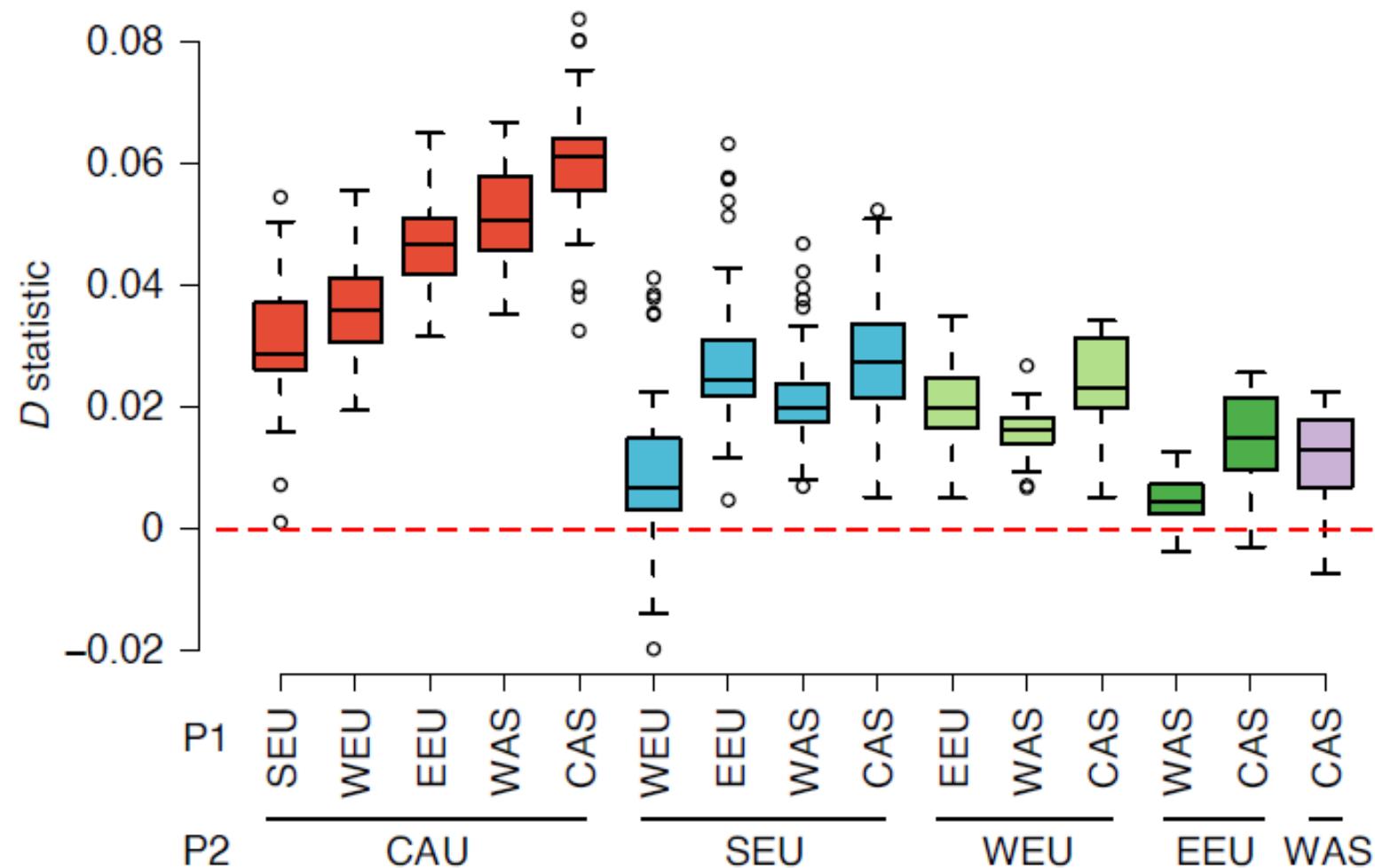
# Origin of the entire-leaf trait



# Origins of domestication traits



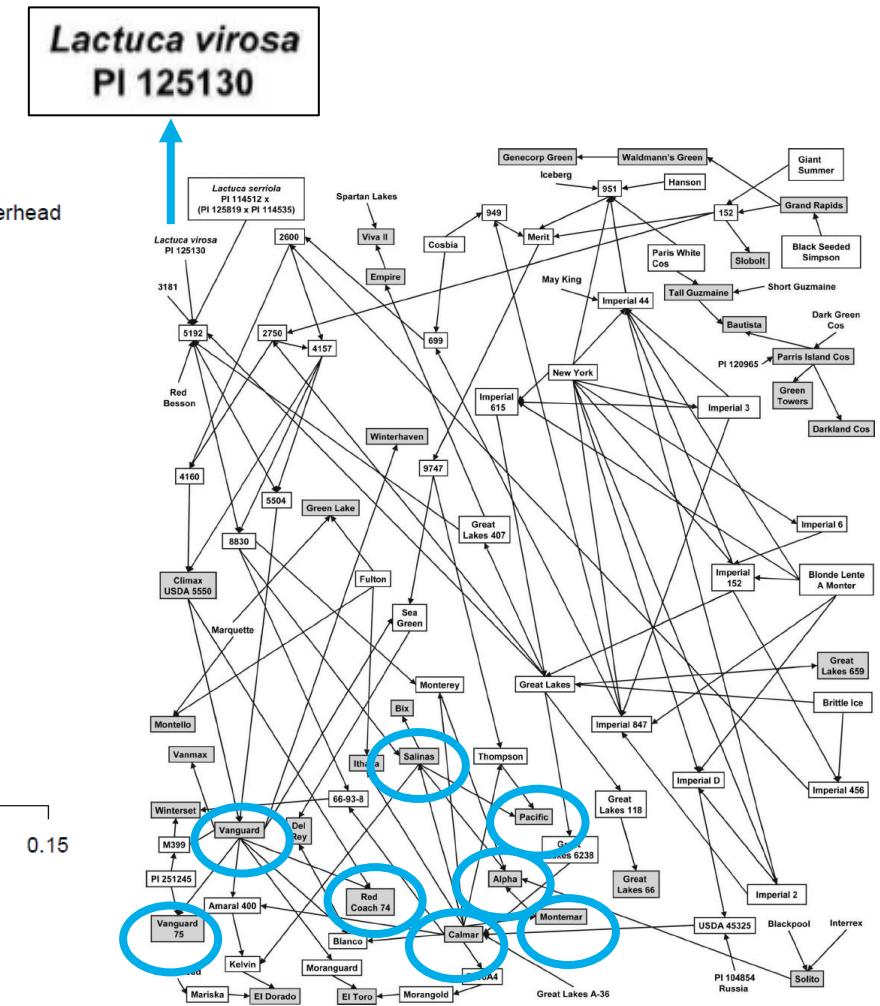
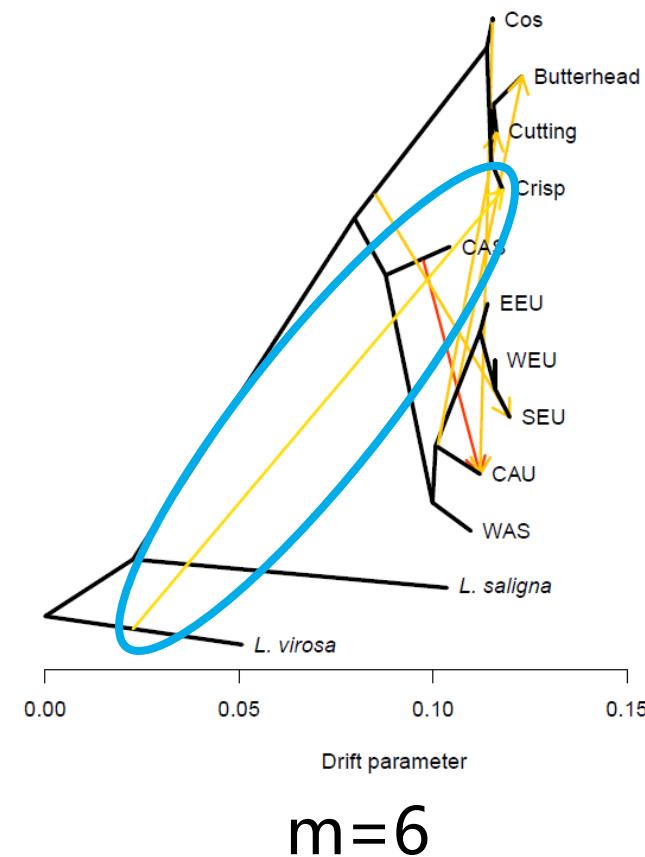
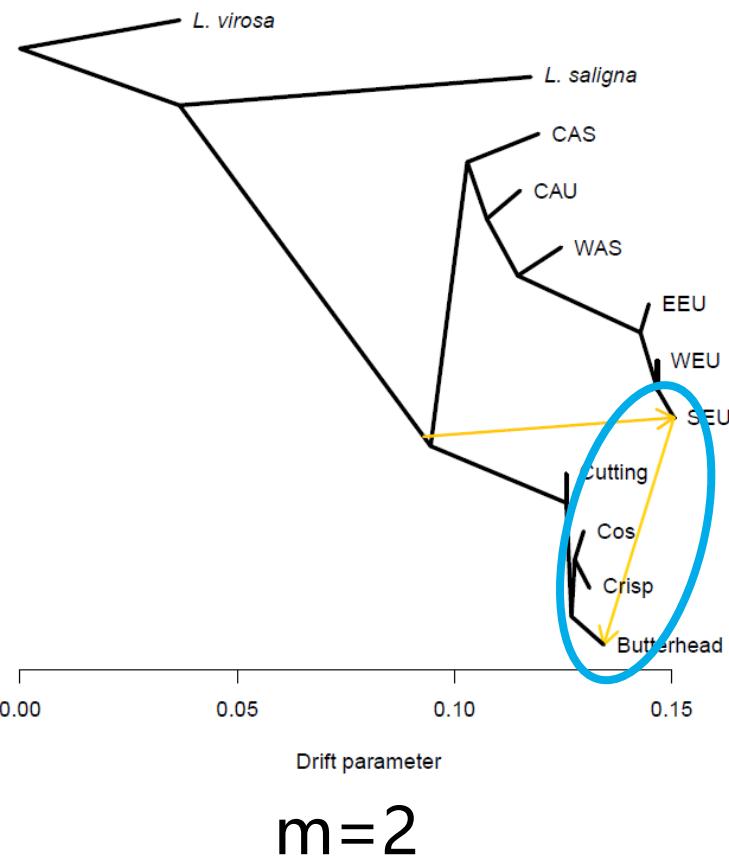
# A close relationship in CAU and SEU



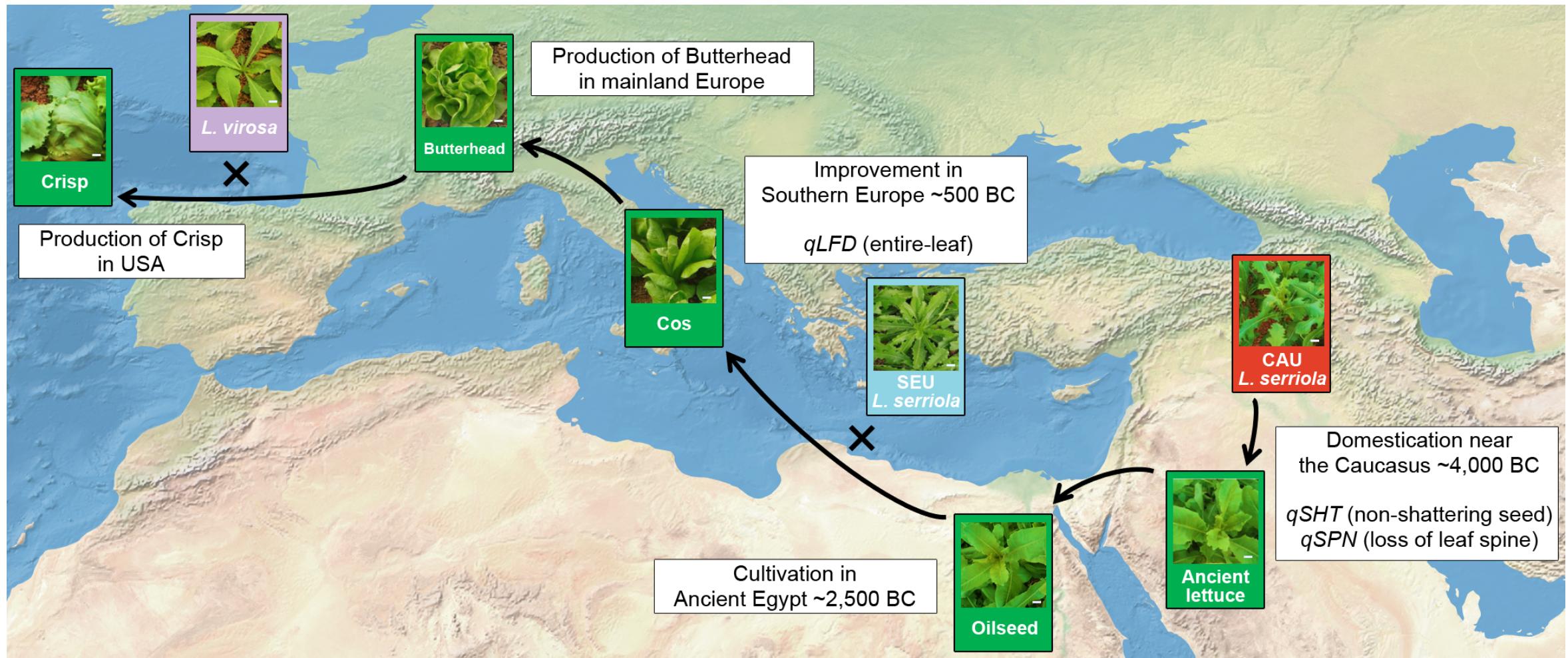
Outgroup: *L. saligna*  
 P3: *L. sativa*  
 P1 & P2: *L. serriola* groups



# Gene flow from SEU and *L. virosa*



# A proposed domestication history



# Answers to the scientific questions

---

- Q4: What are the major events and traits during lettuce domestication and improvement?
- A4: Non-shattering seeds and loss of leaf spines marked lettuce domestication; the entire-leaf trait was introduced from a Southern European wild population.

## GWAS of leaf traits

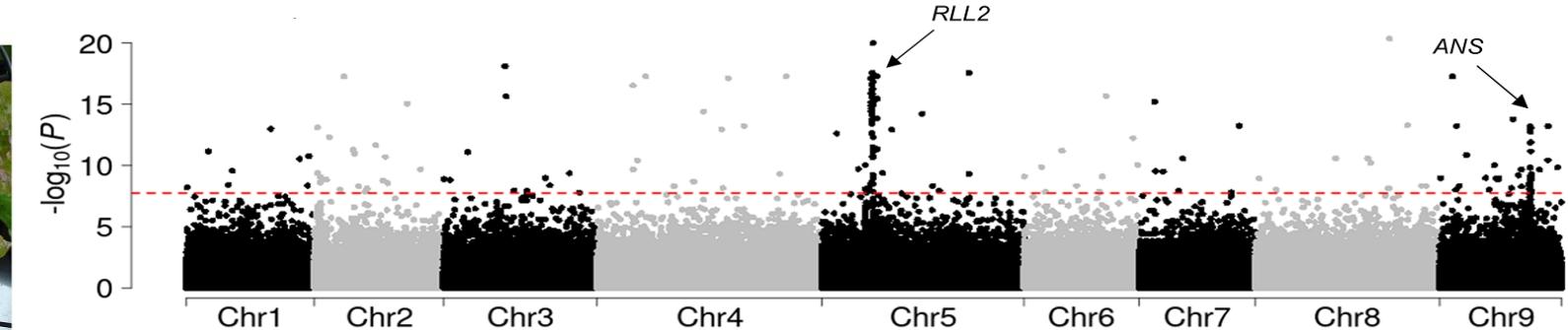
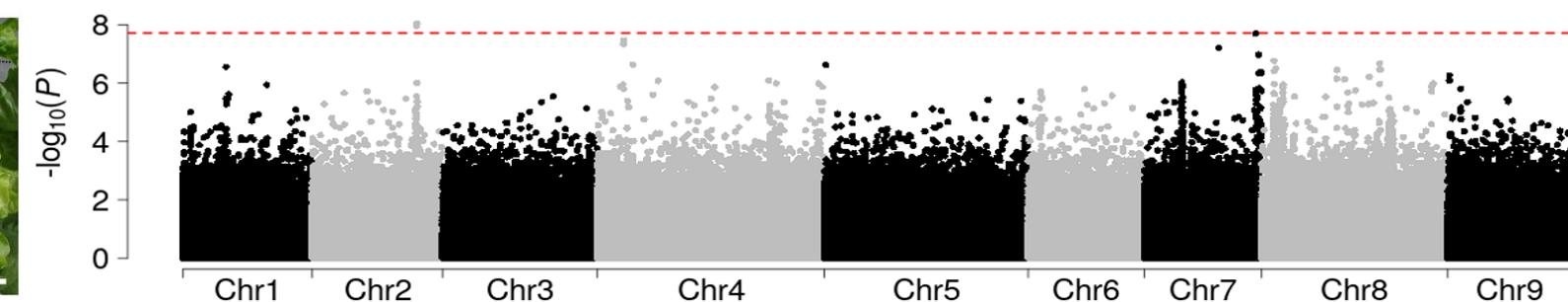
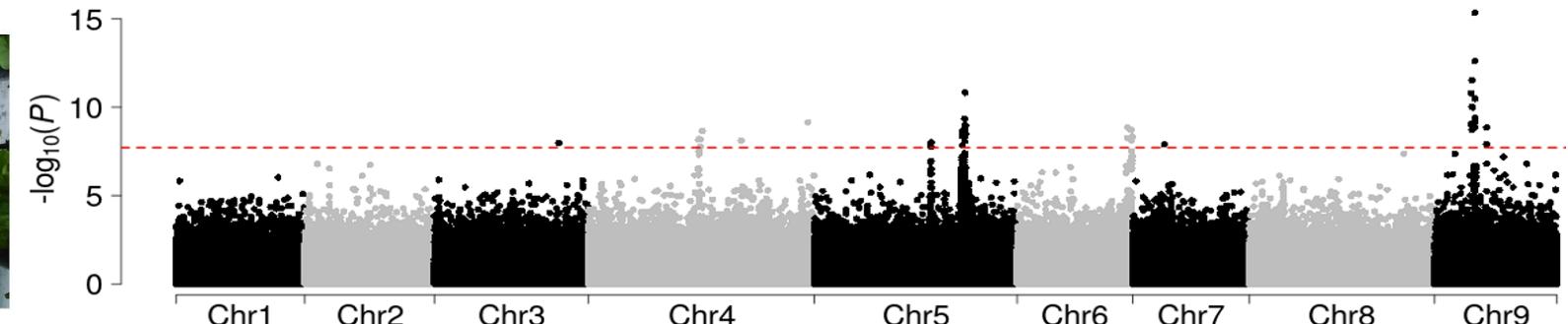
## Leaf margin undulation



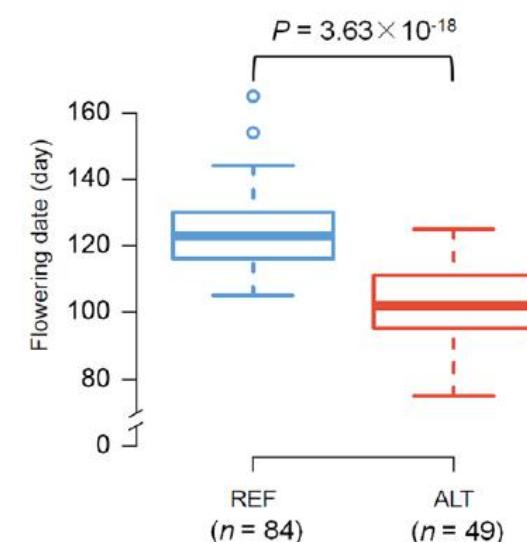
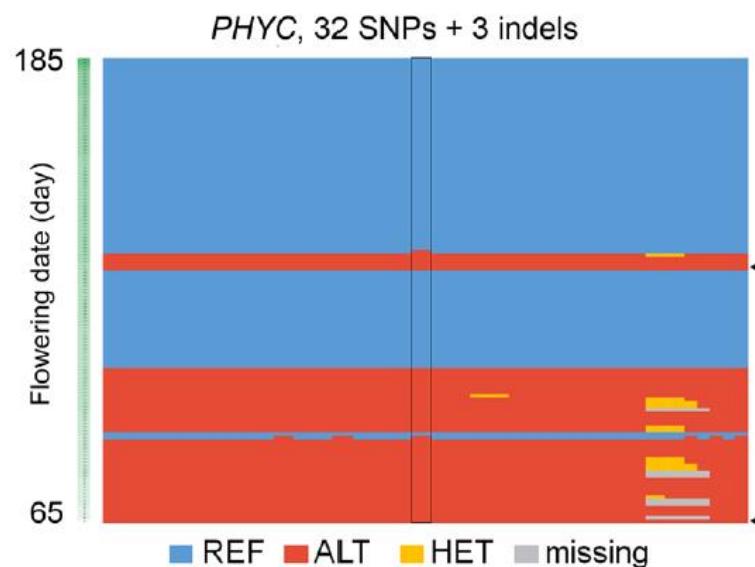
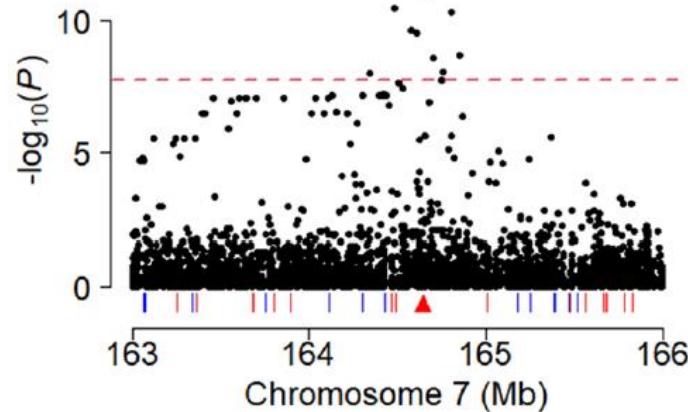
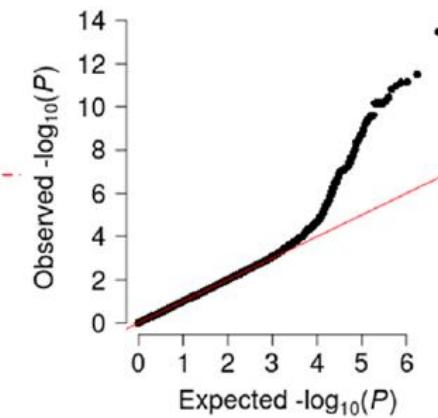
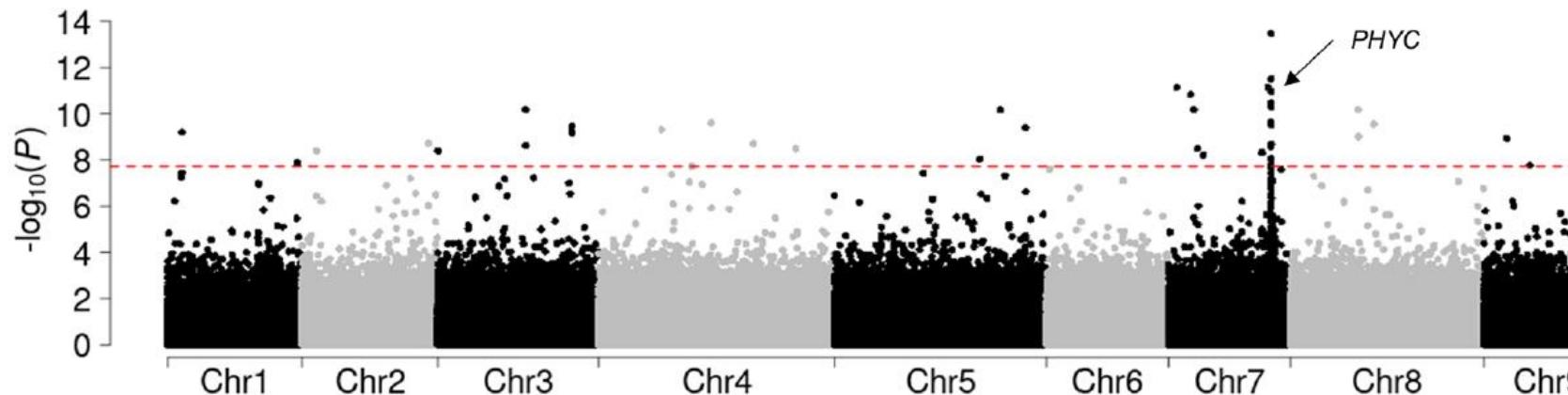
## Leaf venation



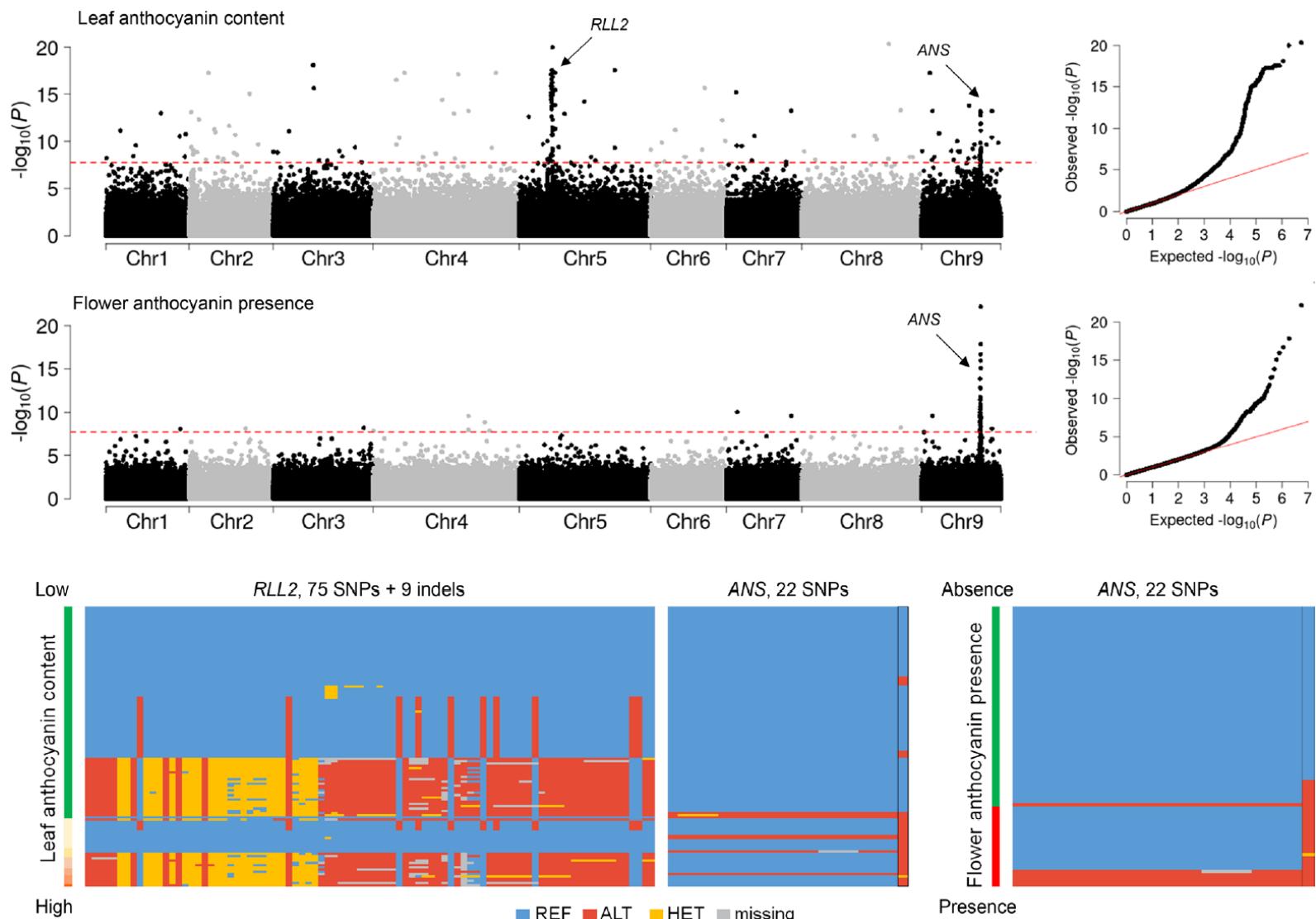
## Leaf anthocyanin contents



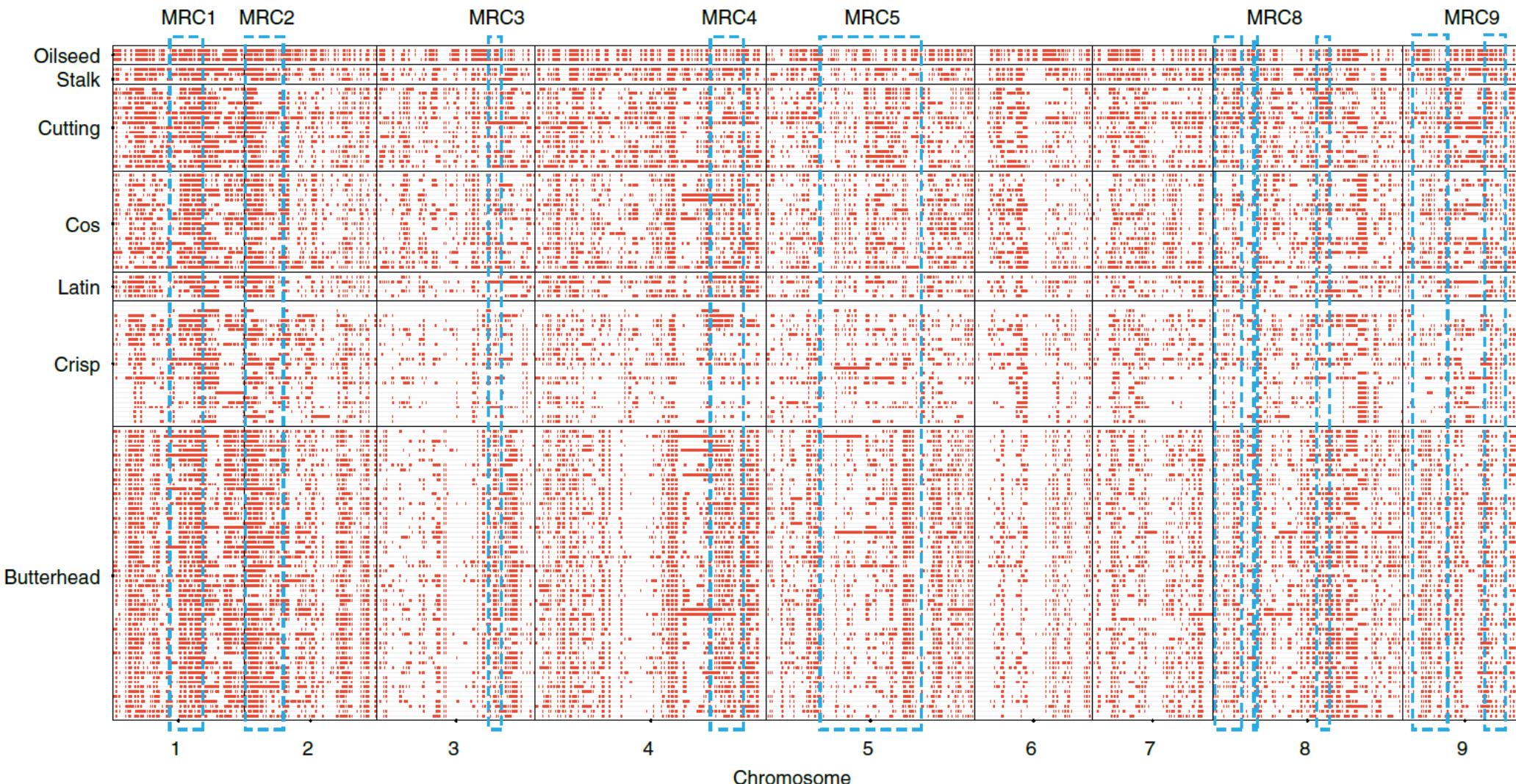
# GWAS of flowering date



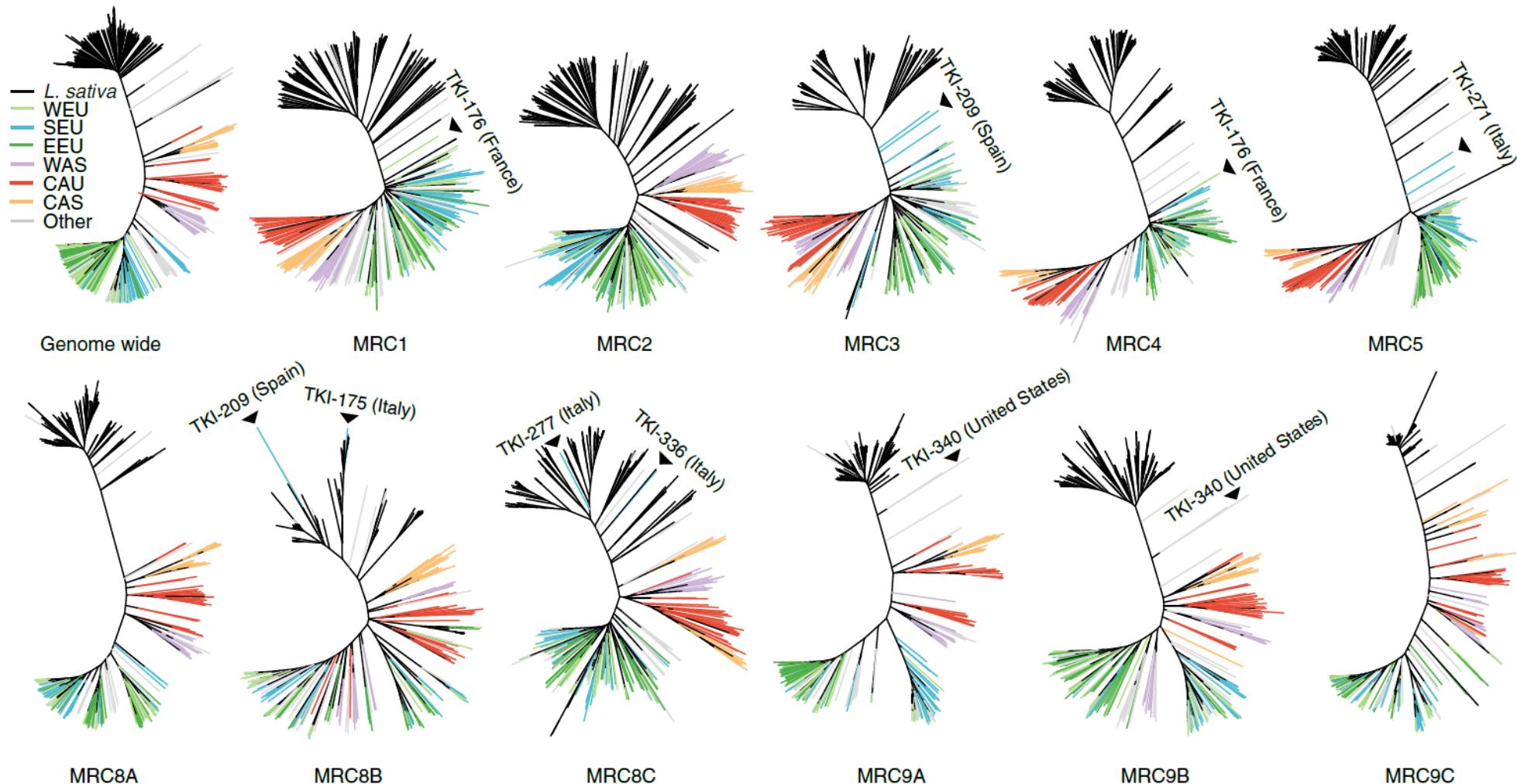
# GWAS of anthocyanin biosynthesis



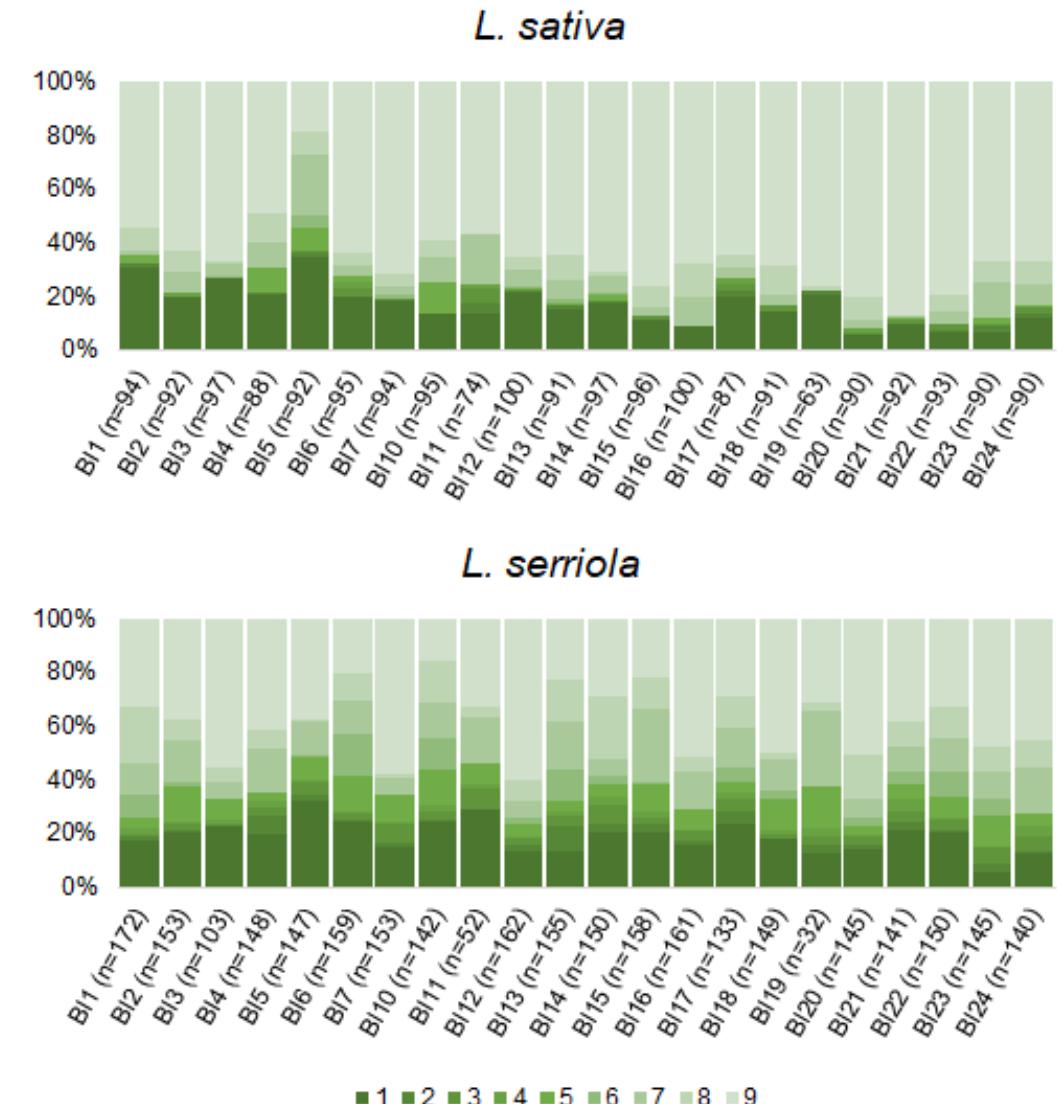
# Wild introgression in resistance clusters



# Origins of the MRCs

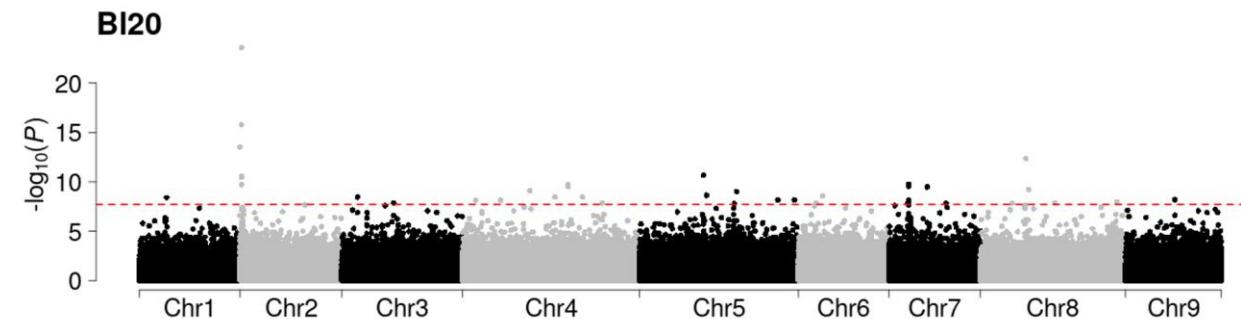
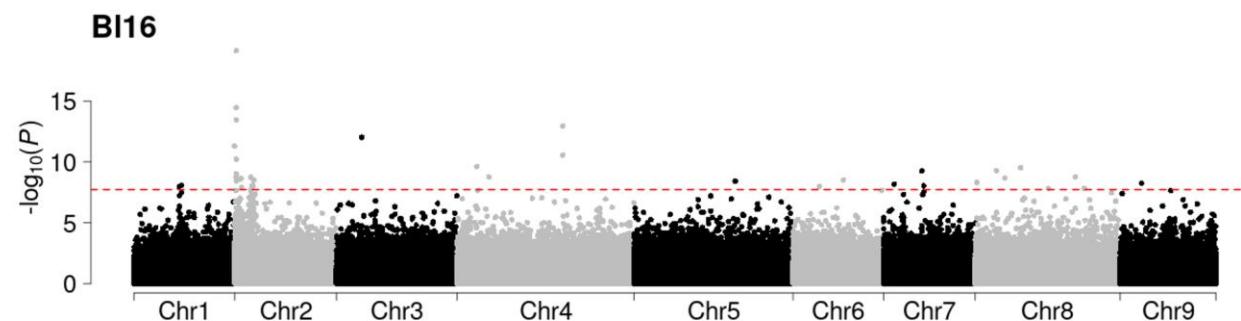


# Bremia resistance in lettuce

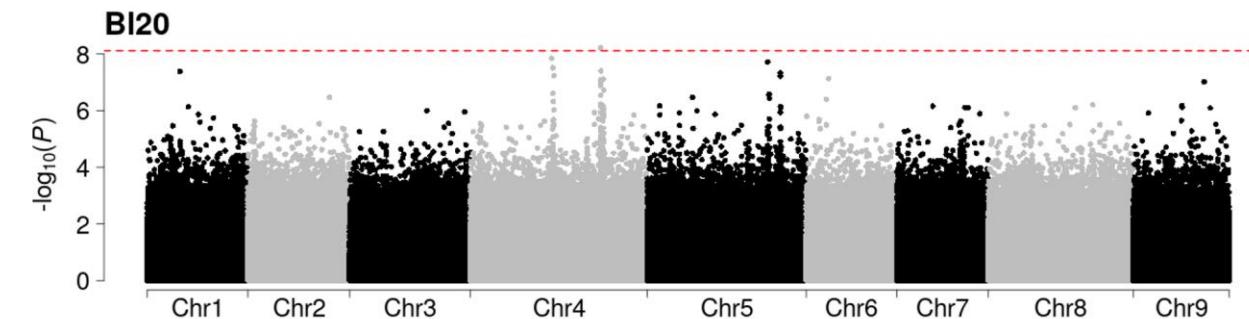
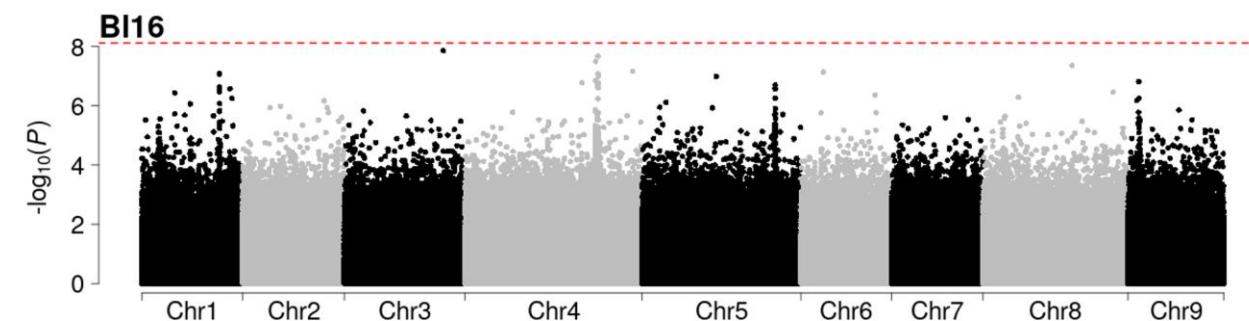


BGI 华大 Multiple *Bremia* resistance loci in *L. serriola*

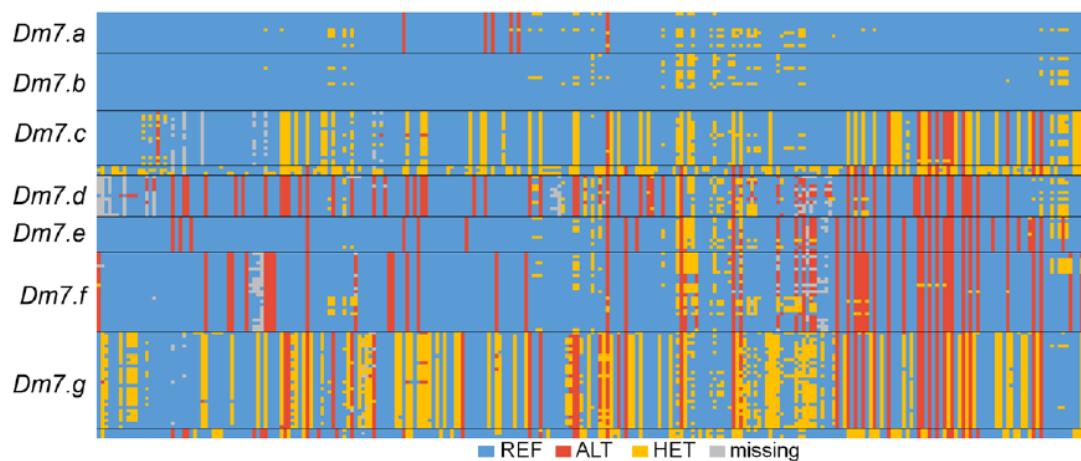
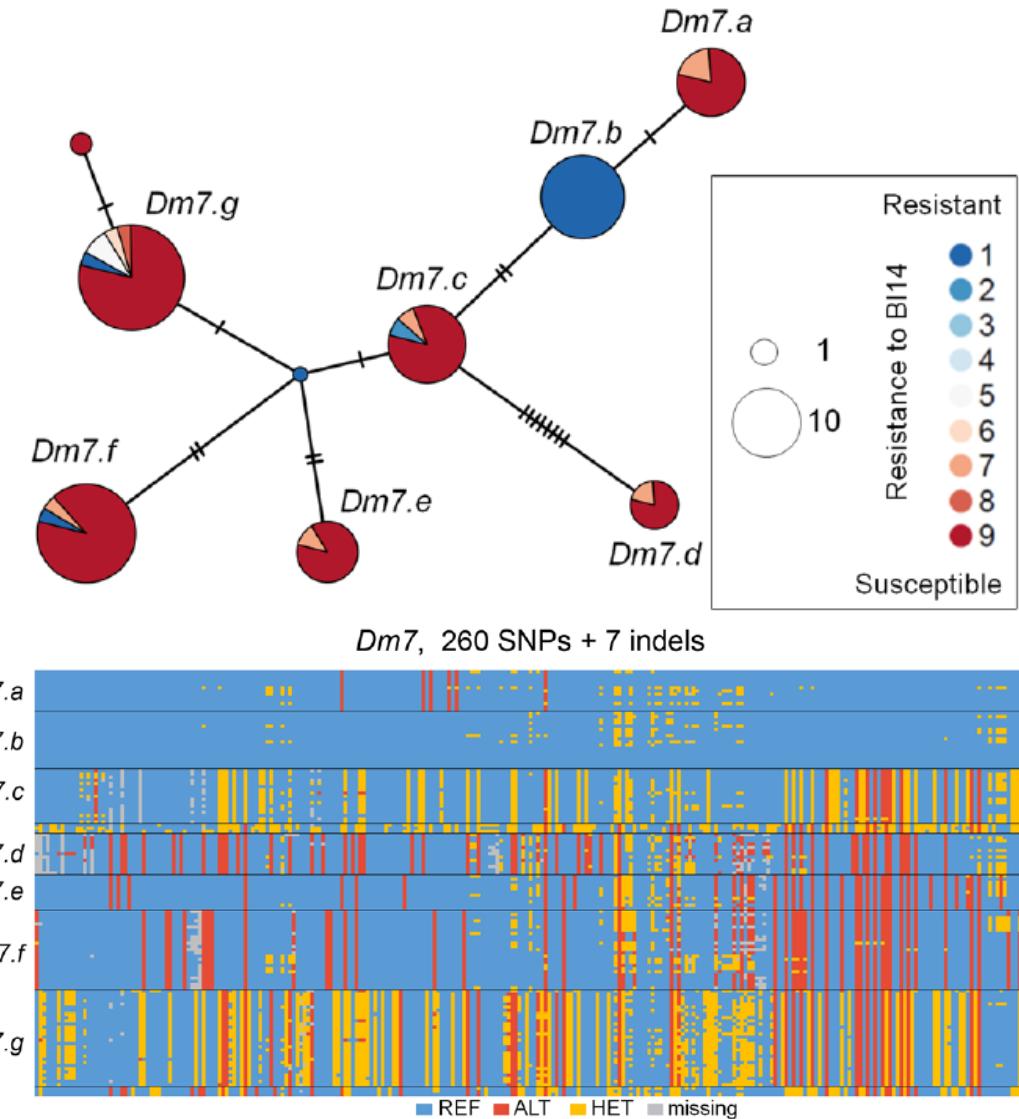
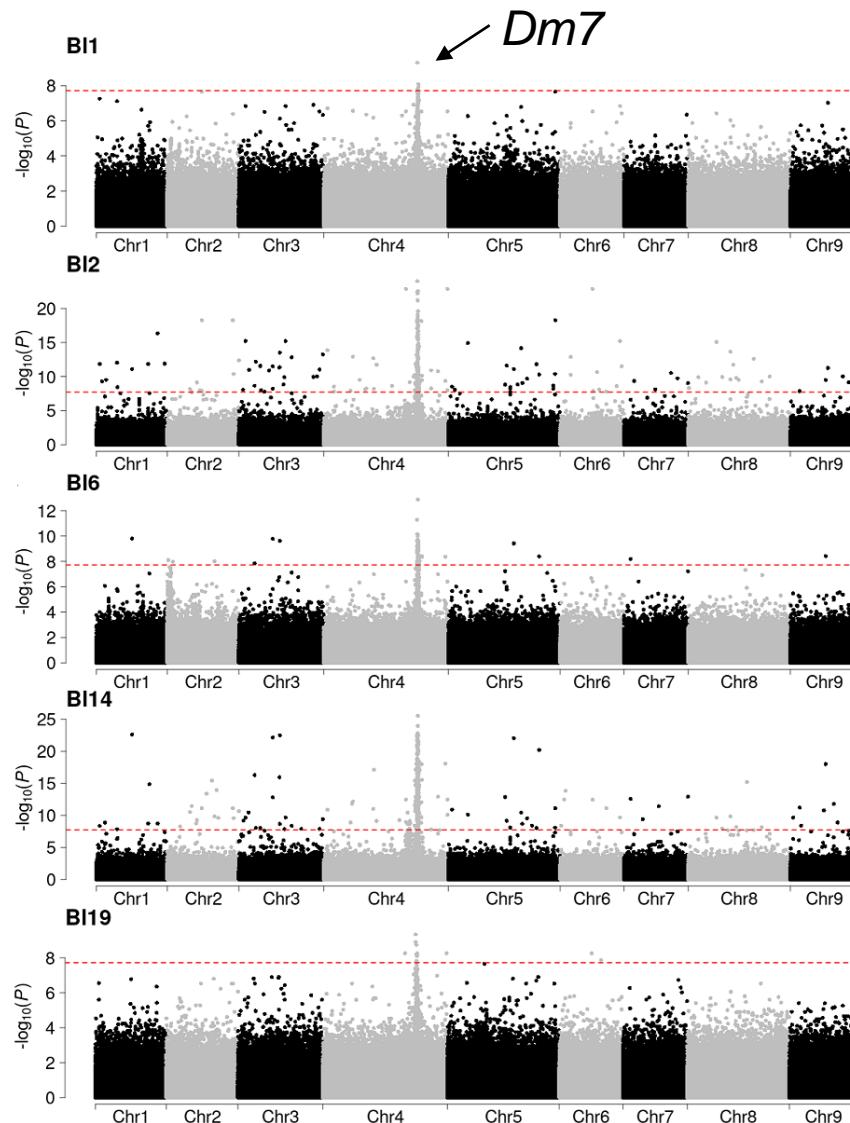
*L. sativa*



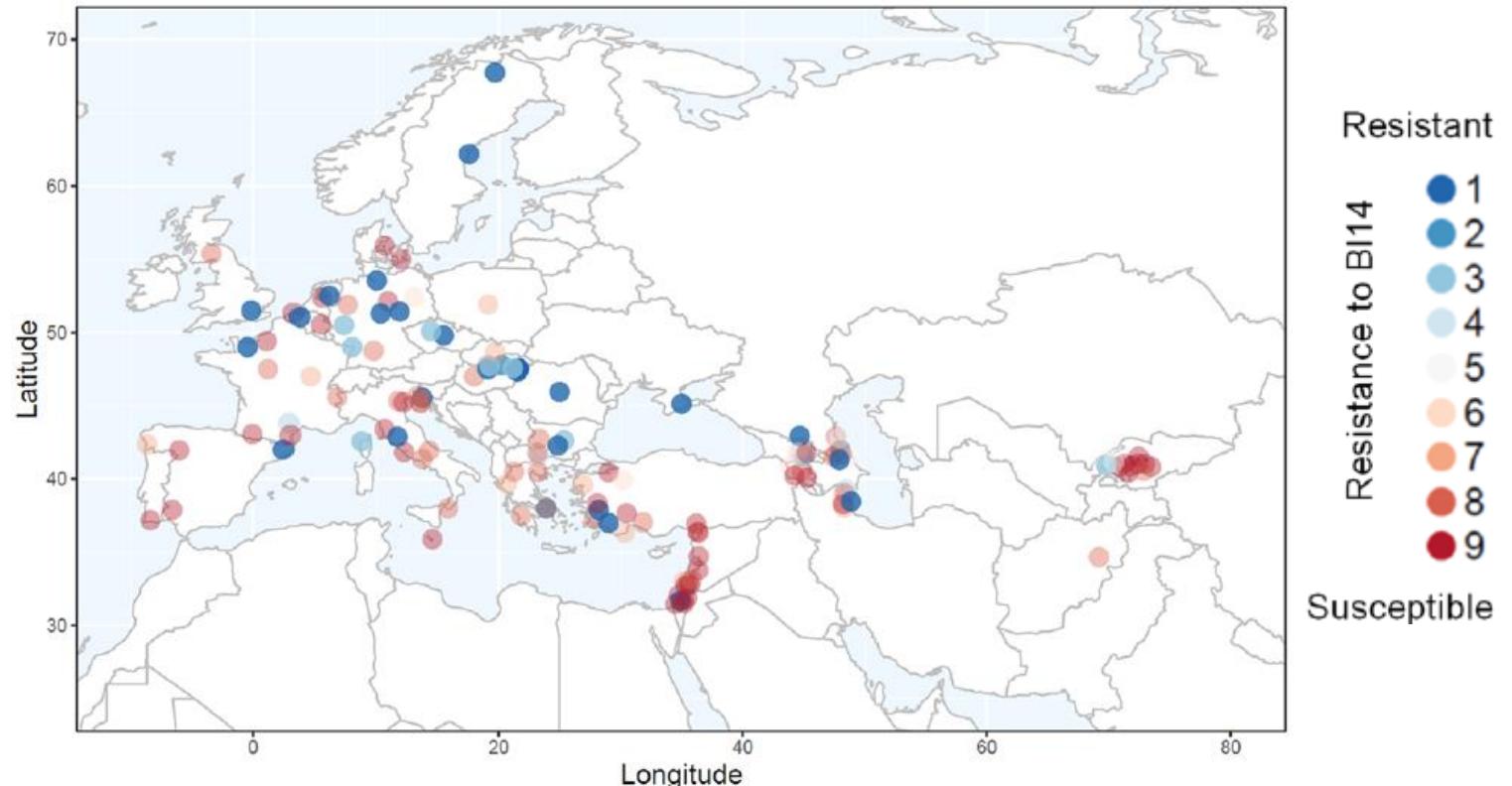
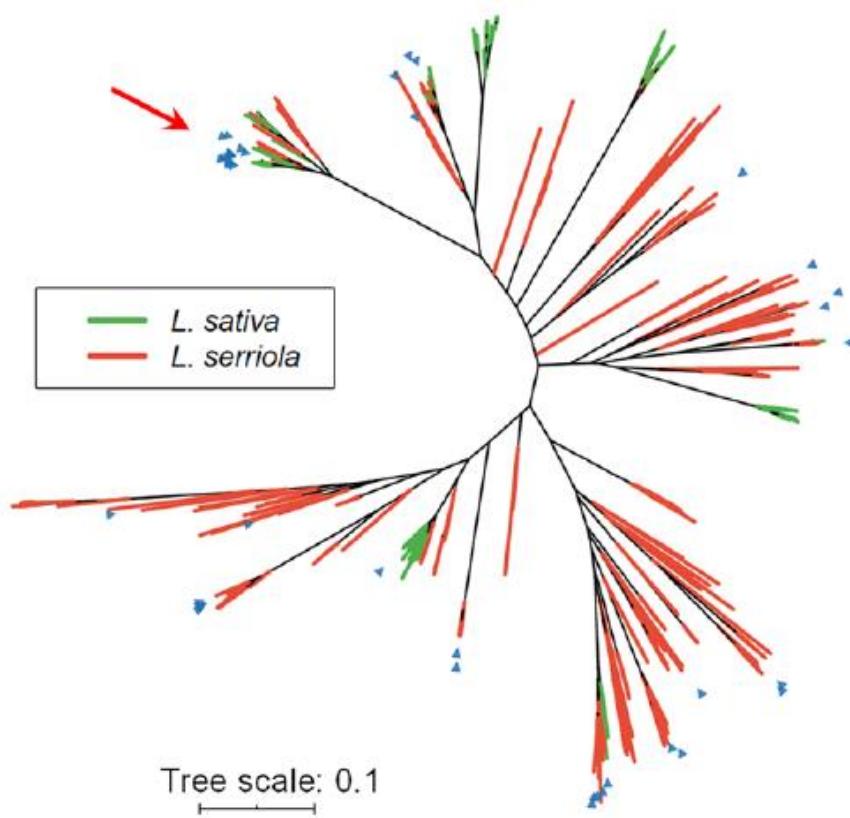
*L. serriola*



# *Dm7.b* confers *Bremia* resistance



# Dm7 alleles in wild population



# Answers to the scientific questions

---

- Q5: What are the genetic determinants for phenotypic differences?
- A5: Nineteen genomic regions were found associated with nine agronomic traits.
- Q6: What are the genetic determinants for *Bremia* resistance?
- A6: Substantial introgressions from *L. serriola* were found in MRCs; more loci are associated with *Bremia* resistance in *L. serriola*.

- Resequencing of 445 accessions identified a comprehensive variation map for cultivated and wild lettuce.
- Population analyses revealed clarified taxonomic issues in GP status.
- Demography identified selective sweeps and major genetic determinants of agronomic traits during lettuce domestication.
- Our study revealed the domestication history of cultivated lettuce.

# Questions?